# Tokenisation and Compression of Java Class Files for Mobile Devices

Shawn Haggett

Bachelor of Information Technology (Honours)

Flinders University of South Australia

A Thesis Submitted for the Degree of Doctor of Philosophy

Flinders University

School of Computer Science, Engineering and Mathematics

Adelaide, South Australia

2012

(Submitted 18th February 2010)

# Contents

# List of Figures

# List of Tables

# Abstract

An object oriented language, such as Java, needs dynamic binding for method calls since the type of the target object will only be known at runtime. Desktop PCs have sufficient processor and memory resources that dynamic binding is not a significant bottleneck to performance. However, smaller devices such as mobile phones have much more limited resources, requiring efficient implementations. C++ makes use of dispatch tables (also called virtual function tables or just vtables) as a way of speeding up this dispatch. A given method call has an offset (or token value) associated with it, which is used as an index into the target object's vtable. The value stored in the vtable will be a pointer to the C++ function to be executed (similar to a function pointer in C). However, the multiple-inheritance support in C++ complicates this, often requiring a class to have a separate vtable for each super-class it extends.

Java Card (a reduced implementation of Java for smart cards) also uses virtual function tables. While Java does not have full multiple-inheritance as C++ does, it has the notion of an interface. Method calls are divided into two categories in Java, those where the declared type of the object is a class, and those where it is an interface. This allows for a form of multiple-inheritance without having multiple implementations for the same method. There are two different bytecode instructions for these, *invokevirtual* and *invokeinterface* respectively. In Java Card, only the *invokevirtual* instruction can be directly dispatched via the vtable. This leads to simpler vtables, but leaves the *invokeinterface* instruction to use a slower and more complicated dispatch mechanism.

This thesis presents a way to allocate tokens to methods such that both *invokevirtual* and *invokeinterface* can be dispatched via the same vtable and avoids the need of multiple tables as in C++. For tokenisation to succeed, all runtime dependencies must be present,

i.e. the class files for all libraries the application uses. These are needed to determine when methods do and do not need the same token values. An initial tokenisation scheme is presented, where the complete system must be tokenised as a single operation, that is, the application, any libraries it uses and the API. Next, this is extended to allow a new, previously unknown, set of class files to be added to an existing tokenisation (incremental tokenisation). For example, the first tokenisation could include the API and base libraries on a device, followed by a third-party library being added in the second pass and then an application can be added in the third pass. During each tokenisation the previous tokenisation does not need to be modified. This allows a device to be released with a tokenised Java system installed and then new applications can be developed, tokenised and released for that platform. The new application will operate as expected even though the tokeniser had no knowledge of the application at the time it tokenised the initial libraries.

The KVM (Kilobyte Virtual Machine) from Sun Microsystems is a reference virtual machine designed for mobile phone and other portable devices. It is shown that the vtable based dispatch can be between 3 to 45 times faster than the equivalent method dispatch in the KVM. The presence of vtables also removes the need for the symbolic information normally used for linking. The tokenisation concept was also carried further to apply to fields and the *getfield/putfield* and *getstatic/putstatic* instructions used to reference them. This allows for similar speedup and simplification when resolving these references. Further, removing the redundant linking information resulted in class files that were between 42 to 72 percent of their original size.

In mobile devices both processing power and memory can be in short supply. These resources are also limited by the amount of battery power available to run them, as faster processors and larger memories can require more power. This thesis therefore makes a significant contribution towards making Java code both faster to execute and smaller, two vital attributes for a language running on small, portable devices.

"I Shawn Haggett, certify that this thesis does not incorporate without acknowledgment any material previously submitted for a degree or diploma in any university; and that to the best of my knowledge and belief it does not contain any material previously published or written by another person except where due reference is made in the text."

Candidate:

_____

Shawn Haggett

# Acknowledgements

Thanks to Professor Greg Knowles, my supervisor, for keeping me on track and for the guidance when I needed it. Also, thanks to (in no particular order): Graham Bignell, Darius Pfitzner, Tanya Bilka and all the staff at Flinders Uni for your help and support over the years. Finally, thanks to my parents, friends and family for putting up with me.