

Entkopplung der Z3 Komponente in ProB mit ZeroMQ

Bachelorarbeit

vorgelegt von

Silas Alexander Kraume

22. Januar 2025

im Studiengang Informatik
zur Erlangung des akademischen Grades

Bachelor of Science (B.Sc.)

Erstgutachter: Prof. Dr. Michael Leuschel
Zweitgutachter: Dr. C. Bolz-Tereick

Selbstständigkeitserklärung

Hiermit versichere ich, die vorliegende Bachelorarbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt zu haben. Alle Stellen, die aus den Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Düsseldorf, den 22. Januar 2025

Silas Alexander Kraume

Zusammenfassung

Fassen Sie hier die Fragestellung, Motivation und Ergebnisse Ihrer Arbeit in wenigen Worten zusammen.

Die Zusammenfassung sollte den Umfang einer Seite nicht überschreiten.

Macht man am Ende ...

Danksagung

Me, Myself and I!

Inhaltsverzeichnis

Tabellenverzeichnis	xi
Abbildungsverzeichnis	xi
Algorithmenverzeichnis	xi
Quellcodeverzeichnis	xi
1 Einführung	1
1.1 Motivation	1
1.2 Ziele	2
2 Grundlagen	3
2.1 ProB	3
2.2 Z3 Solver	4
2.3 ZeroMQ	6
3 Architekturänderung	7
3.1 Planung	8
3.1.1 Prolog Datentypen	8
3.1.2 Struktur der Nachrichten	8
3.1.3 Server Struktur	8
3.2 Implementierung	8
3.2.1 Interfacefunktionen	8
3.2.2 Hilfsfunktionen	8
3.2.3 Optimierungen	8
3.2.4 Serveranbindung	10
3.2.5 Logging	10
4 Exceptions	10
4.1 Kontrollfluss	10
5 Zusätzliche Ergebnisse	10
5.1 Softlock	10
5.2 Versionsinkompatibilität	10

6	Leistungsbewertung	10
6.1	Performance-Overhead	11
7	Zukünftige Arbeiten	17
7.1	Weitere Analyse und Optimierungen	17
7.2	Deinit Hook	17
7.3	Parallelisierung	17
8	Konklusion	18
	Literatur	19

Tabellenverzeichnis

1	Auszug der Daten der Performance-Messung.	11
2	Ausschnitt der gesammelten Messwerte eines Ausreißers.	16

Abbildungsverzeichnis

1	Die geplante Architekturänderung (Die Komponenten-Entkopplung ist in Rot markiert. Die gestrichelten Boxen zeigen die verschiedenen Prozesse an.)	2
2	Durchschnittliche Laufzeiten der Anfragen in den verschiedenen Architekturen.	12
3	Induzierter Overhead der neuen Server-Architektur durch das Serialisieren auf den Socket.	15

Algorithmenverzeichnis

Quellcodeverzeichnis

1	Ein Beispiel zur Verwendung des Z3-Solvers innerhalb der ProB REPL. . . .	6
2	Ein minimales Client Modell mit ZeroMQ.	7
3	Ein minimales Server Modell mit ZeroMQ.	7
4	Ein Ausschnitt einer Interfacefunktion zur Demonstration von Schleifenoptimierung.	9
5	Ein Ausschnitt einer redundanten Objektverwaltung.	9
6	Die Hilfsfunktion <code>escape_string</code>	9

1 Einführung

Die digitale Transformation hat unsere Welt grundlegend verändert und macht Software-Systeme zu einem unverzichtbaren Bestandteil des täglichen Lebens. Von sicherheitskritischen Anwendungen wie der Steuerung autonomer Fahrzeuge bis hin zu Finanzsystemen und medizinischen Geräten sind wir zunehmend auf Software angewiesen, die zuverlässig und fehlerfrei funktioniert. Die Gewährleistung von Korrektheit und Stabilität ist jedoch eine anspruchsvolle Aufgabe insbesondere angesichts der Komplexität moderner Systeme. Ein zentraler Baustein zur Bewältigung dieser Herausforderung ist der Einsatz von Modellierungs- und Verifikationswerkzeugen. Diese ermöglichen es, komplexe Systeme systematisch zu analysieren und sicherzustellen, dass sie den gewünschten Spezifikationen entsprechen. Besonders hervorzuheben ist der Einsatz von SMT¹-Sovern, die sich als leistungsfähige Werkzeuge etabliert haben, um schwierige logische Probleme effizient zu lösen. SMT-Solver wie Z3 bieten durch ihre Fähigkeit zur präzisen und schnellen Verarbeitung logischer Ausdrücke eine wertvolle Unterstützung bei der Verifikation und Validierung. Ein prominentes Beispiel für die Integration eines solchen Solvers ist die Software ProB. Der Animator, Constraint-Solver und Model-Checker ProB nutzt den SMT-Solver Z3, um formale Modelle effizient zu analysieren und zu überprüfen. Dies macht die Software zu einer wichtigen Instanz in der Welt der formalen Methoden, insbesondere im Kontext von Modellierungs- und Verifikationsaufgaben.

1.1 Motivation

Innerhalb von ProB birgt der Einsatz des Z3-Solvers jedoch auch Herausforderungen, die die Effizienz und Zuverlässigkeit der Anwendung beeinträchtigen können. Ein bekanntes Problem besteht in dem sporadischen Auftreten von Speicherlecks und Segmentation Faults, die sowohl die Stabilität als auch die Nutzbarkeit von ProB's Z3-Interface negativ beeinflussen. Diese technischen Mängel erschweren nicht nur die Durchführung formaler Verifikationen, sondern können auch zu einer zeitraubenden Verwendung der Z3-Solver Komponente sowie Unterbrechung von Arbeitsprozessen führen.

Ein weiterer Mangel liegt in der aktuellen sequenziellen Lösung mehrerer Prädikate. Dieser Ansatz, bei dem die Prädikate nacheinander gelöst werden, ist in seiner Natur ressourcenintensiv und zeitaufwendig. Angesichts der steigenden Komplexität formaler Modelle und der wachsenden Nachfrage nach schnellerer Verifikation wird die Limitierung durch die sequenzielle Verarbeitung immer offensichtlicher. Eine Parallelisierung der Lösung von Prädikaten könnte hier erhebliche Leistungsverbesserungen bringen, indem moderne Mehrkernarchitekturen effizienter ausgenutzt werden, um den Anforderungen der Nutzer und der immer komplexer werdenden Modelle gerecht zu werden.

Die Kombination dieser Herausforderungen (sporadische technische Instabilitäten und

¹Satisfiability Modulo Theories

begrenzte Effizienz durch sequenzielle Verarbeitung) macht es notwendig, alternative Ansätze oder Verbesserungen für die Integration des Z3-Solvers in ProB zu erforschen und bildet die Grundlage und Motivation für die vorliegende Arbeit.

1.2 Ziele

Das Hauptziel dieser Arbeit ist es, die Integration des Z3-Solvers in ProB zu verbessern, indem die bestehende Vorgehensweise, die Prädikate direkt im Z3-Interface von ProB zu lösen, verworfen wird. Stattdessen wird eine neue Architektur vorgeschlagen und implementiert, welche eine vollständige Entkopplung der Z3-Solver-Komponente von ProB vorsieht. Hierzu wird also der Z3-Solver in einen eigenständigen, separaten Prozess ausgelagert, wodurch ein System eingeführt wird, bei dem ProB und der Z3-Solver als zwei unabhängige Prozesse agieren, die über eine Kommunikationsschnittstelle miteinander verbunden sind. Prädikate werden hierdurch innerhalb des Z3-Interfaces an den Z3-Solver gesendet, wo diese gelöst und zurückgeschickt werden. Diese geplante Architekturänderung ist in der folgenden Abbildung 1 visualisiert.

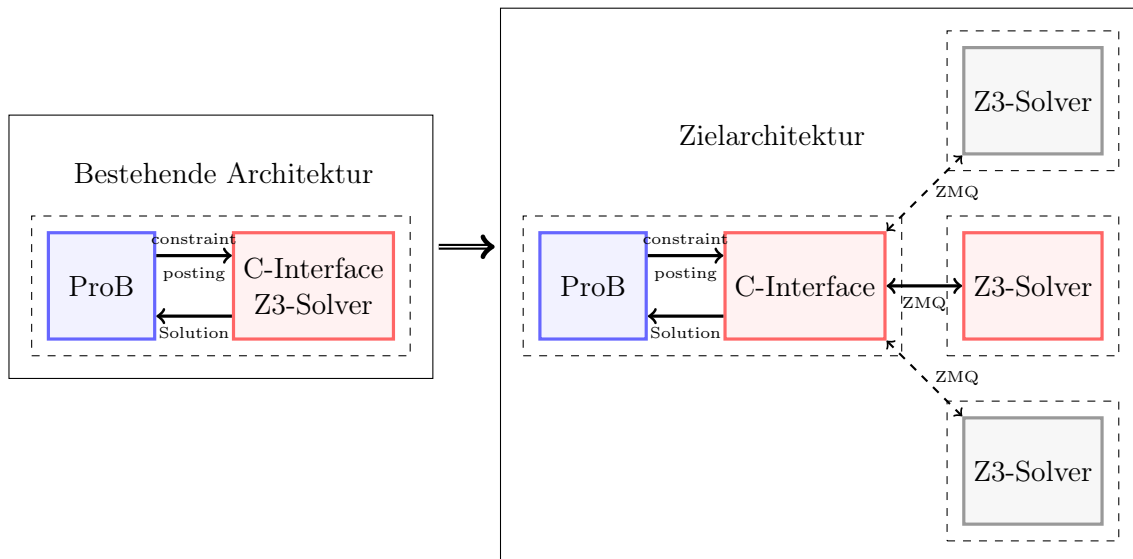


Abbildung 1: Die geplante Architekturänderung (Die Komponenten-Entkopplung ist in Rot markiert. Die gestrichelten Boxen zeigen die verschiedenen Prozesse an.)

Diese Arbeit wird einerseits mit dem Interesse der Erweiterbarkeit verrichtet, sodass zukünftig die Option besteht, gegebenenfalls mehrere Instanzen des Z3-Prozesses zu starten und das Lösen der Prädikate zu parallelisieren. Andererseits dient die Entkopplung selbst bereits zur Verbesserung der Stabilität und Zuverlässigkeit von ProB, da bei eventuellen Fehlern im Z3-Solver-Prozess dieser unabhängig von ProB neu gestartet werden kann, was zu einem robusteren Gesamtsystem führt.

Die genauen Technologien und Konzepte, die hierfür zum Einsatz kommen und Relevanz zeigen, sowie ihre Funktionsweise und Vorteile werden im folgenden Kapitel detailliert erläutert. Daraufhin wird die Planung und Implementierung der neuen Architektur beschrieben und die Leistungsfähigkeit der vorgenommenen Entkopplung anhand von Benchmarks und Tests hinsichtlich des Vergleichs zur vorherigen Systemstruktur evaluiert. Zuletzt wird auf zukünftige Erweiterungen und Verbesserungen eingegangen und eine abschließende Konklusion genannt.

2 Grundlagen

Zur Förderung eines einheitlichen Verständnisses werden in diesem Abschnitt zunächst die erforderlichen Hintergrundinformationen illustriert. Im Folgenden werden die drei zentralen Konzepte behandelt, die für das Verständnis dieser Arbeit von Bedeutung sind: ProB, Z3 und ZeroMQ.

2.1 ProB

Die B-Methode, entwickelt von J.-R. Abrial [Abr96], ist eine formale Methode zur Entwicklung von Softwaresystemen, die auf der Idee der abstrakten Maschinen basiert. Mit abstrakten Maschinen lassen sich Zustände und deren Veränderungen mithilfe mathematischer Konzepte wie Mengen, Relationen und Funktionen modellieren [LB03]. Durch sogenannte Verfeinerungen wird schrittweise von einer abstrakten Beschreibung zu einer konkreten Implementierung übergegangen. Dabei stellt die Methode sicher, dass Invarianten stets eingehalten werden, um die Korrektheit des Systems zu garantieren.

Die an der HHU am Lehrstuhl der Softwaretechnik und Programmiersprachen entwickelte Software ProB [LB03] ist ein Validierungstoolset für Modelle der B-Methode. Als solcher unterstützt ProB mitunter die Modellierung, Animation und Verifikation von B-Modellen, indem Funktionalitäten wie Consistency Checking und Constraint Solving bereitgestellt werden. ProB findet bereits in vielen Systemen Verwendung zur Datenvalidierung und Validierung komplexer Eigenschaften für sicherheitskritische Systeme. Es wird bereits von mehreren Unternehmen eingesetzt und ist November 2022 mit dem „AlainColmerauer Prize“ ausgezeichnet worden.

Der Animator in ProB ermöglicht es, formale Spezifikationen zu visualisieren und zu animieren. Nutzer können durch die Simulation in Echtzeit einen Einblick in die Zustandsübergänge einer Maschine erhalten und schrittweise die Veränderungen nachvollziehen. Der aktuelle Zustand der Maschine wird dabei in einer grafischen Benutzeroberfläche dargestellt.

Ein weiterer Kernbestandteil von ProB ist das Consistency Checking, welches in zwei Ansätzen realisiert wird: Temporal Model Checking und Constraint-Based Checking.

Relevanz für die Arbeit suchen? Wo ist sie? Man weiß es nicht. Aber irgendwas muss man ja schreiben...

Inhaltliche Korrektheit prüfen. Nach A kommt B und was ist das?

Korrekte Referenzierung von Quellen überprüfen.

Beim Temporal Model Checking wird versucht, eine Sequenz von Operationen zu finden, die ausgehend von einem Anfangszustand zu einer Verletzung der Invariante oder einem anderen Fehler führt. Im Gegensatz dazu fokussiert sich das Constraint-based Checking auf die Suche nach einem Zustand des Systems, der die Invariante noch erfüllt. Von dort aus wird geprüft, ob es eine einzelne Operation gibt, welche die Invariante verletzt oder anderweitige Fehler erzeugt.

Während das Model Checking eine umfassende Exploration aller Zustände ermöglicht, ist das Constraint-based Checking spezifischer, da es nur auf Fehler bei einzelnen Operationen fokussiert ist. Zusammen lassen sich so vollständige Fehler und problematische Operationen identifizieren.

Beide Ansätze bieten wertvolle Instrumente für die Konsistenzprüfung von B-Modellen, und sind in der Lage, die Verletzung von Invarianten und daraus folgenden Bedingungen sowie die Abwesenheit von Deadlocks und das Erreichen von spezifizierten Zielpredikaten zu überprüfen [LB08].

Zuletzt bietet ProB auch eine Constraint-Solving-Funktionalität, die es ermöglicht, unter Berücksichtigung von gegebenen Constraints (Einschränkungen) Lösungen für spezifische Prädikate zu finden. Derartige Einschränkungen oder Bedingungen können in Form von logischen Ausdrücken oder Gleichungen gegeben sein, die es zu erfüllen gilt. Ein Constraint-Solver ist ein Algorithmus oder System, welches darauf abzielt, unter Berücksichtigung eben jener Bedingungen eine Belegung aller Variablen zu finden, die die gegebenen Prädikate erfüllt und somit ein Problem auf dessen Erfüllbarkeit zu prüfen. ProB implementiert hierfür verschiedene Constraint-Solving-Strategien, die auf unterschiedlichen Algorithmen basieren und es ermöglichen, Prädikate effizient zu lösen. Einerseits wird CLP(FD)² verwendet, um auf endlichen Domänen beispielsweise Gleichheits- und Ungleichheitsbedingungen sowie arithmetische Relationen zu lösen. Ein weiterer Ansatz ist die Integration des SAT-basierten Kodkod, einem effizienten Constraint-Solver für die Prädikatenlogik erster Ordnung mit Relationen, transitiven Hüllen, Bit-Vektor-Arithmetik und partiellen Modellen [TJ07]. Zuletzt wird auch der SMT-Solver Z3 in ProB integriert, um komplexere Prädikate zu lösen, die über simple boolesche und arithmetische Ausdrücke hinausgehen.

ProB ist im Kern in SICStus Prolog [CWA⁺88] implementiert, bietet jedoch verschiedene Programmiererweiterungen, welche zumeist in C oder C++ geschrieben sind. Einer dieser Erweiterungen ist das Z3-Interface, welches die Integration des Z3-Solvers in ProB ermöglicht.

2.2 Z3 Solver

Das Satisfiability Modulo Theories (SMT)-Problem lässt sich als eine natürliche Erweiterung des klassischen Boolean Satisfiability (SAT)-Problems verstehen. Während SAT sich ausschließlich mit der Erfüllbarkeit von booleschen Formeln beschäftigt, umfasst SMT

²Constraint Logic Programming over Finite Domains

zusätzliche Hintergrundtheorien wie Arithmetik, Bit-Vektoren, Arrays und uninterpretierten Funktionen. Diese Erweiterung ermöglicht eine umfassendere Form der logischen Schlussfolgerung.

Der Z3-Solver [MB08] wurde von Microsoft Research entwickelt. Seit seiner ersten Veröffentlichung hat sich Z3 zu einem der leistungsfähigsten und am weitesten verbreiteten SMT-Solver behauptet. Microsoft hat Z3 als Open-Source-Software frei zugänglich unter der MIT-Lizenz veröffentlicht, was seine Verbreitung und Nutzung in verschiedenen Bereichen weiter gefördert hat.

Er wurde in C++ implementiert, bietet jedoch eine Vielzahl von externen API-Anbindungen, die es ermöglichen, den Solver in verschiedenen Programmiersprachen zu verwenden, wie beispielsweise OCaml, Python, Ruby und Rust. Er ist ein leistungsfähiger SMT-Solver, der sich speziell auf die Lösung von Problemen im Bereich der Softwareverifikation und -analyse spezialisiert. Zum Beispiel wird der Z3-Solver in der Softwareverifikation eingesetzt, um die Korrektheit von Programmen zu überprüfen, indem er formale Spezifikationen testet und beweist. Andere Einsatzbereiche sind die automatische Generierung von Testfällen basierend auf formalen Modellen sowie dem Modellieren von Entscheidungsproblemen und dem Abstrahieren von Prädikaten.

Der Aufbau von Z3 ist modular und basiert auf dem DPLL³(T)-Framework, das die grundlegenden Prinzipien des DPLL-Algorithmus zur Lösung des CNF-SAT-Problems um weitere Theorien (wie SMT) ergänzt. Der DPLL ist ein Algorithmus zur Entscheidung der Erfüllbarkeit einer booleschen Aussagenlogikformel in konjunktiver Normalform (CNF). Er arbeitet rekursiv und kombiniert folgende Schritte:

1. Unit Propagation: Ableitung neuer Wahrheitswerte aus Klauseln mit nur einem Literal (Einheitsklausel).
2. Decision Making: Entscheidung für ein Literal, um die Suchrichtung voranzutreiben.
3. Conflict Detection/Backtracking: Zurücksetzen von Entscheidungen, wenn ein Konflikt entdeckt wird.

hier mache ich bestimmt noch den eigentlichen algorithmus hin, weil ich prinzipiell einfach ein algorithmenverzeichnis haben möchte xD

In Z3 dient der DPLL-Algorithmus als Kern des SAT-Solvers. Weitere Theorien wie lineare und nicht lineare Arithmetik, Bit-Vektoren und Arrays werden durch weitere Theorie-Solver behandelt, welche als spezialisierte Module eng mit dem SAT-Solver integriert sind.

Ein zentrales Merkmal von Z3 ist seine hohe Effizienz und seine Fähigkeit, große und komplexe Probleme in akzeptabler Zeit zu lösen. Diese Leistungsfähigkeit wird durch verschiedene Optimierungen und Heuristiken erreicht, die den Suchraum effizient einschränken und die Lösung von Constraints beschleunigen. Zudem ist ein Simplifier in die System Architektur integriert, der gegebene Constraints vereinfacht. Z3 hat nur wenige Abhängigkeiten an externe Bibliotheken und ist daher leicht zu integrieren und zu verwenden [MR15]. Es

³Davis-Putnam-Logemann-Loveland

Quellcode 1: Ein Beispiel zur Verwendung des Z3-Solvers innerhalb der ProB REPL.

```
1: >>> :z3 X<Y & Y<X & X: INTEGER
2: PREDICATE is FALSE
```

wird die C++ Bibliothek für das Threaden der Anwendung verwendet, um die parallele Verarbeitung zu ermöglichen und damit die Leistungsfähigkeit weiter zu steigern.

Innerhalb von ProB wird Z3 vor allem dann eingesetzt, wenn ProBs interner Constraint-Solver bei der Lösung bestimmter Probleme an seine Grenzen stößt [KL16]. Dies betrifft insbesondere große oder komplexe Constraints mit nicht-linearen Bedingungen. Beispielsweise lässt sich die folgende Formel, welche in Quellcode 1 gezeigt wird, in Z3 aber nicht in ProB lösen:

Insgesamt macht die Leistungsfähigkeit und Flexibilität des Z3-Solvers ihn zu einem wichtigen Bestandteil von ProB, um Verifikationsaufgaben zu bewältigen, die über die Fähigkeiten des internen Constraint-Solvers hinausgehen.

2.3 ZeroMQ

In der geplanten Architektur wird als Technologie zur Kommunikation zwischen den zwei separaten Prozessen die ZeroMQ-Bibliothek [Hin13] verwendet. ZeroMQ⁴ ist eine hochleistungsfähige, asynchrone Nachrichtenaustauschbibliothek, die speziell für verteilte Systeme entwickelt wurde. Da ZeroMQ in der Sprache C implementiert- und ursprünglich für die Börse geschrieben wurde, lag extreme Performanceoptimierung lange Zeit im Fokus [S⁺15], was es zu einer der schnellsten und effizientesten Bibliotheken für den Nachrichtenaustausch macht.

ZeroMQ bietet eine Vielzahl von Kommunikationsmustern, die es ermöglichen, verschiedene Arten von verteilten Systemen zu realisieren. Dazu gehören unter anderem:

- Request-Reply (REQ/REP): Ein einfacher Nachrichtenaustausch, der für synchrone Kommunikation geeignet ist.
- Publish-Subscribe (PUB/SUB): Ermöglicht die Verteilung von Nachrichten an mehrere Empfänger, die sich für bestimmte Themen registrieren.
- Push-Pull (PUSH/PULL): Ein Muster für Lastverteilung, bei dem Nachrichten an Worker-Threads oder Prozesse verteilt werden.
- Dealer-Router (DEALER/ROUTER): Ein erweiterbares und flexibles Muster für komplexere (und asynchrone) Kommunikationsstrukturen.

⁴auch ØMQ, 0MQ oder ZMQ

Quellcode 2: Ein minimales Client Modell mit ZeroMQ.

```

1:  s = socket(REQ);
2:  s.connect("ipc:///tmp/zmq");
3:  s.send("Hello");
4:  reply = s.recv();

```

Quellcode 3: Ein minimales Server Modell mit ZeroMQ.

```

1:  s = socket(REP);
2:  s.connect("ipc:///tmp/zmq");
3:  request = s.recv();
4:  s.send("Hello back");

```

Im Kern stellt ZeroMQ ein API über traditionelle Sockets bereit, welche die Komplexität des Netzwerkprotokollmanagements abstrahiert. Statt sich mit niedrigstufigen Details wie Verbindungen, Paketverwaltung und Fehlerbehebung auseinanderzusetzen, kann so die Implementierung von Anwendungslogik drastisch vereinfacht werden. Derartige API Anbindungen sind für viele Programmiersprachen verfügbar, darunter C++, Python, Java, Go und viele mehr.

ZeroMQ ist flexibel und einfach skalierbar, da es ohne einen dedizierten Message Broker auskommt und direkt zwischen den Prozessen kommuniziert. Es werden verschiedene Transportprotokolle unterstützt, darunter TCP (Transmission Control Protocol), IPC (Inter-Process Communication), (E)PGM ((Encapsuled) Pragmatic General Multicast) und viele mehr, die es ermöglichen, ZeroMQ in verschiedenen Umgebungen zu verwenden.

Insgesamt bietet ZeroMQ also eine hohe Geschwindigkeit, Flexibilität und Skalierbarkeit sowie ein simples Programmierinterface, was es zu einer idealen Wahl für die Kommunikation zwischen ProB und dem Z3-Solver macht. Zu Illustrationszwecken ist in Quellcode 2 und Quellcode 3 beispielhaft ein minimaler Client und Server in C++ mit ZeroMQ dargestellt.

Für die geplante Architekturänderung wird ZeroMQ verwendet, um eine direkte IPC-Protokollverbindung zwischen ProB und dem Z3-Solver herzustellen. Da der entkoppelte Solver auf demselben Rechner ausgeführt wird, bietet sich IPC als Kommunikationsprotokoll an, um eine effiziente und schnelle Kommunikation zu gewährleisten. Es wurde sich für das Request-Reply Muster entschieden, da es im Kontext die am besten passende Lösung darstellt, ohne dem System unnötige Komplexität hinzuzufügen.

3 Architekturänderung

Eigentlicher Kern der Arbeit, mache ich, sobald ich mit Prokrastination fertig bin.

3.1 Planung

3.1.1 Prolog Datentypen

- atoms string - integers longs (laut dokumentation bis version blabla) - floats doubles - typerefs problem, weil kann alles sein

3.1.2 Struktur der Nachrichten

- function identifier - status identifier - message

3.1.3 Server Struktur

long damn switch case threading

3.2 Implementierung

3.2.1 Interfacefunktionen

porting of all 53 interface function

3.2.2 Hilfsfunktionen

insbesondere *mktype* statemachines

3.2.3 Optimierungen

Nach der erfolgreichen Portierung der Schnittstelle und der Implementierung aller notwendigen Funktionen lassen sich zusätzlich kleine Optimierungen und Verbesserungen vornehmen, um Effizienz, Lesbarkeit und Wartbarkeit des Codes zu erhöhen.

Eine der elementarsten Optimierungsmöglichkeiten ist die Vermeidung von Nachrichtenaustausch beider Prozesse innerhalb von Schleifen. Dieses Verhalten tritt insbesondere dann auf, wenn Prolog Datenstrukturen übermittelt werden, die iterierbar sind, wie Listen, Vektoren und Records. In diesen Fällen wird für jedes Element der Struktur eine Nachricht an den Server gesendet, um das aktuelle Element zu übermitteln. Dieses Verhalten kann in einzelnen Fällen durch die Übermittlung der gesamten Struktur in einer einzigen Nachricht vermieden werden. In dem Folgenden Quellcode [4](#) ist gezeigt, wie eine Prolog Liste zunächst

Quellcode 4: Ein Ausschnitt einer Interfacefunktion zur Demonstration von Schleifenoptimierung.

```

1:  std::vector<std::string> string_vec =
2:      term_ref_prolog_list_to_string_vector(element_names);
3:  for (int i = 0; i < cardinality; i++) {
4:      zmsg_addstr(request, string_vec[i].c_str());
5:  }
6:  // send request to server

```

Quellcode 5: Ein Ausschnitt einer redundanten Objektverwaltung.

```

1:  // function: mk_op_comprehension_set_multi
2:  ContextData ctx_dta =
3:      get_translation_representant_ctx_data(translation_type_atom);
4:  prolog_type_list_to_sort_vector(ctx_data, couple_types);
5:  // ...
6:
7:  // function: prolog_type_list_to_sort_vector
8:  mk_sort(ctx_data->get_translation_type_atom());

```

in einem String Vektor akkumuliert wird, um anschließend in einer einzigen Nachricht an den Server gesendet zu werden.

Eine weitere Optimierungsmöglichkeit ist die Vermeidung von unnötigen Nachrichten an den Server. Beispielsweise illustriert Quellcode 5 vermeidbare Komplexität durch die Verwendung unnötiger Datenobjekte. Anhand der Variablen *translation_type_atom* wird das Objekt *ctx_data* ermittelt, welches an die Funktion *prolog_type_list_to_sort_vector* übergeben wird. Diese verwendet das Objekt ausschließlich zur Ermittlung der Variablen *translation_type_atom*. Das Problem hierbei ist, dass *ctx_data* in der neuen Architektur auf der Seite des Server-Prozesses liegt und somit einen Nachrichtenaustausch erfordert. Das Problem lässt sich umgehen, indem die Funktion *prolog_type_list_to_sort_vector* dahingehen refaktorisiert wird, direkt mit *translation_type_atom* aufgerufen zu werden.

Zuletzt wurde das DRY⁵-Prinzip angewendet, um die Wartbarkeit des Codes zu erhöhen, indem gewisse Hilfsfunktionen ausschließlich auf entweder der Prolog- oder der Serverseite implementiert wurden. Der ursprüngliche Kontrollfluss verlangte zum Beispiel die Imple-

⁵Don't Repeat Yourself

Quellcode 6: Die Hilfsfunktion *escape_string*.

```

1:  std::string escape_string(const std::string to_escape) {
2:      return "|" + to_escape + "|";
3:  }

```

mentierung von der in Quellcode 6 gezeigten Funktion in beiden Prozessen. Durch die Anwendung des Programmierprinzips wurde die Funktion ausschließlich an dem Z3-Prozess implementiert.

3.2.4 Serveranbindung

server als subprocess starting as needed

3.2.5 Logging

via sys argv stdout not captureable in sicstus prolog

4 Exceptions

very important. need to be handled properly. need to work with server rep req structure
what if multiple errors? need to notify other process had segfaults

Big Important!

4.1 Kontrollfluss

dependency graph of all functions either throw or catch exceptions

5 Zusätzliche Ergebnisse

5.1 Softlock

endless loop fixed by interrupting on every reset

Wichtig, wird in Performanceanalyse verwiesen.

5.2 Versionsinkompatibilität

makefile hell glibc (OS) incompatible with zlib.so -> darwin12

6 Leistungsbewertung

Nach Abschluss der durchgeführten Architekturänderung ist es elementar, die Auswirkungen auf die Laufzeitperformance zu bewerten. Da sich die grundlegende Funktionsweise des Z3-

Solvers in der neuen Architektur nicht geändert hat, ist es zu erwarten, dass die Laufzeitperformance der ProB-Systemerweiterung durch die Einführung der ZeroMQ-Kommunikation negativ beeinflusst wird. Zusätzlich zu den Aufrufen des Z3-Interfaces müssen zur Lösung eines einzelnen Prädikates nun mehrere Anfragen und Antworten über den Socket serialisiert werden. Diese zusätzliche Kommunikation führt zu einem Performance-Overhead, welcher quantifiziert und evaluiert werden muss. Ebenfalls besteht ein Interesse zum Vergleich verschiedener ZeroMQ-Protokolle und deren Auswirkungen auf die Performance. Hierbei ist zu erwarten, dass das Inter-Process-Communication (IPC) Protokoll schneller ist als das Transmission-Control-Protocol (TCP) Protokoll, da es auf dem gleichen Rechner arbeitet und keine Netzwerkkommunikation benötigt, sondern das Dateisystem verwendet. Im nachfolgenden Abschnitt werden die Methodik der Leistungsbewertung, die erzielten Ergebnisse und ihre Interpretation detailliert beschrieben.

6.1 Performance-Overhead

Um eine Bewertung des Performance-Overheads zu ermöglichen, müssen zunächst empirische Daten erhoben werden. Hierzu werden die Tests zur Verifikation der Funktionsweise des Z3-Interfaces umfunktioniert, um die Laufzeit der einzelnen Anfragen zu messen. Im Code des Z3-Interfaces wird ein Zeitstempel bei Beginn und Ende des Lösungsvorgangs eines Prädikates gesetzt, dessen Differenz berechnet und gespeichert. Insgesamt stehen 53 Tests zur Verfügung, die in der Testumgebung des Z3-Solvers ausgeführt werden können. Diese Tests umfassen insgesamt 679 Prädikate, welche eine ausreichende Grundgesamtheit zur Bewertung der Performance bieten. Ebenfalls wird die Anzahl der Anfragen und Antworten, die über das Netzwerk gesendet werden, gemessen. Ein kleiner Auszug dieser Messdaten sind in Tabelle 1 dargestellt.

Tabelle 1: Auszug der Daten der Performance-Messung.

TestID	QueryID	Old(ns)	New(IPC, ns)	New(TCP, ns)	Req. Count
⋮	⋮	⋮	⋮	⋮	⋮
1510	1	114 815 014	283 392 117	113 503 997	191
1510	2	52 048 678	59 273 375	44 489 012	166
1510	3	30 853 103	24 983 820	25 212 332	286
1511	1	69 240 686	199 325 313	61 823 314	21
1511	2	61 404 523	62 220 124	48 522 297	20
1513	1	80 829 324	202 694 845	75 877 790	25
⋮	⋮	⋮	⋮	⋮	⋮

Die Zeitmessungen sind in Nanosekunden (ns) angegeben und zeigen die Laufzeit der Anfragen in den verschiedenen Konfigurationen. Innerhalb der gegebenen Größenordnung sind die Werte ebenfalls in Millisekunden (ms) zu interpretieren und damit ausschlaggebend für die Performance. Innerhalb der eigentlichen Daten wurden die Messungen mehrfach unabhängig voneinander wiederholt, um eine statistische Aussagekraft zu gewährleisten. Da

dennoch von Messfehlern und Schwankungen auszugehen ist, werden die Daten in einem statistischen Kontext betrachtet. Es wird angenommen, dass die Messfehler normalverteilt sind.

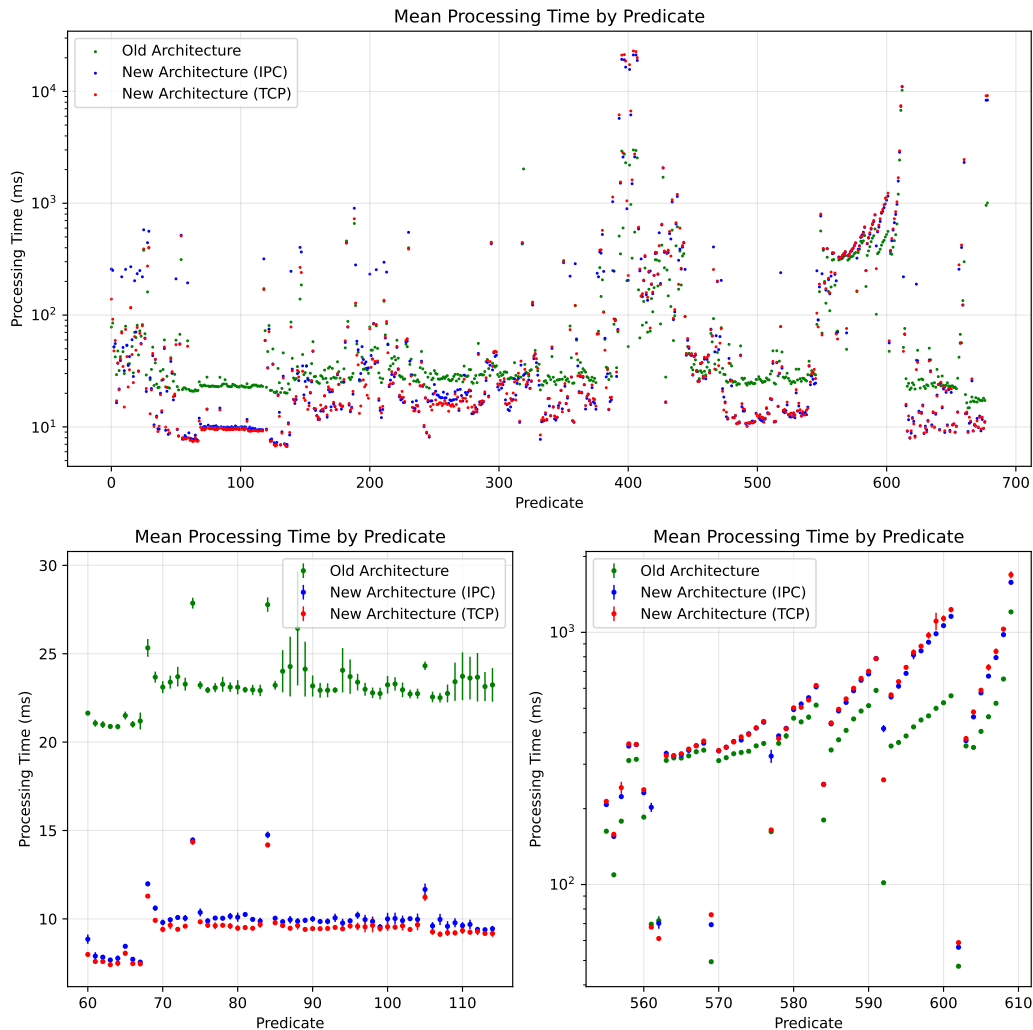


Abbildung 2: Durchschnittliche Laufzeiten der Anfragen in den verschiedenen Architekturen.

Um sich zunächst einen Überblick über die Rohdaten zu verschaffen, werden in [Abbildung 2](#) die durchschnittlichen Laufzeiten in den verschiedenen Konfigurationen dargestellt. Der obere Subgraph zeigt die durchschnittlichen Laufzeiten aller Anfragen in sowohl der alten als auch der neuen Architektur, wobei die neue Architektur in das IPC-Protokoll und TCP-Protokoll unterteilt ist. Wider Erwarten zeigt sich, dass die durchschnittliche Laufzeit in der neuen Architektur bei vielen Prädikaten geringer ausfällt als in der alten Architektur. Insgesamt sind von den 679 Prädikaten nur 172 Prädikate in der alten Architektur schneller gelöst worden. 265 Prädikate wurden in der IPC-Konfiguration und 242 Prädikate in der

TCP-Konfiguration am schnellsten gelöst.

Die unteren beiden Subgraphen stellen interessante Bereiche der Rohdaten erneut dar und zeigen zusätzlich die Standardfehler der Mittelwerte in Form der Fehlerbalken. Diese geben an, wie sehr die Mittelwerte der Laufzeiten von den tatsächlichen Mittelwerten der Grundgesamtheit abweichen. Kalkuliert wird dieser Standardfehler mittels der folgenden Formel:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad (1)$$

Hierbei ist $\sigma_{\bar{x}}$ der Standardfehler des Mittelwertes, σ die Standardabweichung der Grundgesamtheit und n die Anzahl der Messungen. Somit werden ausschließlich statistische Fehler betrachtet und systematische Fehler nicht berücksichtigt.

So zeigt der untere linke Subgraph die Rohdaten im Bereich der Prädikate 60 bis 115 erneut, wobei klar zu erkennen ist, dass hier die neue Server-Architektur konstant schneller ist. Dieser Trend ist über die gesamte Messung hin unterhalb einer gewissen Grenze zu beobachten. Ebenfalls ersichtlich ist der nur minimal ausfallende Unterschied zwischen IPC und TCP.

Auf der anderen Seite zeigt der untere rechte Subgraph die Rohdaten im Bereich der Prädikate 555 bis 610. Hier haben die Laufzeit Messungen einen vergleichsweise hohen Wert und es ist zu erkennen, dass die alte Architektur in diesem Bereich schneller ist.

Entsprechend wurde durch die Protokollkommunikation ein Performance-Overhead eingeführt, der sich in den Rohdaten dahingehend widerspiegelt, dass Prädikate, welche bereits eine hohe Laufzeit aufgewiesen haben, nun in der neuen Architektur noch langsamer gelöst werden. Neben diesem Overhead wurde jedoch eine signifikante Laufzeitverbesserung erworben, die bis zu einer gewissen Gesamtlaufzeitgrenze den Overhead nicht nur annulliert, sondern überkompensiert. Insgesamt scheint die neu eingeführte Server-Architektur hinsichtlich der Laufzeitperformance um einen gewissen Faktor effizienter geworden zu sein.

Ebenfalls ist zu erkennen, dass der Unterschied zwischen IPC und TCP nur minimal ist und keine signifikanten Unterschiede aufweist. Innerhalb der 679 Testprädikaten ist IPC in 365 (53.76%) Fällen schneller als TCP und entsprechend TCP in 314 (46.24%) Fällen schneller als IPC. Da zur TCP-Kommunikation die Adresse `tcp : //127.0.0.1` verwendet wurde, ist eine mögliche Erklärung hierfür, dass die Kommunikation über das Loopback-Interface einen Großteil des Netzwerk-Stacks umgeht und vom Betriebssystem speziell optimiert wird. In der weiteren Analyse wird daher nur das IPC-Protokoll betrachtet.

Um den Overhead der neuen Architektur zu quantifizieren, wird die Differenz der durchschnittlichen Laufzeiten der neuen und alten Architektur berechnet und gegen die Anzahl der ZeroMQ-Anfragen analysiert. Der induzierte Overhead wird also wie folgt definiert:

$$\text{InduzierterOverhead} = \overline{IPC\text{Laufzeit}} - \overline{Alte\text{Laufzeit}} \quad (2)$$

Der oberste Subgraph in Abbildung 3 zeigt den induzierten Overhead der neuen Architektur in Abhängigkeit der Anzahl der ZeroMQ-Anfragen auf einer logarithmischen Skala. Im groben Verlauf ist zu erkennen, dass der Overhead mit steigender Anzahl der Anfragen annähernd linear zunimmt. Zusätzlich scheint jedoch eine ausgeprägte Systematik vorzuliegen, welche sich inhaltlich durch die horizontale Verzerrung der Daten ausprägt. Einige wenige Datenpunkte weichen stark von jeglicher Systematik ab und sind als Ausreißer zu klassifizieren.

Im zweiten Subgraphen wird die lineare Abhängigkeit des Overheads von der Anzahl der Anfragen untersucht, indem der Zusammenhang linear gefittet wird. Dies geschieht unter Berücksichtigung des kombinierten Standardfehlers der Mittelwerte der alten und neuen IPC-Architektur

$$\sigma_{comb} = \sqrt{\sigma_{IPC}^2 + \sigma_{Old}^2} \quad , \quad (3)$$

welche ebenfalls in den Fehlerbalken dargestellt wird. Das Ergebnis des linearen Fits zeigt eine Steigung von $0.03ms$ pro Anfrage, was bedeutet, dass der Overhead pro Anfrage um $0.03ms$ steigt. Zusätzlich gibt es einen y-Achsenabschnitt von $-13.13ms$, was bedeutet, dass der Overhead bei 0 Anfragen $-13.13ms$ beträgt. Dieser Wert beschreibt die zuvor erfasste Laufzeitverbesserung der neuen Architektur. Insgesamt ergibt sich die folgende lineare Funktion zur Beschreibung des Overheads:

$$\text{Overhead} = 0.03ms \cdot \text{Anfragenanzahl} - 13.13ms \quad (4)$$

Durch das Gleichsetzen der Overhead-Funktion mit 0 ergibt sich eine Anzahl an Anfragen von etwa 438, bis der Overhead nicht länger von der Laufzeitverbesserung kompensiert werden kann.

Eine separate Messreihe wurde durchgeführt, um die durchschnittliche Laufzeit einer einzelnen Anfrage an den Z3-Server zu ermitteln. Hierbei wurde die durchschnittliche Laufzeit von $0.0054ms$ bestimmt. Da für jede Anfrage zwei Nachrichten (Anfrage und Antwort) über den Socket gesendet werden, ergibt sich eine Laufzeit von $0.0216ms$ pro Anfrage unter der Annahme, dass das Erhalten einer ZeroMQ-Nachricht die gleiche Laufzeit aufweist wie das Senden einer Nachricht. Die Differenz von $0.0084ms$ zur ermittelten Steigung des linearen Fits ist beispielsweise auf die Größe der Nachrichten zurückzuführen, welche in der kontrollierten Messreihe einer einzelnen Anfrage kleiner ausfällt als in den 52

Testszenarien. Zudem sind Messwerte dieser Größenordnung zunehmend ungenau, weshalb sie annnehmbar plausibel bezüglich der Steigung des linearen Fits sind.

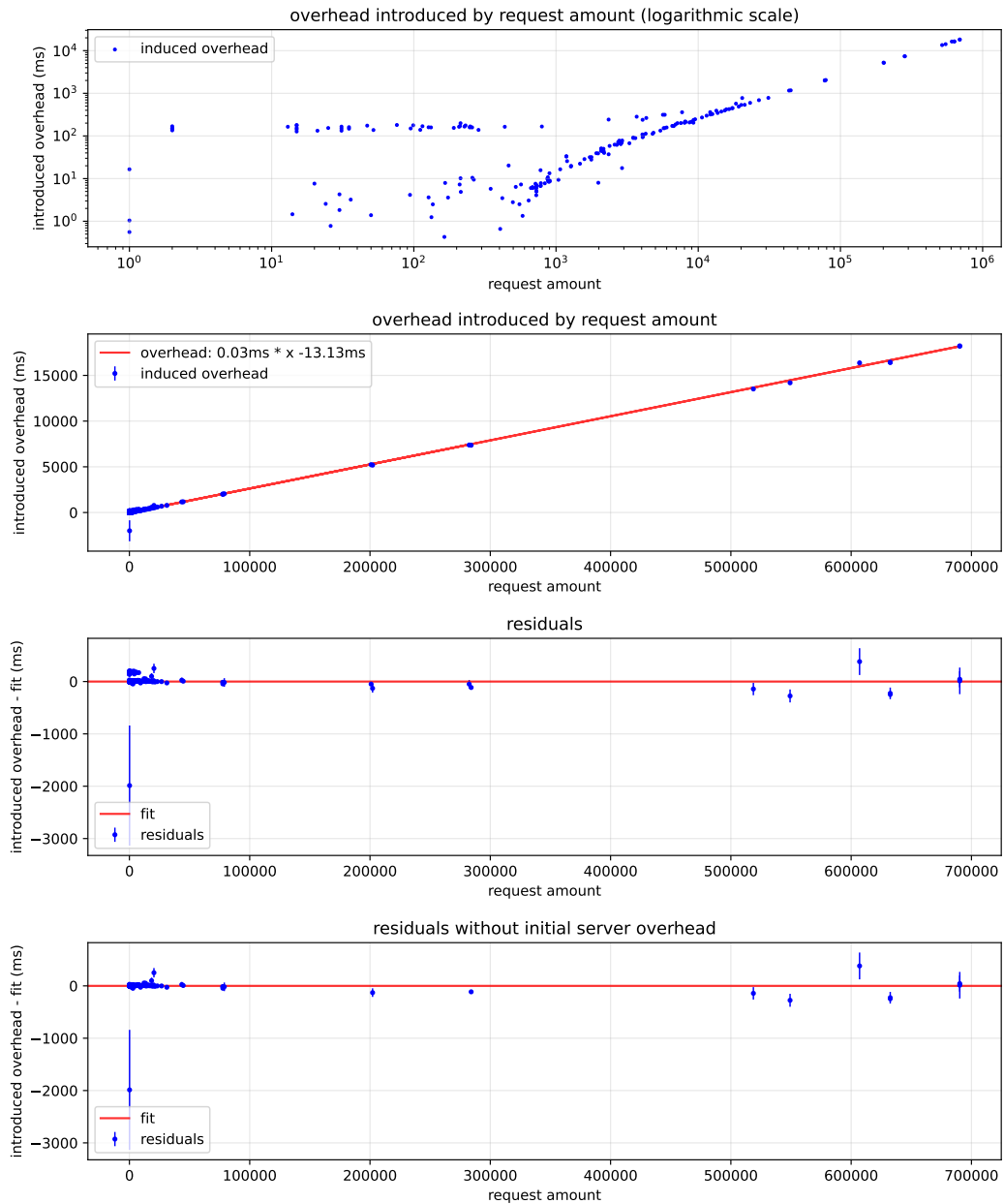


Abbildung 3: Induzierter Overhead der neuen Server-Architektur durch das Serialisieren auf den Socket.

Der lineare Fit zeigt augenscheinlich eine gute Übereinstimmung mit den Daten, welche in dem dritten Subgraphen aus Abbildung 3 verifiziert wird. Hierbei wird die Differenz des

induzierten Overheads und des linearen Fits gegen die Anzahl der Anfragen dargestellt. Es ergibt sich ein Residuenplot. In diesem Plot ist die Systematik des obersten Subgraphen erneut zu erkennen. Sie bildet ein kleines Cluster oberhalb der 0-Linie, bei sehr kleinen Anfragenzahlen. Es stellt sich heraus, dass innerhalb dieses Datenclusters alle Datenpunkte die QueryID 1 aufweisen. Somit lässt sich die Abweichung durch das Initialisieren des Sockets, der Verbindung des Z3-Servers und dessen initialen Overheads zum Aufsetzen der Z3-Konfigurationen erklären. Innerhalb einer TestID wird der Z3-Server nur einmalig initialisiert und folgende Prädikate desselben Tests nutzen denselben Z3-Server, indem sie die Konfiguration zurücksetzen und nicht grundlegend neu initialisieren.

Zur Überprüfung stellt der letzte Subgraph erneut den Residuenplot dar, jedoch ohne diejenigen Datenpunkte mit QueryID 1. Hierbei ist zu erkennen, dass die Systematik des vorherigen Graphen verschwindet und der lineare Fit eine gute Übereinstimmung mit den Daten aufweist. Somit ist gezeigt, dass der Overhead der neuen Architektur durch das Serialisieren auf den Socket annähernd linear mit der Anzahl der Anfragen zunimmt und keine weiteren relevanten Systematiken besitzt. Eine Voraussage des Overheads ist mithilfe des linearen Fits möglich. Es ist zu beachten, dass die absoluten Werte des Overheads nur in der gegebenen Testumgebung gültig sind und nicht auf andere Umgebungen übertragen werden können.

Ein einzelner Datenpunkt scheint eine unwahrscheinlich hohe Laufzeitverbesserung von durchschnittlich 2 Sekunden aufzuweisen, zeichnet jedoch einen hohen Fehlerbalken auf. In Tabelle 2 ist dieser Datenpunkt dargestellt.

Tabelle 2: Ausschnitt der gesammelten Messwerte eines Ausreißers.

TestID-QueryID	$Old_1(\text{ns})$	$Old_2(\text{ns})$	$Old_3(\text{ns})$	$Old_4(\text{ns})$	Req. Count
2122-90	36	4033	3993	43	119

Die verschiedenen Messwerte der alten Architektur in Millisekunden zeigen zwei starke Ausreißer bei der zweiten und dritten Messung. Diese Varianz ist womöglich auf die in Abschnitt 5.1 beschriebene Problematik zurückzuführen. In jedem Fall ist dieser Datenpunkt als individueller Ausreißer zu betrachten und weist keine besondere statistische Relevanz auf.

Die ermittelte Performance-Verbesserung wird an dieser Stelle nicht weiter analysiert, da sie nicht Gegenstand dieser Arbeit ist. Ein möglicher Grund für die Verbesserung könnte die Kompilierung sein, die durch die neue Architektur ermöglicht wird, da der Z3-Solver als eigenständiger Prozess womöglich besser vom Compiler optimiert werden kann. Als eigenständiger Prozess ist es ebenfalls möglich, dass der runtime linker die Z3-Bibliothek optimierter laden und ausführen kann. Zusätzlich könnten die in Abschnitt 3.2.3 beschriebenen Optimierungen eine Rolle spielen.

7 Zukünftige Arbeiten

Es existieren verschiedene Möglichkeiten, um auf die in dieser Arbeit vorgestellte Implementierung aufzubauen und sie weiter zu verbessern. Folgende Abschnitte stellen einige dieser Möglichkeiten vor.

7.1 Weitere Analyse und Optimierungen

Die in Abschnitt 6 vorgestellte Performance-Evaluation ist eine erste Analyse der Implementierung und der Auswirkungen durch die Entkopplung von ProB und Z3. Sie stellt überraschenderweise eine Verbesserung der Performance fest, obwohl die Systemänderung an sich einen Overhead darstellt. Um ein tieferes Verständnis über die Laufzeit zu erlangen, ist es sinnvoll, die Analyse zu erweitern und zu vertiefen. Eine ausreichend tiefgründige Analyse könnte insofern von Vorteil sein, um ein besseres Verständnis über die Performanceverbesserung zu erlangen und somit gegebenenfalls gezielt weitere Optimierungen vorzunehmen.

7.2 Deinit Hook

Nach der erfolgreich abgeschlossenen Systemarchitekturänderung, die in dieser Arbeit vorgestellt wurde, bleibt ein letzter Schritt offen, um die Implementierung vollständig abzurunden. Derzeitig wird der Z3-Solver-Prozess nicht sauber beendet, sondern existiert als sogenannter *child process* von ProB und ist somit an die Existenz des ProB-Prozesses gekoppelt. Dies stellt insofern grundsätzlich kein großes Problem dar, da der Z3-Solver-Prozess nach Beendigung des ProB-Prozesses automatisch beendet wird. Dennoch besteht in der Theorie die Möglichkeit, dass der Z3-Server als sogenannter *dangling process* zurückbleibt. Grundsätzlich ist es zudem eine sauberere Lösung, den Z3-Prozess manuell zu beenden.

Des Weiteren bleibt bei Beendigung des Solvers der von ZeroMQ instanziierte Socket geöffnet und als *dangling socket* im Arbeitsspeicher zurück. Dieser Umstand kann zu Speicherlecks führen und sollte daher vermieden werden.

Um diese Probleme zu beheben, ist es sinnvoll, die Implementierung um einen Deinit-Hook zu erweitern. Dieser Hook wird aufgerufen, wenn der Z3-Solver beendet beziehungsweise deinitialisiert wird, und dient dazu, den Z3-Solver-Prozess und den Z3-Socket sauber zu beenden und zu schließen.

7.3 Parallelisierung

Die Entkopplung von ProB und Z3 dient elementar als Grundlage für die Parallelisierung des Z3-Solvers beziehungsweise die Parallelisierung der Lösung mehrerer Prädikate. Wie in

Abschnitt 1.2 beschrieben, ist die Entkopplung mit eben jenem Interesse der Erweiterbarkeit verrichtet worden.

Es bestehen verschiedene Ansätze, um die Parallelisierung zu realisieren. Eine Möglichkeit wäre, mehrere Instanzen des Z3-Solvers zu starten und die Prädikate auf diese Instanzen zu verteilen. In diesem Fall müsste die Kommunikationsschnittstelle auf der Seite von ProB so erweitert werden, dass sie mehrere Z3-Solver-Prozesse unterstützt. Der Z3-Server selbst müsste in diesem Szenario nicht verändert werden.

Eine andere Möglichkeit besteht darin, auch den Z3-Server zu erweitern, sodass dieser mehrere Prädikate gleichzeitig lösen kann. Hierbei müsste die Kommunikationsschnittstelle auf der Seite des Z3-Servers so erweitert werden, dass sie mehrere Prädikate gleichzeitig empfangen und lösen kann. Es würde sich anbieten, in diesem Fall das Kommunikationsmuster zu wechseln und auf ein asynchrones Modell umzusteigen, wie beispielsweise das Dealer-Router Modell, welches in Abschnitt 2.3 bereits erwähnt wurde.

8 Konklusion

hier wird
konkludiert.
alles supi
mit rundem
ende oder
so...

stuff

Literatur

- [Abr96] ABRIAL, Jean-Raymond: *The B-Book: Assigning Programs to Meanings*. New York, NY, USA : Cambridge University Press, 1996
- [CWA⁺88] CARLSSON, Mats ; WIDEN, Johan ; ANDERSSON, Johan ; ANDERSSON, Stefan ; BOORTZ, Kent ; NILSSON, Hans ; SJÖLAND, Thomas: *SICStus Prolog User's Manual*. Kista, Sweden : Swedish Institute of Computer Science, 1988
- [Hin13] HINTJENS, P: *ZeroMQ: Messaging for Many Applications*. O'Reilly Media, 2013
- [KL16] KRINGS, Sebastian ; LEUSCHEL, Michael: SMT Solvers for Validation of B and Event-B Models. In: ÁBRAHÁM, Erika (Hrsg.) ; HUISMAN, Marieke (Hrsg.): *Integrated Formal Methods*. Cham : Springer International Publishing, 2016. – ISBN 978-3-319-33693-0, S. 361–375
- [LB03] LEUSCHEL, Michael ; BUTLER, Michael: ProB: A model checker for B. In: *FME 2003: Formal Methods: International Symposium of Formal Methods Europe, Pisa, Italy, September 8-14, 2003. Proceedings* Springer, 2003, S. 855–874
- [LB08] LEUSCHEL, Michael ; BUTLER, Michael: ProB: an automated analysis toolset for the B method. In: *International Journal on Software Tools for Technology Transfer* 10 (2008), S. 185–203
- [MB08] MOURA, Leonardo de ; BJØRNER, Nikolaj: Z3: An Efficient SMT Solver. In: RAMAKRISHNAN, C. R. (Hrsg.) ; REHOF, Jakob (Hrsg.): *Tools and Algorithms for the Construction and Analysis of Systems*. Berlin, Heidelberg : Springer Berlin Heidelberg, 2008. – ISBN 978-3-540-78800-3, S. 337–340
- [MR15] MICROSOFT-RESEARCH: *GitHub - Z3Prover/z3: The Z3 Theorem Prover*. 2015
- [S⁺15] SÚSTRIK, Martin u. a.: ZeroMQ. In: *Introduction Amy Brown and Greg Wilson* (2015), S. 16
- [TJ07] TORLAK, Emina ; JACKSON, Daniel: Kodkod: A relational model finder. In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems* Springer, 2007, S. 632–647