

THÈSE DE DOCTORAT DE L'UNIVERSITÉ SORBONNE PARIS CITÉ  
PRÉPARÉE À L'UNIVERSITÉ PARIS DIDEROT AU LABORATOIRE IRIF

Ecole Doctorale 386 — Science Mathématiques De Paris Centre

# Tierless Web programming in ML

par

Gabriel RADANNE

Dirigée par

Roberto DI COSMO et Jérôme VUILLON

Thèse soutenue publiquement le 14 Novembre 2017 devant le jury constitué de

Roberto DI COSMO  
Jérôme VUILLON  
Koen CLAESSEN  
Jacques GARRIGUE  
Coen DE ROOVER  
Xavier LEROY  
Manuel SERRANO  
Jeremy YALLOP

Directeur de thèse  
Co-Directeur de thèse  
Rapporteur  
Rapporteur  
Examineur  
Examineur  
Examineur  
Examineur



**Abstract** ELIOM is a dialect of OCAML for Web programming in which server and client pieces of code can be mixed in the same file using syntactic annotations. This allows to build a whole application as a single distributed program, in which it is possible to define in a composable way reusable widgets with both server and client behaviors.

ELIOM is type-safe, as it ensures that communications are well-behaved through novel language constructs that match the specificity of Web programming. ELIOM is also efficient, it provides static slicing which separates client and server parts at compile time and avoids back-and-forth communications between the client and the server. Finally, ELIOM supports modularity and encapsulation thanks to an extension of the OCAML module system featuring tierless annotations that specify whether some definitions should be on the server, on the client, or both.

This thesis presents the design, the formalization and the implementation of the ELIOM language.

**Résumé** ELIOM est un dialecte d'OCAML pour la programmation Web qui permet, à l'aide d'annotations syntaxiques, de déclarer code client et code serveur dans un même fichier. Ceci permet de construire une application complète comme un unique programme distribué dans lequel il est possible de définir des widgets aisément composables avec des comportements à la fois client et serveur.

ELIOM assure un bon comportement des communications grâce à un système de type et de nouvelles constructions adaptés à la programmation Web. De plus, ELIOM est efficace : un découpage statique sépare les parties client et serveur durant la compilation et évite de trop nombreuses communications entre le client et le serveur. Enfin, ELIOM supporte la modularité et l'encapsulation grâce à une extension du système de module d'OCAML permettant l'ajout d'annotations indiquant si une définition est présente sur le serveur, le client, ou les deux.

Cette thèse présente la conception, la formalisation et l'implémentation du langage ELIOM.



# Remerciements

Ceux qui me connaissent savent que je ne suis pas particulièrement expansif, je serai donc bref: je remercie ma famille, mes amis, mes collègues et l'univers.



## Remerciements, Deuxième essai

Je voudrais tout d’abord remercier mes deux directeurs de thèse, Roberto et Jérôme, qui ont formé une paire remarquablement complémentaire et m’ont beaucoup aidé dans tous les aspects de ma recherche, scientifiques ou non.

Je remercie l’équipe OCSIGEN et les membres de l’Irill pour leur compagnie, conseils et pour avoir pris le temps d’écouter mes idées farfelues sur les prochaines fonctionnalités révolutionnaires qui vont faire d’ELIOM le langage de programmation Web le plus mieux au monde!

Je remercie Jacques, pour cette petite remarque pendant ICFP 2015: “Tu devrais parler avec Oleg, ce que tu fais ressemble à MetaOCaml”. Merci également à Oleg, pour ses remarques et ses encouragements.

Avant de commencer cette thèse, j’ai eu la chance de faire de nombreux stages. Je voudrais remercier tous les gens avec qui j’ai pu interagir à Édimbourg, Gotheborg et Grenoble. Ces stages m’ont beaucoup appris et surtout, je m’y suis énormément amusé! Merci également aux membres d’OCamlabs qui m’ont donné l’occasion de m’échapper de ma thèse pendant 3 mois pour faire joujou avec Mirage et boire du cidre dans les pubs de Cambridge.

En parlant de cidre, merci à tous les potes de promo de Rennes: Benoit, Simon, Nicolas, Clément, Gael et tous les autres. Merci également à David pour son encadrement pendant mes études et pour m’avoir orienté vers d’excellents stages.

Merci aux membres de l’IRIF avec qui j’ai pu interagir, notamment Vincent, Ralf, Jean-Baptiste et Michele avec qui j’ai eu le plaisir de faire mes enseignements et Alexis et Delia, pour leurs discussions et coups de pouce occasionnels. Je remercie particulièrement Odile, sans qui, soyons honnête, je n’aurais jamais réussi à faire quoi que ce soit. Merci à tous les doctorants de PPS, en particulier la folle équipe du bureau 3026: Clément, Yann, Rémi, Pierre, Amina, Charles, Léo, mon quotidien aurait été nettement plus terne sans vous.

Enfin, je remercie ma famille, pour à peu près tout.





# Contents

<b>1</b>	<b>Introduction</b>	<b>13</b>
1.1	On Web programming languages . . . . .	15
1.1.1	The client-server dichotomy . . . . .	15
1.1.2	Tierless programming languages . . . . .	16
1.1.3	Functional programming . . . . .	16
1.1.4	Static type systems . . . . .	17
1.1.5	Modularity and encapsulation . . . . .	17
1.2	ELIOM . . . . .	19
1.2.1	The OCSIGEN project . . . . .	19
1.2.2	Principles of the ELIOM language . . . . .	19
1.3	Plan . . . . .	20
1.4	Contributions . . . . .	21
<b>2</b>	<b>Programming with Eliom</b>	<b>23</b>
2.1	Core concepts . . . . .	23
2.1.1	Sections . . . . .	23
2.1.2	Client fragments . . . . .	24
2.1.3	Injectons . . . . .	24
2.2	Client-server behaviors . . . . .	24
2.2.1	Introducing side effects . . . . .	25
2.2.2	Communication and execution scheme . . . . .	26
2.3	Heterogeneous datatypes . . . . .	26
2.3.1	Remote procedure calls . . . . .	27
2.3.2	Converters . . . . .	28
2.3.3	Client-server reactive broadcasts . . . . .	29
2.4	Modules and signatures . . . . .	30
2.4.1	A primer on OCAML modules . . . . .	30
2.4.2	Modules and locations . . . . .	32
2.4.3	OCAML integration . . . . .	32
2.4.4	Heterogeneous implementations . . . . .	33
2.5	Mixed client-server data structures . . . . .	34
2.5.1	HTML . . . . .	34
2.6	Shared values . . . . .	35
2.6.1	Shared data-structures . . . . .	36
2.6.2	Mixed functors . . . . .	37
2.7	A sophisticated example: accordions . . . . .	37
2.8	Going further . . . . .	40

<b>3</b>	<b>The ML programming language</b>	<b>41</b>
3.1	Syntax	41
3.2	Type system	43
3.2.1	The expression language	43
3.2.2	The module language	45
3.2.3	Inference	49
3.3	Semantics	49
3.3.1	Traces and Printing	52
3.3.2	Modules	53
3.3.3	Notes on Soundness	56
3.4	Related works	56
<b>4</b>	<b>The Eliom programming language</b>	<b>59</b>
4.1	Syntax	60
4.1.1	Locations	60
4.1.2	Expression language	61
4.1.3	Module language	61
4.2	Type system	62
4.2.1	Expressions	62
4.2.2	Modules	65
4.3	Interpreted semantics	72
4.3.1	Generated client programs	72
4.3.2	Base, Client and Server declarations	73
4.3.3	Mixed modules	77
4.4	Results on locations	81
4.4.1	Relation between ML and ELIOM	81
4.4.2	Notes on soundness	83
<b>5</b>	<b>Compilation of Eliom programs</b>	<b>85</b>
5.1	Target languages $ML_s$ and $ML_c$	85
5.1.1	Converters	85
5.1.2	Injections	86
5.1.3	Fragments	87
5.1.4	Modules	88
5.1.5	Type system rules	89
5.1.6	Semantics rules	90
5.2	Compilation	93
5.2.1	Sliceability	95
5.2.2	Slicing rules	96
5.3	Typing preservation	99
5.4	Semantics preservation	101
5.4.1	Hoisting	101
5.4.2	Preliminaries	102
5.4.3	Server expressions and structures	105

5.4.4	Mixed structures . . . . .	107
5.4.5	Proof of the main theorem . . . . .	117
5.5	Discussion around mixed functors . . . . .	119
<b>6</b>	<b>Implementation</b>	<b>121</b>
6.1	The OCAML compiler . . . . .	121
6.2	The ELIOM compiler . . . . .	122
6.2.1	Notes on OCAML compatibility . . . . .	123
6.3	Converters . . . . .	123
6.3.1	Modular implicits . . . . .	124
6.3.2	Wrapping . . . . .	124
6.4	Typechecking and Slicing . . . . .	125
6.4.1	Technical details . . . . .	127
6.5	Runtime and Serialization . . . . .	129
6.5.1	Primitives . . . . .	129
6.5.2	Serialization format . . . . .	131
6.5.3	Optimized placement of sections . . . . .	131
<b>7</b>	<b>State of the art and comparison</b>	<b>133</b>
7.1	Web programming . . . . .	133
7.1.1	Code and data location . . . . .	134
7.1.2	Slicing . . . . .	135
7.1.3	Communications . . . . .	136
7.1.4	Type systems . . . . .	136
7.1.5	Details on some specific approaches . . . . .	137
7.2	Staged meta-programming . . . . .	142
<b>8</b>	<b>Conclusion</b>	<b>145</b>
	<b>Bibliography</b>	<b>148</b>



# 1 Introduction

At the beginning, there was nothing. Suddenly, a burst of light, information and heated discussions: the Internet was born. And for a time, it was enough: Plain text information was exchanged between researchers, computers had 10Mb of RAM and there was no time for frivolous graphical Web pages. In 1989, Tim Berners-Lee disagreed and created the Encyclopedia Of Cats. To make his masterpiece accessible to everyone, he also created the World Wide Web: each page of the encyclopedia was accessed using the HyperText Transfer Protocol (HTTP). Pages were written using the HyperText Markup Language (HTML) and viewed through a Web browser. Related Web pages could be accessed by using Hyperlinks. Web pages could be easily published on the Internet using a Web server. This initiative was so successful that the World Wide Web now counts several billions of Web pages and consolidated cats as the dominant species on Earth.



Figure 1.1: Tim Berners-Lee's first website (<http://info.cern.ch/>) as seen on a modern Web browser. A line-browser simulator is also available, for a more vintage experience.

Early websites such as Tim Berners-Lee’s encyclopedia were static: the Web pages were written by hand and didn’t contain any form of dynamic content. Quickly, people wrote programs to generate Web pages dynamically. In this case, pages are created when the user tries to access them. The programming of Web pages, or Web programming, has been done using a wide array of methods: CGI scripts, dedicated languages such as PHP, but also more generic languages using Web libraries such as Java, C# or Ruby. In all these cases, the idea is the same: a program written in the language in question receives a request from a client, usually a Web browser, and answers it by generating an HTML Web page. The Web page is then sent back to the client. This form of server-side dynamic website was embraced by early dedicated Web programming languages such as PHP (1994) and led to the creation of many popular “Web applications” such as Web search engines (1993), Web forums (around 1994<sup>1</sup>), Wikis (1995), etc. Many modern popular websites still mostly rely on this form of Web programming, for example Wikipedia and Google search.

An important limitation of these websites is that their dynamic behavior is limited to the server. Once the Web page is generated and sent to the client, it doesn’t change until the client requested a new page. To obtain a more dynamic behavior, it was necessary to run programs inside the user’s Web browser. Around 1995, multiple solutions for client-side scripting were developed: Java applets, Flash, and JAVASCRIPT. Numerous websites took advantage of these new capabilities to create interactive experiences available directly in the browser. In particular, Flash and Java allowed the developments of many browser-based games but also more serious applications such as interactive Web maps, music streaming websites, etc.

The combination of both server-side dynamic Web page generation and client-side scripting expanded the capability of websites and allowed to create rich Web applications that rival traditional desktop applications. This paved the way for many modern Web sites such as email clients, browser-based chats, collaborative text editors and more complex websites such as Facebook.

Unfortunately, this situation didn’t come without some complexity. Consider the composition of an email in a Web client such as Gmail: While you type your message, a JAVASCRIPT program provides a rich text editor with bold, italic and emojis. The rich text is rendered in your browser using HTML for the content and CSS for the appearance. When you hit the send button, a dynamic AJAX request, which does not load a new page Web, sends a message containing your email to Google servers. On these servers, a Java program will store the content of the email in a database and send it to the desired recipient. This single task used two programming languages (JAVASCRIPT and Java), at least two description languages (HTML and CSS), all of it spread over two code bases and several distinct communication protocols. While Web applications grew more and more complex, with some subtle interplay of client and server behaviors and very large website with millions of line of code, the way we write Web application stayed similar: a client program, usually called “the frontend”, handles the interactive part and communicates, usually through loosely structured text messages, to a server

---

<sup>1</sup>Although grumpy old-timer will say that Usenet and BBS did it all before.

program called “the backend”. These two programs usually do not share code, not even the definition of the messages communicated, and are usually written by different teams in different programming languages, despite the fact that features used by users usually span client and server side indiscriminately.

## 1.1 On Web programming languages

The core of this issue is, of course, programming languages themselves. While the design of early Web programming languages might have been adapted to how the Web was used at the time, this is not the case anymore. We shall now explore some of the issues prevalent with existing Web languages.

### 1.1.1 The client-server dichotomy

Websites are usually split in two parts, also called “tiers”: a backend that runs on a Web server and a frontend that runs in a browser. Hence, Web programmers must deal with communications between the backend and the frontend. These communications are often static: the backend generates (dynamically or not) a Web page associated with its frontend code and sends it to the browser, which then shows and executes it. Already with this simple scheme, problems can arise: the backend and the frontend should agree about the type and the shape of the generated Web page. For example, if the JAVASCRIPT program tries to find the HTML element named “user”, the backend must have generated it; if the client program is generated dynamically, it should be correct; and so on.

Communications between frontend and backend can also be more complex. Consider the email client example given above: the request to send an email from the client to the server is dynamic and does not generate a new Web page. Communications can be done arbitrarily between client and server through various protocols (AJAXs, Web sockets, Comet, ...). In all these cases, the frontend and the backend must agree on which kind of data is sent. If the backend sends potatoes, the frontend should not expect tomatoes. All these guarantees are difficult to provide statically in today’s Web programming languages and are usually enforced manually by the programmer. This is both time-consuming and error-prone.

This separation of Web applications as separate client and server programs is also problematic from a modularity point of view. Code must now be organized according to where it should run, and not what it does. Indeed, client and server functions can not be mixed arbitrarily, they can not even be in the same file! While this poses significant constraints on how a Web application is organized, it also severely restricts encapsulation. Internal details regarding communications must be accessible to the other side of the application, which often involves exposing them as a protocol to the whole application thus preventing the construction of good abstraction boundaries around individual libraries.

### 1.1.2 Tierless programming languages

Tierless Web programming languages aim to solve these issues by allowing programmers to build dynamic Web pages in a composable way. They allow programmers to define on the server functions that create fragments of a Web page together with their associated client-side behavior. This is done by allowing to freely intersperse client and server expressions with seamless communication in a unique programming language. A tierless program is then sliced in two: a part which runs on the server and a part which is compiled to JAVASCRIPT and runs on the client. This slicing can be done either dynamically, by generating JAVASCRIPT code at runtime, or statically, by cutting the program in two during compilation. Such tierless languages can also leverage both functional programming and static typing. Functional programming provides increased expressiveness and flexibility while static typing gives strong guarantees about client-server separation, in particular ensuring that communications are consistent across tiers.

### 1.1.3 Functional programming

Functional programming languages place the *function* at the center of the art of programming. Such languages allow to manipulate functions as first class values. More concretely, they offer the programmer the ability to create and manipulate functions seamlessly, allowing functions to return or take as argument other functions and enabling functions to be created, composed, stored, copied, etc, in arbitrary manners. Functional programming languages also favor immutable values: values that can not change during the lifetime of the program. To provide new information, new values must be created. Functional programming and immutable data-structures have been proven to provide great benefits in term of expressivity, safety and the ability to create extremely powerful abstractions for a wide variety of use cases as described most elegantly in the seminal essay “Why Functional Programming Matters” by [Hughes \[1989\]](#).

Functional programming was first introduced with the Lisp language by [McCarthy \[1960\]](#). Since then, numerous languages provided support for functional programming: Scheme, OCAML, Haskell, Scala, but also more recently Swift and even Java and C++. Notably, JAVASCRIPT, the language of choice for frontend Web programming, provides support for functional programming which was leveraged to great effect in a large number of libraries such as [Immutable](#) or [Reactjs](#). In a client-server context, functional programming was also used as a tool to structure the control flow across client-server boundaries [[Balat, 2013](#), [Queinnec, 2003, 2004, 2000](#)].

Of particular interest is the ML language. ML, short for Meta Language [[Milner, 1978](#)], aimed to provide a new statically-typed functional programming language tailored for writing theorem provers. This language turned out to be very powerful and useful for a lot more than theorem proving. It inspired numerous languages which now form the larger “ML family”. Some of its direct descendants, such as SML and OCAML, even enjoy large usage today. The main idea common to the numerous members of the ML family is to provide functional, sometimes impure, programming constructs along with a very rich static type system with type inference.



#### 1.1.4 Static type systems

Type systems are an essential component of programming languages. They allow to label each element of the language by its “type” in order to ensure that it is used properly. This can be done dynamically at runtime, statically during a preliminary typechecking phase at compilation, or a mix of both. The choice between dynamically and statically typed languages is a long standing and very heated debate among programmers and numerous essays on the topic can be found online. In the context of this thesis, we will mostly consider the use of statically typed languages as the programming tool of choice. One reason is simply a strong preference from the author of this thesis. More objectively, static typing has been used to great effect in numerous contexts both to prevent bugs, but also to guide and inform the design of programs. Static typing can be further improved through type inference, which allows the programmer to partially or completely omit types annotations. This is notably the case of the ML family of functional programming languages, which combines static typing with a very powerful type inference mechanism to provide both the safety of statically typed languages and the flexible programming style of the more dynamic approaches.

In the context of tierless programming languages, a type system can also provide location information such as “where this code should run” and communication information such as “what am I sending”. Combined with a strong static typing discipline, this ensures that the correctness of communications and library usage is checked statically by the compiler. This can even be done in a fairly non-intrusive way by leveraging type inference, ensuring that client and server expressions can be nested without the need for too many annotations.

#### 1.1.5 Modularity and encapsulation

Mixing tiers directly inside expressions provides a very fined-grained notion of composition. However, programming large-scale software and libraries also requires the capacity for modularity and encapsulation at larger scale.

Modularity is a property ensuring that software can be constructed much like LEGO: by using independent bricks together to build a bigger structure. Each brick, or module, can use other modules but should be independent and self-contained from other non-related modules. Such modularity is also highly desirable in a tierless Web programming contexts. Indeed, parts of a library could be entirely on the server or on the client and programmers should be able to manipulate them freely. Unfortunately, most tierless languages do not support such modular approach to program architecture. Even in traditional Web programming languages such as JAVASCRIPT, support for proper modular programming is very recent (JAVASCRIPT modules were introduced in ES6, in 2015) and does not allow manipulating modules directly nor provides proper encapsulation.

Encapsulation is the complementary aspect of modularity: it allows to hide internal details of a library. Encapsulation improves the flexibility of a complex application: by hiding how each part works internally, the programmer should be able to swap them around as long as their external capabilities are the same. Encapsulation also improves

safety by ensuring that programmers can not misuse the internal details of a library by mistake. Encapsulation is essential for tierless Web programming: the details of how the server and client parts of a widget communicate with each other should not matter, as long as it provides the intended functionality. Communications should only be an internal implementation detail.

To solve these problems, we propose to leverage a well-known tool: ML-style modules.

**ML modules** On top of providing powerful functional programming constructs and static typing, modern ML languages feature a very expressive module system. In these languages, the module language is separate from the expression language. While the language of expression allows to program “in the small”, the module language allows to program “in the large”. In most languages, modules are compilation units: a simple collection of type and value declarations in a file. The SML module language [MacQueen, 1984] uses this notion of collection of declarations (called *structure*) and extends it with types (module specifications, or *signatures*), functions (parametrized modules, or *functors*) and function application, forming a small typed functional language.

In the history of ML languages, ML-style modules have been informally shown to be very expressive tools to architect software. Functors, in particular, allow to write generic implementations by abstracting over a complete module. Furthermore, ML modules provide very good encapsulation through the use of module signatures and abstract datatypes [Leroy, 1995, Crary, 2017]. Module signatures allow to restrict which members of a module are exposed to the outer world while abstract datatypes allow to hide the definition of a given type. Together, they allow to finely control the external interface of a library in order to enforce numerous properties. Another very desirable feature, called separate or incremental compilation, is the ability of the compiler to handle each module separately without needing to recompile all the dependent modules. This allows only the minimal amount of modules to be recompiled after a change which is essential for a fast development cycle. Separate compilation also improves modularity: indeed, checking that a module can indeed be compiled in isolation ensures that compilation can not fail later, in the context of a bigger program. Rich module systems such as ML’s have been shown to provide excellent support for separate compilation [Leroy, 1994, Swasey et al., 2006].

**Modular and tierless** Associating a module language supporting proper modularity, encapsulation and separate compilation in a tierless programming setting is, however, quite delicate. Indeed, while separate compilation is available for languages with rich module systems such as ML, most tierless programming languages rely on whole program compilation to do the slicing. Furthermore, functors are very expressive and raise numerous issues when combined with tierless features that allows to interleave client and server sections.

## 1.2 Eliom

The ELIOM language proposes to solve both fine-grained and large-scale modularity issues. ELIOM is an extension of OCAML, an industrial-strength programming language of the ML family, for tierless web programming. Through the use of new syntactic constructs, it allows to create complete web applications as a single ELIOM program describing both the client and the server behaviors. These new constructions allow statically-checked client-server communications, along with the ability to express rich client-server behaviors in a modular way. Although this increased expressivity could induce significant performance penalties, this is not the case thanks to an efficient compilation and execution scheme. ELIOM programs are compiled statically in two parts: the first part runs on the server while the second part is compiled to JAVASCRIPT and runs on the client. ELIOM also inherits the powerful type system and module system from OCAML, along with a rich ecosystem. This allows us to take advantage of numerous OCAML libraries, such as the rest of the OCSIGEN project.

### 1.2.1 The Ocsigen project

ELIOM is part of the larger OCSIGEN project [Balat et al., 2009]. OCSIGEN provides a comprehensive set of tools and libraries for developing Web applications in OCAML, including the compiler JS\_OF\_OCAML [Vouillon and Balat, 2014], a Web server, and libraries for concurrency [Vouillon, 2008], HTML manipulation [TyXML] and database interaction [Scherer and Vouillon, 2010]. OCSIGEN libraries take deep advantage of the OCAML type system to provide guarantees about various aspects of client- and server-side Web programming. For example, HTML validity is statically guaranteed by the type system [TyXML]. These guarantees are complementary to the ones that ELIOM provides.

Although this thesis will focus on the language extension, ELIOM also comes with a large set of libraries for client and/or server programming that leverages the powerful new constructs introduced by the language. This notably includes RPCs; a functional reactive library for Web programming; a GUI toolkit [Ocsigen Toolkit]; a powerful session mechanism and an advanced service identification mechanism [Balat, 2014]. Some examples of these libraries, and how to implement them, are given in Chapter 2.

### 1.2.2 Principles of the Eliom language

All of the modules and libraries in OCSIGEN, and in particular in the ELIOM framework, are implemented on top of a unique core language. The design of this core language is guided by a set of six properties.

**Explicit communications.** ELIOM uses manual annotations to determine whether a piece of code is to be executed server- or client-side. This design decision stems from our belief that the programmer must be well aware of where the code is to be executed, to avoid unnecessary remote interaction. Explicit annotations also prevent ambiguities in the semantics, allow for more flexibility, and enable the programmer to reason about

where the program is executed and the resulting trade-offs. Programmers can thus ensure that some data stays on the client or on the server, and choose how much communication takes place.

**A simple and efficient execution model.** ELIOM relies on a novel and efficient execution model for client-server communication that avoids back-and-forth communication. This model is simple and predictable. Having a predictable execution model is essential in the context of an impure language, such as OCAML.

**Leveraging the type system.** ELIOM introduces a novel type system that allows composition and modularity of client-server programs while preserving type-safety and abstraction. This ensures, via the type-system, that client functions are not called by mistake inside server code (and conversely) and ensures the correctness of client-server communications.

**Integration with the host language.** ELIOM is an extension of OCAML. Programmers must be able to leverage both the language and the ecosystem of OCAML. OCAML libraries can be useful on the server, on the client or on both. As such, any OCAML file, even when compiled with the regular OCAML compiler, is a valid ELIOM module. Furthermore, we can specify if we want to use a given library on the client, on the server, or everywhere.

**Modularity and encapsulation.** Module and type abstractions are very powerful programming tools. By only exposing part of a library, the programmer can safely hide implementation details and enforce specific properties. ELIOM leverages module abstraction to provide encapsulation and separation of concern for widgets and libraries. By combining module abstraction and tierless features, library authors can provide good APIs that do not expose the fine-grained details of client-server communication to the users.

**Composition.** The ELIOM language allows to define and manipulate *on the server*, as first class values, fragments of code which will be executed *on the client*. This gives us the ability to build reusable widgets that capture both the server and the client behaviors transparently. This makes it possible to define client-server building blocks (and libraries thereof) without further explicit support from the language.

This thesis will explore the consequences of these properties by presenting the design, the formalization and the implementation of the ELIOM language.

## 1.3 Plan

ELIOM aims to be a usable programming language. As such, we introduce ELIOM from a programming perspective in Chapter 2. Tierless programming goes further than just

gluing client and server pieces of code together. Through the various examples, we show that ELIOM enables new programming idioms that allows to build web applications and libraries safely and easily while providing very good encapsulation and modularity properties.

As an extension of OCAML, ELIOM is a compiled language. As for all compiled languages, there is a certain tension between the “intuitive” semantics, explained in term of syntactic reduction over the source language, which is easy to grasp by beginners, and the “real” semantics, which is expressed in term of compilation to a target language and the semantics of this target language. While the interpreted semantics is easier to explain, the implementation of ELIOM uses the compiled semantics. The goal of the formalization of ELIOM was precisely to resolve that tension by ensuring that both semantics agree. For this purpose, we proceed in several steps.

ELIOM is not a language created in isolation, it extends the OCAML programming language. Similarly, our formalization of ELIOM is based on a simpler ML language. The ML family of programming language is quite large and contains many different flavors. Our first step is to describe a simple ML language with modules which forms a reasonable subset of the OCAML language. This is done in Chapter 3.

In Chapter 4, we define the  $\text{ELIOM}_\varepsilon$  language, a subset of ELIOM which is amenable to formalization, and describe its type system and interpreted semantics, along with various properties with relation to the base ML language. ELIOM is compiled by producing two regular OCAML programs for each ELIOM program. Similar, our compilation scheme for  $\text{ELIOM}_\varepsilon$  emits two ML programs with additional primitives. The compilation scheme and the target languages are described in Chapter 5, along with the simulation theorem which ensure that the interpreted semantics and the compiled semantics correspond.

This theorem, however, is not the end of the story. The implementation of programming languages rarely fits directly to the idealized formalization. This certainly applies to OCAML, and thus also to ELIOM. In Chapter 6, we present the implementation of ELIOM as an extension of OCAML. We expose its challenges and the choices we made to make the implementation possible. One particularly important property is the integration between ELIOM and the vanilla OCAML compiler.

Finally, Chapter 7 explores the competitors and inspirations of the ELIOM language, both in the field of tierless Web programming but also in various other domains such as distributed programming and staged meta-programming. Chapter 8 concludes with various remarks on tierless Web programming.

## 1.4 Contributions

ELIOM, and the OCSIGEN project, were initiated by Vincent BALAT and Jérôme VOUIL-LON several years ago. It enjoyed the input of many contributors along the years, who participated to the improvement and the refinement of various ideas you will find in this

thesis. In particular, fragments and sections and their compilation were already present in ELIOM before this thesis started.

My contribution can be summarized as follows:

- A new type system for the ELIOM constructs as an extension of the OCAML type system, including the notion of converters. Previous attempts at typing were limited at best, and unsound and hackish at worst.
- A new module system which fits the programming model of ELIOM while preserving the good properties of the OCAML module system such as abstraction and modularity.
- A formalization of the type system, the interpreted semantics and the compilation scheme.
- A new implementation of the type checker, the compiler and the runtime of ELIOM which more closely reflects the formalization.

More precisely, many examples and libraries presented in Chapter 2 have been developed along the years as part of the OCSIGEN ecosystem, most recently by Vasilis PAPAVALIEOU, Vincent BALAT and Jérôme VOILLON. Furthermore, the compilation of sections, fragments and injections, along with the basic primitives that are presented in Chapter 5, are distilled versions of the original ELIOM implementation due to Pierre CHAMBART, Grégoire HENRY, Benedikt BECKER and Vincent BALAT.

## 2 Programming with Eliom

*"A training course." I look at him. "What in? Windows NT  
system administration?"  
He shakes his head. "Computational demonology for dummies."  
Charles Stross, *The Atrocity Archive**

One of the goal of the ELIOM language is to provide the essential building blocks for type-safe, efficient, tierless web programming. In order to demonstrate this, we first introduce ELIOM's core concept in Section 2.1, then provide numerous examples. Our examples are extracted from code that appears in the OCSIGEN tutorial [Tutorial] and in the ELIOM library [Eliom]. Each example was chosen to illustrate a particular new programming pattern that is used pervasively in the ELIOM ecosystem.

For clarity, we add some type annotations to make the meaning of the code clearer. These annotations are not necessary. While we use the OCAML language, we only assume some familiarity with functional languages.

### 2.1 Core concepts

ELIOM only adds a few new constructions. The aim here is to provide constructions that are sufficiently minimal to be implemented on top of OCAML but also sufficiently expressive to provide convenient implementations of all the Web programming idioms as libraries.

#### 2.1.1 Sections

In ELIOM, we explicitly mark where a definition should be executed through the use of *section* annotations. We can specify whether a declaration is to be performed on the server or on the client as follows:

```
1 let%server s = ...  
2 let%client c = ...
```

A third kind of section, written **shared**, is used for code executed on both sides. Sections allow the programmer to group related code in the same file, regardless of where it is executed.

In the rest of this thesis, we use the following color convention: client is in **yellow**, server is in **blue** and shared is in **green**. Colors are however not mandatory to understand the rest of this thesis.

### 2.1.2 Client fragments

While section annotations allow programmers to gather code across locations, it doesn't allow convenient communication. For this purpose, ELIOM allows to include client-side expression inside a server section: an expression placed inside `[%client ...]` will be computed on the client when it receives the page; but the eventual client-side value of the expression can be passed around immediately as a black box on the server. These expressions are called *client fragments*.

```
1 let%server x : int fragment = [%client 1 + 3]
```

For example, here, the expression `1 + 3` will be evaluated on the client, but it's possible to refer server-side to the future value of this expression (for example, put it in a list). The variable `x` is only usable server-side, and has type `int fragment` which should be read “a fragment containing some integer”. The value inside the client fragment cannot be accessed on the server.

### 2.1.3 Injections

Fragments allow programmers to manipulate client values on the server. We also need the opposite direction. Values that have been computed on the server can be used on the client by prefixing them with the symbol `~%`. We call this an *injection*.

```
1 let%server s : int = 1 + 2
2 let%client c : int = ~%s + 1
```

Here, the expression `1 + 2` is evaluated and bound to variable `s` on the server. The resulting value `3` is transferred to the client together with the Web page. The expression `~%s + 1` is computed client-side.

An injection makes it possible to access client-side a client fragment which has been defined on the server:

```
1 let%server x : int fragment = [%client 1 + 3]
2 let%client c : int = 3 + ~%x
```

The value inside the client fragment is extracted by `~%x`, whose value is `4` here.

Fragments and injections are not novel to this thesis. They draw inspiration from the very rich history of quotations for staging and macros [Bawden, 1999, Sheard and Peyton Jones, 2002] and have been used for several years in various tierless Web programming languages such as HOP [Serrano et al., 2006] and ELIOM. These constructions are quite powerful and allow to express complex client-server interactions in convenient and safe ways, as we see in the rest of this section.

## 2.2 Client-server behaviors

The `hint_button` function below creates a button labeled “Show hint” that pops up a dialog box when activated. The message contained in the dialog box is a server-side



string that is given as argument to the function `hint_button`. One additional property is that the HTML is generated server-side and sent to the client as a regular HTML page.

For this purposes, we use an HTML DSL [**TyXML**] that provides combinators such as `button` and `a onclick` (which respectively create an HTML tag and an HTML attribute). See Section 2.5.1 for more details on this DSL. The `~a` is the OCAML syntax for named arguments. Here, it is used for the list of HTML attributes.

Our function is implemented using a handler for the `onclick` event: since clicks are performed client-side, this handler needs to be a client function inside a fragment. Inside the fragment, an injection is used to access the argument `msg` that contains the string to be showed to the user. The produced HTML fragment is shown in Example 2.1b and the inferred type in Example 2.1a. As we can see, this type does not expose the internal details of the widget's behavior. In particular, the communication between server and client does not leak in the API: This provides proper encapsulation for client-server behaviors. Furthermore, this widget is easily composable: the embedded client state cannot affect nor be affected by any other widget and can be used to build larger widgets.

<pre> 1 let%server hint_button msg = 2   button 3   ~button_type:'Button 4   ~a:[a onclick [%client fun _ -&gt; alert ~%msg] ] 5   [pdata "Show hint"] </pre>	<pre> 1 &lt;button onclick="..."&gt; 2   Show hint 3 &lt;/button&gt; </pre>
<p>(a) Implementation and interface</p>	<p>(b) Emitted HTML</p>

Example 2.1: A button that shows a message

### 2.2.1 Introducing side effects

We now want to generalize our button widget by creating a button that increments a client-side counter and invokes a callback each time it is clicked. This is implemented by the `counter` function, shown below. This client function modifies the widget's state (the client-side reference `state`) and then calls the user-provided client-side callback `action`. This demonstrates that the higher-order nature of OCAML can be used in our client-server setting, and that it is useful for building server-side Web page fragments with parameterized client-side behaviors. In addition, note that the separation between `state` and `action` makes it straightforward to extend this example with a second button that decrements the counter while sharing the associated state.

```

1 let%server counter (action: (int -> unit) fragment) =
2   let state = [%client ref 0 ] in
3   button
4   ~button_type:'Button
5   ~a:[a onclick
6       [%client fun _ -> incr ~%state; ~%action !(~%state) ]]

```

### 2.2.2 Communication and execution scheme

Our counter widget showcases complex patterns of interleaved client and server code, including passing client fragments as arguments to server functions, and subsequently to client code. This would be costly if the communication between the client and the server were done naively.

ELIOM employs an efficient communication mechanism. Specifically, the server only ever sends data along with the initial version of the page. This is made possible by the fact that client fragments are not executed immediately when encountered inside server code. Intuitively, the semantics, presented formally in Section 4.3, is the following. When the server code is executed, the encountered client code is not executed right away; instead it is just registered for later execution once the Web page has been sent to the client. Only then is the client code executed. We also guarantee that client code, be it either client sections or fragments, is executed in the order that it was encountered on the server.

Let us consider the `counter` example above. When calling the `counter` function, we first encounter the `[%client ref 0]` fragment. We generate a fresh identifier that will be used to identify the result of the execution on the client. We store in a queue the fact that this piece of code should be run later. We also return the generated identifier. Later on, we encounter the fragment containing the callback called when the button is pressed: `[%client fun _ -> incr ~%state; ~%action ! (~%state)]`. As before, we generate a fresh identifier and register this piece of code to be run on the client. We also send the content of injections `~%state` and `~%action`. Note that since `action` is a fragment, it is only represented by an identifier that will be used to find the value on the client. Once the server code has been executed, we send the information necessary for the client-side execution to the client. For example, that the identifier previously generated is associated to the result of `ref 0`. We can then execute the client-side part of the ELIOM program. Using the sent information, we can execute client fragments and injections in the order expected by the programmer.

This presentation might make it seem as if we dynamically create the client code during execution of the server code. This is not the case. Like OCAML, ELIOM is statically compiled and separates client and server code at compile time. During compilation, we statically extract the code included inside fragments and compile it as part of the client code to JAVASCRIPT. This allows us to provide both an efficient execution scheme that minimizes communication and preserve side effect orders while still presenting an easy-to-understand semantics. We also benefit from optimizations done by the JS\_OF\_OCAML compiler, thus producing efficient and compact JAVASCRIPT code.

## 2.3 Heterogeneous datatypes

Some datatypes are represented in fundamentally different ways on the server and on the client. This is a consequence of the different nature of the server and the client

environments. ELIOM properly models this heterogeneous aspect by allowing to relate a client and a server datatype that share a similar semantics while having different definitions. We use this feature to present a safe and easy to use API for remote procedure calls (RPCs).

### 2.3.1 Remote procedure calls

When using fragments and injections, the only communication taking place between the client and the server is the original HTTP request and response. However, further communication is sometimes desirable. A remote procedure call is the action of calling, from the client, a function defined on the server. We present here an RPC API implemented using the ELIOM language. The API is shown in Figure 2.1 and an example in Figure 2.2.

```
1 type%server ('i,'o) t
2 type%client ('i,'o) t = 'i -> 'o
3
4 val%server create : ('i -> 'o) -> ('i, 'o) t
```

Figure 2.1: `Rpc` signature

```
1 let%server plus1 : (int, int) Rpc.t =
2   Rpc.create (fun x -> x + 1)
3
4 let%client f x = ~%plus1 x + 1
```

Figure 2.2: Usage of the `Rpc` module

In Figure 2.2, we first create server-side an RPC endpoint using the function `Rpc.create`. Our example RPC adds 1 to its argument. The endpoint is therefore a value of type `(int,int)Rpc.t`, *i.e.*, an RPC whose argument and return values are both of type `int`. The type `Rpc.t` is abstract on the server, but is a synonym for a function type on the client. Of course, this function does not contain the actual implementation of the RPC handler, which only exists server-side.

To use this API, we leverage injections. By using an injection in `~%plus1`, we obtain *on the client* a value of type `Rpc.t`. We describe the underlying machinery that we leverage for converting RPC endpoints into client-side functions in Section 2.3.2. What matters here is that we end up with a function that we can call like any other; calling it performs the remote procedure call.

We can now combine the RPC API with the `counter` widget defined in Section 2.2.1 to create a button that saves the value of the counter on the server. This is presented in Example 2.2. We assume the existence of a `save_counter` function, which saves the counter in a database, and of the `counter` function defined previously. The signature of these functions are shown in Example 2.2a. We then proceed to define `save_counter_rpc` (*i.e.*, the server-side RPC interface for `save_counter`), and inject it into a fragment `f`. This fragment is subsequently used as the user-provided callback for `counter`. This way, each time the counter is incremented, its new value is saved server-side.

The RPC API we proposed is “blocking”: the execution waits for the remote call to finish before pursuing, thus blocking the rest of the client program. Remote procedure calls should, on the contrary, be made asynchronously: the client program keeps running while the call is made and the result is used when the communication is done. In the actual implementation, we use the LWT library [Vouillon, 2008] to express asynchronous calls in a programmer-friendly manner through promises. The use of LWT is pervasive in

the ELIOM ecosystem both on the server and on the client. In this thesis, we will simply omit mentions of the LWT types and operators for pedagogic purposes.

### 2.3.2 Converters

In the RPC API, we associate two types with different implementation on the server and on the client. We rely on injections to transform the datastructure when moving from one side to the other. This ability to transform data before it is sent to the client via an injection is made possible by the use of *converters*. Figure 2.3 broadly presents the converter API. Given a serialization format `serial`, a converter is a pair of a *server* serialization function and a *client* de-serialization function. Note that the client and server types are not necessarily the same. Furthermore, we can arbitrarily manipulate the value before returning it. Several predefined converters are available for fragments, basic OCAML datatypes, and tuples in the module `Conv`. Implementation details about converters can be found in Section 4.2.1.

We can use converters to implement the RPC API (Example 2.3). The server implementation of `Rpc.t` is composed of a handler, which is a server function, and a URL to which the endpoint answers. Our serialization function only sends the URL of the endpoint. The client de-serialization function uses this URL to create a function performing an HTTP request to the endpoint. This way, an RPC endpoint can be accessed simply with an injection. Thus, for the `create` function, we assume that we have a function `serve` of type `string -> (request -> answer) -> unit` that creates an HTTP handler at a specified URL. When `Rpc.create` is called, a unique identifier `id` is created, along with a new HTTP endpoint `"/rpc/id"` that invokes the specified function.

This implementation has the advantage that code using the `Rpc` module is completely independent of the actual URL used. The URL is abstracted away. Converters preserve abstraction by only exposing the needed information.

<pre> 1 val%server save_counter : 2   int -&gt; unit 3 val%server counter : 4   (int -&gt; unit) fragment -&gt; Html.t </pre>	<pre> 1 let%server save_counter_rpc : (int, unit) Rpc.t = 2   Rpc.create save_counter 3 4 let%server widget_with_save : Html.element = 5   let f = [%client ~%save_counter_rpc] in 6   counter f </pre>
---	---

(a) Environment

Example 2.2: Combination of the counter widget and the RPC API.

```

1 type serial (* A serialization format *)
2 type%server ('a, 'b) converter = {
3   serialize : 'a -> serial ;
4   deserialize : (serial -> 'b) fragment
5 }

```

Figure 2.3: Schematized API for converters

### 2.3.3 Client-server reactive broadcasts

In the previous example, we used converters on rather simple datatypes: only a URL was sent, and a closure was created client-side. In this example, we use converters for a more ambitious API: lift *Functional Reactive Programming (FRP)* to be usable across client-server boundaries.

FRP is a paradigm that consists in operating on streams of data, either discrete (events) or continuous (signals). It has been used successfully to program graphical interfaces in a functional fashion, and can also be used to implement Web interfaces. Here, we show how to create an API that allows broadcasting server reactive events to a set of clients.

We assume pre-existing libraries implementing the following two APIs: An *untyped* broadcast API (Figure 2.4) and an FRP event API (Figure 2.5). Both of these APIs are orthogonal to ELIOM's primitives; we can implement broadcast with standard Web techniques, and use the OCAML library **React** for FRP events. The broadcast API operates on messages of type **serial**, the serialization type introduced in Figure 2.3.

Let us now implement the typed broadcast API shown in Figure 2.6. It is quite similar to the RPC API: we have a type **t** with different implementations on the client and the server, and a server function **create** that takes a converter and an event stream as argument and produces a value of type **t**. Here, we use a converter explicitly in order to transfer elements on the broadcast bus.

The implementation of the **Broadcast** module is shown in Figure 2.7. On the server, a **BroadcastEvent.t** is composed of a converter that is used to transfer elements together with a URL. The **create** function starts by creating an untyped broadcast endpoint. We then use **Event.iter** to serialize and then send each occurrence of the provided event.

```
1 type%server ('i,'o) t = {
2   url : string ;
3   handler: 'i -> 'o ;
4 }
5
6 type%client ('i, 'o) t = 'i -> 'o
7
8 let%server serialize t = serialize_string t.url
9 let%client deserialize x =
10   let url = deserialize_string x in
11   fun i -> XmlHttpRequest.get url i
12
13 let conv = {
14   serialize = serialize ;
15   deserialize = [%client deserialize] ;
16 }
17
18 let%server create handler =
19   let url = "/rpc/" ^ generate_new_id () in
20   serve url handler ;
21   { url ; handler }
```

Example 2.3: Simplified RPC implementation corresponding to Figure 2.1.

We now need to create a converter for `BroadcastEvent.t`. We need to transmit two values: the URL of the broadcast endpoint, so that the client can subscribe, and the deserialization part of the provided converter, so that the client can decode the broadcasted messages. `raw_conv` provides a converter for a pair of a URL and a fragment. In addition to receiving this information, the client deserializer creates a new event stream and subscribes to the broadcast endpoint. We connect the broadcast output to the event stream by passing along all the (deserialized) messages.

As we can see in this example, we can use converters explicitly to setup very sophisticated communication schemes in a safe and typed manner. We also use the client deserialization step to execute stateful operations as needed on the client. Note that using a converter here allows effective use of resources: the only clients that subscribe to the broadcast are the ones that really need the event stream, since it has been injected.

## 2.4 Modules and signatures

In the previous examples, we declared simple ELIOM modules containing client and server functions, along with signatures for such modules. We also used “pure” OCAML declarations that are neither client nor server such as the `serial` type (Figure 2.3) or the `Event` module (Figure 2.5). We now give a quick description of the OCAML module system and some of the ELIOM extensions.

### 2.4.1 A primer on OCaml modules

The OCAML module system forms a second language separate from the expression language. While the language of expressions allows to program “in the small”, the module

```
1 val%server t
2 val%server create : url -> t
3 val%server send : t -> serial -> unit
4
5 val%client subscribe :
6   url -> (serial -> unit) -> unit
```

Figure 2.4: Broadcast: Untyped API

```
1 type 'a event
2 (** Events with occurrences of type ['a] *)
3
4 val create : unit -> 'a event * ('a -> unit)
5 (** [create ()] returns an event [e] and a
6     [send] function *)
7
8 val iter : ('a -> unit) -> 'a event -> unit
9 (** [iter f e] applies [f] to
10    [e]'s occurrences *)
```

Figure 2.5: Event: Reactive events API

```
1 type%server ('i, 'o) t
2 type%client ('i, 'o) t = 'o Event.event
3
4 val%server create :
5   ('i, 'o) converter -> 'i event -> ('i, 'o) t
```

Figure 2.6: BroadcastEvent: Shared reactive events API

language allows to program “in the large”. In most languages, modules are compilation units: a simple collection of type and value declarations in a file. The ML module language uses this notion of collection of declarations (called *structure*) and extends it with types (module specifications, or *signatures*), functions (parametrized modules, or *functors*) and function applications, forming a small typed functional language.

In the previous examples, we already implicitly used the module system: each `.ml` file form a structure containing the list of declarations included in the file. It is also possible to specify a signature for such module by adding a `.mli` file. Let us now give a very simple example of functors. For a longer (and better) introduction to modules and functors, please consult the manual [Leroy et al., 2016] or the Real World OCaml

```

1  type%server ('i,'o) t = {
2    conv : ('i, 'o) converter ;
3    url : string ;
4  }
5  type%client ('i,'o) t = 'o Event.event
6
7  let%server create
8    (conv : ('i, 'o) conv) (event : 'i Event.event) =
9    let url = "/broadcast/" ^ generate_new_id () in
10   let t = Broadcast.create url in
11   let send x =
12     Broadcast.send t (conv.serialize x)
13   in
14   let () = Event.iter send event in
15   { conv ; url }
16
17  type%client ('i, 'o) t = 'o Event.event
18
19  let%server raw_conv
20    : (url * 'a fragment, url * 'a) converter
21    = Conv.pair Conv.url Conv.fragment
22
23  let%server serialize t =
24    raw_conv.serialize (t.url, t.conv.deserialize)
25  let%client deserialize s =
26    let url, deserial_msg =
27      ~%raw_conv.deserialize s
28    in
29    let event, send = Event.create () in
30    let handler msg = send (deserial_msg msg) in
31    Broadcast.subscribe url handler ;
32    event
33
34  let%server conv = {
35    serialize ;
36    deserialize = [%client deserialize] ;
37  }

```

Figure 2.7: BroadcastEvent: Shared reactive events. API shown in Figure 2.6.

book [Minsky et al., 2013].

Functors are simply functions that take a module and return another module. They can be used for a large variety of purposes. Here, we use them to build up data structure based on some primitive operations. Let us say we want to create dictionaries with keys of type `t`. One method to implement efficient dictionaries is to use Binary Search Trees, which requires a comparison function for values of type `t` in order to search in the tree. `Map.Make` is a pre-defined functor in the OCAML standard library that takes a module implementing the `COMPARABLE` signature as argument and returns a module that implements dictionaries whose keys are of the type `t` in the provided module. In Figure 2.9, we use this functor to create the `StringMap` module which defines dictionaries with string keys. We then define `d`, a dictionary which associates `"foo"` to `3`.

```
1 module type COMPARABLE = sig
2   type t
3   val compare : t -> t -> int
4 end
5
6 module Make (Key : COMPARABLE) : sig
7   type 'a t
8   val add : Key.t -> 'a -> 'a t -> 'a t
9   (* ... *)
10 end

1 module StringComp = struct
2   type t = string
3   let compare = String.compare
4 end
5 module StringMap = Map.Make(StringComp)
6
7 let d : int StringMap.t =
8   StringMap.add "foo" 3 StringMap.empty
```

Figure 2.9: Dictionaries from strings to ints

Figure 2.8: the `Map` module

## 2.4.2 Modules and locations

Section annotations are also available on module declarations, which allows to define client and server modules. One can also use regular OCAML modules and functors inside client and server code. For example, in Figure 2.10, we use `Map.Make` on the client to define maps whose keys are JAVASCRIPT strings. JAVASCRIPT strings are fairly different from OCAML strings, as they are represented by ropes instead of mutable byte arrays, hence the need for a different type. Note here that a functor from vanilla OCAML is applied to a client module and returns a client module.

As we saw in the previous examples, we can mix declarations from multiple locations inside the same module. Such modules are called “mixed”. An important constraint is that, as you go down inside sub modules, locations should be properly included: A client module can not contain server declarations and conversely, but mixed modules can contain everything.

## 2.4.3 OCaml integration

We consider an additional location, “base”, which can only contain OCAML constructs. Pure OCAML declarations such as `serial`, including pure OCAML modules, are considered of location base and are usable both on the client and on the server. Side-effecting



```

1 module%client JStr = struct
2   type t = Js.string
3   let compare = Js.compare_string
4 end
5
6 module%client JStrMap = Map.Make(JStr)

```

Figure 2.10: Map of js strings

base code will be executed on both the client and the server. ELIOM guarantee that pieces of code inside base locations can only contain pure OCAML code, without any of the additional ELIOM constructs. The semantics of base code is guaranteed to be exactly the same as the vanilla OCAML semantics. Furthermore, OCAML library can be imported and linked inside ELIOM projects freely. It is even possible to use the compilation output of the regular OCAML compiler.

Additionally, it is possible to decide that a given OCAML library should be imported in ELIOM only client or server side. For example, we might want to use a vanilla OCAML database library in our ELIOM project. We can simply specify that this library should only be loaded on the server and the ELIOM type system will prevent its use in client code.

#### 2.4.4 Heterogeneous implementations

Shared sections make it possible to write code for the client and the server at the same time. This provides a convenient way of writing terse shared implementations, without duplicating logic and code. This does not necessarily entail that everything is shared. In particular, base primitives might differ between client and server, though the overall logic is the same. Just as we can implement heterogeneous datatypes with different client- and server-side representations, we can also provide interfaces common to the client and the server, with different client- and server-side implementations. We consider the case of database access. We first assume the existence of a server function `get_age` of type `string -> int` that performs a database query and returns the age of a person.

We can easily create a client version of that function via our RPC API of Figure 2.1.

```

1 let%server get_age_rpc = Rpc.create get_age
2 let%client get_age = %get_age_rpc

```

The API is then:

```

1 val%shared get_age : string -> int

```

We can use this function to write widgets that can be used either on the client or on the server:

```

1 let%shared person_widget name =
2   div ~a:[a_class "person"] [
3     text (name^" : ^string_of_int(get_age name))
4   ]

```

This technique is used pervasively in ELIOM to expose implementations than can be used either on the client or on the server with similar semantics, in a very concise way.

## 2.5 Mixed client-server data structures

We can readily embed client fragments inside server data structures. As a simple example of such a mixed data structure, consider a list of button names (standard server-side strings) and their corresponding client-side actions. Example 2.4 presents a function that takes such a list and builds an unordered HTML list of buttons.

ELIOM makes such mixed data-structure particularly easy to write. Having explicit annotations with the usage of fragments is essential here. This would be quite difficult to achieve if the delimitation between client and server values were implicitly inferred.

```
1 let%server button_list
2   (lst : (string * handler fragment) list) =
3   ul (List.map (fun (name, action) ->
4       li [button
5           ~button_type:'Button
6           ~a:[a_onclick ~%action]
7           [pdata name]])
8       lst)
```

Example 2.4: Function generating a list of buttons

### 2.5.1 HTML

A common idiom in Web programming is to generate the skeleton of a Web page on the server, then fill in the holes on the client with dynamic content, or bind dynamic client-side behaviors on HTML elements. In order to do that, the usual technique is to use the `id` or `class` HTML properties to identify elements, and to manually make sure that these identifiers are used in a coherent manner on the client and the server.

ELIOM simplifies this process by mean of a client-server HTML library that allows injections of HTML elements to the client. Figure 2.11 shows a simplified API, which is uniform across clients and servers. The API provides combinators such as the `div` function shown below, which builds a `div` element with the provided attributes and child elements. We already used this HTML API in several previous examples.

On the server, HTML is implemented as a regular OCAML datatype. When sending the initial HTML document, this datatype is converted to a textual representation. This ensures compatibility with JAVASCRIPT-less clients and preserves the usual behavior of a Web server.

On the client, we represent HTML nodes directly as DOM trees. The mismatch between client and server implementations does not preclude us from providing a uniform API. However, to permit injections of HTML nodes from the server to the client, special care must be taken. In particular, we equip each injected node with an `id`, and `id` is the only piece of data sent by the serialization function. The deserialization function then

```

1 type%shared attribute
2 type%shared element
3
4 val%shared div :
5   ?a:(attribute list) -> element list -> element
6
7 val%server a_onclick :
8   (Event.t -> bool) fragment -> attribute
9
10 module%server Client : sig
11   val node : element fragment -> element
12 end

```

Figure 2.11: Html: The simplified HTML API

finds the element with the appropriate `id` in the page. The `a_onclick` function finds the appropriate HTML element on the client and attaches the specified handler.

The fact that we use a uniform API allows us to abstract the specificities of the DOM and to provide other kinds of representations, such as a virtual DOM approach. A further improvement that fits in our design is nesting client HTML elements inside server HTML documents without any explicit DOM manipulation. This is done by the `Client.node` function (Figure 2.12), which takes a client fragment defining an HTML node and converts it to a server-side HTML node that can be embedded into the page. This function works by including a placeholder element server-side. The placeholder is later replaced by the actual element on the client.

```

1 let%server node (x: element fragment) : element =
2   let placeholder = span [] in
3   let _ = [%client
4     let placeholder = ~%placeholder in
5     let x = ~%x in
6     Option.iter
7       (Dom.parent placeholder)
8       (fun parent ->
9         Dom.replaceChild parent placeholder x)
10  ] in
11   placeholder

```

Figure 2.12: Implementation of `Client.node`

## 2.6 Shared values

We presented in Section 2.4.4 how we can use shared sections to write code that is used both on the client and on the server. Using the ELIOM features we have described, we can also create *shared values*, which have a similar dual meaning, but at the level of expressions. The API is described in Figure 2.13 while the implementation is shown in Figure 2.14. Implementation of the converter for shared values is shown in Figure 2.15.

```

1 type%server ('a, 'b) shared_value =
2   { srv : 'a ; cli : 'b fragment }
3 type%client ('a, 'b) shared_value = 'b
4
5 val%server local : ('a, 'b) shared_value -> 'a
6 val%client local : ('a, 'b) shared_value -> 'b
7
8 val%server cli :
9   ('a, 'b) shared_value -> 'b fragment
10 val%client cli : ('a, 'b) shared_value -> 'b

```

Figure 2.13: Shared values API

```

1 let%server local x = x.srv
2 let%client local x = x
3
4 let%server cli x = x.cli
5 let%client cli x = x

```

Figure 2.14: Shared values Implementation

```

1 let%server shared_conv
2   : (('a, 'b) shared_value, ('a, 'b) shared_value) converter
3   = {
4     serialize = (fun x -> Conv.fragment.serialize x.cli);
5     deserialize = Conv.fragment.deserialize
6   }

```

Figure 2.15: Converter for shared values

The server-side implementation of a shared value clearly needs to contain a fragment that can be injected on the client. On the other hand, the client cannot possibly inject a value on the server, so the client-side representation only consists of a fragment. For injecting a server-side shared value on the client, we use a converter whose server-side portion serializes only the fragment, and whose client-side portion deserializes this fragment.

## 2.6.1 Shared data-structures

Shared values are very useful when a given operation needs to be performed both on the server and on the client, but in a way that matches the specific requirement of each side. As an example, we present a *cached dictionary* API for storing data of interest on both the server and the client. This dictionary API should be well-adapted for ELIOM's client-server style of programming. On the server, the dictionary is to be used while serving a request, e.g., for locally caching data obtained from complex database queries. It is frequently the case that the client needs access to the same data; in that case, it is desirable that we avoid performing multiple RPCs. To achieve this, the semantics of the server-side addition operation (function **add**) is such that the value does not only become available for future server-side lookups, but also for client-side lookups. Of course, additional items may be added client-side, but then there is no expectation of server-side addition; the server-side dictionary may not even exist any longer, given that it was local to the code handling the request.

As with the **Map** module, we want to define our dictionary over arbitrary keys. Let us first review what we need to define such a data-structure: We want two comparison function: one on the server and one on the client. We also need a converter that ensures

we can transmit values from the server to the client. For simplicity, we consider that the type of our keys is an OCAML base type. The type of our shared dictionaries (`'a, 'b`)`table` contains two type variables `'a` and `'b`, corresponding to the server- and client-side contents respectively. The API in Figure 2.16a provides `add` and `find` operations, as is typical for association tables, which are available on both sides. The implementation is shown in Figure 2.16b. A dictionary is implemented as a pair of a server-side dictionary and a client-side one. The server-side `add` implementation stores a new value locally in the expected way, but additionally builds a fragment that has the side-effect of performing a client-side addition. The retrieval operation (`find`) returns a shared value that contains both the server side version and the client side. On the client, however, we can directly use the local values. Note that since the client-side type exactly corresponds to a regular map, we can directly use the usual definitions for the various map operations. This is done by including the `M` module on the client.

Several extensions of this API are possible. For pedagogic purposes, the type of key we used here is a pure OCAML type on the base location (Section 2.4.3). We could also have two different kinds of keys on the server and on the client, for example OCAML strings on the server and JAVASCRIPT strings on the client, which would ensure better efficiency. The function would then need two comparison functions and a converter between the two types. Alternatively, we could also easily create a full blown replicated dictionary: by using the RPC API, the client can require the server dictionary to be updated and by using the broadcast API, we can distribute new additions to all clients.

Going further, shared values empower an approach to reactive programming that is well-adapted for ELIOM's client-server paradigm [Shared reactive programming]. This approach is the subject of ongoing work. One notable possible improvement is the notion of shared fragment, the analogous of shared declarations for expressions, which allow to avoid some code duplication present in the implementation in Figure 2.16b.

## 2.6.2 Mixed functors

In the previous example, we defined a functor which creates a module containing both client and server declarations. Such functors are called *mixed*. As we saw, mixed functors are fairly powerful and can be used for a wide variety of purposes. However, contrary to client and server functors, mixed functors are limited: arguments must be mixed modules, mixed structures and functors can not be nested arbitrarily and injections inside client-side bindings can only reference elements out of the functor. Injections inside client fragments can be used in arbitrary ways. A more precise description of these limitations is provided in Section 4.2.2.

## 2.7 A sophisticated example: accordions

We now demonstrate how to implement the well-known widget *accordion*. An accordion is a kind of application menu that displays collapsible sections in order to present information in a limited amount of space. The section titles are always visible. The content of a section is shown when the user clicks on its title. Only one section is open at a time.

<pre> 1 module type T = sig 2   type t 3   val comparable : t -&gt; t -&gt; int 4   val%server conv : (t, t) conv 5 end 6 7 module Cache (Key : T) : sig 8   module M = Map.Make(Key) 9 10  type%shared ('a, 'b) table = 11    ('a M.t, 'b M.t) Shared.t 12 13  val%shared add : 14    Key.t -&gt; ('a, 'b) Shared.t -&gt; 15    ('a, 'b) table -&gt; ('a, 'b) table 16 17  val%shared find : 18    Key.t -&gt; ('a, 'b) table -&gt; 19    ('a, 'b) Shared.t 20 21  (* ... *) 22 end </pre>	<pre> 1 module Cache (Key : T) = struct 2   module M = Map.Make(Key) 3 4   type%shared ('a, 'b) table = 5     ('a M.t, 'b M.t) Shared.t 6 7   include%client M 8 9   let%server add id v tbl = 10     [%client M.add ~id ~v ~%tbl ]; 11     M.add id v.srv tbl.srv 12 13   let%server find id tbl = 14     { srv = M.find id tbl ; 15       cli = [%client M.find ~id ~%tbl] 16     } 17 18   (* ... *) 19 end </pre>
(a) Signature	(b) Implementation

Figure 2.16: The SharedTable module

This widget can be used in numerous contexts such as the body of news articles, trip details in a train ticket search, etc.

In our example, sections are implemented independently and attached to the accordion given as parameter. The distinctive characteristic of our implementation, made possible by the two-level language, is that a section can be generated freely either on the server or on the client, and attached to an existing accordion. The example contains three sections, two generated server-side and one added dynamically client-side to the same accordion.

The code is shown in [Figure 2.17](#). The data structure representing the accordion contains only a reference to a client-side function that closes the currently open section. Functions `new_accordion` and `accordion_section` are included in both the server and client programs (shared sections). Function `switch_visibility` is implemented client-side only. It just adds or removes an HTML class to the element, which has the effect of hiding or showing the element through CSS rules. Function `my_accordion` builds a server-side HTML page containing an accordion with two sections. It also sends to the client process, together with the page, the request to create the accordion (client fragment in function `new_accordion`) and to append a new section to the accordion. For this purpose, we use function `Client.node` on line [40](#).

```

1 let%client switch_visibility elt =
2   if Class.contain elt "hidden"
3   then Class.remove elt "hidden"
4   else Class.add elt "hidden"
5
6 type%shared toggle =
7   (unit -> unit) ref fragment
8
9 let%shared new_accordion () : toggle =
10  [%client ref (fun () -> ()) ]
11
12 let%shared accordion_section
13   (accordion : toggle) s1 s2 =
14   let contents =
15     div ~a:[a_class ["contents"; "hidden"]]
16     [text s2]
17   in
18   let handler = [%client fun _ ->
19     let toggle = ~%accordion in
20     !toggle (); (*close previous section*)
21     toggle :=
22       (fun () -> switch_visibility ~%contents);
23     switch_visibility ~%contents
24   ]
25   in
26   let title =
27     div
28     ~a:[a_class ["title"]; a_onclick handler]
29     [text s1]
30   in
31   div ~a:[a_class ["section"]]
32     [title; contents]
33
34 let%server my_accordion () =
35   let accordion = new_accordion () in
36   div [
37     accordion_section
38     accordion
39     "Item 1" "Server side generated" ;
40     Client.node [%client
41       accordion_section
42       ~%accordion
43       "Item 2" "Client side generated"
44     ] ;
45     accordion_section
46     accordion
47     "Item 3" "Server side generated" ;
48   ]

```

Figure 2.17: The accordion widget

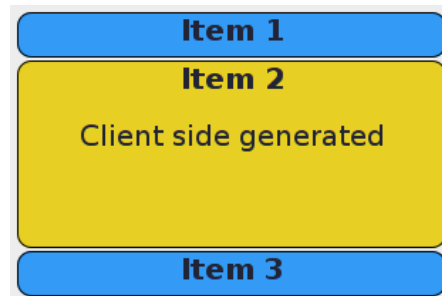


Figure 2.18: Resulting web page

## 2.8 Going further

Our examples demonstrate how the combination of fragments, injections and converters, along with the abstraction and modularity provided by the module system, can be used to build rich Web development libraries that provide convenient programming interfaces. While the module system and the notion of converter presented here are novel to this thesis; fragments, injections and abstraction have been leveraged in the ELIOM library for many years to build very rich structures. In particular, the ELIOM library provides uni- and bi-directional channels, progressive fetching of data, correct-by-construction links, and client-server reactive programming. Interestingly, a common pattern arising across these examples (just like for our RPC and HTML examples of Sections 2.3.1 and 2.5.1) is relating server and client datatypes that differ in their structure and APIs, but that have related intuitive meaning. Of course, the same building blocks and patterns can be used by the programmer to implement additional components outside the ELIOM library, thus catering for their specific use cases.



## 3 The ML programming language

*When I was a child, I was told that Santa Claus came in through the chimney, and that computers were programmed in binary code. Since then, I have learned that programming is better done in higher-level languages, more abstract and more expressive.*

Xavier Leroy, *Polymorphic typing of an algorithmic language*

Programming languages of the ML family are like curry: everyone cook them differently and nobody agrees on the perfect ingredients, but they all end up being delicious. For ML, the original ingredients [Milner, 1978] are first class functions, parametric polymorphism, let bindings and algebraic datatypes. With the years, several new ingredients have been added to the mix; in particular mutable references, exceptions, pattern matching and a module system; which form the base of many languages in the ML family [Milner et al., 1990, Leroy et al., 2016].

Our version of ML contains the minimal amount of ingredients that allows us to describe the ELIOM extensions: a core calculus with polymorphism, let bindings and parametrized datatypes in the style of Wright and Felleisen [1994], accompanied by a fully featured module system in the style of Leroy [1995]. We first present the syntax of the language in Section 3.1 and its type system in Section 3.2. We then present the semantics of this language in Section 3.3. Finally we indicate some relevant work in Section 3.4.

### 3.1 Syntax

Let us first define some notations and meta-syntactic variables. As a general rule, the expression language is in lowercase ( $e$ ) and the module language is in uppercase ( $M$ ). Module types are in calligraphic letters ( $\mathcal{M}$ ). More precisely:  $x$  are variables,  $p$  are module paths,  $X$  are module variables,  $\tau$  are type expressions and  $t$  are type constructors.  $x_i$ ,  $X_i$  and  $t_i$  are identifiers (for values, modules and types). Identifiers (such as  $x_i$ ) have a name part ( $x$ ) and a stamp part ( $i$ ) that distinguish identifiers with the same name. Only the name part of identifiers is exposed in module signature.  $\alpha$ -conversion should keep the name intact and change only the stamp, which allow to preserve module signatures. Lists are noted with a star; for example  $\tau^*$  is a list of type expressions. Indexed lists are noted  $(\tau_i)$ , with an implicit range. Substitution of  $a$  by  $b$  in  $e$  is noted  $e[a \mapsto b]$ . Repeated substitution of each  $a_i$  by the corresponding  $b_i$  is noted  $e[a_i \mapsto b_i]_i$ . The syntax is presented in Figure 3.1.

**Expressions** The expression language is a fairly simple extension of the lambda calculus with a fixpoint combinator  $Y$  and let bindings  $\text{let } x = e_1 \text{ in } e_2$ . The language is parametrized by a set of constants  $Const$ . Variables can be qualified by a module path  $p$ . Paths can be either module identifiers such as  $X_i$ , a submodule access such as  $X_i.Y$ , or a path application such as  $X_i(Y_j.Z)$ . Note that, as said earlier, that fields of modules are only called by their name, without stamp.

**Types** Types are composed of type variables  $\alpha$ , function types  $\tau_1 \rightarrow \tau_2$  and parametrized type constructors  $(\tau_1, \tau_2, \dots, \tau_k)t_i$ . Type constructors can have an arbitrary number of parameters, including zero. Type constructors can be qualified by a module path  $p$ . Type schemes, noted  $\sigma$ , are type expressions that are universally quantified by a list of type

Expressions		Path
$e ::= c$	(Constant)	$p ::= X_i \mid p.X \mid p_1(p_2)$
$\mid x_i \mid p.x$	(Variables)	<b>Type Schemes</b>
$\mid Y$	(Fixpoint)	$\sigma ::= \forall \alpha^*. \tau$
$\mid (e \ e)$	(Application)	<b>Type Expressions</b>
$\mid \lambda x. e$	(Function)	$\tau ::= \alpha$ (Type variables)
$\mid \text{let } x = e \text{ in } e$	(Let binding)	$\mid \tau \rightarrow \tau$ (Function types)
$c \in Const$	(Constants)	$\mid (\tau^*)t_i \mid (\tau^*)p.t$ (Type constructors)

(a) The expression language

Module Expressions		Module types
$M ::= X_i \mid p.X$	(Variables)	$\mathcal{M} ::= \text{sig } S \text{ end}$ (Signature)
$\mid (M : \mathcal{M})$	(Type constraint)	$\mid \text{functor}(X_i : \mathcal{M}_1)\mathcal{M}_2$ (Functor)
$\mid M_1(M_2)$	(Functor application)	<b>Signature body</b>
$\mid \text{functor}(X_i : \mathcal{M})M$	(Functor)	$S ::= \varepsilon \mid \mathcal{D}; S$
$\mid \text{struct } s \text{ end}$	(Structure)	<b>Signature components</b>
<b>Structure body</b>		$\mathcal{D} ::= \text{val } x_i : \tau$ (Values)
$S ::= \varepsilon \mid \mathcal{D}; S$		$\mid \text{type } (\alpha^*)t_i = \tau$ (Types)
<b>Structure components</b>		$\mid \text{type } (\alpha^*)t_i$ (Abstract types)
$D ::= \text{let } x_i = e$	(Values)	$\mid \text{module } X_i : \mathcal{M}$ (Modules)
$\mid \text{type } (\alpha^*)t_i = \tau$	(Types)	<b>Environments</b>
$\mid \text{module } X_i = M$	(Modules)	$\Gamma ::= S$
<b>Programs</b>		
$P ::= \text{prog } S \text{ end}$		

(b) The module language

Figure 3.1: ML grammar

variables. Type schemes can also have free variables. For example:  $\forall\alpha.(\alpha, \beta)t_i$ .

**Modules** The module language is quite similar to a simple lambda calculus: Functors are functions over module (except that arguments are annotated with their types). Module application is noted  $M_1(M_2)$ . Modules can also be constrained by a module type:  $(M : \mathcal{M})$ . Finally, a module can be a structure which contains a list of value, types or module definitions: `struct let  $x_i = 2$  end`. Programs are lists of definitions.

**Module types** Module types can be either the type of a functor or a signature, which contains a list of value, types and module descriptions. Type descriptions can expose their definition or can be left abstract. Typing environments are simply module signatures. We note them  $\Gamma$  for convenience.

## 3.2 Type system

We now present the ML type system. For ease of presentation, we proceed in two steps: we will first forget that the module language exists, and present a self-contained type system for the expression language. We then extend the typing relation to handle modules.

### 3.2.1 The expression language

We introduce the following judgments:

- $\Gamma \triangleright e : \tau$  The expression  $e$  has type  $\tau$  in the environment  $\Gamma$ . See Figure 3.3.
- $\Gamma \triangleright \tau_1 \approx \tau_2$  Types  $\tau_1$  and  $\tau_2$  are equivalent in environment  $\Gamma$ . See Figure 3.4.
- $\Gamma \models \tau$  The type  $\tau$  is well formed in the environment  $\Gamma$ . See Figure 3.2.

We note  $TypeOf(c)$  the type scheme of a given constant  $c$ . The instantiation relation is noted  $\sigma \succ \tau$  for a type scheme  $\sigma$  and a type  $\tau$ . The converse operation which closes a type according to an environment is noted  $Close(\Gamma, \tau)$ . We use  $(D) \in \Gamma$  to test if a given type or value is declared in the environment  $\Gamma$ . Note that for types,  $(\mathbf{type} (\alpha_i)t) \in \Gamma$  holds also if  $t$  is not abstract in  $\Gamma$ .

**Polymorphism** One of the main benefit of programming language of the ML family is the ability to easily define and use functions that operate on values of various types. For example, the `map` function can applies to all lists, regardless of the type of their content. Indeed, the type of `map` is polymorphic:

$$\mathbf{map} : \forall\alpha\beta. (\alpha \rightarrow \beta) \rightarrow (\alpha)\mathbf{list} \rightarrow (\beta)\mathbf{list}$$

From a type checking point of view, this is possible thanks to two operations: instantiation and abstraction. Instantiation takes a type scheme, which is a type where some variables have been universally quantified, and replace all the quantified type variables by

some type. It is used when looking up a variable (rule VAR) or typechecking a constant (rule CONST). For example, the type of map can be instantiated to the following type.

$$(\text{int} \rightarrow \text{bool}) \rightarrow (\text{int})\text{list} \rightarrow (\text{bool})\text{list}$$

Once instantiated, the **map** function can be applied on a list with concrete types. Naturally, we also need the converse operation: constructing a type scheme given a type containing some type variables. Closing a type depends on the current typing environments, we only abstract type variables that have not been introduced by previous binders.  $\text{Close}(\Gamma, \tau)$  returns the type scheme  $\forall \alpha_1 \dots \alpha_n. \tau$  where the  $\alpha_i$  are free variables of  $\tau$  that are not present in  $\Gamma$ . While it is possible to apply the closing operation at any step of a typing derivation, it is only useful at the introduction point of type variables, in let bindings (rule LETIN). In the following example, we derive a polymorphic type for a function that constructs a pair with an element from the environment. We first use the close operation to obtain a type scheme for  $f$ . Note that since  $\alpha$  is present in the environment, it is not universally quantified. We then use the instance operation to apply  $f$  to an integer constant.

$$\frac{\frac{\vdots}{(\text{val } a : \alpha; \text{val } b : \beta) \triangleright (a, b) : \alpha * \beta}}{(\text{val } a : \alpha) \triangleright \lambda b. (a, b) : \beta \rightarrow \alpha * \beta} \quad \frac{\frac{\forall \beta. \beta \rightarrow \alpha * \beta \succ \text{int} \rightarrow \alpha * \text{int}}{\Gamma \triangleright f : \text{int} \rightarrow \alpha * \text{int}} \quad \frac{\text{Const}(3) \succ \text{int}}{\Gamma \triangleright 3 : \text{int}}}{(\text{val } a : \alpha; \text{val } f : \forall \beta. \beta \rightarrow \alpha * \beta) \triangleright f \ 3 : (\alpha * \text{int})} \\ (\text{val } a : \alpha) \triangleright \text{let } f = \lambda b. (a, b) \text{ in } f \ 3 : (\alpha * \text{int})$$

**Parametric datatypes** Parametric polymorphism introduces type variables in type expressions. In the presence of type definitions, it is natural to expect the ability to write type definitions which can contain type variables. This leads us to parametric datatypes: datatypes which are parametrized by a set of variables.  $(\alpha)\text{list}$  is of course an example of such datatype. Note that care must be taken when deciding the equivalence of types. If the type is not abstract, *i.e.*, its definition is available, we can always unfold the definition, as shown in rule DEFTYPEEQ. However, when considering an abstract type, we cannot unfold the type definition. Instead, we check that head symbols are compatible and that parameters are equivalent pairwise<sup>1</sup>. This is done in rule ABSTYPEEQ.

$$\frac{\text{TYPEVAL} \quad (\text{type } (\alpha_i)t) \in \Gamma \quad \forall i, \Gamma \models \tau_i}{\Gamma \models (\tau_i)t} \quad \frac{\text{ARROWVAL} \quad \Gamma \models \tau_1 \quad \Gamma \models \tau_2}{\Gamma \models \tau_1 \rightarrow \tau_2} \quad \frac{\text{VARVAL}}{\Gamma \models \alpha}$$

Figure 3.2: Type validity rules –  $\Gamma \models \tau$

<sup>1</sup>This is similar to the handling of free symbols in the unification literature.

$$\begin{array}{c}
\text{VAR} \\
\frac{(\text{val } x : \sigma) \in \Gamma \quad \sigma \succ \tau}{\Gamma \triangleright x : \tau}
\end{array}
\quad
\begin{array}{c}
\text{LAMBDA} \\
\frac{\Gamma; (\text{val } x : \tau_1) \triangleright e : \tau_2}{\Gamma \triangleright \lambda x. e : \tau_1 \rightarrow \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{CONST} \\
\frac{\text{TypeOf}(c) \succ \tau}{\Gamma \triangleright c : \tau}
\end{array}$$

$$\begin{array}{c}
\text{LETIN} \\
\frac{\Gamma \triangleright e_1 : \tau_1 \quad \Gamma; (\text{val } x : \text{Close}(\tau_1, \Gamma)) \triangleright e_2 : \tau_2}{\Gamma \triangleright \text{let } x = e_1 \text{ in } e_2 : \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{EQUIV} \\
\frac{\Gamma \triangleright e : \tau_1 \quad \Gamma \triangleright \tau_1 \approx \tau_2}{\Gamma \triangleright e : \tau_2}
\end{array}$$

$$\begin{array}{c}
\text{APP} \\
\frac{\Gamma \triangleright e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \triangleright e_2 : \tau_1}{\Gamma \triangleright (e_1 \ e_2) : \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{Y} \\
\frac{}{\Gamma \triangleright \text{Y} : ((\tau_1 \rightarrow \tau_2) \rightarrow \tau_1 \rightarrow \tau_2) \rightarrow \tau_1 \rightarrow \tau_2}
\end{array}$$

$\text{Close}(\tau, \Gamma) = \forall \alpha_0 \dots \alpha_n. \tau$  with  $\{\alpha_0, \dots, \alpha_n\} = \text{FreeTypeVar}(\tau) \setminus \text{FreeTypeVar}(\Gamma)$

Figure 3.3: ML expression typing rules –  $\Gamma \triangleright e : \tau$

$$\begin{array}{c}
\text{REFLEQ} \\
\frac{}{\Gamma \triangleright \tau \approx \tau}
\end{array}
\quad
\begin{array}{c}
\text{TRANSEQ} \\
\frac{\Gamma \triangleright \tau_1 \approx \tau_2 \quad \Gamma \triangleright \tau_2 \approx \tau_3}{\Gamma \triangleright \tau_1 \approx \tau_3}
\end{array}
\quad
\begin{array}{c}
\text{COMMEQ} \\
\frac{\Gamma \triangleright \tau_2 \approx \tau_1}{\Gamma \triangleright \tau_1 \approx \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{FUNEQ} \\
\frac{\Gamma \triangleright \tau_1 \approx \tau'_1 \quad \Gamma \triangleright \tau_2 \approx \tau'_2}{\Gamma \triangleright \tau_1 \rightarrow \tau_2 \approx \tau'_1 \rightarrow \tau'_2}
\end{array}$$

$$\begin{array}{c}
\text{DEFTYPEEQ} \\
\frac{(\text{type } (\alpha_i) t = \tau) \in \Gamma}{\Gamma \triangleright (\tau_i) t \approx \tau[\alpha_i \mapsto \tau_i]_i}
\end{array}
\quad
\begin{array}{c}
\text{ABSTYPEEQ} \\
\frac{(\text{type } (\alpha_i) t) \in \Gamma \quad \forall i, \Gamma \triangleright \tau_i \approx \tau'_i}{\Gamma \triangleright (\tau_i) t \approx (\tau'_i) t}
\end{array}$$

Figure 3.4: Type equivalence rules –  $\Gamma \triangleright \tau \approx \tau'$

### 3.2.2 The module language

We introduce the following judgments:

- $\Gamma \blacktriangleright M : \mathcal{M}$  The module  $M$  is of type  $\mathcal{M}$  in  $\Gamma$ . See Figure 3.5.
- $\Gamma \blacktriangleright \mathcal{M} <: \mathcal{M}'$  The module type  $\mathcal{M}$  is a subtype of  $\mathcal{M}'$  in  $\Gamma$ . See Figure 3.6.
- $\Gamma \models \mathcal{M}$  The module type  $\mathcal{M}$  is well-formed in  $\Gamma$ . See Figure 3.7.

The typing rules for OCAML-style modules are quite complex. In particular, the inner details of the rules are not well known, even by OCAML programmers. Before presenting the typing rules in details, we will attempt to give insight on why some features are present in the languages and what are their advantages. For this purpose, we present two examples illustrating the need for applicative functors and strengthening, respectively. We assume that readers are familiar with simpler usages of ML modules.

**Applicative Functors** Let us consider the following scenario: we are given a module  $G$  implementing a graph data-structure and would like to implement a simple graph algorithm which takes a vertex and returns all the accessible vertices. We would like the

returned module to contain a function of type  $G.\text{graph} \rightarrow G.\text{vertex} \rightarrow \text{set\_of\_vertices}$ . How to implement `set_of_vertices`? An easy but inefficient way would be to use lists. A better way is to use proper sets (implemented with balanced binary tree, for example). In OCAML, this is provided in the standard library by the functor *Set.Make*, presented in Section 2.4.1, which takes a module implementing comparison functions for the given type. We would obtain a signature similar to the one below.

```
module Access(G : Graph) : sig
  module VerticesSet : sig ... end
  val run : G.graph → G.vertex → VerticesSet.set
end
```

However, this means we need to expose a complete module implementing set of vertices that is independent from any other set module. This prevents modularity, since any usage of our new function must use this specific set implementation. Furthermore, this make the signature bigger than strictly necessary. What we really want to expose is that the return type comes from an application of *Set.Make*. Fortunately, we can do so by using the following signature.

```
module Access(G : Graph) : sig
  val run : G.graph → G.vertex → Set.Make(G.Vertex).set
end
```

Here, we export the fact that the set type must be the result of a functor application on a module that is compatible with *G.Vertex*. The type system guarantees that any such functor application will produces types that are equivalent. In particular, if multiple libraries uses the *Access* functor, their sets will be of the same types, which make composition of libraries easier. This behavior of functors is usually called *applicative*.

**Strengthening** Let us now consider the program presented in Example 3.1. We assume the existence of two modules, presented in Example 3.1a. The module *Showable* exposes the abstract type *t*, along with a *show* function that turns it into a string. The module *Elt* exposes a type *t* equal to *Showable.t* and a value that inhabits this type. The program is presented in Example 3.1b. We define a functor *F* taking two arguments *E* and *S* whose signature are similar to *Elt* and *Showable*, respectively. The main difference is that *E* comes first and *S.t* is defined as an alias of *E.t*. The functor uses the *show* function on the element in *E* to create a string. It is natural to expect the functor application *F(Elt)(Showable)* to type check, since *Elt.t* = *Showable.t*. We must, however, check for module inclusion. While *Elt* is clearly included in the signature of the argument *E*, the same is not clear for *Showable*. We first need to enrich its type signature with additional type equalities. We give *Showable* the type `sig type t = Showable.t ... end`. It makes sense to enrich the signature in such a manner since *Showable* is already in

<pre> module Showable : sig   type t   val show : t → string end module Elt : sig   type t = Showable.t   val v : elt end </pre> <p>(a) Typing environment</p>	<pre> module F   (E : sig type t val v : t end)   (S : sig type t = E.t val show : t → string end) = struct   let s = (S.show E.v) end module X = F(Elt)(Showable) </pre> <p>(b) Application of multi-argument functor using manifests</p>
--	--

Example 3.1: Program using functors and manifest types

the environment. Given this enriched signature, we can now deduce that  $\blacktriangleright(\text{type } t = \text{Showable.t}) <: (\text{type } t = E.t)$  since  $E.t = \text{Elt.t} = \text{Showable.t}$ .

The operation that consists in enriching type signatures of module identifiers with new equalities by using elements in the environment is called *strengthening* [Leroy, 1994].

### Typing rules

In the previous two examples, we showcased some delicate interactions between functor, type equalities and modularity in the context of an ML module system. We now see in details how the rules presented in Figures 3.5 to 3.7 produce these behaviors.

**Qualified access** Unqualified module variables are typechecked in a similar manner than regular variables in the expression language, with the MODVAR typing rule. Qualified access (of the form  $X.a$ ), both for the core and the module language, need more work. As with the expression language, the typing environment is simply a list of declaration. In particular, typing environments do not store paths. This means that in order to prove  $\Gamma \triangleright p.a : \tau$ , we must first verify that the module  $p$  typechecks in  $\Gamma$ :  $\Gamma \blacktriangleright p : \mathcal{M}$ . We then need to verify that the module type  $\mathcal{M}$  contains a declaration  $(\text{val } a : \tau)$ . This is done in the QUALMODVAR rule for the module language. The rules for the expression language are given in Figure 3.8.

Let us now try to apply these rules to the module  $X$  with the following module type.  $X$  contains a type  $t$  and a value  $a$  of that type. We note that module type  $\mathcal{X}$ .

$$X : \text{sig type } t; \text{val } a : t \text{ end}$$

We wish to typecheck  $X.a$ . One expected type for this expression is  $X.t$ . However, the binding of  $v$  in  $\mathcal{X}$  gives the type  $t$ , with no mention of  $X$ . We need to prefix the type variable  $t$  by the access path  $X$ . This is done in the rule QUALMODVAR by the substitution  $\mathcal{M}[n_i \mapsto p.n \mid n_i \in \text{BV}(\mathcal{S}_1)]$  which prefixes all the bound variables of  $\mathcal{S}_1$ , noted  $\text{BV}(\mathcal{S}_1)$ , by the path  $p$ . Note here that we substitute only by the names declared *before* the variable  $a$ . Indeed, a variable or a type can only reference names declared

previously in ML. To prove that  $X.a$  has the type  $X.t$ , we can write the following type derivation.

$$\text{QUALVAR} \frac{\text{MODVAR} \frac{(\text{module } X : \mathcal{X}) \in (\text{module } X : \mathcal{X})}{(\text{module } X : \mathcal{X}) \blacktriangleright X : \text{sig type } t; \text{ val } a : t \text{ end}}}{(\text{module } X : \mathcal{X}) \triangleright X.a : X.t} \text{ WITH } X.t = t[t \mapsto p.t]$$

**Strengthening** The strengthening operation, noted  $\mathcal{M}/p$ , is defined in Figure 3.9 and is used in the STRENGTH rule. It takes a module type  $\mathcal{M}$  and a path  $p$  and returns a module type  $\mathcal{M}'$  where all the type declarations, abstract or not, have been replaced by type aliases pointing to the path  $p$ . These type aliases are usually called “manifest types”. This operator relies on the following idea: if  $p$  is of type  $\mathcal{M}$ , then  $p$  is available in the environment. In order to expose as many type equalities as possible, it suffices to give  $p$  a type where all the type definition point to definitions available in the environment. This way, we preserve type equalities even for abstract types. This also mean that type equalities can be deduced by only looking at the path and the module type. In particular, we do not need to look at the implementation of  $p$ , which is important for the purpose of separate compilation.

**Applicative functors** Let us consider a functor  $F$  with the following type. It takes a module containing a single type  $t$  and return a module containing an abstract type  $t'$  and a conversion function.

$$F : \text{functor}(X : \text{sig type } t \text{ end})(\text{sig type } t'; \text{ val } \text{make} : X.t \rightarrow t' \text{ end})$$

If we consider two modules  $X_1$  and  $X_2$ , does  $X_1 = X_2$  imply  $F(X_1).t = F(X_2).t$ ? If that is the case, we say that functors are *applicative*. Otherwise, they are *generative*<sup>2</sup>. Here, we consider the applicative behavior of functors. This is implemented with the last strengthening rule which ensures that the body of functors is also strengthened. For example, if  $\mathcal{M}$  is the type of the functor above,  $\mathcal{M}/F$  is the following module type:

$$\text{functor}(X : \text{sig type } t \text{ end})(\text{sig type } t' = F(X).t; \text{ val } \text{make} : X.t \rightarrow t' \text{ end})$$

This justifies the presence of application inside paths. Otherwise, such type manifests inside functors could not be represented. A more type-theoretic description of generative and applicative functors can be found in Leroy [1996].

**Separate compilation** Separate compilation is an important properties of programming languages. In fact, almost all so-called “mainstream” languages support it. We can distinguish two aspects of this property: separate typechecking and separate code generation. In both cases, it means that in order to process the file (either to type check

<sup>2</sup>SML only supports generative functors. OCAML originally only supported applicative functors, but also supports the generative behavior since version 4.03.



it or to transform it into another representation), we only need to look at the type of its dependencies, not their implementation.

It turns out that the ML module system with manifest types lends itself very well to separate typechecking [Leroy, 1994]. Indeed, let us consider a program as a list of modules. Each module represents a compilation unit (*i.e.*, a file). Since module bindings in the typing environment only contains module types, and not the actual module, typechecking a file only needs the module type of the previous files, which ensure that we can typecheck each file separately, as long as all its dependencies have been typechecked before. This is expressed more formally in Theorem 1.

**Theorem 1** (Separate Typechecking). Given a list of module declarations that form a typed program, there exists an order such that each module can be typechecked with only knowledge of the type of the previous modules.

More formally, given a list of  $n$  declarations  $D_i$  and a signature  $\mathcal{S}$  such that

$$\blacktriangleright (D_1; \dots; D_n) : \mathcal{S}$$

then there exists  $n$  definitions  $\mathcal{D}_i$  and a permutation  $\pi$  such that

$$\forall i < n, \mathcal{D}_1; \dots; \mathcal{D}_i \blacktriangleright D_{i+1} : \mathcal{D}_{i+1} \quad \blacktriangleright \mathcal{D}_{\pi(1)}; \dots; \mathcal{D}_{\pi(n)} <: \mathcal{S}$$

*Proof.* It is always possible to reorder declarations in a signature using the SUBSTRUCT rule. This means we can choose the appropriate permutation of definitions that matches the order of declarations. The rest follows by definition of the typing relation.  $\square$

### 3.2.3 Inference

Full inference is one of the greatest strength of the ML programming language. While we do not address inference formally in this thesis, here are some remarks. Inference is of course decidable for the core language using the well known  $\mathcal{W}$  algorithm. It is “efficient”, which means here that it is fast for usual programs, though pathological cases can be constructed. The typechecking rules for modules are not syntax directed. The STRENGTH rule, in particular, is free floating. Leroy [1994] presents how to turn this into a syntax-directed type system, which allows inference as long as functor arguments are annotated.

## 3.3 Semantics

We now define the semantics of our ML language. We use a rule-based big step semantics with traces. Traces allows us to reason about execution order in a way that is compatible with modules, as we will see in Section 3.3.1.

We note  $v$  for values in the expression language and  $V$  for values in the module language. Values are defined in Figure 3.11. Values in the expression language can be either constants or lambdas. Module values are either structures, which are list of bindings of

$$\begin{array}{c}
\text{MODVAR} \quad \frac{(\text{module } X_i : \mathcal{M}) \in \Gamma}{\Gamma \blacktriangleright X_i : \mathcal{M}} \quad \text{QUALMODVAR} \quad \frac{\Gamma \blacktriangleright p : (\text{sig } \mathcal{S}_1; \text{module } X_i : \mathcal{M}; \mathcal{S}_2 \text{ end})}{\Gamma \blacktriangleright p.X : \mathcal{M}[n_i \mapsto p.n \mid n_i \in \text{BV}(\mathcal{S}_1)]} \quad \text{STRENGTH} \quad \frac{\Gamma \blacktriangleright p : \mathcal{M}}{\Gamma \blacktriangleright p : \mathcal{M}/p} \\[10pt]
\frac{\Gamma \blacktriangleright M : \mathcal{M}' \quad \Gamma \blacktriangleright \mathcal{M}' <: \mathcal{M}}{\Gamma \blacktriangleright M : \mathcal{M}} \quad \frac{\Gamma \blacktriangleright M_1 : \text{functor}(X_i : \mathcal{M})\mathcal{M}' \quad \Gamma \blacktriangleright M_2 : \mathcal{M}}{\Gamma \blacktriangleright M_1(M_2) : \mathcal{M}'[X_i \mapsto M_2]} \\[10pt]
\frac{\Gamma \models \mathcal{M} \quad X_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{module } X_i : \mathcal{M}) \blacktriangleright M : \mathcal{M}'}{\Gamma \blacktriangleright \text{functor}(X_i : \mathcal{M})M : \text{functor}(X_i : \mathcal{M})\mathcal{M}'} \quad \frac{\Gamma \models \mathcal{M} \quad \Gamma \blacktriangleright M : \mathcal{M}}{\Gamma \blacktriangleright (M : \mathcal{M}) : \mathcal{M}} \\[10pt]
\frac{\Gamma \triangleright e : \tau \quad x_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{val } x_i : \text{Close}(\tau, \Gamma)) \blacktriangleright S : \mathcal{S}}{\Gamma \blacktriangleright (\text{let } x_i = e; s) : (\text{val } x_i : \tau; \mathcal{S})} \\[10pt]
\frac{\Gamma \models \tau \quad t_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{type } (\alpha^*)t_i = \tau) \blacktriangleright S : \mathcal{S}}{\Gamma \blacktriangleright (\text{type } (\alpha^*)t_i = \tau; s) : (\text{type } (\alpha^*)t_i = \tau; \mathcal{S})} \\[10pt]
\frac{\Gamma \blacktriangleright M : \mathcal{M} \quad X_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{module } X_i : \mathcal{M}) \blacktriangleright S : \mathcal{S}}{\Gamma \blacktriangleright (\text{module } X_i = M; s) : (\text{module } X_i : \mathcal{M}; \mathcal{S})} \\[10pt]
\frac{\Gamma \blacktriangleright S : \mathcal{S}}{\Gamma \blacktriangleright \text{struct } S \text{ end} : \text{sig } \mathcal{S} \text{ end}} \quad \frac{}{\Gamma \blacktriangleright \varepsilon : \varepsilon}
\end{array}$$

Figure 3.5: Module typing rules –  $\Gamma \blacktriangleright m : \mathcal{M}$

$$\begin{array}{c}
\text{SUBSTRUCT} \quad \frac{\pi : [1; m] \rightarrow [1; n] \quad \forall i \in [1; m], \Gamma; \mathcal{D}_1; \dots; \mathcal{D}_n \blacktriangleright \mathcal{D}_{\pi(i)} <: \mathcal{D}'_i}{\Gamma \blacktriangleright (\text{sig } \mathcal{D}_1; \dots; \mathcal{D}_n \text{ end}) <: (\text{sig } \mathcal{D}'_1; \dots; \mathcal{D}'_m \text{ end})} \\[10pt]
\frac{\Gamma \triangleright \tau_1 \approx \tau_2}{\Gamma \blacktriangleright (\text{val } x_i : \tau_1) <: (\text{val } x_i : \tau_2)} \quad \frac{\Gamma \blacktriangleright \mathcal{M}_1 <: \mathcal{M}_2}{\Gamma \blacktriangleright (\text{module } X_i : \mathcal{M}_1) <: (\text{module } X_i = \mathcal{M}_2)} \\[10pt]
\frac{\Gamma \blacktriangleright \mathcal{M}'_a <: \mathcal{M}_a \quad \Gamma, (\text{module } X : \mathcal{M}'_a) \blacktriangleright \mathcal{M}_r <: \mathcal{M}'_r}{\Gamma \blacktriangleright \text{functor}(X : \mathcal{M}_a)\mathcal{M}_r <: \text{functor}(X : \mathcal{M}'_a)\mathcal{M}'_r} \\[10pt]
\frac{\Gamma \triangleright \tau_1 \approx \tau_2}{\Gamma \blacktriangleright (\text{type } (\alpha^*)t_i = \tau_1) <: (\text{type } (\alpha^*)t_i = \tau_2)} \quad \frac{}{\Gamma \blacktriangleright (\text{type } (\alpha^*)t_i) <: (\text{type } (\alpha^*)t_i)} \\[10pt]
\frac{\Gamma \triangleright (\alpha^*)t_i \approx \tau}{\Gamma \blacktriangleright (\text{type } (\alpha^*)t_i) <: (\text{type } (\alpha^*)t_i = \tau)} \quad \frac{}{\Gamma \blacktriangleright (\text{type } (\alpha^*)t_i = \tau_1) <: (\text{type } (\alpha^*)t_i)}
\end{array}$$

Figure 3.6: Module subtyping rules –  $\Gamma \blacktriangleright \mathcal{M} <: \mathcal{M}'$

$$\begin{array}{c}
\frac{\Gamma \models \mathcal{S}}{\Gamma \models \text{sig } \mathcal{S} \text{ end}} \quad \frac{}{\Gamma \models \varepsilon} \quad \frac{\Gamma \models \mathcal{M}_a \quad x_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{module } X_i : \mathcal{M}_a) \models \mathcal{M}_r}{\Gamma \models \text{functor}(X_i : \mathcal{M}_a) \mathcal{M}_r} \\
\\
\frac{t_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{type } (\alpha^*) t_i) \models \mathcal{S}}{\Gamma \models \text{type } (\alpha^*) t_i; \mathcal{S}} \quad \frac{\Gamma \models \mathcal{M} \quad x_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{module } X_i : \mathcal{M}) \models \mathcal{S}}{\Gamma \models \text{module } X_i : \mathcal{M}; \mathcal{S}} \\
\\
\frac{\Gamma \models \tau \quad t_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{type } (\alpha^*) t_i = \tau) \models \mathcal{S}}{\Gamma \models \text{type } (\alpha^*) t_i = \tau; \mathcal{S}} \\
\\
\frac{\Gamma \models \tau \quad x_i \notin \text{BV}(\Gamma) \quad \Gamma; (\text{val } x_i : \tau) \models \mathcal{S}}{\Gamma \models \text{val } x_i : \tau; \mathcal{S}}
\end{array}$$

Figure 3.7: Module type validity rules –  $\Gamma \models \mathcal{M}$

$$\begin{array}{c}
\text{QUALVAR} \quad \frac{\Gamma \blacktriangleright p : (\text{sig } \mathcal{S}_1; \text{val } x_i : \tau; \mathcal{S}_2 \text{ end})}{\Gamma \triangleright p.x : \tau[n_i \mapsto p.n \mid n_i \in \text{BV}(\mathcal{S}_1)]} \quad \text{QUALDEFTYPEEQ} \quad \frac{\Gamma \blacktriangleright p : (\text{sig } \mathcal{S}_1; \text{type } (\alpha^*) t_i = \tau; \mathcal{S}_2 \text{ end})}{\Gamma \triangleright (\tau_i) p.t \approx \tau[n_i \mapsto p.n \mid n_i \in \text{BV}(\mathcal{S}_1)][\alpha_i \mapsto \tau_i]_i} \\
\\
\text{QUALABSTYPEEQ} \quad \frac{\Gamma \blacktriangleright p : (\text{sig } \mathcal{S}_1; \text{type } (\alpha^*) t_i; \mathcal{S}_2 \text{ end}) \quad \forall i, \Gamma \triangleright \tau_i \approx \tau'_i}{\Gamma \triangleright (\tau_i) p.t \approx (\tau'_i) p.t}
\end{array}$$

Figure 3.8: Additional typing rules for the expression language

$$\begin{array}{l}
\varepsilon/p = \varepsilon \\
(\text{sig } \mathcal{S} \text{ end})/p = \text{sig } \mathcal{S}/p \text{ end} \\
(\text{module } X_i = \mathcal{M}; \mathcal{S})/p = \text{module } X_i = \mathcal{M}/p; \mathcal{S}/p \\
(\text{type } (\alpha^*) t_i = \tau; \mathcal{S})/p = \text{type } (\alpha^*) t_i = (\alpha^*) p.t; \mathcal{S}/p \\
(\text{type } (\alpha^*) t_i; \mathcal{S})/p = \text{type } (\alpha^*) t_i = (\alpha^*) p.t; \mathcal{S}/p \\
(\text{val } x_i : \tau; \mathcal{S})/p = \text{val } x_i : \tau; \mathcal{S}/p \\
(\text{functor}(X_i : \mathcal{M}) \mathcal{M}')/p = \text{functor}(X_i : \mathcal{M})(\mathcal{M}'/p(X_i))
\end{array}$$

Figure 3.9: Module strengthening operation –  $\mathcal{M}/p$

values, or functors. We note  $\rho$  the execution environment. Execution environments are a list of value bindings. Note here that the execution environment is not mutable, since reference cells are not in the language. We note the concatenation of environment  $+$ . Environment access is noted  $\rho(x) = v$  where  $x$  has value  $v$  in  $\rho$ . The same notation is also used for structures. Traces are lists of messages. For now, we consider messages that are values and are emitted with a `print` operation. The empty trace is noted  $\langle \rangle$ . Concatenation of traces is noted  $@$ .

Given an expression  $e$  (resp. a module  $m$ ), an execution environment  $\rho$ , a value  $v$  (resp.  $V$ ) and a trace  $\theta$ ,

$$e \xRightarrow{\rho} v, \theta$$

means that  $e$  reduces to  $v$  in  $\rho$  and prints  $\theta$ . The reduction rules are given in Figure 3.12. The rules for the expression language are fairly traditional. Variables and paths must be resolved using the VAR and QUALVAR rules. Applications are done in two steps: first, we reduce both the function and the argument with the APP rule, then we apply the appropriate reduction rule for the application: BETA for lambda expressions, Y for fixpoints and DELTA for constants. The  $\delta$  operation gives meaning to application of a constant to a value.  $\delta(c, v) = v', \theta$  means that  $c$  applied to  $v$  returns  $v'$  and emits the trace  $\theta$ . Let bindings are treated in a similar manner than lambda expressions: the left hand side is executed, added to the environment, then the right hand side is executed.

The module language has similar rules for identifiers and application. In this case, the BETA and APP rule have been combined in MODBETA. Additional rules for declarations are also present. Type declarations are ignored (TYPEDECL). Values and module declarations (VALDECL and MODDECL) are treated similarly to let bindings: the body of the binding is executed, added to the environment and then the rest of the structure is executed.

### 3.3.1 Traces and Printing

Traces allow us to visualize the execution order of programs. In particular, if we prove that code transformation preserves traces, it ensures that the execution order is preserved. Traces allow us to reason about execution without introducing references and other side-effecting operations in our language, which would make the presentation significantly more complex.

One example of operation using traces is the `print` constant. Typing and semantics of `print` are provided in Figure 3.10. `print` accepts any value, prints it, and returns it. From a typing point of view, `print` has the same type as the identity: a polymorphic function which returns its input. We make use of the fact that the CONST typing rule also uses the instantiation for type schemes. The semantics of `print` is provided via the DELTA rule: it returns its argument directly but also emits a trace containing the given argument.

We now present an example using `print`. We assume the existence of the type `int`, a set of constant corresponding to the integers and an associated operation  $+$ . We wish to type and execute the expression  $e$  defined as `let  $x = (\text{print } 3)$  in  $(\text{print } (x + 1))$`

$$\begin{array}{ll}
\text{PRINTTY} & \text{PRINTEXEC} \\
\text{TypeOf}(\text{print}) = \forall \alpha. (\alpha \rightarrow \alpha) & \delta(\text{print}, v) = v, \langle v \rangle
\end{array}$$

Figure 3.10: Typing and execution rules for `print`

Let us first show that  $e$  is of type `int`. The type derivation is provided in Example 3.2. The typing derivation is fairly direct: we use the `CONST` rule to type `print` as `int`  $\rightarrow$  `int` and apply it to integers with the rule `APP`. We can now look at the execution of  $e$ , which returns 4 with a trace  $\langle 3; 4 \rangle$ . The execution derivation is shown in Example 3.3. The first step is to decompose the let-binding. We first reduce  $(\text{print } 3)$ , which can be directly done with the `DELTA` rule. This gives us 3 with a trace  $\langle 3 \rangle$ . We then reduce  $(\text{print } (x + 1))$  in the environment where  $x$  is associated to 3. Before resolving the application of `print` with the `DELTA` rule, we need to reduce its argument with the `APP` rule. We obtain 4 with a trace  $\langle 4 \rangle$ . We return the result of the right hand side of the left and the concatenation of both traces by usage of the `LETIN` rule, which gives us 4 with a trace  $\langle 3; 4 \rangle$ .

$$\begin{array}{c}
\text{CONST} \frac{\text{TypeOf}(\text{print}) \succ \text{int} \rightarrow \text{int}}{\triangleright \text{print} : \text{int} \rightarrow \text{int}} \quad \frac{\text{TypeOf}(3) = \text{int}}{\triangleright 3 : \text{int}} \quad \frac{\vdots}{\triangleright \text{print} : \text{int} \rightarrow \text{int}} \quad \frac{\vdots}{\triangleright x + 1 : \text{int}} \\
\text{APP} \frac{\triangleright \text{print} : \text{int} \rightarrow \text{int} \quad \triangleright 3 : \text{int}}{\triangleright (\text{print } 3) : \text{int}} \quad \frac{\triangleright \text{print} : \text{int} \rightarrow \text{int} \quad \triangleright x + 1 : \text{int}}{(\text{val } x : \text{int}) \triangleright (\text{print } (x + 1)) : \text{int}} \\
\text{LETIN} \frac{\triangleright (\text{print } 3) : \text{int} \quad (\text{val } x : \text{int}) \triangleright (\text{print } (x + 1)) : \text{int}}{\triangleright \text{let } x = (\text{print } 3) \text{ in } (\text{print } (x + 1)) : \text{int}}
\end{array}$$

Example 3.2: Typing derivation for  $e - \triangleright e : \text{int}$

$$\begin{array}{c}
\text{DELTA} \frac{\delta(\text{print}, 3) = 3, \langle 3 \rangle}{(\text{print } 3) \Rightarrow 3, \langle 3 \rangle} \quad \frac{\text{print} \Rightarrow \text{print}, \langle \rangle \quad x + 1 \xrightarrow{\{x \mapsto 3\}} 4, \langle \rangle}{(\text{print } (x + 1)) \xrightarrow{\{x \mapsto 3\}} 4, \langle 4 \rangle} \quad \frac{\delta(\text{print}, 4) = 4, \langle 4 \rangle}{(\text{print } 4) \Rightarrow 4, \langle 4 \rangle} \\
\text{APP} \frac{(\text{print } 3) \Rightarrow 3, \langle 3 \rangle \quad (\text{print } (x + 1)) \xrightarrow{\{x \mapsto 3\}} 4, \langle 4 \rangle}{\text{let } x = (\text{print } 3) \text{ in } (\text{print } (x + 1)) \Rightarrow 4, \langle 3; 4 \rangle} \\
\text{LET} \frac{\text{let } x = (\text{print } 3) \text{ in } (\text{print } (x + 1)) \Rightarrow 4, \langle 3; 4 \rangle}{\text{let } x = (\text{print } 3) \text{ in } (\text{print } (x + 1)) \Rightarrow 4, \langle 3; 4 \rangle}
\end{array}$$

Example 3.3: Execution derivation for  $e - e \Rightarrow 4, \langle 3; 4 \rangle$

### 3.3.2 Modules

We now present an example of reduction involving modules. Our example program  $P$  is presented in Example 3.4a. It consists of two declarations: a module declaration  $X$  which contains a single declaration  $a$ , and the return value of the program, which is equal to  $X.a$ . It is fairly easy to see that the program  $P$  return a value of type `int`, hence we focus on the execution of  $P$ , which is presented in Example 3.4b. The derivation is slightly simplified for clarity. In particular, rules such as `EMPTYSTRUCT` are elided. The first step is to apply the `PROGRAM` and `MODULEDECL` rules in order to execute the content of

each declaration. The declaration of  $X$ , on the left side, can be reduced by first applying the **STRUCT** rule in order to extract the content of the module structure, then **VALDECL**, to reduce the declaration of  $a$ . These reductions give us the structure value  $\{a \mapsto 3\}$ . We now execute the declaration of **return**. According to the **MODULEDECL** rule, we must do so in a new environment containing  $X$ :  $\{X \mapsto \{a \mapsto 3\}\}$ . In order to reduce  $X.a$ , we must use the **QUALMODVAR** rule, which reduces qualified variables. This means we first reduce  $X$ , which according to the environment gives us  $\{a \mapsto 3\}$ , noted  $V$ . We then look up  $a$  in  $V$ , which returns 3. To return, we first compose the resulting structure value from both declaration:  $\{X \mapsto \{a \mapsto 3\}\} + \{\text{return} \mapsto 3\}$ . We then lookup **return** in this structure, which gives us 3.

```

prog
  module X = struct let a = 3 end
  let return = X.a
end

```

(a) The program  $P$

$$\begin{array}{c}
\text{VALDECL} \frac{3 \Rightarrow 3, \langle \rangle}{\text{let } a = 3 \Rightarrow \{a \mapsto 3\}, \langle \rangle} \quad \text{MODVAR} \frac{\rho(X) = V}{X \xRightarrow{\rho} V \equiv \{a \mapsto 3\}, \langle \rangle} \quad V(a) = 3 \\
\text{STRUCT} \frac{\left( \begin{array}{l} \text{struct} \\ \text{let } a = 3 \\ \text{end} \end{array} \right) \Rightarrow \{a \mapsto 3\}, \langle \rangle}{\text{let return} = X.a \xRightarrow{\{X \mapsto \{a \mapsto 3\}\}} \{\text{return} \mapsto 3\}, \langle \rangle} \quad \text{QUALMODVAR} \frac{X.a \xRightarrow{\{X \mapsto \{a \mapsto 3\}\}} 3, \langle \rangle}{\text{let return} = X.a \xRightarrow{\{X \mapsto \{a \mapsto 3\}\}} \{\text{return} \mapsto 3\}, \langle \rangle} \\
\text{MODULEDECL} \frac{\left( \begin{array}{l} \text{module } X = \text{struct let } a = 3 \text{ end} \\ \text{let return} = X.a \end{array} \right) \Rightarrow \{X \mapsto \{a \mapsto 3\}\} + \{\text{return} \mapsto 3\}, \langle \rangle}{P \Rightarrow 3, \langle \rangle} \quad \text{MODULEDECL} \\
\text{PROGRAM} \frac{\left( \begin{array}{l} \text{module } X = \text{struct let } a = 3 \text{ end} \\ \text{let return} = X.a \end{array} \right) \Rightarrow \{X \mapsto \{a \mapsto 3\}\} + \{\text{return} \mapsto 3\}, \langle \rangle}{P \Rightarrow 3, \langle \rangle} \quad \text{PROGRAM}
\end{array}$$

(b) Execution of  $P$

Example 3.4: Example of execution with modules

**Why big steps?** One might wonder why use a big step semantics with traces, instead of a small step semantics. Indeed, small step usually make proofs easier, especially for simulations which we will use on ELIOM later. A first remark is that modules are not stable by substitution since  $(\text{struct } \dots \text{end}).x$  is not valid (and is problematic to type). Hence we need to use a semantics with environments and closures. Let us now consider doing one step deep inside a structure using a small step semantics. The previously evaluated declarations in the structure should be available in the local environment which mean we would need to rebuild environments as we explore a context to execute a small step. We would also need to manipulate partially evaluated structures as we execute declarations. Furthermore, typing preservation for small steps would be difficult to express in the presence of abstract types. While this is all possible, big steps semantics with environments is, by comparison, fairly straightforward.

Expressions	Modules	Bindings
$v ::= c$ (Constant)	$V ::= \{ V_b^* \}$ (Structure)	$V_b ::= \{ x_i \mapsto v \}$ (Values)
$  \lambda x. \rho. e$ (Function)	$  \mathbf{functor}(\rho)(X_i : \mathcal{M})M$	$  \{ X_i \mapsto V \}$ (Modules)

Figure 3.11: ML values

$\frac{\text{VAR} \quad \rho(x) = v}{x \xRightarrow{\rho} v, \langle \rangle}$	$\frac{\text{QUALVAR} \quad p \xRightarrow{\rho} V, \theta \quad V(x) = v}{p.x \xRightarrow{\rho} v, \theta}$	$\frac{\text{CONSTANT}}{c \xRightarrow{\rho} c, \langle \rangle}$	$\frac{\text{CLOSURE}}{\lambda x. e \xRightarrow{\rho} \lambda x. \rho. e, \langle \rangle}$
$\frac{\text{LETIN} \quad e' \xRightarrow{\rho} v', \theta \quad e \xRightarrow{\rho + \{x \mapsto v'\}} v, \theta'}{(\mathbf{let} \ x = e' \ \mathbf{in} \ e) \xRightarrow{\rho} v, \theta @ \theta'}$	$\frac{\text{APP} \quad e \xRightarrow{\rho} v, \theta \quad e' \xRightarrow{\rho} v', \theta' \quad (v \ v') \xRightarrow{\rho} v'', \theta''}{(e \ e') \xRightarrow{\rho} v'', \theta @ \theta' @ \theta''}$		
$\frac{\text{BETA} \quad e \xRightarrow{\rho' + \{x \mapsto v'\}} v', \theta}{(\lambda x. \rho'. e \ v) \xRightarrow{\rho} v', \theta}$	$\frac{\text{Y} \quad (v \ \lambda x. (\mathbf{Y} \ v \ x)) \xRightarrow{\rho} v', \theta}{(\mathbf{Y} \ v) \xRightarrow{\rho} v', \theta}$	$\frac{\text{DELTA} \quad \delta(c, v) = v', \theta}{(c \ v) \xRightarrow{\rho} v', \theta}$	
$\frac{\text{MODVAR} \quad \rho(X) = V}{X \xRightarrow{\rho} V, \langle \rangle}$	$\frac{\text{QUALMODVAR} \quad p \xRightarrow{\rho} V', \theta \quad V'(X) = V}{p.X \xRightarrow{\rho} V, \theta}$	$\frac{\text{STRUCT} \quad S \xRightarrow{\rho} V_s, \theta}{(\mathbf{struct} \ S \ \mathbf{end}) \xRightarrow{\rho} V_s, \theta}$	$\frac{\text{EMPTYSTRUCT}}{\varepsilon \xRightarrow{\rho} \{ \}, \langle \rangle}$
$\frac{\text{MODCLOSURE}}{\mathbf{functor}(X : \mathcal{M})M \xRightarrow{\rho} \mathbf{functor}(\rho)(X : \mathcal{M})M, \langle \rangle}$	$\frac{\text{MODCONSTR} \quad M \xRightarrow{\rho} V, \theta}{(M : \mathcal{M}) \xRightarrow{\rho} V, \theta}$		
$\frac{\text{MODBETA} \quad M \xRightarrow{\rho} \mathbf{functor}(\rho')(X : \mathcal{M})M_f, \theta \quad M' \xRightarrow{\rho} V', \theta' \quad M_f \xRightarrow{\rho' + \{X \mapsto V'\}} V'', \theta''}{M(M') \xRightarrow{\rho} V'', \theta @ \theta' @ \theta''}$			
$\frac{\text{TYPEDECL} \quad S \xRightarrow{\rho} V_s, \theta}{(\mathbf{type} \ (\alpha^*)t_i = \tau; S) \xRightarrow{\rho} V_s, \theta}$	$\frac{\text{MODULEDECL} \quad M \xRightarrow{\rho} V, \theta \quad S \xRightarrow{\rho + \{X_i \mapsto V\}} V_s, \theta'}{(\mathbf{module} \ X_i = M; S) \xRightarrow{\rho} \{X \mapsto V\} + V_s, \theta @ \theta'}$		
$\frac{\text{VALDECL} \quad e \xRightarrow{\rho} v, \theta \quad S \xRightarrow{\rho + \{x \mapsto v\}} V_s, \theta'}{(\mathbf{let} \ x_i = e; S) \xRightarrow{\rho} \{x \mapsto v\} + V_s, \theta'}$	$\frac{\text{PROGRAM} \quad S \xRightarrow{\rho} V_s, \theta}{\mathbf{prog} \ S \ \mathbf{end} \xRightarrow{\rho} V_s(\mathbf{return}), \theta}$		

Figure 3.12: Big step semantics –  $e \xRightarrow{\rho} v, \theta$

### 3.3.3 Notes on Soundness

Soundness properties, which correspond to the often misquoted “Well typed programs cannot go wrong.”, have been proven for many variants of the ML language. Unfortunately, stating and proving the soundness property for big step semantics and ML modules requires a fairly large amount of machinery which we do not attempt to provide in this thesis. Instead, we give pointers to various relevant work containing such proofs.

Soundness for a small step semantics of our expression language is provided in [Wright and Felleisen \[1994\]](#). At a larger scale, [Owens \[2008\]](#) proves the soundness of a small step semantics for a very large portion of the OCAML expression language using the Locally Nameless Coq framework [[Aydemir et al., 2008](#)]. Soundness of a big step semantics has been proved and mechanized for several richer languages [[Amin and Rompf, 2017](#), [Tofte, 1988](#), [Owens et al., 2016](#), [Lee et al., 2007](#), [Garrigue, 2009](#)].

Unfortunately, as far as we are aware, soundness of Leroy’s module language with higher order applicative functors has not been proved directly and is a fairly delicate subject. The most recent work of interest is [Rossberg et al. \[2014\]](#), which presents an elaboration scheme from ML modules, including applicative OCAML-style modules, into System  $F_\omega$ . Soundness then relies on soundness of the elaboration (provided in the article) and soundness of System  $F_\omega$ . In this work, the applicative/generative behavior of functors is decided depending on its purity, which is much more precise than what is done in OCAML. Notes that our language does not contain side-effects, which means that our language has a chance to be sound. The same cannot be said about OCAML in general.

## 3.4 Related works

Our formalization of the expression language is inspired by [Wright and Felleisen \[1994\]](#) which contains a small step semantics for the core language and extensions for references and exceptions, along with their soundness proofs. A discussion around the various styles of semantics in the context of traces is also provided by [Nakata and Uustalu \[2009\]](#). The big step semantics takes inspiration from [Owens et al. \[2016\]](#) and [Amin and Rompf \[2017\]](#). [Tofte \[1988\]](#) also presents a big-step semantics for ML and proves the soundness in the context of polymorphic type inference. A gentle introduction to soundness for big step semantics can be found in [Siek \[2013\]](#).

The landscape of ML modules is extremely rich, we only point to several key work that are directly relevant. Our module language is almost directly taken from [Leroy \[1994, 1995\]](#) which is the basis of the OCAML module language [[Leroy et al., 2016](#)]. The SML module system employs a different mechanism for propagating type equalities based on a **sharing** annotation [[Milner et al., 1990](#)]. Such annotation does not lend itself immediately to separate compilation [[Swasey et al., 2006](#)]. However, the SML module system has been (partially) mechanized [[Lee et al., 2007](#)]. A thorough and very instructive comparison between the various modules systems was done by [Dreyer \[2005, Chapter 1 and 2\]](#).

The module system we presented here is expressed directly in term of the syntax of the language. This is on purpose for two reasons. First, this is how it is implemented



in OCAML, making it possible to reason both about the formal system and the eventual implementation. Second, it makes the various extensions of ELIOM fairly easy to model. In particular, we shall see that our compilation scheme for ELIOM preserves the typing relation. Such results would be more delicate to express if the module system were expressed by elaboration. However, if such constraints were not considered, [Rossberg et al. \[2014\]](#) define an ML module system, including first class modules and applicative functors, in term of System  $F_\omega$ . This yields a simple yet feature-full system that could prove easier to extend.



## 4 The Eliom programming language

*Write a paper promising salvation, make it a 'structured' something or a 'virtual' something, or 'abstract', 'distributed' or 'higher-order' or 'applicative' and you can almost be certain of having started a new cult.*

Edsger W. Dijkstra, *My hopes of computing science*

We now present the formalization of ELIOM, a high-order applicative language for distributed client-server programming with strong support for abstraction and structured programming through modules<sup>1</sup>. The goal of this formalization is not to capture the complete ELIOM language, nor the larger OCSIGEN framework. Instead, we simply aim to capture the specific additions to the OCAML language from a typing and execution perspective. As such, we propose a new tierless calculus: ELIOM<sub>ε</sub>. In the ELIOM<sub>ε</sub> language, programs are series of bindings returning a single client value which symbolizes the web page showed to the user. In particular, we do not try to capture the non-terminating aspect of a complete web server, nor the back-and-forth interactions of a web browser making HTTP requests to a server.

This might seem like a very limited formalism. However, when a browser requests a web page, each request follows a similar pattern: an HTTP request is made, the server executes a specific handler to answer this HTTP request and to send a web page and a client program to the browser, this client program is then executed. The cycle starts again when the user clicks on another link. Each handler can then be considered as its own individual program. In the practical context of ELIOM and OCSIGEN, this little program is executed in the context of a larger program: the web server, which can contain some state and be non-terminating. Our formalization, however, focuses on modeling the typing, execution and compilation of each little program that are run in handlers<sup>2</sup>. ELIOM<sub>ε</sub> nevertheless captures several difficult points of ELIOM, namely the handling of distinct type universes, the transmission of values between server and client stages, the detailed semantics, its interaction with side effects, the module system in the presence of stages, and a compilation process that supports separate compilation. This should hopefully keep most readers entertained and give us a stable foothold before attempting an implementation. While ELIOM is an extension of OCAML, ELIOM<sub>ε</sub> is an extension of the simple ML calculus with modules presented in Chapter 3. We first present the ELIOM<sub>ε</sub> language in Section 4.1 and its type system in Section 4.2. We then present the semantics in Section 4.3 before stating various useful properties in Section 4.4

---

<sup>1</sup>Donations to the Cult of the Oxygenated Camel can be made by joining the OCAML consortium.

<sup>2</sup>Or, said in another way, ELIOM<sub>ε</sub> is the calculus of modular tierless CGI scripts.

## 4.1 Syntax

We now present  $\text{ELIOM}_\varepsilon$ , a core language for ELIOM.  $\text{ELIOM}_\varepsilon$  is an extension of the minimal ML language presented in Chapter 3, with both an expression and a module language. To emphasize the new elements introduced by ELIOM, these additional elements will be colored in blue. This is only for ease of reading and is not essential for understanding the formalization.

The syntax is presented in Figure 4.4. It extends the syntax presented in Figure 3.1.

### 4.1.1 Locations

Before describing the syntax of  $\text{ELIOM}_\varepsilon$ , let us introduce the notation of *locations*. The grammar of locations is given in Figure 4.1. A location is “a place where the code runs”. There are three core locations: **server**, **client** or **base**. The base side represents expressions that are “location-less”, that is, which can be used everywhere. We use the meta-variable  $\ell$  for an unspecified core location. There is a forth location that is only available for modules: **mixed**. A mixed module can have client, server and base components. We use the meta-variable  $\varsigma$  for locations that are either **m** or one of the core locations. In most contexts, locations are annotated with subscripts.

We also introduce two relations.  $\varsigma \succ \varsigma'$ , defined in Figure 4.2, means that a variable (either a value, a type or a module) defined on location  $\varsigma$  can be used on a location  $\varsigma'$ . For example, a base type can be used in a client context. Base declarations are usable everywhere, while mixed declarations are not usable in base code.  $\varsigma <: \varsigma'$ , defined in Figure 4.3, means that a module defined on location  $\varsigma$  can contain component on location  $\varsigma'$ . In particular, the mixed location  $m$  can contain any component, while other location can contain only component declared on the same location. Note that both relations are reflexive: For instance, it is always possible to use client declarations when you are on the client.

$$\ell ::= s \mid c \mid b \qquad \varsigma ::= m \mid \ell$$

Figure 4.1: Grammar of locations –  $\ell$  and  $\varsigma$

$$m \succ s \quad m \succ c \quad b \succ s \quad b \succ c \quad b \succ m \quad \forall \varsigma \in \{s, c, m, b\} \varsigma \succ \varsigma$$

Figure 4.2: “can be used in” relations on locations –  $\ell \succ \ell'$

$$m <: s \quad m <: c \quad m <: b \quad \forall \varsigma \in \{s, c, m, b\} \varsigma <: \varsigma$$

Figure 4.3: “can contain” relation on locations –  $\varsigma <: \varsigma'$

Expressions		Type Expressions	
$e ::= \dots$		$\tau ::= \dots$	
$\mid \{\{ e \}\}$	(Fragment)	$\mid \alpha_\ell$	(Type variables)
$\mid f\%v$	(Injection)	$\mid \{\tau\}$	(Fragment types)
$f ::= p.x \mid x_i \mid c$	(Converter)	$\mid \tau \rightsquigarrow \tau$	(Converter types)
Module Expressions		Module types	
$m ::= \dots$		$M ::= \dots$	
$\mid \text{functor}_m(X_i : \mathcal{M})M$	(Mixed functor)	$\mid \text{functor}_m(X_i : \mathcal{M}_1)\mathcal{M}_2$	(Mixed functor)
Structure components		Signature components	
$d ::= \text{let}_\ell x_i = e$		$D ::= \text{val}_\ell x_i : \tau$	
$\mid \text{type}_\ell (\alpha_{\ell_j}^*)t_i = \tau$		$\mid \text{type}_\ell (\alpha_{\ell_j}^*)t_i = \tau$	
$\mid \text{module}_\varsigma X_i = M$		$\mid \text{type}_\ell (\alpha_{\ell_j}^*)t_i$	
		$\mid \text{module}_\varsigma X_i : \mathcal{M}$	

Figure 4.4: ELIOM<sub>ε</sub>'s grammar

#### 4.1.2 Expression language

The expression language is extended with two new constructs: fragments and injections. A *client fragment*  $\{\{ e \}\}$  can be used on the server to represent an expression that will be computed on the client, but whose future value can be manipulated on the server. An *injection*  $f\%v$  can be used on the client to access values defined on the server. An injection must make explicit use of a converter  $f$  that specifies how to send the value. In particular, this should involve a serialization step, executed on the server, followed by a deserialization step executed on the client. For ease of presentation, injections are only done on variables and constants. In the implementation, this restriction is removed by adding a lifting transformation. For clarity, we sometimes distinguish injections *per se*, which occur outside of fragments, and escaped values, which occur inside fragments.

The syntax of types is also extended with two constructs. A *fragment type*  $\{\tau\}$  is the type of a fragment. A *converter type*  $\tau_s \rightsquigarrow \tau_c$  is the type of a converter taking a server value of type  $\tau_s$  and returning a client value of type  $\tau_c$ . All type variables  $\alpha_\ell$  are annotated with a core location  $\ell$ . There are now three sets of constants: client, server and base.

#### 4.1.3 Module language

The main change in the module language of ELIOM<sub>ε</sub> is that structure and signature components are annotated with locations. Value and type declarations can be annotated with a core location  $\ell$  which is either  $b$ ,  $s$  or  $c$ . Module declarations can also have one additional possible location: the mixed location  $m$ . Only modules on location

$m$  can have subfields on different locations. We also introduce mixed functors, noted  $\text{functor}_m(X:\mathcal{M})\mathcal{M}$ , which body can contain both client and server declarations. A program is a list of declarations including a client value declaration **return** which is the result of the program.

## 4.2 Type system

Judgments are now annotated with a location that specifies where the given code should be typechecked. Judgments on the expression language can accept any core location  $\ell$  while module judgments accept mixed locations  $\varsigma$ . We note  $\text{TypeOf}_\ell(c)$  the type of a given constant  $c$  on the location  $\ell$ . Binding in typing environments, just like in signatures, are annotated with a location. The first three kind of bindings, corresponding to the core language, can only appear on core locations:  $s$ ,  $c$  or  $b$ . Modules can also be of mixed location  $m$ . Names are namespaced by locations, which means it is valid to have both a client and a server value with both the same name.

### 4.2.1 Expressions

The new typing rules for expressions are presented in Figures 4.5 to 4.7. We introduce two typing rules for the new constructions. Rule **FRAGMENT** is for the construction of client fragments and can only be applied on the server. If  $e$  is of type  $\tau$  *on the client*, then  $\{\{ e \}\}$  is of type  $\{\tau\}$  *on the server*. Rule **INJECTION** is for the communication from the server to the client and can only be applied on the client. If  $e$  is of type  $\tau_s$  on the server and  $f$  is of type  $\tau_s \rightsquigarrow \tau_c$  on the server, then  $f\%e$  is of type  $\tau_c$  on the client. Since no other typing rules involves client fragments, it is impossible to deconstruct them.

The last difference with usual ML rules are the visibility of variable. As described earlier, bindings in  $\text{ELIOM}_m$  are located. Since access across sides are explicit, we want to prevent the use of client variables on the server, for example. In rule **VAR**, to use on location  $\ell$  the variable  $v$  which is bound on location  $\ell'$ , we check that  $\ell' \succ \ell$ , defined in Figure 4.2, which means that the definition on  $\ell$  can be used in  $\ell'$ . Base elements  $b$  are usable everywhere. Mixed elements  $m$  are usable in both client and server. Type variables are also annotated with a location and follow the same rules. Using type variables from the client on the server, for example, is disallowed.

The validity judgement on types presented in Figure 4.6 is extended to check that locations are respected both for type constructors and type variables. This judgement is used in type declaration, which are presented in the module system.

### Converters

To transmit values from the server to the client, we need a serialization format. We assume the existence of a type **serial** in  $\text{Const}_b$  which represents the serialization format. The actual format is irrelevant. For instance, one could use JSON or XML.

Converters are special values that describe how to move a value from the server to the client. A converter can be understood as a pair of functions. A converter  $f$  of type

### Common rules

$$\begin{array}{c}
\text{VAR} \\
\frac{(\text{val}_{\ell'} x : \sigma) \in \Gamma \quad \ell' \succ \ell \quad \sigma \succ \tau}{\Gamma \triangleright_{\ell} x : \tau}
\end{array}
\quad
\begin{array}{c}
\text{LAM} \\
\frac{\Gamma; (\text{val}_{\ell} x : \tau_1) \triangleright_{\ell} e : \tau_2}{\Gamma \triangleright_{\ell} \lambda x. e : \tau_1 \rightarrow \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{CONST} \\
\frac{\text{TypeOf}_{\ell}(c) \succ \tau}{\Gamma \triangleright_{\ell} c : \tau}
\end{array}$$

$$\begin{array}{c}
\text{LETIN} \\
\frac{\Gamma \triangleright_{\ell} e_1 : \tau_1 \quad \Gamma; (\text{val}_{\ell} x : \text{Close}(\tau_1, \Gamma)) \triangleright_{\ell} e_2 : \tau_2}{\Gamma \triangleright_{\ell} \text{let } x = e_1 \text{ in } e_2 : \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{EQUIV} \\
\frac{\Gamma \triangleright_{\ell} e : \tau_1 \quad \Gamma \triangleright_{\ell} \tau_1 \approx \tau_2}{\Gamma \triangleright_{\ell} e : \tau_2}
\end{array}$$

$$\begin{array}{c}
\text{APP} \\
\frac{\Gamma \triangleright_{\ell} e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \triangleright_{\ell} e_2 : \tau_1}{\Gamma \triangleright_{\ell} (e_1 \ e_2) : \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{Y} \\
\frac{}{\Gamma \triangleright_{\ell} \text{Y} : ((\tau_1 \rightarrow \tau_2) \rightarrow \tau_1 \rightarrow \tau_2) \rightarrow \tau_1 \rightarrow \tau_2}
\end{array}$$

$$\begin{array}{c}
\text{QUALVAR} \\
\frac{\Gamma \blacktriangleright_{\ell} p : (\text{sig } S_1; \text{val}_{\ell'} x_i : \tau; S_2 \text{ end}) \quad \ell' \succ \ell}{\Gamma \triangleright_{\ell} p.v : \tau[n_i \mapsto_{\ell} p.n \mid n_i \in \text{BV}_{\ell}(S_1)]}
\end{array}$$

### Server rules

$$\begin{array}{c}
\text{FRAGMENT} \\
\frac{\Gamma \triangleright_c e : \tau}{\Gamma \triangleright_s \{\{ e \}\} : \{\tau\}}
\end{array}$$

### Client rules

$$\begin{array}{c}
\text{INJECTION} \\
\frac{\Gamma \triangleright_s f : \tau_s \rightsquigarrow \tau_c \quad \Gamma \triangleright_s e : \tau_s}{\Gamma \triangleright_c f \% e : \tau_c}
\end{array}$$

$\text{Close}(\tau, \Gamma) = \forall \alpha_0 \dots \alpha_n. \tau$  with  $\{\alpha_0, \dots, \alpha_n\} = \text{FreeTypeVar}(\tau) \setminus \text{FreeTypeVar}(\Gamma)$

Figure 4.5: ELIOM<sub>ε</sub> expression typing rules –  $\Gamma \triangleright_{\ell} e : \tau$

$$\begin{array}{c}
\text{TYPEVAL} \\
\frac{(\text{type}_{\ell}(\alpha_{\ell_i})t) \in \Gamma \quad \forall i, \Gamma \models_{\ell_i} \tau_i}{\Gamma \models_{\ell} (\tau_i)t}
\end{array}
\quad
\begin{array}{c}
\text{ARROWVAL} \\
\frac{\Gamma \models_{\ell} \tau_1 \quad \Gamma \models_{\ell} \tau_2}{\Gamma \models_{\ell} \tau_1 \rightarrow \tau_2}
\end{array}
\quad
\begin{array}{c}
\text{VARVAL} \\
\frac{}{\Gamma \models_{\ell} \alpha_{\ell}}
\end{array}
\quad
\begin{array}{c}
\text{FRAGVAL} \\
\frac{}{\Gamma \models_s \{\tau\}}
\end{array}$$

$$\begin{array}{c}
\text{QUALIFIEDVAL} \\
\frac{\Gamma \blacktriangleright_{\ell} p : (\text{sig } S_1; \text{type}_{\ell'}(\alpha_{\ell_i})t; S_2 \text{ end}) \quad \forall i, \Gamma \models_{\ell_i} \tau_i \quad \ell' \succ \ell}{\Gamma \models_{\ell} (\tau_i)p.t}
\end{array}
\quad
\begin{array}{c}
\text{CONVVAL} \\
\frac{\Gamma \models_s \tau_1 \quad \Gamma \models_c \tau_2}{\Gamma \models_s \tau_1 \rightsquigarrow \tau_2}
\end{array}$$

Figure 4.6: Type validity rules –  $\Gamma \models_{\ell} \tau$

$$\begin{array}{c}
\text{REFLEQ} \quad \frac{}{\Gamma \triangleright_{\ell} \tau \approx \tau} \quad \text{TRANSEQ} \quad \frac{\Gamma \triangleright_{\ell} \tau_1 \approx \tau_2 \quad \Gamma \triangleright_{\ell} \tau_2 \approx \tau_3}{\Gamma \triangleright_{\ell} \tau_1 \approx \tau_3} \quad \text{COMMEQ} \quad \frac{\Gamma \triangleright_{\ell} \tau_2 \approx \tau_1}{\Gamma \triangleright_{\ell} \tau_1 \approx \tau_2} \quad \text{FUNEQ} \quad \frac{\Gamma \triangleright_{\ell} \tau_1 \approx \tau'_1 \quad \Gamma \triangleright_{\ell} \tau_2 \approx \tau'_2}{\Gamma \triangleright_{\ell} \tau_1 \rightarrow \tau_2 \approx \tau'_1 \rightarrow \tau'_2} \\
\\
\text{DEFTYPEEQ} \quad \frac{(\text{type}_{\ell'}(\alpha_{\ell_i})t = \tau) \in \Gamma \quad \ell' \succ_{\ell}}{\Gamma \triangleright_{\ell}(\tau_i)t \approx \tau[\alpha_i \mapsto_{\ell_i} \tau_i]_i} \quad \text{ABSTYPEEQ} \quad \frac{(\text{type}_{\ell'}(\alpha_{\ell_i})t) \in \Gamma \quad \forall i, \Gamma \triangleright_{\ell_i} \tau_i \approx \tau'_i \quad \ell' \succ_{\ell}}{\Gamma \triangleright_{\ell}(\tau_i)t \approx (\tau'_i)t} \\
\\
\text{QUALDEFTYPEEQ} \quad \frac{\Gamma \blacktriangleright_{\ell} p : (\text{sig } \mathcal{S}_1; \text{type}_{\ell'}(\alpha_{\ell_i})t = \tau; \mathcal{S}_2 \text{ end}) \quad \forall i, \Gamma \triangleright_{\ell_i} \tau_i \approx \tau'_i \quad \ell' \succ_{\ell}}{\Gamma \triangleright_{\ell}(\tau_i)p.t \approx \tau[n_i \mapsto_{\ell} p.n \mid n_i \in \text{BV}_{\ell}(\mathcal{S}_1)][\alpha_i \mapsto_{\ell_i} \tau_i]_i} \\
\\
\text{QUALABSTYPEEQ} \quad \frac{\Gamma \blacktriangleright_{\ell} p : (\text{sig } \mathcal{S}_1; \text{type}_{\ell'}(\alpha_{\ell_i})t; \mathcal{S}_2 \text{ end}) \quad \forall i, \Gamma \triangleright_{\ell_i} \tau_i \approx \tau'_i \quad \ell' \succ_{\ell}}{\Gamma \triangleright_{\ell}(\tau_i)p.t \approx (\tau'_i)p.t} \\
\\
\text{FRAGMENTEQ} \quad \frac{\Gamma \triangleright_{\mathbf{c}} \tau \approx \tau'}{\Gamma \triangleright_{\mathbf{s}} \{\tau\} \approx \{\tau'\}} \quad \text{CONVEQ} \quad \frac{\Gamma \triangleright_{\mathbf{s}} \tau_s \approx \tau'_s \quad \Gamma \triangleright_{\mathbf{c}} \tau_c \approx \tau'_c}{\Gamma \triangleright_{\mathbf{s}} \tau_s \rightsquigarrow \tau_c \approx \tau'_s \rightsquigarrow \tau'_c}
\end{array}$$

Figure 4.7: Type equivalence rules –  $\Gamma \triangleright_{\ell} \tau \approx \tau'$

$\tau_s \rightsquigarrow \tau_c$  is composed of a server-side encoding function of type  $\tau_s \rightarrow \mathbf{serial}$ , and a client-side decoding function of type  $\mathbf{serial} \rightarrow \tau_c$ . We assume the existence of two built-in converters:

- The **serial** converter of type  $\mathbf{serial} \rightsquigarrow \mathbf{serial}$ . Both sides are the identity.
- The **frag** converter of type  $\forall \alpha_c. (\{\alpha_c\} \rightsquigarrow \alpha_c)$ .

## Type universes

It is important to note that there is no identity converter (of type  $\forall \alpha. (\alpha \rightsquigarrow \alpha)$ ). Indeed the client and server type universes are distinct and we cannot translate arbitrary types from one to the other. Some types are only available on one side: database handles, system types, JAVASCRIPT API types. Some types, while available on both sides (because they are in base for example), are simply not transferable. For example, functions cannot be serialized in general. Another example is file handles: they are available both on the server and on the client, but moving a file handle from server to client seems adventurous.

Finally, some types may share a semantic meaning, but not their actual representation. This is the case where converters are used, as demonstrated in Section 2.3.



## Mixed datatypes

In Sections 3.2.1 and 3.2.2, we saw that the version of ML we consider supports an interesting combination of three features: abstract datatypes, parametrized datatypes and separate compilation at the module level.  $\text{ELIOM}_\varepsilon$ , as an extension of ML, also supports these features. These three features have non-trivial interactions that need to be accounted for, in particular when introducing properties on types, such as locations.

Let us consider the module shown in Example 4.1. We declare a server datatype  $\mathbf{t}$  with two parameters and we hide the definition in the signature. We now want to check that  $(\mathbf{t1}, \mathbf{t2})\mathbf{t}$  is a correct type expressions. However, without the type definition, we don't know if  $\mathbf{t1}$  and  $\mathbf{t2}$  are base, client or server types. In order to type check the type sub-expressions, we need more information about the definition of  $\mathbf{t}$ . The solution, much like variance, is to annotate type variables in datatypes with extra information. This is done in the syntax for type declarations given in Figure 4.4. Each type parameters is annotated with a location. Type variables can only be used on the right location. This ensures proper separation of client and server type variables and their proper usage.

```
1 module M : sig
2   type%server ('a, 'b) t
3 end = struct
4   type%server ('a, 'b) t =
5     'a fragment * 'b
6 end
```

(a) Incorrect abstract datatype

```
1 module M : sig
2   type%server ('a[@client], 'b) t
3 end = struct
4   type%server ('a[@client], 'b) t =
5     'a fragment * 'b
6 end
```

(b) Correct abstract datatype

Example 4.1: A module with an abstract datatype.

Example 4.1a does not exposes information about acceptable sides for  $\mathbf{'a}$  and  $\mathbf{'b}$ . In Example 4.1b, annotations specifying the side of type variables are exposed in the interface.

## 4.2.2 Modules

We now detail ELIOM's module system. In the expression language, location transitions can only happen in very specific constructions: fragments and injections. This allow us to keep most of the ML type system unchanged. This is not the case anymore for modules: we allow users to create base, client and server modules, but also mixed modules that can contain base, client and server declarations, including other modules. This means we need to track locations quite precisely.

We first introduce the various feature of our module system along with some motivating examples. We then detail how those features are enforced by the typing rules.

### Base location and specialization

In Section 2.4, we presented an example where a base functor  $\mathbf{Map.Make}$ , is applied to a client module to obtain a new client module. As  $\mathbf{Map.Make}$  is a module provided by the

standard library of OCAML, it is defined on location  $b$ . In particular, its input signature has components on location  $b$ , thus it would seem a module whose components are on the client or the server should not be accepted. We would nevertheless like to create maps of elements that are only available on the client. To do so, we introduce a specialization operation, defined in Figure 4.8, that allows to use a base module in a client or server scope by replacing instances of the base location with the current location.

The situation is quite similar to the application of a function of type  $\forall\alpha.\alpha \rightarrow \alpha$  to an argument of type *int*: we need to instantiate the function before being able to use it. Mixed modules only offer a limited version of polymorphism for locations: there is only one “location variable” at a time, and it’s always called  $b$ . The specialization operation simply rewrites a module signature by substituting all instances of the location  $b$  or  $m$  by the specified  $c$  or  $s$  location. Note that before being specialized, a module should be accessible according to the “can be used” relation defined Figure 4.2. This means that we never have to specialize a server module on the client (or conversely). Specialization towards location  $b$  has no effect since only base modules are accessible on location base. Specialization towards the location  $m$  has no effect either: since all locations are allowed inside the mixed location, no specialization is needed. Mixed functors are handled in a specific way, as we see in the next section.

## Mixed Functors

Mixed functors are functors declared in a mixed scope. We note  $\mathbf{functor}_m(X_i : \mathcal{M})M$  the mixed functor that takes an argument  $X_i$  of type  $\mathcal{M}$  and return a module  $M$ . They can contain both client and server declarations (or mixed submodules). Mixed functors and regular functors have different types that are not compatible. We saw in Section 2.6.2 an example of usage for mixed functors. Mixed functors have several restrictions compared to regular functors which we now detail using various examples.

**Specialization** A naive implementation of specialization of mixed functors would be to specialize on both side of the arrow and apply the resulting functor. Let us see on an example why this solution does not work. In Example 4.2, the functor  $\mathbf{F}$  takes as argument a module containing a base declaration and uses it on both sides. If the type of the functor parameter were specialized, the functor application in Example 4.2b would be well-typed. However, this makes no sense:  $\mathbf{M.y}$  is supposed to represent a fragment whose content is the client value of  $b$ , but this value doesn’t exist, since  $b$  was declared on the server. There would be no value available to inject in the declaration of  $\mathbf{y'}$ .

The solution here is that specialization on mixed functors should only specialize the return type, not the argument.

**Injections** Injections inside client sections (as opposed to escaped values inside client fragments) are fairly static: the value might be dynamic, but the position of the injection and its use sites are statistically known and does not depend on the execution of the program. In particular, injections are independent of the control flow. We can just give a unique identifier to each injection, and use that unique name for lookup on the client.

```

1 module%mixed F (A : sig val b : int end)
2 = struct
3   let%server x = A.b
4   let%server y = [%client A.b]
5 end

```

(a) A mixed functor using a base declaration

```

1 module%server M =
2   F(struct let%server b = 2 end)
3 let%client y' = ~%M.y

```

(b) An ill-typed application of F

Example 4.2: A mixed functor using base declaration polymorphically

This property comes from the fact that injected server identifiers cannot be bound in a client section.

Unfortunately, this property does not hold in the presence of mixed functor when we assume the language can apply functor at arbitrary positions, which is the case in OCAML. Let us consider Example 4.3. The functor  $F$  takes a structure containing a server declaration  $x$  holding an integer and returns a structure containing the same integer, injected in the client. In Example 4.3b, the functor is used on  $A$  or  $B$  conditionally. The issue is that the client integer depends both on the server integer and on the local *client* control flow. Lifting the functor application at toplevel would not preserve the semantics of the language, due to side effects. Thus, we avoid this kind of situation by forbidding injections that access dynamic names inside mixed functors.

```

1 module%mixed F
2   (A : sig val%server x : int end)
3 = struct
4   let%client x' = ~%A.x
5 end

```

(a) An problematic mixed functor with an injection

```

1 module%mixed A = struct let%server x = 2 end
2 module%mixed B = struct let%server x = 4 end
3 let%client a =
4   if Random.bool ()
5   then let module M = F(A) in M.x'
6   else let module M = F(B) in M.x'

```

(b) A pathological functor application

Example 4.3: Problematic example of injection inside a mixed functor

In order to avoid this situation, we add the constraints that injections inside the body of a mixed functors can only refer to outside of the functor. Escaped values, which are injections inside client fragments, are still allowed. The functor presented in Example 4.4a is not allowed while the one in Example 4.4b is allowed. Formally, this is guaranteed by the MIXEDFUNCTOR rule, where each injection is typechecked in the outer typing environment.

**Functor application** Mixed functors can only be applied to mixed structures. This means that in a functor application  $F(M)$ ,  $M$  must be a structure defined by a `modulem` declaration. Note that this breaks the property that the current location of an expression or a module can be determined syntactically: The location inside `F(struct ... end)` can be either mixed or not, depending on  $F$ . This could be mitigated by using a different syntax for the application of mixed functor. The justification for this restriction is detailed in Section 4.3.

```

1
2 module%mixed F
3   (A:sig val%server x : int end)
4 = struct
5   let%client y = ~%A.x + 2
6 end
7

```

(a) An ill-typed mixed functor using an injection

```

1 let%server x = 3
2 module%mixed F
3   (A:sig val%server y : int end)
4 = struct
5   let%client z = ~%x
6   let%server z' = [%client ~%A.y + 1]
7 end

```

(b) A well-typed mixed functor using an injection

Example 4.4: Mixed functor and injections

## Type rules

We now review how these various language constructs are reflected in the rules of our type system. As before, the ELIOM module system is built on the ML module system. We extend the typing, validity and subtyping judgments by a location annotation that specifies the location of the current scope. The program typing judgments don't have a location, since a program is always considered mixed. Most rules are fairly straightforward adaptations of the ML rules, annotated with locations.

The typing rules MODVAR and QUALMODVAR follow the usual rules of ML modules with two modifications: We first check that the module we are looking up can indeed be used on the current location. This is done by the side condition  $\varsigma' \succ \varsigma$  where  $\varsigma$  is the current location and  $\varsigma'$  is the location where the identifier is defined. This allows, for instance, to use *base* identifiers in a *client* scope. We also specialize the module type of the identifier towards the current location  $\varsigma$ . The specialization operation, which was described in Section 4.2.2, is noted  $[M]_\varsigma$  and is defined in Figure 4.8.

There are two new typing rules compared to ML: the rules MIXEDFUNCTOR and MIXEDAPPLICATION define mixed functor definition and application. We use INJS( $\cdot$ ) which returns the set of all injections in client declarations.

## Subtyping and equivalence of modules

Subtyping rules are given in Figure 4.11. For brevity, we note  $\varsigma <: (\varsigma_1 \succ \varsigma_2)$  as a shorthand for  $\varsigma <: \varsigma_1 \wedge \varsigma <: \varsigma_2 \wedge \varsigma_1 \succ \varsigma_2$ , that is, both  $\varsigma_1$  and  $\varsigma_2$  are valid locations for components of a module on location  $\varsigma$  and location  $\varsigma_1$  encompasses location  $\varsigma_2$ . Note that the following holds:

$$\Gamma \triangleright_\varsigma \text{struct val}_b t_i : \text{int end} <: \text{struct val}_c t_i : \text{int end}$$

This is perfectly safe, since for any identifier  $x_i$  on base,  $\text{let}_c x'_j = x_i$  is always valid. This allows programmers to declare some code on base (and get the guarantee that the code is only using usual OCAML constructs) but to expose it as client or server in the module type.

$$\begin{aligned}
\llbracket \mathcal{M} \rrbracket_b &= \mathcal{M} & \llbracket \mathcal{M} \rrbracket_m &= \mathcal{M} \\
\llbracket \text{sig } \mathcal{S} \text{ end} \rrbracket_\iota &= \text{sig } \llbracket \mathcal{S} \rrbracket_\iota \text{ end} & \llbracket \text{functor}_m(X_i : \mathcal{M}) \mathcal{M}' \rrbracket_\iota &= \text{functor}_m(X_i : \mathcal{M}) \llbracket \mathcal{M}' \rrbracket_\iota \\
\llbracket \varepsilon \rrbracket_\iota &= \varepsilon & \llbracket \text{functor}(X_i : \mathcal{M}) \mathcal{M}' \rrbracket_\iota &= \text{functor}(X_i : \llbracket \mathcal{M} \rrbracket_\iota) \llbracket \mathcal{M}' \rrbracket_\iota \\
\llbracket \text{val}_\ell x_i : \tau; \mathcal{S} \rrbracket_\iota &= \begin{cases} \text{val}_\iota x_i : \tau; \llbracket \mathcal{S} \rrbracket_\iota & \text{when } \ell \succ \iota \\ \llbracket \mathcal{S} \rrbracket_\iota & \text{otherwise} \end{cases} \\
\llbracket \text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i = \tau; \mathcal{S} \rrbracket_\iota &= \begin{cases} \text{type}_\iota (\alpha_{\ell_j}^*) \text{t}_i = \tau; \llbracket \mathcal{S} \rrbracket_\iota & \text{when } \ell \succ \iota \\ \llbracket \mathcal{S} \rrbracket_\iota & \text{otherwise} \end{cases} \\
\llbracket \text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i; \mathcal{S} \rrbracket_\iota &= \begin{cases} \text{type}_\iota (\alpha_{\ell_j}^*) \text{t}_i; \llbracket \mathcal{S} \rrbracket_\iota & \text{when } \ell \succ \iota \\ \llbracket \mathcal{S} \rrbracket_\iota & \text{otherwise} \end{cases} \\
\llbracket \text{module}_\varsigma X_i : \mathcal{M}; \mathcal{S} \rrbracket_\iota &= \begin{cases} \text{module}_\iota X_i : \llbracket \mathcal{M} \rrbracket_\iota; \llbracket \mathcal{S} \rrbracket_\iota & \text{when } \varsigma \succ \iota \\ \llbracket \mathcal{S} \rrbracket_\iota & \text{otherwise} \end{cases}
\end{aligned}$$

Where  $\iota$  is either  $c$  or  $s$ .

Figure 4.8: Module specialization operation –  $\llbracket M \rrbracket_\varsigma$

$$\begin{aligned}
\varepsilon/p &= \varepsilon \\
(\text{sig } \mathcal{S} \text{ end})/p &= \text{sig } \mathcal{S}/p \text{ end} \\
(\text{module}_\varsigma X_i = \mathcal{M}; \mathcal{S})/p &= \text{module}_\varsigma X_i = \mathcal{M}/p; \mathcal{S}/p \\
(\text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i = \tau; \mathcal{S})/p &= \text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i = (\alpha^*) \text{p.t}; \mathcal{S}/p \\
(\text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i; \mathcal{S})/p &= \text{type}_\ell (\alpha_{\ell_j}^*) \text{t}_i = (\alpha^*) \text{p.t}; \mathcal{S}/p \\
(\text{val}_\ell x_i : \tau; \mathcal{S})/p &= \text{val}_\ell x_i : \tau; \mathcal{S}/p \\
(\text{functor}(X_i : \mathcal{M}) \mathcal{M}')/p &= \text{functor}(X_i : \mathcal{M}) (\mathcal{M}'/p(X_i)) \\
(\text{functor}_m(X_i : \mathcal{M}) \mathcal{M}')/p &= \text{functor}_m(X_i : \mathcal{M}) (\mathcal{M}'/p(X_i))
\end{aligned}$$

Figure 4.9: Module strengthening operation –  $M/p$

$$\begin{array}{c}
\text{MODVAR} \\
\frac{(\text{module}_{\zeta'} X_i : \mathcal{M}) \in \Gamma \quad \zeta' \succ \zeta}{\Gamma \blacktriangleright_{\zeta} X_i : \lfloor \mathcal{M} \rfloor_{\zeta}} \\
\\
\text{QUALMODVAR} \qquad \text{STRENGTH} \\
\frac{\Gamma \blacktriangleright_{\zeta} p : (\text{sig } \mathcal{S}_1; \text{module}_{\zeta'} X_i : \mathcal{M}; \mathcal{S}_2 \text{ end}) \quad \zeta' \succ \zeta}{\Gamma \blacktriangleright_{\zeta} p.X : \lfloor \mathcal{M}[n_i \mapsto_{\zeta'} p.n \mid n_i \in \text{BV}_{\zeta'}(\mathcal{S}_1)] \rfloor_{\zeta}} \quad \frac{\Gamma \blacktriangleright_{\zeta} p : \mathcal{M}}{\Gamma \blacktriangleright_{\zeta} p : \mathcal{M}/p} \\
\\
\frac{\Gamma \blacktriangleright_{\zeta} M : \mathcal{M}' \quad \Gamma \blacktriangleright_{\zeta} \mathcal{M}' <: \mathcal{M}}{\Gamma \blacktriangleright_{\zeta} M : \mathcal{M}} \quad \frac{\Gamma \blacktriangleright_{\zeta} M_1 : \text{functor}(X_i : \mathcal{M})\mathcal{M}' \quad \Gamma \blacktriangleright_{\zeta} M_2 : \mathcal{M}}{\Gamma \blacktriangleright_{\zeta} M_1(M_2) : \mathcal{M}'[X_i \mapsto_{\zeta} M_2]} \\
\\
\frac{\Gamma \models_{\ell} \mathcal{M} \quad X_i \notin \text{BV}_{\ell}(\Gamma) \quad \Gamma; (\text{module}_{\ell} X_i : \mathcal{M}) \blacktriangleright_{\ell} M : \mathcal{M}'}{\Gamma \blacktriangleright_{\ell} \text{functor}(X_i : \mathcal{M})M : \text{functor}(X_i : \mathcal{M})\mathcal{M}'} \\
\\
\frac{\Gamma \models_{\zeta} \mathcal{M} \quad \Gamma \blacktriangleright_{\zeta} M : \mathcal{M}}{\Gamma \blacktriangleright_{\zeta} (M : \mathcal{M}) : \mathcal{M}} \\
\\
\text{MIXEDFUNCTOR} \\
\frac{\Gamma \models_m \text{sig } \mathcal{S} \text{ end} \quad x_i \notin \text{BV}_m(\Gamma) \quad \Gamma; (\text{module}_m X_i : \text{sig } \mathcal{S} \text{ end}) \blacktriangleright_m M : \mathcal{M}' \quad \forall f_j \% X_j \in \text{INJS}(M), \Gamma \triangleright_c f_j \% X_j : \tau_j}{\Gamma \blacktriangleright_m \text{functor}_m(X_i : \text{sig } \mathcal{S} \text{ end})M : \text{functor}_m(X_i : \text{sig } \mathcal{S} \text{ end})\mathcal{M}'} \\
\\
\text{MIXEDAPPLICATION} \\
\frac{\Gamma \blacktriangleright_{\zeta} M_1 : \text{functor}_m(X_i : \mathcal{M})\mathcal{M}' \quad \Gamma \blacktriangleright_m M_2 : \mathcal{M} \quad m \succ \zeta}{\Gamma \blacktriangleright_{\zeta} M_1(M_2) : \mathcal{M}'[X_i \mapsto_{\zeta} M_2]} \\
\\
\frac{\Gamma \triangleright_{\ell} e : \tau \quad x_i \notin \text{BV}_{\ell}(\Gamma) \quad \Gamma; (\text{val}_{\ell} x_i : \text{Close}(\tau, \Gamma)) \blacktriangleright_{\zeta} S : \mathcal{S} \quad \zeta <: \ell}{\Gamma \blacktriangleright_{\zeta} (\text{let}_{\ell} x_i = e; s) : (\text{val}_{\ell} x_i : \tau; \mathcal{S})} \\
\\
\frac{\Gamma \models_{\ell} \tau \quad t_i \notin \text{BV}_{\ell}(\Gamma) \quad \Gamma; (\text{type}_{\ell} (\alpha_{\ell_j}^*) t_i = \tau) \blacktriangleright_{\zeta} S : \mathcal{S} \quad \zeta <: \ell}{\Gamma \blacktriangleright_{\zeta} (\text{type}_{\ell} (\alpha_{\ell_j}^*) t_i = \tau; s) : (\text{type}_{\ell} (\alpha_{\ell_j}^*) t_i = \tau; \mathcal{S})} \\
\\
\frac{\Gamma \blacktriangleright_{\zeta} M : \mathcal{M} \quad X_i \notin \text{BV}_{\zeta}(\Gamma) \quad \Gamma; (\text{module}_{\zeta} X_i : \mathcal{M}) \blacktriangleright_{\zeta'} S : \mathcal{S} \quad \zeta' <: \zeta \quad \forall \zeta'' \in \text{locations}(\mathcal{M}). \zeta'' \succ \zeta}{\Gamma \blacktriangleright_{\zeta'} (\text{module}_{\zeta} X_i = M; s) : (\text{module}_{\zeta} X_i : \mathcal{M}; \mathcal{S})} \\
\\
\frac{\Gamma \blacktriangleright_{\zeta} S : \mathcal{S}}{\Gamma \blacktriangleright_{\zeta} \text{struct } S \text{ end} : \text{sig } \mathcal{S} \text{ end}} \quad \frac{}{\Gamma \blacktriangleright_{\zeta} \varepsilon : \varepsilon}
\end{array}$$

Figure 4.10: Module typing rules –  $\Gamma \blacktriangleright_{\zeta} m : M$

$$\begin{array}{c}
\text{SUBSTRUCT} \\
\frac{\pi : [1; m] \rightarrow [1; n] \quad \forall i \in [1; m], \Gamma; \mathcal{D}_1; \dots; \mathcal{D}_n \blacktriangleright_{\varsigma} \mathcal{D}_{\pi(i)} <: \mathcal{D}'_i}{\Gamma \blacktriangleright_{\varsigma} (\text{sig } \mathcal{D}_1; \dots; \mathcal{D}_n \text{ end}) <: (\text{sig } \mathcal{D}'_1; \dots; \mathcal{D}'_m \text{ end})} \\
\\
\frac{\Gamma \triangleright_{\ell_2} \tau_1 \approx \tau_2 \quad \varsigma <: (\ell_1 \succ \ell_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{val}_{\ell_1} x_i : \tau_1) <: (\text{val}_{\ell_2} x_i : \tau_2)} \\
\\
\frac{\Gamma \blacktriangleright_{\varsigma_2} \mathcal{M}_1 <: \mathcal{M}_2 \quad \varsigma <: (\varsigma_1 \succ \varsigma_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{module}_{\varsigma_1} X_i : \mathcal{M}_1) <: (\text{module}_{\varsigma_2} X_i = \mathcal{M}_2)} \\
\\
\frac{\Gamma \blacktriangleright_{\ell} \mathcal{M}'_a <: \mathcal{M}_a \quad \Gamma, (\text{module}_{\ell} X : \mathcal{M}'_a) \blacktriangleright_{\ell} \mathcal{M}_r <: \mathcal{M}'_r}{\Gamma \blacktriangleright_{\ell} \text{functor}(X : \mathcal{M}_a) \mathcal{M}_r <: \text{functor}(X : \mathcal{M}'_a) \mathcal{M}'_r} \\
\\
\frac{\Gamma \blacktriangleright_m \mathcal{M}'_a <: \mathcal{M}_a \quad \Gamma, (\text{module}_m X : \mathcal{M}'_a) \blacktriangleright_m \mathcal{M}_r <: \mathcal{M}'_r}{\Gamma \blacktriangleright_m \text{functor}_m(X : \mathcal{M}_a) \mathcal{M}_r <: \text{functor}_m(X : \mathcal{M}'_a) \mathcal{M}'_r} \\
\\
\frac{\Gamma \triangleright_{\ell_2} \tau_1 \approx \tau_2 \quad \varsigma <: (\ell_1 \succ \ell_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{type}_{\ell_1} (\alpha_{\ell_j}^*) t_i = \tau_1) <: (\text{type}_{\ell_2} (\alpha_{\ell_j}^*) t_i = \tau_2)} \\
\\
\frac{\varsigma <: (\ell_1 \succ \ell_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{type}_{\ell_1} (\alpha_{\ell_j}^*) t_i) <: (\text{type}_{\ell_2} (\alpha_{\ell_j}^*) t_i)} \quad \frac{\Gamma \triangleright_{\ell_2} (\alpha_{\ell_j}^*) t_i \approx \tau \quad \varsigma <: (\ell_1 \succ \ell_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{type}_{\ell_1} (\alpha_{\ell_j}^*) t_i) <: (\text{type}_{\ell_2} (\alpha_{\ell_j}^*) t_i = \tau)} \\
\\
\frac{\varsigma <: (\ell_1 \succ \ell_2)}{\Gamma \blacktriangleright_{\varsigma} (\text{type}_{\ell_1} (\alpha_{\ell_j}^*) t_i = \tau_1) <: (\text{type}_{\ell_2} (\alpha_{\ell_j}^*) t_i)}
\end{array}$$

Figure 4.11: Module subtyping rules –  $\Gamma \blacktriangleright_{\varsigma} M <: M'$

### 4.3 Interpreted semantics

While ELIOM, just like OCAML, is a compiled language, it is desirable to present a semantics that does not involve complex program transformation. The reason is two-fold: First, this simple semantics should be reasonably easy to explain to users. Indeed, this semantics is the one used to present ELIOM in Chapter 2. However, we must also show that this semantics is correct, in that it does actually corresponds to our compilation scheme. This is done in Section 5.4. As presented in Chapter 2, ELIOM execution proceeds in two steps: The server part of the program is executed first. This creates a client program, which is then executed.

Let us first introduce a few notations. Generated client programs are noted  $\mu$ . Server expressions (resp. declarations) that do not contain injections are noted  $\bar{e}$  (resp.  $\bar{D}$ ). Values are the same as for ML: constants, closures, structures and functor closures. We consider a new class of identifiers called “references” and noted in bold, such as **r** or **R**. We assume the existence of a name generator that can create arbitrary new fresh **r** identifiers at any point of the execution. References are used as global identifiers that ignore scoping rules. References can also be qualified as “reference paths”, noted **X.r**. This is used for mixed functors, in particular. We use  $\gamma$  to note the global environment where such references are stored.

We now introduce a new reduction relation,  $\Rightarrow_\varsigma$ , which is the reduction over ELIOM constructs on side  $\varsigma$ . The notation  $\Rightarrow_\varsigma$  actually represents several reduction relations which are presented in Figures 4.13, 4.16 and 4.17. Four of these relations reduce the server part of the code and emit a client program. We note  $e \xRightarrow{\rho}_\iota v, \mu, \theta$  the reduction of a server expression  $e$  inside a context  $\iota$  in the environment  $\rho$ . It returns the value  $v$ , the client program  $\mu$  and emits the trace  $\theta$ . The context  $\iota$  can be either base ( $b$ ), server ( $s$ ), server code inside client contexts ( $c/s$ ) or server code inside mixed contexts ( $m$ ). We also have a client reduction, noted  $e \xRightarrow{\rho | \gamma \rightarrow \gamma'}_c v, \theta$  which reduces a client expression  $e$  inside an environment  $\rho$ , returns a value  $v$  and emits a trace  $\theta$ . It also updates a global environment from  $\gamma$  to  $\gamma'$ .

Note that the first family of relation executes *only* the server part of a program and returns a client program, which is then executed by  $\Rightarrow_c$ . This is represented formally by the PROGRAM rule. In order to reduce an ELIOM program  $P$ , we first reduce the server part using  $\Rightarrow_m$ . This returns no value and a client program  $\mu$  which we execute. We now look into each specific feature in greater detail

#### 4.3.1 Generated client programs

Let us first describe evaluation rules for generated client programs. Generated client programs are ML programs with some additional constructions which are described in Figure 4.12. The new evaluation rules are presented in Figure 4.13. The construction **bind env f** binds the current accessible environment to **f** in the global environment  $\gamma$ . This is implemented by the BINDENV rule. **bind r = e with f** computes  $e$  in the environment previous associated to **f**. The results is then stored as **r** in  $\gamma$ . This construction is also usable for module expressions and is implemented by the BIND and



$\text{BIND}_m$  rules. All these constructions also accept paths of references such as  $\mathbf{R.f}$ .

The new **bind** constructs are similar to the ones used in languages with continuations in the catch/throw style. Instead of storing both an environment and the future computation, we store only the environment. This will allow us to implement closures across locations, in particular the case where fragments are used inside a server closure.

The client reduction relation also inherits the ML rules (rule **CLIENTCODE**). In such a case, the global environment is passed around following an order compatible with traces. For example, the **LETIN** rule for let expression would be modified like so:

$$\frac{e' \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\mathbf{c}} v', \theta \quad e \xRightarrow{\rho + \{x \mapsto v'\}|\gamma' \rightarrow \gamma''}_{\mathbf{c}} v, \theta'}{(\text{let } x = e' \text{ in } e) \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\mathbf{c}} v, \theta @ \theta'}$$

Here,  $e'$  is evaluated first (since  $\theta$  is present first in the resulting traces), hence it uses the initial environment  $\gamma$  and returns the environment  $\gamma'$ , which is then passed along.

In the rest of this thesis, we use  $\mathbf{f}$  to denote the reference associated to fragments closures and  $\mathbf{r}$  to denote the reference associated to a specific value of a fragment.

$$\begin{aligned} \mathbf{p} &::= (\mathbf{X}.)^* \mathbf{x} && \text{(Reference path)} \\ D_c &::= D_{\text{ML}} \\ &| \text{bind env } \mathbf{p} && \text{(Env binding)} \\ &| \text{bind } \mathbf{p} = e \text{ with } \mathbf{p}' \mid \text{bind } \mathbf{p} = M \text{ with } \mathbf{p}' && \text{(Global binding)} \end{aligned}$$

Figure 4.12: Grammar of client programs

$$\begin{aligned} \text{BIND} \quad & \frac{e \xRightarrow{\gamma(\mathbf{p_f})|\gamma \rightarrow \gamma'}_{\mathbf{c}} v, \theta \quad S \xRightarrow{\rho|(\gamma' + \{\mathbf{p} \mapsto v\}) \rightarrow \gamma''}_{\mathbf{c}} V, \theta'}{(\text{bind } \mathbf{p} = e \text{ with } \mathbf{p_f}; S) \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\mathbf{c}} V, \theta @ \theta'} & \text{BINDENV} \quad \frac{S \xRightarrow{\rho|(\gamma + \{\mathbf{p_f} \mapsto \rho\}) \rightarrow \gamma'}_{\mathbf{c}} V, \theta}{(\text{bind env } \mathbf{p_f}; S) \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\mathbf{c}} V, \theta} \\ \text{BIND}_m \quad & \frac{M \xRightarrow{\gamma(\mathbf{p_f})|\gamma \rightarrow \gamma'}_{\mathbf{c}} V, \theta \quad S \xRightarrow{\rho|(\gamma' + \{\mathbf{p} \mapsto V\}) \rightarrow \gamma''}_{\mathbf{c}} V', \theta'}{(\text{bind } \mathbf{p} = M \text{ with } \mathbf{p_f}; S) \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\mathbf{c}} V', \theta @ \theta'} & \text{CLIENTCODE} \quad \text{Inherit the rules from ML} \end{aligned}$$

Figure 4.13: Semantics for client generated programs –  $e \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\mathbf{c}} v, \theta$

### 4.3.2 Base, Client and Server declarations

We now consider the case of base, client and server declarations. The rules are presented in Figure 4.16. Let us first describe the execution of complete  $\text{ELIOM}_\epsilon$  programs (rule **PROGRAM**). A program  $P$  reduces to a client value  $v$  if and only if we can first create a *server* reduction of  $P$  that produces no value, emits a client program  $\mu$  and a trace  $\theta_s$ . We can then create a reduction of  $\mu$  that reduces in  $v$  with a trace  $\theta_c$ . The trace of the program is the concatenation of the traces. We see that the execution of  $\text{ELIOM}_\epsilon$

program is split in two as described earlier. Let us now look in more details at various construction of the  $\text{ELIOM}_\varepsilon$  language.

**Base** The base reduction relation corresponds exactly to the ML reduction relation, and always returns empty programs (rule  $\text{BASECODE}$ ). When reducing a base declaration in a mixed context, we both reduce the declaration using  $\Rightarrow_b$ , but also add the declaration to the emitted client program (rule  $\text{BASEDECL}$ ). As we can see, base declarations are executed twice: once on the server and once on the client.

**Client contexts and injections** The goal of the client reduction relation  $\Rightarrow_{c/s}$  is *not* to reduce client programs. It only reduces server code contained by injections inside client code. It returns a client expression without injections, a client program and a trace. Since we don't want to execute client code, it does not inherit the reduction rules for ML. Given an injection  $f\%e$ , the rule  $\text{INJECTION}$  reduces the server side expression  $(f^s e)$  to a value  $v$ . We then transform the server value  $v$  into a client value using the  $\downarrow$  operator presented Figure 4.15. We then returns the client expression  $(f^c \downarrow v)$  without executing it. This expression will be executed on the client side, to deserialize the value. The value injection operator, noted  $\downarrow$  represents the serialization of values from the server to the client and is the identity over constants in  $\text{Const}_b$  and references, and fail on any other values. According to the definition of converters, if  $f$  is a converter  $\tau_s \rightsquigarrow \tau_c$ , then  $f^s$  is the server side function of type  $\tau_s \rightarrow \text{serial}$  and  $v$  should be of type  $\text{serial}$ . Since  $\text{serial}$  is defined on  $b$ , the injection of values should be the identity.

The rule  $\text{CLIENTCONTEXT}$  defines the evaluation of server expression up to client contexts. Client contexts are noted  $E[e_1, \dots, e_n]$  and are defined in Figure 4.14. A client context can have any number of holes which must all contain injections. The rest of the context can contain arbitrary ML syntax that are not injections. Evaluation under a multi-holed context proceed from left to right. The resulting programs and traces are constructed by concatenation of each program and trace.

In order to evaluate client declarations, the rule  $\text{CLIENTDECL}$  uses  $\Rightarrow_{c/s}$  to evaluate the server code present in the declaration  $D_c$  which returns a declaration without injections  $\overline{D}_c$  and a client program  $\mu$ . We then return the client program composed by the concatenation of the two. We demonstrate this in Example 4.5. The  $\text{ELIOM}_\varepsilon$  program is presented on the left side. It first declares the integer  $a$  on the server then inject it on the client and returns the result. The emitted code, shown in the middle, contains an explicit call to the  $\text{int}^c$  deserializer while the rest of the client code is unchanged. The returned value is shown on the right.

$$\begin{array}{l} \text{let}_s a = 3 \\ \text{let}_c \text{return} = \text{int}\%x + 1 \end{array} \quad \Longrightarrow_m \quad \text{let return} = (\text{int}^c 3) + 1 \quad \Longrightarrow_c \quad 4, \langle \rangle$$

Example 4.5: Execution of a client declaration

**Server code and fragments** The server reduction relation reduces server code and emits the appropriate client program associated to client fragments. Since client program are

$$\begin{array}{ll}
E_e ::= [f\%e] \mid e \mid (E_e E_e) \mid \lambda x. E_e \mid \text{let } x = E_e \text{ in } E_e & \downarrow c = c \quad \text{when } c \in \text{Const}_b \\
E_M ::= M \mid (E_M : \mathcal{M}) \mid E_M(E_M) & \downarrow \mathbf{p} = \mathbf{p} \\
& \mid \text{functor}(X_i : \mathcal{M}) E_M \mid \text{struct } (E_D)^* \text{ end} & \downarrow v = \perp \quad \text{otherwise} \\
E_D ::= D \mid \text{let}_c x_i = E_e \mid \text{module}_c X_i = E_M & 
\end{array}$$

Figure 4.15: Injections

Figure 4.14: Execution contexts for injections –  $E[\cdot]$

of values –  $\downarrow v$

mostly ML programs, it inherits the ML reduction rules (rule `SERVERCODE`) where client programs are concatenated in the same order as traces. Client fragments are handled by the rule `FRAGMENT`. Let us consider a fragment  $\{\{ e \}\}$ , this evaluation proceeds in two steps: first, we evaluate all the injections inside the client expression  $e$  using the relation  $\Rightarrow_{c/s}$  described in the previous section. We thus obtain an expression without injection  $\bar{e}$  and a client program  $\mu$ .

The second step is to register  $\bar{e}$  to be evaluated in the client program. One could propose to simply consider client fragments as values. This is however quite problematic, as it could lead to duplicated side effects. Consider the program presented on the left side of Example 4.6. If we were simply to pass fragments along, the `print` statement would be evaluated twice. Instead, we create a fresh identifier  $\mathbf{r}$  that will be globally bound to  $\bar{e}$  in the client program, as shown in rule `FRAGMENT`. This way, the client expression contained inside the fragment will be executed once, in a timely manner. The execution rule for fragment is demonstrated in Example 4.6. As before, the `ELIOMε` program is presented on the left, the emitted client program is shown in the middle and the returned value is on the right. Note that both `frags` and `fragc` are the identity function.

$$\begin{array}{lll}
\text{let}_s x = \{\{ (\text{print } 3) \}\} & & \text{bind env } \mathbf{f} \\
\text{let}_c \text{return} = & \Longrightarrow_m & \text{bind } \mathbf{r} = (\text{print } 3) \text{ with } \mathbf{f} \\
\text{frag}\%x + \text{frag}\%x & & \text{let return} = & \Longrightarrow_c & 6, \langle 3 \rangle \\
& & (\text{frag}^c \mathbf{r}) + (\text{frag}^c \mathbf{r})
\end{array}$$

Example 4.6: Execution of a fragment containing side-effects

**Closures and fragments** In the client program above, we also use a reference  $\mathbf{f}$  and the `bind env` construct. To see why this is necessary, we now consider a case where fragments are used inside closures. This is presented in Example 4.7. The `ELIOMε` program, presented on the left, computes  $1 + 3 + 2$  on the client (although in a fairly obfuscated way). We first define the client variable  $a$  as 1. We then define a server closure  $f$  containing a client fragment capturing  $a$ . We then define a new variable also named  $a$  and call  $(f \ 3)$ , inject the results and returns. When evaluating the definitions of  $f$ , since it contains syntactically a client fragment, we will emit the client instruction `bind env f`, where  $\mathbf{f}$  is a fresh identifier. This will capture the local environment, which is  $\{a \mapsto 1\}$  at this point of the client program. When we execute  $(f \ 3)$ , we will finally reduce the client fragment and emit the `(bind r = (intc 3) + a with f)` instruction. On the client,

this will be executed in the **f** environment, hence  $a$  is 1 and the result is 4. Once this is executed, we move back to the regular environment, where  $a$  is 2, and proceed with the execution.

Thanks to this construction, the capturing behavior of closures is preserved across location boundaries. The **bind env** construct is generated by the SERVERDECL rule. FRAGS( $D_s$ ) returns the fragments syntactically present in  $D_s$ . For each fragment, the local environment is bound to the associated reference.

$$\begin{array}{ll}
\text{let}_c a = 1 & \text{let } a = 1 \\
\text{let}_s f \ x = \{\{ \text{int}^c x + a \}\} & \text{bind env } \mathbf{f} \\
\text{let}_c a = 2 & \text{let } a = 2 \\
\text{let}_s y = (f \ 3) & \text{bind } \mathbf{r} = (\text{int}^c 3) + a \text{ with } \mathbf{f} \\
\text{let}_c \text{return} = \text{frag}^c y + a & \text{let } \text{return} = (\text{frag}^c \mathbf{r}) + a
\end{array}
\quad \Longrightarrow_m \quad \Longrightarrow_c \quad 6, \langle \rangle$$

Example 4.7: Execution of a fragment inside a closure

### Server code inside client contexts

$$\begin{array}{ll}
\text{INJECTION} & \text{CLIENTCONTEXT} \\
\frac{(f^s \ e) \xRightarrow{\rho}_s v, \mu, \theta}{f^c e \xRightarrow{\rho}_{c/s} (f^c \downarrow v), \mu, \theta} & \frac{\forall i, \ e_i \xRightarrow{\rho}_{c/s} v_i, \mu_i, \theta_i}{E[e_1, \dots, e_n] \xRightarrow{\rho}_{c/s} E[v_1, \dots, v_n], \mu_1; \dots; \mu_n, @_i \theta_i}
\end{array}$$

### Server code

$$\begin{array}{l}
\text{FRAGMENT} \\
\frac{e \xRightarrow{\rho}_{c/s} \bar{e}, \mu, \theta \quad \mathbf{r} \text{ fresh}}{\{\{ e \}\} \mathbf{f} \xRightarrow{\rho}_s \mathbf{r}, (\mu; \text{bind } \mathbf{r} = \bar{e} \text{ with } \mathbf{f}), \theta}
\end{array}$$

### Base code

$$\begin{array}{ll}
\text{SERVERCODE} & \text{BASECODE} \\
\text{Inherit the rules} & \xRightarrow{\rho} \equiv \xRightarrow{\rho}_b \\
\text{from ML} &
\end{array}$$

### Declarations

$$\begin{array}{ll}
\text{BASEDECL} & \text{CLIENTDECL} \\
\frac{D_b \xRightarrow{\rho}_b V, \varepsilon, \theta \quad S \xRightarrow{\rho+V}_m V', \mu', \theta'}{D_b; S \xRightarrow{\rho}_m V + V', (D_b; \mu'), \theta @ \theta'} & \frac{D_c \xRightarrow{\rho}_{c/s} \overline{D_c}, \mu, \theta \quad S \xRightarrow{\rho}_m V, \mu', \theta'}{D_c; S \xRightarrow{\rho}_m V, (\mu; \overline{D_c}; \mu'), \theta} \\
\\
\text{SERVERDECL} & \\
\frac{\text{FRAGS}(D_s) = \{\{ e_i \}\}_{\mathbf{f}_i} \quad D_s \xRightarrow{\rho}_s V, \mu, \theta \quad S \xRightarrow{\rho+V}_m V', \mu', \theta'}{D_s; S \xRightarrow{\rho}_m V + V', (\text{bind env } \mathbf{f}_i; \mu; \mu'), \theta @ \theta'} & \\
\\
\text{PROGRAM} & \\
\frac{P \xRightarrow{\rho}_m (), \mu, \theta_s \quad \mu \xRightarrow{\rho | \varepsilon \rightarrow \gamma}_c v, \theta_c}{P \xRightarrow{\rho}_s v, \theta_s @ \theta_c} &
\end{array}$$

Figure 4.16: Semantics for base, client and server sections –  $e \xRightarrow{\rho}_s v, \mu, \theta$

**Fragment annotations** In the previous examples, we presented the server reduction rules where, for each syntactic fragment, a fresh reference  $\mathbf{f}$  is generated and bound to the environment. In the rest of this thesis, we will simply assume that all fragments syntactically present in the program are annotated with a unique reference. Such annotation is purely syntactic and can be done by walking the syntax tree of the program. Annotated fragments are noted  $\{\{ \dots \}\}_{\mathbf{f}}$ .

Mixed structures syntactically present in the program are also annotated in a similar manner with a unique module reference. Annotated mixed structures are noted `struct ... end $\mathbf{F}$` .

### 4.3.3 Mixed modules

Let us now describe the reduction relation for mixed modules. The mixed reduction relation is presented in Figure 4.17 and, just like the server relation, has for goal to evaluate all the server code and emit a client program to be later evaluated by the client relation. Mixed modules can be composed of either mixed functors, functor applications or structures. The mixed relation contains various rules that are similar to the ML reduction rules for modules. The notable novel aspect of mixed functor is that they both have a client part and a server part. This is different from client fragments, which only have a client part that can be manipulated on the server via an identifier. The server part of mixed modules also need to indicate its client part. In order to do this, each mixed structure will contains an additional field called `Dyn` which contains a module identifier. The identifier points to a globally bound module on the client which is the result of the client-side evaluation.

Let us first demonstrate these features in Example 4.8. In this example, we declare a mixed module  $X$  containing a fragment  $x$  and an integer  $y$ . We then declare another mixed module  $Y$  containing a submodule. The structure of the emitted client code mimics closely the structure of the server code. In particular, the `bind` operation is nested inside the mixed module  $X$  that is emitted on the client. The exact same names are reused on the client. We also register each structure in the global environment using the annotated identifier of the structure. Here, we use the `bind` construct as a shorthand for `bind with` that doesn't change the environment. The shape of the program is kept intact thanks to the `MIXEDMODVAR`, `MIXEDQUALMODVAR` and `MIXEDSTRUCT` rules. The first two are similar to the non mixed version, but the last one deserves some explanation. First, it prefixes all the fragment references inside the body of the structure. This is for consistency with functors, as we will see later. It then adds the `Dyn` field to the returned structure, as discussed before. Finally, it emits a `bind` on the client and returns the module reference. Each structure is thus bound appropriately, even when nested.

Module identifiers are not used in the present program, but they are used in the case of mixed functors, as we will see now.

**Mixed functors, injections and client side application** Before exposing the complex interaction of mixed functors and fragment, let us illustrate various details about mixed functors in Example 4.9. The server code proceed in the following way: we first define a

<pre> module<sub>m</sub> X = struct   let<sub>s</sub> x = {{ 1 }}   let<sub>c</sub> y = 2 + frag%<sub>c</sub>x end<sub>X</sub> module<sub>m</sub> Y = struct   module<sub>m</sub> A = X end<sub>Y</sub> let<sub>c</sub> return = Y.A.y </pre>	$\Longrightarrow_m$	<pre> bind X = struct   bind env X.f   bind r = 1 with X.f   let y = 2 + (frag<sup>c</sup> r) end module X = X bind Y = struct   module A = X end module Y = Y let return = Y.A.y </pre>	$\Longrightarrow_c$	$3, \langle \rangle$
---	---------------------	--	---------------------	----------------------

Example 4.8: Execution of mixed modules

server variable  $x$  followed by a mixed functor  $F$  containing an injection. We then define a mixed module  $Y$  and executes *on the client* the functor application  $F(Y)$ .

First, let us recall that injections inside mixed functors can only refer to elements outside the functor. This means that injections inside functors can be reduced as soon as we consider a functor. In particular, we do not wait for functor application. This can be seen in the MODCLOSURE rule which returns a functor closure on the server side and emit the client part of the functor on the client side. We then take the client part of the body of the functor (noted  $M|_c$ ) and applies the  $\Rightarrow_{c/s}$  reduction relation, which executes injections inside client code. In this example, it results in the injection `int%cx` being resolved immediately in the client-side version of the functor.

Mixed functor application can be done in client and server contexts. When it is done in a client context, we simply call the client-side definition and omits the server-side execution completely. Hence we can simply emit the client-code  $F(Y)$ . Execution is done through the usual rules for client sections. This is always valid since each mixed declaration emits a client declaration with the same name and the same shape.

<pre> let<sub>s</sub> x = 1 module<sub>m</sub> F(X : M) = struct   let<sub>c</sub> b = X.a + int%<sub>c</sub>x end<sub>Y</sub> module<sub>m</sub> Y = struct   let<sub>c</sub> a = 2 end<sub>Y</sub> module<sub>c</sub> Z = F(Y) let<sub>c</sub> return = Z.b </pre>	$\Longrightarrow_m$	<pre> module F(X : M) = struct   let b = X.a + (int<sup>c</sup> 1) end<sub>Y</sub> bind Y = struct   let a = 2 end module Y = Y module Z = F(Y) let return = Z.b </pre>	$\Longrightarrow_c$	$3, \langle \rangle$
--	---------------------	---	---------------------	----------------------

Example 4.9: Execution of mixed functors with injections

**Mixed functors and fragments** The difficulty of the reduction of mixed functor containing fragments is that the server-side application of a mixed functor should result in both server and client effects. This makes the reduction rules for mixed functor application quite delicate. We illustrate this with Example 4.10. In this example, we define a functor  $F$  contains only the server declaration  $x$ . The argument of the functor simply contains two integers, one on the server and one on the client. In the fragment bound to  $x$ , we add the two integers (using an escaped value). The interesting aspect here is that the body of the client fragment depends on both the client and the server side of the argument, even if there is no actual client side for the functor  $F$ . The rest of the program is composed of a simple mixed module  $Y$  and the mixed functor application  $F(Y)$ .

The first step of the execution is to define the client side part of  $F$  and  $Y$ , as demonstrated in the previous example. In this case, since  $F$  only contains a server side declaration, the client part of the functor returns an empty structure. We then have to execute  $F(Y)$ . This is done with the STRUCTBETA rule. When reducing a mixed functor application, we first generate a fresh identifier ( $\mathbf{R_Z}$  here) and prefix all the fragment closure identifiers. We then evaluate the body of the functor on the server, which gives us both the server module value and the generated client code. In this case, we simply obtain the binding of  $\mathbf{r_x}$ . Note that this reference is not prefixed by  $\mathbf{R_Z}$  since it is freshly generated at runtime. If the functor was applied again, we would simply generated a new one. In order for functor arguments to be properly available on the client, we need to introduce additional bindings. For this purpose, we lookup the `Dyn` field for each module argument and insert the additional binding. In this case, `module  $X = Y$` . This gives us a complete client structure which we can bind to  $\mathbf{R_Z}$ .

<pre> module<sub>m</sub> F(X : M) = struct   let<sub>s</sub> x = { { X.a + int%<i>X.b</i> } }f<sub>x</sub> end<sub>F</sub> module<sub>m</sub> Y = struct   let<sub>c</sub> a = 4   let<sub>s</sub> b = 2 end<sub>Y</sub>  module<sub>m</sub> Z = F(Y)  let<sub>c</sub> return = frag%<i>Z.x</i> </pre>	$\Rightarrow_m$	<pre> bind env F module F(X : M) = struct end bind Y = struct   let a = 4 end module Y = Y  bind R<sub>Z</sub> = struct   module X = Y   bind env R<sub>Z</sub>.f<sub>x</sub>   bind r<sub>x</sub> = Y.a + (int<sup>c</sup> 2)   with R<sub>Z</sub>.f<sub>x</sub> end with F module Z = F(Y)  let return = (frag<sup>c</sup> r<sub>x</sub>) </pre>	$\Rightarrow_c$ 6, $\langle \rangle$
--	-----------------	--	--------------------------------------

Example 4.10: Execution of mixed functors with fragments

We see here that the body of functors allows to emit client code in a dynamic but controlled way. Generated module references used on the client are remembered on the server using the `Dyn` field while closure identifiers ensure that the proper environment is used. One problematic aspect of this method is that it leads to two executions of the client side. We shall discuss this in Section 5.5.

### Mixed module expressions

$$\begin{array}{c}
\text{MIXEDSTRUCT} \\
\frac{S[\mathbf{f}_i \mapsto \mathbf{X.f}_i]_i \xRightarrow{\rho}_m V, \mu, \theta \quad V' = V + \{\text{Dyn} \mapsto \mathbf{X}\}}{\text{struct } S \text{ end}_{\mathbf{X}} \xRightarrow{\rho}_m V', \mathbf{X}, \text{bind } \mathbf{X} = \text{struct } \mu \text{ end}, \theta} \\
\\
\text{MIXEDMODVAR} \\
\frac{\rho(X) = V}{X \xRightarrow{\rho}_m V, X, \varepsilon, \langle \rangle} \\
\\
\text{APP} \\
\frac{M \xRightarrow{\rho}_m V, M_c, \mu, \theta \quad M' \xRightarrow{\rho}_m V', M'_c, \mu', \theta' \quad V(V') \xRightarrow{\rho}_m V'', \mu'', \theta''}{M(M') \xRightarrow{\rho}_m V'', M_c(M'_c), \mu; \mu'; \mu'', \theta @ \theta' @ \theta''} \\
\\
\text{STRUCTBETA} \\
\frac{\mathbf{R} \text{ fresh} \quad V_f = \text{functor}_m(\rho')(X_i : \mathcal{M}_i)_i \text{struct } S \text{ end}_{\mathbf{F}} \quad V_i(\text{Dyn}) = \mathbf{R}_i \quad S[\mathbf{f}_i \mapsto \mathbf{R.f}_i]_i \xRightarrow{\rho' + \{X_i \mapsto V_i\}_i}_m V, \mu, \theta}{V_f(V_1) \dots (V_n) \xRightarrow{\rho}_m V + \{\text{Dyn} \mapsto \mathbf{R}\}, \left( \begin{array}{c} \text{bind } \mathbf{R} = \text{struct} \\ (\text{module } X_i = \mathbf{R}_i;)_i \\ \mu \\ \text{end with } \mathbf{F} \end{array} \right), \theta} \\
\\
\text{NOTSTRUCTBETA} \\
\frac{V = \text{functor}_m(\rho')(X_i : \mathcal{M}_i)_i M \quad M \xRightarrow{\rho' + \{X_i \mapsto V_i\}_i}_m V_r, \mu, \theta}{V(V_1) \dots (V_n) \xRightarrow{\rho}_m V_r, \mu, \theta,} \quad \text{MIXEDQUALMODVAR} \quad \frac{p \xRightarrow{\rho}_m V, \mu, \theta}{p.X \xRightarrow{\rho}_m V(X), p.X, \mu, \theta} \quad \text{EMPTY} \quad \frac{}{\varepsilon \xRightarrow{\rho}_m \{\}, \varepsilon, \langle \rangle} \\
\\
\text{MODCLOSURE} \\
\frac{M|_c \xRightarrow{\rho}_{c/s} \overline{M}, \mu, \theta}{\text{functor}_m(X : \mathcal{M})M \xRightarrow{\rho}_m \text{functor}_m(\rho)(X : \mathcal{M})M, \text{functor}(X : \mathcal{M}|_c)\overline{M}, \mu, \theta}
\end{array}$$

### Mixed declarations

$$\begin{array}{c}
\text{MIXEDMODDECL} \\
\frac{\text{MixedStructIds}(M) = \overline{\mathbf{F}_i} \quad M \xRightarrow{\rho}_m V, M^c, \mu, \theta \quad S \xRightarrow{\rho + \{X \mapsto V\}}_m V', \mu', \theta'}{\text{module}_m X = M; S \xRightarrow{\rho}_m \{X \mapsto V\} + V', (\overline{\text{bind env } \mathbf{F}_i}; \mu; \text{module } X = M^c; \mu'), \theta @ \theta'}
\end{array}$$

Figure 4.17: Semantics for mixed modules –  $M \xRightarrow{\rho}_s V, \mu, \theta$



## 4.4 Results on locations

The behavior of locations and specialization in the presence of the various language constructs is not obvious, even with various examples. We now present various elementary results related to locations that should make the behavior of some constructions easier to grasp and inform the design of our compilation scheme. Let us note  $\mathcal{M}_{[\ell \mapsto \ell']}$  the substitution on locations in an ELIOM module type  $\mathcal{M}$ .

**Proposition 1.** Given an ELIOM module type  $\mathcal{M}$  and a location  $\ell \in \{b, c, s\}$ , if  $\Gamma \models_{\ell} \mathcal{M}$ , then  $\lfloor \mathcal{M} \rfloor_{\ell} = \mathcal{M}$ .

*Proof.* By definition of  $<:$ ,  $\mathcal{M}$  can only contain declarations on  $\ell$ . This means that, by reflexivity of  $\succ$ , only specialization rules that leave the declaration unchanged are involved.  $\square$

**Proposition 2.** Given an ELIOM module type  $\mathcal{M}$  and a location  $\ell \in \{c, s\}$ , if  $\Gamma \models_b \mathcal{M}$ , then  $\lfloor \mathcal{M} \rfloor_{\ell} = \mathcal{M}_{[b \mapsto \ell]}$ .

*Proof.* We remark that for all  $\ell \in \{c, s\}$ ,  $b \succ \ell$ . Additionally, mixed functors cannot appear on base (since  $m \not\succ b$ ). We can then proceed by induction over the rules for specialization.  $\square$

### 4.4.1 Relation between ML and Eliom

ELIOM is an extension of ML. However, we also want ELIOM to be well integrated in its host language. We now detail the relation between ML and ELIOM. In particular, we show that ML can be completely embedded in ELIOM without changes.

Given an ML module  $m$ , we note  $m_{[\text{ML} \mapsto \ell]}$  the ELIOM module where all the module components have been annotated with location  $\ell$ . Given an ELIOM module  $m$ , we note  $m_{[\mapsto \text{ML}]}$  the ML module where all the location have been erased. We extend these notations to module types and environments.

**Proposition 3.** Given ML type  $\tau$ , expression  $e$ , module  $m$  and module type  $\mathcal{M}$  and locations  $\ell, \ell'$ :

$$\begin{array}{ll}
 \Gamma \models_{\text{ML}} \tau \implies \Gamma_{[\text{ML} \mapsto \ell']} \models_{\ell} \tau & \text{Where } \ell' \succ \ell \\
 \Gamma \triangleright_{\text{ML}} e : \tau \implies \Gamma_{[\text{ML} \mapsto \ell']} \triangleright_{\ell} e : \tau & \text{Where } \ell' \succ \ell \\
 \Gamma \models_{\text{ML}} \mathcal{M} \implies \Gamma_{[\text{ML} \mapsto \ell']} \models_{\ell} \mathcal{M}_{[\text{ML} \mapsto \ell]} & \text{Where } \ell' \succ \ell \\
 \Gamma \blacktriangleright_{\text{ML}} M : \mathcal{M} \implies \Gamma_{[\text{ML} \mapsto \ell']} \blacktriangleright_{\ell} M_{[\text{ML} \mapsto \ell]} : \mathcal{M}_{[\text{ML} \mapsto \ell]} & \text{Where } \ell' \succ \ell
 \end{array}$$

*Proof.* We remark that each syntax, typing rule or well formedness rule for ML has a direct equivalent rule in ELIOM. We can then simply rewrite the proof tree of the hypothesis to use the ELIOM type and well-formedness rules. We consider only some specific cases:

- By Proposition 1 and since the modules are of uniform location, the specialization operation in VAR and MODVAR are the identity.

- The side conditions  $\ell' <: \ell$  are always respected since the modules are of uniform location and by reflexivity of  $<:.$
- The side conditions  $\ell' \succ \ell$  are respected by hypothesis.  $\square$

**Proposition 4.** Given ML type  $\tau$ , expression  $e$ , module  $m$  and module type  $\mathcal{M}$ :

$$\begin{aligned} \Gamma \models_b \tau &\implies \Gamma_{[\mapsto \text{ML}]} \models_{\text{ML}} \tau & \Gamma \models_b \mathcal{M} &\implies \Gamma_{[\mapsto \text{ML}]} \models_{\text{ML}} \mathcal{M}_{[\mapsto \text{ML}]} \\ \Gamma \triangleright_b e : \tau &\implies \Gamma_{[\mapsto \text{ML}]} \triangleright_{\text{ML}} e : \tau & \Gamma \blacktriangleright_b M : \mathcal{M} &\implies \Gamma_{[\mapsto \text{ML}]} \blacktriangleright_{\text{ML}} M_{[\mapsto \text{ML}]} : \mathcal{M}_{[\mapsto \text{ML}]} \end{aligned}$$

*Proof.* We first remark that the following features are forbidden in the base part of the language: injections, fragments, mixed functors and any other location than base. The rest of the language contains no tierless features and coincides with ML. We can then proceed by induction over the proof trees.  $\square$

**Proposition 5.** Given an ML module  $M$  (resp. expression  $e$ ), an execution environment  $\rho$ , a location  $\ell$ , a value  $V$  (resp.  $v$ ) and a trace  $\theta$ :

$$\begin{aligned} M \xRightarrow{\rho} V, \theta &\iff M_{[\text{ML} \mapsto \ell]} \xRightarrow{\rho_{[\text{ML} \mapsto \ell]}}_{\ell} V_{[\text{ML} \mapsto \ell]}, \varepsilon, \theta \\ e \xRightarrow{\rho} v, \theta &\iff e_{[\text{ML} \mapsto \ell]} \xRightarrow{\rho_{[\text{ML} \mapsto \ell]}}_{\ell} v_{[\text{ML} \mapsto \ell]}, \varepsilon, \theta \end{aligned}$$

Furthermore, given an ML program  $P$ , an execution environment  $\rho$ , a value  $v$  and a trace  $\theta$ :

$$\begin{aligned} P \xRightarrow{\rho} v, \theta &\iff P_{[\text{ML} \mapsto \iota]} \xRightarrow{\rho_{[\text{ML} \mapsto \iota]}} v_{[\text{ML} \mapsto \iota]}, \theta & \iota \in \{c, s\} \\ P \xRightarrow{\rho} v, \theta &\iff P_{[\text{ML} \mapsto b]} \xRightarrow{\rho_{[\text{ML} \mapsto b]}} v_{[\text{ML} \mapsto b]}, \theta @ \theta \end{aligned}$$

*Proof.* Let us first note that the ML reduction relation is included in the base, the server and the client-only relations. Additionally, the considered programs, modules or expressions can not contain fragments, injections or binds. The additional rules in the server and client-only relations are only used for these additional syntactic constructs. For the first three statements, we can then proceed by induction. For the last statement, we remark that base code is completely copied to the client during server execution. Using rule PROGRAM, we execute the program twice, which returns the same value but duplicates the trace.  $\square$

**Theorem 2** (Base/ML correspondance). ELIOM modules, expressions and types on base location  $b$  correspond exactly to the ML language.

*Proof.* By Propositions 3 to 5.  $\square$

Thanks to Theorem 2, we can completely identify the language ML and the part of ELIOM on base location. This is of course by design: the base location allows us to reason about the host language, OCAML, inside the new language ELIOM. It also provides the guarantee that anything written in the base location does not contain any communication

between client and server. In the rest of the thesis, we omit location substitutions of the form  $[ML \mapsto b]$  and  $[b \mapsto ML]$ .

Proposition 3 also has practical consequences: Given a file previously typechecked by an ML typechecker, we can directly use the module types either on base, but also on the client or on the server, by simply annotating all the signature components. This give us the possibility, in the implementation, to completely reuse compiled objects from the OCAML typechecker and load them on an arbitrary location. In particular, it guarantees that we can reuse pure OCAML libraries safely and freely.

#### 4.4.2 Notes on soundness

We do not give a complete proof of soundness for  $ELIOM_\varepsilon$ . One of the reason is that a soundness proof for our version of ML is already out of the scope of this thesis (See Section 3.3.3). The  $ELIOM_\varepsilon$  semantics is mostly implemented in terms of the ML one. Here we will assume that the ML semantics is sound, and give some arguments towards the soundness of  $ELIOM_\varepsilon$ .

**Server and Base semantics** Since base code exactly corresponds to ML code, and can not depend on server and client identifiers, the soundness of the base reduction relation is equivalent to the soundness of ML. The server reduction relation only adds the `FRAGMENT` rule. If a fragment is well typed and the various injection can be reduced, it is easy to see that this rule always applies. Since references  $\mathbf{r}$  are opaque values on the server and can only be used in injections, they do not compromise the soundness of the server reduction relation.

**Client emitted code** In our presentation, we do not give typing rules for the `bind env` and the `bind with` constructs. Such typing rules can simply be given by typing  $\mathbf{f}_i$  references similarly to environments: with a module type. With this addition, it should be possible to show that client emitted programs are always well typed. We can then rely on the soundness of ML equipped with `bind`. The difficult point here are fragments, due to the delayed nature of the `bind` constructs. However, the content of a fragment is always executed in an environment with a type compatible with its typing environment, which should ensure correctness.

**Mixed modules** The client-side behavior of a mixed module is fairly easy to model, since it is equivalent to a client module. The server-side behavior is mostly the same as the one of a pair composed by a server-side module and a fragment containing a client-side module. As such, the soundness arguments should be the same that the one used for fragments.



## 5 Compilation of Eliom programs

In Chapter 4, we gave a tour of the  $\text{ELIOM}_\varepsilon$  language from a formal perspective, providing a type system and an interpreted semantics.  $\text{ELIOM}$ , however, is not an interpreted language. The interpreted semantics is here both as a formal tool and to make the semantics of the language more approachable to users, as demonstrated in Chapter 2. In the implementation,  $\text{ELIOM}$  programs are compiled to two programs: one server program (which is linked and executed on the server) and a client program (which is compiled to `JAVASCRIPT` and executed in a browser). The resulting programs are efficient and avoid unnecessary back-and-forth communications between the client and the server.

Description of the complete compilation toolchain, including emission of `JAVASCRIPT` code, is out of scope of this thesis (see [Vouillon and Balat \[2014\]](#)). Instead, we describe the compilation process in term of emission of client and server programs in an ML-like language equipped with additional primitives. Hence, we present the typing and execution of compiled programs, in Section 5.1 and the compilation process, in Section 5.2. However, we also want to ensure that the interpreted semantics, which is explained to users, corresponds to the compiled semantics<sup>1</sup>. This is done in Section 5.4. Finally, we discuss the design of mixed functors from a compilation perspective in Section 5.5.

### 5.1 Target languages $\text{ML}_s$ and $\text{ML}_c$

We introduce the two target languages  $\text{ML}_c$  and  $\text{ML}_s$  as extensions of  $\text{ML}$ . The additions in these two new languages are highlighted in Figure 5.1. Typing is provided in Section 5.1.5. The semantics is provided in Section 5.1.6. As before, we use globally bound identifiers, which we call “references” and note in bold: **r**. References can also be paths, such as **X.r**. In some contexts, we accept a special form of reference path, noted **Dyn.x** which we explain in Section 5.1.4. In practice, these references are implemented with uniquely generated names and associative tables. Contrary to the interpreted semantics, references are also used to transfer values from the server to the client and can appear in expressions. A reference used inside an expression is always of type **serial**.

#### 5.1.1 Converters

For each converter  $f$ , we note  $f^s$  and  $f^c$  the server side encoding function and the client side decoding function. If  $f$  is of type  $\tau_s \rightsquigarrow \tau_c$ , then  $f^s$  is of type  $\tau_s \rightarrow \text{serial}$  and  $f^c$  is of type  $\text{serial} \rightarrow \tau_c$ . We will generally assume that if the converter  $f$  is available in the environment, then  $f^c$  and  $f^s$  are available in the client and server environment respectively.

---

<sup>1</sup>Thus, the main result of this chapter is that you do not need to read it to understand  $\text{ELIOM}_\varepsilon$  programs!

ML <sub>s</sub> grammar	ML <sub>c</sub> grammar
$\mathbf{p} ::= \text{Dyn.f} \mid (\mathbf{X}.)^*\mathbf{f}$ (Reference path)	$\mathbf{p} ::= \text{Dyn.f} \mid (\mathbf{X}.)^*\mathbf{f}$ (Reference path)
$\tau ::= \dots \mid \text{frag}$ (Fragment type)	$e ::= \dots \mid \mathbf{x}$ (Reference)
$v ::= \dots \mid \mathbf{r}$ (Reference)	$M ::= \dots \mid \mathbf{X}$ (Module reference)
$e ::= \dots$	$D ::= \dots$
$\mid \text{fragment } \mathbf{p} \ e^*$ (Fragment call)	$\mid \text{bind } \mathbf{p} = e$ (Fragment closure)
$M ::= \dots$	$\mid \text{bind}_m \mathbf{p} = M$ (Functor fragment)
$\mid p.\text{Dyn}$ (Dynamic field)	$\mid \text{exec } ()$ (Fragment execution)
$\mid \text{fragment}_m \mathbf{p} (M^*)$ (Fragment)	
$D ::= \dots$	
$\mid \text{injection } \mathbf{x} \ e$ (Injection)	
$\mid \text{end } ()$ (End token)	

Figure 5.1: Grammar for ML<sub>s</sub> and ML<sub>c</sub> as extensions of ML<sub>ε</sub>

### 5.1.2 Injections

For injections, we associate server-side the injected value  $e$  to a reference  $\mathbf{v}$  using the construction `injection  $\mathbf{v}$   $e$` , where  $e$  is of type `serial`. When the server execution is over, a mapping from references to injected values is sent to the client.  $\mathbf{v}$  is then used client-side to access the value.

An example is given in Example 5.1. In this example, two integers are sent from the server to the client and add them on the client. We suppose the existence of a base abstract type `int`, a converter `int` of type  $(\text{int} \rightsquigarrow \text{int})$  and the associated encoding and decoding functions. The server program, in Example 5.1a, creates two injections,  $\mathbf{v}_1$  and  $\mathbf{v}_2$  and does not expose any bindings nor return any values. These injections hold the serialized integers 4 and 2. The client program, in Example 5.1b, uses these two injections, deserialize their values, adds them, and returns the result. Note that `injection` is *not* a network operation. It simply stores a mapping between references (*i.e.*, names) and serialized values. The mapping generated at the end of the server execution is shown in Example 5.1c. After the execution of the server code, this mapping is sent to the client, and used to execute the client code.

<code>injection <math>\mathbf{v}_1</math> (<math>\text{int}^s</math> 4);</code>	<code>let return =</code>	$\mathbf{v}_1 \mapsto 4$
<code>injection <math>\mathbf{v}_2</math> (<math>\text{int}^s</math> 2);</code>	<code>(<math>\text{int}^c \mathbf{v}_1</math>) + (<math>\text{int}^c \mathbf{v}_2</math>);</code>	$\mathbf{v}_2 \mapsto 2$
(a) Server program	(b) Client program	(c) Mapping of injections

Example 5.1: Client-server programs calling and using injections

### 5.1.3 Fragments

The primitive related to fragments also relies on shared references between the server program and the client program. However, these references allow to uniquely identify functions that are defined on the client but are called on the server. To implement this, we use the following primitives:

- In  $ML_C$  structures, **bind**  $\mathbf{p} = e$  declares a new client function bound to the reference  $\mathbf{p}$ . The function  $e$  takes an arbitrary amount of argument of type **serial** and returns any type.
- In  $ML_S$  expressions, **fragment**  $\mathbf{p} \ e_1 \dots e_n$  is a delayed function application which registers that, on the client, the function associated to  $\mathbf{p}$  will be applied to the arguments  $e_i$ . All the arguments must be of type **serial**. It returns a value of type **frag**, which holds a unique identifier referring to the result of this application.

Here again, none of these primitives are network communication primitives. While the API is similar to Remote Procedure Calls, the execution is very different: **fragment** only accumulates the function call in a list, to be executed later. When the server execution is over, the list of calls is sent to the client, and used during the client execution. OCAML, and consequently ELIOM, are impure languages: the order of execution of statement is important. In order to control the order of execution, we introduce two additional statements: **end** (), on the server, introduces an **end** marker in the list of calls. **exec** (), on the client, executes all the calls until the next **end** token.

Example 5.2 presents a pair of programs which emit the client trace  $\langle 2; 3; 3 \rangle$ , but in such a way that, while the client does the printing, the values and the execution order are completely determined by the server. The server code (Example 5.2a) calls **f** with 2 as argument, injects the result and then calls **f** with 3 as argument. The client code, in Example 5.2b, declares a fragment closure **f**, which simply adds one to its arguments, and **exec** both fragments. In-between both executions, it prints the content of the injection **v**. During the execution of the server, the list of calls (Example 5.2c) and the mapping of injections (Example 5.2d) are built. First, when **fragment f (int<sup>s</sup> 2)** is executed, a fresh reference **r<sub>1</sub>** is generated, the call to the fragment is added to the list and **r<sub>1</sub>** is returned. The injection adds the association **v<sub>1</sub> ↦ r<sub>1</sub>** to the mapping of injections. The call to **end** () then adds the token **end** to the list of fragments. The second fragment proceeds similarly to the first, with a fresh identifiers **r<sub>2</sub>**. Once server execution is over, the newly generated list of fragments and mapping of injections are sent to the client. During the client execution, the execution of the list is controlled by the **exec** calls. First, **(f 2)** emits  $\langle 2 \rangle$  and is evaluated to 3, and the mapping **r<sub>1</sub> ↦ 3** is added to a global environment. Then **v<sub>1</sub>** is resolved to **r<sub>1</sub>** and printed (which shows  $\langle 3 \rangle$ ). Finally **(f 3)** emits  $\langle 3 \rangle$  and is evaluated to 4.

The important thing to note here is that both the injection mapping and the list of fragments are completely dynamic. We could add complicated control flow to the server program that would drive the client execution according to some dynamic values. The only static elements are the names **f** and **v<sub>1</sub>**, the behavior of **f** and the number of call

to `exec ()`. We cannot, however, make the server-side control flow depend on client-side values, since we obey a strict phase separation between server and client execution.

Finally, remark that we do not need the `bind env` construct introduced in Section 4.3.1. Instead, we directly capture the environment using closures that are extracted in advance. We will see how this extraction works in more details while studying the compilation scheme, in Section 5.2.

<code>let <math>x_1</math> = fragment <math>f</math> (<math>\text{int}^s 2</math>);</code>	<code>bind <math>f</math> = <math>\lambda x.((\text{print } (\text{int}^c x)) + 1)</math>;</code>	<code>{<math>r_1 \mapsto (f 2)</math>}; end;</code>
<code>injection <math>v_1</math> (<math>\text{frag}^s x_1</math>);</code>	<code>exec ();</code>	<code>{<math>r_2 \mapsto (f 3)</math>}; end;</code>
<code>end ();</code>	<code>let <math>a</math> = (<math>\text{print } (\text{frag}^c v_1)</math>);</code>	(c) List of fragments
<code>let <math>x_2</math> = fragment <math>f</math> (<math>\text{int}^s 3</math>);</code>	<code>exec ();</code>	
<code>end ();</code>	<code>let return = <math>a</math></code>	$v_1 \mapsto r_1$
(a) Server program	(b) Client program	(d) Mapping of injections

Example 5.2: Client-server program defining and calling fragments

#### 5.1.4 Modules

We introduce three new module-related construction that are quite similar to fragment primitives:

- `bindm p = M` is equivalent to `bind` for modules. It is a client instruction that associates the module or functor  $M$  to the reference  $p$ .
- `fragmentm p ( $R_1$ ) ... ( $R_n$ )` is analogous to `fragment p e` for modules. It is a delayed functor application that is used on the server to register that the functor associated to  $p$  will have to be applied to the modules associated to  $R_i$ . It returns a fresh reference that represents the resulting module. Contrary to `fragment`, it can only be applied to module references.
- `p.Dyn` returns a reference that represents the client part of a server module  $p$ . This is used for  $\text{ELIOM}_\varepsilon$  mixed structure that have both a server and a client part.

The first argument of `fragment`, `fragmentm`, `bind` and `bindm` can also be a reference path `Dyn.f`, where `Dyn` is the locally bound `Dyn` field inside a module. This allows us to isolate some bound references inside a fresh module reference. This is useful for functors, as we will now demonstrate in Example 5.3.

In this example, we again add integers<sup>2</sup> on the client while controlling the values and the control flow on the server. We want to define server modules that contain server values but also trigger some evaluation on the client, in a similar way to fragments. The first step is to define a module  $X$  on the server and to bind a corresponding module  $\mathbf{X}$  on the client. Similarly to the interpreted semantics presented in Section 4.3, we add a

<sup>2</sup>But better! or at very least, more obfuscated.



**Dyn** field to the server module that points to the client module. Plain structures such as  $X$  are fairly straightforward, as we only need to declare each part statically and add the needed reference. `bindm` allows to declare modules globally.

We then declare the functor  $F$  on the server and bind the functor  $\mathbf{F}$  on the client. The server-side functor contains a call to a fragment defined in the client-side functor. The difficulty here is that we should take care of differentiating between fragment closures produced by different functor applications. For this purpose, we use a similar technique than the one presented in Section 4.3.3, which is to prefix the fragment closure identifier  $\mathbf{f}$  with the reference of the client-side module. This reference is available on the server side as the **Dyn** field and is generated by a call to `fragmentm`. When  $F$  is applied to  $X$  on the server, we generate a fresh reference  $\mathbf{R}$  and add  $\{\mathbf{R}_1 \mapsto \mathbf{F} \mathbf{X}_0\}$  to the execution queue. When `exec ()` is called, we introduce the additional binding  $\{\mathbf{Dyn} \mapsto \mathbf{R}_1\}$  in the environment and apply  $\mathbf{F}$  to  $\mathbf{X}_0$ , which will register the  $\mathbf{R}_1.\mathbf{f}$  fragment closure. Since it is the result of this specific functor application, the closure  $\mathbf{R}_1.\mathbf{f}$  will always add 4 to its argument. The rest of the execution proceed as shown in the previous section: we call a new fragment, which triggers the client-side addition  $2 + 4$  and use an injection to pass the results around.

<pre> module X = struct   module Dyn = <math>\mathbf{X}_0</math>;   let a = 2 end; module F(Y : <math>\mathcal{M}_s</math>) = struct   module Dyn = fragment<sub>m</sub> <math>\mathbf{F}</math> (Y.Dyn);   let b = fragment Dyn.f (int<sup>s</sup> Y.a); end; module Z = F(X); end (); injection <math>\mathbf{v}_1</math> (frag<sup>s</sup> Z.b); </pre>	<pre> bind<sub>m</sub> <math>\mathbf{X}_0</math> = struct   let c = 4 end; bind<sub>m</sub> <math>\mathbf{F}(Y : \mathcal{M}_c) = struct   bind Dyn.f =     <math>\lambda a. ((\text{int}^c a) + Y.c)</math>; end; exec () let return = (frag<sup>c</sup> <math>\mathbf{v}_1</math>); </math></pre>	<pre> {<math>\mathbf{R}_1 \mapsto \mathbf{F} \mathbf{X}_0</math>}; {<math>\mathbf{r}_2 \mapsto \mathbf{R}_1.\mathbf{f} 2</math>}; end </pre>
(a) Server program	(b) Client program	(c) List of fragments  $\mathbf{v}_1 \mapsto \mathbf{r}_2$ (d) Mapping of injections

Example 5.3: Client-server program using module fragments

### 5.1.5 Type system rules

The  $\text{ML}_s$  and  $\text{ML}_c$  typing rules are presented in Figures 5.2 and 5.3 as a small extension over the ML typing rules presented in Section 3.2. Note that the typing rules for the new primitives are weakly typed and are certainly not sound with respect to serialization and deserialization. Given arbitrary  $\text{ML}_s$  and  $\text{ML}_c$  programs, there is no guarantee that (de)serialization will not fail at runtime. This is on purpose. Indeed, all these guarantees are provided by ELIOM itself.  $\text{ML}_s$  and  $\text{ML}_c$  are target languages that are very liberal by design, so that all patterns permitted by ELIOM are expressible with them. Furthermore, from an implementation perspective,  $\text{ML}_s$  and  $\text{ML}_c$  are simply OCAML libraries and do not rely on further compiler support. Note that **Dyn** fields are not reflected in signatures.

The fragment  $\text{FRAGMENT}_m$  rule does not enforce that the **Dyn** field is present in all the arguments. This is enforced by construction during compilation.

### 5.1.6 Semantics rules

We define two reduction relations as extensions of the ML reduction rules (see Section 3.3). The  $\Rightarrow_{\text{ML}_s}$  reduction for  $\text{ML}_s$  server programs is presented in Figure 5.4. The  $\Rightarrow_{\text{ML}_c}$  reduction for  $\text{ML}_c$  client programs is presented in Figure 5.5. Let us consider a server structure  $S_s$  and a client structure  $S_c$ . A paired execution of the two structures is presented below:

$$S_s \xRightarrow{\rho_s}_{\text{ML}_s} V_s, \xi, \zeta, \theta_s \qquad S_c, \xi \xRightarrow{\rho_c \mid \gamma \cup \zeta \rightarrow \gamma'}_{\text{ML}_c} V_c, \xi', \theta_c$$

Let us now detail these executions rules. As with the ML reduction,  $\rho_s$  and  $\rho_c$  are the local environments of values while  $\theta_s$  and  $\theta_c$  are the traces for server and client executions respectively.  $V_s$  and  $V_c$  are the returned values. Similarly to the interpreted semantics for  $\text{ELIOM}_\varepsilon$ , the client reduction uses global environments noted  $\gamma$ .

As we saw in the previous examples, server executions emits two sets of information during execution: a queue of fragments and a map of injections. Mapping of injection is a traditional (global) environment where bindings are noted  $\{\mathbf{v} \mapsto \dots\}$ . The queue of fragments is noted  $\xi$  and contains end tokens **end** and fragment calls  $\{\mathbf{r} \mapsto \mathbf{f} \ v_1 \dots v_n\}$ . Concatenation of fragment queues is noted  $++$ . We now see the various rules in more details.

**Injections** Injection bindings are collected on the server through the **INJECTION** rule. When creating a new injection binding, we inject the server-side value using the injection of value operator, noted  $\downarrow v$  and presented in Figure 4.15. This models the serialization

$$\begin{array}{c} \text{FRAGMENT} \\ \frac{\forall i, \Gamma \triangleright_{\text{ML}_s} e_i : \text{serial}}{\Gamma \triangleright_{\text{ML}_s} \text{fragment } \mathbf{p} \ e_1 \ \dots \ e_n : \text{frag}} \qquad \text{INJECTION} \qquad \frac{\Gamma \triangleright_{\text{ML}_s} e : \text{serial}}{\Gamma \blacktriangleright_{\text{ML}_s} \text{injection } \mathbf{v} \ e : \varepsilon} \qquad \text{END} \qquad \frac{}{\Gamma \blacktriangleright_{\text{ML}_s} \text{end } () : \varepsilon} \\[10pt] \text{FRAGMENT}_m \\ \frac{\forall i, \Gamma \triangleright_{\text{ML}_s} p_i : \mathcal{M}_i}{\Gamma \blacktriangleright_{\text{ML}_s} \text{module Dyn} = \text{fragment}_m \ \mathbf{p} \ p_1.\text{Dyn} \ \dots \ p_n.\text{Dyn} : \varepsilon} \end{array}$$

Figure 5.2: Typing rules for  $\text{ML}_s$

$$\begin{array}{c} \text{BIND} \qquad \frac{}{\Gamma \triangleright_{\text{ML}_c} e : \tau} \qquad \text{BIND}_m \qquad \frac{}{\Gamma \blacktriangleright_{\text{ML}_c} M : \mathcal{M}} \qquad \text{REFERENCE} \qquad \frac{}{\Gamma \triangleright_{\text{ML}_c} \mathbf{x} : \text{serial}} \qquad \text{EXEC} \qquad \frac{}{\Gamma \blacktriangleright_{\text{ML}_c} \text{exec } () : \varepsilon} \\[10pt] \Gamma \blacktriangleright_{\text{ML}_c} \text{bind } \mathbf{p} = e : \varepsilon \quad \Gamma \blacktriangleright_{\text{ML}_c} \text{bind}_m \ \mathbf{p} = M : \varepsilon \quad \Gamma \triangleright_{\text{ML}_c} \mathbf{x} : \text{serial} \quad \Gamma \blacktriangleright_{\text{ML}_c} \text{exec } () : \varepsilon \end{array}$$

Figure 5.3: Typing rules for  $\text{ML}_c$

of values before transmission from server to client by ensuring that only base values and references are injected. Other kinds of values should be handled using converters explicitly.

The injection environment  $\zeta$  forms a valid client-side global environment. When executing the client-side program, we simply assume that  $\zeta$  is included in the initial global environment  $\gamma$ .

**Fragments and functors** On the server, fragments and functors calls are added to the queue through the **FRAGMENT** and **FRAGMENT<sub>m</sub>** rules. In both rules, the reference of the associated closure or functor is provided, along with a list of arguments. A fresh reference symbolizing the fragment is generated and the call is added to the queue  $\xi$ . Note that in the case of regular fragments, the arguments are expressions which can themselves contains fragment calls. The module rule **FRAGMENT<sub>m</sub>** is similar, the main difference being that it only accept module references as arguments of the call.

Fragment closures and functors are bound on the client through the **BIND** and **BIND<sub>m</sub>** rules, which simply binds a reference to a value or a module value in the global environment  $\gamma$ . Since **bind** accepts references of the form **Dyn.f**, it must first resolves **Dyn** to the actual reference. This is done through the **DYN** rule.

**Segmented execution** In  $ML_s$  and  $ML_c$  programs, the execution of fragments is segmented through the use of the **exec ()/end ()** instructions. On the server, **end** instructions are handled through the **END** rule, which simply adds an **end** token to the execution queue  $\xi$ . On the client, we use the **EXEC** rule, associated to the **FRAG** and **FRAG<sub>m</sub>** rules. When **exec** is called, the **EXEC** rule triggers the execution of the segment of the queue until the next **end** token. Each token  $\{\mathbf{r} \mapsto \mathbf{f} \ v_1 \dots v_n\}$  is executed with the **FRAG** rule as the function call  $(\mathbf{f} \ v_1 \dots v_n)$ . The result of this function call is bound to **r** in the global environment  $\gamma$ . Similarly, functor calls are executed using the **FRAG<sub>m</sub>** rule. Note that for functors we also introduce the **Dyn** field in the local environment, which allows local **bind** definitions. Once all the tokens have been executed in the considered fragment queue, we resume the usual execution.

$$\begin{array}{c}
\text{INJECTION} \\
\frac{e \xRightarrow{\rho}_{\text{ML}_s} v, \xi, \zeta, \theta \quad S \xRightarrow{\rho}_{\text{ML}_s} V, \xi', \zeta', \theta'}{\text{injection } \mathbf{x} \ e; S \xRightarrow{\rho}_{\text{ML}_s} V, \xi ++ \xi', (\zeta \cup \{\mathbf{x} \mapsto \downarrow v\} \cup \zeta'), \theta @ \theta'} \\
\\
\text{FRAGMENT} \\
\frac{\mathbf{p} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{p}' \quad \forall i, e_i \xRightarrow{\rho}_{\text{ML}_s} v_i, \xi_i, \zeta_i, \theta_i \quad \mathbf{r} \text{ fresh}}{\text{fragment } \mathbf{p} \ e_1 \ \dots \ e_n \xRightarrow{\rho}_{\text{ML}_s} \mathbf{r}, (\xi_1 ++ \dots \xi_n ++ \{\mathbf{r} \mapsto \mathbf{p}' \downarrow v_1 \dots \downarrow v_n\}), \cup_i \zeta_i, @_i \ \theta_i} \\
\\
\text{FRAGMENT}_m \\
\frac{\mathbf{p} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{p}' \quad \forall i, p_i.\text{Dyn} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{R}_i, \varepsilon, \varepsilon, \langle \rangle \quad \mathbf{R} \text{ fresh}}{\text{fragment}_m \ \mathbf{p} \ p_1.\text{Dyn} \ \dots \ p_n.\text{Dyn} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{R}, \{\mathbf{R} \mapsto \mathbf{p}' \ \mathbf{R}_1 \dots \mathbf{R}_n\}, \{\}, \langle \rangle} \\
\\
\text{END} \qquad \text{DYN} \qquad \text{SERVERCODE} \\
\frac{S \xRightarrow{\rho}_{\text{ML}_s} V, \xi, \zeta, \theta}{\text{end } (); S \xRightarrow{\rho}_{\text{ML}_s} V, \text{end} ++ \xi, \zeta, \theta} \quad \frac{\rho(\text{Dyn}) = \mathbf{R}}{\text{Dyn.f} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{R.r}} \quad \begin{array}{l} \text{Inherit the rules} \\ \text{from ML} \end{array}
\end{array}$$

Figure 5.4: Semantics rules for  $\text{ML}_s - S \xRightarrow{\rho}_{\text{ML}_s} V, \xi, \zeta, \theta$

$$\begin{array}{c}
\text{BIND} \\
\frac{\mathbf{p} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{p}' \quad e, \xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} v, \xi', \theta \quad S, \xi' \xRightarrow{\rho|\gamma' + \{\mathbf{p}' \mapsto v\} \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta'}{(\text{bind } \mathbf{p} = e; S), \xi \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta @ \theta'} \\
\\
\text{BIND}_m \\
\frac{\mathbf{p} \xRightarrow{\rho}_{\text{ML}_s} \mathbf{p}' \quad M, \xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} V_{\mathbf{p}}, \xi', \theta \quad S, \xi' \xRightarrow{\rho|\gamma' + \{\mathbf{p}' \mapsto V_{\mathbf{p}}\} \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta'}{(\text{bind}_m \ \mathbf{p} = M; S), \xi \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta @ \theta'} \\
\\
\text{REFERENCE} \qquad \text{EXEC} \\
\frac{\gamma(\mathbf{r}) = v}{\mathbf{r}, \xi \xRightarrow{\rho|\gamma \rightarrow \gamma}_{\text{ML}_c} v, \xi, \langle \rangle} \quad \frac{\xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} \{\}, \theta \quad S, \xi' \xRightarrow{\rho|\gamma' \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta'}{(\text{exec } (); S), \xi ++ \text{end} ++ \xi' \xRightarrow{\rho|\gamma \rightarrow \gamma''}_{\text{ML}_c} V, \xi'', \theta @ \theta'} \\
\\
\text{FRAG} \qquad \text{DYN} \\
\frac{(\mathbf{f} \ v_1 \dots v_n) \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} v, \theta \quad \xi \xRightarrow{\rho|\gamma' + \{\mathbf{r} \mapsto v\} \rightarrow \gamma''}_{\text{ML}_c} \{\}, \theta'}{\{\mathbf{r} \mapsto \mathbf{f} \ v_1 \dots v_n\} ++ \xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} \{\}, \theta @ \theta'} \quad \frac{\rho(\text{Dyn}) = \mathbf{R}}{\text{Dyn.f} \xRightarrow{\rho}_{\text{ML}_c} \mathbf{R.r}} \\
\\
\text{FRAG}_m \qquad \text{CLIENTCODE} \\
\frac{\mathbf{F}(\mathbf{R}_1) \dots (\mathbf{R}_n) \xRightarrow{\rho \cup \{\text{Dyn} \mapsto \mathbf{R}\}|\gamma \rightarrow \gamma'}_{\text{ML}_c} V, \theta \quad \xi \xRightarrow{\rho|\gamma' + \{\mathbf{R} \mapsto V\} \rightarrow \gamma''}_{\text{ML}_c} \{\}, \theta'}{\{\mathbf{R} \mapsto \mathbf{F} \ \mathbf{R}_1 \dots \mathbf{R}_n\} ++ \xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} \{\}, \theta @ \theta'} \quad \begin{array}{l} \text{Inherit the} \\ \text{ML rules} \end{array}
\end{array}$$

Figure 5.5: Semantics rules for  $\text{ML}_c - S, \xi \xRightarrow{\rho|\gamma \rightarrow \gamma'}_{\text{ML}_c} V, \xi', \theta$

## 5.2 Compilation

In Section 5.1, we presented the two target languages  $ML_s$  and  $ML_c$ . We now present the compilation process transforming *one*  $ELIOM_\epsilon$  program into two distinct  $ML_s$  and  $ML_c$  programs. Before giving a more formal description in Section 5.2.2, we present the compilation process through three examples of increasing complexity.

**Injections and fragments** Example 5.4 presents an  $ELIOM_\epsilon$  program containing only simple declarations involving fragments and injections without modules. The  $ELIOM_\epsilon$  program is presented on the left, while the compiled  $ML_s$  and  $ML_c$  programs are presented on the right. In this example, a first fragment is created. It only contains an integer and is bound to  $a$ . A second fragment that uses  $a$  is created and bound on the server to  $b$ . Finally,  $b$  is used on the client via an injection. The program returns 4.

For each fragment, we emit a **bind** declaration on the client. The client expression contained in the fragment is abstracted and transformed in a closure that is bound to a fresh reference. The number of arguments of the closure corresponds to the number of injections inside the fragment. Similarly to the interpreted semantics, we use the client part of the converter on the client. In this case,  $\{\{ 1 \}\}$  is turned into  $\lambda().1$  and  $\{\{ \text{frag}\%a + 1 \}\}$  is turned into  $\lambda v.((\text{frag}^c v) + 1)$ . On the server, each fragment is replaced by a call to the primitive **fragment**. The arguments of the call are the identifier of the closure and all the injections that are contained in the fragment. The **fragment** primitive, which was presented in Section 5.1.3, registers that the closure declared on the client should be executed later on. Since all the arguments of **fragment** should be of type **serial**, we apply the client and server parts of the converters at the appropriate places. The **exec** and **end** primitives synchronize the execution so that the order of side effects is preserved. When **exec** is encountered, it executes queued fragment up to an **end** token which was pushed by an **end** primitive. We place an **exec/end** pair at each server section. This is enough to ensure that client code inside server fragment and client code in regular client declaration is executed in the expected order.

Note that injections, which occur outside of fragments, and escaped values, which occur inside fragments, are compiled in a very different way. Injections have the useful property that the use site and number of injections is completely static: we can collect all the injections on the server, independently of the client control flow and send them to the client. This is the property that allows us to avoid communications from the client to the server.

$ELIOM_\epsilon$	$ML_s$	$ML_c$
$\text{let}_s a = \{\{ 1 \}\};$	$\text{let } a = \text{fragment } f_0 ();$ $\text{end } ();$	$\text{bind } f_0 = \lambda().1;$ $\text{exec } ();$
$\text{let}_s b = \{\{ \text{frag}\%a + 1 \}\};$	$\text{let } b = \text{fragment } f_1 (\text{frag}^s a);$ $\text{end } ();$	$\text{bind } f_1 = \lambda v.((\text{frag}^c v) + 1);$ $\text{exec } ();$
$\text{let}_c \text{return} = \text{frag}\%b + 2;$	$\text{injection } v (\text{frag}^s b);$	$\text{let } \text{return} = (\text{frag}^c v) + 2;$

Example 5.4: Compilation of expressions

**Sections and modules** We now present an example with client and server modules in Example 5.5. The lines of the various programs have been laid out in a way that should highlight their correspondence.

We declare a server module  $X$  containing a client fragment, a client functor  $F$  containing an injection, a client functor application  $Z$  and finally the client return value, with another injection. The compilation of the server module  $X$  proceeds in the following way: on the server, we emit a module  $X$  similar to the original declaration but where fragments have been replaced by a call to the **fragment** primitive. On the client, we only emit the call to **bind**, without any of the server-side code structure. Compilation for the rest of the code proceeds in a similar manner.

This compilation scheme is correct thanks to the following insight: In client and server modules or functors, the special instructions for fragments and injection can be freely lifted to the outer scope. Indeed, the fragment closure bound in  $\mathbf{f_0}$  can only reference client elements. Since the server  $X$  can only introduce server-side variables, the body of the fragment closure is independent from the server-side code. A similar remark can be made about the client functor  $F$ : the functor argument must be on the client, hence it cannot introduce new server binding. The server identifier that is injected must have been introduced outside of functor and the injection can be lifted outside the functor.

Using this remark, the structure of the  $\text{ML}_s$  and  $\text{ML}_c$  programs is fairly straightforward: Code on the appropriate side has the same shape as in the original  $\text{ELIOM}_\varepsilon$  program and code on the other side only contains calls to the appropriate primitives.

$\text{ELIOM}_\varepsilon$	$\text{ML}_s$	$\text{ML}_c$
$\text{module}_s X = \text{struct}$ $\quad \text{let}_s a = \{\{ 2 \}\}$ $\quad \text{let}_s b = 4$ $\text{end};$  $\text{module}_c F(Y : \mathcal{M}) = \text{struct}$ $\quad \text{let}_c a = Y.b + \text{int}\%X.b$ $\text{end};$ $\text{module}_c Z =$ $\quad F(\text{struct let}_c b = 2 \text{ end});$ $\text{let}_c \text{return} =$ $\quad \text{frag}\%X.a + Z.a;$	$\text{module } X = \text{struct}$ $\quad \text{let } a = \text{fragment } \mathbf{f_0} ()$ $\quad \text{let } b = 4$ $\text{end};$ $\text{end } ();$  $\text{injection } \mathbf{v_0} (\text{int}^s X.b);$   $\text{injection } \mathbf{v_1} (\text{frag}^s X.a);$	$\text{bind } \mathbf{f_0} = \lambda().2;$   $\text{exec } ();$ $\text{module } F(Y : \mathcal{M}) = \text{struct}$ $\quad \text{let } a = Y.b + (\text{int}^c \mathbf{v_0})$ $\text{end};$ $\text{module } Z =$ $\quad F(\text{struct let } b = 2 \text{ end});$ $\text{let return} =$ $\quad (\text{frag}^c \mathbf{v_1}) + Z.a;$

Example 5.5: Compilation of client and server modules and functors

**Mixed modules** Finally, Example 5.6 presents the compilation of mixed modules. In this example, we create a mixed structure  $X$  containing a server declaration and a client declaration with an injection. We define a functor  $F$  that takes a module containing a client integer and use it both inside a client fragment, and inside a client declaration. We then apply  $F$  to  $X$  and use an injection to compute the final result of the program.

The compilation of the mixed module  $X$  is similar to the procedure for programs: we compile each declaration and use the **injection** primitive as needed. Additionally, we

add a **Dyn** field on the server-side version of the module. The content of the **Dyn** field is determined statically for simple structures (here, it is  $\mathbf{X}_0$ ). The client-side version of the module is first bound to  $\mathbf{X}_0$  using the  $\text{bind}_m$  primitive. We then declare  $X$  as a simple alias. This alias ensures that  $X$  is also usable in client sections transparently.

For functors, the process is similar. One additional complexity is that the **Dyn** field should be dynamically generated. For this purpose, we add a call to the  $\text{fragment}_m$  primitive. Each call to  $\text{fragment}_m$  generates a new, fresh identifier. We also prefix each call to **fragment** by the **Dyn** field. On the client, we emit two different functors. The first one is called  $F$  and contains only the client declarations to be used inside the rest of the client code. It is used for client-side usage of mixed functors. An example with the interpreted semantics was presented in Section 4.3.3. The other one is bound to a new reference (here  $\mathbf{F}_1$ ) and contains both client declaration, along with calls to the **bind** and **exec** primitives. This function is used to perform client side effects: when the server version of  $F$  is applied, a call to  $\mathbf{F}_1$  is registered and will be executed when the client reaches the associated **exec** call (here, the last one).

$\text{ELIOM}_\varepsilon$	$\text{ML}_s$	$\text{ML}_c$
$\text{module}_m X = \text{struct}$ $\quad \text{let}_s a = 2$ $\quad \text{let}_c b = 4 + \text{int}\%a$ $\text{end};$	$\text{module } X = \text{struct}$ $\quad \text{module Dyn} = \mathbf{X}_0;$ $\quad \text{let } a = 2; \text{end } ();$ $\quad \text{injection } \mathbf{v}_0 (\text{int}^s a);$ $\text{end};$	$\text{bind}_m \mathbf{X}_0 = \text{struct}$ $\quad \text{exec } ();$ $\quad \text{let } b = 4 + (\text{int}^c \mathbf{v}_0)$ $\text{end};$ $\text{module } X = \mathbf{X}_0;$
$\text{module}_m F(Y : \mathcal{M}) = \text{struct}$ $\quad \text{let}_s c = \{\{ Y.b \}\}$ $\quad \text{let}_c d = 2 * Y.b$ $\text{end};$	$\text{module } F(Y : \mathcal{M}) = \text{struct}$ $\quad \text{module Dyn} =$ $\quad \quad \text{fragment}_m \mathbf{F}_1 (Y.\text{Dyn});$ $\quad \text{let } c = \text{fragment Dyn.f}_0 ();$ $\quad \text{end } ();$ $\text{end};$	$\text{bind}_m \mathbf{F}_1(Y : \mathcal{M}) = \text{struct}$ $\quad \text{bind}_m \text{Dyn.f}_0 = \lambda().(Y.b);$ $\quad \text{exec } ();$ $\quad \text{let } d = 2 * Y.b$ $\text{end};$ $\text{module } F(Y : \mathcal{M}) = \text{struct}$ $\quad \text{let } d = 2 * Y.b$ $\text{end};$
$\text{module}_m Z = F(X);$	$\text{module } Z = F(X);$ $\text{end } ();$	$\text{module } Z = F(X);$ $\text{exec } ()$
$\text{let}_c \text{return} =$ $\quad \text{frag}\%Z.c + Z.d;$	$\text{injection } \mathbf{v}_1 (\text{frag}^s Z.c);$	$\text{let return} =$ $\quad (\text{frag}^c \mathbf{v}_1) + Z.d;$

Example 5.6: Compilation of a mixed functor

### 5.2.1 Sliceability

In order to simplify our presentation of mixed functors, both in the slicing rules and in the simulation proofs, we consider a sliceability constraint which dictates which programs can be sliced.

A program is said sliceable if mixed structures are only defined at top level, or directly inside a toplevel mixed functor. We demonstrate this constraint in Example 5.7. The

program presented on the left is not sliceable, since it contains a structure which is nested inside a structure in a functor. The semantically equivalent program on the right is sliceable, since structures are not nested.

We discuss this restriction and explain how to partially relax this restriction in Section 5.5.

<pre> module<sub>m</sub> F(X : <math>\mathcal{M}</math>) = struct   module<sub>m</sub> Y = struct     ...   end end </pre>	<pre> module<sub>m</sub> Y'(X : <math>\mathcal{M}</math>) = struct   ... end module<sub>m</sub> F(X : <math>\mathcal{M}</math>) = struct   module<sub>m</sub> Y = Y'(X) end </pre>
(a) An unsliceable functor	(b) A sliceable functor

Example 5.7: The sliceability constraint

### 5.2.2 Slicing rules

Given an  $\text{ELIOM}_\varepsilon$  module  $M$  (resp. module type  $\mathcal{M}$ , structure  $S$ , ...) and a location  $\iota$  that is either client or server, we note  $\langle M \rangle_\iota$  the result of the compilation of  $M$  to the location  $\iota$ . The result of  $\langle M \rangle_s$  is a module of  $\text{ML}_s$  and the result of  $\langle M \rangle_c$  is a module in  $\text{ML}_c$ .

Let us defines a few notation. As before, we use  $e[a \mapsto b]$  to denote the substitution of  $a$  by  $b$  in  $e$ .  $e[a_i \mapsto b_i]_i$  denotes the repeated substitution of  $a_i$  by  $b_i$  in  $e$ . We note  $\text{FRAGS}(e)$  (resp.  $\text{INJS}(e)$ ) the fragments (resp. injections) present in the expression  $e$ . We note  $(e_i)_i$  the list of elements  $e_i$ . For ease of presentation, we use  $D_\varsigma$  (resp.  $\mathcal{D}_\varsigma$ ) for definitions (resp. declarations) located on location  $\varsigma$ .

In the rest of this section, we assume that all the considered programs are well typed and sliceable.

We now describe how to slice the various constructions of our language. The slicing rules for modules and expressions are defined in Figure 5.6. The slicing rules for structures and declarations are presented in Figure 5.7.

Base structure and signature components are left untouched. Indeed, according to Proposition 4, base elements are valid  $\text{ML}_\varepsilon$  elements. We do not need to modify them in any way. Signature components that are not valid on the target location are simply omitted. Signature components that are valid on the target have their type expressions translated. The translation of a type expression to the client is the identity: indeed, there are no new  $\text{ELIOM}_\varepsilon$  type constructs that are valid on the client. Server types, on the other hand, can contains pieces of client types inside fragments  $\{\tau_c\}$  and inside converters  $\tau_s \rightsquigarrow \tau_c$ . Fragments in  $\text{ML}_s$  are represented by a primitive type, **fragment**, without parameters. The type of converters is represented by the type of their server part, which is  $\tau_s \rightarrow \text{serial}$ . Module and module type expressions are traversed recursively.



### Type expressions

$$\langle \{\tau\} \rangle_s = \mathbf{frag} \qquad \langle \tau_s \rightsquigarrow \tau_c \rangle_s = \langle \tau_s \rangle_s \rightarrow \mathbf{serial}$$

### Signature components

$$\begin{aligned} \langle \mathcal{D}_b; \mathcal{S} \rangle_\iota &= \mathcal{D}_b; \langle \mathcal{S} \rangle_\iota & \langle \mathcal{D}_\varsigma; \mathcal{S} \rangle_\iota &= \langle \mathcal{S} \rangle_\iota & \text{when } \varsigma \not\succ \iota \\ & & &= \langle \mathcal{D}_\varsigma \rangle_\iota; \langle \mathcal{S} \rangle_\iota & \text{when } \varsigma \succ \iota \end{aligned}$$

### Declarations and Definitions

$$\begin{aligned} \langle \mathbf{type}_\iota t_i \rangle_\iota &= \mathbf{type} \ t_i & \langle \mathbf{type}_\iota t_i = \tau \rangle_\iota &= \mathbf{type} \ t_i = \langle \tau \rangle_\iota \\ \langle \mathbf{val}_\iota x_i : \tau \rangle_\iota &= \mathbf{val} \ x_i : \langle \tau \rangle_\iota & \langle \mathbf{module}_\varsigma X_i : \mathcal{M} \rangle_\iota &= \mathbf{module} \ X_i : \langle \mathcal{M} \rangle_\iota \\ \langle \mathbf{let}_\iota x_i = e \rangle_\iota &= \mathbf{let} \ x_i = e & \langle \mathbf{module}_\iota X_i = M \rangle_\iota &= \mathbf{module} \ X_i = \langle M \rangle_\iota \end{aligned}$$

### Module Expressions

$$\begin{aligned} \langle \mathbf{struct} \ S \ \mathbf{end} \rangle_\iota &= \mathbf{struct} \ \langle S \rangle_\iota \ \mathbf{end} \\ \langle M(M') \rangle_\iota &= \langle M \rangle_\iota (\langle M' \rangle_\iota) \\ \langle \mathbf{functor}(X_i : \mathcal{M}) M \rangle_\iota &= \mathbf{functor}(X_i : \langle \mathcal{M} \rangle_\iota) \langle M \rangle_\iota \\ \langle \mathbf{functor}_m(X_i : \mathcal{M}) M \rangle_\iota &= \mathbf{functor}(X_i : \langle \mathcal{M} \rangle_\iota) \langle M \rangle_\iota \end{aligned}$$

### Module Type Expressions

$$\begin{aligned} \langle \mathbf{sig} \ S \ \mathbf{end} \rangle_\iota &= \mathbf{sig} \ \langle S \rangle_\iota \ \mathbf{end} \\ \langle \mathbf{functor}(X_i : \mathcal{M}) \mathcal{M}' \rangle_\iota &= \mathbf{functor}(X_i : \langle \mathcal{M} \rangle_\iota) \langle \mathcal{M}' \rangle_\iota \\ \langle \mathbf{functor}_m(X_i : \mathcal{M}) \mathcal{M}' \rangle_\iota &= \mathbf{functor}(X_i : \langle \mathcal{M} \rangle_\iota) \langle \mathcal{M}' \rangle_\iota \end{aligned}$$

Figure 5.6: Slicing –  $\langle \cdot \rangle_\iota$

Functors and functor applications have each part sliced. Mixed functors are turned into normal functors.

Slicing of structure components inserts additional primitives that were described in Section 5.1. In client structure components, we need to handle injections. We associate each injection to a new fresh reference noted **v**. In  $\text{ML}_s$ , we use the **injection** primitive to register the fact that the given server value should be associated to a given reference. In  $\text{ML}_c$ , we replace each injection by its associated reference. This substitution is applied both inside expressions and structures. Note that for each injection  $f\%x$ , we use the encoding part  $f^s$  and decoding part  $f^c$  for the server and client code, respectively. For server structure components, we apply a similar process to handle fragments. For each fragment, we introduce a reference noted **f**. In  $\text{ML}_s$ , we replace each fragment by a call to **fragment** with argument the associated reference and each escaped value inside the fragment (with the encoding part of the converters). We also add, after the translated component, a call to **end** which indicates that the execution of the component is finished. In  $\text{ML}_c$ , we use the **bind** primitives to associate to each reference a closure where all the escaped values are abstracted. We also introduce the decoding part of each converter for escaped values. We then call **exec**, which executes all the pending fragments until the

$$\langle D_b; S \rangle_\iota = D_b; \langle S \rangle_\iota$$

$$\langle D_m; S \rangle_\iota = \langle D_m \rangle_\iota; \langle S \rangle_\iota$$

$$\langle D_c; S \rangle_s = (\text{injection } \mathbf{x}_i \ (f_i^s \ x_i);)_i \langle S \rangle_s$$

$$\langle D_c; S \rangle_c = \frac{D_c \left[ \overline{f_i \% x_i} \mapsto \overline{(f_i^c \ \mathbf{x}_i)} \right]}{\langle S \rangle_c};$$

$$\text{Where} \begin{cases} \overline{f_i \% x_i} = \text{INJS}(D_c) \\ \overline{\mathbf{x}_i} \text{ is a list of fresh variables.} \end{cases}$$

$$\langle D_s; S \rangle_s = \frac{\langle D_s \rangle_s \left[ \{ \{ e_i \} \}_{\mathbf{f}_i} \mapsto \text{fragment } \mathbf{f}_i \ \overline{(f_{i,j}^s \ a_{i,j})} \right]_i}{\text{end } (); \langle S \rangle_s};$$

$$\langle D_s; S \rangle_c = \frac{(\text{bind Dyn. } \mathbf{f}_i = \lambda \overline{x_{i,j}}. e_i [f_{i,j} \% a_{i,j} \mapsto (f_{i,j}^c \ x_{i,j})];)_i}{\text{exec } (); \langle S \rangle_c};$$

$$\text{Where} \begin{cases} \{ \{ e_i \} \}_{\mathbf{f}_i} = \text{FRAGS}(D_s) \\ \forall i, \overline{f_{i,j} \% a_{i,j}} = \text{INJS}(e_i) \\ \forall i, \overline{x_{i,j}} \text{ are fresh variables} \end{cases}$$

$$\left\langle \frac{\text{module}_m \ \overline{F_i(X_{i_k} : \mathcal{M}_k)} = \text{struct}}{S} \right\rangle_s = \frac{\text{module } \overline{F_i(X_{i_k} : \langle \mathcal{M}_k \rangle_s)} = \text{struct}}{\text{end}} \frac{\text{module Dyn = fragment}_m \ \mathbf{F} \ (\overline{X_{i_k} \text{.Dyn}})}{\langle S \rangle_s};$$

$$\left\langle \frac{\text{module}_m \ \overline{F_i(X_{i_k} : \mathcal{M}_k)} = \text{struct}}{S} \right\rangle_c = \frac{\text{bind}_m \ \mathbf{F}(\overline{X_{i_k} : \langle \mathcal{M}_k \rangle_c}) = \text{struct} \ \langle S \rangle_c \ \text{end};}{\text{module } \overline{F_i(X_{i_k} : \langle \mathcal{M}_k \rangle_c)} = \text{struct} \ S|_c \ \text{end};}$$

$$\langle \text{module}_m \ X_i = \text{struct } S \ \text{end}_{\mathbf{X}} \rangle_s = \frac{\text{module } X_i = \text{struct}}{\text{end}_{\mathbf{F}}} \frac{\text{module Dyn = } \mathbf{X};}{\langle S \rangle_s};$$

$$\langle \text{module}_m \ X_i = \text{struct } S \ \text{end}_{\mathbf{X}} \rangle_c = \frac{\text{bind}_m \ \mathbf{X} = \text{struct} \ \langle S \rangle_c \ \text{end};}{\text{module } X_i = \mathbf{X};}$$

$$\langle \text{module } X_i = M \rangle_s = \frac{\text{module } X_i = \langle M \rangle_s;}{\text{end } ()};$$

$$\langle \text{module } X_i = M \rangle_c = \frac{\text{module } X_i = \langle M \rangle_c;}{\text{exec } ()};$$

Figure 5.7: Slicing of declarations –  $\langle D \rangle_\iota$

next **end** (). This allows to synchronize interleaved side effects between fragments and client components.

Given the constraint of sliceability, a mixed module is either a multi-argument functor returning a structure, or it does not contain any structure at all. For each structure, we use the reference annotated on structures, as described in Section 4.3.2. Mixed modules without structures can simply be sliced by leaving the module expression unchanged. Mixed module types are also straightforward to slice. Mixed structures (with an arbitrary number of arguments) need special care. In  $ML_s$ , we add a **Dyn** field to the structure. The value of this field is the result of a call to the primitive **fragment**<sub>*m*</sub> with arguments **F** and all the **Dyn** fields of the arguments of the functor. In  $ML_c$ , we create two structures for each mixed structure. One is simply a client functor where all the server parts have been removed. Note here that we don't use the slicing operation. The resulting structure does not contain any call to **bind** and **exec**. We also create another structure that uses the regular slicing operation. This structure is associated to **F** with the **bind**<sub>*m*</sub> primitive

### 5.3 Typing preservation

One desirable property is that the introduction of new elements in the language and the compilation operation does not compromise the guarantees provided by the host language. To ensure this, we show that slicing a well typed  $ELIOM_\varepsilon$  program provides two well typed  $ML_c$  and  $ML_s$  programs.

We only consider typing environments  $\Gamma$  containing the primitive types **frag** and **serial** *i.e.*,  $(\text{type}_s \text{ frag}) \in \Gamma$  and  $(\text{type}_b \text{ serial}) \in \Gamma$ . We also extend the slicing operation to typing environments. Slicing a typing environment is equivalent to slicing a signature with additional rules for converters. Converters, in  $ELIOM_\varepsilon$ , are not completely first class: they are only usable in injections and not manipulable in the expression language. As such, they must be directly provided by the environment. We add the two following slicing rules that ensures that converters are properly present in the sliced environment:

$$\begin{aligned} \langle \text{val } f : \tau_s \rightsquigarrow \tau_c \rangle_s &= \text{val } f^s : \langle \tau_s \rangle_s \rightsquigarrow \text{serial} \\ \langle \text{val } f : \tau_s \rightsquigarrow \tau_c \rangle_c &= \text{val } f^c : \text{serial} \rightsquigarrow \tau_c \end{aligned}$$

**Theorem 3** (Compilation preserves typing). Let us consider  $M$  and  $\mathcal{M}$  such that  $\Gamma \blacktriangleright_m M : \mathcal{M}$ . Then  $\langle \Gamma \rangle_s \blacktriangleright \langle M \rangle_s : \langle \mathcal{M} \rangle_s$  and  $\langle \Gamma \rangle_c \blacktriangleright \langle M \rangle_c : \langle \mathcal{M} \rangle_c$

*Proof.* We proceed by induction over the proof tree of  $\Gamma \blacktriangleright_m M : \mathcal{M}$ . The only difficult cases are client and server structure components and mixed structures. For brevity, we only detail the case of client structure components with one injection.

Let us consider  $D_c$  such that  $\Gamma \blacktriangleright_m (D_c; S) : \mathcal{S}$  and  $\text{INJS}(D_c) = f\%x$ . We note **x** the fresh reference. By definition of the typing relation on  $ELIOM_\varepsilon$ , there exists  $\Gamma'$  and  $\tau_c$ ,  $\tau_s$  such that  $\Gamma \subset \Gamma'$ ,  $\Gamma' \triangleright_s f : \tau_s \rightsquigarrow \tau_c$  and  $\Gamma' \triangleright_s x : \tau_s$ . We observe that there cannot be any server bindings in  $D_c$ , Hence we can assume  $\Gamma' = \Gamma$ . This is illustrated on the proof tree

below.

$$\frac{\frac{\Gamma \triangleright_s f : \tau_s \rightsquigarrow \tau_c \quad \overline{\Gamma \triangleright_s x : \tau_s}}{\Gamma \triangleright_c f \% x : \tau_c} \quad \vdots}{\Gamma \blacktriangleright_m(D_c; S) : \mathcal{S}}$$

By definition of slicing on typing environments,  $(\text{val } f^s : \langle \tau_s \rangle_s \rightarrow \text{serial}) \in \langle \Gamma \rangle_s$  and  $(\text{val } f^c : \text{serial} \rightarrow \tau_c) \in \langle \Gamma \rangle_c$ . By definition of  $\text{ML}_c$  and  $\text{ML}_s$  typing rules, we have  $\langle \Gamma \rangle_s \triangleright_{\text{ML}_s}(f^s x) : \text{serial}$  and  $\langle \Gamma \rangle_c \triangleright_{\text{ML}_c}(f^c \mathbf{x}) : \tau_c$ .

We easily have that  $\langle \Gamma \rangle_s \blacktriangleright_{\text{ML}_s} \text{injection } \mathbf{x} (f^s x) : \varepsilon$ , as seen on the proof tree below.

$$\frac{\frac{(\text{val } f^s : \langle \tau_s \rangle_s \rightarrow \text{serial}) \in \langle \Gamma \rangle_s}{\langle \Gamma \rangle_s \triangleright_{\text{ML}_s} f^c : \langle \tau_s \rangle_s \rightsquigarrow \text{serial}} \quad \overline{\langle \Gamma \rangle_s \triangleright_{\text{ML}_s} x : \langle \tau_s \rangle_s}}{\langle \Gamma \rangle_s \blacktriangleright_{\text{ML}_s} \text{injection } \mathbf{x} (f^s x) : \varepsilon}$$

By induction hypothesis on  $\Gamma$ ,  $(\text{val}_c x_j : \tau_c) \triangleright_m d_c[f \% x \mapsto x_j] : \varepsilon$  where  $v_j$  is fresh, we have

$\langle \Gamma \rangle_c, (\text{val } x_j : \tau_c) \triangleright_{\text{ML}_c} d_c[f \% x \mapsto x_j] : \varepsilon$ . We can then replace the proof tree of  $v_j$  by the one of  $(f^c \mathbf{x})$ . We simply need to ensure that the environments coincide. This is the case since  $f^c$  cannot be introduced by new bindings. We can then remove the binding of  $v_j$  from the environment, since it is unused. We obtain that  $\langle \Gamma \rangle_c \blacktriangleright_{\text{ML}_c} d_c[f \% x \mapsto (f^c \mathbf{x})] : \varepsilon$  which allows us to conclude.

$$\frac{\frac{(\text{val } f^c : \text{serial} \rightarrow \tau_c) \in \langle \Gamma \rangle_c}{\langle \Gamma \rangle_c \triangleright_{\text{ML}_c} f^c : \text{serial} \rightsquigarrow \tau_c} \quad \overline{\langle \Gamma \rangle_c \triangleright_{\text{ML}_c} \mathbf{x} : \text{serial}}}{\langle \Gamma \rangle_c \triangleright_{\text{ML}_c} (f^c \mathbf{x}) : \tau_c} \quad \vdots \quad \frac{}{\langle \Gamma \rangle_c \blacktriangleright_{\text{ML}_c} D_c[f \% x \mapsto (f^c \mathbf{x})] ; S : \mathcal{S}}$$

□

## 5.4 Semantics preservation

We now prove that the compilation process preserves the semantics of  $\text{ELIOM}_\varepsilon$  programs. In order to do that, we show that, given an  $\text{ELIOM}_\varepsilon$  program  $P$ , the trace of its execution is the same as the concatenation of the traces of  $\langle P \rangle_s$  and  $\langle P \rangle_c$ .

First, let us put some constraints on the constants of the  $\text{ELIOM}_\varepsilon$ ,  $\text{ML}_s$  and  $\text{ML}_c$  language:

**Hypothesis 1** (Well-behaved converters). Converters are said to be well-behaved if for each constant  $c$  in  $\text{Const}$  such that  $\text{TypeOf}(c) = \tau_s \rightsquigarrow \tau_c$ , then  $c^s \in \text{Const}_s$  and  $c^c \in \text{Const}_c$ .

We now assume that converters in  $\text{ELIOM}_\varepsilon$ ,  $\text{ML}_s$  and  $\text{ML}_c$  are *well-behaved*. We can then state the following theorem.

**Theorem 4** (Compilation preserves semantics). Given sets of constants where converters are well-behaved, given an  $\text{ELIOM}_\varepsilon$  program  $P$  respecting the slicability hypothesis and such that  $P \xRightarrow{\{\}} v, \theta$  then

$$\langle P \rangle_s \xRightarrow{\{\}}_{\text{ML}_s} (), \xi, \zeta, \theta_s \quad \langle P \rangle_c, \xi \xRightarrow{\{\} | \zeta \rightarrow \gamma}_{\text{ML}_c} v, \xi', \theta_c \quad \theta = \theta_s @ \theta_c$$

### 5.4.1 Hoisting

In Section 5.2, we mentioned that a useful property of injections and fragments is that they can be partially lifted outside sections. This property can be used to simplify the simulation proofs. We consider the code transformation that hoists the content of injections out of fragments, client declarations and mixed functors in a way that preserve semantics. This transformation can be decomposed in two parts.

**Injections** We decompose injections inside fragments and client declarations into simpler components. For example, the  $\text{ELIOM}_\varepsilon$  piece of code presented in Example 5.8a is decomposed in Example 5.8b by moving out the application of the converter and leaving only a call to the **serial** converter. All injections using a converter than is not **serial** nor **frag** can be decomposed in such a way.

Since injections can only be used on variables or constants and that no server bindings can be introduced inside a fragment, scoping is preserved. Furthermore, by definition of converters and their client and server components, this transformation preserves typing. It also preserves the dynamic semantics as long as the order of hoisting correspond to the order of evaluation. This can be seen by inspecting the reduction relation for server code under client contexts  $\Rightarrow_{c/s}$ . Finally, it trivially preserves the semantics of the compiled program since it corresponds exactly to how converters are decomposed during compilation.

This allows us to assume that reduction of server code in client context only uses variable lookup and never leads to any evaluation. In particular, this will avoid having to deal with the case of fragments being executed inside the reduction of another fragment (to see why this could happen, consider the case of a converter of type  $\forall \alpha_c. (\text{unit} \rightarrow \{\alpha_c\}) \rightsquigarrow \alpha_c$ ).

```

let a = 1 + 2 in
{{ 3 + int%a }}

```

(a) Fragment with injections

```

let a = 1 + 2 in
let a' = (ints a) in
{{ 3 + (intc serial%a') }}

```

(b) Fragment with hoisted injections

Example 5.8: Hoisting on fragments

In the rest of this section, we assume that reductions of the  $\text{ELIOM}_\varepsilon$  rule **FRAGMENT** are always of the following shape:

$$\frac{e \xrightarrow{\rho_c}_{c/s} \bar{e}, \varepsilon, \langle \rangle}{\{\{ e \} \} \xrightarrow{\rho_s}_s \mathbf{r}, (\mathbf{bind} \ \mathbf{r} = \bar{e}), \langle \rangle} \text{ where } \bar{e} = e[f_i \% x_i \mapsto \rho_s(x_i)]_i \text{ and } f_i \in \{\mathbf{serial}, \mathbf{frag}\}$$

and that reductions of the  $\text{ELIOM}_\varepsilon$  rule **CLIENTDECL** are always of the following shape:

$$\frac{}{D_c \xrightarrow{\rho_s}_{c/s} \overline{D_c}, \varepsilon, \langle \rangle} \text{ where } \overline{D_c} = D_c[f_i \% x_i \mapsto \rho_s(x_i)]_i \text{ and } f_i \in \{\mathbf{serial}, \mathbf{frag}\}$$

**Injections inside mixed modules** We also hoist injections completely out of mixed contexts to the outer englobing scope. For example in the functor presented in Example 5.9a, we can lift the injection out of the functor, as show in Example 5.9b. This is valid since injections can only reference content outside of the functor, by typing. Semantics is similarly preserved since injections inside functors are reduced immediately when encountering a functor, as per rule **MODCLOSURE** in Figure 4.16.

This allows us to assume that the reduction of a mixed module will never lead to the reduction of an injection.

```

lets x = ...
modulem F(X : M) = struct
  letc y = f%x
end

```

(a) Mixed functor with injections

```

lets x = ...
letc y' = f%x
modulem F(X : M) = struct
  letc y = y'
end

```

(b) Mixed functor with hoisted injections

Example 5.9: Hoisting on mixed modules

## 5.4.2 Preliminaries

Let us start with some naming conventions. Identifiers with a hat, such as  $\hat{\gamma}$ , are related to the compiled semantics. For example, while the server environment for the interpreted

semantics is noted  $\rho_s$ , the environment for the execution of the target language  $ML_s$  is noted  $\widehat{\rho}_s$ . This naming convention is only for ease of reading and does not apply a formal relation between the objects with and without hats, unless indicated explicitly.

### Remarks about global environments

Let us make some preliminary remarks about global environments in the  $ELIOM_\varepsilon$  client generated programs and in  $ML_C$ .

Given a global environment  $\gamma$  resulting of  $\Rightarrow_c$ , it contains only two kinds of references:

- Closure fragments, noted **f**, which come from the execution of **bind env**. The associated value is always a environment (*i.e.*, a signature).
- Fragment values, noted **r**, which come from the execution of **bind with**.

In the rest of this section, we consider that we can always decompose global environments  $\gamma$  in two parts: a fragment value environment  $\gamma_r$  containing all the references **r** that were produced by **bind with** and a fragment closure environment  $\gamma_f$  containing only binding of the form  $\{f \mapsto \rho\}$  that were produced by **bind env**.

Similarly, given a global environment  $\widehat{\gamma}$  used in  $ML_C$ . There are only three kind of references:

- Closure fragments, noted **f**, which come from the slicing of syntactic fragments in the source program. The associated value is always a closure.
- Fragment values, noted **r** come from the execution of fragments in the fragment queue.
- Injections, noted **x**. The associated values must be serializable, and hence can only be references or constants in  $Const_b$ .

In the rest of this section, we consider that we can always decompose global compiled environment  $\widehat{\gamma}$  into a fragment closure environment  $\widehat{\gamma}_f$ , a fragment value environment  $\widehat{\gamma}_r$  and an injection environment  $\zeta$ .

### Client equivalence

**Definition 1** (Client values equivalence). Given  $v$  an  $ELIOM$  client value,  $v'$  an  $ML_C$  value and  $\zeta$  an environment of references,  $v$  and  $v'$  are equivalent under  $\zeta$ , noted  $v \simeq_\zeta^c v'$ , if and only if they are equals after substitution by  $\zeta$ :  $v[\zeta] = v'[\zeta]$ .

We extend this notation to environments and traces.

**Definition 2** (Global environment equivalence). We say that an  $ELIOM_\varepsilon$  global environment  $\gamma = \gamma_f \cup \gamma_r$  and an  $ML_C$  global environment  $\widehat{\gamma} = \widehat{\gamma}_f \cup \widehat{\gamma}_r \cup \zeta$  are synchronized if and only if the following conditions hold.

- The reference environments are equivalent:  $\gamma_r \simeq_\zeta^c \widehat{\gamma}_r$

- The domains of  $\gamma_f$  and  $\hat{\gamma}_f$  coincides, and:
  - For each  $\mathbf{f}$  in these environments such that  $\gamma_f(\mathbf{f}) = \rho$  and that  $\hat{\gamma}_f(\mathbf{f}) = \lambda x_0 \dots x_n. \hat{\rho}.e$ , then the following property must hold.

We must have that  $\rho \simeq_\zeta^c \hat{\rho}$  and that for all  $v_0, \dots, v_n, \hat{v}_0, \dots, \hat{v}_n$  such that for all  $i$ ,  $v_i \simeq_\zeta^c \hat{v}_i$ ; then:

$$e[x_i \mapsto v_i]_i \xrightarrow{\rho \mid \gamma \rightarrow \gamma}_c v, \theta \implies (\lambda \bar{x}_i. \hat{\rho}.e \hat{v}_0 \dots \hat{v}_n) \xrightarrow{|\hat{\gamma} \rightarrow \hat{\gamma}}_{\text{ML}_c} \hat{v}, \hat{\theta}$$

with  $v \simeq_\zeta^c \hat{v}$  and  $\theta \simeq_\zeta^c \hat{\theta}$

- For each  $\mathbf{F}$  in these environments such that  $\gamma_f(\mathbf{F}) = \rho$  and that  $\hat{\gamma}_f(\mathbf{F}) = \text{functor}(\hat{\rho})(Y_i : \mathcal{M}_i)_i M$ , we have  $\rho \simeq_\zeta^c \hat{\rho}$ .

**Definition 3** (Fragment closure environment). We consider that  $\hat{\gamma}_f$  is a fragment closure environment for the  $\text{ELIOM}_\varepsilon$  server expression  $e_s$ , noted  $FCE(\hat{\gamma}_f, e_s)$ , if for each  $\{\mathbf{f} \mapsto \lambda \bar{x}_i. \hat{\rho}.e'\}$  in  $\hat{\gamma}_f$ , for each  $\{\{e\}\}_{\mathbf{f}}$  in  $\text{FRAGS}(e_s)$  we have  $e' = e[f_i \% x_i \mapsto x_i]_i$ .

**Definition 4** (Functor closure environment). We consider that  $\hat{\gamma}_f$  is a functor closure environment for the  $\text{ELIOM}_\varepsilon$  module expression  $M$ , noted  $FCE(\hat{\gamma}_f, M)$ , if for each  $\{\mathbf{F} \mapsto \text{functor}(\hat{\rho})(Y_i : \mathcal{M}_i)_i \hat{S}\}$  in  $\hat{\gamma}_f$ , for each  $(\text{struct } S \text{ end}_{\mathbf{F}})$  in  $M_s$  we have  $\hat{S} = \langle S \rangle_c$ . Additionally, we require that  $\hat{\gamma}_f$  be a fragment closure environment for each expression contained in  $S$ .

In the rest of this section, we use the same notation for both properties. We extend this notation to server declarations, server values (by looking under closures) and server environments.

**Lemma 1** (Reduction up to equivalence). Given  $\rho, \hat{\rho}, \gamma = \gamma_f \cup \gamma_r, \hat{\gamma} = \hat{\gamma}_f \cup \hat{\gamma}_r \cup \zeta, e$  and  $\hat{e}$  such that:

$$\rho \simeq_\zeta^c \hat{\rho} \qquad \gamma \simeq_\zeta^c \hat{\gamma} \qquad e[\zeta] = \hat{e}[\zeta] \qquad e \xrightarrow{\rho \mid \gamma \rightarrow \gamma}_c v, \theta$$

Then we have:

$$\hat{e} \xrightarrow{\hat{\rho} \mid \hat{\gamma} \rightarrow \hat{\gamma}}_{\text{ML}_c} \hat{v}, \hat{\theta} \qquad v \simeq_\zeta^c \hat{v} \qquad \theta \simeq_\zeta^c \hat{\theta}$$

*Proof.* The only difference between  $\text{ELIOM}_\varepsilon$  client expressions and  $\text{ML}_c$  expressions are the presence of extra references for injections in  $\text{ML}_c$ . Indeed, syntactic injections have been removed either by the server execution or by compilation and **bind** constructs are only accessible at the module level. Since we assume that the original expression  $e$  and the compiled expression  $\hat{e}$  are the same up to the injection environment  $\zeta$ , we can trivially mimic the execution of  $e$  in  $\hat{e}$  by induction.  $\square$



## Server equivalence

**Definition 5** (Server value equivalence). Given  $v$  an ELIOM server value,  $\hat{v}$  an  $\text{ML}_s$  value. We say they are equivalent, noted  $v \simeq^s \hat{v}$  if and only if

$$v[\{\{ e_i \}\}_{\mathbf{f}} \mapsto \mathbf{fragment} \mathbf{f} \overline{x_{i,j}}]_i = \hat{v} \quad \text{where } \overline{x_{i,j}} = \text{INJS}(e_i)$$

We extend this notation to environments and traces.

### 5.4.3 Server expressions and structures

We first look at server expressions and structures. By definition of the server reduction relation for  $\text{ELIOM}_\varepsilon$ , the emitted program is a series of binds.

**Lemma 2** (Server expressions are simulable). We consider an  $\text{ELIOM}_\varepsilon$  server expression  $e$ ; the  $\text{ELIOM}_\varepsilon$  environments  $\rho_s, \rho_c$  and  $\gamma = \gamma_f \cup \gamma_r$ ; the target environment  $\hat{\rho}_s, \hat{\rho}_c$  and  $\hat{\gamma} = \hat{\gamma}_f \cup \hat{\gamma}_r \cup \zeta$ .

If the expression  $e$  has valid server and client executions:

$$e \xrightarrow{\rho_s}_s v, \mu, \theta_s \quad \mu \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma'}_c \{\}, \theta_c$$

and the following invariants hold:

$$\hat{\rho}_c \simeq_\zeta^c \rho_c \quad \hat{\rho}_s \simeq^s \rho_s \quad \hat{\gamma} \simeq^c \gamma \quad FCE(\hat{\gamma}_f, e) \quad FCE(\hat{\gamma}_f, \rho_s)$$

Then  $\hat{e} = e[\{\{ e_i \}\}_{\mathbf{f}} \mapsto \mathbf{fragment} \mathbf{f} \overline{x_{i,j}}]_i$  has an equivalent execution.

$$\frac{}{\hat{e} \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \hat{v}, \xi_\bullet, \{\}, \hat{\theta}_s} \quad \frac{}{\text{exec } (), \xi_\bullet \text{ ++ end} \xrightarrow{\hat{\rho}_c \mid \hat{\gamma} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \varepsilon, [\ ], \hat{\theta}_c}$$

with the following invariants:

$$\hat{\gamma}' \simeq^c \gamma' \quad FCE(\hat{\gamma}'_f, v) \quad \hat{v} \simeq^s v \quad \hat{\theta}_s \simeq^s \theta_s \quad \hat{\theta}_c \simeq_\zeta^c \theta_c$$

*Proof.* We consider an expression  $e$ ; the  $\text{ELIOM}_\varepsilon$  environments  $\rho_s, \rho_c$  and  $\gamma = \gamma_f \cup \gamma_r$ ; the target environment  $\hat{\rho}_s, \hat{\rho}_c$  and  $\hat{\gamma} = \hat{\gamma}_f \cup \hat{\gamma}_r \cup \zeta$ . such that

$$\hat{\gamma} \simeq^c \gamma \quad \hat{\rho}_c \simeq_\zeta^c \rho_c \quad \hat{\rho}_s \simeq^s \rho_s \quad FCE(\hat{\gamma}_f, e) \quad FCE(\hat{\gamma}_f, \rho_s)$$

We will proceed by induction over the executions of  $e$  and  $\mu$ . The only case of interest is when the server expression is a fragment.

- *Case*  $\{\{ e \}\}_{\mathbf{f}}$ .

We assume that the following executions hold:

$$\frac{\rho_s(x_i) = v_i}{\{\{ e \}\}_{\mathbf{f}} \xrightarrow{\rho_s}_m \mathbf{r}, \text{bind } \mathbf{r} = \bar{e} \text{ with } \mathbf{f}, \langle \rangle} \quad \frac{\gamma(\mathbf{f}) = \rho \quad \bar{e} \xrightarrow{\rho \mid \gamma \rightarrow \gamma'}_c v_c, \theta_c}{(\text{bind } \mathbf{r} = \bar{e} \text{ with } \mathbf{f}) \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma'}_c \{\}, \theta_c}$$

where  $\bar{e} = e[f_i \% x_i \mapsto \downarrow v_i]_i$  and  $\gamma' = \gamma \cup \{\mathbf{r} \mapsto v_c\}$ . We have  $\hat{e}$  equal to **fragment**  $\mathbf{f} \ x_1 \dots x_n$ .

We first consider the execution of  $\hat{e}$ . We can easily construct the following execution.

$$\frac{\hat{\rho}_s(x_i) = \hat{v}_i}{\mathbf{fragment} \ \mathbf{f} \ x_1 \dots x_n \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \mathbf{r}, \{\mathbf{r} \mapsto \mathbf{f} \ \downarrow \hat{v}_1 \dots \downarrow \hat{v}_n\}, \{\}, \langle \rangle}$$

By hypothesis, for each  $i$ ,  $v_i \simeq_{\hat{\gamma}}^s \hat{v}_i$ . This gives us that  $\downarrow v_i \simeq_{\hat{\gamma}}^c \downarrow \hat{v}_i$ . We trivially have that  $\mathbf{r} \simeq_{\hat{\gamma}}^s \mathbf{r}$

Let us now look at the client execution. By client execution of  $\mu$ ,  $\gamma(\mathbf{f}) = \rho$ . Since  $\gamma \simeq_{\zeta}^c \hat{\gamma}$ , we have  $\{\mathbf{f} \mapsto \lambda x_0 \dots x_n. \hat{\rho}. e'\} \in \hat{\gamma}$  and  $\rho \simeq_{\hat{\gamma}}^c \hat{\rho}$ . Furthermore, since  $FCE(\hat{\gamma}_f, \{\{e\}\}\mathbf{f})$ , we know that that  $e' = e[f_i \% x_i \mapsto x_i]_i$ . We have by hypothesis that  $\bar{e} \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma'}_c v_c, \theta_c$ . Since  $\bar{e} = e'[x_i \mapsto \downarrow v_i]$  and since for all  $i$ ,  $\downarrow v_i \simeq_{\hat{\gamma}}^c \downarrow \hat{v}_i$ , we can use Lemma 1 to build the following reduction:

$$\frac{\frac{\frac{e' \xrightarrow{\hat{\rho}_c \cup \{x_i \mapsto \hat{v}_i\}_i \mid \hat{\gamma} \rightarrow \hat{\gamma}}_{\text{ML}_c} \hat{v}, \hat{\theta}_c}{\lambda x_1 \dots x_n. \hat{\rho}. e' \ \downarrow \hat{v}_1 \dots \downarrow \hat{v}_n \xrightarrow{\hat{\rho}_c \mid \hat{\gamma} \rightarrow \hat{\gamma}}_{\text{ML}_c} \hat{v}, \hat{\theta}_c}}{\mathbf{f} \ \downarrow \hat{v}_1 \dots \downarrow \hat{v}_n \xrightarrow{\hat{\rho}_c \mid \hat{\gamma} \rightarrow \hat{\gamma}}_{\text{ML}_c} \hat{v}, \hat{\theta}_c}}{\mathbf{exec} \ (\cdot), \{\mathbf{r} \mapsto \mathbf{f} \ \downarrow \hat{v}_1 \dots \downarrow \hat{v}_n\} \xrightarrow{\hat{\rho}_c \mid \hat{\gamma} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \varepsilon, [\cdot], \hat{\theta}_c}$$

Where  $\hat{\gamma}' = \hat{\gamma} \cup \{\mathbf{r} \mapsto \hat{v}\}$ . By Lemma 1, we have that  $v \simeq_{\hat{\gamma}}^c \hat{v}$  and  $\theta_c \simeq_{\hat{\gamma}}^c \hat{\theta}_c$ . The only part that is changed in  $\gamma'$  and  $\hat{\gamma}'$  is the fragment reference environment, hence we easily have that  $\hat{\gamma}' \simeq^c \gamma'$ , which concludes.  $\blacksquare$

• *Other cases.*

In other cases, we first note that references manipulated inside server code can only fragment references  $\mathbf{r}$ . By hypothesis, the same references are considered before and after compilation. Since the fragment closure environment hypothesis ranges over all server expressions, including the one in closures, it is easy to preserve it during execution. The rest is a very simple induction.  $\blacksquare$

□

**Corollary 1** (Server module declarations are simulable). We consider an  $\text{ELIOM}_\varepsilon$  server declaration  $D_s$ ; the  $\text{ELIOM}_\varepsilon$  environments  $\rho_s, \rho_c$  and  $\gamma = \gamma_f \cup \gamma_r$ ; the target environment  $\hat{\rho}_s, \hat{\rho}_c$  and  $\hat{\gamma} = \hat{\gamma}_f \cup \hat{\gamma}_r \cup \zeta$ .

If the expression  $e$  has valid server and client executions:

$$D \xrightarrow{\rho_s}_s V, \mu, \theta_s \qquad \mu \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma'}_c \{\}, \theta_c$$

and the following invariants hold:

$$\hat{\rho}_c \simeq_{\zeta}^c \rho_c \qquad \hat{\rho}_s \simeq^s \rho_s \qquad \hat{\gamma} \simeq^c \gamma \qquad FCE(\hat{\gamma}_f, D) \qquad FCE(\hat{\gamma}_f, \rho_s)$$

Then  $\widehat{D} = D[\{\{ e_i \}\}_{\mathbf{f}} \mapsto \mathbf{fragment} \mathbf{f} \overline{x_{i,j}}]_i$  have an equivalent execution.

$$\frac{}{\widehat{D} \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \widehat{V}, \xi_{\bullet}, \{\}, \widehat{\theta}_s} \quad \frac{}{\mathbf{exec} \ () , \xi_{\bullet} ++ \mathbf{end} \xrightarrow{\widehat{\rho}_c | \widehat{\gamma} \rightarrow \widehat{\gamma}'}_{\text{ML}_c} \varepsilon, [\ ], \widehat{\theta}_c}$$

with the following invariants:

$$\widehat{\gamma}' \simeq^c \gamma' \quad FCE(\widehat{\gamma}'_f, V) \quad \widehat{V} \simeq^s V \quad \widehat{\theta}_s \simeq^s \theta_s \quad \widehat{\theta}_c \simeq^{\varepsilon}_{\zeta'} \theta_c$$

#### 5.4.4 Mixed structures

**Lemma 3** (Structures are simulable). We consider a slicable structure  $S$ ; the  $\text{ELIOM}_{\varepsilon}$  environments  $\rho_s, \rho_c$  and  $\gamma = \gamma_f \cup \gamma_r$ ; the target environment  $\widehat{\rho}_s, \widehat{\rho}_c$  and  $\widehat{\gamma} = \widehat{\gamma}_f \cup \widehat{\gamma}_r \cup \zeta$ .

If the structure has valid server and client executions:

$$S \xrightarrow{\rho_s}_m V_s, \mu, \theta_s \quad \mu \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c, \theta_c$$

and the following invariants hold:

$$\widehat{\rho}_c \simeq^{\varepsilon}_{\zeta} \rho_c \quad \widehat{\rho}_s \simeq^s \rho_s \quad \widehat{\gamma} \simeq^c \gamma \quad FCE(\widehat{\gamma}_f, S) \quad FCE(\widehat{\gamma}_f, \rho_s)$$

then for any  $\xi$ , the compiled structures have equivalent executions

$$\langle S \rangle_s \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \widehat{V}_s, \xi_{\bullet}, \zeta_{\bullet}, \widehat{\theta}_s \quad \langle S \rangle_c, \xi_{\bullet} ++ \xi \xrightarrow{\widehat{\rho}_c | \zeta_{\bullet} \cup \widehat{\gamma} \rightarrow \widehat{\gamma}'}_{\text{ML}_c} \widehat{V}_c, \xi', \widehat{\theta}_c$$

with the following invariants:

$$\begin{array}{ccc} \widehat{\gamma}' \simeq^c \gamma' & \widehat{V}_s \simeq^s V_s & \widehat{\theta}_s \simeq^s \theta_s \\ FCE(\widehat{\gamma}'_f, V_s) & \widehat{V}_c \simeq^{\varepsilon}_{\zeta'} V_c & \widehat{\theta}_c \simeq^{\varepsilon}_{\zeta'} \theta_c \end{array}$$

*Proof.* We consider a slicable structure  $S$ ; the  $\text{ELIOM}_{\varepsilon}$  environments  $\rho_s, \rho_c$  and  $\gamma = \gamma_f \cup \gamma_r$ ; the target environment  $\widehat{\rho}_s, \widehat{\rho}_c$  and  $\widehat{\gamma} = \widehat{\gamma}_f \cup \widehat{\gamma}_r \cup \zeta$ . such that

$$\widehat{\rho}_c \simeq^{\varepsilon}_{\zeta} \rho_c \quad \widehat{\rho}_s \simeq^s \rho_s \quad \widehat{\gamma} \simeq^c \gamma \quad FCE(\widehat{\gamma}_f, S) \quad FCE(\widehat{\gamma}_f, \rho_s)$$

We will now proceed by induction over the execution of  $S$ .

- *Case  $S = D_b; S'$  – Base declaration.*

We assume that the following executions hold:

$$\frac{D_b \xrightarrow{\rho_s}_b V_s, \varepsilon, \theta_s \quad S' \xrightarrow{\rho_s + V_s}_m V'_s, \mu, \theta'_s}{D_b; S' \xrightarrow{\rho_s}_m V_s + V'_s, (D_b; \mu), \theta_s @ \theta'_s} \quad \frac{D_b \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c, \theta_c \quad \mu \xrightarrow{\rho_c + V_c | \gamma \rightarrow \gamma'}_c V'_c, \theta'_c}{D_b; \mu \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c + V'_c, \theta_c @ \theta'_c}$$

Let us consider the executions of  $D_b$ . By definition of base, it contains neither injections nor fragments. By Proposition 5,  $D_b \xrightarrow{\rho}_b V_s, \varepsilon, \theta_s$  and  $D_b \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_b, \theta_c$  both

correspond to  $D_b \xRightarrow{\rho_s} V_s, \varepsilon, \theta_s$  and  $D_b \xRightarrow{\rho_c} V_c, \varepsilon, \theta_c$  respectively. By definition, base fragments can't be present, hence we also have  $FCE(\hat{\gamma}_f, V_s)$

Additionally, the compilation functions are the identity on base, which mean that  $\langle D_b \rangle_s$  and  $\langle D_b \rangle_c$  contains only ML constructs. The reduction relation over  $ML_s$  and  $ML_c$  coincide with the ML one on the ML fragment of the language. Hence, for any  $\xi$ , we have  $\langle D_b \rangle_s \xRightarrow{\rho_s}_{ML_s} \hat{V}_s, [ ], \{ \}, \hat{\theta}_s$  and  $\langle D_b \rangle_c, \xi \xRightarrow{\rho_c | \hat{\gamma} \rightarrow \hat{\gamma}}_{ML_c} \hat{V}_c, \xi, \hat{\theta}_c$  with

$$\hat{V}_s \simeq^s V_s \quad \hat{V}_c \simeq_{\hat{\gamma}}^c V_c \quad \hat{\theta}_s \simeq^s \theta_s \quad \hat{\theta}_c \simeq_{\hat{\gamma}}^c \theta_c$$

Let us consider the execution of  $S'$  and  $\mu$ . We easily have the following properties:

$$\hat{\rho}_c + \hat{V}_c \simeq_{\hat{\gamma}}^c \rho_c + V_c \quad \hat{\rho}_s + \hat{V}_s \simeq^s \rho_s + V_s \quad \hat{\gamma} \simeq^c \gamma \quad FCE(\hat{\gamma}_f, S) \quad FCE(\hat{\gamma}_f, \rho_s + V_s)$$

hence, by induction on the execution of  $S'$  and  $\mu'$ , we have  $\langle S' \rangle_s \xRightarrow{\hat{\rho}_s + \hat{V}_s}_{ML_s} \hat{V}'_s, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}'_s$  and  $\langle S' \rangle_c, \xi_{\bullet} ++ \xi \xRightarrow{\hat{\rho}_c + \hat{V}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{ML_c} \hat{V}'_c, \xi, \hat{\theta}'_c$  for any  $\xi$ , with

$$\begin{array}{ccc} \hat{\gamma}' \simeq^c \gamma' & \hat{V}'_s \simeq^s V'_s & \hat{\theta}'_s \simeq^s \theta'_s \\ FCE(\hat{\gamma}'_f, V'_s) & \hat{V}'_c \simeq_{\hat{\gamma}'}^c V'_c & \hat{\theta}'_c \simeq_{\hat{\gamma}'}^c \theta'_c \end{array}$$

We can then build the following derivations:

$$\begin{array}{c} \frac{\langle D_b \rangle_s \xRightarrow{\hat{\rho}_s}_{ML_s} \hat{V}_s, [ ], \{ \}, \theta_s \quad \langle S' \rangle_s \xRightarrow{\hat{\rho}_s + \hat{V}_s}_{ML_s} \hat{V}'_s, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}'_s}{\langle D_b; S' \rangle_s \xRightarrow{\hat{\rho}_s}_{ML_s} \hat{V}_s + \hat{V}'_s, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}_s @ \hat{\theta}'_s} \\ \frac{\langle D_b \rangle_c, \xi_{\bullet} ++ \xi \xRightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{ML_c} \hat{V}_c, \xi_{\bullet} ++ \xi, \hat{\theta}_c \quad \langle S' \rangle_c, \xi_{\bullet} ++ \xi \xRightarrow{\hat{\rho}_c + \hat{V}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{ML_c} \hat{V}'_c, \xi, \hat{\theta}'_c}{\langle D_b; S' \rangle_c, \xi_{\bullet} ++ \xi \xRightarrow{\rho_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{ML_c} \hat{V}_c + \hat{V}'_c, \xi, \hat{\theta}_c @ \hat{\theta}'_c} \end{array}$$

and the following invariants are easily verified:

$$\begin{array}{ccc} \hat{\gamma}' \simeq^c \gamma' & \hat{V}_s + \hat{V}'_s \simeq^s V_s + V'_s & \hat{\theta}_s @ \hat{\theta}'_s \simeq^s \theta_s @ \theta'_s \\ FCE(\hat{\gamma}'_f, V_s + V'_s) & \hat{V}_c + \hat{V}'_c \simeq_{\hat{\gamma}'}^c V_c + V'_c & \hat{\theta}_c @ \hat{\theta}'_c \simeq_{\hat{\gamma}'}^c \theta_c @ \theta'_c \end{array} \quad \blacksquare$$

- *Case  $S = D_s; S'$  – Server declaration.*

We assume that the following executions hold:

$$\begin{array}{c} \frac{D_s \xRightarrow{\rho_s}_s V_s, \mu, \theta_s \quad S' \xRightarrow{\rho_s + V_s}_m V'_s, \mu', \theta'_s}{D_s; S' \xRightarrow{\rho_s}_m V_s + V'_s, (\text{bind env } \mathbf{f}_i)_i; \mu; \mu', \theta_s @ \theta'_s} \\ \frac{(\text{bind env } \mathbf{f}_i)_i \xRightarrow{\rho_c | \gamma \rightarrow \gamma'}_c \{ \}, \langle \rangle \quad \mu \xRightarrow{\rho_c | \gamma' \rightarrow \gamma''}_c V_c, \theta_c \quad \mu' \xRightarrow{\rho_c | \gamma'' \rightarrow \gamma'''}_c V'_c, \theta'_c}{(\text{bind env } \mathbf{f}_i)_i; \mu; \mu' \xRightarrow{\rho_c | \gamma \rightarrow \gamma''}_c V_c + V'_c, \theta_c @ \theta'_c} \end{array}$$

Let us note  $\{\{ e_i \}\}_{\mathbf{f}_i}$  the fragments syntactically present in  $D_s$ . let us note  $\{\{ e_j \}\}_{\mathbf{f}_j}$  the fragments *executed* during the reduction of  $D_s$  and  $\mathbf{r}_j$  the associated fresh variables.

We have the following compilations:

$$\begin{aligned}\langle D_s; S' \rangle_s &= \langle D_s \rangle_s[\{\{ e_i \}\}_{\mathbf{f}_i} \mapsto \mathbf{fragment} \mathbf{f}_i \overline{x_{i,k}}_i; \mathbf{end} (); \langle S' \rangle_s] \\ \langle D_s; S' \rangle_c &= (\mathbf{bind} \mathbf{f}_i = \lambda \overline{x_{i,k}}. e_i[f_{i,k} \% x_{i,k} \mapsto x_{i,k}]; )_i; \mathbf{exec} (); \langle S' \rangle_c\end{aligned}$$

After hoisting, converters can only be the **serial** or **frag**. Its server and client parts are the identity, hence we simply omit them. We also note that  $\langle D_s \rangle_s$  differs with  $D_s$  only on type annotations and type declarations which are ignored by reduction relations. We note  $\widehat{D}_s = \langle D_s \rangle_s[\{\{ e_i \}\} \mapsto \mathbf{fragment} \mathbf{f}_i \overline{x_{i,k}}_i]$ .

Let us consider the reduction of  $(\mathbf{bind} \mathbf{f}_i = \lambda \overline{x_i}. \widehat{e_i})_i$ . Let us note  $e'_i = e_i[f_{i,j} \% x_{i,j} \mapsto x_{i,j}]$ . For any queue  $\xi$ , we have the following reduction:

$$\begin{aligned}& \frac{\lambda x_1 \dots x_m. e'_i, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma}_i \rightarrow \widehat{\gamma}_i} \text{ML}_c \lambda \overline{x_i}. \widehat{\rho}_c. e'_i, \xi, \langle \rangle}{\forall i, \quad \mathbf{bind} \mathbf{f}_i = \lambda \overline{x_i}. e'_i, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma}_i \rightarrow \widehat{\gamma}_{i+1}} \text{ML}_c \{\}, \xi, \langle \rangle} \\ & \frac{}{(\mathbf{bind} \mathbf{f}_i = \lambda \overline{x_i}. e'_i)_i, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma}_1 \rightarrow \widehat{\gamma}_{n+1}} \text{ML}_c \{\}, \xi, \langle \rangle}\end{aligned}$$

where  $\widehat{\gamma}_1 = \widehat{\gamma}$  and  $\widehat{\gamma}_{i+1} = \widehat{\gamma}_i \cup \{\mathbf{f}_i \mapsto \lambda x_1 \dots x_m. \widehat{\rho}_c. e'_i\}$ . Let  $\gamma_{bind}$  be  $\{\mathbf{f}_i \mapsto \lambda x_1 \dots x_m. \widehat{\rho}_c. e'_i\}_i$ . We note  $\widehat{\gamma}' = \widehat{\gamma}_{n+1} = \widehat{\gamma} \cup \gamma_{bind}$  and  $\widehat{\gamma}'_f = \widehat{\gamma}_f \cup \gamma_{bind}$ .

Since  $(\mathbf{bind} \mathbf{env} \mathbf{f}_i)_i \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c \{\}, \langle \rangle$ , we have  $\gamma' = \gamma \cup \{\mathbf{f}_i \mapsto \rho_c\}_i$ . and  $\rho_c \simeq^c \widehat{\rho}_c$ , we have that  $\gamma' \simeq^c \widehat{\gamma}'$ . Furthermore, given one of the  $\mathbf{f}_i$  in  $\gamma_{bind}$ , each fragment annotated with this  $\mathbf{f}_i$  syntactically appear in  $D_s$  by uniqueness of the annotation function. This also holds inside functors, since each  $\mathbf{f}_i$  will be prefixed by a unique module reference. Hence  $FCE(\widehat{\gamma}'_f, D_s)$  and  $FCE(\widehat{\gamma}'_f, \rho_s)$ .

We now have all the ingredients to uses Corollary 1 on the execution of  $D_s$  and  $\mu$ . This gives us the following reductions:

$$\frac{}{\widehat{D}_s \xrightarrow{\widehat{\rho}_s} \text{ML}_s \widehat{V}_s, \xi_\bullet, \{\}, \widehat{\theta}_s} \quad \frac{}{\mathbf{exec} (), \xi_\bullet ++ \mathbf{end} \xrightarrow{\widehat{\rho}_c | \widehat{\gamma}' \rightarrow \widehat{\gamma}''} \text{ML}_c \varepsilon, [], \widehat{\theta}_c}$$

with the following invariants:

$$\widehat{\gamma}'' \simeq^c \gamma'' \quad FCE(\widehat{\gamma}''_f, V) \quad \widehat{V} \simeq^s V \quad \widehat{\theta}_s \simeq^s \theta_s \quad \widehat{\theta}_c \simeq^c \theta_c$$

We remark that  $\zeta'' = \zeta$  since no injection took place during a server section and that  $\widehat{\gamma}''_f = \widehat{\gamma}'_f$ , by definition of the reduction for **exec**.

We now consider the execution of  $S'$ . The following invariants holds:

$$\widehat{\rho}_c \simeq^c \rho_c \quad \widehat{\rho}_s + \widehat{V}_s \simeq^s \rho_s + V_s \quad \widehat{\gamma}'' \simeq^c \gamma'' \quad FCE(\widehat{\gamma}''_f, S') \quad FCE(\widehat{\gamma}''_f, \rho'_s + V_s)$$

By induction on the execution of  $S'$  and  $\mu'$ , we have  $\langle S' \rangle_s \xrightarrow{\widehat{\rho}_s + \widehat{V}_s}_{\text{ML}_s} \widehat{V}'_s, \xi'_s, \zeta'_s, \widehat{\theta}'_c$  and  $\langle S' \rangle_c, \xi'_s \text{++} \xi'_c \xrightarrow{\widehat{\rho}_c | \widehat{\gamma}'' \cup \zeta'_s \rightarrow \widehat{\gamma}'''}_{\text{ML}_c} \widehat{V}'_c, \xi'_c, \widehat{\theta}'_c$  where

$$\begin{array}{lll} \widehat{\gamma}''' \simeq^c \gamma''' & \widehat{V}'_s \simeq^s V'_s & \widehat{\theta}'_s \simeq^s \theta'_s \\ FCE(\widehat{\gamma}'''_f, V'_s) & \widehat{V}'_c \simeq^c_{\widehat{\gamma}'} V'_c & \widehat{\theta}'_c \simeq^c_{\widehat{\gamma}'} \theta'_c \end{array}$$

Finally, we can construct the following executions:

$$\frac{\widehat{D}_s \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \widehat{V}_s, \xi_\bullet, \{\}, \widehat{\theta}_s \quad \langle S' \rangle_s \xrightarrow{\rho_s + V_s}_{\text{ML}_s} \widehat{V}'_s, \xi'_s, \zeta'_s, \widehat{\theta}'_s}{\widehat{D}_s; \text{end } (); \langle S' \rangle_s \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \widehat{V}_s + \widehat{V}'_s, \xi_\bullet \text{++} \text{end++} \xi'_s, \zeta'_s, \widehat{\theta}_s @ \widehat{\theta}'_s} \\ \frac{\widehat{\mu}, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma} \cup \zeta'_s \rightarrow \widehat{\gamma}'''}_{\text{ML}_c} \{\}, \xi, \langle \rangle \quad \text{exec } (); \langle S' \rangle_c, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma} \cup \zeta'_s \rightarrow \widehat{\gamma}'''}_{\text{ML}_c} \widehat{V}'_c, \xi'_c, \widehat{\theta}'_c @ \widehat{\theta}'_c}{\widehat{\mu}; \text{exec } (); \langle S' \rangle_c, \xi \xrightarrow{\widehat{\rho}_c | \widehat{\gamma} \cup \zeta'_s \rightarrow \widehat{\gamma}'''}_{\text{ML}_c} \widehat{V}'_c, \xi'_c, \widehat{\theta}'_c @ \widehat{\theta}'_c}$$

where  $\xi = \xi_\bullet \text{++} \text{end++} \xi'_s \text{++} \xi'_c$  and  $\widehat{\mu} = (\text{bind } \mathbf{f}_i = \lambda \bar{x}_i. e'_i)_i$ . We verify the following invariants:

$$\begin{array}{lll} \widehat{\gamma}''' \simeq^c \gamma''' & \widehat{V}_s + \widehat{V}'_s \simeq^s V_s + V'_s & \widehat{\theta}_s + \widehat{\theta}'_s \simeq^s \theta_s + \theta'_s \\ FCE(\widehat{\gamma}'''_f, V_s + V'_s) & \widehat{V}'_c \simeq^c_{\widehat{\gamma}'''} V'_c & \widehat{\theta}'_c + \widehat{\theta}'_c \simeq^c_{\widehat{\gamma}'''} \theta_c + \theta'_c \end{array} \quad \blacksquare$$

• *Case  $S = D_c; S'$  – Client declaration.*

We assume that the following executions hold:

$$\frac{D_c \xrightarrow{\rho_s}_{c/s} \overline{D}_c, \varepsilon, \langle \rangle \quad S' \xrightarrow{\rho_s}_m V'_s, \mu', \theta'_s}{D_c; S' \xrightarrow{\rho_s}_m V'_s, (\overline{D}_c; \mu'), \theta'_s} \quad \frac{\overline{D}_c \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c, \theta_c \quad \mu' \xrightarrow{\rho_c + V_c | \gamma \rightarrow \gamma'}_c V'_c, \theta'_c}{\overline{D}_c; \mu' \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c + V'_c, \theta_c @ \theta'_c}$$

Let us note  $f_i \% x_i$  the injections in  $D_c$  and  $\mathbf{x}_i$  the associated fresh variables. Since hoisting has been applied, all the  $f_i$  are either **serial** or **frag**. Furthermore, no fragments are executed due to the execution of injections and  $\overline{D}_c = D_c[f_i \% x_i \mapsto \rho_s(x_i)]_i$ .

We have the following compilations:

$$\begin{array}{l} \langle D_c; S \rangle_s = (\text{injection } \mathbf{x}_i \ x_i;)_i \text{end } (); \langle S \rangle_s \\ \langle D_c; S \rangle_c = \text{exec } (); D_c[f_i \% x_i \mapsto \mathbf{x}_i]_i; \langle S \rangle_c \end{array}$$

In the rest of this proof, we note  $\widehat{D}_c = D_c[f_i \% x_i \mapsto \mathbf{x}_i]_i$ .

We consider the server reduction  $D_c \Rightarrow_{c/s} \overline{D}_c$ . We know that  $\overline{D}_c = D_c[f_i \% x_i \mapsto \downarrow \rho_s(x_i)]_i$ . Let us note  $v_i = \rho_s(x_i)$ . We can build the following  $\text{ML}_s$  reduction:

$$\frac{\widehat{\rho}_s(x_i) = \widehat{v}_i}{\forall i. \text{injection } \mathbf{x}_i \ x_i \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \varepsilon, [\ ], \{\mathbf{x}_i \mapsto \downarrow \widehat{v}_i\}, \langle \rangle} \\ (\text{injection } \mathbf{x}_i \ x_i;)_i; \text{end } () \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \varepsilon, \text{end}, \{\mathbf{x}_i \mapsto \downarrow \widehat{v}_i\}_i, \langle \rangle$$

Since  $\hat{\rho}_s \simeq^s \rho_s$ , we also have that  $\hat{v}_i \simeq^s v_i$  and  $\downarrow \hat{v}_i \simeq_{\zeta}^c \downarrow v_i$  for each  $i$ . We note  $\zeta_{\bullet} = \{\mathbf{x}_i \mapsto \downarrow \hat{v}_i\}_i$ . By definition of the slicing relation, the  $\mathbf{x}_i$  are fresh, hence they are not bound in  $\hat{\gamma}$ . We can thus construct the global environment  $\hat{\gamma}' = \hat{\gamma} \cup \zeta_{\bullet}$ . Since we only extend the part with injection references, we still have that  $\gamma \simeq^c \hat{\gamma}'$ .

We now consider the client reduction  $\overline{D}_c \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma}_c V_c, \theta_c$ . We know that  $\overline{D}_c$  is equal to  $D_c[f_i \% x_i \mapsto \downarrow v_i]_i$ , hence the reduction tree contains for each  $i$  a reduction  $\downarrow v_i \xrightarrow{\mid \gamma \rightarrow \gamma}_c \downarrow v_i, \langle \rangle$ . To obtain a reduction of  $\hat{D}_c = D_c[f_i \% x_i \mapsto \mathbf{x}_i]_i$ , we simply substitute each of these sub-reduction by one of the form  $\mathbf{x}_i, \xi \xrightarrow{\mid \hat{\gamma}' \rightarrow \hat{\gamma}'}_{\text{ML}_c} \downarrow \hat{v}_c, \xi, \langle \rangle$ . for any queue  $\xi$ . By Lemma 1, we can build the following reduction:

$$\overline{\hat{D}_c, \xi \xrightarrow{\hat{\rho}_c \mid \hat{\gamma}' \rightarrow \hat{\gamma}'}_{\text{ML}_c} \hat{V}_c, \xi, \hat{\theta}_c}$$

where  $\hat{V}_c \simeq_{\zeta'}^c V_c$  and  $\hat{\theta}_c \simeq_{\zeta'}^c \theta_c$ , for any queue  $\xi$ .

We now consider the execution of  $S'$ . We have the following properties:

$$\hat{\rho}_c + \hat{V}_c \simeq_{\zeta'}^c \rho_c + V_c \quad \hat{\rho}_s \simeq^s \rho_s \quad \hat{\gamma}' \simeq^c \gamma \quad FCE(\hat{\gamma}_f, S') \quad FCE(\hat{\gamma}_f, \rho_s)$$

By induction on the execution of  $S'$  and  $\mu'$ , we have  $\langle S' \rangle_c \xrightarrow{\hat{\rho}_c + \hat{V}_c}_{\text{ML}_s} \hat{V}'_s, \xi'_s, \zeta'_s, \hat{\theta}'_c$  and  $\langle S' \rangle_s, \xi'_s \mapsto \xi' \xrightarrow{\hat{\rho}_s \mid \zeta'_s \cup \hat{\gamma}' \rightarrow \hat{\gamma}''}_{\text{ML}_c} \hat{V}'_c, \xi', \hat{\theta}'_c$  where

$$\begin{array}{lll} \hat{\gamma}'' \simeq^c \gamma' & \hat{V}'_s \simeq_{\zeta''}^s V'_s & \hat{\theta}'_s \simeq_{\zeta''}^s \theta'_s \\ FCE(\hat{\gamma}'_f, V'_s) & \hat{V}'_c \simeq_{\zeta''}^c V'_c & \hat{\theta}'_c \simeq_{\zeta''}^c \theta'_c \end{array}$$

Finally, we can build the following derivations:

$$\begin{array}{c} \overline{\forall i, \text{ injection } \mathbf{x}_i \ x_i \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \varepsilon, [\ ], \zeta_{\bullet}, \langle \rangle} \quad \overline{\langle S' \rangle_s \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \hat{V}'_s, \xi'_s, \zeta'_s, \hat{\theta}'_s} \\ \hline (\text{injection } \mathbf{x}_i \ (f_i^s \ x_i);)_i \text{ end } (); \langle S' \rangle_s \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \hat{V}'_s, \text{end} \mapsto \xi'_s, \zeta_{\bullet} \cup \zeta'_s, \hat{\theta}'_s \\ \hline \hat{D}_c, \xi'_s \mapsto \xi' \xrightarrow{\hat{\rho}_c \mid \hat{\gamma}' \cup \zeta'_s \rightarrow \hat{\gamma}''}_{\text{ML}_c} \hat{V}_c, \xi'_s \mapsto \xi', \hat{\theta}'_c \quad \langle S' \rangle_c, \xi'_s \mapsto \xi' \xrightarrow{\hat{\rho}_c + \hat{V}_c \mid \hat{\gamma}' \cup \zeta'_s \rightarrow \hat{\gamma}''}_{\text{ML}_c} \hat{V}'_c, \xi', \hat{\theta}'_c \\ \hline \hat{D}_c; \langle S' \rangle_c, \xi'_s \mapsto \xi' \xrightarrow{\hat{\rho}_c \mid \hat{\gamma}' \cup \zeta'_s \rightarrow \hat{\gamma}''}_{\text{ML}_c} \hat{V}_c + \hat{V}'_c, \xi', \hat{\theta}_c @ \hat{\theta}'_c \\ \hline \text{exec } (); \hat{D}_c; \langle S' \rangle_c, \text{end} \mapsto \xi'_s \mapsto \xi' \xrightarrow{\hat{\rho}_c \mid \hat{\gamma}' \cup \zeta_{\bullet} \cup \zeta'_s \rightarrow \hat{\gamma}''}_{\text{ML}_c} \hat{V}_c + \hat{V}'_c, \xi', \hat{\theta}_c @ \hat{\theta}'_c \end{array}$$

We verify the following invariants:

$$\begin{array}{lll} \hat{\gamma}'' \simeq^c \gamma' & \hat{V}_s \simeq^s V_s & \hat{\theta}_s \simeq^s \theta_s \\ FCE(\hat{\gamma}'_f, V_s + V'_s) & \hat{V}_c + \hat{V}'_c \simeq_{\hat{\gamma}''}^c V_c + V'_c & \hat{\theta}_c + \hat{\theta}'_c \simeq_{\hat{\gamma}''}^c \theta_c + \theta'_c \quad \blacksquare \end{array}$$

- *Case  $\text{module}_m X = M; S'$  – Declaration of a mixed module.*

We assume that the following executions hold:

$$\begin{array}{c}
\frac{M \xrightarrow{\rho_s}_m V_s, M^c, \mu, \theta_s}{\text{module}_m X = M \xrightarrow{\rho_s}_m \{X \mapsto V_s\}, \text{module } X = M^c; \mu, \theta_s} \quad S' \xrightarrow{\rho_s + \{X \mapsto V_s\}}_m V'_s, \mu', \theta'_s \\
\hline
\text{module}_m X = M; S' \xrightarrow{\rho_s}_m \{X \mapsto V_s\} + V'_s, (\mu; \text{module } X = M^c; \mu'), \theta_s @ \theta'_s \\
\\
\frac{M^c \xrightarrow{\rho_c | \gamma' \rightarrow \gamma''}_c V_c, \theta'_c}{\mu \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c \{\}, \theta_c} \quad \frac{\text{module } X = M^c \xrightarrow{\rho_c | \gamma' \rightarrow \gamma''}_c \{X \mapsto V_c\}, \theta'_c \quad \mu' \xrightarrow{\rho_c + \{X \mapsto V_c\} | \gamma'' \rightarrow \gamma''}_c V'_c, \theta''_c}{\mu; \text{module } X = M^c; \mu' \xrightarrow{\rho_c | \gamma \rightarrow \gamma''}_c \{X \mapsto V_c\} + V'_c, \theta_c @ \theta'_c @ \theta''_c}
\end{array}$$

Let us assume that we can build the following reductions

$$\begin{array}{c}
\frac{}{\langle \text{module}_m X = M \rangle_s \xrightarrow{\widehat{\rho}_s}_{\text{ML}_s} \{X \mapsto \widehat{V}_s\}, \xi_\bullet, \zeta, \widehat{\theta}_s} \\
\\
\frac{}{\langle \text{module}_m X = M \rangle_c, \xi_\bullet ++ \xi \xrightarrow{\widehat{\rho}_c | \zeta \cup \widehat{\gamma} \rightarrow \widehat{\gamma}'}_{\text{ML}_c} \{X \mapsto \widehat{V}_c\}, \xi, \widehat{\theta}_c}
\end{array}$$

for any  $\xi$ , and that the following invariants hold:

$$\begin{array}{ccc}
\widehat{\gamma}' \simeq^c \gamma'' & \widehat{V}_s \simeq^s V_s & \widehat{\theta}_s \simeq^s \theta_s \\
FCE(\widehat{\gamma}'_f, V_s) & \widehat{V}_c \simeq^{\xi_\bullet} V_c & \widehat{\theta}_c \simeq^{\xi_\bullet} \theta_c + \theta'_c
\end{array}$$

By induction on the execution of  $S'$  and  $\mu'$ , we can build the following reduction:

$$\langle S' \rangle_s \xrightarrow{\widehat{\rho}_s + \{X \mapsto \widehat{V}_s\}}_{\text{ML}_s} V'_s, \xi'_\bullet, \zeta', \theta'_s \text{ and } \langle S' \rangle_c, \xi'_\bullet ++ \xi \xrightarrow{\widehat{\rho}_c + \{X \mapsto \widehat{V}_c\} | \zeta' \cup \widehat{\gamma}' \rightarrow \widehat{\gamma}''}_{\text{ML}_c} V'_c, \xi, \theta'_c, \text{ which}$$

allows us to conclude.

To build the compiled reduction, we will operate by case analysis over  $M$ .

- *Subcase  $M = \text{struct } S \text{ end}_{\mathbf{X}}$  – Declaration of a mixed structure.*

We have  $\mu = \text{bind } \mathbf{X} = (\text{struct } \mu_0 \text{ end})$  and  $M^c = \mathbf{X}$  with the following reductions:

$$\begin{array}{c}
\frac{S \xrightarrow{\rho_s}_m V_s, \mu, \theta_s}{\text{struct } S \text{ end} \xrightarrow{\rho_s}_m V_s + \{\text{Dyn} \mapsto \mathbf{X}\}, \mathbf{X}, \mu, \theta_s} \\
\\
\frac{\frac{\mu_0 \xrightarrow{\rho_c | \gamma \rightarrow \gamma'}_c V_c, \theta_c}{\text{bind } \mathbf{X} = \text{struct } \mu_0 \text{ end; module } X = \mathbf{X}} \quad \frac{\gamma''(\mathbf{X}) = V_c}{\mathbf{X} \xrightarrow{\rho_c | \gamma'' \rightarrow \gamma''}_c V_c, \langle \rangle}}{\text{bind } \mathbf{X} = \text{struct } \mu_0 \text{ end; module } X = \mathbf{X} \xrightarrow{\rho_c | \gamma \rightarrow \gamma''}_c, \theta_c}
\end{array}$$

where  $\gamma'' = \gamma' \cup \{\mathbf{X} \mapsto V_c\}$ .



We have the following compilations:

$$\begin{aligned}
\langle \text{module}_m X = \text{struct } S \text{ end}_{\mathbf{X}} \rangle_s &= \begin{array}{l} \text{module } X = \text{struct} \\ \text{module Dyn} = \mathbf{X}; \\ \langle S \rangle_s \\ \text{end} \end{array} \\
\langle \text{module}_m X = \text{struct } S \text{ end}_{\mathbf{X}} \rangle_c &= \begin{array}{l} \text{bind}_m \mathbf{X} = \text{struct } \langle S \rangle_c \text{ end;} \\ \text{module } X = \mathbf{X}; \end{array}
\end{aligned}$$

By induction on the execution of  $S$  and  $\mu$ , we have  $\langle S \rangle_s \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \hat{V}_s, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}_s$  and  $\langle S \rangle_c, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \hat{V}_c, \xi, \hat{\theta}_c$  with the following invariants:

$$\begin{array}{lll}
\hat{\gamma}' \simeq^c \gamma' & \hat{V}_s \simeq^s V_s & \hat{\theta}_s \simeq^s \theta_s \\
FCE(\hat{\gamma}'_f, V_s) & \hat{V}_c \simeq_{\hat{\gamma}'}^c V_c & \hat{\theta}_c \simeq_{\hat{\gamma}'}^c \theta_c
\end{array}$$

We can then build the following executions:

$$\begin{array}{c}
\frac{\langle S \rangle_s \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \hat{V}_s, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}_s}{\langle \text{module}_m X = \text{struct } S \text{ end}_{\mathbf{X}} \rangle_s \xrightarrow{\hat{\rho}_s}_{\text{ML}_s} \left\{ X \mapsto \{\text{Dyn} \mapsto \mathbf{X}\} + \hat{V}_s \right\}, \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}_s} \\
\frac{\langle S \rangle_c, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \hat{V}_c, \xi, \hat{\theta}_c}{\text{module } X = \text{struct } \langle S \rangle_c \text{ end}, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \left\{ X \mapsto \hat{V}_c \right\}, \xi, \hat{\theta}_c} \\
\frac{\text{module } X = \text{struct } \langle S \rangle_c \text{ end}, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \left\{ X \mapsto \hat{V}_c \right\}, \xi, \hat{\theta}_c}{\langle \text{module}_m X = \text{struct } S \text{ end} \rangle_c, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c | \hat{\gamma} \cup \zeta_{\bullet} \rightarrow \hat{\gamma}''}_{\text{ML}_c} \left\{ X \mapsto \hat{V}_c \right\}, \xi, \hat{\theta}_c}
\end{array}$$

Where  $\hat{\gamma}'' = \hat{\gamma}' \cup \left\{ \mathbf{X} \mapsto \hat{V}_c \right\}$ . We verify the following invariants:

$$\begin{array}{lll}
\hat{\gamma}'' \simeq^c \gamma' & \hat{V}_s + \{\text{Dyn} \mapsto \mathbf{X}\} \simeq^s V_s + \{\text{Dyn} \mapsto \mathbf{X}\} & \hat{\theta}_s \simeq^s \theta_s \\
FCE(\hat{\gamma}''_f, V_s) & \hat{V}_c \simeq_{\hat{\gamma}'}^c V_c & \hat{\theta}_c \simeq_{\hat{\gamma}'}^c \theta_c
\end{array}$$

which concludes. ■

- *Subcase  $M = \text{functor}(X_i : \mathcal{M}_i)_i \text{struct } S \text{ end}_{\mathbf{F}}$  – Declaration of a mixed functor.*

In this case, we have the client program  $\mu = \text{bind env } \mathbf{F}$  and the module expression

$M^c = \text{functor}(X_i : \mathcal{M}_i)_i \text{struct } S|_c \text{end}$ . The following reductions hold:

$$\frac{\left( \begin{array}{c} \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S \\ \text{end}_F \end{array} \right) \xRightarrow{\rho_s}_m \{F \mapsto V_s\}, \left( \begin{array}{c} \text{bind env } \mathbf{F} \\ \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S|_c \\ \text{end} \end{array} \right), \langle \rangle}{\left( \begin{array}{c} \text{bind env } \mathbf{F} \\ \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S|_c \\ \text{end} \end{array} \right) \xRightarrow{\rho_c | \gamma \rightarrow \gamma'}_c \{F \mapsto V_c\}, \langle \rangle}$$

Where  $\gamma' = \gamma \cup \{\mathbf{F} \mapsto \rho_c\}$  and the following values:

$$\begin{aligned} V_s &= \text{functor}(\rho_s)(X_i : \mathcal{M}_i)_i \text{struct } S \text{end} \\ V_c &= \text{functor}(\rho_c)(X_i : \mathcal{M}_i)_i \text{struct } S|_c \text{end} \end{aligned}$$

We recall that by hoisting, the body of the functors contains no injection, hence we don't need to evaluate server code in the client part.

We have the following compilations:

$$\begin{aligned} \left\langle \begin{array}{c} \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S \\ \text{end}_F \end{array} \right\rangle_s &= \begin{array}{c} \text{module } F(X_i : \langle \mathcal{M}_i \rangle_s)_i = \text{struct} \\ \text{module Dyn} = \text{fragment}_m \mathbf{F} (X_i.\text{Dyn})_i; \\ \langle S \rangle_s \\ \text{end} \end{array} \\ \left\langle \begin{array}{c} \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S \\ \text{end}_F \end{array} \right\rangle_c &= \begin{array}{c} \text{bind}_m \mathbf{F}(X_i : \langle \mathcal{M}_i \rangle_c)_i = \text{struct } \langle S \rangle_c \text{end}; \\ \text{module } F(X_i : \langle \mathcal{M}_i \rangle_c)_i = \text{struct } S|_c \text{end}; \end{array} \end{aligned}$$

We trivially have the following execution:

$$\begin{aligned} \left\langle \begin{array}{c} \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S \\ \text{end}_F \end{array} \right\rangle_s &\xRightarrow{\hat{\rho}_s}_{\text{ML}_s} \{F \mapsto \hat{V}_s\}, \xi_\bullet, \{\}, \langle \rangle \\ \left\langle \begin{array}{c} \text{module}_m F(X_i : \mathcal{M}_i)_i = \text{struct} \\ S \\ \text{end}_F \end{array} \right\rangle_c &\xRightarrow{\hat{\rho}_c | \hat{\gamma} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \{F \mapsto \hat{V}_c\}, \xi, \langle \rangle \end{aligned}$$

Where  $\gamma' = \gamma \cup \{\mathbf{F} \mapsto \widehat{V}_{\mathbf{F}}\}$  with the following values:

$$\begin{aligned}\widehat{V}_s &= \text{functor}(\widehat{\rho}_s)(X_i : \mathcal{M}_i)_i \text{struct } \langle S \rangle_s \text{ end} \\ \widehat{V}_c &= \text{functor}(\widehat{\rho}_c)(X_i : \mathcal{M}_i)_i \text{struct } S|_c \text{ end} \\ \widehat{V}_{\mathbf{F}} &= \text{functor}(\widehat{\rho}_c)(X_i : \mathcal{M}_i)_i \text{struct } \langle S \rangle_c \text{ end}\end{aligned}$$

We now need to show that the invariants still hold. We easily have that  $\widehat{V}_c \simeq_{\widehat{\gamma}}^c V_c$ . By definition of equivalence over mixed functors, we have  $\widehat{V}_s \simeq^s V_s$ . Indeed, the body of the functor in  $\widehat{V}_s$  is the server compilation of the body of the mixed functor  $V_s$  and the captured environments corresponds. Finally, we have that the body of  $\widehat{V}_{\mathbf{F}}$  is the client compilation of the body of the mixed functors and that the capture environment corresponds to  $\gamma'(\mathbf{F})$ . Thus we get that  $FCE(\widehat{\gamma}'_f, V_s)$ . By definition of the annotation function, the reference  $\mathbf{F}$  could not have appeared on a previously executed structure, hence we still have that  $FCE(\widehat{\gamma}'_f, \widehat{\rho}_s)$ .

Hence, all the following invariants are respected, which concludes.

$$\begin{array}{ll}\widehat{\gamma}' \simeq^c \gamma' & \widehat{V}_s \simeq^s V_s \\ FCE(\widehat{\gamma}'_f, V_s) & \widehat{V}_c \simeq_{\widehat{\gamma}''}^c V_c\end{array} \quad \blacksquare$$

Otherwise,  $M$  is a module expression. By definition of slicability,  $M$  does not syntactically contain any structure. In the general case, we should proceed by induction over module expressions. We will simply present the case of a mixed functor application where the functor returns a mixed structure.

We consider  $M = F(X_1) \dots (X_n)$ . We have  $M^c = F(X_1) \dots (X_n)$  with the following executions:

$$\frac{\begin{array}{l} F \xRightarrow{\rho_s}_m \text{functor}_m(\rho'_s)(Y_i : \mathcal{M}_i)_i \text{struct } S \text{ end}_{\mathbf{F}, \varepsilon, \langle \rangle} \quad X_i \xRightarrow{\rho_s}_m V_i^s, \varepsilon, \langle \rangle \\ V_i^s(\text{Dyn}) = \mathbf{R}_i \quad \mathbf{R} \text{ fresh} \quad S[\mathbf{f}_i \mapsto \mathbf{R}.\mathbf{f}_i]_i \xRightarrow{\rho'_s + \{Y_i \mapsto V_i^s\}_i}_m V, \mu_0, \theta \end{array}}{\begin{array}{l} F(X_1) \dots (X_n) \xRightarrow{\rho_s}_m V_s + \{\text{Dyn} \mapsto \mathbf{R}\}, \mu = \left( \begin{array}{l} \text{bind } \mathbf{R} = \text{struct} \\ (\text{module } Y_i = \mathbf{R}_i; )_i \\ \mu_0 \\ \text{end with } \mathbf{F} \end{array} \right), \theta \end{array}} \\ \text{module}_m X = F(X_1) \dots (X_n) \xRightarrow{\rho}_m V_s, (\mu; \text{module } X = F(X_1) \dots (X_n)), \theta_s\end{array}$$

$$\begin{array}{c}
\frac{\gamma(\mathbf{F}) = \rho_{\mathbf{F}} \quad \gamma(\mathbf{R}_i) = V_i^c \quad \mu_0 \xrightarrow{\rho_{\mathbf{F}} + \{Y_i \mapsto V_i^c\}_i \mid \gamma \rightarrow \gamma'}_c V_c, \theta_c}{\left( \begin{array}{l} \text{bind } \mathbf{R} = \text{struct} \\ \quad (\text{module } Y_i = \mathbf{R}_i;)_i \\ \quad \mu_0 \\ \text{end with } \mathbf{F} \end{array} \right) \xrightarrow{\rho_c \mid \gamma \rightarrow \gamma' \cup \{\mathbf{R} \mapsto V_c\}}_c \varepsilon, \theta_c} \\
\\
\frac{\begin{array}{l} F \xrightarrow{\rho_s \mid \dots \rightarrow \dots}_c \text{functor}(\rho'_c)(Y_i : \mathcal{M}_i)_i \text{struct } S_c \text{ end}, \langle \rangle \\ X_i \xrightarrow{\rho_c \mid \dots \rightarrow \dots}_c V_i, \langle \rangle \quad S_c \xrightarrow{\rho'_c + \{Y_i \mapsto V_i\}_i \mid \gamma' \cup \{\mathbf{R} \mapsto V_c\} \rightarrow \gamma''}_c V'_c, \theta'_c \end{array}}{\text{module } X = F(X_1) \dots (X_n) \xrightarrow{\rho_c \mid \gamma' \cup \{\mathbf{R} \mapsto V_c\} \rightarrow \gamma''}_c \{X \mapsto V'_c\}, \theta'_c}
\end{array}$$

We note  $V_F$  the value of  $F$  in  $\gamma_s$ , which is  $\text{functor}(\rho'_s)(Y_i : \mathcal{M}_i)_i \text{struct } S \text{ end}$ .

We have the following compilations:

$$\begin{aligned}
\langle \text{module}_m X = F(X_1) \dots (X_n) \rangle_s &= \text{module } X = F(X_1) \dots (X_n); \\
&\quad \text{end } (); \\
\langle \text{module}_m X = F(X_1) \dots (X_n) \rangle_c &= \text{module } X = F(X_1) \dots (X_n); \\
&\quad \text{exec } ();
\end{aligned}$$

Let us now look at the execution of the server application  $F(X_1) \dots (X_n)$ . By hypothesis,  $V_F$  is a mixed functor. By equivalence, we know that  $\hat{\rho}_s(F) = \hat{V}_F$  where  $\hat{V}_F \simeq^s V_F$ . By definition of the equivalence on server values,  $\hat{V}_F$  is of the following shape:

$$\hat{V}_F = \left( \begin{array}{l} \text{functor}(\hat{\rho}'_s)(Y_i : \mathcal{M}_i)_i \text{struct} \\ \quad \text{module Dyn} = \text{fragment}_m \mathbf{F} (Y_i.\text{Dyn})_i; \\ \quad \langle S \rangle_s \\ \text{end} \end{array} \right)$$

where  $\hat{\rho}'_s \simeq_{\hat{\gamma}}^s \rho'_s$ . For each  $i$  we note  $\hat{V}_i^s = \hat{\rho}_s(X_i)$ . By equivalence, we have  $\hat{V}_i^s \simeq^s V_i^s$ .

Furthermore, since  $\hat{\gamma} \simeq^c \gamma$  and by hypothesis,  $\mathbf{F}$  is also bound in  $\hat{\gamma}$ . We note  $V_{\mathbf{F}}$  the corresponding value. Since  $FCE(\hat{\gamma}, V_c)$  (via  $\rho_s$ ), then  $\hat{V}_{\mathbf{F}}$  is of the following shape:

$$\hat{V}_{\mathbf{F}} = \left( \begin{array}{l} \text{functor}(\hat{\rho}_{\mathbf{F}})(Y_i : \mathcal{M}_i)_i \text{struct} \\ \quad \langle S \rangle_c \\ \text{end} \end{array} \right)$$

where  $\hat{\rho}_{\mathbf{F}} \simeq_{\hat{\gamma}}^c \rho_{\mathbf{F}} = \gamma(\mathbf{F})$ . Additionally, for each  $i$  we note  $V_i^c = \hat{\gamma}(\mathbf{R}_i)$ . By equivalence, we have  $\hat{V}_i^c \simeq^s V_i^c$ .

We can now proceed by induction on  $S$  and  $\mu_0$  in the environment  $\hat{\gamma}$ ,  $\hat{\rho}_c \cup \{Y_i \mapsto \mathbf{R}_i\}$  and  $\hat{\rho}_s \cup \{\text{Dyn} \mapsto \mathbf{R}\} \cup \{Y_i \mapsto \hat{V}_i^s\}_i$ . We obtain the following reductions:

$$\langle S \rangle_s \xrightarrow{\hat{\rho}'_s \cup \{Y_i \mapsto \hat{V}_i^s\}_i}_{\text{ML}_s} \xi_{\bullet}, \zeta_{\bullet}, \hat{\theta}_s, \text{ and } \langle S \rangle_c, \xi_{\bullet} ++ \xi \xrightarrow{\hat{\rho}_c \cup \{Y_i \mapsto \mathbf{R}_i\}_i \mid \zeta_{\bullet} \cup \hat{\gamma} \rightarrow \hat{\gamma}'}_{\text{ML}_c} \hat{V}_c, \xi, \hat{\theta}_c$$

with the usual invariants. We can now build the following executions:

$$\begin{array}{c}
\text{fragment}_m \mathbf{F} (Y_i.\text{Dyn})_i \xrightarrow{\widehat{\rho}_s \cup \{Y_i \mapsto \widehat{V}_i^s\}_i} \text{ML}_s \mathbf{R}, \xi_{\mathbf{R}}, \varepsilon, \langle \rangle \\
\widehat{\rho}_s(F) = \widehat{V}_f \quad \widehat{\rho}_s(X_i) = \widehat{V}_i^s \quad \langle S \rangle_s \xrightarrow{\widehat{\rho}_s \cup \{\text{Dyn} \mapsto \mathbf{R}\} \cup \{Y_i \mapsto \widehat{V}_i^s\}_i} \text{ML}_s \xi_{\bullet}, \zeta_{\bullet}, \widehat{\theta}_s, \\
\text{module } X = F(X_1) \dots (X_n); \text{end } () \xRightarrow{\widehat{\rho}_s} \text{ML}_s \left\{ X \mapsto \widehat{V}_s \right\}, \xi_{\mathbf{R}} ++ \xi_{\bullet} ++ \text{end}, \zeta_{\bullet}, \widehat{\theta}_s \\
\widehat{\gamma}(\mathbf{F}) = \widehat{V}_{\mathbf{F}} \quad \langle S \rangle_c, \xi_{\bullet} ++ \text{end} ++ \xi \xrightarrow{\widehat{\rho}_{\mathbf{F}} \cup \{Y_i \mapsto \mathbf{R}_i\}_i \mid \zeta_{\bullet} \cup \widehat{\gamma} \rightarrow \widehat{\gamma}'} \text{ML}_c \widehat{V}_c, \text{end} ++ \xi, \widehat{\theta}_c \\
\text{exec } () , \xi_{\mathbf{R}} ++ \xi_{\bullet} ++ \text{end} ++ \xi \xrightarrow{\widehat{\rho}_c \mid \zeta_{\bullet} \cup \widehat{\gamma} \rightarrow \widehat{\gamma}' \cup \{\mathbf{R} \mapsto \widehat{V}_c\}} \text{ML}_c \varepsilon, \xi, \widehat{\theta}_c
\end{array}$$

where  $\xi_{\mathbf{R}} = \{\mathbf{R} \mapsto \mathbf{F}(\mathbf{R}_1) \dots (\mathbf{R}_n)\}$ . We respect the following invariants:

$$\begin{array}{ccc}
\widehat{\gamma}' \simeq^c \gamma' & \widehat{V}_s \simeq^s V_s & \widehat{\theta}_s \simeq^s \theta_s \\
FCE(\widehat{\gamma}', V_s) & \widehat{V}_c \simeq_{\widehat{\gamma}'}^c V_c & \widehat{\theta}_c \simeq_{\widehat{\gamma}'}^c \theta_c
\end{array}$$

Let us now consider the client application. Since  $\rho_c \simeq_{\widehat{\gamma}}^c \widehat{\rho}_c$ , we have that the body of the functor  $F$  is equivalent. We can thus build the following reduction:

$$\begin{array}{c}
\widehat{\rho}_c(F) = \text{functor}(\widehat{\rho}'_c)(Y_i : \mathcal{M}_i)_i \text{struct } \widehat{S}_c \text{end} \\
\widehat{\rho}_c(X_i) = \widehat{V}_i^c \quad \widehat{S}_c \xrightarrow{\widehat{\rho}_c \cup \{Y_i \mapsto \widehat{V}_i^c\}_i \mid \widehat{\gamma}' \cup \{\mathbf{R} \mapsto \widehat{V}_c\} \rightarrow \widehat{\gamma}''} \text{ML}_c \widehat{V}_c', \widehat{\theta}_c' \\
\text{module } X = F(X_1) \dots (X_n), \xi \xrightarrow{\widehat{\rho}_c \mid \widehat{\gamma}' \cup \{\mathbf{R} \mapsto \widehat{V}_c\} \rightarrow \widehat{\gamma}''} \text{ML}_c \left\{ X \mapsto \widehat{V}_c' \right\}, \xi, \widehat{\theta}_c'
\end{array}$$

By equivalence of  $F$  and  $X_i$  in  $\rho_c$  and  $\widehat{\rho}_c$ , we have that  $\widehat{\gamma}' \simeq^c \gamma'$ ,  $\widehat{V}_c' \simeq_{\widehat{\gamma}'}^c V_c'$  and  $\widehat{\theta}_c \simeq_{\widehat{\gamma}'}^c \theta_c$ .

We already built the reduction for the compiled server program. We can now build the compiled client program:

$$\begin{array}{c}
\text{exec } () , \{\mathbf{R} \mapsto \mathbf{F}(\mathbf{R}_1) \dots (\mathbf{R}_n)\} ++ \xi_{\bullet} ++ \text{end} ++ \xi \xrightarrow{\widehat{\rho}_c \mid \zeta_{\bullet} \cup \widehat{\gamma} \rightarrow \widehat{\gamma}' \cup \{\mathbf{R} \mapsto \widehat{V}_c\}} \text{ML}_c \varepsilon, \xi, \widehat{\theta}_c \\
\text{module } X = F(X_1) \dots (X_n), \xi \xrightarrow{\widehat{\rho}_c \mid \widehat{\gamma}' \cup \{\mathbf{R} \mapsto \widehat{V}_c\} \rightarrow \widehat{\gamma}''} \text{ML}_c \left\{ X \mapsto \widehat{V}_c' \right\}, \xi, \widehat{\theta}_c' \\
\langle \text{module}_m X = F(X_1) \dots (X_n) \rangle_c, \xi_{\mathbf{R}} ++ \xi_{\bullet} ++ \text{end} ++ \xi \xrightarrow{\widehat{\rho}_c \mid \zeta_{\bullet} \cup \widehat{\gamma} \rightarrow \widehat{\gamma}''} \text{ML}_c \left\{ X \mapsto \widehat{V}_c' \right\}, \xi, \widehat{\theta}_c @ \widehat{\theta}_c'
\end{array}$$

where the invariants still hold. This concludes. ■

□

#### 5.4.5 Proof of the main theorem

Finally, we prove Theorem 4. This is a direct consequence of Lemma 3.

*Proof of Theorem 4.* We have that  $P \xRightarrow{\{\}} v, \theta$ . By definition of an  $\text{ELIOM}_\varepsilon$  program execution, we can decompose this rule as following:

$$\frac{P \xRightarrow{\{\}}_m(), \mu, \theta_s \quad \mu \xRightarrow{\{\} | \varepsilon \rightarrow \gamma}_c v, \theta_c}{P \xRightarrow{\{\}} v, \theta_s @ \theta_c}$$

We trivially have the following invariants:

$$\{\} \simeq_{\{\}}^c \{\} \quad \{\} \simeq^s \{\} \quad \hat{\gamma} \simeq^c \gamma \quad FCE(\{\}, P) \quad FCE(\{\}, \{\})$$

which allow us to apply Lemma 3 and conclude.  $\square$

## 5.5 Discussion around mixed functors

In this thesis, we presented the notion of “mixed functors”, which are functors that take as argument and return mixed module composed of both client and server code. As we saw, those functors, while quite expressive, have several limitations. In this section, we will try to explore a little bit the design space around mixed functors and which limitations we think could be lifted.

First, let us recall the design constraints of ELIOM: typing and slicing are done statically and separate compilation is supported. This prevents us from “erasing” functors by inlining them and also prevents us from dynamically generating the code contained in functors. Furthermore, we want to support languages extensions such as OCAML, where functor application can depend on the control flow (notably, first class modules).

One alternative solution would be *fully separable* functors: mixed functors such that client and server execution are completely independent. While this would be easy to implement, it would also mean preventing any meaningful usage of fragments inside functors. Our version of mixed functors is slightly more expressive: the client part of the functor is indeed independent from the server part, but the server part is not. The cost is that we must do some extra book-keeping to ensure that for each server-side application of a mixed functors, all the client side effects are performed. We believe this expressive power is sufficient for most use cases. There are however several limitations to our approach, which we shall discuss point-by-point now.

**Mixed functors arguments** An important restriction of mixed functors is that their arguments can only be mixed structures. The reason for this restriction is that the `Dyn` field is used in order to keep track of client applications from the server. In the semantics presented in this thesis, the `Dyn` field is represented simply as a regular structure field, which means only structures can be passed as argument. A first step would be to allow a similar field to be added to functors. While it is a bit delicate to formalize, it should be possible to implement it in OCAML by simply adding an additional field to the closure. A second step would be to allow base, client and server modules as arguments. As with usual mixed functors, care must be taken during typechecking to not specialize eagerly. Given that this constraint is respected, one possibility would be to introduce a new module-level operation that can take a base, server or client structure and turn it into a mixed structure. Given a client module `A`, this could be done by inserting a structure of the form `(struct include%client A end)`. This could even be done transparently, since the location of modules is always known.

**Injections** Injections inside mixed functors can only refer to identifiers outside of the functor. This restriction seems particularly difficult to lift, as demonstrated by the example in Section 4.2.2. One would need to evaluate the usefulness of mixed functors in first-class contexts. It might be possible to rewrite usages of injections to escaped values in fragments. In all cases, these changes would be invasive and of limited usefulness.

**Sliceability constraints** The sliceability requirement presented in Section 5.2.1 is quite restrictive. Its goal is to ensure that `bindm` expressions are not nested and that module references are unique. This requirement can be relaxed in different manners. For mixed structures without functors, this requirements can be completely removed trivially, since internal modules are known.

For functors, there are two possibilities. A first idea would be to apply lambda-lifting to functors. By simply lifting mixed functors to the outer scope, we ensure that `bind` operations are not nested. Another possibility would be to use a similar technique to the one used for fragments inside functors: by prefixing each statically generated reference with the locally available `Dyn` field, we ensure uniqueness while allowing arbitrary nesting.

**Double application of the client side of mixed functors** As mentioned in Section 4.3.3, the client-side part of a mixed functor, when applied in a mixed context, might be called twice. While this is not so problematic in the context of applicative and pure functors, it might prove more troublesome with generative or impure functors. This problem is difficult to solve in general. In particular, it is in fairly direct conflict with the design decision, justified by Section 4.2.2, to make the client and server side of a mixed functor application independent. One potential solution would be to provide a special interpretation of mixed application in mixed contexts that would ensure that the result of the client-side functor application is properly reused. Notably, this might be doable by generating at compile time an identifier for each known returned client module.



## 6 Implementation

*Its black gates are guarded by more than just orcs. There is evil there that does not sleep. The great eye is ever watchful. It is a barren wasteland, riddled with fire, ash, and dust. The very air you breathe is a poisonous fume.*

J. R. R. Tolkien, *The Fellowship of the Ring*

The main goal of the formalization of  $\text{ELIOM}_\varepsilon$  was to inform the design and the implementation of ELIOM as a real extension of OCAML usable to develop real Web applications. As a consequence, everything presented in Chapter 4 except for mixed functors is also implemented in a fork of the OCAML compiler. This fork supports the complete OCAML language. OCAML is quite a large language and a lot of aspects of a real programming language are not covered by the formal specification.

This chapter presents how ELIOM is implemented, the various design decisions that were made related to compilation, interaction with OCAML and tooling, along with the few compromises that were occasionally made. Some sections are quite technical and assume knowledge of the inner workings of the OCAML compiler.

Everything developed in relation to this thesis is published as free software<sup>1</sup> under the GPL and LGPL licence (with the OCAML linking exception):

- <https://github.com/ocsigen/ocaml-eliom> is the fork of the compiler;
- <https://github.com/ocsigen/eliomlang> contains the ELIOM runtime and a collection of associated tools.

### 6.1 The OCaml compiler

Let us start by a quick reminder of the architecture of the OCAML compiler. The compiler pipeline is presented in Figure 6.1a. Figure 6.1b summarizes the numerous files involved in the OCAML compilation. The C files are also presented, for comparison.

Due to separate compilation and typechecking, the OCAML compiler uses many files. `.cmi` files are compiled versions of interfaces, or module types. They are used for the purpose of separate compilation: when subsequent compilations need an already compiled module, the compiler loads the file in the typing environment. Bytecode (resp. Native) compilation of a module produces a `.cmo` (resp. `.cmx`) object file. Several of these files can be grouped together in libraries as `.cma` (resp. `.cmxa`) files. Compiled objects and libraries can be linked to produce a bytecode or native executable. The compiler also

---

<sup>1</sup>I would certainly not deserve Roberto as my advisor otherwise!

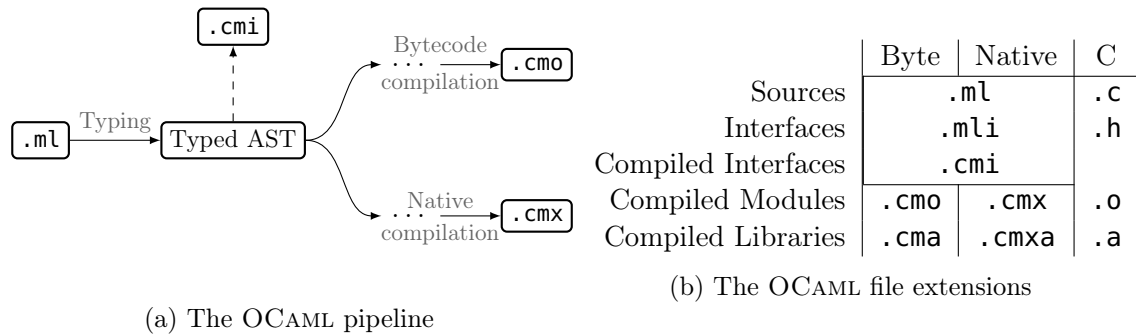


Figure 6.1: The OCAML compilation

produces several other files such as `.cmt` and `.cmti` for documentation and `.cmxs` for inlining.

The pipeline is composed of three parts: bytecode compilation, native compilation and a common part, which we simply call “frontend”, composed of parsing and typechecking. An interface source file (`.mli`) is the text format for module types. After parsing and verification of well-formedness, it is transformed into a compiled interface and produces a standalone `.cmi` file. A source file (`.ml`) is treated in several steps. Parsing produces an abstract syntax tree (or AST) with no typing information. Typechecking produces a Typed AST and a compiled interface. The typed AST is an AST where all the types have been made completely explicit: each node of the AST is annotated with its type and its local environment. The interface is the inferred module type of the file. If a corresponding module type is present (as a `.mli` file), inclusion of the two interfaces is checked. If no module type is present, a `.cmi` file is produced. After typechecking, the compiler proceeds to bytecode or native code generation. For our purpose, everything after typechecking is a black box that produces the desired compiled files.

## 6.2 The Eliom compiler

We now look at our modified ELIOM compiler. The pipeline for OCAML file is identical to the regular pipeline. The pipeline for ELIOM file is presented in Figure 6.2.

ELIOM files are either `.eliom` or `.eliomi`, which corresponds to `.ml` and `.mli` respectively. The pipeline is similar to the regular OCAML one. Each ELIOM file is first typechecked using a modified typechecker presented in Section 6.4. This creates an ELIOM typed AST, which still contains both side of the program. A unique compiled interface (`.cmi`) is created. We then slice the typed AST according to the procedure presented in Section 5.2, with returns two untyped AST, one for the client part and one for the server part. After slicing, we obtain two pure OCAML programs that can be compiled with the regular OCAML compiler which return two compiled objects `.server.cm[ox]` and `.server.cm[ox]`. The  $ML_s$  and  $ML_c$  primitives introduced in Section 5.1 are implemented in two libraries which are presented in Section 6.5.

One might wonder why the slicing is implemented in such a convoluted way: we first

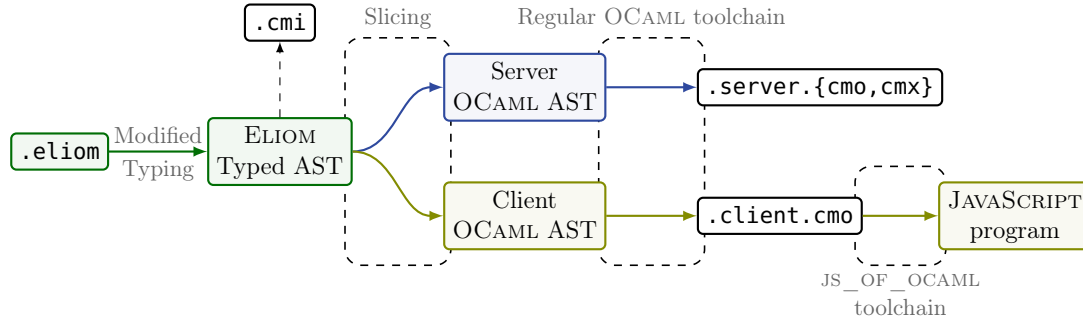


Figure 6.2: The ELIOM pipeline

type using a modified typechecker, then slices on the typed program and return two *untyped* programs which are typed again using the regular typechecker. The first reason is safety: By using the regular typechecker, we increase trust in our modified typechecker. The second reason is more pragmatic: The OCAML typed AST is a complex datastructure which is highly optimized for fast type inference and contains non trivial invariants. Doing transformations on such datastructure while preserving its invariants is a difficult task. Emission of untyped AST, on the other hand, is well supported. Thanks to the PPX ecosystem, numerous tools are available and the compiler is robust against ill-formed ASTs which makes developing such transformations easier.

### 6.2.1 Notes on OCaml compatibility

One of the important results on the ELIOM module system given Section 4.4.1 is that any vanilla OCAML module is also an  $\text{ELIOM}_\varepsilon$  module, located either on base, client or server. This is also true for the implementation. We can even go further.

Compiled interfaces created by the regular compiler can be loaded by the ELIOM compiler. A compiler flag allows to specify on which location to load the given interfaces. Compiler interfaces created by the ELIOM compiler, which contains extra location information, cannot be loaded by the OCAML compiler (as it would not be able to make sense of them). Server and Client object files are compatible with the regular compilers (and can be linked together with object files emitted by the regular compiler). The consequence is that the OCAML and ELIOM compilers can live graciously side by side and any OCAML library can be used directly in ELIOM.

The various files used by external tools (such as `.cmt` files) are also similarly compatible.

## 6.3 Converters

In the formalization, injections are modeled as the application of a converter; *i.e.*, a pair of a serialization and a deserialization function; to a value available on the server (See Section 4.2.1). Furthermore, converters are global constants introduced by the typing

environment. This formalization, however, is inconvenient to program with: it would mean the programmer needs to specify which converter to use for each injection, even the most trivial ones. We also need a way to define new converters easily. One promising idea would be to use ad-hoc polymorphism to infer which converter to use based on type information. We explore this lead in Section 6.3.1 and how it interacts with our module system. Unfortunately, ad-hoc polymorphism is not yet available in OCAML. In Section 6.3.2 we present what is currently implemented in ELIOM instead.

### 6.3.1 Modular implicits

Modular implicits [White et al., 2014] are an extension of OCAML that adds ad-hoc polymorphism in a way that is compatible with OCAML’s module system. An introduction to modular implicits can be found in White et al. [2014]. The main idea is that we can use modules themselves as implicit arguments. Converters would then be mixed modules satisfying the signature **CONV** presented in Figure 6.3. The **String** module satisfies the **CONV** signature and exports that the client and the server type are equal to the base type **string**. We can also specify converters for parametrized datatypes using functors, as presented by the **List** mixed functor. The injection operator  $\sim\%$  will then take an implicit module **C** satisfying the **CONV** signature, take an element of type **C.t** *on the server* and returns an element of type **C.t** *on the client*. Finally, we can also easily implement the **frag** converter by using a functor taking a module with just a datatype. Note that we only need a client datatype and no serialization method.

This implementation has several advantages. First, it means that injections are properly parametric: Inference will work on implicit modules in parametric functions and it plays well with polymorphism. Additionally, defining a new converter is very easy: it is simply a module respecting a fairly simple signature. Finally, the projection functions  $f^s$  and  $f^c$  used in the target languages  $ML_s$  and  $ML_c$  (Section 5.1.1) are very easy to define: given a converter module **M**, they are **M.serialize** and **M.deserialize**. Note that implementing such a signature is not completely trivial. We discuss serialization techniques in Section 6.5.2.

### 6.3.2 Wrapping

Unfortunately, at the time of writing of this thesis, modular implicits are not yet integrated in OCAML. ELIOM uses instead a method that was originally devised by Pierre Chambart and Gregoire Henry and relies on value introspection. The main idea is the following: most values can be serialized simply by using the module **Marshal** from the OCAML standard library [marshal]. As shown on multiple occasions in Chapter 2, it can be useful to transform data before serializing it. In this case, a transformation function is attached to the value. Before serializing, the function is applied to the value and the result is serialized. For customized deserialization, a similar technique is applied using tags and a client-side association table from tag to functions.

While this has the advantage of leveraging **Marshal**, which works for any data structure, is very fast and preserves sharing, it also has significant disadvantages. First, it is

```

1 module type CONV = sig
2   type%server t
3   type%client t
4   val%server serialize : t -> serial
5   val%client deserialize : serial -> t
6 end
7
8 implicit%mixed String : CONV
9   with type%server t = string
10    and type%client t = string
11
12 implicit%mixed List {M : CONV} : CONV
13   with type%server t = M.t list
14    and type%client t = M.t list
15
16 implicit%mixed Fragment {M : sig type%client t end} : CONV
17   with type%server t = M.t fragment
18    and type%client t = M.t
19
20 val%client (~%) : {C : CONV} -> C.t(*server*) -> C.t(*client*)

```

Figure 6.3: A signature for converters

not type-directed. Closures for example cannot be serialized, which will not be detected by the type system. Furthermore, defining converters is quite delicate: the programmer must create a function operating on values in a fairly untyped way and attach it to the value being converted. This is an expert task that is very error-prone.

## 6.4 Typechecking and Slicing

Most of the differences between the ELIOM compiler and the OCAML compiler are modifications of the typechecker. A formal account of the difference between the two type systems was given in Section 4.2. We now give a practical account of these differences and explain how the OCAML type checker was modified.

The OCAML language is much larger than our ML language. In particular it contains intrusive features such as first class modules and Generalized Algebraic Datatypes which affect many parts of the language. Our modified typechecker handles most of these features correctly thanks to the fact that most additional typing rules for ELIOM are at the boundaries between locations. By being conservative in the type conversions across location boundaries (in particular, by prohibiting existential and universal types), we avoid difficulties related to complex features of the OCaml type system. In other aspects, the implementation follows the formalization closely: there are now distinct namespaces for each location, with visibility rules that correspond to the relations presented in Section 4.1.1. The use of explicit locations allows to provide good error messages for missing identifiers and wrongly located declarations such as client declarations inside server modules.

One notable new feature of the implementation is the inference of some location in-

formation. In particular declarations inside base, client and server modules do not need any annotations. Similarly, most mixed annotations on modules can be elided. These elisions are based on the fact that in those cases, declaration can only have one location, which make it superfluous.

**Shared declarations** Another significant addition is the **shared** annotation. This annotation has no equivalent in the formalization. In particular, it does *not* correspond to mixed modules. Shared sections were presented in Section 2.4.4 and are implemented simply by duplicating the code between client and server, as demonstrated in Example 6.1. Note that this is a purely syntactic transformation done before typechecking. From a typechecking perspective, **y** is not a singular variable to which can be assigned a type. There are two variables, one client and one server, with two different types. This is different from **mixed** annotation where a structure contains both client and server parts. This slightly schizophrenic nature of shared declarations does not lead itself to being internalized in the type system. Nevertheless, it avoids repetitive code patterns where a similar tasks is done both client and server side with a few distinct values, which makes it very useful in practice.

<pre> 1 let%client x = 2 2 let%server x = "foo" 3 let%shared y = x 4 </pre>	<pre> 1 let%client x = 2 2 let%server x = "foo" 3 let%client y = x 4 let%server y = x </pre>
(a) Original code	(b) Desugared code

Example 6.1: A shared declaration

**Slicing** Slicing is also identical to the formalization. An example of sliced programs is given in Example 6.2. References are represented by unique generated strings.

A notably difficult point is the generation of fresh identifiers. Indeed, in our slicing formalization, we rely on generating fresh identifiers for each injection and each fragment closure. These identifiers should be globally unique in the complete program. However,

	Server	Client
let%client c = ~%s + 1	push_injection "A.s1" s	let c = get_injection "A.s1" + 1
let%server y = [%client 2 + ~%x ]	let y = fragment "B1" (x) push_fragments "B"	register_closure "B1" (fun x -> 2 + x) execute_fragments "B"

Example 6.2: Client-server code slicing

compilation units are sliced separately. The list of identifiers generated by slicing other modules is not available. To overcome this limitation, we rely on the property that, in a given OCAML program, top level module names are unique. This property is enforced by the OCAML compiler at link time. Consequently, we can generate locally unique identifiers prefixed by the current module path, which produces globally unique name in a way that is compatible with separate compilation.

### 6.4.1 Technical details

We now present a technical account of the modifications made to the compiler. This section is targeted at ambitious readers who are willing to gaze into the typechecker<sup>2</sup> in order to modify or improve the ELIOM language. We assume that the reader is familiar with the implementation of the OCAML typechecker. For a more gentle introduction to the typechecker, please consult the file `typing/HACKING.adoc` in the standard OCAML distribution.

The typechecker is both complex and frequently updated. Consequently, changes to the implementation must be made with the smallest footprint possible, in order to avoid conflicts with future modifications. One of the important early design decision for the ELIOM compiler was complete compatibility with standard OCAML tools. OCAML files format such as `.cmi` correspond to serializations of the internal data structure of the compiler (notably, the `typedtree`). Any changes to the `typedtree`, type environment or representation of type expressions would lead to changes in the compiler files which would break binary compatibility. This leads to the following programming rule: “Thou shall not change data structures”, of which we will now explore the consequences.<sup>3</sup> Some techniques used by this implementation are partially inspired by the METAOCAML patches<sup>4</sup>. METAOCAML [Kiselyov, 2014] is an extension of OCAML for meta-programming which employs quotations and slicing constructions that are similar to ELIOM’s fragments and injections. Notably, METAOCAML manages to add a quotation mechanism to the OCAML type system in a remarkably lean way (less than 1000 lines of changes on existing files).

**Locations, fragments and injections** The first task the modified typechecker need to do is to keep track of the local location of the considered code. In order to avoid passing this information around by additional arguments, we simply use a global reference. While this is not particularly elegant, there are various precedents in the OCAML typechecker (levels, notably). One additional difficulty is that the typechecker occasionally uses exception for control flow, which means we need to ensure that our global reference is put back in order. For this purpose, we use the traditional functional pattern of closure-based handler. While this approach sacrifices tail-recursion, it has not proved problematic on concrete programs.

```
1 val in_location : location -> (unit -> 'a) -> 'a
```

<sup>2</sup>No, it does not gaze back. Typecheckers do not gaze. They do however inflict 2d6 sanity damage on failed bootstraps.

<sup>3</sup>Naturally, as demonstrated by the [Oulipo](http://okmij.org/ftp/ML/MetaOCaml.html), programming constraints only encourage creativity.

<sup>4</sup><http://okmij.org/ftp/ML/MetaOCaml.html>

In order to typecheck injections and fragments, we simply leverage `in_location` to typecheck the inner expression and use unification to reflect the inferred type upward. Unfortunately, it is not possible to introduce additional nodes in the typedtree, but it is possible to add extra PPX annotations. Thus we simply return the typedtree of the inner expression with custom annotations (`[@eliom.fragment]` and `[@eliom.injection]`, notably). We then use these annotations to drive the slicing phase.

Note that `in_location` is not only used for typechecking. Indeed, locations must also be tracked when walking type expressions for validation and unification due to the presence of mixed datatypes (Section 4.2.1).

**Identifiers and bindings** In the formalization, bindings are annotated with location information. This is, unfortunately, difficult to achieve in the current implementation of the typechecker without large changes to the `Env` module. An alternative solution is used: instead of annotating binders, we annotate identifiers. Identifiers are represented by the `Ident.t` datatype, whose definition is given in Figure 6.4. The name field is the name of the identifier given in the source file. The stamp is a unique identifier. The flag field is a bitfield collecting certain flags necessary for the compiler. Only two bits are used in this integer (one to recognize top level module names, one for built in exceptions) which leaves at least 29 free bits. Fortunately, we can use this free space to store locations information! The ELIOM compilers then uses two additional bits which we will call “c” and “s” in order to determine the location of a given identifier according to the table presented Figure 6.5. The typechecker maintains a set of association tables of identifiers to values. This is implemented by `Ident.tbl`, which implements lookup by identifiers, stamp but also lookup by name. We re-purpose these lookup functions by making them respect the “can be used” relation (Figure 4.2). This can be achieved easily (although hackily) by simply accessing the current location via the global reference during lookup. Thanks to this addition to the `Ident.t` datatype and the changes in the lookup function, almost all the ELIOM location mechanisms can be implemented without touching the typechecker’s code itself.

**Specialization** Specialization is the action of taking a base type expression and bringing it to the current location. In the formalization, specialization only acts on binders, which is where location information is stored. Due to the fact that location is stored inside identifiers, the implementation is significantly more complex: we need to change the location annotations in the identifiers of the specialized expression. This is made

```

1 type t = {
2   stamp: int;
3   name: string;
4   mutable flags: int
5 }

```

Figure 6.4: The `Ident.t` datatype

		flag “c”	
		1	0
flag “s”	1	mixed	server
	0	client	base

Figure 6.5: Encoding of locations



even more delicate by the use of physical equality and mutations in the typechecker (the typechecker implements unification using a union-find algorithm embedded in the typedtree). The solution used for this purpose is to simulate copy of type expressions (as already done in the typechecker), except each identifier is copied with its flag corrected.

**Binary compatibility** Thanks to the various modifications made, the typedtree for ELIOM is compatible with the OCAML one, but with some additional information embedded in PPX attributes and identifiers flags. In order to prevent the original typechecker to misuse this additional information, we use a different magic number for cmi files.

We hope that such binary compatibility will mean that the numerous OCAML tools, notably ocp-index, ocp-browser, odoc and merlin, can eventually be adapted to present ELIOM's additional typing information.

## 6.5 Runtime and Serialization

We now present some practical consideration regarding the ELIOM runtime. There are mainly two points of interest: the implementation of the ELIOM primitives, Section 6.5.1, and the serialization format used for client-server communication, Section 6.5.2. We also discuss a small improvement to the slicing scheme in Section 6.5.3.

### 6.5.1 Primitives

The primitives used by the slicing scheme are presented formally in Section 5.1. The external signature for the primitives is shown in Figure 6.6. The main difference is that the notion of references has been further divided in `closure_id` and `inj_id` for fragment closure and injections references.

```

1 type closure_id
2 type inj_id
3
4 val%server fragment :
5   closure_id -> 'injs -> 'a fragment
6 val%server push_fragments :
7   id -> unit
8 val%server push_injection :
9   inj_id -> 'a -> unit
10
11 val%client get_injection :
12   inj_id -> 'a
13 val%client register_closure :
14   closure_id -> ('a -> 'b) -> unit
15 val%client execute_fragments :
16   id -> unit

```

Figure 6.6: API for ELIOM primitives

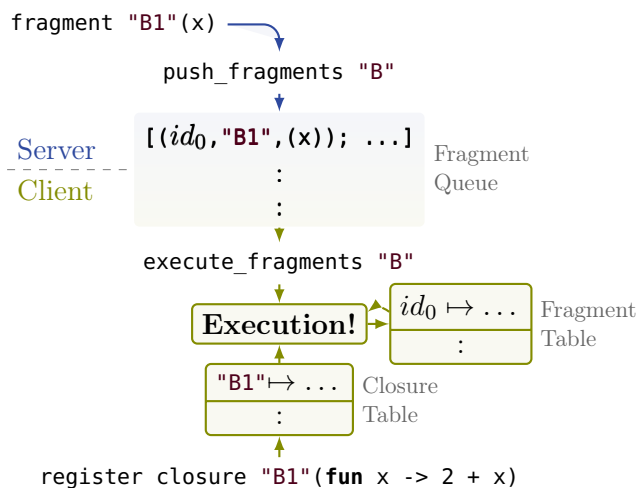


Figure 6.7: Execution of fragments

Transmitted data is represented by internal data structures which are shown in Figure 6.8. First, we use the type `poly` to (unsafely) represent arbitrary values. Similarly to closures and injections, fragments have their own kind of identifiers. Fragments on the server are represented by their identifier. The queue of fragments contains records of type `fragment`, which contains the unique identifier, the identifier of the associated closure and a tuple containing all the arguments. Injections are represented by the `injection` type, which is a pair of an identifier and the transmitted value. Finally, the transmitted information, which corresponds to the injection mapping  $\zeta$  and the fragment queue  $\xi$  are represented by `global_data` and `request_data`.

The mapping of injection is represented as an array of id/values pairs. Instead of using an explicit end token in the queue of injection, we segment the queue for each section by using an array of array. The usage is very similar. Note that we use fairly basic data structures on purpose here, as they result in more compact serialized representations.

Global data is sent for all requests while request data is local to the request currently considered. More precisely, we can consider that the execution of a web server proceeds in two steps. First, we start the application: services are launched, database is accessed and so on. Fragments and injections executed during this start-up period are considered global. Then, the request handler is installed and start listening and responding to requests. Fragments executed during a request are of course local to it. This scheme allows to avoid recomputing things (global data is kept around and can be sent along the client program when needed) and also allows to avoid replaying some side effects. For example, the ELIOM framework will reload only the minimum necessary when changing page inside a web application and thus only send new request data. Request data can only contain fragments (since injections are always global) and do not need end token (since everything is executed in one shot).

After executing the server program, we extract the resulting data, serialize it and add it to the generated HTML page. The client runtime will deserialize it, load it and use it to run the client program.

```

1  type poly (** Arbitrary values *)
2  val inject: 'a -> poly
3  val project: poly -> 'a
4
5  type fragment_id
6
7  type fragment = {
8    closure : closure_id;
9    args : poly;
10   value : fragment_id
11 }
12
13 type injection = {
14   id : inj_id;
15   value : poly;
16 }
17
18 type global_data = {
19   frags : fragment array array;
20   injs : injection array;
21 }
22
23 type request_data = fragment array

```

Figure 6.8: Internal data structures for the ELIOM runtime

### 6.5.2 Serialization format

ELIOM typing and semantics are agnostic with regard to the serialization format, which means we can choose any of the numerous serialization methods available. Serialization methods usually are a compromise between several aspects: safety, composition, size and speed.

In our case, serialized messages are transmitted on the network, so we need a very compact serialization method. Furthermore, it should also be very fast, since serialization is one of the contention points of our system. Safety properties, on the other hand, are less important: type safety of serialization and deserialization is guaranteed by the ELIOM type system. Messages could also be modified en-route. This is mitigated by two facts: HTTPS should always be used and, in case modifications are made, the worst outcome is a failure in the JAVASCRIPT client program, which is far less problematic than a crash in the server program. Finally, proper composition is not strictly needed. Indeed, Given our execution scheme for compiled program in Section 5.2, the complete arrays of injections can be serialized at once, after the execution of the server part of the program. The definition of converters, as discussed in Section 6.3, should be modular.

**Marshal** [marshal] is very fast and produces compact messages. The downside is that incompatible types at deserialization point might cause runtime errors. As detailed above, this is mitigated by ELIOM’s type system. Its distinct feature, compared to most other serialization formats, is that it preserves the sharing of OCAML values. This is extremely important to properly handle HTML across client-server boundaries (pointers to HTML elements on the server should point to the relevant DOM elements on the client). A safer and more modular alternative could be to use the “generic” library<sup>5</sup> by Balestrieri and Mauny [2016]. It leverages **Marshal** but provides improved safety checks by combining generic programming as popularized by approaches such as “Scrap Your Boilerplate” [Lämmel and Peyton Jones, 2003] with typed unmarshalling [Henry et al., 2012].

### 6.5.3 Optimized placement of sections

In Section 5.2, we present the slicing scheme for  $\text{ELIOM}_\varepsilon$  programs. In particular, we present how to place *section* annotations **exec** () and **end** (). While this scheme is correct, it is also very wasteful.

Let us consider Figure 6.9. We see that the first **exec** is not necessary: all the fragments can be executed at once at the boundary between the set of server section and the following client section. Furthermore, since  $x$  is a value, it could not even lead to a fragment execution to begin with! While this has little impact on efficiency, the high number of sections can lead to issues in some browsers<sup>6</sup>. Additionally, it can produces bigger messages than necessary by introducing many **end** tokens.

Fortunately, this can be solved in a principled way with the following remarks:

---

<sup>5</sup><https://github.com/balez/generic>

<sup>6</sup>Notably Safari. See <https://github.com/ocsigen/eliom/pull/387>.

- If a server declaration does not lead to any evaluation of client fragment, a **end/exec** pair is not needed.
- In a succession of server declarations, only one pair of **end/exec** is needed, after the last successive server declaration.

The first remark can be approximated by looking if the declaration contains something that is not a value. If that is the case, we need a section, otherwise we do not. Given that OCAML and ELIOM code mostly consist of function and type declarations, this is sufficient to eliminate more than two thirds of **end/exec** calls in medium-sized websites<sup>6</sup>. The second remark can easily be implemented by operating slicing only on the longest sequences of declarations with a common location.

1 <b>let</b> <b>%server</b> x = 1	1 execute_fragment "A"	1
2 <b>let</b> <b>%server</b> y =	2 register_closure "A1"	2 register_closure "A1"
3 <b>%client</b> ~%x + 1]	3   ( <b>fun</b> x -> x + 1)	3   ( <b>fun</b> x -> x + 1)
4	4 execute_fragment "A"	4 execute_fragment "A"
5 <b>let</b> <b>%client</b> res =	5 <b>let</b> res =	5 <b>let</b> res =
6   ~%y	6   get_injection "A.z1"	6   get_injection "A.z1"
(a) Original ELIOM code	(b) Naive client compilation	(c) Optimized version

Figure 6.9: Optimized placement of sections for a simple ELIOM program

## 7 State of the art and comparison

*Togusa: How great is the sum of thy thoughts? If I should count them, they are more in number than the sand.*

*Batou: Psalms 139, Old Testament. The way you spout these spontaneous exotic references, I'd say your own external memory's pretty twisted.*

Mamoru Oshii, *Ghost in the Shell 2: Innocence* (2004)

ELIOM takes inspiration from many sources. The two main influences are, naturally, the extremely diverse ecosystem of web programming languages and frameworks, which we explore in Section 7.1, and the long lineage of ML programming languages, which we described in Section 3.4. One of the important contributions of ELIOM is the use of a programming model similar to languages for distributed systems (Section 7.1.5) while using an execution model inspired by staged meta-programming (Section 7.2).

### 7.1 Web programming

Various directions have been explored to simplify Web development and to adapt it to current needs. ELIOM places itself in one of these directions, which is to use the same language on the server and the client. Several unified client-server languages have been proposed. They can be split in two categories depending on their usage of JAVASCRIPT. JAVASCRIPT can either be used on the server, with NODE.JS, or as a compilation target, for example with GOOGLE WEB TOOLKIT for Java or EMSCRIPTEN for C. The approach of compiling to JAVASCRIPT was also used to develop new client languages aiming to address the shortcomings of JAVASCRIPT. Some of them are new languages, such as HAXE, ELM or DART. Others are only JAVASCRIPT extensions, such as TYPESCRIPT or COFFEESCRIPT.<sup>1</sup>

However, these proposals only address the fact that JAVASCRIPT is an inadequate language for Web programming. They do not address the fact that the model of Web programming itself – server and client aspects of web applications are split in two distinct programs with untyped communication – raises usability, correctness and efficiency issues. A first attempt at tackling these concerns is to specify the communication between client and server. Such examples includes SOAP<sup>2</sup> (mostly used for RPCs) and REST<sup>3</sup> (for Web APIs). A more recent attempt is the GraphQL [GraphQL] query language which

---

<sup>1</sup>A fairly exhaustive list of languages compiling to JAVASCRIPT can be found in <https://github.com/jashkenas/coffeescript/wiki/List-of-languages-that-compile-to-JS>

<sup>2</sup><https://en.wikipedia.org/wiki/SOAP>

<sup>3</sup>[https://en.wikipedia.org/wiki/Representational\\_state\\_transfer](https://en.wikipedia.org/wiki/Representational_state_transfer)

attempts to describe, with a type system, the communications between the client and server parts of the application. These proposals are very powerful and convenient ways to check and document Web-based APIs. However, while making the contract between the client and the server more explicit, they further separate web applications into distinct tiers.

Tierless languages attempt to go in the opposite direction: by removing tiers and allowing web applications to be expressed in one single program, they make the development process easier and restore the possibility of encapsulation and abstraction without compromising correctness. In the remainder of this section, we attempt to give a fairly exhaustive taxonomy of tierless programming languages. We first give a high-level overview of the various trade-offs involved, then we give a detailed description of each language.

### 7.1.1 Code and data location

In ELIOM, code and data locations are specified through syntactic annotations. Other approaches for determining locations have been proposed. The first approach is to infer locations based on known elements through a control flow analysis (Stip.js, OPA, UR/WEB): database access is on the server, dynamic DOM interaction is done on the client, etc. Another approach is to extend the type system with locations information (LINKS, ML5). Locations can then be determined by relying on simple type inference and checking.

These various approaches present a different set of compromises:

- We believe that the semantics of a language should be easy to predict by looking at the code, which is why ELIOM uses syntactic annotations to specify locations. This fits well within the OCAML language, which is specifically designed to have a predictable behavior. On the other end of the spectrum, languages with inferred location sacrifice predictability for a very light-weight syntax which provides very little disruption over the rest of the program. Typed based approaches sit somewhere in the middle: locations are not visible in the code but are still accessible through types. Such approaches benefit greatly from IDEs allowing exploring inferred types interactively.
- Naturally, explicit approaches are usually more expressive than implicit approaches. Specifying locations manually gives programmers greater control over the performance of their applications. Furthermore, it allows to express mixed *data structures*, *i.e.*, data structures that contain both server and client parts as presented in Section 2.5. Such idioms are difficult to express when code locations are inferred. We demonstrate this with an example in the UR/WEB description.
- Type-directed approaches to infer code location is an extremely elegant approach. It can be employed either in an algebraic effect setting (LINKS) or as modal logic annotations (ML5). By its type-directed nature, error messages can be expressed in term of the source language and it should lend itself naturally to separate compilation (although this has not yet been achieved). However, such novel type systems significantly extend traditional general purpose type systems (the ML one, in this case) to

the point where it seems difficult to retrofit them on an existing languages. One lead would be to provide a tight integration using a form of Foreign Function Interface. Such integration has yet to be proposed.

- ELIOM, as a language based on OCAML, is an effectful language. Marrying inference of locations and a side-effecting semantics is delicate. The Stip.js library [Philips et al., 2014] attempts to solve this by automatically providing replication and eventual consistency on shared references. This might cause many more communications that necessary if not done carefully. We believe that such decision is better left in the hands of the programmer.

### 7.1.2 Slicing

Once code location has been determined, the tierless program must be sliced in (at least) two components. In ELIOM, slicing is done statically at compile time in a modular manner: each module is sliced independently. Another common approach is to use a static whole-program slicing transformation (UR/WEB, Stip.js). This is most common for languages where code location is inferred, simply due to the fact that such inference is often non-modular. This allows precise analysis of location that can benefit from useful code transformations such as CPS transformation [Philips et al., 2016], inlining and defunctionalization. However, this can make it difficult for the users to know where each piece of code is executed and hinder error messages,. It also prevents any form of separate compilation.

Finally, slicing can be done at runtime simply by generating client JAVASCRIPT code “on the fly” during server execution (LINKS, HOP, PHP). Such solution has several advantages: it is easier to implement and provides a very flexible programming style by allowing programmers to compose the client program in arbitrary ways. The downside is that it provides less guarantees to the users. Furthermore, it prevents generating and optimizing a single JAVASCRIPT file in advance, which is beneficial for caching and execution purposes.

**Separate and incremental compilation** Most current mainstream compiled language support some form of incremental compilation. Indeed, incremental compilation avoids recompiling files of which no dependency has changed. This accelerates the feedback loop between development and testing greatly and allow very fast recompilation times. In the case of statically typed languages, it also allows immediate checking of the modified file thus providing developers very fast iteration cycles. The easiest way to implement incremental compilation is through separate compilation, where each file can be compiled completely independently. Furthermore, separate compilation is compatible with link-time optimization and thus does not prevent generation of heavily optimized code, as demonstrated by nearly every C compiler. As a consequence, we consider languages that do not support incremental compilation completely unusable for practical usages.

### 7.1.3 Communications

ELIOM uses asymmetric communication between client and server (see Section 2.2.2): everything needed to execute the client code is sent during the initial communication that also sends the Web page. It also exposes a convenient API for symmetric communications using RPC (Section 2.3.1) and broadcasts (Section 2.3.3), which must be used manually.

We thus distinguish several kind of communications. First, manual communications are exposed through normal APIs and are performed explicitly by programmers. Of course, the convenience and safety of such functions vary a lot depending on the framework. Then, we consider automatic communications, that are inserted automatically by the language at appropriate points, as determined by code locations and slicing. We can further decompose automatic communications further in two categories. In static asymmetric communications, information is sent from the server to the client automatically, when sending the page. In dynamic symmetric communications, information is sent back and forth between the client and the server dynamically through some form of channel (AJAX, websockets, Comet, ...).

While symmetric communications are very expressive, they impose a significant efficiency overhead: a permanent connection must be established and each communication imposes a round trip between client and server. Furthermore, such communication channel must be very reliable. On the other hand, asymmetric communications are virtually free: data is sent with the web page (and is usually much smaller). Only a thin instrumentation is needed. Of course, the various communication methods can be mixed in arbitrary manner. ELIOM, for example, uses both automatic asymmetric and manual communications.

**Offline usage** Many web applications are also used on Mobile phones, where connection is intermittent at best. As such, we must consider the case where the web application produced by a tierless language is used offline. In this context, asymmetric communication offer a significant advantage: given the initially transmitted information by the server, the client program can run perfectly fine without connection. This guarantee, however, does not extend to dynamic manual communications done by the use of RPCs and channels. Philips et al. [2014] explore this question for symmetric communications through their R5 requirement.

### 7.1.4 Type systems

Type safety in the context of tierless languages can encompass several notions. The first notion is the traditional distinction between weakly and strongly typed languages. In the interest of avoiding a troll war among the jury, we will not comment further. A more interesting question is whether communication errors between client and server are caught by the typechecker. This is, surprisingly, not the case of UR/WEB since location inference and slicing is done very late in the compilation process, far after type checking. One consequence of this is that slicing errors are fairly difficult to understand [Chlipala,



2015a, page 10].<sup>4</sup> While the ELIOM formalization is type safe, the ELIOM implementation is not, due to the use of wrapping and Marshall (Section 6.3.2), which will fail at runtime on functional values.

Another remark is the distinction between client and server universes.<sup>5</sup> ELIOM has separate type universes for client and server types (see Section 4.2.1). Most tierless languages do not provide such distinction, notably for the purpose of convenience. Distributed systems such as ACUTE, however, do make such distinction to provide a solution for API versionning and dynamic reloading of code. In this case, there are numerous distinct type universes.

**Module systems** The notion of module system varies significantly depending on the language. In ELIOM we consider an ML-style module system composed of a small typed language with structures and functors. We believe modules are essential for building medium to large sized programs: this has been demonstrated for general purpose languages but also holds for web programming languages, as demonstrated by the size of large modern websites (the web frontend of facebook alone is over 9 *millions* lines of code). Even JAVASCRIPT recently obtained a module system in ES6. In the context of tierless languages, an interesting question is the interaction between locations and modules. In particular, can modules contain elements of different locations and, for statically typed languages, are locations reflected in signatures?

**Types and documentation** Type systems are indisputably very useful for correctness purposes, but they also serve significant documentation purposes. Indeed, given a function, its type signature provides many properties. In traditional languages, this can range from very loose (arguments and return types) to very precise (with dependent types and parametricity [Wadler, 1989]). In the context of tierless languages, important questions we might want to consider are “Where can I call this function?” and “Where should this argument come from?”. The various languages exposes this information in different ways: ELIOM does not expose location in the types, but it is present in the signature. ML5 exposes this information directly in the types. UR/WEB and LINKS do not expose that information at all.

### 7.1.5 Details on some specific approaches

We now provide an in-depth comparison with the most relevant approaches. A summary in Figure 7.1 classifies each approach according to the main distinctive features described in the previous paragraphs. Each language or framework is also described below.

UR/WEB [Chlipala, 2015a,b] is a new statically typed language especially designed for Web programming. It features a rich ML-like type and module system and a fairly

---

<sup>4</sup> “However, the approach we adopted instead, with ad-hoc static analysis on whole programs at compile time, leads to error messages that confuse even experienced Ur/Web programmers.”

<sup>5</sup> Or, more philosophically: Is your favorite language platonist or nominalist ?

	Locations	Slicing	Communications	Type safe	Host language
ELIOM	Syntactic	Modular	Asymmetric	✓	OCAML
LINKS	Type-based*	Dynamic*	Symmetric	✓	-
UR/WEB	Inferred	Global	(A)symmetric~	✓*	-
Haste	Type-based	Modular	Symmetric	✓	HASKELL
HOP	Syntactic	Dynamic*	(A)symmetric~	×	JAVASCRIPT*
Meteor.js	Syntactic	Dynamic	Manual	×	JAVASCRIPT
Stip.js	Inferred	Global	Symmetric*	×	JAVASCRIPT
ML5	Type-based	Global*	Symmetric	✓	-
Acute	Syntactic	Modular	Distributed	✓	OCAML

Figure 7.1: Summary of the various tierless languages

See previous sections for a description of each headline. A star \* indicates that details are available in the description of the associated language. A tilde ~ indicates that we are unsure, either because the information was not specified, or because we simply missed it.

original execution model where programs only execute as part of a web-server request and do not have any state (the language is completely pure). While similar in scope to ELIOM, it follows a very different approach: Location inference and slicing are done through a whole-program transformation operated on a fairly low level representation. Notably, this transformation relies on inlining and removal of high-order functions (which are not supported by the runtime). The advantages of this approach are twofold: It makes UR/WEB applications extremely fast (in particular because it doesn't use a GC: memory is trashed after each request) and it requires very little syntactic overheads, allowing programs to be written in a very elegant manner.

The downsides, however, are fairly significant. UR/WEB's approach is incompatible with any form of separate compilation. Many constructs are hard-coded into the language, such as RPCs and reactive signals and it does not seem possible to implement them as libraries. The language is clearly not general and has a limited expressivity, in particular when trying to use mixed data-structures (see Section 2.5). For example, Example 7.1. presents the server function `button_list` which takes a list of labels and client functions and generates a list of buttons. We show the ELIOM implementation and a tentative UR/WEB implementation. The UR/WEB version typechecks but slicing fails. We are unable to write a working version and do not believe it to be possible: indeed, in the ELIOM version we use a client fragment to build the list `l` as a mixed data-structure. This annotation is essential to make the desired semantics explicit. Other examples, such as the accordion widget (Section 2.7) are expressible only using reactive signals, which present a very different semantics.

HOP [Serrano et al., 2006] is a dialect of Scheme for programming Web applications. Its successor, HOP.js [Serrano and Prunet, 2016], takes the same concepts and brings them to JAVASCRIPT. The implementation of HOP.js is very complete and allow them to run both the JAVASCRIPT and the scheme dialect while leveraging the complete node.js ecosystem. HOP uses very similar language constructions to the one provided by ELIOM:

```

1 let%client handler _ = alert "clicked!"
2 let%server l =
3   [ ("Click!", [%client handler]) ]
4
5 let%server button_list lst =
6   ul (List.map (fun (name, action) ->
7     li [button
8       ~button_type:'Button
9       ~a:[a_onclick action]
10      [pdata name]])
11   lst)
12
13 let main () =
14   body (button_list l)

```

(a) ELIOM version

```

1 fun main () : transaction page =
2   let
3     fun handler _ = alert "clicked!"
4     val l = Cons (("Click!", handler), Nil)
5
6     fun button_list lst =
7       case lst of
8         Nil => <xml/>
9         | Cons ((name, action), r) =>
10          <xml>
11            <xml><button value={name}
12              onclick={action}/>
13            {buttons r}
14          </xml>
15   in
16     return <xml>
17       <body>{button_list l}</body>
18     </xml>
19   end

```

(b) Tentative UR/WEB version. Typechecks but does not compile.

Example 7.1: Programs building a list of buttons from a list of client side actions

~-expressions are fragments and \$-expressions are injections. All functions seem to be **shared** by default. Communications are asymmetric when possible and use channels otherwise. However, contrary to ELIOM, slicing is done dynamically during server execution [Loitsch and Serrano, 2007]. In the tradition of Scheme, HOP only uses a minimal type system for optimizations and does not have a notion of location. In particular HOP does not provide static type checking and does not statically enforce the separation of client and server universes (such as preventing the use of database code inside the client). The semantics of HOP has been formalized [Serrano and Queinnec, 2010, Boudol et al., 2012] and does present similarities to the interpreted ELIOM semantics (Section 4.3). HOP is however significantly more dynamic than ELIOM: it allows dynamic communication patterns through the use of channels and allows nested fragments in the style of Lisp quotations which allows to generate client code inside client code.

For dynamically-typed inclined programmers, HOP currently presents the most convincing approach to tierless Web programming. In particular given its solid implementation, great flexibility and support for the JAVASCRIPT ecosystem.

LINKS [Cooper et al., 2006] is an experimental functional language for client-server Web programming with a syntax close to JAVASCRIPT and an ML-like type system. Its type system is extended with a notion of *effects*, allowing a clean integration of database queries in the language [Lindley and Cheney, 2012]. In Example 7.2, we highlight two notable points of LINKS: the function **adults** takes as argument a list *l* and returns the name of all the person over 18. This function has no effect and can thus run on the

client, the server, but can also be transformed into SQL to run in a database query. On the other hand, the `print` function has an effect called “wild” which indicates it can’t be run inside a query. Effects are also used to provide type-safe channel-based concurrency.

LINKS also allows to annotate functions by indicating on which location they should run. Those annotations, however, are not reflected in the type system. Communications are symmetric and completely dynamic through the use of AJAX. Client-server slicing is dynamic (although some progress has been made towards static *query* slicing [Cheney et al., 2014]) and can introduce “code motion”, which can move closures from the server to the client. This can be extremely problematic in practice, both from an efficiency and a security point of view. The current implementation of LINKS is interpreted but a compilation scheme leveraging the Multicore-OCAML efforts has been recently added.

Although LINKS is very seducing, the current implementation presents many shortcomings given its statically typed nature: slicing is dynamic and produces fairly large JAVASCRIPT code and the type system does not really track client-server locations.

```
1 links> fun adults(l) { for (x <- l) where (x.age >= 18) [(name = x.name)] } ;;
2 adults = fun : ((age:Int,name:a[_]) -> [(name:a)])
3
4 links> print ;;
5 print : (String) {wild}-> ()
```

#### Example 7.2: Small pieces of LINKS code

METEOR.JS [Meteor.js] is a framework where both the client and the server sides of an application are written in JAVASCRIPT. It has no built-in mechanism for sections and fragments but relies on conditional `if` statements on the `Meteor.isClient` and `Meteor.isServer` constants. It does not perform any slicing. This means that there are no static guarantees over the respective execution of server and client code. Besides, it provides no facilities for client-server communication such as fragments and injections. Compared to ELIOM, this solution only provides coarse-grained composition.

STIP.JS [Philips et al., 2014] allows to slice tierless JAVASCRIPT programs with a minimal amount of annotations. It emits METEOR.JS programs with explicit communications. Annotations are optionally provided through the use of comments, which means that STIP.JS are actually perfectly valid JAVASCRIPT programs. Location inference and slicing are whole-program static transformations. Communications are symmetric, through the use of fairly elaborate consistency and replication mechanisms for shared references. This approach allows the programmer to write code with very little annotations. As opposed to UR/WEB, manual annotations are possible, which might allow to express delicate patterns such as mixed data-structures and prevents security issues.

## Distributed programming

Tierless languages in general are very inspired by distributed programming languages. The main difference being that distributed programs contain an arbitrary number of locations while tierless web programs only have two: client and server. Communications are generally symmetric and dynamic, due to the multi-headed aspect of distributed systems. There are of course numerous programming languages dedicated to distributed programming. We present here two relevant approaches that put greater emphasis on the typing and tierless aspects.

[Ekblad \[2017\]](#) proposes an EDSL of HASKELL for distributed programming. This DSL allows to express complex orchestrations of multiple nodes and external components (for example databases and IoT components), with handling of distinct type universes when necessary. Instead of using syntactic annotations, locations are determined through typing. This approach works particularly well in the context of HASKELL, thanks to the advanced type system and the syntactic support for monads and functors. Multiple binaries are produced from one program. Slicing relies on type information and dead code elimination, as provided by the GHC compiler. Explicit slicing markers similar to ELIOM’s section annotations are the subject of future work. Communications are dynamic and symmetric through the use of websockets. One notable feature of this DSL is that it offers a client-centric view: The control flow is decided by the client which pilots the other nodes. This is the opposite of ELIOM where the server can assemble pieces of client code through fragments. This work also inherits the HASKELL and GHC features in term of modules, data abstraction and separate-compilation. A module language has been developed for HASKELL by [Kilpatrick et al. \[2014\]](#).

An earlier version, HASTE.APP [\[Ekblad and Claessen, 2014\]](#), was limited to only one client and one server and used a monadic approach to structure tierless programs.

ML5 [\[Murphy VII et al., 2007\]](#) is an ML language that introduces new constructs for type-safe communication between distributed actors through the use of location annotations inside the types called “modal types”. It is geared towards a situation where all actors have similar capabilities. It uses dynamic communication, which makes the execution model very different from ELIOM. ML5 provides a very rich type system that allows to precisely export the capabilities of the various locations. For example, it is possible to talk about addresses on distant locations and pass them around arbitrary. ELIOM only supports such feature through the use of fragments, for client code.

Unfortunately, ML5 does not have a module system. However, we believe that ML5’s modal types can play a role similar to ELIOM’s location annotations on declarations, including location polymorphism. ML5 uses a global transformation for slicing. Given the rich typing information present in ML5’s types, it should lend itself fairly well to a modular slicing approach, but this has not been done.

ACUTE [\[Sewell et al., 2007\]](#) is an extension of OCAML for distributed programming. It provides typesafe serialization and deserialization and also allows arbitrary loading of

modules at runtime. Like ELIOM, it provides a full-blown module system. However, it takes an opposite stance on the execution model: each actor runs independent programs and communications are completely dynamic.

Handling of multiple type universes is done by providing a description of the type with each message and by versioning APIs. In particular, great care is taken to provide type safe serialization by also transmitting the type of messages alongside each message. This gives ACUTE very interesting capabilities, such as reloading only part of the distributed system in a type-safe way.

## 7.2 Staged meta-programming

An important insight regarding ELIOM is that, while it is a tierless programming language and tries to disguise itself as a distributed programming language, ELIOM corresponds exactly to a staged meta-programming language. ELIOM simply provides only two stages: stage 0 is the server, stage 1 is the client. ELIOM's client fragments are the equivalent of stage quotations.

Most approaches to partial evaluation are done implicitly (not unlike tierless languages with implicit locations). We take inspiration from several approaches that combine staged meta-programming with explicit stages annotations that are reflected in the type system, which we describe here.

METAOOCAML [Kiselyov, 2014] is an extension of OCAML for meta programming. It introduces a quotation annotation for staged expressions, whose execution is delayed. Quotations and antiquotations corresponds exactly to fragments and injections. The main difference is that METAOOCAML is much more dynamic: quoted code does not have to be completely closed when produced and well-scopedness is checked dynamically, just before running the quoted code. This allows very dynamic behaviors such as automatic insertion of let-bindings [Kiselyov, 2015] and dynamically determining staged stream pipelines [Kiselyov et al., 2017]. One difference is the choice of universes: ELIOM has two universes, client and server, which are distinct. METAOOCAML has a single type universe but a series of scopes, for each stage, included in one another.

METAOOCAML itself provides no support for modules and only leverages the OCAML module system. Staging annotations are only on expressions, not on declarations.

Modular macros [Nicole, 2016, Yallop and White, 2015] are another extension of OCAML. It uses staging to implement macros. It provides both a quotation-based expression language along with staging annotations on declarations. It also aims to support modules and functors. The slicing can be seen as dynamic (since code is executed at compile time to produce pieces of programs). In particular, this allows to lift most of the restriction imposed on multi-stage functors. They also use a notion similar to converters, except that the serialization format here is simply the OCAML AST.

The main difference compared to ELIOM is how the asymmetry between stage 0 and stage 1 is treated. Only one type universe is used and there is no notion of slicing that

would allow a distant execution.

[Feltman et al. \[2016\]](#) presents a slicing technique for a two-staged simply typed lambda calculus. Their technique is similar to the one used in ELIOM. They distinguish their language in three parts:  $1\mathcal{G}$ , which corresponds to base code;  $1\mathcal{M}$ , which corresponds to server code; and  $2\mathcal{M}$ , which corresponds to client code. They also provide a proof of equivalence between the dynamic semantics and the slicing techniques. This proof has been mechanized in Twelf. While their work is done in a more general setting, they do not specify how to transfer rich data types across stages (which is solved in ELIOM using converters). They also do not propose a module system.





## 8 Conclusion

*To the designer of programming languages, I say: unless you can support the paradigms I use when I program, or at least support my extending your language into one that does support my programming methods, I don't need your shiny new languages; [..] To persuade me of the merit of your language, you must show me how to construct programs in it.*

Robert W. Floyd, *The paradigms of programming*

In this thesis, I presented the ELIOM language, its design, formalization and implementation. Through ELIOM, I also presented several related notions, such as ML languages, tierless web programming and staged meta-programming. At its core, ELIOM combines the various insights made by staged meta-programming languages and the very powerful OCAML language in order to provide a safe and efficient programming language that allows to write client-server tierless programs in a convenient way.

One might note that, for most of this thesis, we do not talk about the Web all that much. Indeed, while OCSIGEN is a Web programming framework, the language constructs we presented are not specific to Web programming. In fact, the minimal runtime developed during this thesis only needs the OCAML standard library to run, and can be compiled to any targets! Given the static nature of the ELIOM programming language, the server and the client part are separated and compiled statically, so ELIOM could be used to write more general client-server applications. Similarly, we believe that several techniques developed in the context of this thesis are of more general use. Notably, converters are a generic solution to make cross-stage persistence manipulable in a first class manner by programmers. From this perspective, ELIOM can be seen as a general approach for type-safe client-server communications. Indeed ELIOM does not force a specific programming style on the user. The additional primitives can be used for traditional Web programming techniques, as well as Model-View-Controller architectures or a more modern Functional Reactive approach. It simply turns out that combining those primitives with a programming language that supports modularity and encapsulation such as OCAML yields a powerful Web programming language.

Another focus of this thesis is to make ELIOM usable in practice. Chapter 2 was dedicated to presenting the ELIOM programming language from a user's perspective, with numerous examples demonstrating its usefulness for building Web applications and libraries. Chapter 6 presents how ELIOM can be implemented as an actual extension of the OCAML compiler. In general, many design choices of ELIOM were inspired by practical concerns. Indeed, in order to be useful, a language must have an ecosystem. The simplest way to have an ecosystem is to reuse the one of an existing language, hence the development of ELIOM as an extension of OCAML. Since incremental compilation is

indispensable for any non-trivial programming projects, our design must be compatible with it. Given the existence of the OCaml module system, we needed good interaction between modularity, abstraction and the tierless annotations. All this led us to develop the various features presented in this thesis such as fragments, converters and location annotations inside modules.

Of course, there is still significant work to be done on the ELIOM language. From a practical point of view, the new implementation needs to be improved and tested by more users. On the more theoretical aspects, mixed functors are still very much work in progress, as noted in Section 5.5. Instead of dwelling on these unfinished tasks, I would like, after three and a half years working on the subject, to offer my personal vision for the perfect tierless statically typed Web programming language.

**My ideal tierless programming language** ELIOM is, as said above, very practically-minded. By extending OCAML, it leverages its many strengths but also inherits its design choices and weaknesses. As such, it represents a local optimum in the design space of tierless Web programming languages. I believe we can do better, and I will attempt here to give my ideas on how.

Before considering the tierless aspects, let me settle the kind of language I will consider here. As should be apparent by now, I strongly prefer statically-typed strict-by-default languages. Dynamic languages do have numerous advantages, for example in the context of distributed systems like Erlang or in very dynamic settings where the Scheme family excels. Indeed, HOP already provides a very convincing solution for dynamically-typed tierless Web programming. My personal taste<sup>1</sup>, however, goes to static typing, both for the numerous benefits it provides in term of soundness and convenience, but also because the discipline it offers help me organize my own thoughts on how and what to do while exploring the complicated maze of design decisions that is programming.

Given this context, I believe that the modularity and encapsulation properties of a strong module system with abstract datatypes are one of the best programming tools available in modern programming languages. Modules give the programmer the ability to manipulate structured programs directly inside the language and allows to enforce invariants inside module boundaries, as demonstrated on multiple occasions in Chapter 2. These capabilities are essential for any kind of modular programming, including Web programming. A module system providing these capabilities can take many forms: the “classical” ML module system, the more practical OCAML one; more experimental module systems such as 1ML [Rossberg et al., 2014] or even something quite different, such as Backpack [Kilpatrick et al., 2014]. All these module systems provide the ability to hide internal details through abstraction and to manipulate modules in very powerful ways, thus providing the necessary tools for modular large scale programming.

On the tierless aspects, syntactic annotations are essential to provide finer control, as argued in several occasions in this thesis. However, I believe a type-and-effect system in the style of Eff [Bauer and Pretnar, 2012] could provide a tighter integration between

---

<sup>1</sup>For the willing, I can propose other lively debates: Vim vs. Emacs; Zelazny vs. Sanderson, K. Dick vs. Asimov; Tomatoes vs. Potatoes.

the type system and the tierless annotations. Indeed, we would have two effects: “client” and “server”. Fragments would only contain code that has no “server” effect, but would produce a server effect themselves. This could allow to elide some of the most obvious annotations, hence allowing a programming style similar to Ur/Web’s implicitness while still providing the necessary control when needed. Furthermore, location information would be exposed at the module level directly through the effects systems, thus replacing the need for annotations in module types. Effects would then play a similar role than modalities in ML5. Links already demonstrated that such an effect-based language can be compiled down to quotations [Cheney et al., 2014], which would allow to keep the efficient execution model of ELIOM.

Finally, while there is a rather large tooling ecosystem surrounding Web languages, the same cannot be said about tierless languages (except maybe Hop.js). In particular, a REPL for a statically sliced tierless language would both be technically challenging and provide a more exploratory style of programming.

This design can be decomposed in a combination of powerful general-purpose features (here, modularity, abstraction and effect systems) and a few selected specific-purpose primitives (our two new effects). This combination allows to build up new idioms for a chosen domain and provides good integration into existing languages. By exposing orthogonal features that interact well instead of baking complex constructions in the language, we can provide tools that are more flexible and easier to extend by programmers. I believe this is a nice way of designing programming languages, and I hope to pursue it in other domains in the future.



# Bibliography

- N. Amin and T. Rompf. Type soundness proofs with definitional interpreters. In G. Castagna and A. D. Gordon, editors, *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*, pages 666–679. ACM, 2017. ISBN 978-1-4503-4660-3. doi: 10.1145/3009837. URL <http://dl.acm.org/citation.cfm?id=3009866>.
- B. E. Aydemir, A. Charguéraud, B. C. Pierce, R. Pollack, and S. Weirich. Engineering formal metatheory. In G. C. Necula and P. Wadler, editors, *Proceeding of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, pages 3–15. ACM, 2008. URL <http://arthur.chargueraud.org/research/2007/binders/>.
- V. Balat. Rethinking Traditional Web Interaction. In *Proceedings of the Eighth International Conference on Internet and Web Applications and Services*, Rome, Italy, 2013.
- V. Balat. Rethinking traditional web interaction: Theory and implementation. *International Journal on Advances in Internet Technology*, 2014. ISSN 1942-2652. URL [http://www.iariajournals.org/internet\\_technology/](http://www.iariajournals.org/internet_technology/).
- V. Balat, J. Vouillon, and B. Yakobowski. Experience report: Ocsigen, a Web programming framework. In G. Hutton and A. P. Tolmach, editors, *ICFP*, pages 311–316. ACM, 2009. ISBN 978-1-60558-332-7.
- F. Balestrieri and M. Mauny. Generic programming in OCaml. *OCaml Workshop*, 2016.
- A. Bauer and M. Pretnar. Programming with algebraic effects and handlers. *CoRR*, abs/1203.1539, 2012. URL <http://arxiv.org/abs/1203.1539>.
- A. Bawden. Quasiquotation in lisp. In *PEPM*, pages 4–12. University of Aarhus, 1999.
- G. Boudol, Z. Luo, T. Rezk, and M. Serrano. Reasoning about Web applications: An operational semantics for HOP. *ACM Trans. Program. Lang. Syst.*, 34(2):10, 2012.
- J. Cheney, S. Lindley, G. Radanne, and P. Wadler. Effective quotation: relating approaches to language-integrated query. In W. Chin and J. Hage, editors, *Proceedings of the ACM SIGPLAN 2014 workshop on Partial evaluation and program manipulation, PEPM 2014, January 20-21, 2014, San Diego, California, USA*, pages 15–26. ACM, 2014. doi: 10.1145/2543728.2543738. URL <http://doi.acm.org/10.1145/2543728.2543738>.

- A. Chlipala. Ur/Web: A simple model for programming the Web. In *Proceedings of the 42Nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '15, pages 153–165, New York, NY, USA, 2015a. ACM. ISBN 978-1-4503-3300-9. doi: 10.1145/2676726.2677004. URL <http://doi.acm.org/10.1145/2676726.2677004>.
- A. Chlipala. An optimizing compiler for a purely functional Web-application language. In *ICFP*, 2015b.
- E. Cooper, S. Lindley, P. Wadler, and J. Yallop. Links: Web programming without tiers. In *FMCO*, pages 266–296, 2006.
- K. Crary. Modules, abstraction, and parametric polymorphism. In G. Castagna and A. D. Gordon, editors, *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*, pages 100–113. ACM, 2017. URL <http://dl.acm.org/citation.cfm?id=3009892>.
- D. Dreyer. *Understanding and Evolving the ML Module System*. PhD thesis, CMU, May 2005. URL <https://people.mpi-sws.org/~dreyer/thesis/main.pdf>.
- A. Ekblad. A meta-edsl for distributed web applications. In I. S. Diatchki, editor, *Proceedings of the 10th ACM SIGPLAN International Symposium on Haskell, Oxford, United Kingdom, September 7-8, 2017*, pages 75–85. ACM, 2017. doi: 10.1145/3122955.3122969. URL <http://doi.acm.org/10.1145/3122955.3122969>.
- A. Ekblad and K. Claessen. A seamless, client-centric programming model for type safe web applications. In *Proceedings of the 2014 ACM SIGPLAN Symposium on Haskell, Haskell '14*, pages 79–89, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-3041-1. doi: 10.1145/2633357.2633367. URL <http://doi.acm.org/10.1145/2633357.2633367>.
- Eliom. *Eliom web site*. <https://ocsigen.org/eliom>, 2017.
- Immutable. *Immutable*. Facebook, <https://facebook.github.io/immutable-js/>, 2017.
- Reactjs. *Reactjs*. Facebook, <https://reactjs.org/>, 2017.
- N. Feltman, C. Angiuli, U. A. Acar, and K. Fatahalian. Automatically splitting a two-stage lambda calculus. In Thiemann [2016], pages 255–281. ISBN 978-3-662-49497-4. doi: 10.1007/978-3-662-49498-1\_11. URL [http://dx.doi.org/10.1007/978-3-662-49498-1\\_11](http://dx.doi.org/10.1007/978-3-662-49498-1_11).
- J. Garrigue. A Certified Interpreter for ML with Structural Polymorphism, 2009.
- GraphQL. GraphQL, 2016. URL <http://graphql.org/>.

- G. Henry, M. Mauny, E. Chailloux, and P. Manoury. Typing unmarshalling without marshalling types. In P. Thiemann and R. B. Findler, editors, *ACM SIGPLAN International Conference on Functional Programming, ICFP'12, Copenhagen, Denmark, September 9-15, 2012*, pages 287–298. ACM, 2012. ISBN 978-1-4503-1054-3. doi: 10.1145/2364527.2364569. URL <http://doi.acm.org/10.1145/2364527.2364569>.
- J. Hughes. Why functional programming matters. *Comput. J.*, 32(2):98–107, 1989. doi: 10.1093/comjnl/32.2.98. URL <https://doi.org/10.1093/comjnl/32.2.98>.
- S. Kilpatrick, D. Dreyer, S. L. Peyton Jones, and S. Marlow. Backpack: retrofitting haskell with interfaces. In S. Jagannathan and P. Sewell, editors, *The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '14, San Diego, CA, USA, January 20-21, 2014*, pages 19–32. ACM, 2014. doi: 10.1145/2535838.2535884. URL <http://doi.acm.org/10.1145/2535838.2535884>.
- O. Kiselyov. The design and implementation of BER metaocaml - system description. In M. Codish and E. Sumii, editors, *Functional and Logic Programming - 12th International Symposium, FLOPS 2014, Kanazawa, Japan, June 4-6, 2014. Proceedings*, volume 8475 of *Lecture Notes in Computer Science*, pages 86–102. Springer, 2014. ISBN 978-3-319-07150-3. doi: 10.1007/978-3-319-07151-0\_6. URL [http://dx.doi.org/10.1007/978-3-319-07151-0\\_6](http://dx.doi.org/10.1007/978-3-319-07151-0_6).
- O. Kiselyov. Generating code with polymorphic let: A ballad of value restriction, copying and sharing. In J. Yallop and D. Doligez, editors, *Proceedings ML Family / OCaml Users and Developers workshops, ML Family/OCaml 2015, Vancouver, Canada, 3rd & 4th September 2015.*, volume 241 of *EPTCS*, pages 1–22, 2015. doi: 10.4204/EPTCS.241.1. URL <https://doi.org/10.4204/EPTCS.241.1>.
- O. Kiselyov, A. Biboudis, N. Palladinos, and Y. Smaragdakis. Stream fusion, to completeness. In G. Castagna and A. D. Gordon, editors, *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*, pages 285–299. ACM, 2017. URL <http://dl.acm.org/citation.cfm?id=3009880>.
- R. Lämmel and S. L. Peyton Jones. Scrap your boilerplate: a practical design pattern for generic programming. In Z. Shao and P. Lee, editors, *Proceedings of TLDI'03: 2003 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation, New Orleans, Louisiana, USA, January 18, 2003*, pages 26–37. ACM, 2003. ISBN 1-58113-649-8. doi: 10.1145/604174.604179. URL <http://doi.acm.org/10.1145/604174.604179>.
- D. K. Lee, K. Crary, and R. Harper. Towards a mechanized metatheory of standard ML. In M. Hofmann and M. Felleisen, editors, *Proceedings of the 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2007, Nice, France, January 17-19, 2007*, pages 173–184. ACM, 2007. ISBN 1-59593-575-4. doi: 10.1145/1190216.1190245. URL <http://doi.acm.org/10.1145/1190216.1190245>.

- X. Leroy. Manifest types, modules, and separate compilation. In H. Boehm, B. Lang, and D. M. Yellin, editors, *Conference Record of POPL'94: 21st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Portland, Oregon, USA, January 17-21, 1994*, pages 109–122. ACM Press, 1994. doi: 10.1145/174675.176926. URL <http://doi.acm.org/10.1145/174675.176926>.
- X. Leroy. Applicative functors and fully transparent higher-order modules. In R. K. Cytron and P. Lee, editors, *Conference Record of POPL'95: 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, San Francisco, California, USA, January 23-25, 1995*, pages 142–153. ACM Press, 1995. doi: 10.1145/199448.199476. URL <http://doi.acm.org/10.1145/199448.199476>.
- X. Leroy. A syntactic theory of type generativity and sharing. *J. Funct. Program.*, 6(5):667–698, 1996. doi: 10.1017/S0956796800001933. URL <http://dx.doi.org/10.1017/S0956796800001933>.
- X. Leroy, D. Doligez, A. Frisch, J. Garrigue, D. Rémy, and J. Vouillon. *The OCaml system release 4.04, Documentation and user's manual*. Projet Gallium, INRIA, Nov. 2016.
- S. Lindley and J. Cheney. Row-based effect types for database integration. In B. C. Pierce, editor, *Proceedings of TLDI 2012: The Seventh ACM SIGPLAN Workshop on Types in Languages Design and Implementation, Philadelphia, PA, USA, Saturday, January 28, 2012*, pages 91–102. ACM, 2012. doi: 10.1145/2103786.2103798. URL <http://doi.acm.org/10.1145/2103786.2103798>.
- F. Loitsch and M. Serrano. Hop client-side compilation. In M. T. Morazán, editor, *Proceedings of the Eighth Symposium on Trends in Functional Programming, TFP 2007, New York City, New York, USA, April 2-4, 2007.*, volume 8 of *Trends in Functional Programming*, pages 141–158. Intellect, 2007.
- D. B. MacQueen. Modules for standard ML. In *LISP and Functional Programming*, pages 198–207, 1984.
- marshal. *OCaml Standard Library – Marshal*. <http://caml.inria.fr/pub/docs/manual-ocaml/libref/Marshal.html>, 2016.
- J. McCarthy. Recursive functions of symbolic expressions and their computation by machine, part I. *Commun. ACM*, 3(4):184–195, 1960. doi: 10.1145/367177.367199. URL <http://doi.acm.org/10.1145/367177.367199>.
- Meteor.js. *Meteor.js*. <http://meteor.com>, 2017.
- R. Milner. A theory of type polymorphism in programming. *J. Comput. Syst. Sci.*, 17(3):348–375, 1978. doi: 10.1016/0022-0000(78)90014-4. URL [http://dx.doi.org/10.1016/0022-0000\(78\)90014-4](http://dx.doi.org/10.1016/0022-0000(78)90014-4).



- R. Milner, M. Tofte, and R. Harper. *Definition of standard ML*. MIT Press, 1990. ISBN 978-0-262-63132-7.
- Y. Minsky, A. Madhavapeddy, and J. Hickey. *Real World OCaml - Functional Programming for the Masses*. O'Reilly, 2013. ISBN 978-1-4493-2391-2. URL <https://realworldocaml.org/>.
- T. Murphy VII, K. Crary, and R. Harper. Type-safe distributed programming with ML5. In G. Barthe and C. Fournet, editors, *TGC*, volume 4912 of *Lecture Notes in Computer Science*, pages 108–123. Springer, 2007. ISBN 978-3-540-78662-7.
- K. Nakata and T. Uustalu. Trace-based coinductive operational semantics for while. In S. Berghofer, T. Nipkow, C. Urban, and M. Wenzel, editors, *Theorem Proving in Higher Order Logics, 22nd International Conference, TPHOLs 2009, Munich, Germany, August 17-20, 2009. Proceedings*, volume 5674 of *Lecture Notes in Computer Science*, pages 375–390. Springer, 2009. doi: 10.1007/978-3-642-03359-9\_26. URL [https://doi.org/10.1007/978-3-642-03359-9\\_26](https://doi.org/10.1007/978-3-642-03359-9_26).
- O. Nicole. Bringing typed, modular macros to ocaml, 2016. URL [https://oliviernicole.github.io/about\\_macros.html](https://oliviernicole.github.io/about_macros.html).
- Ocsigen Toolkit. *Ocsigen Toolkit*. <http://ocsigen.org/ocsigen-toolkit/>, 2017.
- S. Owens. A sound semantics for ocamlight. In S. Drossopoulou, editor, *Programming Languages and Systems, 17th European Symposium on Programming, ESOP 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*, volume 4960 of *Lecture Notes in Computer Science*, pages 1–15. Springer, 2008. doi: 10.1007/978-3-540-78739-6\_1. URL [https://doi.org/10.1007/978-3-540-78739-6\\_1](https://doi.org/10.1007/978-3-540-78739-6_1).
- S. Owens, M. O. Myreen, R. Kumar, and Y. K. Tan. Functional big-step semantics. In Thiemann [2016], pages 589–615. ISBN 978-3-662-49497-4. doi: 10.1007/978-3-662-49498-1\_23. URL [https://doi.org/10.1007/978-3-662-49498-1\\_23](https://doi.org/10.1007/978-3-662-49498-1_23).
- L. Philips, C. De Roover, T. Van Cutsem, and W. De Meuter. Towards tierless Web development without tierless languages. In *Proceedings of the 2014 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming & Software, Onward! 2014*, pages 69–81, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-3210-1. doi: 10.1145/2661136.2661146. URL <http://doi.acm.org/10.1145/2661136.2661146>.
- L. Philips, J. D. Koster, W. D. Meuter, and C. D. Roover. Dependence-driven delimited CPS transformation for javascript. In B. Fischer and I. Schaefer, editors, *Proceedings of the 2016 ACM SIGPLAN International Conference on Generative Programming: Concepts and Experiences, GPCE 2016, Amsterdam, The Netherlands, October 31 - November 1, 2016*, pages 59–69. ACM, 2016. doi: 10.1145/2993236.2993243. URL <http://doi.acm.org/10.1145/2993236.2993243>.

- C. Queinnec. The influence of browsers on evaluators or, continuations to program web servers. In M. Odersky and P. Wadler, editors, *Proceedings of the Fifth ACM SIGPLAN International Conference on Functional Programming (ICFP '00), Montreal, Canada, September 18-21, 2000.*, pages 23–33. ACM, 2000. doi: 10.1145/351240.351243. URL <http://doi.acm.org/10.1145/351240.351243>.
- C. Queinnec. Inverting back the inversion of control or, continuations versus page-centric programming. *SIGPLAN Notices*, 38(2):57–64, 2003. doi: 10.1145/772970.772977. URL <http://doi.acm.org/10.1145/772970.772977>.
- C. Queinnec. Continuations and web servers. *Higher-Order and Symbolic Computation*, 17(4):277–295, 2004. doi: 10.1007/s10990-004-4866-z. URL <https://doi.org/10.1007/s10990-004-4866-z>.
- React. *React*. <http://erratique.ch/software/react>, 2017.
- A. Rossberg, C. V. Russo, and D. Dreyer. F-ing modules. *J. Funct. Program.*, 24(5): 529–607, 2014. doi: 10.1017/S0956796814000264. URL <https://doi.org/10.1017/S0956796814000264>.
- G. Scherer and J. Vouillon. Macaque : Interrogation sûre et flexible de base de données depuis OCaml. In *Ving et unième journées francophones des langages applicatifs*, Studia Informatica Universalis, pages –, La Ciotat, France, Jan. 2010. Hermann. URL <https://hal.archives-ouvertes.fr/hal-00495977>.
- M. Serrano and V. Prunet. A glimpse of hopjs. In J. Garrigue, G. Keller, and E. Sumii, editors, *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming, ICFP 2016, Nara, Japan, September 18-22, 2016*, pages 180–192. ACM, 2016. ISBN 978-1-4503-4219-3. doi: 10.1145/2951913.2951916. URL <http://doi.acm.org/10.1145/2951913.2951916>.
- M. Serrano and C. Queinnec. A multi-tier semantics for Hop. *Higher-Order and Symbolic Computation*, 23(4):409–431, 2010.
- M. Serrano, E. Gallesio, and F. Loitsch. Hop: a language for programming the Web 2.0. In *OOPSLA Companion*, pages 975–985, 2006.
- P. Sewell, J. J. Leifer, K. Wansbrough, F. Z. Nardelli, M. Allen-Williams, P. Habouzit, and V. Vafeiadis. Acute: High-level programming language design for distributed computation. *J. Funct. Program.*, 17(4-5):547–612, 2007. doi: 10.1017/S0956796807006442. URL <http://dx.doi.org/10.1017/S0956796807006442>.
- Shared reactive programming. *Shared React*. <https://ocsigen.org/eliom/5.0/manual/clientserver-react>, 2017.
- T. Sheard and S. L. Peyton Jones. Template meta-programming for haskell. *SIGPLAN Notices*, 37(12):60–75, 2002.

- J. Siek. Type safety in three easy lemmas, May 2013. URL <https://siek.blogspot.fr/2013/05/type-safety-in-three-easy-lemmas.html>.
- D. Swasey, T. Murphy VII, K. Crary, and R. Harper. A separate compilation extension to standard ML. In A. Kennedy and F. Pottier, editors, *Proceedings of the ACM Workshop on ML, 2006, Portland, Oregon, USA, September 16, 2006*, pages 32–42. ACM, 2006. doi: 10.1145/1159876.1159883. URL <http://doi.acm.org/10.1145/1159876.1159883>.
- P. Thiemann, editor. *Programming Languages and Systems - 25th European Symposium on Programming, ESOP 2016, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2016, Eindhoven, The Netherlands, April 2-8, 2016, Proceedings*, volume 9632 of *Lecture Notes in Computer Science*, 2016. Springer. ISBN 978-3-662-49497-4. doi: 10.1007/978-3-662-49498-1. URL <http://dx.doi.org/10.1007/978-3-662-49498-1>.
- M. Tofte. *Operational Semantics and Polymorphic Type Inference*. PhD thesis, University of Edinburgh, 1988.
- Tutorial. *Ocsigen Tutorial*. <https://ocsigen.org/tuto/manual>, 2017.
- TyXML. *TyXML*. <http://ocsigen.org/tyxml/>, 2017.
- J. Vouillon. Lwt: a cooperative thread library. In E. Sumii, editor, *Proceedings of the ACM Workshop on ML, 2008, Victoria, BC, Canada, September 21, 2008*, pages 3–12. ACM, 2008. ISBN 978-1-60558-062-3. doi: 10.1145/1411304.1411307. URL <http://doi.acm.org/10.1145/1411304.1411307>.
- J. Vouillon and V. Balat. From bytecode to JavaScript: the Js\_of\_ocaml compiler. *Software: Practice and Experience*, 44(8):951–972, 2014. ISSN 1097-024X. doi: 10.1002/spe.2187. URL <http://dx.doi.org/10.1002/spe.2187>.
- P. Wadler. Theorems for free! In J. E. Stoy, editor, *Proceedings of the fourth international conference on Functional programming languages and computer architecture, FPCA 1989, London, UK, September 11-13, 1989*, pages 347–359. ACM, 1989. doi: 10.1145/99370.99404. URL <http://doi.acm.org/10.1145/99370.99404>.
- L. White, F. Bour, and J. Yallop. Modular implicits. In O. Kiselyov and J. Garrigue, editors, *Proceedings ML Family/OCaml Users and Developers workshops, ML/OCaml 2014, Gothenburg, Sweden, September 4-5, 2014.*, volume 198 of *EPTCS*, pages 22–63, 2014. doi: 10.4204/EPTCS.198.2. URL <http://dx.doi.org/10.4204/EPTCS.198.2>.
- A. K. Wright and M. Felleisen. A syntactic approach to type soundness. *Information and computation*, 115(1):38–94, 1994.
- J. Yallop and L. White. Modular macros. *OCaml Workshop*, 2015. URL <http://www.lpw25.net/ocaml2015-abs1.pdf>.