

Hardwarebeschleunigung von paralleler Ablaufplanung für Multicoresysteme

Zwischenbericht zum 1. Jahr

Nina Engelhardt

12. Juni 2012

Ziel meiner Arbeit ist es, durch speziell angepasste Hardware die Synchronisation und Aufgabenverteilung für parallele Programme auf Multicoresystemen zu verbessern. Als Grundlage dient dafür das VMS-System, ein Framework das es ermöglicht, Laufzeitumgebungen für verschiedenste parallele Programmiermodelle mit nur wenig Aufwand zu implementieren.

In den ersten 6 Monaten habe ich wie geplant das VMS-System analysiert. VMS umfasst die Grundlegende Struktur, die allen Laufzeitumgebungen gemein ist: sie beinhaltet den Wechsel zwischen Programm und Laufzeitumgebung, und die Möglichkeit, eine Zeitliche Ordnung zwischen zwei Codepunkten in zwei virtuellen Prozessoren zu garantieren (semantikfreie Synchronisation). Auf Basis dieser Elemente muss der/die Laufzeitumgebungsentwickler in nur noch die jeweilige Semantik der parallelen Umgebung hinzufügen, in Form eines Plugins, das Handlers für jedes parallele Konstrukt und einen Assigner, der entscheidet, welche Aufgaben als nächstes bearbeitet werden, bereitstellt.

Diese zweiteilige Struktur ermöglicht zwei Ansatzpunkte für die Beschleunigung: den VMS-Kern selbst – Verbesserungen an dieser Stelle kommen allen Laufzeitumgebungen zu Gute – und das jeweilige Plugin. VMS ist minimal gehalten und bietet nur wenige Möglichkeiten zur Verbesserung, jedoch haben Änderungen hier den größten Effekt. VMS ist außerdem, da es die grundlegendsten Elemente zusammenfasst, bereits sehr nah an der existierenden Hardware und die häufigen Aufgaben sind gut unterstützt. Zusätzliche Beschleunigung kann hier nur durch sehr ressourcenintensive Lösungen (z.B. zusätzliche Registersätze für schnellen Kontextwechsel) erwartet werden. Beschleuniger für Pluginfunktionen haben deutlich mehr Spielraum, besonders für Sprachen mit komplexen, hardwarefernen Konstrukten (und dies sind die besonders interessanten, weil sie sich der Struktur des Problems nähern und damit dem Programmierer die Arbeit abnehmen, das Problem in Hardwarenahe Konzepte umzudenken bzw. übersetzen zu müssen.) Da die Pluginfunktionen aber Sprachspezifisch sind, kommen Beschleunigungen an dieser Stelle nur dem Teil der Programme zugute, die in dieser Sprache geschrieben sind. Die Herausforderung wird also sein, Hardwareelemente zu finden, die möglichst flexible und vielfältige Verwendung finden können. Hauptsächlich wird aber die Programmiersprache StarS anvisiert werden, die im Fachgebiet AES häufig Verwendung findet.

VMS ist die praktische Ausformulierung eines Modells des Parallelismus, die von Dr. Sean Halle unter dem Namen “Holistic Model of Parallel Computation” [] vorgeschlagen wurde. Dieses Modell isoliert die schedulingrelevanten Teile eines Programms und legt das Augenmerk besonders auf die Effekte von Scheduling Decisions auf die Performance.

Während meiner Analyse von VMS stieß ich auf mehrere unklare Stellen in der Theorie, und habe diese in Zusammenarbeit mit Dr. Halle ausgearbeitet. Insbesondere habe ich das Modell um das Konzept der Schichten erweitert, das dem Umstand Rechnung trägt, dass in einem System meistens mehrere hierarchisch untergeordnete Scheduler existieren. So werden

die Arbeitseinheiten, über die die Laufzeitumgebung als Eins entscheidet, möglicherweise vom Betriebssystem noch einmal unterbrochen und verteilt, spätestens aber vom Prozessor in kleinere Einheiten – Assemblerbefehle – zerteilt, und auf unterschiedliche Funktionseinheiten aufgeteilt. Im Fall von Out-of-order Prozessoren wird deutlich, dass auch diese parallele Aufgabenverteilung ist, und das Holistische Modell auch auf diese anwendbar ist. Zur Zeit kollaborieren wir an einem Journal Article der das Holistische Modell detailliert darlegt.

Zu den Errungenschaften des Holistischen Modells gehören auch zwei grafische Darstellungen, die Unit Constraint Collection und der Scheduling Consequence Graph, deren Ziel es ist, die Gründe der (guten oder schlechten) Performance eines parallelen Programms ersichtlich zu machen. Die Unit Constraint Collection (kurz UCC) zeigt die Arbeitseinheiten (units) die im Programm definiert werden und die Randbedingungen (constraints) ihrer Ausführbarkeit, die vom Programm selbst ausgehen. Der Consequence Graph (CG) repräsentiert eine realisierte Ausführung des Programms und zeigt die räumliche und zeitliche Platzierung der Arbeitseinheiten, die der Scheduler gewählt hat, sowie die unterschiedlichen Randbedingungen, die diese Wahl eingeschränkt haben. Diese können vom Programm selbst stammen, sich aus Ressourcenbegrenzungen ergeben, oder durch Beschränkungen der Laufzeitumgebung entstehen.

Nachdem ich VMS bereits zum Zweck der Analyse des Laufzeitverhaltens instrumentalisiert hatte, habe ich diese Infrastruktur wiederverwendet, um aus den gewonnenen Daten die räumliche und zeitliche Platzierung der Arbeitseinheiten im Consequence Graph abzuleiten. Zusätzlich habe ich ein VMS-Plugin, das synchrones Senden und Empfangen von Nachrichten zwischen Threads in einem Programm unterstützt (synchronous send-recv, kurz SSR), instrumentalisiert, um die Constraints zu erfassen, die für die Konstruktion der UCC und des CG notwendig sind. Schließlich erstellte ich ein kleines Programm, das diese Aufzeichnungen entgegennehmen und UCC & CG darstellen kann.

Quasi mit dem ersten Blick stellte sich heraus, dass die Matrixmultiplikation, die ich seit Monaten als Testapplikation verwendete, einige Performanceprobleme hatte. So hatte z.B. der applikationseigene Lastenverteiler einen Bug, sodass von 40 Prozessorkernen nur 3 den Großteil der Arbeit zugeteilt bekamen, und 10 weitere einen kleinen Teil, die 27 restlichen jedoch gar keine. Des Weiteren wurde die Arbeit so verteilt, dass der Kern, der den Verteiler ausführt, als erster Arbeit zugeteilt bekam, und diese Arbeit dann den Verteiler für längere Zeit unterbricht, bevor dieser weiteren Kernen Arbeit zuteilen kann. All dies war in den Statistiken nicht aufgefallen, in der Visualisierung jedoch offensichtlich. Als ich noch eine detailliertere Aufteilung des Overheads, der für eine bestimmte Arbeitseinheit aufgewendet wird, und mehrere Metriken (Takte, Befehle, Cache Misses) hinzugefügt hatte, wurden auch subtilere Effekte sichtbar. Dabei stellte sich z.B. heraus, dass die zentralisierte Architektur von VMS, die alle Informationen die die Laufzeitumgebung über den Zustand des Programms bereithalten muss in einem gemeinsamen Pool speichert, auf den nur ein Kern gleichzeitig zugreifen kann, einen größeren Engpass als erwartet darstellte. Da mir keine äquivalenten Performance-debugging Tools bekannt waren, und mehrere Mitarbeiter des Fachgebiets bereits händeringend nach solchen gesucht hatten, habe ich über dieses Tool einen weiteren Artikel verfasst [Anhang?], den ich voraussichtlich im August für die PPOP-2013 Konferenz einreichen werde.

Aufgrund der gewonnenen Erkenntnisse verändert sich mein Plan wie folgt:

Da die Vermutung, dass ein zentralisiertes System, in dem ein Master-Kern die Aufgabenverteilung für das gesamte System übernimmt, die beste Lösung sein würde, nicht so gesichert wie am Anfang erscheint, werde ich mich stärker für die verschiedenen Möglichkeiten der Parallelisierung interessieren.