

MLton

A Whole-Program Optimizing Compiler for Standard ML

Matthew Fluet

`mtf@cs.rit.edu`

January 5, 2011

Who am I?

- ▶ Matthew Fluet
- ▶ 2nd year faculty member (Dept. of CS)
- ▶ Hooked by sophomore Programming Languages course
- ▶ Studied Programming Languages in graduate school and beyond
 - ▶ theory, design, implementation
 - ▶ **MLton** (a Standard ML compiler)
 - ▶ Cyclone (a safe dialect of C)
 - ▶ Transactional Events (a message-passing + transactions abstraction)
 - ▶ Manticore (a heterogeneous parallel functional language)
 - ▶ Delta ML (a language for self-adjusting computation)

Programming Languages = “Languages” for “Programs”?

I like programming languages because there is nothing like a good language to help us express computations precisely, in ways that we can reason about them, while still keeping things at a high level.

–Norman Ramsey (Tufts University)

Programming Languages = “Languages” for “Programs”?

I like programming languages because there is nothing like a good language to help us express computations precisely, in ways that we can reason about them, while still keeping things at a high level.

–Norman Ramsey (Tufts University)

Language

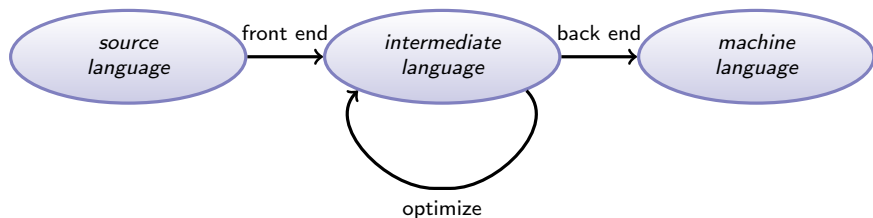
- ▶ Agreed upon medium for communication

Program

- ▶ Executed/interpreted by machine

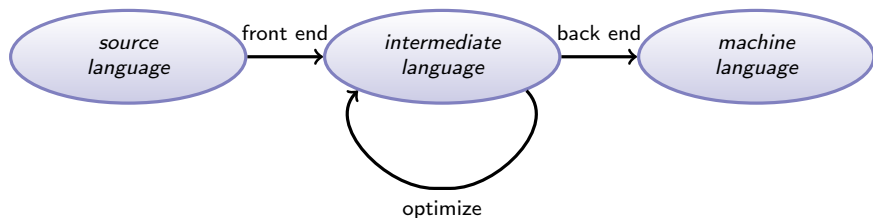
Compilers

Translate high-level program (source language)
into an equivalent executable program (machine language).



Compilers

Translate high-level program (source language)
into an equivalent executable program (machine language).



Compilers are themselves fascinating artifacts:

- ▶ theory meets practice
- ▶ large software system
- ▶ programming language design (programmer-side)
- ▶ programming language implementation (machine-side)

Language Design vs. Language Implementation

Fast programs are better than slow programs.

Short, understandable programs are better than long, confusing programs.

Advanced features make a language more attractive to programmers:

- ▶ simplify the development and maintenance of programs
- ▶ programs are shorter, easier to write and read, less likely to contain errors

Advanced features make a language more difficult for compiler writers:

- ▶ complicate the implementation of the programming language
- ▶ compilers supporting such features are often larger, more difficult to write and extend, and less likely to provide superior performance

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

Design vs. Implementation

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

- ▶ the Right tradeoff

Unfortunate if compiler writers expended great effort to implement advanced features only to have programmers avoid using such features.

Design vs. Implementation

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

- ▶ the Right tradeoff

Unfortunate if compiler writers expended great effort to implement advanced features only to have programmers avoid using such features.

If moving a function definition from one module to another will prevent the function from being inlined and incurs a cost of an extra procedure call at run time, then programmers are less likely to move the function.

Design vs. Implementation

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

- ▶ the Right tradeoff

Unfortunate if compiler writers expended great effort to implement advanced features only to have programmers avoid using such features.

If using a generic type incurs a cost of an extra indirection or an extra word of space, then programmers are less likely to adopt this kind of abstract programming style.

Design vs. Implementation

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

- ▶ the Right tradeoff

Unfortunate if compiler writers expended great effort to implement advanced features only to have programmers avoid using such features.

If using a higher-order function to abstract out an iterator incurs a cost of 20% in the program's running time, then programmers are less likely to use the abstraction.

Design vs. Implementation

The lives of the programmers (the majority) are made easier, while the lives of the compiler writers (the minority) are made harder.

- ▶ the Right tradeoff

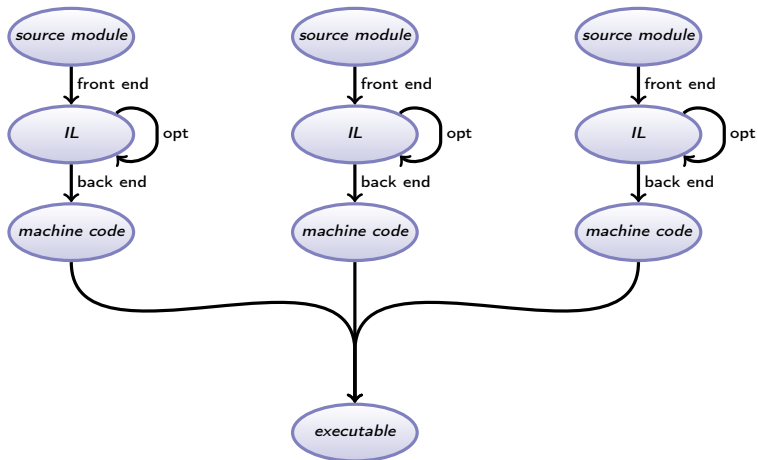
Unfortunate if compiler writers expended great effort to implement advanced features only to have programmers avoid using such features.

A compiler writer's implementation decision forces the programmer to choose between performance and clarity.

Choosing performance (by not choosing to use advanced features) leads to programs that are longer, harder to write, read, and maintain, and more likely to contain errors.

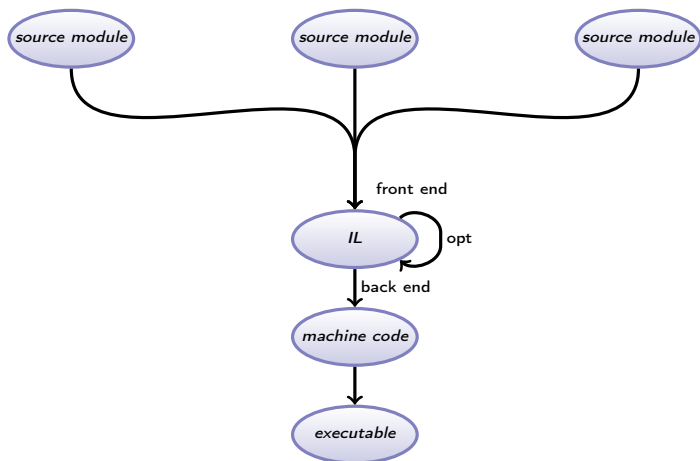
- ▶ the Wrong tradeoff.

Separate Compilation



Compiler has only partial information about the program being compiled.

Whole-Program Compilation



Compiler has maximal information about the program being compiled.
And simplifies the compiler itself (and makes experimentation easier).

A whole-program, optimizing compiler for Standard ML
www.mlton.org

A whole-program, optimizing compiler for Standard ML

www.mlton.org

Robustness

- ▶ Supports the full SML 97 language.
 - ▶ (as given in The Definition of Standard ML (Revised))
- ▶ A complete implementation of the Basis Library.
- ▶ Generates standalone executables.
- ▶ Compiles large programs.
- ▶ Support for large amounts of memory.
 - ▶ (up to 4G on 32-bit systems; more on 64-bit systems)
- ▶ Support for large array lengths.
 - ▶ (up to $2^{31} - 1$ on 32-bit systems; up to $2^{63} - 1$ on 64-bit systems)
- ▶ Support for large files, using 64-bit file positions.

A whole-program, optimizing compiler for Standard ML
`www.mlton.org`

Performance

- ▶ Executables have excellent running times.
- ▶ Generates small executables.
- ▶ Untagged and unboxed native integers, reals, and words.
- ▶ Unboxed native arrays.
- ▶ Multiple garbage collection strategies.
- ▶ Fast arbitrary-precision arithmetic based on the GnuMP.

A whole-program, optimizing compiler for Standard ML

`www.mlton.org`

Tools

- ▶ Source-level profiling for both time and allocation.
- ▶ ML-Lex lexer generator.
- ▶ ML-Yacc parser generator.
- ▶ ML-NLFFIGEN foreign-function-interface generator.
- ▶ ML-ULex lexer generator (next release)
- ▶ ML-Antlr parser generator (next release)

A whole-program, optimizing compiler for Standard ML

www.mlton.org

Extensions

- ▶ A simple and fast C FFI.
 - ▶ (supports calling from SML to C and from C to SML)
- ▶ The ML Basis system for programming in the very large.
- ▶ Libraries for continuations, finalization, interval timers, random numbers, resource limits, resource usage, signal handlers, object hashing, object size, system logging, threads, weak pointers, and world save and restore.

A whole-program, optimizing compiler for Standard ML

www.mlton.org

Portability

- ▶ ARM[†]: Linux (Debian).
- ▶ Alpha[†]: Linux (Debian).
- ▶ AMD64[‡]: Darwin (Mac OS X), FreeBSD, Linux (Debian, Fedora, ...), Solaris (10 and above).
- ▶ HPPA[†]: HP-UX (11.11 and above), Linux (Debian).
- ▶ IA64[†]: HP-UX (11.11 and above), Linux (Debian).
- ▶ PowerPC[†]: AIX (5.2 and above), Darwin (Mac OS X), Linux (Debian, Fedora).
- ▶ PowerPC64[†]: AIX (5.2 and above).
- ▶ S390[†]: Linux (Debian).
- ▶ Sparc[†]: Linux (Debian), Solaris (8 and above).
- ▶ X86[‡]: Cygwin/Windows, Darwin (Mac OS X), FreeBSD, Linux (Debian, Fedora, ...), MinGW/Windows, NetBSD, OpenBSD, Solaris (10 and above).

[†] C backend; [‡] native backend

A whole-program, optimizing compiler for Standard ML

`www.mlton.org`

Statistics

- ▶ Developed since 1997
- ▶ Code
 - ▶ 160k lines SML for the compiler
 - ▶ 23k lines C for the runtime system (incl. garbage collector)
 - ▶ 39k lines SML for the Basis Library
- ▶ Current release: 20100608

MLton: Practical Programming in SML

Efficiency:

- ▶ raw speed
- ▶ eliminate performance disincentives of advanced features

Robustness:

- ▶ adherence to standards, completeness
- ▶ bugs and correctness are a priority
- ▶ support long runs and large inputs

Usability:

- ▶ good type error messages
- ▶ command-line interface, standalone executables
- ▶ large programs (> 100k lines)
- ▶ short enough compile times (< 5 minute self compile)

Commercial Users

Intel Research

- ▶ a top-secret compiler for Tim Sweeney (Epic Games)
- ▶ <http://dl.acm.org/citation.cfm?id=1900175>

PolySpace (MathWorks)

- ▶ a code verifier to detect (or prove the absence of) run-time errors in src code
- ▶ <http://www.mathworks.com/products/polyspace/>

Reactive Systems

- ▶ a testing and validation tool supporting model-based design
- ▶ <http://www.reactive-systems.com/>
- ▶ <http://dl.acm.org/citation.cfm?id=1291172>

SSH Communications Security

- ▶ a collection of internal-use software
- ▶ <http://dl.acm.org/citation.cfm?id=1362714>

Academic Users

Delta ML

- ▶ a language for self-adjusting computation
- ▶ <http://www.mpi-sws.org/~umut/projects/delta-ml/>
- ▶ <http://dl.acm.org/citation.cfm?id=1411249>

Metis

- ▶ an automatic theorem prover for first order logic with equality
- ▶ <http://www.gilith.com/software/metis/>

OpenTheory

- ▶ a tool for processing higher order logic theory packages
- ▶ <http://www.gilith.com/software/opentheory/>

Twelf

- ▶ an implementation of the LF logical framework
- ▶ http://twelf.plparty.org/wiki/Main_Page

A Type Error Slicer for SML

- ▶ a tool for explaining type errors in programs
- ▶ <http://www2.macs.hw.ac.uk/~rahli/cgi-bin/slicer/index.html>

Less Academic Users

SML3d

- ▶ graphics programming using SML and OpenGL
- ▶ <http://sml3d.cs.uchicago.edu/>

HaMLet

- ▶ a reference implementation of SML
- ▶ <http://www.mpi-sws.org/~rossberg/hamlet/>

<http://mlton.org/Users>

Performance: SML Compilers

<http://mlton.org/Performance>

- ▶ (a little out of date)

Compares all of the active SML compilers

- ▶ ML Kit, Moscow ML, MLton, Poly/ML, SML/NJ

with 40+ benchmarks ranging from 15 to 23k lines.

Run-time ratios over all benchmarks:

	MLton	ML Kit	MoscowML	Poly/ML	SML/NJ
median	1.0	2.4	30.6	4.6	3.1
geo. mean	1.0	3.3	25.9	6.2	3.9

Microbenchmarks helpful to compiler writers,
but miss the point for users and whole-program optimization.

- ▶ Every compiler is “whole-program” if the program is 15 lines.

Performance: Large Programs

Large successes:

- ▶ Twelf (67k): “significantly faster” than SML/NJ
- ▶ HaMLet (22k): 3x faster than SML/NJ
- ▶ HOL (120k): 10.3x faster than Moscow ML
- ▶ ML Kit (120k): “significantly faster” than SML/NJ
- ▶ MLton (145k): 81x faster than SML/NJ
- ▶ RML (22k): 2x faster than SML/NJ
- ▶ SML.NET (80k): 3x faster than SML/NJ
- ▶ PolySpace (>100k): commercial, speedup not public

Large failures:

- ▶ HOL (400k)

What is Standard ML?

SML is a general-purpose functional programming language with

- ▶ strict evaluation
- ▶ strong and static typing
- ▶ polymorphic types
- ▶ type inference
- ▶ datatypes and pattern matching
- ▶ functional impurities (mutable objects, side-effects)
- ▶ exceptions
- ▶ a sophisticated module system
- ▶ a rigorous formal definition

Standard ML = Core Language + Module Language

SML is made up of two sub-languages

- ▶ core language:
 - ▶ expressing types and computations
- ▶ module language
 - ▶ packaging elements of core language into units for modularity and reuse
 - ▶ Structures
 - ▶ an encapsulated, named, collection of (type and value) declarations
 - ▶ Signatures
 - ▶ an encapsulated, named, collection of specifications
 - ▶ classify structures
 - ▶ Functors
 - ▶ an encapsulated, named, function from structures to structures

The module language is a *language*:

- ▶ It has non-trivial static and dynamic semantics.
- ▶ It is not simply a namespace management veneer.

SML Implementations

- ▶ Poly/ML (1983)
 - ▶ <http://www.polyml.org>
 - ▶ very fast compilation; REPL
- ▶ Standard ML of New Jersey (1986)
 - ▶ <http://www.smlnj.org>
 - ▶ continuation-passing style; incremental compilation and REPL
- ▶ ML Kit (1989)
 - ▶ http://www.itu.dk/research/mlkit/index.php/Main_Page
 - ▶ region-based memory management
- ▶ MoscowML (1994)
 - ▶ <http://www.itu.dk/people/sestoft/mosml.html>
 - ▶ bytecode compiler (using Caml Light runtime system)
- ▶ MLton (1997)
 - ▶ <http://www.mlton.org>
 - ▶ whole-program optimization
- ▶ HaMLet (1999)
 - ▶ <http://www.mpi-sws.org/~rossberg/hamlet/>
 - ▶ reference interpreter (written in SML, following the Definition)

Compiling SML Efficiently is Hard

Advanced features lead to missing information.

- ▶ higher-order functions \Rightarrow missing control-flow info
- ▶ polymorphism \Rightarrow missing type info
- ▶ functors \Rightarrow missing control-flow and type info

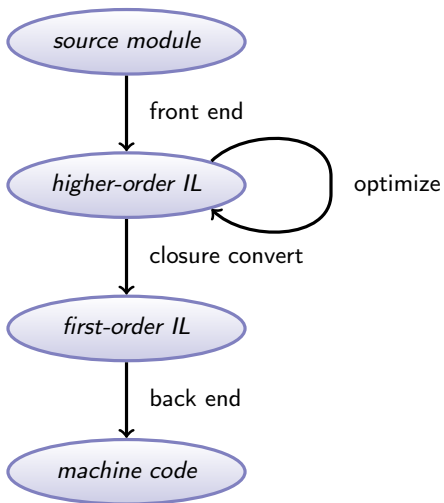
Missing information leads to bad code.

- ▶ inefficient data representations:
 - ▶ tagged integers, boxing, no packing, extra variant tags
- ▶ missed control-flow optimizations:
 - ▶ inlining, loop optimizations, dead-code elimination

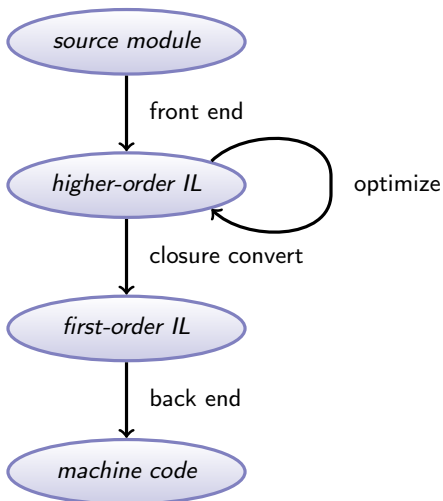
Compiler writers tie both hands behind their back.

- ▶ separate compilation
- ▶ complex, nonstandard intermediate languages for optimization

Traditional Approach to Compiling SML

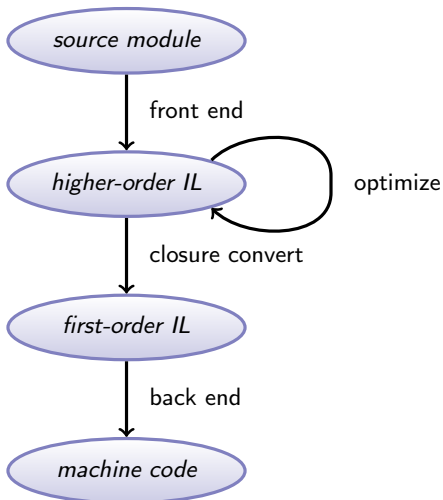


Traditional Approach to Compiling SML

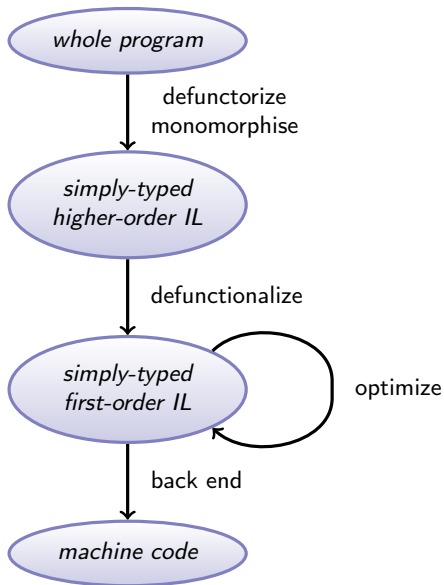


- ▶ Separate compilation:
 - ▶ bad type info
 - ▶ bad control-flow info
- ▶ Polymorphic or untyped IL:
 - ▶ bad data representations
 - ▶ bad dataflow analyses
- ▶ Higher-order IL:
 - ▶ can't use traditional optimizations
 - ▶ optimizations do their own control-flow analysis
- ▶ Poor closure optimization

Traditional Approach

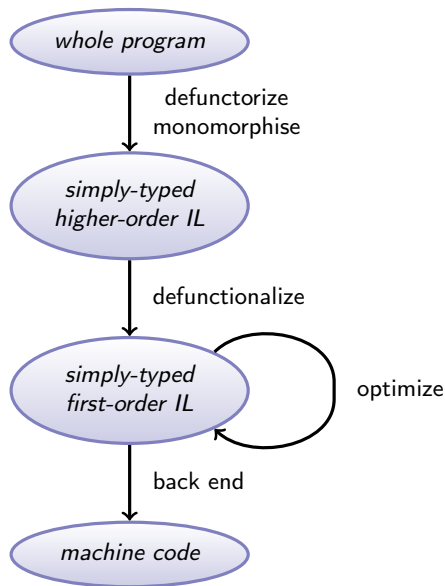


MLton's Approach



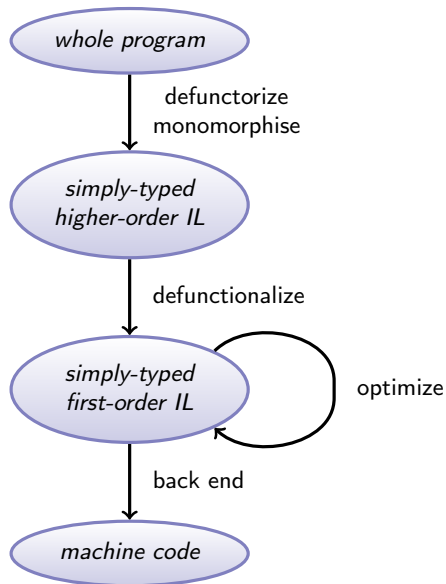
Benefits of MLton's Approach

- ▶ Absolute efficiency:
 - ▶ massive optimization
 - ▶ good control-flow info
 - ▶ good data representations
- ▶ Relative efficiency:
 - ▶ zero-cost or low-cost advanced features
- ▶ Simplicity:
 - ▶ simple, typed IL
 - ▶ traditional optimizations
 - ▶ optimizations don't have to do their own CFA

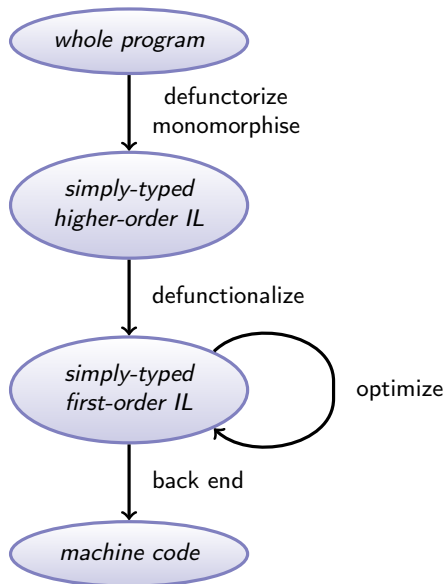


Drawbacks of MLton's Approach

- ▶ Compile time.
- ▶ Compile memory.
- ▶ Executable size.



MLton's Approach



Defunctorization

Goals:

- ▶ turn full SML into a polymorphic, higher-order IL
- ▶ expose types hidden by functors and signatures
- ▶ expose function calls across modules
- ▶ zero-cost modules for programmer

Method:

- ▶ eliminate structures and signatures
- ▶ duplicate each functor at every use

Code explosion in theory, but not in practice.

Defunctorization: Examples

Structure declarations are eliminated,
with all declarations moved to the top level.
Long identifiers are renamed.

```
structure S =  
  struct  
    type t = int  
    val x : t = 13  
  end  
val y : S.t = S.x
```

⇒

```
val x_0 : int = 13  
val y_0 : int = x_0
```

Open declarations are eliminated.

```
val x = 13  
val y = 14  
structure S =  
  struct  
    val x = 15  
  end  
open S  
val z = x + y
```

⇒

```
val x_0 = 13  
val y_0 = 14  
val x_1 = 15  
val z_0 = x_1 + y_0
```

Defunctorization: Examples

Functor declarations are eliminated,
and the body of a functor is duplicated wherever the functor is applied.

```
functor F(val x : int) =  
  struct  
    val y = x  
  end  
structure F1 = F(val x = 13)  
structure F2 = F(val x = 14)  
val z = F1.y + F2.y
```



```
val x_0 = 13  
val y_0 = x_0  
val x_1 = 14  
val y_1 = x_1  
val z_0 = y_0 + y_1
```

Defunctorization

Goals:

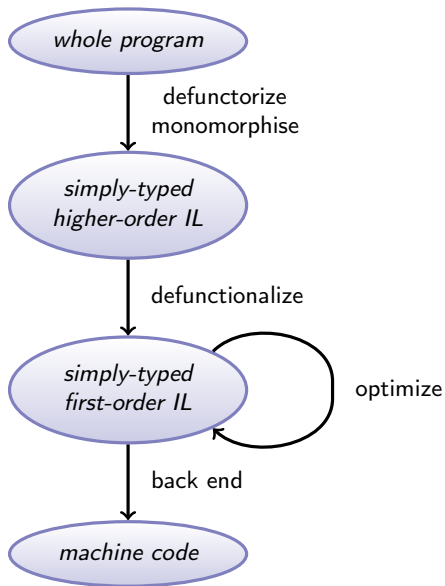
- ▶ turn full SML into a polymorphic, higher-order IL
- ▶ expose types hidden by functors and signatures
- ▶ expose function calls across modules
- ▶ **zero-cost modules for programmer**

Method:

- ▶ eliminate structures and signatures
- ▶ duplicate each functor at every use

Code explosion in theory, but not in practice.

MLton's Approach



Monomorphisation

Goals:

- ▶ eliminate polymorphism, producing a simply-typed IL
- ▶ enable good data representations
- ▶ zero-cost polymorphism for programmers

Method:

- ▶ duplicate type declarations at each type used
- ▶ duplicate function declarations at each type used
- ▶ rely on properties of SML for termination

Code explosion in theory, but manageable in practice.

- ▶ (Max increase seen is 30% in MLton)

Subtleties: non-uniform datatypes, phantom types.

Monomorphisation: Example

Polymorphic program:

```
datatype 'a t = T of 'a
fun 'a f (x: 'a) = T x
val a = f 1
val b = f 2
val z = f (3, 4)
```

Monomorphic program:

```
datatype t1 = T1 of int
datatype t2 = T2 of int * int
fun f1 (x: t1) = T1 x
fun f2 (x: t2) = T2 x
val a = f1 1
val b = f1 2
val z = f2 (3, 4)
```

Monomorphisation

Goals:

- ▶ eliminate polymorphism, producing a simply-typed IL
- ▶ enable good data representations
- ▶ **zero-cost polymorphism for programmers**

Method:

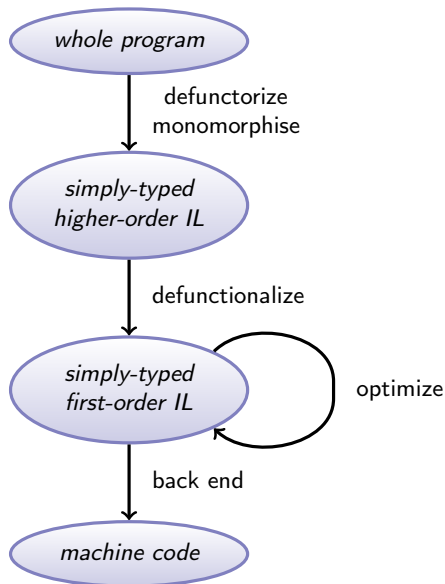
- ▶ duplicate type declarations at each type used
- ▶ duplicate function declarations at each type used
- ▶ rely on properties of SML for termination

Code explosion in theory, but manageable in practice.

- ▶ (Max increase seen is 30% in MLton)

Subtleties: non-uniform datatypes, phantom types.

MLton's Approach



Defunctionalization

Goals:

- ▶ eliminate higher-order functions, producing a first-order IL
- ▶ make direct top-level calls, which are easy to optimize
- ▶ make control-flow info available to rest of optimizer
- ▶ optimize closures just like other data structures

Method:

- ▶ moves nested functions to top level
- ▶ function = tagged record of free variables
- ▶ call = dispatch on tag followed by top-level call
- ▶ control-flow analysis to minimize dispatches

Control-Flow Analysis (OCFA)

OCFA: whole-program dataflow analysis

- ▶ computes set of functions at each call site

Imprecise in theory, but precise in practice.

- ▶ almost no calls require case dispatch

Cubic time in theory, but very fast in MLton.

- ▶ less than 2s to analyze MLton itself
- ▶ preprocessing based on types
- ▶ ignore first-order values
- ▶ hash cons sets and cache binary operations
- ▶ use union-find for equality constraints

Prior code duplication helps speed and precision.

Defunctionalization: Example

Higher-order program:

```
val f = fn a => fn b => a
val g = fn c => fn d => d
val h = if x < y then f else g
val m = h 13
val z = m 7
```

Flow analysis:

$$\begin{aligned} F(f) &= \{\mathbf{fn} \ a\} \\ F(g) &= \{\mathbf{fn} \ c\} \\ F(h) &= \{\mathbf{fn} \ a, \mathbf{fn} \ c\} \\ F(m) &= \{\mathbf{fn} \ b, \mathbf{fn} \ d\} \end{aligned}$$

Defunctionalization: Example

First-order program:

```
datatype t1 = C1 of unit (* fn a *)
datatype t2 = C2 of int  (* fn b *)
datatype t3 = C3 of unit (* fn c *)
datatype t4 = C4 of unit (* fn d *)
datatype t5 = C5 of unit (* fn a  $\in F(h)$  *)
      | C6 of unit (* fn c  $\in F(h)$  *)
datatype t6 = C7 of int  (* fn b  $\in F(m)$  *)
      | C8 of unit (* fn d  $\in F(m)$  *)

fun F0 (r, a) = C2 a      (* code for fn a *)
fun F2 (r as (a), b) = a (* code for fn b *)
fun G0 (r, c) = C4 ()    (* code for fn c *)
fun G2 (r, d) = C5 ()    (* code for fn d *)

val f = C1()
val g = C3()
val h = if x < y
      then (case f of C1 r => C5 r)
      else (case g of C3 r => C6 r)

val m = case h of
      C5 r => (case F0 (r, 13) of C2 r => C7 r)
      | C6 r => (case G0 (r, 13) of C4 r => C8 r)

val z = case m of
      C7 r => F1 (r, 7)
      | C8 r => G1 (r, 7)
```

Defunctionalization

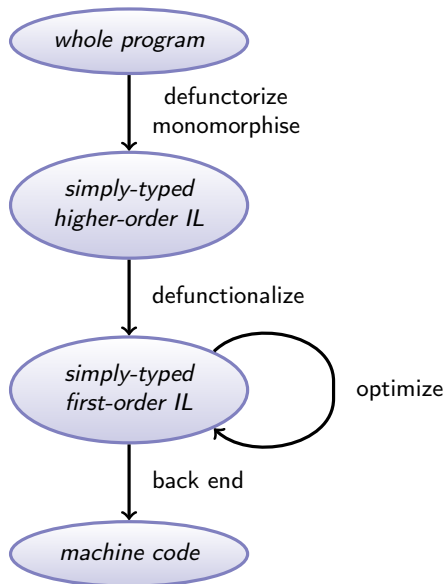
Goals:

- ▶ eliminate higher-order functions, producing a first-order IL
- ▶ make direct top-level calls, which are easy to optimize
- ▶ **make control-flow info available to rest of optimizer**
- ▶ **optimize closures just like other data structures**

Method:

- ▶ moves nested functions to top level
- ▶ function = tagged record of free variables
- ▶ call = dispatch on tag followed by top-level call
- ▶ control-flow analysis to minimize dispatches

MLton's Approach



SSA Intermediate Language

Traditional, simple IL.

- ▶ simply-typed, first-order
- ▶ program = datatypes + functions
- ▶ function = type + arguments + control-flow graph
- ▶ usual SSA conditions: def once, def dominates use

300 line interface, 2K line implementation

- ▶ immutable IL
- ▶ pretty printer, CFG visualizer, DFS, utilities

MLton options:

- ▶ `-drop-pass`, `-diag-pass`, `-keep-pass`

SSA Type Checker

Verifies:

- ▶ uniqueness of names
- ▶ variable definitions dominate uses
- ▶ control-flow graphs are well formed
- ▶ types at primitive applications and calls

Runs at beginning and end of optimizer.

Can be run after each pass (`-type-check true`).

- ▶ slows down optimizer by 50%
- ▶ catches bugs in optimizer

700 lines of code.

SSA Example: Nontail Fib

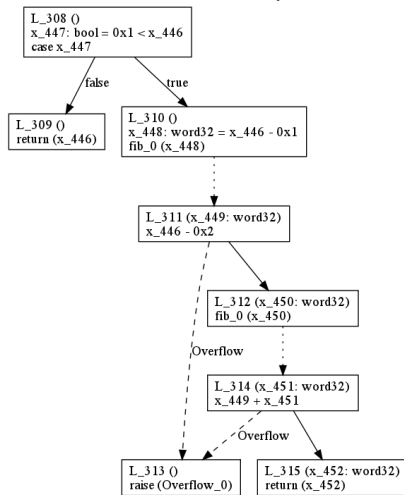
SML Source Function

```
fun fib n =  
  if n <= 1 then  
    n  
  else  
    fib (n - 1) +  
    fib (n - 2)
```

SSA Top-Level Function

```
fun fib_0 (x_446: word32):  
  {raises = Some (exn),  
   returns = Some (word32)} =  
  L_308 ()
```

SSA Control-Flow Graph



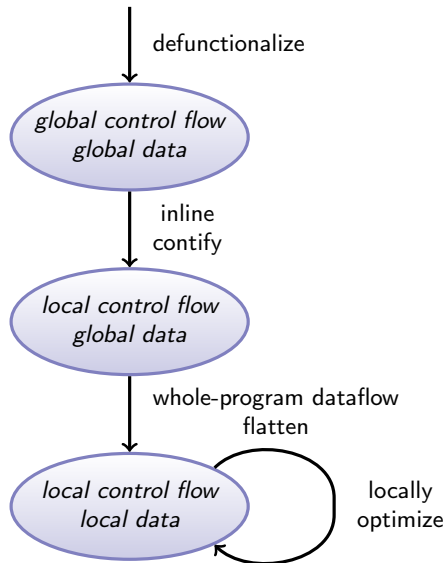
SSA Optimizer

Goals:

- ▶ turn function calls into control-flow graphs
- ▶ expose interprocedural data
- ▶ reduce tuple allocation
- ▶ traditional local optimizations

Method:

- ▶ 24 small, independent, SSA-to-SSA rewrite passes
- ▶ each pass: analyze, transform, shrink



SSA Shrinker

Goals:

- ▶ perform “obvious” local simplification
- ▶ let other optimizations focus on what they do best
- ▶ keep SSA IL programs small

Method:

- ▶ depth-first search of control-flow graph for each function
- ▶ reduce: primapps, case of variant, select of tuple, ...
- ▶ Appel-Jim shrinker applied to SSA.

One of the largest SSA pass, 1400 lines.

Inlining and Contification

Goals:

- ▶ turn function calls into control-flow graphs
- ▶ eliminate call overhead

Leaf inlining.

- ▶ uncurrying for free

Call-graph inlining.

- ▶ inline if: $(numCalls - 1) * (size - c) \leq limit$

Contification.

- ▶ turns functions used as continuations into jumps

Relies on OCFA and first-order whole program.

Whole-Program Dataflow Optimizations

Goals:

- ▶ expose data to shrinker and later optimizations
- ▶ clean up across modules

Constant propagation.

- ▶ analyze: forwards from constants with a flat lattice
- ▶ transform: replace variables with constants

Useless-component removal.

- ▶ analyze: backwards from primitives, tests, FFI, ...
- ▶ transform: eliminate useless component

Flattening Optimizations

Goals:

- ▶ eliminate indirection (save space and time)
- ▶ pack tuples
- ▶ reduce allocation

Method:

- ▶ flatten function arguments and results
- ▶ flatten constructor applications
- ▶ flatten ref cells into data structures and stack frames
- ▶ flatten array components
- ▶ flatten basic-block arguments

Caveat: space safety

SSA Example: List Fold Becomes a Loop

SML Source Functions

```
fun fold (l, b, f) = let
  fun loop (l, b) =
    case l of
      [] => b
    | x :: l => loop (l, f (x, b))
  in loop (l, b) end

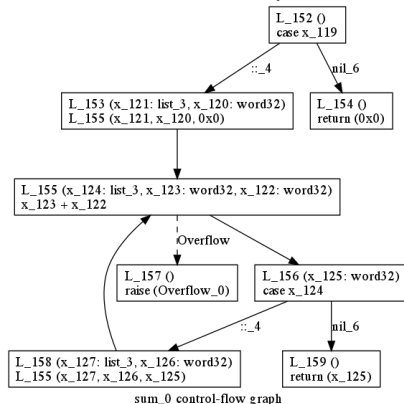
fun sum l =
  fold (l, 0, fn (x, y) => x + y)
```

SSA Top-Level Function

```
list_3 = nil_6 | ::_4 of (list_3, word32)

fun sum_0 (x_119: list_3):
  {raises = Some (exn),
   returns = Some (word32)} =
  L_152 ()
```

SSA Control-Flow Graph



Dominator-based Local Optimizations

Goals:

- ▶ apply traditional intraprocedural optimizations
- ▶ take advantage of prior whole-program optimization

Method:

- ▶ compute dominator tree for each function's CFG
- ▶ recursively walk tree, use known facts in subtrees

Examples:

- ▶ common-subexpression elimination
- ▶ known-case elimination
- ▶ redundant-test elimination (includes bounds checks)
- ▶ overflow-detection elimination

SSA Optimizer During a Self Compile

	#functions	#statements	size
start	32,576	457,944	248,319,824 bytes
leaf inline	18,019	453,670	177,866,544 bytes
contify	6,940	429,630	162,863,152 bytes
constant prop	6,940	271,492	124,894,336 bytes
inline	1,893	373,038	141,519,200 bytes
flatten	1,889	260,916	133,347,432 bytes
local opts	1,889	239,228	131,298,528 bytes

Data Representations

Goals:

- ▶ choose efficient representation for each IL type
- ▶ save space and allocation
- ▶ make GC fast and easy

Method:

- ▶ pack tuples and array elements
- ▶ unbox datatype variants (including lists)
- ▶ reorder fields
- ▶ use untagged integers and words
- ▶ fast card marking

Simply-typed whole program is essential.

Technical Lessons

Whole-program compilation is feasible.

- ▶ compile a 100k line program in minutes with 1G RAM
- ▶ myths: defunctorization, monomorphisation, OCFA

Whole-program compilation is effective.

- ▶ fast code and compact data representations
- ▶ total information \Rightarrow optimizations rewrite at will

Whole-program compilation is simple.

- ▶ simplifies compiler
- ▶ simplifies optimizations
- ▶ simplifies intermediate languages

Technical Lessons

Simply-typed, first-order SSA is an excellent compiler IL, even for advanced languages.

- ▶ complete type information
- ▶ all passes benefit from CFA, which is only done once
- ▶ traditional optimizations

Structuring an optimizer as small, independent rewrite passes on an immutable IL makes life easy.

- ▶ easy to develop new passes
- ▶ easy to debug old passes
- ▶ easy to experiment with phase ordering
- ▶ easy for passes to help each other

Future Work

New Optimizations:

- ▶ Improved closure representations
- ▶ Array bounds check elimination
- ▶ Overflow check elimination
- ▶ Loop-invariant code motion
- ▶ Partial redundancy elimination
- ▶ Loop unrolling
- ▶ Type splitting

Full employment theorem for compiler writers.

Future Work

New Optimizations Frameworks:

- ▶ Hoopl: A Modular, Reusable Library for Dataflow Analysis and Transformation
 - ▶ <http://portal.acm.org/citation.cfm?id=1863539>
- ▶ Equality Saturation: A New Approach To Optimization
 - ▶ <http://portal.acm.org/citation.cfm?id=1480915>
- ▶ Supercompilation
 - ▶ <http://portal.acm.org/citation.cfm?id=1480916>
 - ▶ <http://portal.acm.org/citation.cfm?id=1863588>
 - ▶ <http://portal.acm.org/citation.cfm?id=1863540>

New Optimization Verifications:

- ▶ Safe-for-space

Future Work

New Analyses:

- ▶ Uncaught exception analysis
- ▶ Alternative control-flow analyses for defunctionalization

CFA research has had a renaissance;
an opportunity to test new “theory” in “practice”.

Future Work

New Backends:

- ▶ LLVM
- ▶ MLRISC

MLton wins because of its SSA optimizer,
and in spite of its simple native codegen.

Future Work

New Tools:

- ▶ Heap profiler
- ▶ Debugger
- ▶ Statistically sound benchmarking

Future Work

New Language Features (easy):

- ▶ Record punning, Record extension, Record update
- ▶ Or-patterns
- ▶ Vector literals
- ▶ Higher-order functors
- ▶ Nested signatures
- ▶ Local modules

Future Work

New Intermediate-Language Features (moderate):

- ▶ Multi-return function calls
- ▶ Multi-entry function calls
- ▶ Heap-allocated stack frames

Future Work

New Runtime System Features (moderate to difficult):

- ▶ Garbage collector improvements
- ▶ Multicore

Future Work

New Language Features (hard):

- ▶ Polymorphic recursion
- ▶ First-class polymorphism
- ▶ Generalized abstract datatypes

These features cannot be *completely* monomorphised.

Future Work

New Language Features (hard):

- ▶ Polymorphic recursion
- ▶ First-class polymorphism
- ▶ Generalized abstract datatypes

These features cannot be *completely* monomorphised.

“Type-Flow Analysis for Partial Monomorphisation”

- ▶ monomorphise “as much as possible, and no more”

“Representation Monomorphisation”

- ▶ monomorphise with respect to back-end data representations

Future Work

???

Languages and Compilers provide lots of room to explore.

MLton Credits

Design:

- ▶ Henry Cejtin, Matthew Fluet, Suresh Jagannathan, Stephen Weeks

Implementation:

- ▶ Matthew Fluet, Stephen Weeks

Code:

- ▶ `gdt oa`, GnuMP, ML Kit, Moscow ML, SML/NJ

Support:

- ▶ NEC, InterTrust, PolySpace, Reactive Systems

People:

- ▶ Jesper Louis Andersen, Johnny Andersen, Alain Deutsch, Martin Elsman, Brent Fulgham, Adam Goode, Simon Helsen, Joe Hurd, Vesa Karvonen, Richard Kelsey, Ville Laurikari, Geoffrey Mainland, Eric McCorkle, Tom Murphy, Michael Neumann, Barak Pearlmutter, Filip Pizlo, John Reppy, Sam Rushing, Jeffrey Mark Siskind, Wesley Terpstra, Luke Ziarek

Questions?