

Erhard Rahm

Hochleistungs- Transaktionssysteme

Konzepte und Entwicklungen
moderner Datenbankarchitekturen

erschienen im:

Vieweg-Verlag, Reihe Datenbanksysteme, 1993,
286 Seiten, ISBN 3-528-05343-7

Vorwort

Transaktionssysteme sind in der kommerziellen Datenverarbeitung weitverbreitet und von enormer ökonomischer Bedeutung. Die Leistungsanforderungen vieler Anwendungen steigen sehr stark und können nur von sogenannten **Hochleistungs-Transaktionssystemen** erfüllt werden. Dies betrifft insbesondere die zu bewältigenden Transaktionsraten für vorgeplante Anwendungsfunktionen. Daneben sind auch zunehmend ungeplante Ad-Hoc-Anfragen auf derselben Datenbank auszuführen, welche oft den Zugriff auf große Datenmengen verlangen und für die die Gewährleistung kurzer Antwortzeiten ein Hauptproblem darstellt. Weitere Anforderungen an Hochleistungs-Transaktionssysteme betreffen die Gewährleistung einer hohen Verfügbarkeit, die Möglichkeit eines modularen Wachstums sowie die Unterstützung einer hohen Kosteneffektivität. Für die Erzielung einer hohen Leistungsfähigkeit müssen u.a. die enormen technologischen Fortschritte bei der CPU-Kapazität für die Transaktionsverarbeitung umgesetzt werden können. Die Kosteneffektivität leidet darunter, daß Transaktionssysteme herkömmlicherweise auf zentralen Großrechnern (Mainframes) laufen. Zur Steigerung der Kosteneffektivität gilt es daher, zunehmend Mikroprozessoren für die Transaktionsverarbeitung zu nutzen.

Die erwähnten Anforderungen können mit herkömmlichen Systemarchitekturen zur Transaktions- und Datenbankverarbeitung meist nicht erfüllt werden. Ein zunehmendes Problem stellt das E/A-Verhalten für die Externspeicherzugriffe dar. Denn die Leistungsmerkmale von Magnetplatten konnten in der Vergangenheit im Gegensatz zur Prozessorkapazität kaum verbessert werden, so daß v.a. für Hochleistungssysteme zunehmend die Gefahr von E/A-Engpässen besteht. Weiterhin können viele der genannten Anforderungen nur mit verteilten Transaktionssystemen erfüllt werden, während Transaktionssysteme traditionellerweise zentralisiert realisiert wurden.

Im Mittelpunkt dieses Buchs steht der Einsatz neuerer Systemarchitekturen zur Realisierung von Hochleistungs-Transaktionssystemen. Dabei werden vor allem Ansätze zur Optimierung des E/A-Verhaltens sowie verteilte Transaktionssysteme behandelt, mit denen die Beschränkungen herkömmlicher Systeme behoben werden können. Wie auch im Titel des Buches zum Ausdruck kommt, konzentriert

sich die Darstellung dabei vor allem auf die Datenbankaspekte bei der Transaktionsverarbeitung.

Zur **Optimierung des E/A-Verhaltens** werden neben allgemeinen Implementierungstechniken des Datenbanksystems (Pufferverwaltung, Logging) folgende Ansätze behandelt:

- Hauptspeicher-Datenbanksysteme
- Verwendung von Disk-Arrays (Plattenfeldern)
- Nutzung von seitenadressierbaren Halbleiterspeichern wie Platten-Caches, Solid-State-Disks und erweiterte Hauptspeicher im Rahmen einer erweiterten Speicherhierarchie.

Besonders ausführlich wird letzterer Ansatz dargestellt. Eine detaillierte Leistungsanalyse zeigt, daß selbst mit einem begrenzten Einsatz der Speichertypen signifikante Leistungssteigerungen gegenüber konventionellen E/A-Architekturen erreicht werden.

Einen weiteren Schwerpunkt bildet die Untersuchung **verteilter Transaktionssysteme**. Anhand einer allgemeinen Klassifikation werden die wichtigsten Systemarchitekturen zur verteilten Transaktionsverarbeitung vorgestellt und hinsichtlich der zu erfüllenden Anforderungen bewertet. Zur Realisierung von Hochleistungs-Transaktionssystemen besonders geeignet sind zwei Klassen lokal verteilter Mehrrechner-Datenbanksysteme, nämlich der DB-Sharing (Shared-Disk)- und der Shared-Nothing-Ansatz. Diese Mehrrechner-DBS können auch innerhalb von Workstation-/Server-Architekturen zur Realisierung eines verteilten Server-Systems eingesetzt werden. Workstation-/Server-Systeme erlauben eine Steigerung der Kosteneffektivität, da dabei Mikroprozessoren relativ einfach zur Transaktionsverarbeitung genutzt werden können (für die Workstations sowie auf Server-Seite).

Die effektive Nutzung eines Mehrrechner-DBS verlangt vor allem eine weitgehende Reduzierung von Kommunikationsverzögerungen und -Overhead. Da dieses Ziel mit lose gekoppelten Rechnerarchitekturen häufig nur begrenzt erreicht werden kann, wird der Einsatz einer sogenannten nahen Rechnerkopplung untersucht. Dabei sollen v.a. gemeinsame Halbleiterspeicher für eine effiziente Kommunikation sowie zur Realisierung globaler Kontrollaufgaben genutzt werden. Vor allem DB-Sharing-Systeme können von einer nahen Kopplung profitieren, da eine einfache und effiziente Realisierung kritischer Funktionen (z.B. globale Sperrbehandlung) ermöglicht wird. Ein Leistungsvergleich belegt, daß damit vor allem für reale Lasten hohe Leistungsgewinne gegenüber loser Kopplung erreicht werden.

Um auch für komplexe Datenbankanfragen kurze Bearbeitungszeiten zu ermöglichen, wird der Einsatz von Intra-Transaktionsparallelität auf einem Multipro-

zessor bzw. Mehrrechner-DBS immer wichtiger. Hierzu diskutieren wir die Realisierung verschiedener Parallelisierungsarten und betrachten den Einfluß der Datenverteilung und Systemarchitektur. Weitere Kapitel behandeln die Realisierung einer automatischen Lastkontrolle in Transaktionssystemen sowie die Unterstützung einer schnellen Katastrophen-Recovery.

Das Buch richtet sich an Informatiker in Studium, Lehre, Forschung und Entwicklung, die an neueren Entwicklungen im Bereich von Transaktions- und Datenbanksystemen interessiert sind. Es entspricht einer überarbeiteten Version meiner im Februar 1993 vom Fachbereich Informatik der Universität Kaiserslautern angenommenen Habilitationsschrift. Neben der Präsentation neuer Forschungsergebnisse erfolgen eine breite Einführung in die Thematik sowie überblicksartige Behandlung verschiedener Realisierungsansätze, wobei auf eine möglichst allgemeinverständliche Darstellung Wert gelegt wurde. Der Text wurde durchgehend mit Marginalien versehen, welche den Aufbau der Kapitel zusätzlich verdeutlichen und eine schnelle Lokalisierung bestimmter Inhalte unterstützen sollen.

Die vorgestellten Forschungsergebnisse entstanden zum Teil innerhalb des DFG-Projektes "Architektur künftiger Transaktionssysteme", das seit 1990 unter meiner Leitung am Fachbereich Informatik der Universität Kaiserslautern durchgeführt wird. Teile des Stoffes gingen auch aus meinen Vorlesungen über Datenbanksysteme an der Universität Kaiserslautern hervor. Die Untersuchungen zur Lastkontrolle begannen während eines einjährigen Forschungsaufenthaltes am IBM T.J. Watson Research Center in Yorktown Heights, N.Y., USA.

Mein besonderer Dank geht an Herrn Prof. Dr. Theo Härder für die langjährige Betreuung und Förderung meines wissenschaftlichen Werdegangs. Herrn Prof. Dr. Gerhard Weikum von der ETH Zürich danke ich für die Übernahme des Koreferats sowie detaillierte und wertvolle Verbesserungsvorschläge. Hilfreiche Anmerkungen erhielt ich darüber hinaus von Herrn Prof. Dr. Andreas Reuter (Univ. Stuttgart). Für die geleistete Arbeit im Rahmen des DFG-Projektes danke ich meinen Mitarbeitern sowie Studenten, insbesondere den Herren K. Butsch, R. Marek, T. Stöhr, P. Webel und G. Wollenhaupt. Für die fruchtbare Zusammenarbeit während des Aufenthaltes bei IBM sei vor allem Herrn Dr. C. Nikolaou, Herrn Dr. D. Ferguson und Herrn Dr. A. Thomasian gedankt.

Kaiserslautern, im März 1993

Erhard Rahm

Inhaltsverzeichnis

Vorwort

Inhaltsverzeichnis

1	Einführung	1
1.1	Transaktionssysteme	1
1.2	Das Transaktionskonzept	4
1.3	Transaktionslasten.....	5
1.4	Anforderungen an künftige Transaktionssysteme	10
1.5	Aufbau des Buchs	13
2	Aufbau und Funktionsweise von Transaktionssystemen	15
2.1	Modellbildung	15
2.1.1	Schichtenmodell	15
2.1.2	Modell von Bernstein	18
2.1.3	X/Open-Modell.....	19
2.2	Implementierungsaspekte	21
2.2.1	Prozeß- und Taskverwaltung.....	21
2.2.2	Betriebssystem-Einbettung von TP-Monitor und DBS.....	22
2.2.3	Speicher- und Programmverwaltung	24
2.2.4	Transaktionsverwaltung.....	25
3	Verteilte Transaktionssysteme.....	27
3.1	Allokationsprobleme in verteilten Transaktionssystemen	28
3.2	Horizontal verteilte Transaktionssysteme.....	31
3.2.1	Verteilung nur im DC-System.....	31
3.2.2	Verteilung im DBS (horizontal verteilte Mehrrechner-DBS)	37
3.2.3	Verteilung im DC-System und im DBS	41
3.3	Vertikal verteilte Transaktionssysteme.....	42
3.3.1	Verteilung im DC-System.....	42
3.3.2	Verteilung im DBS (vertikal verteilte Mehrrechner-DBS)	44
3.3.3	Vertikale Verteilung über mehr als zwei Ebenen	45
3.4	Kombination von horizontaler und vertikaler Verteilung	46
3.5	Bewertung verschiedener Verteilformen	48
3.6	Vergleich mit anderen Klassifikationsansätzen.....	50

4	Technologische Entwicklungstrends	53
4.1	Speedup und Scaleup	53
4.2	CPU	55
4.3	Halbleiterspeicher	58
4.4	Magnetplatten	61
4.5	Weitere Entwicklungen	66
4.6	Konsequenzen für die Transaktionsverarbeitung.....	67
5	Optimierung des E/A-Verhaltens	71
5.1	DBS-Implementierungstechniken zur E/A-Optimierung	72
5.1.1	Systempufferverwaltung	73
5.1.2	Logging	78
5.2	Hauptspeicher-DBS	83
5.3	Disk-Arrays	88
5.3.1	Datenverteilung	89
5.3.2	Fehlertoleranz.....	94
5.3.2.1	Fehlerkorrektur über Paritätsbits.....	96
5.3.2.2	Datenreplikation.....	103
5.3.3	Zusammenfassende Bewertung	105
5.4	Sonstige Ansätze	107
6	Einsatz erweiterter Speicherhierarchien	109
6.1	Speicherhierarchien	110
6.2	Seitenadressierbare Halbleiterspeicher.....	113
6.3	Einsatzformen zur E/A-Optimierung.....	117
6.4	Realisierung einer mehrstufigen Pufferverwaltung	121
6.4.1	Schreibstrategien in Speicherhierarchien	124
6.4.2	Verwaltung von Platten-Caches.....	127
6.4.2.1	Flüchtige Platten-Caches	128
6.4.2.2	Nicht-flüchtige Platten-Caches.....	130
6.4.3	Verwaltung eines EH-Puffers	133
6.5	Leistungsbewertung.....	138
6.5.1	Simulationsmodell	138
6.5.1.1	Datenbank- und Lastmodell	139
6.5.1.2	Modellierung der Transaktionsverarbeitung.....	143
6.5.1.3	Externspeichermodellierung.....	145

6.5.2	Simulationsergebnisse	146
6.5.2.1	Parametereinstellungen für die Debit-Credit-Experimente.....	147
6.5.2.2	Allokation der Log-Datei	148
6.5.2.3	Allokation von DB-Partitionen.....	150
6.5.2.4	FORCE vs. NOFORCE	152
6.5.2.5	Mehrstufige DB-Pufferung bei Debit-Credit	154
6.5.2.6	Mehrstufige DB-Pufferung bei der realen Last.....	159
6.5.2.7	Einfluß von Sperrkonflikten.....	161
6.5.3	Zusammenfassung der Simulationsstudie.....	164
7	Mehrrechner-Datenbanksysteme mit naher Kopplung	167
7.1	Alternativen zur Rechnerkopplung	169
7.2	Generelle Realisierungsalternativen zur nahen Kopplung.....	174
7.2.1	Nutzung von Spezialprozessoren.....	174
7.2.2	Verwendung gemeinsamer Halbleiterspeicher.....	175
7.3	Globaler Erweiterter Hauptspeicher (GEH)	178
7.4	Nachrichtenaustausch über einen GEH.....	180
7.5	Nutzung eines GEH bei DB-Sharing.....	183
7.5.1	Synchronisation und Kohärenzkontrolle bei lose gekoppelten DB-Sharing-Systemen.....	183
7.5.2	Synchronisation über den GEH.....	188
7.5.3	Weitere Einsatzmöglichkeiten des GEH.....	193
7.6	Leistungsanalyse	196
7.6.1	Simulationsmodell.....	196
7.6.2	Simulationsergebnisse	199
7.6.2.1	Parametereinstellungen für die Debit-Credit-Experimente.....	199
7.6.2.2	Leistungsverhalten bei naher Kopplung (GEH-Synchr.)	200
7.6.2.3	Lose vs. nahe Kopplung.....	206
7.6.2.4	Ergebnisse für die reale Last	209
7.6.3	Zusammenfassung der Simulationsstudie.....	214
8	Weitere Entwicklungsrichtungen.....	217
8.1	Lastkontrolle in Transaktionssystemen.....	218
8.1.1	Anforderungen.....	221
8.1.2	Architekturvorschlag für eine mehrstufige Lastkontrolle	223
8.1.3	Realisierung der globalen Lastkontrolle	226
8.1.3.1	Transaktions-Routing.....	226
8.1.3.2	Globale Kontrollentscheidungen.....	229
8.1.4	Realisierung der lokalen Lastkontrolle.....	230

8.2	Parallele DB-Verarbeitung.....	234
8.2.1	Arten der Parallelverarbeitung.....	235
8.2.2	Alternativen zur Datenverteilung.....	238
8.2.3	Realisierung von Intra-DML-Parallelität.....	242
8.2.3.1	Intra-Operatorparallelität	243
8.2.3.2	Inter-Operatorparallelität.....	248
8.2.4	Transaktionsverwaltung	251
8.2.5	Probleme der parallelen DB-Verarbeitung.....	253
8.3	Katastrophen-Recovery.....	256
9	Zusammenfassung und Ausblick.....	261
10	Literatur	269
	Index	281