Lucrative Late Lambda Lifting

Sebastian Graf

September 17, 2018

1 Introduction

2 Transformation

Lambda lifting is a well-known technique [2]. Although Johnsson's original algorithm runs in wort-case cubic time relative to the size of the input program, Morazán and Schultz [4] gave an algorithm that runs in $\mathcal{O}(n^2)$.

Our lambda lifting transformation is unique in that it operates on terms of the *spineless tagless G-machine* (STG) [1] as currently implemented [3] in GHC. This means we can assume that the nesting structure of bindings corresponds to the condensation (the directed acyclic graph of strongly connected components) of the dependency graph. TODO: less detail? less language? Additionally, every binding in a (recursive) let expression is annotated with the free variables it closes over. The combination of both properties allows efficient construction of the set of required TODO: any better names? former free variables, abstraction variables... variables for a total complexity of $\mathcal{O}(n^2)$, as we shall see.

2.1 Syntax

Although STG is but a tiny language compared to typical surface languages such as Haskell, its definition [3] still contains much detail irrelevant to lambda lifting. As can be seen in fig. 1, we therefore adopt a simple lambda calculus with **let** bindings as in Johnsson [2], with a few STG-inspired features:

- 1. **let** bindings are annotated with the free variables they close over
- 2. Arguments in an application expression are all atomic (e.g. variable references)
- 3. Every lambda abstraction is the right-hand side of a **let** binding
- 4. All applications are fully saturated

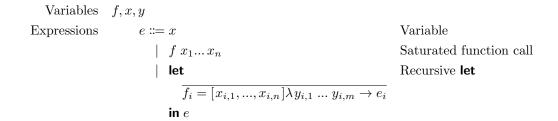


Figure 1: An STG-like untyped lambda calculus

3 When to lift

Lambda lifting a binding to top-level is always TODO: except when we would replace a parameter occurrence by an application a sound transformation. The challenge is in identifying *when* it is beneficial to do so. This section will discuss operational consequences of lambda lifting, introducing multiple criteria based on a cost model for estimating impact on heap allocations.

3.1 Syntactic consequences

Deciding to lift a binding let $f = [x \ y \ z] \lambda a \ b \ c \rightarrow e_1$ in e_2 to top-level has the following consequences:

- (S1) It eliminates the **let** binding.
- (S2) It creates a new top-level definition.
- (S3) It replaces all occurrences of f in e_2 by an application of the lifted toplevel binding to its former free variables, replacing the whole **let** binding by the term $[f \mapsto f_{\uparrow} \times y \ z] \ e_2$. TODO: Maybe less detail here
- (S4) All non-top-level variables that occurred in the **let** binding's right-hand side become parameter occurrences.

Consider what happens if f occurred in e_2 as an argument in an application, as in g 5 \times f. (S3) demands that the argument occurrence of f is replaced by an application expression. This, however, would yield a syntactically invalid expression, because the STG language only allows trivial arguments in an application.

An easy fix would be to bind the complex expression to an auxiliary **let** binding, thereby re-introducing the very allocation we wanted to eliminate through lambda lifting TODO: Move this further down?. Therefore, we can identify a first criterion for non-beneficial lambda lifts:

(C1) Don't lift binders that occur as arguments

3.2 Operational consequences

We now ascribe operational symptoms to combinations of syntactic effects. These symptoms justify the derivation of heuristics which will decide when *not* to lift.

Closure growth. (S1) means we don't allocate a closure on the heap for the **let** binding. On the other hand, (S3) might increase or decrease heap allocation. Consider this example:

```
let f = [x \ y] \lambda a \ b \rightarrow ...

g = [f \ x] \lambda d \rightarrow f \ d \ d + x

in g \ 5

Should f be lifted? It's hard to say without actually seeing the lifted version:

f_{\uparrow} = \lambda x \ y \ a \ b \rightarrow ...;

let g = [x \ y] \lambda d \rightarrow f_{\uparrow} x \ y \ d \ d + x

in g \ 5
```

Just counting the number of variables occurring in closures, the effect of (S1) saved us two slots. At the same time, (S3) removes f from g's closure (no need to close over the top-level f_{\uparrow}), while simultaneously enlarging it with f's former free variable g. The new occurrence of g doesn't contribute to closure growth, because it already occurred in g prior to lifting. The net result is a reduction of two slots, so lifting g seems worthwhile. In general:

(C2) Don't lift a binding when doing so would increase closure allocation

Estimation of closure growth is crucial to identifying beneficial lifting opportunities. We discuss this further in 3.3.

Calling Convention. (S4) means that more arguments have to be passed. Depending on the target architecture, this means more stack accesses and/or higher register pressure. Thus

(C3) Don't lift a binding when the arity of the resulting top-level definition exceeds the number of available hardware registers (e.g. 5 arguments on x86_64)

Turning known calls into unknown calls. There's another aspect related to (S4), relevant in programs with higher-order functions:

Here, there is a *known call* to f in mapF that can be lowered as a direct jump to a static address [3]. Lifting mapF (but not f) yields the following program:

$$\begin{aligned} & \textit{mapF}_{\uparrow} = \lambda f \; \textit{xs} \rightarrow ...f \; \textit{x...}; \\ & \textbf{let} \; f = [\,] \lambda x \rightarrow 2 * x \\ & \textbf{in} \; \textit{mapF}_{\uparrow} \; f \; [1,2,3] \end{aligned}$$

(C4) Don't lift a binding when doing so would turn known calls into unknown calls

Undersaturated calls. When GHC spots an undersaturated call, it arranges allocation of a partial application that closes over the supplied arguments. Pay attention to the call to f in the following example:

let
$$f = [x]\lambda y \ z \rightarrow x + y + z;$$

in map $(f \ x) \ [1, 2, 3]$

Here, the undersaturated (e.g. curried) call to f leads to the allocation of a partial application, carrying two pointers, to f and x, respectively. What happens when f is lambda lifted?

$$f_{\uparrow} = \lambda x \ y \ z \rightarrow x + y + z;$$

 $map (f_{\uparrow} x x) [1, 2, 3]$

The call to f_{\uparrow} will still allocate a partial application, with the only difference that it now also closes over f's free variable x, canceling out the beneficial effects of (S1). Hence

(C5) Don't lift a binding that has undersaturated calls

Sharing. Let's finish with a no-brainer: Lambda lifting updatable bindings (e.g. thunks) or constructor bindings is a bad idea, because it destroys sharing, thus possibly duplicating work in each call to the lifted binding.

(C6) Don't lift a binding that is updatable or a constructor application

3.3 Estimating Closure Growth

Of the criterions above, (C2) is the most important for reliable performance gains. It's also the most sophisticated, because it entails estimating closure growth.

Let's revisit the example from above:

We concluded that lifting f would be beneficial, saving us allocation of one free variable slot. There are two effects at play here. Not having to allocate the closure of f due to (S1) always leads to a one-time benefit. Simultaneously, each closure occurrence of f would be replaced by its referenced free variables. Removing f leads to a saving of one slot per closure, but the free variables x and y each occupy a closure slots in turn. Of these, only y really contributes to closure growth, because x already occurred in the single remaining closure of g.

This phenomenon is amplified whenever allocation happens under a multishot lambda, as the following example demonstrates:

let
$$f = [x \ y] \lambda a \ b \rightarrow ...$$

 $g = [f \ x] \lambda d \rightarrow$
let $h = [f] \lambda e \rightarrow f \ e \ e$

```
\begin{array}{c} \text{in } h \ d \\ \text{in } g \ 1+g \ 2+g \ 3 \end{array}
```

Is it still beneficial to lift f? Following our reasoning, we still save two slots from f's closure, the closure of g doesn't grow and the closure h grows by one. We conclude that lifting f saves us one closure slot. But that's nonsense! Since g is called thrice, the closure for h also gets allocated three times relative to single allocations for the closures of f and g.

In general, h might be occurring inside a recursive function, for which we can't reliably estimate how many times its closure will be allocated. Disallowing to lift any binding which is called inside a closure under such a multi-shot lambda is conservative, but rules out worthwhile cases like this:

```
\begin{aligned} \textbf{let } f &= [x\ y] \lambda a\ b \rightarrow ... \\ g &= [f\ x\ y] \lambda d \rightarrow \\ \textbf{let } h_1 &= [f] \lambda e \rightarrow f\ e\ e \\ h_2 &= [f\ x\ y] \lambda e \rightarrow f\ e\ e + x + y \\ \textbf{in } h_1\ d + h_2\ d \\ \textbf{in } g\ 1 + g\ 2 + g\ 3 \end{aligned}
```

Here, the closure of h_1 grows by one, whereas that of h_2 shrinks by one, cancelling each other out. We express this in our cost model by an infinite closure growth whenever there was any positive closure growth under a multishot lambda.

One final remark regarding analysis performance. TODO: equation operates directly on STG expressions. This means the cost function has to traverse whole syntax trees for every lifting decision.

Instead, our implementation first abstracts the syntax tree into a *skeleton*, retaining only the information necessary for our analysis. In particular, this includes allocated closures and their free variables, but also occurrences of multishot lambda abstractions. Additionally, there are the usual "glue operators", such as sequence (e.g. the case scrutinee is evaluated whenever one of the case alternatives is), choice (e.g. one of the case alternatives is evaluated *mutually exclusively*) and an identity (e.g. literals don't allocate).

References

- [1] Olivier Danvy and Ulrik P. Schultz. "Lambda-lifting in quadratic time". In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). Vol. 2441. 2002, pp. 134–151. ISBN: 3540442332. DOI: 10.1007/3-540-45788-7.
- [2] Thomas Johnsson. "Lambda lifting: Transforming programs to recursive equations". In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). Vol. 201 LNCS. 1985, pp. 190–203. ISBN: 9783540159759. DOI: 10.1007/3-540-15975-4_37.

- [3] Simon Marlow and Simon Peyton Jones. "Making a fast curry". In: Proceedings of the ninth ACM SIGPLAN international conference on Functional programming ICFP '04. 2004, p. 4. ISBN: 1581139055. DOI: 10.1145/1016850.1016856. URL: http://portal.acm.org/citation.cfm?doid=1016850.1016856.
- [4] Marco T. Morazán and Ulrik P. Schultz. "Optimal lambda lifting in quadratic time". In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). Vol. 5083 LNCS. 2008, pp. 37–56. ISBN: 3540853723. DOI: 10.1007/978-3-540-85373-2_3.