

Besondere Lernleistung - Dokumentation

Inokentiy Babushkin

November 18, 2014

1 Grundlagen

Die vorliegende Besondere Lernleistung hat das Ziel, ein Softwareprojekt umfassend zu planen und umzusetzen. Dieses Dokument beschreibt Zielsetzung, Aufbau und konkrete Aspekte der Umsetzung, wobei der Hauptaugenmerk auf Algorithmen, Vorgehensweisen und Designentscheidungen liegen wird.

2 Problemantik

Da moderne Software immer komplexer wird, wird es entsprechend schwieriger, in binärer Form vorliegende Softwareprodukte so umzuformen, dass sie ein tiefgehendes Verständnis, als auch die Möglichkeit zur Modifikation erlauben. Aus diesem Grund werden bereits seit einiger Zeit Ansätze unternommen, diesen Vorgang so weit wie möglich zu automatisieren, da die Wiederherstellung von Kontrollfluss, Variablen, Funktionen und weiteren HLL¹-spezifischen Strukturen ohnehin nicht komplett möglich ist, das Ergebnis also bloß eine Annäherung an den Originalcode sein kann, obwohl der Vorgang sehr zeitraubend ist. Diese Programme beruhen auf verschiedenen theoretischen Ansätzen, Programme zu beschreiben und Muster in ihnen zu erkennen, die in ihrer Komplexität stark variieren. Dabei sind mehrer Programmkategorien zu unterscheiden: Decompiler, deren Ziel es ist, ein Programm komplett wiederherzustellen, die meistens eine non-interaktive Arbeitsweise erfordern, und Programme, die bestimmte Aspekte des zu analysierenden Codes wiederherstellen oder eine Art Anleitung für den Benutzer generieren, wie er dabei am besten verfährt. Das Resultat der vor-

liegenden Besonderen Lernleistung ist in die zweite Kategorie einzuordnen, wobei es aus drei Modulen besteht, die unabhängig arbeiten und verschiedene Aspekte des Analysetargets abdecken.

3 Zielsetzung

Wie bereits im letzten Abschnitt erwähnt, handelt es sich bei diesem Projekt um ein Analyseframework zur (teilweisen) Decompilierung von Software, wobei auch eine Vielzahl weiterer Reverse-Engineering²-Aufgaben mit seiner Hilfe vereinfacht werden kann. Dabei wird die Ausgabe von IDA Pro als Grundlage der Analyse verwendet. Dies mag zuerst als Nachteil erscheinen, da es den Arbeitsverlauf des Nutzers einschränkt, wobei bedacht werden muss, dass der Analyse- und Input-Parsing-Code strikt getrennt wird, was es ermöglicht, andere Eingabeformate mithilfe weiterer Module zu akzeptieren, ohne Änderungen in der Programmmechanik umsetzen zu müssen. Die Umsetzung des Moduls für IDA wurde zuerst durchgeführt, da diese keine eigentliche Analyse der Eingabe erfordert, um die für die Mechanik notwendigen Informationen bereitzustellen. Auch ist die Implementierung eines Moduls, welches selbstständig Binärdateien aufbereitet, durchaus denkbar.

3.1 Kontrollflussanalyse

Einer der zentralen Aspekte moderner wie historischer Programme ist die Möglichkeit, Code in Abhängigkeit von Bedingungen ausführen zu können. Diesen wiederherstellen zu können ist eine der Hauptaufgaben dieses Projekts. Dabei ist diese Zielsetzung keineswegs trivial, zumal Strukturen lin-

¹High Level Language, also Programmiersprache mit hohem Abstraktionsgrad

²Der Prozess, Software oder andere technische Systeme in ihrem Aufbau und Verhalten zu analysieren, um sie zu reproduzieren, zu modifizieren oder das dabei gewonnene Wissen weiterzuverwenden, z.B. um Anti-Malware Routinen zu entwickeln

³Strukturen in einem CFG bestehen in der Regel aus zwei oder Mehr Teilen, wobei sogenannte primitive Teile einfache Knoten des Graphen sind. In der Praxis können jedoch auch nich-primitive Teile, also andere Strukturen an dieser Stelle auftreten. Wenn eine Struktur komplett ein Teil ersetzt, wird der Begriff "linear verschachtelt" verwendet. Gleichzeitig ist es möglich, dass Strukturen verschiedener Ebenen den gleichen Knoten verwenden, was als nichtlineare Verschachtelung bezeichnet wird. Beide Begriffe sind vom Autor selbst geprägt.

ear un nicht linear verschachtelt³ sein können etc. Als Grundlage des verwendeten Algorithmus diente die in (1) beschriebene Vorgehensweise.

3.2 Analyse von Variablen und anderen Daten

Daten bilden den zweiten Hauptbestandteil eines Programms, weswegen das Layout der Variablen, Argumente/Parameter, Rückgabewerte von großem Interesse für den Benutzer sein sollte. Gleichzeitig sind Arrays und Variablen nicht immer trivial zu erkennen, weswegen sich einige einfache, dennoch effektive Heuristiken anbieten.

3.3 Behandlung von Integerdivision

Die Integerdivision ist, ebenso wie die Modulo-division, eine verhältnismäßig häufig verwendete Rechenarten in Programmen, bzw. ist meistens von zentraler Bedeutung für die Resultate der Ausführung. Allerdings sind diese Operationen auch sehr ressourcenintensiv, weswegen seit den späten 1990er Jahren die entsprechenden Prozessorinstruktionen kaum noch verwendet werden, zumal die Division durch Konstanten durch Multiplikation und Shifts ersetzt werden kann, was deutlich schnellere Operationen sind⁽²⁾, wodurch die Ausführung gerade in Schleifen beschleunigt werden kann. Gleichzeitig wird dadurch die eigentliche Bedeutung des Codes verschleiert, was zu Hürden bei der Analyse durch Menschen führt. Auch hier ist eine automatisierte Lösung also naheliegend, zudem noch nicht allgemein verbreitet⁴.

4 Umsetzung

Im Folgenden werden die verwendeten Algorithmen im Detail beschrieben.

4.1 Kontrollflussanalyse

Wie bereits in Abschnitt 3.1 beschrieben, ist die Kontrollflussanalyse eine verhältnismäßig komplexe Aufgabe, da das zu untersuhende Objekt eine Vielzahl an Formen annehmen kann, beispielsweise arrangieren Compiler den generierten Assemblercode nicht immer auf eine direkt erwartete Weise. Aus diesem Grund wurde im Laufe der Planung vorliegender Software sehr schnell klar, dass die Analyse des Programms auf Codeebene nicht erfolgre-

ich sein kann, da die Strukturen kaum noch erkannt werden, wenn verschachtelte Strukturen vorliegen. Eine Untersuchung der Fachliteratur ließ erkennen, dass die Analyse des Kontrollflusses deutlich effektiver ist, wenn dieser als Graph betrachtet wird. Per Definition ist ein Kontrollflussgraph (CFG) ein gerichteter Graph G , der über eine Menge Knoten V , eine Menge gerichteter Kanten E und einen Wurzelknoten r mit $r \in V$ verfügt (3).

$$G(V, E, r)$$

Aufgrund der Konzeption von Codeflusssteuerung auf Machinenenebene kann zudem davon ausgegangen werden, dass jeder Knoten 0 bis 2 Nachfolger hat, allerdings eine beliebige Anzahl Vorgänger, wobei nur der Wurzelknoten keinen Vorgänger haben darf. Die Verwendung einer solchen Datenstruktur zur Representation des Programms hat eine Reihe von Vorteilen, wobei die wichtigsten eine deutlich einfachere Analyse möglicher Pfade vom Start- zum Endknoten und die Unabhängigkeit vom Codelayout durch den Compiler im Speicher sind. Die Software geht bei der Analyse folgendermaßen vor: Die Knoten werden nacheinander bearbeitet, wobei die Reihenfolge der Postorder-Traversierung des DFS-Trees entspricht. Ist der aktuelle Knoten der Anfang einer Struktur, die es zu erkennen gilt, werden die dazugehörigen Knoten und Kanten extrahiert und durch einen Knoten ersetzt, der Informationen über die enthaltene Struktur speichert, und welcher als nächstes betrachtet wird. Danach können die Informationen über erkannte Schleifen, Verzweigungen etc. in strukturierter Form dargestellt werden. Das aktuelle Featureset umfasst vor- und nachprüfende Schleifen, Verzweigungen mit beliebig vielen alternativen Pfaden und alle dem Autor bekannten Möglichkeiten der Verschachtelung. Einzige Einschränkungen bilden derzeit komplexere Bedingungen, die, falls sie mehrere logische *OR* Ausdrücke enthalten, aufgrund nicht-linearer Verschachtelung nicht richtig erkannt werden, und Sprungbefehle wie *goto*, *break*, *continue* und in manchen Fällen *return*. Auch dürfen Schleifen derzeit nur singuläre⁵ Bedingungen haben, wobei diese Einschränkungen im Laufe des Projekts beseitigt werden sollen. Für eine detaillierte Beschreibung des theoretischen Hintergrunds des verwendeten Algorithmus sei an dieser Stelle auf (1) und den Quellcode verwiesen,

⁴Eine Recherche des Autors ergab, dass z.B. der Hex-Rays Decompiler "ohne Probleme" mit solchen Konstrukten umgehen kann, da dieser allerdings ein kommerzielles und noch dazu recht teures Produkt darstellt, kann vermutlich nicht von allgemeiner Verbreitung gesprochen werden.

⁵atomare

⁶Objektorientierte Programmierung

da dieser den Konventionen der OOP⁶ angepasst wurde, um Modularität und Effizienz zu maximieren und die Änderungen in Hochsprachen seit der Veröffentlichung von (1) zu reflektieren.

4.2 Datenanalyse

Bei diesem Modul liegt der Hauptaugenmerk auf der Möglichkeit, nicht nur adressierte Speicherbereiche aus einer Analyse des Codes zu extrahieren, wie es z.B. in IDA Pro geschieht, sondern auch aus dem Kontext der Speicheraufrufe auf die Funktion der entsprechenden Speicherabschnitte zu schließen. Dabei sei erwähnt, dass solche Verfahren grundsätzlich nur Heuristiken darstellen, da nicht jeder Speicherabschnitt, aus dem z.B. ein DWORD⁷ gelesen wird, automatisch einer Integervariable gleichgesetzt werden kann. So verwenden moderne Compiler die push⁸- und pop⁹-Instruktionen nicht mehr für Funktionsaufrufe, da diese zu viele Ressourcen verbrauchen. Stattdessen wird am Anfang der Funktion mehr Speicher alloziert als für die lokalen Variablen notwendig, so dass Parameter direkt in den Speicher geschrieben werden, bevor eine Funktion aufgerufen wird. Darüber hinaus werden vermehrt sogenannte Canaries auf dem Stack hinterlegt, die sicherstellen sollen, dass kein Stackoverflow stattfindet. Dabei handelt sich um "freien" Speicher, der zwischen den lokalen Variablen und dem für Parameter reservierten Platz alloziert wird und auf seine Integrität überprüft wird, wenn die Funktion zurückkehrt. Auf diese Weise wird es viel schwieriger, die Rückkehradresse durch Nutzereingaben zu verändern. Solche Canary-Werte können auch zwischen Arrays und anderen Variablen eingefügt werden, was z.B. die

Längenbestimmung von Arrays erschwert, zumal bei diesen selten direkte Addressierungen der Elemente stattfinden. Diese und andere Aspekte der Kompilierung, bei denen Informationen verlorengehen, erschweren die Analyse und Wiederherstellung mitunter sehr, was gleichermaßen für Menschen und Programme gilt, wobei letztere weniger trugschlussanfällig zu sein scheinen, da sie das Analyseergebnis schrittweise aufbauen, so dass mitunter Informationen nur unvollständig sind, jedoch nur selten objektiv falsch.

4.3 Optimierte Intergerdivision

Die derzeitige Umsetzung dieses Moduls basiert auf praktischen Erfahrungen des Autors und von Nutzern der Reverse-Engineering-Stackexchange Website(4).

Bibliotheksreferenz

- [1] M. Sharir, "Structural analysis: A new approach to flow analysis in optimizing compilers,"
- [2] T. Granlund and P. L. Montgomery, "Division by invariant integers using multiplication,"
- [3] <http://de.wikipedia.org/wiki/Kontrollflussgraph>. Eingesehen am 4.11.2014.
- [4] <http://reverseengineering.stackexchange.com/questions/1397/how-can-i-reverse-optimized-integer-division-modulo-by-constant-operations>. Eingesehen am 15.11.2014.

⁷4 byte (32 bit) großer Speicherabschnitt, meist ein Integer

⁸push alloziert Speicher auf dem Stack und speichert dort den übergebenen Operanden

⁹pop entfernt den zuletzt auf den Stack geschobenen Wert und speichert ihn im übergebenen Operanden