



Trinity College Dublin

Coláiste na Tríonóide, Baile Átha Cliath

The University of Dublin

Eliminating Small, Transient Memory Allocations

Dário Tavares Antunes

B. A. (Mod.) Computer Science

Final Year Project May 2018

Supervisor: Dr. David Abrahamson

School of Computer Science and Statistics
O'Reilly Institute, Trinity College, Dublin 2, Ireland

Declaration

I hereby declare that this project is entirely my own work and that it has not been submitted as an exercise for a degree at this or any other university.

Dário Tavares Antunes, January 1, 1970

Permission To Lend

I agree that the Library and other agents of the College may lend or copy this report upon request.

Dário Tavares Antunes, January 1, 1970

Abstract

This is where I would put the abstract, IF I HAD ONE

Acknowledgements

This is where I would put my acknowledgements, IF I HAD ANY

Contents

1	Introduction	1
2	Background and Objectives	2
3	Implementation	3
4	Case Studies	4
5	Conclusion	5
6	Future Work	6
6.1	Detecting Arbitrary Memory Allocations	6
6.2	Automatically Performing Fixes	7
6.3	Studying Other Architectures	7
	Bibliography	8

Chapter 1

Introduction

640K ought to be enough for
anybody.

Not Bill Gates

Despite the often misattributed epigraph above often being used to mock past beliefs that some amount of memory should be enough for any reasonable purposes, the mentality behind it is still pervasive.

With the broad availability of large amounts of computational power, memory and storage, conservation or efficient use of the same is often overlooked in programming. This is largely perpetuated by the (often valid) view that programmer time is more valuable than the benefits that more efficient but more complex code brings.

However, there remain situations where these benefits are in fact worth the effort required. One of these such cases is in code intended to be deployed in embedded or mobile devices, where resources are limited and preservation of power is essential.

A blog post [Ste17] by Daniel Stenberg, original author of the `curl` command line tool and ubiquitous URL data transfer tool, is a retrospective on an attempt to reduce unnecessary heap allocations.

Inspired by that post, the aim of this project is to produce a tool to identify cases where similar changes could be made in order to potentially reduce a program's energy and processing power footprint, and at the same time improve its performance.

Chapter 2

Background and Objectives

This chapter will contain:

2.1 Background on the Patch

Background on the patch as applied by Daniel Stenberg

2.2 Hypothesis for the Results of the Patch

Hypothesis as to why the patch caused a performance increase (why that much vs more/less)

2.3 Similarity of the Patch to Concepts in Generational Garbage Collection

Background based on similarity to generational GC assumptions (Appel, Shao)

2.4 Predictions on Results of Generalisation of the Patch

Predictions on the results of a generalised application of the patch (vs proebsting's law for compilers)

Chapter 3

Implementation

This chapter will contain descriptions of:

3.1 Goals of the Plugin

What are the initial goals of the plugin, including those that weren't achieved? How do these goals directly relate to the patch?

3.2 The Frama-C Platform

What is Frama-C? Why choose Frama-C (over bespoke solutions or other platforms)? Background on Frama-C, its goals.

3.3 Difficulties Encountered

What difficulties were encountered (new language, installation, documentation, usefulness of results from other plugins [no location data in Base], any future issues)

3.4 Final State of the Plugin

What state is the plugin currently in? How well has it achieved its goals? Refer to future work section.

Chapter 4

Case Studies

This chapter will include, for each chosen case study:

4.1 Case Study 1

4.1.1 Program Reasoning

Why choose Case Study 1 to attempt to apply the patch?

4.1.2 Predictions

What results are expected as a result of the patch being applied to Case Study 1?

4.1.3 Patch Code

Either a listing of the patch inlined here, referred to the appendix, or referred to the attached source code or GH. Also needs a commitish for the version of the code to apply the patch to.

4.1.4 Results

What were the results before and after the patch? (number of mallocs, performance, speed, power usage if measurable)

4.1.5 Comparison to Predictions

How did the predictions line up? Better/worse?

4.1.6 Hypothesis

Why were the results what they were?

Chapter 5

Conclusion

This chapter will include high level summaries and conclusions on:

5.1 Results

Did it match the predictions, if not, hypothesize why not

5.2 State of the Plugin

What state is the plugin left in, is there much work to be done, what would be done with more time, what benefits would those changes have/what are the priorities

5.3 Benefit of Further Work

Would further work in this space (not specifically the plugin) be beneficial? If so, why/what should be done first?

Chapter 6

Future Work

A program is never less than 90% complete, and never more than 95% complete.

Terry Baker

Any non-trivial work is never complete. To that end, listed below are some ideas for potential improvements on the *forgetful* plugin or related works on the same principle.

6.1 Detecting Arbitrary Memory Allocations

The current implementation only finds allocations based on uses of `malloc` and `free`. Other ways to allocate memory exist (`calloc`, `realloc`, `alloca`, direct uses of `mmap` and `sbrk`), and platforms that stand to gain the most from this optimisation may have their own implementations.

An extension to this work could involve allowing an arbitrary list of functions declared to allocate or deallocate memory, potentially with fully annotated files specifying their behaviour so that `frama-c` can be used to its full potential (particularly for value analysis, which relies on these specifications).

Alternatively, if there is a willingness to assume a unix-like platform, the depth of analysis could be extended to attempt to automatically determine which functions might allocate memory by searching for `mmap` or `sbrk` calls and propagating annotations indicating functions that directly or indirectly allocate memory.

The approach propagating allocation information already exists in some form in Facebook's Infer [Fac13] static analyser, so future work could also involve extending that platform instead.

6.2 Automatically Performing Fixes

Ideally, fixes would be automatically generated and patched into the code at compile time, avoiding added complexity from the programmer's point of view while still taking advantage of the performance and memory benefits.

Potential intermediate steps toward that goal could involve generation of patches that could be applied to code before compilation, introducing the optimisation. Fortunately, `frama-c` already has code generation capabilities which could be taken advantage of for this purpose.

6.3 Studying Other Architectures

For the sake of practicality and convenience, the analysis presented was only performed on a single x86 machine. To be certain whether these results apply more generally, the analysis should also be performed on other architectures (for example: ARM, x64, embedded systems without address translation/paging, SPARC).

Bibliography

- [Fac13] Facebook. *Facebook Infer*. 2013. URL: <http://fbinfer.com/> (visited on 03/01/2018).
- [Ste17] Daniel Stenberg. *Fewer mallocs in curl*. 2017. URL: <https://daniel.haxx.se/blog/2017/04/22/fewer-mallocs-in-curl/> (visited on 02/01/2018).