

Adaptive Zugriffskontrollverfahren

Ein Entscheidungsmodell für die Kontrolle des Zugriffs auf gemeinsam genutzte IT-Infrastrukturen

Die Autoren

Thomas Setzer
Martin Bichler
Oliver Hühn

Dipl.-Wi.-Ing. Thomas Setzer
Prof. Dr. Martin Bichler
Dipl.-Inf. Oliver Hühn
Lehrstuhl für Internetbasierte Geschäftssysteme
Fakultät für Informatik
TU München
Boltzmannstr. 3
85748 Garching
{setzer | bichler | huehno}@in.tum.de
<http://ibis.in.tum.de>

tungsvereinbarungen beantworten zu können.

Das Service-Level-Management (SLM), ein zentraler Bereich des IT-Service-Management (ITSM), beschäftigt sich mit der Bereitstellung von IT-Dienstleistungen in vertraglich festgelegter Qualität zu vereinbarten Zeitpunkten. Ein wichtiger Bestandteil hiervon ist die automatisierte Behandlung von Überlastsituationen (overload control). In der Praxis konkurrieren Dienstanfragen mit unterschiedlichen Erträgen bzw. Geschäftswerten, Qualitätsvereinbarungen und Ressourcenanforderungen um gemeinsam genutzte, in ihrer Kapazität beschränkte Ressourcen (CPU, I/O, Hauptspeicher, Bandbreite, Lizenzen etc.). Nachdem Anwendungen individuelle Dienstleistungsvereinbarungen (Service-

Level-Agreements, SLAs) zugeordnet sind, ist eine Verletzung der SLAs mit unterschiedlichen Konsequenzen für den IT-Dienstleister verbunden. In Extremfällen, wie beispielsweise bei Stapelverarbeitungsaufgaben ohne zeitliche Restriktionen, ist eine spätere Ausführung eines Dienstes mit keinen Kosten oder negativen Folgen für den weiteren Geschäftsverlauf verbunden. In anderen Fällen, beispielsweise bei der Abwicklung einer Online-Buchung, kann eine spätere Ausführung den Verlust eines Auftrags bedeuten. Gemeinsam genutzte IT-Ressourcen erfordern somit ausgereifte, adaptive Konzepte für die Überlastkontrolle.

Knappe IT-Ressourcen sind beispielsweise CPU, Hauptspeicher, I/O, Bandbreite oder die Anzahl parallel ausführbarer

1 Einleitung

Betriebliche Infrastrukturkomponenten wie Datenbanksysteme, Web-Application-Server oder Content- bzw. Media-Server werden aufgrund von Effizienzvorteilen bei der Bereitstellung von IT-Dienstleistungen zunehmend von mehreren Diensten (bzw. Anwendungen) gemeinsam genutzt. Oftmals ist es aufgrund von Nachfrageschwankungen, unerwartetem Laufzeitverhalten einzelner Dienste oder auch durch den Ausfall von Hard- oder Softwarekomponenten unwirtschaftlich oder gar unmöglich, zu jedem Zeitpunkt ausreichend Ressourcen vorhalten zu können, um alle eintreffenden Anfragen entsprechend ihren Dienstleis-

Kernpunkte

Der Artikel stellt ein Entscheidungsmodell für die Zugriffskontrolle von Diensten vor, die auf einer gemeinsam genutzten IT-Infrastruktur (shared infrastructure) erbracht werden. Ziel ist hierbei die Maximierung des Ertrags eines IT-Dienstleisters bzw. die effiziente Erfüllung von Dienstleistungsverträgen durch dynamische Priorisierung von Diensten in Phasen von Ressourcenengpässen.

- Opportunitätskostenüberlegungen dienen als Kriterium zur Entscheidung über direkte Annahme, Pufferung oder Ablehnung einer Dienstanfrage.
- Das Modell berücksichtigt kombinatorische Effekte, die sich aufgrund der Belegung mehrerer Ressourcen durch Dienstauführungen und deren Interdependenzen ergeben.
- Das darauf aufbauende Zugriffskontrollverfahren antizipiert zukünftige Nachfrage nach Diensten und korrespondierenden Ressourcenanforderungen und optimiert daher vorausschauend.
- Die Evaluierung des Verfahrens mittels Simulationen realitätsnaher Szenarien eines Anbieters von Media-On-Demand-Diensten zeigt eine wesentliche Verbesserung des erzielten Ertrags gegenüber Standardverfahren.

Stichworte: Zugriffskontrolle, Überlastkontrolle, Ertragsmanagement, Service-Level-Management

Prozesse auf verwendeten Servern. Allgemein gilt es zu entscheiden, wie viele Einheiten einer oder mehrerer Ressourcen für bevorzugte Dienste reserviert werden sollen, um die Dienstleistungsvereinbarungen weitgehend erfüllen zu können. In der Regel besteht das Service-Portfolio eines Diensteanbieters aus mehreren Diensten mit unterschiedlichen Erträgen, Nachfrageverteilungen und Anforderungen an verschiedene IT-Ressourcen, was zu komplexen Zugriffskontrollproblemen führt. Zugriffskontroll- und Priorisierungsverfahren wurden in der Vergangenheit vor allem im Bereich der Netzwerktechnik eingesetzt. Auf Anwendungsebene bieten Hersteller von Datenbank-, Application- oder OLTP-Servern oft einfache statische Regeln an. Beispielsweise werden ab bestimmten Auslastungen der CPU niedrig priorisierte Dienstanfragen abgelehnt. Die a priori Fixierung dieser Schwellwerte stellt allerdings ein grundsätzliches Problem dar und wird meist aufgrund von Erfahrungswerten oder basierend auf Schätzungen vorgenommen. Regeln dieser Art sind oft suboptimal und vernachlässigen wichtige Aspekte des Zugriffskontrollproblems.

Dieser Beitrag beschreibt einen neuen Ansatz zur dynamischen Priorisierung von Anfragen basierend auf Erträgen von Diensten, Dienstleistungsvereinbarungen, erwarteten Ressourcenanforderungen von Diensten sowie aktuellen Systemauslastungen. Antizipation zukünftigen Bedarfs an IT-Ressourcen findet in der Modellformulierung ebenso Berücksichtigung wie die kombinatorischen Effekte aufgrund temporärer Belegung verschiedener (unterschiedlich knapper) Ressourcen durch mehrere Dienste. Im Kern wird anhand eines Vergleichs von Dienstpreisen bzw. -werten, SLAs und den mit Dienstannahmen verbundenen Opportunitätskosten über Akzeptanz, Pufferung oder Ablehnung von Dienstanfragen entschieden. Eine Pufferung oder Ablehnung einer Dienstanfrage entspricht einer Reservierung von Kapazität für zeitnah erwartete Anfragen nach Diensten höherer Priorität. Ziel ist die Maximierung des Gesamtertrags eines IT-Dienstleisters in Überlastphasen.

Der Artikel ist wie folgt aufgebaut: Kapitel 2 gibt einen Überblick über Zugriffskontroll- und Priorisierungsverfahren im Kontext von Informationssystemen sowie über Ertragsmanagementmodelle in verschiedenen Dienstleistungsbranchen. Kapitel 3 beschreibt die Formulierung des Optimierungsproblems. In Kapitel 4 wird die Effizienz des Modells mittels Simulation eines Praxisfalls analysiert. Als Anwendungsbei-

spiel wird in diesem Artikel die Zugriffskontrolle für Media-Streaming-Dienste dargestellt. Kapitel 5 fasst die Arbeit zusammen und gibt einen Ausblick auf weitere Forschungsfragen.

2 Literaturübersicht

Zugriffskontroll- und Priorisierungsverfahren zur Vermeidung bzw. zur effizienten Behandlung von Überlastsituationen in IT-Infrastrukturen wurden bereits intensiv im Bereich der Netzwerktechnik, insbesondere in Telekommunikationsnetzwerken, analysiert [KiMP05; DeGo04; DeGo01]. Diese Verfahren sind allerdings für die Ressourcenallokation auf Anwendungsebene wenig geeignet, da sie ausschließlich die Charakteristika der Daten- und Sprachübertragung über Netzwerkkomponenten modellieren, sich meist nur auf eine knappe Ressource beschränken und Wirtschaftlichkeits- bzw. Ertragsüberlegungen nicht berücksichtigen.

2.1 Bisherige Ansätze

Einige Arbeiten haben sich mit Zugriffskontrollverfahren für Webserver beschäftigt. So gibt es eine Reihe von Vorschlägen wie bei Überlast ganze HTTP-Sessions einzelner Nutzer zugelassen oder abgelehnt werden können (session-based admission control) anstatt einzelner HTTP-Anfragen aller Nutzer [ChPh02; ChMC01; RuCT05]. Dadurch kann auch bei Überlastsituationen die Leistung für bereits akzeptierte Nutzer aufrechterhalten werden. Eggert und Heidemann [EgHe99] schlugen Zugriffskontrollverfahren für Webserver mit zwei verschiedenen Nutzerklassen vor. Dabei wird zwischen Benutzer- und Proxy-initiierten (proxy prefetching) HTTP-Anfragen unterschieden. Diese und verwandte Arbeiten konzentrieren sich auf Server mit statischen Inhalten wie einfache Webseiten.

Singhmar et al. [SMAM02] unterteilen Anfragen auf E-Commerce-Webserver in Browsing- und Transaktions-Anfragen, wobei letztere mit höherer Wahrscheinlichkeit Umsatz generieren und daher zu bevorzugen sind. In Überlastphasen werden deshalb sämtliche Transaktions-Anfragen vor Browsing-Anfragen beantwortet, um die Gesamtrate der erfolgreichen Sessions zu erhöhen.

Differenzierende Zugriffskontrollalgorithmen für Multimedia-Dienste wurden von Chen et al. [ChCh96] untersucht. Die Autoren gruppieren Dienste in unter-

schiedliche Prioritätsklassen ein, die bei der Überschreitung definierter Grenzwerte bezüglich der Systemauslastung sukzessive abgelehnt werden. Statistische Zugriffskontrollverfahren für Multimedia-Server mit dem Ziel einer höheren Auslastung der Festplattensysteme unter mit hoher statistischer Sicherheit eingehaltenen Dienstgütern wurden beispielsweise von Vin et al. sowie von Kwon und Yeom analysiert [ViGG94a; KwYe00; ViGG94b].

Chen und Mohapatra betrachten Verfahren zur Priorisierung von Anfragen, die Transaktionen bzw. Sessions mit guten Aussichten bezüglich einer SLA-konformer Abarbeitung zugeordnet sind. Die Autoren ordnen Anfragen unterschiedlich gewichteten Warteschlangen zu (weighted fair sharing), wobei die einzelnen Gewichte dynamisch angepasst werden [ChMo02]. Zugriffskontrollmodelle in allgemeinen Warteschlangennetzwerken wurden unter anderem auch in [Stid99; LeAF02] analysiert.

Verma und Ghosal [VeGh03] stellen ein allgemeines Zugriffskontrollverfahren für IT-Dienstleister mit einer knappen Ressource vor. Hierbei werden die prognostizierte Dauer einer Dienstanfrage, der Ertrag bei Einhaltung einer definierten Antwortzeit sowie zu entrichtende Vertragsstrafen bei Nichteinhaltung der Antwortzeit betrachtet.

Die Ertragsmaximierung von Anwendungen mit unterschiedlichen Antwortzeitgarantien, die auf einer Menge von Servern installiert sind, wurde eingehend von Liu, Squillante und Wolf [LiSW01] analysiert. Die Autoren verwenden ein Warteschlangennetzwerk zur Modellierung der Systemarchitektur. Die Optimierungsmodelle maximieren den Ertrag unter Berücksichtigung der zu entrichtenden Vertragsstrafen bei Nichteinhaltung der vereinbarten Antwortzeiten. Es wird allerdings davon ausgegangen, dass der durchschnittliche Ressourcenbedarf für alle Anwendungen geringer ist als die zur Verfügung stehenden Ressourcen, bei optimaler Allokation also keine Überlastsituation vorliegt.

Die aufgeführten Ansätze treffen Zugriffskontroll-Entscheidungen basierend auf einer Auslastungs-Kennzahl des Gesamtsystems bzw. einer Systemressource, welche die maximale Leistung des Systems begrenzt (Flaschenhals). In Infrastrukturen, bei denen je nach Nachfrageintensität nach den verschiedenen, angebotenen Diensten unterschiedliche Ressourcen die Gesamtleistung des Systems limitieren, kann der Gesamtertrag eines IT-Dienstleisters in Überlastphasen jedoch nur maximiert werden, falls diese unterschiedlichen

Ressourcen in der Entscheidung auch berücksichtigt werden.

2.2 Ertragsmanagement

Um diese Lücke zu schließen stellen wir in diesem Artikel ein neues Verfahren vor, welches kombinatorische Effekte mit in die Entscheidungsmodellierung einbezieht. Einige Überlegungen des in nachfolgenden Kapiteln vorgestellten Entscheidungsmodells stammen aus der ausgedehnten Literatur im Bereich des Ertragsmanagements [Simp89; Will92; TaRy99].

Ertragsmanagementmodelle behandeln den optimalen Einsatz stochastisch nachgefragter, knapper Ressourcen. Ein Großteil der Modelle konzentriert sich auf die speziellen Gegebenheiten im Bereich der Fluglinien und der Hotelbranche, wo sie heute in vielen Fällen erfolgreich eingesetzt werden [BoBi03]. Die Angebotsseite von Hotels und Fluglinien ist typischerweise durch hohe Fixkosten und geringe marginalen Kosten gekennzeichnet. Die Nachfrager sind heterogen in Bezug auf ihre Kaufpräferenzen und ihre Preissensitivität und lassen sich dementsprechend segmentieren. Voraussetzung für die Anwendung solcher Modelle ist, dass unterschiedliche Kunden mit denselben Ressourcen bedient werden können.

Die benötigten Entscheidungsmodelle lassen sich nach ihrer Fristigkeit in taktische bzw. strategische und in operative Modelle unterscheiden. Zu den taktisch/strategischen Entscheidungen gehören Fragen nach der Segmentierung der Kundenbasis, der Preissetzung und den benötigten Kapazitäten der Ressourcen. Wir konzentrieren uns nachfolgend auf Ertragsmanagement auf operativer Ebene. Diese reservieren Kapazität für Dienstnehmer mit hoher Zahlungsbereitschaft und vergeben die restliche Kapazität an die übrigen Kundensegmente [Belo87; FKSR02; NaBa01]. Die Differenzierung erfolgt dabei nach qualitativen Merkmalen des Dienstes (Flexibilität bei der Buchung etc.).

„Network Revenue Management“ erweitert die klassischen Ansätze um die Betrachtung zusammenhängender Ressourcenanforderungen (beispielsweise Sitze in Flügen einer bestimmten Flugroute), also der von Kunden nachgefragten Güterbündel [McRy99]. Einige Überlegungen hieraus lassen sich auch für den Anwendungsbereich eines IT-Dienstleisters übernehmen. IT-Ressourcen können zur Befriedigung der Nachfrage in unterschiedlichen Kundensegmenten bzw. Dienstklassen eingesetzt werden – Dienste konsumieren also

temporär verschiedene Ressourcen (Ressourcenbündel) und ungenutzte Kapazität entspricht ebenso entgangenem Erlöspotenzial. Auch hier steht der Dienstanbieter in Phasen hoher Last auf einer oder mehreren seiner Infrastrukturkomponenten vor der Aufgabe, beschränkte Ressourcen verschiedenen Dienststanfragen, bzw. verschiedenen Klassen von Dienststanfragen, optimal zuzuweisen.

In der neueren Literatur schlagen Wissenschaftler auch Modelle für Internet-Zugangsanbieter vor, bei denen der Zugriff auf eine Menge homogener Ressourcen in Form von Modems geregelt wird [NaBa01]. Die Bereitstellung von IT-Dienstleistungen auf gemeinsam genutzten Infrastrukturen führt allerdings zu einer Reihe weiterer Annahmen, die nach neuen Modellierungsansätzen verlangen. So werden bei IT-Dienstleistungen ganze Bündel heterogener IT-Ressourcen (CPU, Hauptspeicher, I/O etc.) nachgefragt. Im Gegensatz zum Ertragsmanagement bei Fluglinien erfolgt die Nachfrage nach IT-Ressourcen meist kontinuierlich. Bei Fluglinien werden die Ressourcen (Sitze in bestimmten Flügen) ausschließlich innerhalb fixer Zeitintervalle verwendet (Start- bis Ankunftszeit des Fluges). IT-Dienstleistungen werden hingegen zu beliebigen Zeitpunkten in Anspruch genommen ohne vorher reserviert zu werden. Dieser Sachverhalt erfordert kurzfristige Entscheidungen über die Ressourcenallokation.

3 Problemformulierung

Wie bereits in dem vorangegangenen Kapitel beschrieben, entstehen je nach Anwendungssituation sehr unterschiedliche Problemstellungen die nach speziell angepassten Lösungen verlangen. Nachfolgend schlagen wir ein einfaches Klassifikationsschema vor, mit welchem die wichtigsten Problemklassen charakterisiert werden können, und entwickeln Entscheidungsmodelle für ausgewählte Klassen.

3.1 Klassifikation von Zugriffskontrollproblemen

Das Klassifikationsschema besteht aus einem 3-Tupel, welches den Modus der Abarbeitung, die Dauer der Jobs bzw. Dienstausführung sowie die Anzahl verfügbarer Ressourcen beschreibt:

– $d/e/1$ -Probleme beschreiben die Klasse von Problemen, bei denen Anfragen zu

diskreten Zeitpunkten ((d) iscrete) bearbeitet werden, die Bearbeitungsdauern aller Anfragen identisch sind ((e) qual length) und nur einen (1) Ressourcentyp belegen. Klassische Ertragsmanagementmodelle wie EMSR [Belo87] stellen eine Möglichkeit dar, um diese Probleme zu lösen.

– $c/u/n$ -Probleme charakterisieren die Problemklasse, die durch kontinuierliche Abarbeitung ((c) ontinuous), unterschiedliche Ressourcenanforderungen (unterschiedlich lange Bearbeitungsdauern ((u) nequal length)) und durch n heterogene Ressourcen charakterisiert wird.

Diese Notation kann noch um einige relevante Parameter (z. B. Verteilungsangaben und Zielfunktionen) erweitert werden, genügt aber, um die in diesem Artikel beschriebenen und bisher bekannten Methoden zu kategorisieren. Viele der in Abschnitt 2 beschriebenen technischen Lösungen konzentrieren sich auf $c/u/1$ -Probleme.

Die beschriebenen $d/e/n$ -Probleme können durch das nachfolgend beschriebene Modell DLP (Deterministisches Lineares Programm) gelöst werden. Im Anschluss daran stellen wir ein hierauf aufbauendes, erweitertes Modell DLPc (DLP continuous) vor, das bei $c/u/n$ -Problemen, wie sie bei IT-Dienstleistungen mit gemeinsam genutzten Ressourcen vorliegen, in Zugriffskontrollmethoden verwendet werden kann.

3.2 Basismodell und Opportunitätskosten

Nachfolgend beschreiben wir ein grundlegendes Zugriffskontrollmodell für gemeinsam genutzte IT-Infrastrukturen, welches einer Reihe restriktiver ($d/e/n$ -) Annahmen unterworfen ist. Ausgehend von diesem Basismodell werden wir Modellerweiterungen und Heuristiken vorstellen, um relevante Realweltbedingungen mit abzubilden.

Der Dienst-Katalog eines IT-Dienstleisters bestehe aus I Diensten i ($i = 1, \dots, I$), die zu diskreten Zeitpunkten t_k ($k = 0, \dots, \infty$) stochastisch nachgefragt werden. Eine Nachfrageverteilung D_i habe hierbei den Erwartungswert μ_i . Der Ertrag, der mit einer erfolgreichen Erbringung eines angefragten Dienstes einhergeht, sei r_i . Die Belegungsdauer der während der Dienstleistung verwendeten Ressourcen und die Dauer der Dienstleistung selbst seien von konstanter Länge Δt ($\Delta t = t_{k+1} - t_k$) und somit zum jeweils nächstmöglichen

diskreten Nachfragezeitpunkt t_{k+1} abgeschlossen. Die Ressourcenbelegungskoeffizienten a_{ei} geben die Anforderungs- bzw. Verwendungsmenge an Einheiten der Ressource e ($e = 1, \dots, E$) für die Zeitdauer Δt durch die Ausführung von Dienst i an. Eine Ressource e hat eine begrenzte Kapazität C_e . Die Ermittlung von Ressourcenbelegungskoeffizienten a_{ei} einzelner Dienste (CPU-Zyklen, Hauptspeicher in Bytes, I/O in Blocks etc.) kann über Messungen in isolierten Testumgebungen, wie diese bei Lasttests und Softwareabnahmen zum Einsatz kommen, mit hinreichend hoher Genauigkeit erfolgen [Bran06]. Bild 1 veranschaulicht die genannten Zusammenhänge.

Die zur Verfügung stehenden Ressourcenkapazitäten sollen nun so für die Beantwortung von Dienstanfragen eingesetzt werden, damit der Gesamtertrag maximiert wird. Unter den gegebenen Annahmen kann das Problem durch folgendes Ganzzahliges Lineares Programm (Integer Program, IP) modelliert werden:

$$\begin{aligned} \max \quad & \sum_{i \in I} r_i \cdot x_i \\ \text{s.t.} \quad & \sum_{i \in I} a_{ei} x_i \leq C_e \quad \forall e \in E \\ & x_i \leq D_i \quad \forall i \in I \\ & x_i \in \mathbb{Z}_+ \quad \forall i \in I \end{aligned} \quad (\text{IP})$$

Die ganzzahlige Variable x_i beschreibt die Anzahl der zu akzeptierenden Anfragen nach Dienst i für eine Zeitperiode Δt . Für die Zufallsvariable D_i wird in dieser Formulierung deren Erwartungswert als deterministische Größe verwendet. Nachdem IP ein NP-vollständiges Problem darstellt, kann hier auch die LP-Relaxation des IP, das sog. Deterministischen Linearen Programm (DLP), gelöst werden [Will92].

Die dualen Variablen λ_e der Kapazitätsrestriktionen der LP-Relaxation können ökonomisch als Schattenpreise oder Opportunitätskosten der Verwendung einer Ressourceneinheit interpretiert werden. Die gesamten Opportunitätskosten einer Anfrage nach Dienst i können durch Addition der Produkte aus Ressourcenbelegungskoeffizient und Opportunitätskosten pro Ressourceneinheit ($\sum_e a_{ei} \lambda_e$) berechnet werden [TaRy99]. Es werden ausschließlich solche Dienstanfragen akzeptiert, bei denen der Ertrag die Opportunitätskosten der Dienstanfrage übersteigt. Wie auch Analysen vergleichbarer Modelle in der Ertragsmanagementliteratur für Fluglinien zeigten, erzielte das DLP-Modell bei IT-Dienstleistungen in Simulationen unter oben getroffenen Modellannahmen bei zuverlässigen Nachfrageprognosen sehr gute Ergebnisse [TaRy99; BoFP99].

3.3 Berücksichtigung kontinuierlicher Dienstanfrage

Die im vorhergehenden Abschnitt formulierte Basismodellvariante DLP unterstellt Dienstanfragen zu diskreten Zeitpunkten, deren Bearbeitung zum nächsten diskreten Nachfragezeitpunkt abgeschlossen ist. Diese Annahmen sind nur in speziellen Fällen, wie beispielsweise bei Stapelverarbeitungsaufgaben (batch jobs) mit etwa gleicher Länge oder dem Bedienen verschiedener, abonnierbarer Kanäle im Bezahlfernsehen, gegeben. Sind diese Annahmen nicht gegeben, sinkt, wie auch Simulationen zeigen, die Effizienz des Modells. IT-Dienstleister sehen sich in der Realität meist mit kontinuierlicher Dienstanfrage konfrontiert und verschiedene Dienste haben unterschiedliche Ressourcenanforderungen sowie unterschiedlich lange Bearbeitungsdauern. Diese Art von Zugriffskontrollproblemen stellt somit eine Verallgemeinerung des oben beschriebenen Netzwerk-Ertragsmanagementproblems dar.

Für $c/u/n$ -Probleme ist die Ressourcenverbrauchsmatrix um eine zeitliche Dimension t_{ei} zu erweitern (vgl. Bild 2). Diese gibt die Anzahl an Zeiteinheiten an, die ein

Dienst i für a_{ei} Einheiten einer Ressource e belegt (bei manchen Anwendungen ist darüber hinaus relevant, in welcher Reihenfolge die Ressourcen angefordert werden). Ressourcen können in einem Moment nahezu ausgelastet sein und im nächsten Moment, nach Beendigung aktiver Dienstanfragen, wieder fast in maximaler Kapazität zur Verfügung stehen. Ziel der Zugriffskontrollmodelle ist es, abhängig von verfügbaren Ressourcenmengen, Anfragen auf Dienste niedriger Priorität vorausschauend abzuweisen um Ressourcen für Anfragen nach Diensten höherer Priorität zu reservieren.

Die Planungsintervalle werden dabei möglichst kurz gehalten um einer möglichst exakten Bestimmung aktueller Ressourcenauslastungen Rechnung zu tragen. Die Konsequenz für die praktische Umsetzung ist, dass die Neuberechnung der Schattenpreise im DLPc-Modell zu jedem Dienstanfragezeitpunkt vorgenommen wird. Die Planungshorizonte werden dabei jeweils auf die prognostizierten Endzeitpunkte der Ressourcennutzung durch die angefragten Dienste gelegt. Die Berechnungsdauer der Opportunitätskosten lag bei Szenarien mit 100 Dienstypen und 10 Ressourcen auf einem Pentium-III-2 GHz-Prozessor bei unter 10 ms.

Für die Neuberechnung der Opportunitätskosten sind sowohl die aktuelle Auslastung der vom Dienstyp verwendeten Ressourcen zu bestimmen, als auch die Anzahl der weiteren im Planungsintervall noch zu erwartenden Dienstanfragen. DLP verwendet als Kapazitätsrestriktionen die gesamten Ressourcenkapazitäten C_e , da zu den diskreten Berechnungszeitpunkten für das nächste Planungsintervall keine Dienste aktiv, und somit sämtliche Kapazitäten verfügbar sind. Bei kontinuierlicher Dienstanfrage mit unterschiedlichen Bearbeitungsdauern können Ressourcen jedoch im Moment einer Schattenpreisberechnung durch aktive Dienste teilweise belegt sein. Diese belegten Ressourcenkapazitäten stehen für die weiteren im Planungsintervall angefragten Dienste erst wieder nach Erbringung der aktiven Dienste, und damit nach Freigabe der belegten Ressourcen, zur Verfügung.

Die Ermittlung der Ressourcenauslastungen zum Dienstanfragezeitpunkt, beispielsweise über System-Monitoring-Werkzeuge, erbringt die aktuell verfügbaren Kapazitäten. Die Parametrisierung des Optimierungsprogramms mit diesen Kapazitätsrestriktionen liefert jedoch zu hohe Schattenpreise, da dies unterstellt, dass die derzeit verwendeten Ressourcen für den gesamten

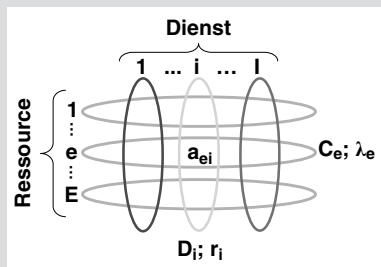


Bild 1 Ressourcenbelegungsmatrix

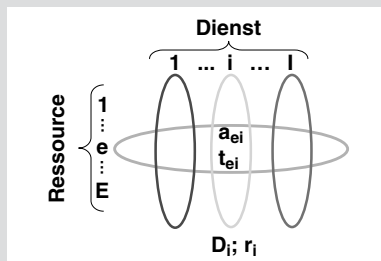


Bild 2 Zeitdauer und Menge der belegten Ressourcen durch Dienst i

Planungszeitraum belegt bleiben. Werden Dienstanfragen noch während eines Planungsintervalls beendet, so stehen die frei werdenden Ressourcen in demselben Planungsintervall wieder zur Verfügung.

Diese grundlegende Problematik wird in DLPC durch eine Heuristik adressiert, die in Simulationen zu guten Ergebnissen geführt hat. Diese Heuristik besteht aus mehreren Berechnungsschritten, bei denen die tatsächlich im Planungsintervall zur Verfügung stehenden Kapazitäten und der Ressourcenbedarf im Planungsintervall näherungsweise bestimmt werden. Dieses Approximationsverfahren wird im Folgenden beschrieben.

Während des gesamten analysierten Zeitraums, also der Sequenz aller betrachteten Planungsintervalle, treffen insgesamt K Dienstanfragen ein, wobei jede einzelne Dienstanfrage k mit $k = 1, \dots, K$ einem bestimmten Dienst bzw. einer Dienstklasse i zugeordnet ist. Zum Zeitpunkt t_k einer Anfrage k nach Dienst i setzt man für alle von i benötigten Ressourcen e deren Planungszeiträume auf die voraussichtliche Dauer der Ressourcenbelegung durch Anfrage k ($t_{ek} = t_{ei}$) (vgl. Bild 3). Danach werden folgende Schritte für alle Ressourcen e mit $a_{ei} > 0$ durchgeführt:

- Die prognostizierten Restlaufzeiten $l_{ek'}$, der Belegung der Ressource e ab t_k durch derzeit aktive Dienstanfragen k' werden bestimmt. Die Ressourcenbelegungsdauern aktiver Dienste werden so um den bereits abgearbeiteten und daher nicht entscheidungsrelevanten Zeitanteil vor t_k gekürzt.
- Für alle k' , deren prognostizierte Endzeitpunkte der Verwendung der Ressource e nach dem Zeitpunkt $t_k + t_{ek}$ liegen, wird deren in das Modell eingehende Restlaufdauer $l_{ek'}$ auf das entscheidungsrelevante Intervall $l_{ek} = [t_k, t_k + t_{ek}]$ begrenzt.
- Die Summe $\sum_{k'} a_{ek'} l_{ek'}$ aller aktiven Dienste k' entspricht den im Betrachtungszeitraum l_{ek} nicht mehr für die aktuelle Dienstanfrage und neue Dienstanfragen zur Verfügung stehenden Kapazitätseinheiten von e . Diese Summe wird von der in l_{ek} theoretisch maximal verfügbaren Kapazität $l_{ek} C_e$ (C_e entspricht der Kapazität von e pro Zeiteinheit) subtrahiert und ergibt die im gesamten Intervall noch verfügbare Kapazitätseinheiten C_{ek} . Dieser Wert wird als Annäherung für die im Planungszeitraum verfügbare Kapazität verwendet.

C_{ek} gibt die in den Ressourcenplanungsintervallen l_{ek} verfügbare Kapazitätseinheiten an, die für verschiedene Ressourcen

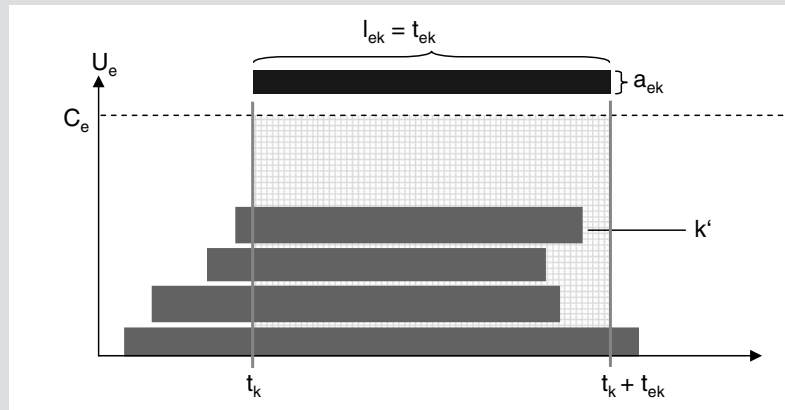


Bild 3 Approximation zukünftiger Ressourcenkapazität

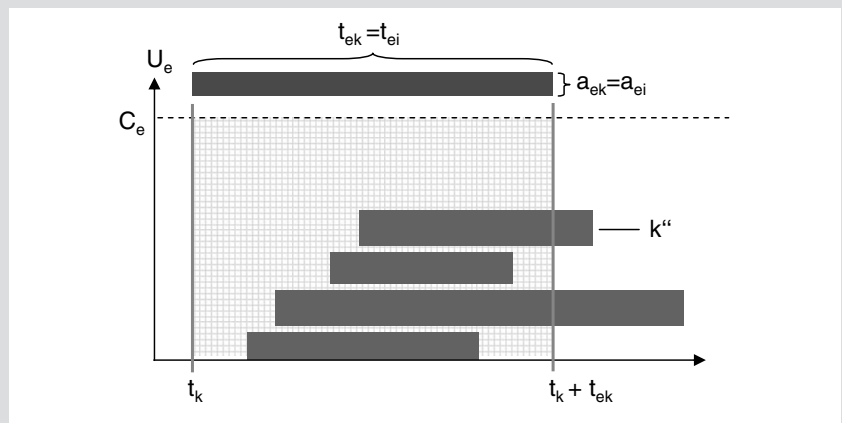


Bild 4 Approximation zukünftiger Ressourcenbelegung

e unterschiedlich sein können. Zur Ableitung der Schattenpreise λ_e aus den dualen Variablen der Kapazitätsrestriktionen sind auch hier wieder die zu akzeptierenden Anfragen x_i zu ermitteln, nun in Abhängigkeit von der, zur Länge des gesamten Planungsintervalls l_{ek} proportionalen, erwarteten Nachfrage nach Diensten. Wählt man als Gesamtintervall das längste Intervall $\max(l_{ek})$, so sind die für die kürzeren Intervalle verfügbaren Kapazitäten entsprechend hochzurechnen:

$$C_{ek}^{\max(l_{ek})} = \frac{\max(l_{ek})}{l_{ek}} C_{ek}$$

Analog zu der beschriebenen Beschränkung durch bereits akzeptierte Dienstanfragen sind auch die Anforderungen der im Intervall l_{ek} noch erwarteten Dienstanfragen k'' zu berücksichtigen. Nimmt man beispielsweise an, dass mit hoher Wahrscheinlichkeit im nächsten Zeitintervall weitere 100 MB Hauptspeicher durch eine Dienstanfrage beansprucht werden, so

sollte dies bei der Kapazitätsabschätzung mit berücksichtigt werden. Auch diese Anfragen werden auf das Ende der Planungsperiode $t_k + t_{ek}$ beschränkt (vgl. Bild 4), da Ressourcenbelegungen nach diesem Zeitpunkt für die Annahme-/Ablehnungsentscheidung der Anfrage k nicht mehr relevant sind. Es sind somit für alle Dienste und Ressourcen Korrekturfaktoren q_{ei} festzulegen, welche die relevanten Anteile der Belegungsdauern einer Ressource durch einen Dienst im Planungsintervall angeben. Die Ressourcenanforderungen durch erwartete Dienste k'' werden in analoger Weise approximiert wie die Anforderungen der aktuell bearbeiteten Dienste k' .

DLPC erweitert DLP um die beschriebenen Heuristiken zur Bestimmung verfügbarer Kapazitäten sowie zur Bestimmung der Ressourcenanforderungen im Planungsintervall. Die Prognose erwarteter Dienstanfragen in einem Planungsintervall kann hierbei entweder über Zeitreihenanalysen vergangener Nachfrage oder dyna-

misch, beispielsweise über die Bildung eines gleitenden Durchschnitts aus den letzten Planungsperioden, oder aus einer Kombination der beiden Ansätze vorgenommen werden.

$$\begin{aligned} \max \quad & \sum_{i \leq I} r_i \cdot x_i \\ \text{s.t.} \quad & \sum_{i \leq I} q_{ei} a_{ei} t_{ei} x_i \leq C_{ek}^{\max(l_{ek})} \quad \forall e \leq E \\ & x_i \leq \max(l_{ek}) D_i \quad \forall i \leq I \\ & x_i \in \mathbb{Z}_+ \quad \forall i \leq I \end{aligned} \quad (\text{DLPc})$$

DLPc betrachtet Ressourcenbelegungskoeffizienten sowohl hinsichtlich Zeitdauer als auch Quantität pro Zeiteinheit als deterministische Größen. Die von uns im Labor durchgeführten Messungen der Ressourcenanforderungen (CPU-Zeit, Hauptspeicherbelegung und I/O) von Web- und Media-Streaming-Diensten wiesen, auch bei starker Veränderung der Arbeitslast, eine geringe Varianz auf. Auch Elnikety et al. [ENTZ04] zeigten durch Messungen, dass die Varianzen der Ressourcenbelegungen auf einem Server in Phasen unvollständiger Auslastung des Servers gering sind.

Kriterien für den Einsatz entsprechender Modelle sind zum einen die hinreichend genaue Prognostizierbarkeit zukünftiger Dienstanfrage sowie Kenntnisse über aktuelle Ressourcenauslastungen, zum anderen Kenntnisse über Ressourcenanforderungen der einzelnen Dienste. Die Prognose zukünftiger Nachfrage basiert in der Regel auf Kenntnissen historischer Nachfragemuster bzw. dem Vorhandensein entsprechender Transaktionsdaten. Das Monitoring der IT-Infrastruktur über System-Monitoring-Werkzeuge wird meist standardmäßig durchgeführt. Ressourcenanforderungen einzelner Dienste können durch Messungen beispielsweise bei Abnahmetests ermittelt werden. Sollte dies nicht möglich sein, können oft anhand von Log-Dateien Rückschlüsse auf die Ressourcenanforderungen einzelner Dienste ermittelt werden. Generell müssen die Gegebenheiten des jeweiligen Anwendungskontexts beim Einsatz des Modells berücksichtigt werden. Beispiele hierfür sind Optimierungen (Indexierung, Cache etc.), die von Hard- und Softwarekomponenten selbständig und fortlaufend vorgenommen werden.

4 Evaluierung

c/u/n-Zugriffskontrollprobleme finden sich überall dort, wo unterschiedliche Dienste

um gemeinsame Ressourcen konkurrieren. Anwendungsmöglichkeiten ergeben sich beispielsweise beim Zugriff auf gemeinsam genutzte Datenbanken, in Call-Centern, bei virtuellen Servern, beim Zugriff auf Web-Applikations-Server sowie bei Media-Streaming-Servern im Falle von Video- oder Audio-On-Demand. Letztgenanntes Szenario verwenden wir zur Evaluierung des vorgestellten Entscheidungsmodells.

Nach einer im März 2006 im Auftrag des Bundesministeriums für Wirtschaft und Technologie herausgegebenen Studie „Gesamtwirtschaftliche Auswirkungen der Breitbandnutzung“ stieg die Anzahl existierender Breitband-Internetzugänge in Deutschland von 1,9 Millionen im Jahr 2001 auf 10,7 Millionen Ende des Jahres 2005. Es wird von einer zukünftigen Fortsetzung dieses Trends ausgegangen, wobei zudem die Bandbreiten pro Internetzugang ebenfalls noch weiter steigen werden [FoOB06].

Die zunehmende Verfügbarkeit hoher Bandbreiten sowie das Vorhandensein effizienter Audio- und Videokomprimierungstechnologien wie MP3-Pro oder H.264 ermöglichen die Bereitstellung von interaktiven Medien-Diensten wie Audio- und Video-On-Demand (AoD, VoD), Online-Spiele, Video-Konferenzen oder auch Internet-Telephonie. Zahlreiche DSL-Internet-Anbieter und Betreiber von Kabelnetzen bieten bereits heute einen Teil oben genannter Zusatzdienstleistungen an. Bei großen internationalen Hotelketten sind solche Dienste meist bereits Standard, wobei hierfür großteils auf Standard-Plattformen wie *Microsoft TV*, *Siemens Surpass Home Entertainment*, *NXTV* oder *Deuro-media* zurückgegriffen wird.

Die Nachfrage nach solchen Diensten im Internet entwickelt sich dynamisch. Ebenso sind Produkte in der Unterhaltungsindustrie von starken Nachfrageschwankungen gekennzeichnet (Erscheinen neuer Filme oder Musikalben). Dementsprechend anspruchsvoll gestalten sich das Kapazitäts- und das Service-Level-Management in solchen Infrastrukturen, in denen man einerseits möglichst wenig unwirtschaftliche Überkapazität bereithalten will, andererseits die hohen Echtzeitanforderungen der Dienste soweit wie möglich einhalten will. Verzögerungen und Stockungen aufgrund von Überlast bei Video-On-Demand sind beispielsweise vom Kunden unmittelbar wahrnehmbar und wirken sich direkt auf die Kundenzufriedenheit aus. Aus diesem Grund kommt der effizienten Behandlung bzw. Vermeidung von Überlast eine besondere Bedeutung zu. Darüber

hinaus stellen die verschiedenen Dienstypen sehr unterschiedliche Anforderungen an verschiedenen IT-Ressourcen wie Netzwerkbandbreite, Durchsatz und Zugriffszeit auf die Festplatte bzw. die Speichersysteme, die CPUs des Multimedia-Servers sowie dessen Hauptspeicher und Cache. Im Gegensatz zur Historie in der Multimediakommunikation, bei der auf Grund wesentlich schmalere Netzwerkbandbreiten Engpässe in der Regel in den Netzwerkkomponenten zu finden waren, können heutzutage je nach den jeweiligen Dienst-Nutzungsintensitäten (Spielfilme, Video-Clips, Audio-Dateien) unterschiedliche Ressourcen Engpässe darstellen [ChTV04; ChSt03].

Hierzu folgendes Beispiel: Ein IT-Dienstleister stellt neben längeren Videos mit hohen Bandbreiten kurze Videoclips sowie Audio-Dateien mit geringeren Bandbreiten bereit. Während eines Tages wird ein Mix der verschiedenen Medien nachgefragt. Stößt man mit zunehmender Anzahl paralleler Breitband-Videos an die Grenze des maximalen Durchsatzes des Festplattensystems, stellt dieses den Systemengpass dar. Am frühen Morgen wird verstärkt auf Nachrichten und Änderungsmeldungen zugegriffen, also auf die Medien mit den geringeren Anforderungen bzgl. Bandbreite. Da es sich hierbei um kleinere Dateien handelt, die teilweise oder komplett in den Cache bzw. den Hauptspeicher des Medien-Servers geladen werden können, wird seltener auf die Festplattensysteme zugegriffen. Dies führt dazu, dass auf Serverseite CPU und Hauptspeicher die primären Flaschenhälse darstellen. Je nach Nachfrage-Mix können somit unterschiedliche Ressourcen zum Systemflaschenhals werden. Diese Problematik beschreibt [ChTV04] ausführlich bei Lasttests mit Media-Streaming-Servern.

4.1 Simulation

Um das in der Simulation verwendete Szenario möglichst einfach und anschaulich zu halten, wurde eine Streaming-Infrastruktur zur ausschließlichen Bereitstellung von Audio- und Video-On-Demand-Diensten verwendet. Das Angebot des Dienstleisters bestand aus Filmen der Preiskategorien VoD-Premium (neue Filme und Blockbuster) und VoD-Standard (ältere Filme). Die Audio-Dateien entsprachen Musik-Alben und wurden in die beiden Preiskategorien AoD-Premium (Album aus den aktuellen Top 100) oder AoD-Standard (sonstige Titel) unterteilt.

Als Media-Streaming-Server wurde in den Simulationen ein 3,8 GHz Intel Pentium 64-Bit-Server mit 2 GB DDR2 SDRAM Hauptspeicher verwendet. Ausgestattet war der Server mit einem RAID-5-Verbund aus Fast-ATA-Festplatten mit 7200 rpm und jeweils 16 MB Cache sowie einer 1 Gbit/s-Fast Ethernet-Netzwerkverbindung. Als Software kam Windows 2003 64 Bit Enterprise Edition inklusive des im Umfang des Betriebssystems enthaltenen Windows Media Servers zum Einsatz.

Um die Ressourcenanforderungen der Audio- und Video-Streams zu ermitteln und Flaschenhals-Kandidaten des Servers zu identifizieren, wurden zunächst Lasttests und Messungen durchgeführt. Hierfür griffen Client-Computer auf Inhalte des Media-Servers zu, um die maximale Anzahl paralleler Streams zu ermitteln, mit welchen die Medien ohne Qualitätsverluste übertragen werden konnten. Der Cache des Media Servers wurde nicht aktiviert.

Als Medien-Abspielsoftware wurde auf den Client-Rechnern der Microsoft Windows Media Player (Version 10) eingesetzt. Die Abspielsoftware war so konfiguriert, dass zunächst eine Sekunde der Streams in den Puffer der Abspielsoftware geladen wurde, bevor das Abspielen begann. Die Abspieldauern der Videos betragen jeweils 90 Minuten, die von Musik-Alben 45 Minuten.

Bei ausschließlichem Zugriff auf Videos konnten von dem Festplattensystem maximal 269 Streams ohne Qualitätsverluste gelesen werden. Dies entspricht einem Durchsatz des Festplattenverbunds von 1.147.188 Kbit/s. Aufgrund dieses hohen Durchsatzwertes (zum Vergleich liegt der maximale Durchsatz eines modernen PCI-Bus 32 Bit/33 Mhz bei 1.064.000 Kbit/s) trat der Systemflaschenhals in der beschriebenen Konfiguration jedoch bereits zuvor bei der Netzwerkverbindung auf. So war die Anzahl paralleler Video-Streams auf 222 limitiert, was einer Allokation von Netzwerkbandbreite von 950.000 Kbit/s entspricht. Ein noch höherer Durchsatz als 950.000 Kbit/s (0,96 Gbit/s) ging in unseren Lasttests einher mit steigenden Fehleraten und Qualitätsproblemen einzelner Streams, sodass die maximal verfügbare Kapazität für dieses Media-Streaming-Szenario auf diesen Wert festzusetzen war. Die Festplattensysteme, der Hauptspeicher sowie die CPU des Media-Servers waren beim Streaming von 222 Video-On-Demand-Dateien in der gewählten Konfiguration keine Flaschenhälse.

Bei ausschließlichem Zugriff auf Audio-Inhalte betrug die maximale Anzahl paral-

Dienst	Preis (€)	Dauer (min)	Netzwerk (%)	Festplatten (%)
VoD-Premium	3,-	90	4262 Kbit/s 950000 Kbit/s	1 269
VoD-Standard	1,-	90	4262 Kbit/s 950000 Kbit/s	1 269
AoD-Premium	0,03	45	129 Kbit/s 950000 Kbit/s	1 6333
AoD-Standard	0,01	45	129 Kbit/s 950000 Kbit/s	1 6333

ler Streams 6333; Systemflaschenhals war hier jedoch nicht die Netzwerkanbindung, sondern der Durchsatz der Festplattensysteme. Aufgrund des Zugriffs auf viele kleinere Streams und der Notwendigkeit einer hochfrequenten, zeitintensiven Repositionierung des Magnetschreib- und -lesekopfs der Festplatte reduzierte sich deren Durchsatz auf 817.000 Kbit/s. Hier ist zu erwähnen, dass je nach bereitgestellten Medientypen, Systemkonfigurationen und eingesetzter Softwareplattform auch weitere Ressourcen zu Engpässen werden können.

Im Folgenden bezeichnen wir das Verhältnis der Summen der Bandbreitenanforderungen von Audio-Streams zu Video-Streams, welches sich durch unterschiedliche Nachfrage nach Audio- und nach Video-Diensten jeweils ergibt, als Audio-Video-Mix (AVM).

Zur Evaluation der Effizienz des vorgestellten Zugriffskontrollmodells DLPc wurden mehrere Simulationsrunden mit unterschiedlich hoher Dienst-Nachfrage und unterschiedlichen AVMs analysiert. Effizienzkriterium war jeweils der während der Überlastphase generierte Gesamtertrag unter Verwendung des vorgestellten Modells im Vergleich zu dem Ertrag, der sich durch Verwendung traditioneller Zugriffskontrollregeln ergab. Solche Regeln setzen Schwellwerte bezogen auf die Auslastung einer knappen Ressource, ab welchen nur noch bestimmte Dienste zugelassen werden (bei Auslastungswerten gegen 100 % werden sämtliche Anfragen abgelehnt).

Zur Parametrisierung, Durchführung und Auswertung der Simulationen wurde ein Simulationstool entwickelt, welches während der Simulationen die Ressourcenauslastungen, die aktuellen Schattenpreise sowie die Entscheidungen der Zugriffskontrollverfahren visualisiert. Das Programm generiert Last entsprechend Nachfrageverteilungsparametern, in dem es über entfernte Methodenaufrufe Streaming-Prozesse auf den eingesetzten Client-

Rechnern initiiert. Das Programm ermöglicht zudem das Abspeichern und Wiederherstellen von Simulationsparametern, das Aufzeichnen und Abspielen von Simulationsabläufen sowie den direkten Vergleich unterschiedlicher Zugriffskontrollmethoden unter identischen Bedingungen.

Eine Übersicht der Dienste, deren Preise und Ressourcenanforderungen zeigt Tabelle 1. Ressourcenanforderungen werden in Prozent bezüglich der Gesamtkapazität einer Ressource dargestellt. So beträgt die relative Last eines Video-Streams bzgl. des Festplattensystems beispielsweise 1/269, da von der Festplatte maximal 269 Videos gleichzeitig übertragen werden können. Der gemessene Verbrauch an Netzwerkbandbreite pro Dienst ist in Spalte Netzwerk (%) beschrieben.

In einer Simulation wurden jeweils 36 Zeitperioden (jeweils 10 Minuten) analysiert. Zu Beginn jeder Simulation war die gesamte Kapazität verfügbar; es waren keine Media-Streams bereits aktiv. Ausgehend von einer Nachfragehöhe, die gerade noch ohne Qualitätsprobleme von dem Server befriedigt werden konnte, wurde die Nachfrage in jeder Simulation um 10 % erhöht, um damit erhöhte Nachfrage (Übernachfrage) zu generieren, die ohne Zugriffskontrollverfahren zu steigenden Überlastsituationen auf dem Server führen würde. Die Simulationsrunden wurden für verschiedene AVMs durchgeführt. In den Simulationen wurden jeweils gleiche Mengen von Premium- und Standard-Diensten nachgefragt.

4.2 Ergebnisse

Die durchschnittlichen Ergebnisse der numerischen Experimente aus 36 Zeitperioden werden in den Bildern 5–8 dargestellt. Die Dienstinachfrage war während einer Simulation konstant und damit sehr gut prognostizierbar. Unterschiede in den

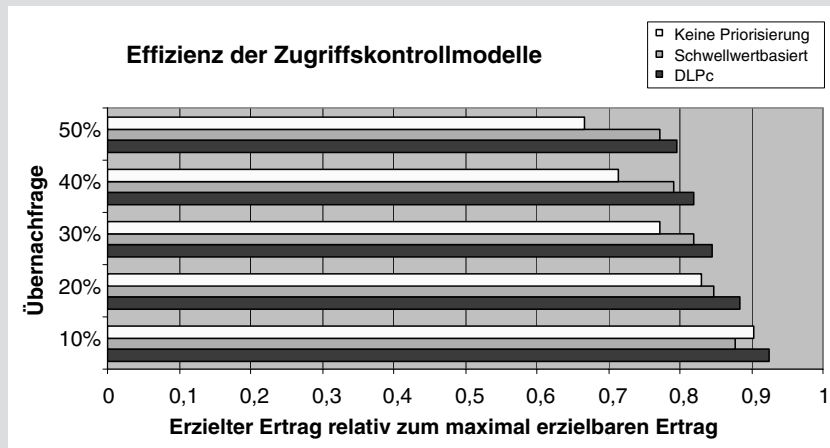


Bild 5 Ausschließlicher Zugriff auf Audio-Dateien

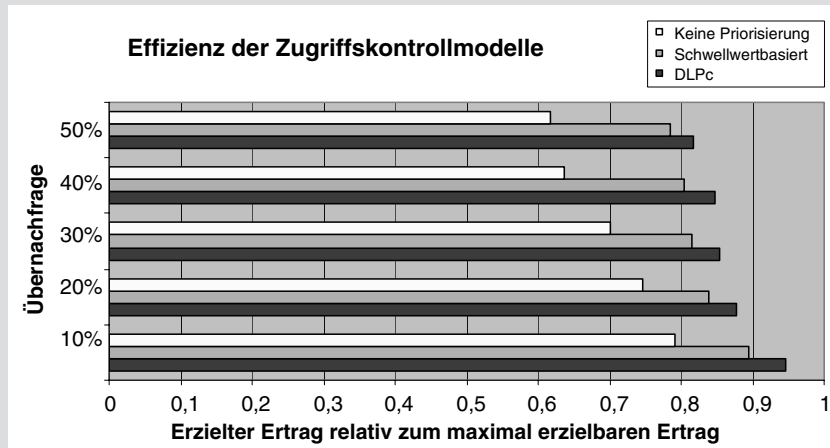


Bild 6 Audio-/Video-Nachfrageverhältnis = 3 : 1

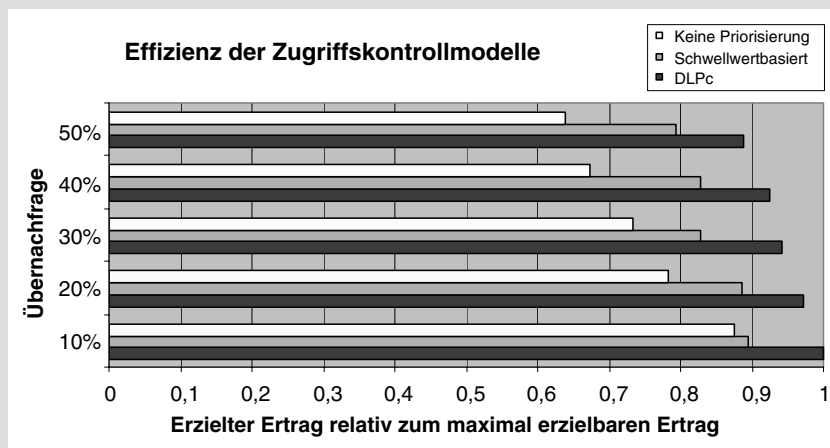


Bild 7 Audio-/Video-Nachfrageverhältnis = 1 : 3

Ergebnissen einzelner Zeitperioden ergeben sich daher vor allem aus systembedingten Varianzen im Eintreffen der Dienst-anfragen. Alle Unterschiede zwischen den Ergebnissen waren signifikant (basierend auf einem gepaarten *t*-Test mit $\alpha = 0,05$).

Auf den Bildern sind jeweils die mit verschiedenen Zugriffskontrollmodellen erzielten Erträge relativ zu dem Ertrag dargestellt, der maximal möglich gewesen wäre, falls ausreichend Kapazität zur Befriedigung aller Anfragen vorhanden gewesen wäre. Das Modell *Keine Priorisierung* erlaubt die Annahme aller eingehenden Anfragen solange ausreichend Kapazität zur Verfügung steht und lehnt ansonsten Anfragen ab. Modell *Schwellwertbasiert* erlaubt die Annahme sämtlicher eingehender Anfragen bis zu einer Systemauslastung von 80% und lässt ab diesem Wert ausschließlich Premium-Dienste zu. Modell *DLPc* trifft die Entscheidung wie in den vorangegangenen Abschnitten beschrieben anhand eines Vergleichs zwischen Dienstpreisen und Schattenpreisen. Als Kennzahl der Systemauslastung wurde in den ersten beiden Modellen der Durchsatz gewählt, wie dies auch in den von Windows Media Server angebotenen Überlast-Kontrollmöglichkeiten angeboten wird. Um Systemüberlast zu vermeiden war die Gesamtkapazität auf einen Wert festzulegen, dessen Einhaltung bei jedem beliebigen AVM eine reibungslose Übertragung aller aktuellen Media-Streams sicherstellt. Dieser Wert entspricht in dem gewählten Szenario wie oben beschrieben 817.000 Kbit/s. Das Modell *DLPc* konnte mit den tatsächlichen Kapazitätsgrenzen des Netzwerks und des Festplattensystems parametrisiert werden, da *DLPc* auch im Falle mehrerer potenzieller Flaschenhalse Systemüberlast verhindert.

DLPc erwies sich bei fast allen AVMs und unterschiedlicher Übernachfrage gegenüber den beiden statischen Modellen als überlegen. Beispielsweise lag bei einem AMV von 1:3 und einer Übernachfrage von 30% der Ertrag auf Basis der Preise in Tabelle 1 nach 6 Stunden bei *DLPc* bei 212 Euro im Vergleich zu 186 Euro bei der schwellwertbasierten Zugriffskontrolle und 165 Euro bei dem Zugriffsmodell ohne Priorisierung. Der Grund hierfür liegt zum einen in der Berücksichtigung kombinatorischer Effekte durch die Verwendung mehrerer Ressourcen und deren Abhängigkeiten, zum anderen an der vorausschauenden Reservierung von Kapazität für zeitnah erwartete, lukrativere Dienste.

So stieg die Überlegenheit des *DLPc* mit zunehmendem Anteil an Video-On-De-

mand-Diensten, da hierdurch die relative Auslastung monoton von den Festplattensystemen zur Netzwerkkomponente verschoben wurde. Hierdurch konnten mehr Daten übertragen und Dienste zugelassen werden als dies bei den statischen Modellen der Fall ist. So wurden bei einem Übernachfragewert von 10% bei hohem bzw. ausschließlichem Anteil an Video-On-Demand-Diensten, wie in Bild 7 und Bild 8 dargestellt, sämtliche Dienste akzeptiert. DLPC war aber auch, wie in Bild 5 dargestellt, bei ausschließlicher Nachfrage nach Audio-On-Demand überlegen, obwohl hier auch bei den statischen Verfahren die Kapazität bezüglich des Festplatten-Flaschenhalses korrekt angenommen wurde. Aufgrund der Antizipation zukünftigen Ressourcenbedarfs von Premium-Diensten und der einhergehenden vorausschauenden Reservierung von Kapazität wurden bei DLPC allerdings nur so viele Standard-Dienste zugelassen, um weitgehend die Annahme aller erwarteten Premium-Dienste realisieren zu können. Ähnliche Ergebnisse zeigten zahlreiche Simulationen von komplexeren Szenarien mit mehreren knappen Ressourcen und einer Vielzahl von Media-On-Demand-Diensten mit unterschiedlichen Übertragungsraten.

5 Fazit und Ausblick

IT-Service-Management bezeichnet Aufgaben und Prozesse, um IT-Dienstleistungen in vereinbarter Qualität zu möglichst geringen Kosten zu erbringen. Die automatisierte Behandlung von Überlastsituationen (overload control) stellt hierbei eine zentrale Aufgabe dar. Neue Möglichkeiten durch gemeinsame Nutzung von IT-Ressourcen erfordern neue Verfahren zur effektiven Behandlung von Überlastsituationen.

Hersteller von Servern bieten bereits heute Priorisierungsmöglichkeiten über statische Regeln an. Diese Ansätze sind jedoch nur bedingt in der Lage, auf spontane Nachfrageschwankungen, Änderungen des Service-Katalogs oder gar Veränderungen einzelner Dienstleistungsvereinbarungen dynamisch reagieren zu können. Zusätzlich ignorieren diese Ansätze die unterschiedlichen Anforderungen einzelner Anfragetypen an verschiedene Ressourcen, deren Interdependenzen und die damit verbundene Kombinatorik.

In diesem Artikel wurde ein neues Entscheidungsmodell namens DLPC zur effi-

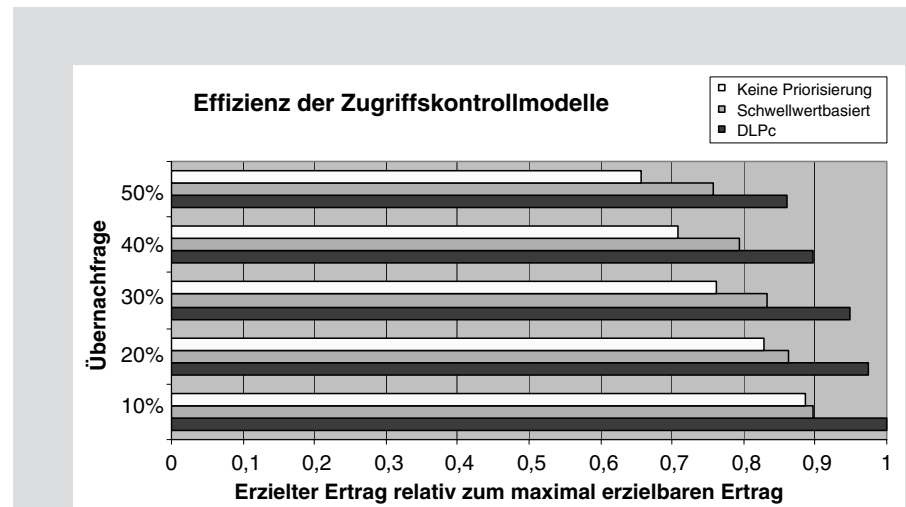


Bild 8 Ausschließlicher Zugriff auf Video-Dateien

zienten Behandlung von Überlastsituationen auf Anwendungsebene vorgestellt und anhand von Simulationen evaluiert. DLPC entscheidet dynamisch, welche Dienstklassen bei bestimmten Auslastungswerten zugelassen werden und welche im Hinblick auf Ertragsmaximierung gepuffert oder abgelehnt werden. Das Entscheidungsmodell ermöglicht somit, dass Dienstleistungsvereinbarungen über die Erreichbarkeit und Antwortzeit (availability) auch bei Überlast bestmöglich erfüllt werden. Die Simulationsresultate zeigen deutliche Verbesserungen im Vergleich zur Ertragssituation mit statischen Priorisierungsrichtlinien.

Der Vorteil liegt zum einen in der dynamischen Anpassung des DLPC an sich ändernde Nachfrage und zum anderen in der Berücksichtigung mehrerer Ressourcen und deren Abhängigkeiten bei den Zugriffskontrollentscheidungen. In unserer aktuellen Forschung gehen wir verstärkt auf die praktischen Anforderungen konkreter Anwendungen ein und analysieren die Güte verschiedener Prognoseverfahren bei stark schwankender Nachfrage. Ziel der Forschung ist es, für verschiedene Arten von Zugriffskontrollproblemen in Zukunft angepasste, effiziente Lösungen bereitstellen zu können.

Abstract

Revenue Maximizing Admission Control Policies for Shared IT Infrastructures

IT service providers are increasingly hosting different services of different customers on a shared IT infrastructure. While this fosters utilization of hardware infrastructure, system malfunctions, unexpected service behaviour or peak demands for one or more services may exploit resource pools (CPU, I/O, main memory, bandwidth etc.), entailing rejection of service requests. In this paper we describe models for dynamic admission control on shared infrastructures.

The admission control model decides whether to accept, buffer or reject a service request based on the revenue, Service Level Agreements (SLAs) and its resource demand in comparison to the actual workload to maximize overall revenue. Simulations of a media streaming infrastructure have been used for evaluation and comparison with traditional admission control policies.

Keywords: Admission Control, Demand Management, Overload Control, Revenue Management, Service Level Management

Notation

- i Dienststyp ($i = 1, \dots, I$) im Dienstleistungskatalog
- D_i Zufallsvariable der Nachfragemenge nach Dienst i während eines Planungsintervalls
- μ_i Erwartungswert der Nachfragemenge nach Dienst i während eines Planungsintervalls
- r_i Ertrag pro Erbringung eines Dienstes vom Typ i
- e Gemeinsam genutzte Ressource ($e = 1, \dots, E$)
- C_e Vorhandene Gesamtkapazität einer Ressource e
- U_e Auslastungsgrad der Ressource e
- a_{ei} Quantität der Einheiten einer Ressource e , die bei Ausführung eines Dienstes i belegt wird (Ressourcenbelegungskoeffizient)
- t_{ei} Dauer der Belegung einer Ressource e durch Dienst i
- Nummer einer Dienstanfrage ($k = 1, \dots, 8$)
- t_k Zeitpunkt einer Dienstanfrage k
- l_{ek} Planungszeitraum für Dienstanfrage k ab Anfragezeitpunkt t_k für Ressource e
- k' Dienstanfrage in Abarbeitung
- k'' Im Planungsintervall erwartete Dienstanfrage
- $l_{ek'}$ Prognostizierte Restlaufzeit der Verwendungen von e ab t_k durch einen sich gerade in Abarbeitung befindlichen Dienst
- q_{ei} Korrekturfaktor, der die relevanten Anteile der Belegungsdauern einer Ressource e durch einen Dienst i im Planungsintervall angibt
- x_i Anzahl der in der Planungsperiode zu akzeptierenden Anfragen nach Dienst i
- λ_e Erwartete Gesamtertragssteigerung im Planungsintervall durch die Bereitstellung einer zusätzlichen Einheit der Ressource e (Schattenpreis, Opportunitätskosten)

Danksagung

Die Ergebnisse des Artikels entstammen dem Gemeinschaftsprojekt „Dynamic Value Webs for IT Services“ zwischen der TU München und Siemens Business Services GmbH & Co. OHG. Wir bedanken uns für die freundliche, fachliche und finanzielle Unterstützung durch Siemens Business Services.

Literatur

- [Belo87] *Belobaba, P.*: Airline Yield Management – An Overview of Seat Inventory Control. In: *Transportation Science* 21 (1987) 2, S. 63–73.
- [BoBi03] *Boyd, E. A.; Bilegan, I. C.*: Revenue Management and E-Commerce. In: *Management Science* 49 (2003) 10, S. 1363–1386.
- [Bran06] *Brandl, R.*: Reinraum-Messungen zur Verrechnung von IT-Anwendungen. In: *Proceedings of Multikonferenz Wirtschaftsinformatik 2006*.
- [ChMo02] *Chen, H.; Mohapatra, P.*: Session-based overload control in QoS-aware Web servers. In: *IEEE INFOCOM* (2002).
- [ChCh96] *Chen, I.-R.; Chen, C.-M.*: Threshold-Based Admission Control Policies for Multimedia Servers. In: *The Computer Journal* 39 (1996) 9, S. 757–766.
- [ChMC01] *Chen, X.; Mohapatra, P.; Chen, H.*: An Admission Control Scheme for Predictable Server Response Time for Web Accesses. In: *Proceedings of the World Wide Web Conference 2001*.
- [ChPh02] *Cherkasova, L.; Phaal, P.*: Session-based admission control: A mechanism for peak load management of commercial Web sites. In: *IEEE Transactions on Computers* 51 (2002) 6, S. 669–685.
- [ChSt03] *Cherkasova, L.; Staley, L.*: Building a Performance Model of Streaming Media Applications in Utility Data Center Environments. In: *Proceedings of the International Symposium on Cluster Computing and the Grid 2003*.
- [ChTV04] *Cherkasova, L.; Tang, W.; Vahdat, A.*: MediaGuard: A Model-Based Framework for Building QoS-aware Streaming Media Services. HP Labs Report No. HPL-2004-25 (2004).
- [BoFP99] *de Boer, S.; Freling, R.; Piersma, N.*: Stochastic Programming for Multiple-Leg Network Revenue Management. In: *Tech. Rep. EI-9935/A* (1999).
- [FKSR02] *de Farias, D. P.; King, A. J.; Squillante, M. S.; Van Roy, B.*: Dynamic Control of Web Server Farms. In: *Proceedings of the 2nd Annual INFORMS Revenue Management Section Conference*. Columbia University, New York (2002).
- [DeGo01] *Debasis, M. L. T. I.; Gopalas, R. K. L. T. I.*: (2001). Band Allocating Method. European Patent Office, Patent Nr. US6721270 (2001).
- [DeGo04] *Debasis, M. L. T. I.; Gopalas, R. K. L. T. I.*: Multicommodity flow method for designing traffic distribution on a multiple-service packetized network. European Patent Office, Patent Nr. EP1076472 (2002).
- [EgHe99] *Eggert; Heidemann, J.*: Application-Level differentiated service from an internet server. In: *World Wide Web Journal* 3 (1999) 2, S. 133–142.
- [ENTZ04] *Elnikety, S.; Nabum, E.; Tracey, J.; Zwaenepoel, W.*: A Method for Transparent Admission Control and Request Scheduling in E-Commerce Web Sites. In: *Proceedings of the World Wide Web Conference*. New York 2004.
- [FoOB06] *Fornfeld, M.; Oefinger, P.; Braulke, T.*: Gesamtwirtschaftliche Auswirkungen der Breitbandnutzung. http://www.bitkom.org/files/documents/ BITKOM_Studie_Breitbandnutzung.pdf, Abruf am 2006-04-12.
- [KiMP05] *Kim, R. Y.; Manas, T.; Pramod, K. S. U.*: (2005). Policy-Based Admission Control And Bandwidth Reservation For Future Sessions. European Patent Office, Patent Nr. WO2005072321 (2005).
- [KwYe00] *Kwon, J.; Yeom, H.*: An Admission Control Scheme for Continuous Media Servers Using Caching. In: *Proceedings of the Int'l Performance, Computing and Communication Conference (IPCCC)*. Phoenix 2000, S. 456–462.
- [LeAF02] *Lewis, M.; Ayhan, H.; Foley, R.*: Bias optimal admission control policies for a multi-class non-stationary queuing system. In: *Journal of Applied Probability* 39 (2002) 1, S. 20–37.
- [LiSW01] *Liu, Z.; Squillante, M. S.; Wolf, J. L.*: On Maximizing Service-Level-Agreement Profits. In: *Proceedings of the 3rd ACM Conference on Electronic Commerce*. Orlando 2001.
- [McRy99] *McGill, J. I.; van Ryzin, G. J.*: Revenue management: Research overview and prospects. In: *Transportation Science* 33 (1999), S. 233–256.
- [NaBa01] *Nair, S. K.; Bapna, R.*: An Application of Yield Management for Internet Service Providers. In: *Naval Research Logistics* 48 (2001) 5, S. 348–362.
- [RuCT05] *Rumsewicz, M.; Castro, M.; Tai Le, M.*: Eddie Admission Control Scheme: Algorithm Description, Prototype Design Details and Capacity Benchmarking. <http://eddie.sourceforge.net/techrep.html>, Abruf am 2005-08-08.
- [Simp89] *Simpson, R. W.*: Using Network Flow Techniques for Origin-Destination Seat Inventory control. Report at the Department of Aeronautics and Astronautics: MIT 1998.
- [SMAM02] *Singhmar, N.; Mathur, V.; Apte, V.; Manjunath, D.*: A Combined LIFO-Priority Scheme for Overload Control of E-commerce Web-Servers. In: *Proceedings of the IEEE RTSS International Infrastructure Survivability Workshop* (2002).
- [Stid99] *Stidham, S.*: Optimal control of admission to a queuing system. In: *IEEE Trans Automat Control* 30 (1999) 8, S. 705–713.
- [TaRy99] *Talluri, K. T.; van Ryzin, G. J.*: An Analysis of Bid-Price Controls for Network Revenue Management. In: *Management Science* 44 (1999) 1577–1593.
- [VeGh03] *Verma, A.; Ghosal, S.*: On Admission Control for Profit Maximization of Networked Service Providers. In: *Proceedings of WWW2003*. Budapest 2003.
- [ViGG94a] *Vin, H.; Goyal, A.; Goyal, P.*: A statistical admission control algorithm for multimedia servers. *International Multimedia Conference*. San Francisco, California, USA 1994, S. 33–40.
- [ViGG94b] *Vin, H. M.; Goyal, A.; Goyal, P.*: An observation-based admission control algorithm for multimedia servers. *Multimedia Computing and Systems*, 1994. Boston, MA, USA 1994, S. 234–243.
- [Will92] *Williamson.*: Airline Network Seat Inventory Control – Methodologies and Revenue Impacts. Doctorial Thesis at the Department of Aeronautics and Astronautics: MIT 1992.