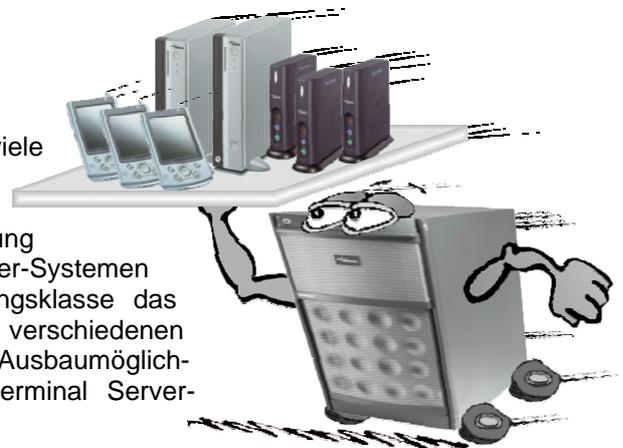


Terminal Server Sizing Guide

Abstract

In diesem technischen Papier wird die Frage diskutiert, wie viele Benutzer eine als Terminal Server eingesetzte PRIMERGY mit adäquater Performance bedienen kann. Das Papier richtet sich in erster Linie an Personen, die sich mit der Planung und Konfektionierung von Microsoft Terminal Server-Systemen beschäftigen. Es soll helfen, für eine geforderte Leistungsklasse das passende PRIMERGY Modell zu finden. Dabei wird auf die verschiedenen Leistungsklassen der PRIMERGY Modellpalette und deren Ausbaumöglichkeiten mit Prozessoren, Arbeitsspeicher und anderen Terminal Server-relevanten Komponenten eingegangen.



Inhalt

PRIMERGY	2	Netzwerk	49
Windows Terminal Server	3	Benutzerverhalten	51
Microsoft Terminal Services 2003.....	4	Eingabegeschwindigkeit	51
Citrix Presentation Server	5	Betriebssystem	52
Skalierung.....	8	Windows Server 2003 R2	52
Scale-Up	8	64-bit.....	53
Scale-Out.....	9	Nutzbarer Speicher.....	54
Dimensionierung.....	11	Physical Address Extension (PAE).....	56
Benutzer	11	Anzahl Prozesse.....	57
Benutzersimulation	12	Terminal Server Version	58
Vergleichbarkeit.....	14	Microsoft Terminal Server vs. Citrix Presentation	
»Tool for User Simulation«	15	Server	58
Messumgebung	17	Citrix Presentation Server Version.....	58
Lastprofil	18	Anwendungen	59
Messmethode	19	Microsoft Office Version.....	59
Ressourcenbedarf	23	Einstellungen für Microsoft Office in einer	
Rechenleistung	24	Terminal Server Umgebung.....	59
Prozessortyp	24	Infrastruktur	60
Taktfrequenz	25	Vergleich der Messwerkzeuge	62
Front-Side-Bus.....	32	Microsoft Testwerkzeuge und -Skripte.....	62
Caches.....	33	Ergebnisse von Fujitsu Siemens Computers und	
Hyper-Threading	35	Microsoft	64
Anzahl Prozessoren.....	36	Resümee	66
Verhalten bei hoher CPU-Last	37	Literatur	68
Arbeitsspeicher	39	Kontakt	68
Disk-Subsystem.....	46		

PRIMERGY

All den Lesern, denen der Name PRIMERGY noch kein Begriff ist, sei hier zunächst ein kleiner Überblick gegeben: PRIMERGY Server ist seit 1995 der Markenname für die sehr erfolgreiche Industrie-Standard-Server-Familie aus dem Hause Fujitsu Siemens Computers. Es handelt sich dabei um eine bei Fujitsu Siemens Computers entwickelte und produzierte Produktlinie mit Systemen für kleine Arbeitsgruppen bis hin zu Lösungen für Großunternehmen.

Scalability, Flexibility & Expandability

Vom kleinen Monoprocessorsystem bis hin zu Systemen mit 16 Prozessoren kommen in der PRIMERGY-Familie die neuesten Technologien zum Einsatz. Als Herzstück werden Intel Prozessoren der obersten Leistungsklasse verwendet. Mehrere 64-bit PCI-X-I/O- und Memory-Busse, schnelles RAM und performante Komponenten, wie SCSI-Technologie und Fibre-Channel, sorgen für hohen Datendurchsatz. Dies bedeutet Leistung satt, gleich ob für Scaling-Out oder Scaling-Up. Bei der

Methode des Scaling-Out, die nach dem Ameisenstaat-Modell mehr Leistung durch eine Vielzahl von Einzelindividuen erzielt, können idealerweise Blade-Server und kompakte Compu-Node Systeme platziert werden. Für die Methode des Scale-Ups, d.h. Hochrüsten eines vorhandenen Systems, sorgen die umfangreichen Ausbaumöglichkeiten der PRIMERGY Systeme, auf bis zu 16 Prozessoren und 128 Gigabyte Arbeitsspeicher. PCI- und PCI-X-Slots sorgen für die notwendige Erweiterbarkeit von I/O-Komponenten. Eine Langzeitplanung in enger Zusammenarbeit mit namhaften Zulieferern von Komponenten, wie z.B. Intel, LSI, ServerWorks, sichert kontinuierliche und bestmögliche Kompatibilität von einer zur nächsten Server-Generation. Die PRIMERGY-Planung reicht zwei Jahre in die Zukunft und garantiert eine möglichst frühe Einbeziehung neuester Technologien.

Reliability & Availability

Neben der Leistung steht die Qualität im Vordergrund. Dazu zählen nicht nur eine exzellente Verarbeitungsqualität und der Einsatz qualitativ hochwertiger Einzelkomponenten, sondern auch Vorkehrungen zur Ausfallsicherheit, frühzeitiger Fehlerdiagnose und Datenschutz. Wichtige Systemkomponenten sind redundant ausgelegt, und werden vom System auf Funktionalität überwacht. Viele Teile können problemlos im laufenden Betrieb ausgetauscht werden, so dass Ausfallzeiten minimiert werden und die Verfügbarkeit gewährleistet wird.

Security

Ihre Daten sind der PRIMERGY heilig. Schutz vor Datenverlusten bieten leistungsfähige Disk-Subsysteme der PRIMERGY und FibreCAT Produktlinie. Eine noch höhere, größtmögliche Verfügbarkeit bieten PRIMERGY Cluster-Konfigurationen, bei denen nicht nur die Server, sondern auch die Disk-Subsysteme sowie die gesamte Verkabelung redundant ausgelegt werden können.

Manageability

Umfassende Management-Software für alle Phasen des Server-Lebenszyklus sorgt für einen reibungslosen Betrieb und erleichtert Wartung und Fehlerdiagnose der PRIMERGY.



ServerStart, ein benutzerfreundliches, menübasiertes Software-Paket für die optimale Installation und Konfigurierung des Systems mit automatischer Hardware-Erkennung und Installation aller notwendigen Treiber.



ServerView zur Serverüberwachung mit Alarm-, Schwellen-, Berichts- und Basis-Management, Prefailure Detection and Analyzing, Alarm-Service und Versionsmanagement.



RemoteView zur von der Hardware und dem Betriebssystem unabhängigen Fern-Wartung und -Diagnose via LAN oder Modem.



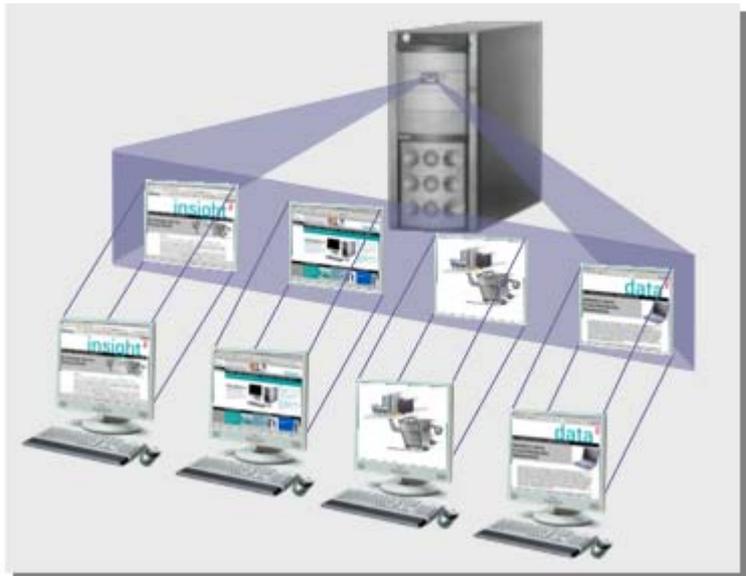
Weitere detaillierte Informationen zu den PRIMERGY Systemen finden Sie im Internet unter <http://www.primergy.de/>.

Windows Terminal Server

Terminal Server steht als Oberbegriff für Server-based Computing auf Basis von Microsoft® Windows® Server Betriebssystemen.

Server-based Computing ist eine Systemarchitektur, bei der Microsoft Windows Client-Anwendungen zu 100 Prozent auf dem Server installiert und ausgeführt werden. Von dort erfolgt nicht nur deren Einsatz, sondern auch deren Wartung, Verwaltung und Support finden direkt auf dem Server statt. Lediglich die Benutzeroberfläche, d.h. die Bildschirm-, Maus- und Tastatur-Informationen werden zwischen Client und Server übertragen. Der Benutzer kann so von fast beliebigen Clients aus, auch nicht Windows basierten, über einen solchen Terminal Server sofort auf Windows-Anwendungen zugreifen,

ohne dass die jeweiligen Applikationen erst zum Client übertragen, dort gestartet, oder gar auf lokalen Massenspeichern vorgehalten werden müssten. Wird ein Client ausschließlich in diesem Server-based Szenario eingesetzt, so hat er hinsichtlich Speicher- und Plattenausstattung wesentlich geringere Anforderungen als ein traditioneller Client, man spricht daher auch von so genannten Thin-Clients.



Senkung der TCO durch Rezentralisierung

Rasant ansteigende Betriebskosten (Total Cost of Ownership) zählen heute zu den größten Problemen in den IT-Umgebungen der Unternehmen. Früher achtete man bei der Einrichtung eines unternehmensweiten IT-Systems vorrangig auf die Anschaffungskosten und weniger auf die Folgekosten. Nach Angaben von Analysten haben jedoch die Anschaffungskosten, die zweifellos eine beträchtliche einmalige Investition darstellen, nur einen Anteil von 15 Prozent an den Gesamtkosten einer unternehmensweiten IT-Lösung. Daher richtet sich heute das Augenmerk mehr auf die laufenden Kosten.

Das Konzept des Server-based Computing hilft durch Rezentralisierung von Anwendungen und Daten, diese Kosten zu reduzieren. Man hat erkannt, dass es effektiver ist, in einer Client-Server-Architektur die Bereitstellung der Applikationen sowie Hardware- und Softwarepflege von einer zentralen Stelle aus im gesamten Unternehmen durchzuführen statt an jedem einzelnen Arbeitsplatz. Server-based Computing kann sowohl für die Endanwender als auch für die Systemadministratoren die Produktivität und Effizienz erheblich verbessern. Nach Meinung von Analysten kann das Server-based Computing die IT-Betriebskosten um 30 bis 50% senken.

Einsatzgebiet

Ein Terminal Server kann prinzipiell für alle Arten von Applikationen eingesetzt werden. Wo bislang kleine Rechner oder Terminals für einfache Dateneingabe bzw. -abfrage Verwendung fanden, können mit dem Terminal Server moderne Anwendungen in ein bestehendes Umfeld integriert werden. Aber auch in Umgebungen, in denen ein einzelner Benutzer bereits eine höhere Rechen- oder Grafikleistung braucht, bietet der Terminal Server den Vorteil der zentralen Bereitstellung der Anwendungen.

Historie

Das von Mainframes seit langem bekannte Konzept des Server-based Computing hielt 1994 Einzug in die Windows-Welt. Als erstes entwickelte das US-amerikanische Softwarehaus Citrix eine Multiuser-Erweiterung für Microsoft Windows NT 3.51, die unter dem Namen »WinFrame« als Gesamtprodukt aus Windows NT und Multiuser-Erweiterung vertrieben wurde. 1997 hat Microsoft die so genannte »Citrix MultiWin Technologie« und damit einen Teil der Citrix-Betriebssystemerweiterung für NT von Citrix lizenziert und in das Produkt »Windows NT 4.0 Terminal Server« einfließen lassen. Seit 2000 hat diese Technologie unter dem Namen »Terminal Services« einen festen Platz in allen Server-Produkten der Microsoft Windows 2000 Server™ und Windows Server 2003 Produktlinie. Selbst in dem Client-Betriebssystem Windows XP Professional steht in begrenztem Umfang der Terminal Service unter dem Namen »Remote Desktop« zur Verfügung. Von einem beliebigen Windows-Client kann somit auf das entfernte System zugegriffen werden, wobei die Anwendungen komplett auf dem Remote-System ablaufen. Das unterliegende Protokoll wird als »Remote Desktop Protocol« (RDP) bezeichnet.

Bereits seit Windows 2000 Server bietet Terminal Server umfangreiche Funktionalitäten, hier nur einige der wichtigsten:

Unterstützte Clients

- 32-bit Clients für Windows-basierte PCs
- 16-bit Client für Windows for Workgroups
- Windows CE-basierter Thin-Client
- Windows XP Embedded-basierter Thin-Client
- Microsoft ActiveX® Control

Client-Eigenschaften

- Datenaustausch von Text und Grafiken über die Zwischenablage zwischen Client und Server
- Bereitstellung der lokalen parallelen und seriellen Schnittstellen des Client innerhalb der Server-basierten Anwendung
- Druck auf lokalen Druckern am Thin-Client
- Zugriff auf lokale Laufwerke des Clients innerhalb der Server-basierten Anwendung
- Bitmap-Caching zur Performance-Steigerung

Kommunikation

- Client-Anbindung über lokales Netzwerk (LAN), »wide area network« (WAN), Dial-up, »Integrated Services Digital Network« (ISDN), »x digital subscriber line« (xDSL), und »virtual private network« (VPN)
- Verschlüsselung der Client-Kommunikation

Microsoft Terminal Services 2003

Windows Server 2003 Terminal Services bieten weitere Neuerungen gegenüber der Vorgängerversion in Windows 2000 Server. Hier einige der wichtigsten:

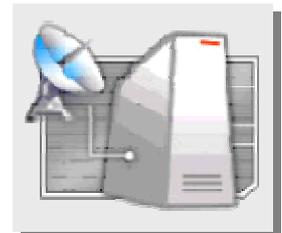
Server-Eigenschaften

- Optimiertes Ressourcen-Management, so dass Windows Server 2003 auf gleicher Hardware nun mehr Benutzer unterstützen kann als unter Windows 2000 Server.
- Windows Server 2003 Enterprise Edition und Datacenter Edition bieten ein Session Directory und somit die Möglichkeit, eine Terminal Server-Farm mit Load Balancing aufzubauen.
- Verbesserte Manageability durch Group Policies, Windows Management Instrumentation (WMI).
- Nutzung von Windows Server 2003 Erweiterungen wie Software Restriction Policies, Erweiterung von Roaming Profiles und neuen Windows-Anwendungscompatibilitätsmodi.

Client-Eigenschaften

- Unterstützung weiterer lokaler Geräte, wie Smart Cards, und Audio-Ausgabe innerhalb der Server-basierten Anwendung
- Farbtiefe bis zu True Color (24-bit)
- Bildschirmauflösung bis zu 1600 x 1200
- Individuelle Zeitzonen je Benutzer

Die Terminal Services sind Bestandteil des jeweiligen Betriebssystems und liegen daher beim Windows 2000 Server und beim Windows Server 2003 als 32-bit Version und bei der Windows Server 2003 x64 Edition als 64-bit Version vor.



Citrix Presentation Server

Als Erweiterung der Basis Terminal Services in Windows bietet Citrix mit der Produktfamilie »Citrix Presentation Server for Windows« (im Weiteren abgekürzt als »Citrix Presentation Server«) sinnvolle Ergänzungen. Die aktuelle Version ist »Citrix Presentation Server 4.0 for Windows«. Die vorangegangenen Versionen sind unter dem Produktnamen »Citrix MetaFrame« bekannt geworden.



Um den Ansprüchen verschiedener Unternehmen gerecht zu werden, gibt es drei verschiedene Produktvarianten von Citrix Presentation Server:

- Citrix Presentation Server, Standard Edition
Applikationsbereitstellung für kleinere Unternehmen
- Citrix Presentation Server, Advanced Edition
Applikationsbereitstellung mit Load Balancing für mittlere Unternehmen
- Citrix Presentation Server, Enterprise Edition
Applikationsbereitstellung mit Load Balancing für größere Unternehmen mit mehreren Standorten

Die wichtigsten Erweiterungen von Citrix Presentation Server gegenüber Windows Terminal Services sind:

Server

- Veröffentlichte Anwendungen (Published Applications), randlose Fenster (Seamless Windows)
Direkter Start einer Server-seitigen Anwendung, ohne einen Windows Desktop zu starten.
Direkter Zugriff auf einzelne Anwendungen.
- Load Balancing
Ein integriertes Load Balancing sorgt für die automatische, lastspezifische Verteilung der Benutzer auf die einzelnen Terminal Server einer Terminal Server-Farm. Das Load Balancing ist in Citrix Presentation Server Advanced Edition und Enterprise Edition enthalten, es wird dafür keine zusätzliche Software benötigt. Es kann eine sehr feine und individuelle Einstellung der Load Balancing-Kriterien vorgenommen werden, insbesondere können auch Terminal Server-spezifische Parameter gewählt werden.

Unterstützte Clients

- Unterstützung heterogener Clients
Auch nicht Windows-basierte Clients können durch das Betriebssystem-unabhängige »Independent Computing Architecture« (ICA)-Protokoll auf vom Server bereitgestellte Anwendungen zugreifen. Zusätzlich werden bei Citrix 16-bit Clients für ältere Windows-Versionen und für Microsoft MS-DOS® sowie Clients für UNIX, Macintosh, Java und ein Browser Client angeboten.

Die Versionen »Citrix MetaFrame XP«, »Citrix MetaFrame Presentation Server 3.0« und »Citrix Presentation Server 4.0« laufen nur unter einem 32-bit Windows Betriebssystem. Für Windows x64 ist der »Citrix Presentation Server for Windows Server 2003 x64« einzusetzen, der funktional der 32-bit Version »Citrix Presentation Server 4.0« entspricht.

Im Folgenden werden die Performance-relevanten Funktionen von »Citrix Presentation Server 4.0« im Vergleich zur Vorgängerversion 3.0 kurz charakterisiert.

Citrix Presentation Server 4.0

Citrix Presentation Server 4.0 enthält einige Neuerungen, die Einfluss auf die Performance haben können. Die ersten beiden Features dienen im Wesentlichen der Applikationskompatibilität.

Application Isolation

»Application Isolation« ermöglicht eine isolierte Installations- und Ablaufumgebung für Applikationen, mit dem Ziel gegenseitige Störungen durch z.B. Registry Settings, Konfigurationsdateien etc. auszuschließen. Das ist vorteilhaft, um verschiedene Versionen einer Applikation auf einem Terminal Server vorzuhalten, oder auch »Alt-Last« Anwendungen, die in einer Multi-User-Umgebung sonst nicht ablauffähig wären. Realisiert wird dieses Feature durch eine Virtualisierungsschicht für Registry Settings, Dateisystem und Named Objects z.B. Semaphoren, Sections etc., auf der die Anwendung installiert wird. Der Citrix Präsentation Server übernimmt dann das Mapping von den virtuellen Ressourcen auf die physikalischen Ressourcen des Betriebssystems. Dabei ist es möglich, die Anwendung bereits in eine isolierte Umgebung zu installieren oder auch nur eine veröffentlichte Anwendung in einer isolierten Umgebung ablaufen zu lassen. Letzteres ist sinnvoll, wenn die Applikation nicht in einer Mehrbenutzerumgebung ablaufen kann.

Virtual Address Support

»Virtual Address Support« unterstützt Anwendungen, die eine eindeutige IP-Adresse pro Session benötigen.

Virtual Memory Optimization

Das Ziel der »Virtual Memory Optimization« ist es, Speicherplatz zu sparen, indem Adress Konflikte beim Laden von DLLs nicht durch »relocation« sondern durch »rebasing« gelöst werden. »Relocation« bedeutet, die DLL wird nicht an die im Image stehende Basisadresse geladen, sondern in einen anderen Teil des virtuellen Speichers und alle in der DLL benutzten Adressen müssen relativ zur Basisadresse umgerechnet werden. Eine Benutzung der DLL durch mehrere Applikationen hat damit auch ein mehrfaches Laden in den virtuellen Speicher zur Folge. »Rebasing« bedeutet, es wird eine Schattenkopie der DLL-Datei angelegt, die eine (konfliktfreie) optimale virtuelle Basisadresse enthält. Dadurch braucht ein »rebased« Objekt nur einmal geladen werden, auch wenn es von mehreren Applikationen genutzt wird.

Das »rebasing« von DLLs führt damit also zu Einsparungen beim virtuellen Speicherplatzverbrauch. In welchem Maße die Speichieranforderungen verringert werden, ist stark applikationsabhängig. Außerdem funktioniert dieses Feature auch nicht bei allen Anwendungen, z.B. können Applikationen, deren DLLs geschützt sind durch »Windows Rights Managements« oder die »digitally signed components« haben, nicht rebased werden. Solche Applikationen können durch eine Ausschlussliste vom Prozess des Optimierens ausgeschlossen werden.

Die »Virtual Memory Optimization« wird durch einen Monitorprozess realisiert, der feststellt, wo »relocation« von DLLs erfolgt und dies in deiner Datei protokolliert.

Zu vom Administrator festzulegenden Zeiten wird dann ein Prozess tätig, der diese Datei liest und das »Rebasing« der entsprechenden DLLs durchführt.

Zusätzlich zum »Rebasing« erfolgt auch das »Binding« von DLLs, d.h. in der DLL »import Section« wird die virtuelle Ladeadresse der importierten Funktionen gleich eingetragen. Auf diese Weise wird CPU-Zeit beim Initialisieren der Applikation gespart.

CPU Utilization Management

Das »CPU Utilization Management« dient dazu, die CPU-Leistung gleichmäßig auf die vorhandenen Benutzer aufzuteilen, um dadurch Spitzen in der CPU-Auslastung eines Servers auszugleichen, und daher mehr User pro Server zu ermöglichen.

Der tatsächliche Performancegewinn ist sehr abhängig von der Anwendungslast.

Wenn das »CPU Utilization Management« aktiviert ist, bekommt jeder Benutzer den gleichen Anteil CPU-Zeit zugeteilt. Auf diese Weise soll verhindert werden, dass ein Benutzer, der besonders intensiv die CPU nutzt, andere Benutzer benachteiligt. Schöpft ein Benutzer seinen ihm zur Verfügung stehenden Anteil an CPU-Leistung nicht voll aus, so können die anderen Benutzer diese Leistung nutzen. Es ist auch möglich, einem Benutzer einen relativ zu den anderen Benutzern höheren Anteil an CPU-Leistung zuzusprechen. Ebenso ist es möglich, einem Benutzer einen festen Anteil der Rechenleistung, unabhängig von der Gesamtlast, zuzuteilen.

Leider gibt es für die Verwaltung der benutzerspezifischen CPU-Zuteilung kein grafisches Benutzerinterface und die Einstellungen erfolgen über entsprechende »Registry Settings«.

Realisiert wird die entsprechende CPU-Zuteilung über Prozesse, die den CPU-Verbrauch der Benutzer monitoren und dann über Betriebssystemaufrufe den Scheduler beeinflussen.

Support für Windows Server 2003 x64 Edition

Citrix unterstützt Windows Server 2003 x64. Mit der 64-bit Version von Windows können verschiedene Limitierungen von 32-bit überwunden werden. So sind z.B. unter dem 64-bit Betriebssystem die Kernel-Ressourcen ausreichend dimensioniert; ein Beispiel ist der Non-Paged Pool mit 128 GB, der unter dem 32-bit Windows durch die geringe Größe (256 MB) schon zum Engpass werden konnte (vgl. auch Kapitel [Nutzbarer Speicher](#)). Auch unter dem 32-bit Betriebssystem bietet Citrix Presentation Server Version 4.0 Vorteile gegenüber der Vorgängerversion, da bestimmte Seiten, die früher dem »Non-Paged Pool« zugeordnet waren, jetzt dem »Paged Pool« zugeordnet sind.

Verbessertes Printing

Ein neuer »Universal Print Driver« verspricht laut Dokumentation eine bis zu vierfach verbesserte Druckgeschwindigkeit bei verminderter Speicher- und Bandbreitennutzung.

Skalierung

Bei der Skalierung – das ist der Prozess, das System an die benötigte Leistung anzupassen – werden zwei Methoden unterschieden:

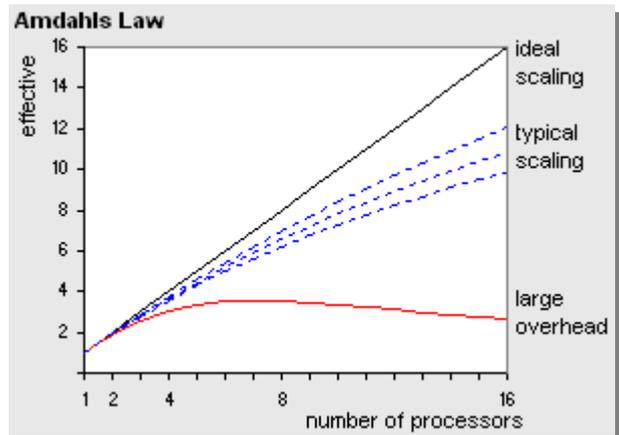


Beide Szenarien, sowohl Scale-Up als auch Scale-Out, werden von Terminal Server unterstützt.

Scale-Up

Beim Scale-Up wird die Leistung eines Terminal Servers durch den Einsatz leistungsfähigerer Hardware, also insbesondere Rechenleistung und Arbeitsspeicher, erhöht. Diesem Skalierungsprozess sind Grenzen durch die maximale Größe eines Server-Systems gesetzt.

Theoretisch benötigt man »nur« eine beliebig leistungsfähige Hardware und würde im Scale-Up-Szenario einen beliebig leistungsfähigen Terminal Server erhalten. Dies ist jedoch leider nur Theorie. So ist die Skalierung mit wachsender Anzahl Prozessoren nur im Idealfall einer optimal parallelisierbaren Anwendung linear. Je mehr Zugriffe jedoch auf gemeinsame Ressourcen, wie Arbeitsspeicher, Festplatten oder Netzwerk erfolgen, und somit eine Koordination zwischen den Prozessoren bedingen, umso mehr flacht die Skalierungskurve ab. Im Extremfall kann es bei einer sehr großen Anzahl Prozessoren und sehr hohem Koordinationsanteil der Prozessoren untereinander sogar zu einem »Umkippen« der Skalierung kommen. Man bezeichnet diesen Sachverhalt auch als »Amdahls Gesetz«, nach Gene Amdahl, der dieses 1967 untersuchte und in ein mathematisches Modell fasste.



Designer von großen Multiprozessorsystemen wirken dem entgegen, indem sie den Prozessoren große Caches beiseite stellen oder Gruppen von Prozessoren bilden und diesen eigenen Arbeitsspeicher und I/O-Komponenten zuordnen.

In der Praxis setzt heute oft nicht die Hardware die Grenzen, sondern die Software-Architektur. Die heute zumeist eingesetzte Software im 32-bit Design kann die zur Verfügung stehende Hardware häufig nicht mehr voll nutzen. Im speziellen ergeben sich Limitierungen bei der Adressierung des Arbeitsspeichers, durch die 32-bit Anwendungen auf 4 GB virtuellen Adressraum begrenzt sind. Ist der Server physikalisch mit mehr als 4 GB Arbeitsspeicher ausgestattet, so kann dieser Speicher zumeist nicht effektiv genutzt werden. Durch die Abhängigkeit zwischen dem Bedarf an Arbeitsspeicher und Rechenleistung können viele Anwendungen auch die Rechenleistung, die moderne Systeme mit 8 oder 16 CPU-Sockeln bereitstellen, nicht ausschöpfen.

Auch für Terminal Server ergibt sich eine Grenze, ab der ein Scale-Up nicht mehr die gewünschte Leistungssteigerung zeigt. Diese ist bei dem heutigen 32-bit Windows Server 2003 bei einem 4-way System mit 4 GB Arbeitsspeicher zu sehen. Daher waren Terminal Server-Umgebungen bisher klassische Scale-Out-Szenarien. Mit 64-bit-Betriebssystemen und 64-bit-Anwendungen werden diese Grenzen überwunden, so dass viele Kunden heute vor der Frage stehen, ob die neue 64-bit-Welt eine Lösung für die bisherigen Engpässe darstellt. Siehe hierzu die Kapitel »[Rechenleistung](#)«, »[Arbeitsspeicher](#)« und »[Betriebssystem](#)«).

Scale-Up ist eine adäquate Skalierungsmethode, wenn eine überschaubare Anzahl an Benutzern zu bedienen ist (vgl. Kapitel »[Resümee](#)«). Ist eine größere Anzahl an Benutzern von mehreren hundert oder tausenden mit Terminal Server zu bedienen, so kann das Scale-Up-Szenario nicht mehr verwendet werden und es bedarf anderer Mechanismen um die Leistung eines Terminal Servers zu steigern.

Scale-Out

Das Scale-Out-Szenario verfolgt einen anderen Weg als das Scale-Up. Anstatt einen Server immer größer zu dimensionieren, werden beim Scale-Out viele Server zu einer Gruppe zusammengefasst. Man spricht auch von Server-Farmen.

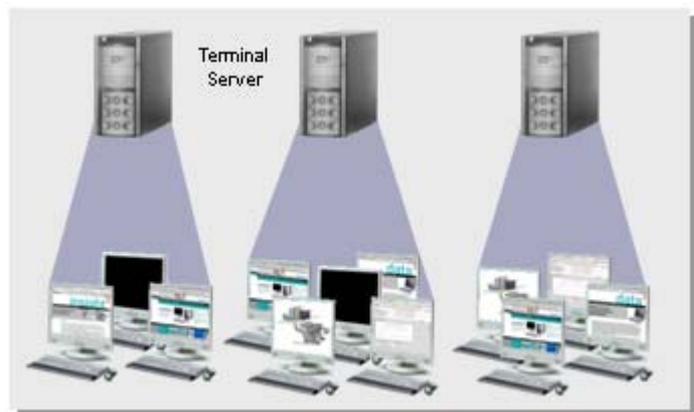
Mit diesem Konzept kann leicht die Grenze überwunden werden, die ein einzelner Terminal Server aufgrund seiner Software-Architektur bedingt. Die Skalierung ist aber auch bei einer Server-Farm nicht ideal linear, denn analog zu Amdahls Gesetz bei Multiprozessorsystemen gibt es auch in einer Server-Farm Overhead durch interne Kommunikation. Allerdings fällt dieser meist geringer aus als bei großen Multiprozessorsystemen.

Beim Scale-Out kann man drei Varianten unterscheiden:

Just a Bunch of Servers

»Just a Bunch of Servers« ist eine lose Ansammlung von Servern, in unserem Fall Terminal Servern. Diesen Terminal Servern sind dediziert Benutzergruppen oder Applikationen zugeordnet, es findet jedoch unter den Terminal Servern kein Informationsaustausch und kein Lastausgleich statt.

Der Vorteil dieser Architektur ist die sehr einfache Erweiterbarkeit. Nachteilig ist, dass kein automatischer Lastausgleich zwischen den einzelnen Servern stattfindet, so dass je nach Zuordnung der Benutzer zu den Servern Rechenleistung ungenutzt bleibt. Der administrative Aufwand ist recht hoch, da jedes System separat verwaltet werden muss.



Dennoch wird diese Variante des Scale-Outs in der Praxis in kleineren Konfigurationen eingesetzt.

Server-Farm

Eine Terminal Server-Farm ist ein Zusammenschluß von Terminal Servern, die eine gemeinsame Verwaltungseinheit, Data Store genannt, besitzen. Diese werden gemeinsam administriert. Die Zuordnung von Benutzern zu Servern und Applikationen erfolgt meist statisch, ein Load Balancing wird nicht notwendigerweise verwendet. Vorteil gegenüber der »Just a Bunch of Servers« Variante ist die vereinfachte Administration. Redundanz und ein automatischer Lastausgleich sind jedoch nicht gegeben.

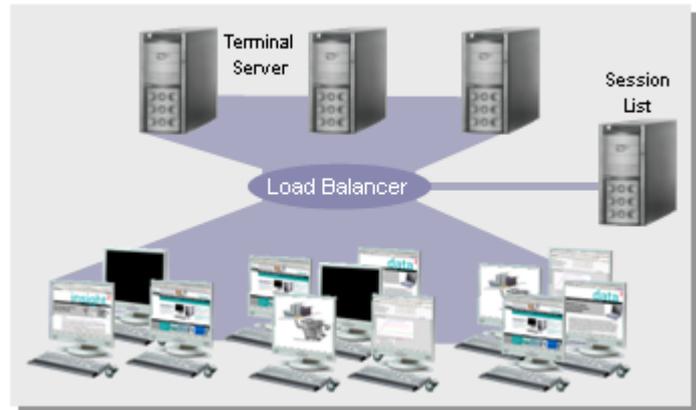
Auch diese Variante des Scale-Outs wird in der Praxis sehr häufig eingesetzt.



Load-balanced Server-Farm

Bei einer »load-balanced Server-Farm« werden die einzelnen Terminal Server zu einer logischen Einheit zusammengefasst. Wird von einem Client eine Session initiiert, so wird diese von einem Load Balancer nach bestimmten Mechanismen an den Server mit der momentan geringsten Auslastung delegiert.

Wesentliche Basis für ein Load Balancing von Terminal Servern ist das Führen einer Session-Liste. Terminal Server erlaubt es eine Verbindung zwischen Client und Server zu trennen (disconnect), wobei die Session auf dem Terminal Server jedoch weiterläuft. Baut der Client erneut eine Verbindung zu der Terminal Server-Farm auf, so muss anhand dieser Session-Liste sichergestellt werden, dass er wieder zu »seiner« bestehenden Session verbunden wird und nicht aufgrund des Load Balancing zu einem anderen Terminal Server der Farm, der momentan die geringste Auslastung zeigt. Eine Verschiebung von Sessions zwischen den einzelnen Mitgliedern der Farm wird von Terminal Server nicht unterstützt.



Neben dem Verteilen der Benutzer-Verbindungen in Abhängigkeit der Auslastung bietet die Methode der load-balanced Server-Farm auch eine gewisse Redundanz. Fällt ein Terminal Server aus, so können die Benutzer von den anderen Mitgliedern in der Server-Farm bedient werden. Bei dediziert zugeordneten Terminal Servern im »Just a Bunch of Server« Szenario stehen bei einem Ausfall eines Terminal Servers den zugeordneten Clients keine Terminal Server-Dienste mehr zur Verfügung. Allerdings bietet eine load-balanced Server-Farm keine Ausfallsicherheit für die einzelnen Client-Sessions. Fällt ein Terminal Server im laufenden Betrieb aus, so gehen alle auf diesem Terminal Server aktiven Sessions verloren.

Große Terminal Server Farmen, die sich auch über mehrere Standorte erstrecken können, findet man in der Praxis häufig in Enterprise-Umgebungen.

Scale-Out mit Terminal Server

Die einzelnen Versionen von Terminal Server unterscheiden sich hinsichtlich Ihrer Scale-Out Fähigkeiten.

Windows 2000 Server Terminal Services unterstützt keine Session-Liste. Eine Server-Farm mit Load Balancing kann somit ohne eine Zusatzsoftware wie Citrix Presentation Server nicht realisiert werden.

Windows Server 2003 Terminal Services unterstützt in den Enterprise und Datacenter Editionen eine Session-Liste. Das Load Balancing kann wahlweise mit Windows »Network Load Balancing« (NLB) oder mit dedizierten 3rd-Party Load Balancern, wie z.B. F5 Network BIG-IP, realisiert werden.

Citrix Presentation Server Advanced Edition und Enterprise Edition unterstützen ein Session-Directory und bieten ein eigenes sehr flexibles Load Balancing, das speziell auf die Bedürfnisse von Terminal Server zugeschnitten ist.

Dimensionierung

Aus Gründen der zentralen Administration werden Terminal Server heute in einem breiten Aufgabenspektrum eingesetzt. Nicht nur für Aufgaben, wo bislang kleine Rechner oder Terminals für einfache Dateneingabe bzw. -abfrage Verwendung fanden, sondern auch in Umgebungen, in denen ein einzelner Benutzer durchaus die Rechen- oder Grafikleistung eines dedizierten PCs benötigt.

Vor jeder Implementierung eines Applikationsservers steht immer die gleiche Frage: Welches ist die passende Hardware für die geforderte Aufgabe? Natürlich möchte man dabei ein möglichst optimales System, welches weder für die Anforderungen zu klein noch (aus Kostengründen) total überdimensioniert ist. Die Frage ist also: Wie findet man ein wohl dimensioniertes System?

Die einzige meist vorliegende Kenngröße ist die Anzahl Benutzer, die mit dem System arbeiten sollen. Die typische Frage, die also zumeist auftritt, ist: »*Welches PRIMERGY Modell benötigt man für einen Terminal Server zur Unterstützung einer bestimmten Anzahl von Benutzern?*«. Optimalerweise würde man als Antwort eine handliche Tabelle erwarten, aus der anhand der Benutzerzahl in der einen Spalte unmittelbar aus der zweiten Spalte das ideale PRIMERGY System abgelesen werden kann. Leider gibt es eine solche Tabelle nicht – auch wenn mancher Mitbewerber dies dem Kunden mit bunten Web-Seiten suggeriert. Die Antwort auf die scheinbar so einfache Frage ist doch wesentlich komplexer, denn sie enthält eine große Unbekannte, und die ist der *Benutzer*. Ein Benutzer ist, auch wenn dies viele vielleicht wünschen, keine standardisierte und berechenbare Größe, sondern ein Individuum mit unterschiedlichem Arbeitstempo und Arbeitsverhalten. Hinzu kommen unterschiedliche Arbeitsaufgaben, die in unterschiedlichen Anforderungen an ein Computersystem resultieren. Ein Benutzer, dessen Aufgabe aus Abfragen an ein Lagerhaltungssystem besteht, wird eine andere Last auf einem Computersystem erzeugen, als ein Benutzer, dessen Aufgabe es ist, eine grafische Werbebroschüre zu entwerfen.

Benutzer

Um den unterschiedlichen Einsatzszenarien und Anwendern gerecht zu werden und dennoch eine Vereinheitlichung zu erreichen, definiert man Benutzergruppen. Dabei befassen sich neben den Autoren von Sizing Guides auch Marktforschungsinstitute mit diesem Thema. Die von der Gartner Group getroffene Einteilung dürfte eine der gebräuchlichsten in der IT-Branche sein (Quelle »TCO: A Critical Tool for Managing IT« Gartner Group, 12.10.1998). Darin wird eine Vielzahl von Benutzergruppen definiert:

High-Performance Worker	verwendet EDV zur Erstellung von Produkten nutzt sehr spezialisierte Anwendungen Ingenieure, Graphiker und Programmierer
Knowledge Worker	verwendet EDV zur Sammlung von Daten aus unterschiedlichen Quellen nutzt ein Mix aus Office- und spezialisierten entscheidungsunterstützenden Anwendungen Analysten, Berater und Projekt-Manager
Mobile Worker	Im wesentlichen ein Knowledge Worker, jedoch Standort unabhängig nutzt ein Mix aus Office-Anwendungen Analysten, Berater und Projekt-Manager
Process Worker	verwendet EDV zur Bearbeitung immer wiederkehrender Aufgaben in einem Produktionsprozess nutzt einen Mix aus Office- und Enterprise-Anwendungen Sachbearbeiter, Kundendienst und Helpdesk
Data Entry Worker	Nutzt die EDV zur Eingabe von Daten verwendet zumeist nur eine Anwendung Bestellwesen, Wareneingang und Verwaltungsaufgaben

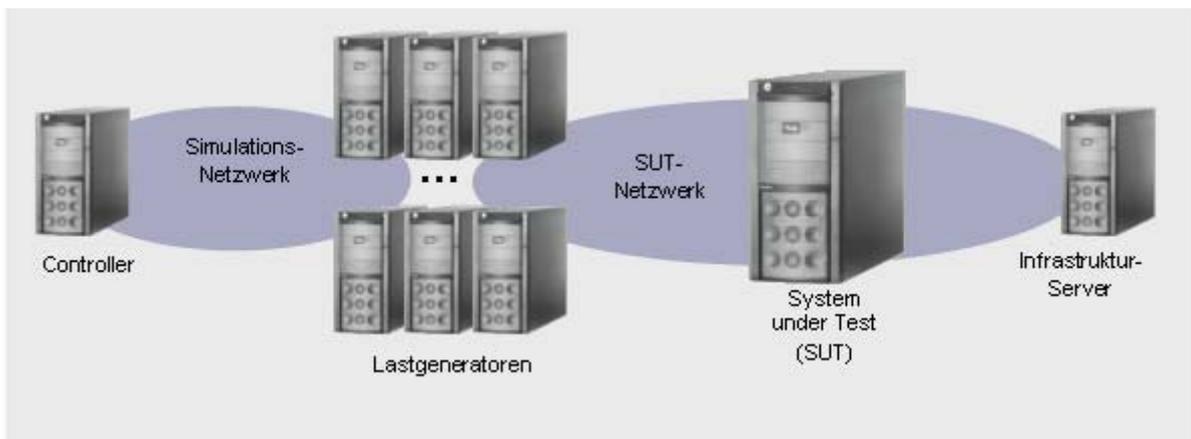
Für Anwendungen, die auf Terminal Server basieren, sind nicht alle Benutzergruppen von Belang. Die Gruppe der High-Performance Worker nutzt typischerweise dedizierte Workstations; und Mobile Worker nutzen Anwendungen lokal auf ihren mobilen Arbeitsplätzen. Benötigen diese Benutzergruppen in Ergänzung zu ihren lokalen Arbeitsplätzen Terminal Server Anwendungen, so können sie der Gruppe der Knowledge Worker zugeordnet werden.

Die Klassifizierung der Benutzer und Einteilung in Gruppen nach Gartner sagt jedoch noch nichts über die tatsächliche Aktivität aus, also konkret, welche Anwendung Benutzer einer bestimmten Gruppe nutzen und mit welcher Intensität. Insbesondere die Arbeitsgeschwindigkeit, also wie schnell ein Benutzer Text eingibt oder auf Dialoge der Anwendung reagiert, spielt eine große Rolle. Berücksichtigt man dieses, so kann man basierend auf Gartner drei Benutzerklassen definieren, die zumeist mit Heavy, Medium und Light bezeichnet werden:

Heavy	Knowledge Worker	nutzt gleichzeitig mehrere verschiedene Anwendungen gibt Daten mäßig schnell ein führt komplexere Operationen aus
Medium	Process Worker	arbeitet zu einer Zeit intensiv mit einer Anwendung gibt schnell viele Daten ein arbeitet kontinuierlich
Light	Data Entry Worker	arbeitet zu einer Zeit nur mit einer Anwendung gibt wenig Daten ein längere Pausen zwischen den Eingaben

Benutzersimulation

Bei Leistungsmessungen werden generell keine realen Benutzer verwendet, sondern die Benutzer werden mit Hilfe von Computern, so genannten Lastgeneratoren, und einer speziellen Software simuliert. Dabei wird von einem physikalischen Lastgenerator zumeist eine Vielzahl von logischen Benutzern simuliert, so dass je nach Lastgenerator einige zig oder hundert Benutzer simuliert werden können. Die folgende Abbildung zeigt eine typische Simulationsumgebung.



Der Controller ist die zentrale Steuerkonsole, die die Simulation steuert und überwacht. Über ein Simulations-Netzwerk ist dieser mit den Lastgeneratoren verbunden. Jeder Lastgenerator kann eine Vielzahl von Benutzern simulieren. Die Lastgeneratoren erreichen das Testsystem (System unter Test (SUT)) über ein zweites Netzwerk, in dem sich auch noch ein Infrastruktur-Server befindet. Dieser liefert dem SUT die notwendigen Dienste, aber er wird selbst nicht vermessen.

Unterschiedliche Benutzergruppen, wie oben diskutiert, werden von Lastsimulatoren zumeist durch verschiedene Lastprofile, im Terminal Server Umfeld auch Skript genannt, berücksichtigt.

Bei der Benutzersimulation unterscheidet man die Begriffe »Lastgenerator«, »Client« und »Benutzer«. Im weiteren Verlauf wird als »Lastgenerator« die Hardware bezeichnet. Ein »Client« ist der Terminal Server-Client, von dem einer oder mehrere auf dem Lastgenerator ausgeführt werden. Ein simulierter »Benutzer« arbeitet innerhalb einer Terminal Server Sitzung.

Für Terminal Server gibt es verschiedene Werkzeuge zur Simulation von Last. Einige der bekanntesten sind:

Terminal Server Scalability Planning Tool	<ul style="list-style-type: none"> • Eine Lastsimulator-Suite von Microsoft (http://www.microsoft.com), die Bestandteil des Windows Server 2003 Resource Kit ist. • Arbeitet nur mit Microsoft Terminal Services. • Modifiziert den RDP-Client zur Simulation der Benutzereingaben. • Simulation von Maus-Eingaben nur bedingt möglich. • Es werden modifizierbare Skripte für Lastprofile mitgeliefert, die ein aufwändiges Testszenario bedingen und nicht nur Terminal Server, sondern auch BackOffice Dienste, wie Exchange und SQL in den Test mit einbeziehen. • Es können maximal 20 Benutzer mit einem Lastgenerator simuliert werden. • Zur Vermessung von load-balanced Terminal Server-Farmen nur bedingt geeignet.
CSTK	<ul style="list-style-type: none"> • Ein kostenloser Lastsimulator von Citrix (http://www.citrix.com/cdn). • Arbeitet nur mit Citrix Presentation Server zusammen, da nur das ICA-Protokoll unterstützt wird. • Die eigentliche Simulation der Eingabe-Daten erfolgt nicht auf dem Lastgenerator (Client), sondern auf dem Server. • Nur Simulation von Tastatureingaben möglich. • Die Ausgaben der Terminal Server Session werden nicht auf Korrektheit überprüft. • Es werden Skripte für Lastprofile mitgeliefert, die aber nicht veränderbar und einsehbar sind. • Kundenspezifische Skripte müssen manuell und mit Hilfe eines kostenpflichtigen 3rd-Party Tools erstellt werden. • Produziert Last, aber ermittelt keine Antwortzeiten von einzelnen Aktionen, nur die Gesamtlaufzeit für ein komplettes Skript ist messbar. • CSTK ist instabil und für leistungsfähige Server-Systeme nicht nutzbar. Die Messergebnisse sind nicht reproduzierbar. Nutzbar als Test-Tools aber nicht zur Messung geeignet. • Zur Vermessung von load-balanced Terminal Server-Farmen ungeeignet.
CitraTest	<ul style="list-style-type: none"> • Ein kommerzielles Produkt der gehobenen Preisklasse von Tevron (http://www.tevron.com). • Kann sowohl für Microsoft Terminal Server als auch für Citrix Presentation Server verwendet werden. • Das Simulationswerkzeug läuft ausschließlich auf dem Client, ohne Client und Server zu modifizieren. • Es werden Tastatur- und Maus-Eingaben simuliert. • Die Ausgaben der Terminal Server Session können auf Korrektheit überprüft werden. • Es werden keine Skripte zur Benutzersimulation mitgeliefert. • Tool-gestützte Erstellung kundenspezifischer Skripte möglich. • Etwa nur 5 bis 10 Benutzer auf einem Lastgenerator simulierbar. • Zur Vermessung von load-balanced Terminal Server-Farmen geeignet.
LoadRunner for Citrix	<ul style="list-style-type: none"> • Ist ein kommerzielles Produkt der gehobenen Preisklasse von Mercury Interactive (http://www.mercuryinteractive.com). • Unterstützt nur Citrix Presentation Server, jedoch keinen Microsoft Terminal Server. • Die Simulation der Eingaben erfolgt mittels eines modifizierten ICA-Clients. • Es werden Text- und Maus-Eingaben simuliert. • Die Ausgaben der Terminal Server Session werden auf Korrektheit überprüft. • Tool-gestützte Erstellung kundenspezifischer Skripte möglich. • Zur Vermessung von load-balanced Terminal Server-Farmen ungeeignet.

Wie die Auflistung zeigt, sind viele dieser Lastsimulatoren leider sehr spezialisiert und nicht universell und uneingeschränkt nutzbar. Einige können nur mit einer speziellen Version des Terminal Servers zusammenarbeiten, bei anderen kann der zu simulierende Benutzer nicht modelliert werden, wiederum andere verfälschen möglicherweise das Messergebnis durch zusätzliche Komponenten der Simulationssoftware auf Client oder Server.

Vergleichbarkeit

Anders als bei anderen Benchmarks, wo es vom Hersteller der Applikation oder einem unabhängigen Gremium einen Benchmark und ein entsprechendes Reglement zur Durchführung gibt, wie z.B. bei Microsoft Exchange Server, SAP R/3, SPECweb oder TPC-C, gibt es für Terminal Server bis heute keinen standardisierten und akzeptierten Benchmark.

Es gibt zwar verschiedene Lastgeneratoren, die auch mit vordefinierten Lastprofilen bereitgestellt werden, wie es beim Terminal Server Scalability Planning Tool von Microsoft und dem CSTK von Citrix der Fall ist, jedoch mangelt es an einem Reglement bezüglich der Messumgebung, der Durchführung der Messung und standardisierten Lastprofilen, so dass durch diesen Spielraum jeder Hersteller andere Resultate ermittelt hat. Ein weiterer Mangel ist die Tatsache, dass es kein standardisiertes Werkzeug gibt, mit dem sowohl Microsoft als auch Citrix Terminal Server vermessen und verglichen werden können.

Die Ergebnisse solcher Performance-Messungen verschiedener Hersteller oder Benchmark-Labore sind unter diesen Bedingungen natürlich nicht untereinander vergleichbar. Nur Messungen, die in gleicher Umgebung und mit gleichem Lastprofil durchgeführt wurden, können auch sinnvoll verglichen werden. Daher haben Microsoft und Fujitsu Siemens Computers zusammengearbeitet, um die Ergebnisse ihrer zwei Messwerkzeuge auf der gleichen PRIMERGY Hardware zu vergleichen. Die Unterschiede in den Ergebnissen des »Microsoft Terminal Server Capacity and Scaling« Werkzeugs und des Fujitsu Siemens Computers »T4US« Werkzeugs werden detailliert im Kapitel »[Vergleich der Messwerkzeuge](#)« diskutiert.

Des Weiteren ist zu beachten, dass Performance-Messungen nicht in realen Produktivumgebungen, sondern in idealisierten Laborumgebungen durchgeführt werden. Zwar wird versucht, diese möglichst realitätsnah nachzubilden, es können jedoch nicht alle kundenspezifischen Gegebenheiten berücksichtigt werden.

Obleich die Einheit vieler Performance-Messungen »Anzahl Benutzer pro Server« ist, sollte man die Ergebnisse von Performance-Messungen in erster Linie relativ betrachten, also beispielsweise »ein System A ist doppelt so leistungsfähig wie ein System B« oder »die Verdopplung des Arbeitsspeichers resultiert in x% Leistungssteigerung«. Denn wie bereits im Kapitel »[Benutzer](#)« erläutert, ist ein Benutzer schwer zu quantifizieren, und ein synthetischer Benutzer muss nicht in allen Fällen mit einem realen Benutzer korrelieren.

Mit dieser Ausgabe liegt Version 3.x des PRIMERGY Terminal Server Sizing Guides vor. Zwischen jeder Ausgabe haben sich die Randbedingungen grundlegend geändert, so dass sich die in den bisherigen Dokumenten genannten absoluten Benutzerzahlen leider nicht miteinander vergleichen lassen. Zum einen hat jeweils ein Generationswechsel in den Betriebssystemen und Terminal Servern vorgelegen, zum anderen musste leider auch die Messmethodik den wechselnden Voraussetzungen in der IT-Landschaft angepasst werden. So zeigte z.B. das für die Version 2.0 des Sizing Guides eingesetzte Lastsimulationstool gravierende Schwächen hinsichtlich der Stabilität und Reproduzierbarkeit der Ergebnisse, so dass bei heutiger Leistungsfähigkeit der PRIMERGY Server lediglich ein Monoprozessorsystem vermessen werden könnte (vgl. Kapitel »[Benutzersimulation](#)«). Um dennoch einen Vergleich zu älteren PRIMERGY Systemen vornehmen zu können, wurden in den Messreihen für diese Ausgabe auch einige Prozessoren älterer PRIMERGY Server einbezogen, so dass Rückschlüsse auf die Leistungsfähigkeit zwischen älteren und aktuellen PRIMERGY Server gezogen werden können (vgl. Kapitel »[Prozessortyp](#)«).

Es hat sich gezeigt, dass viele Messtools, z.B. CSTK, im Vergleich zur Realität zu hohe Benutzerzahlen liefern. Ein Grund hierfür ist, dass sich während der gesamten Messphase kein Benutzer an- oder abmeldet. In unseren neuen Messreihen haben wir dem Rechnung getragen und können daher davon ausgehen, dass die ermittelten Benutzerzahlen denen aus realen Produktionsumgebungen nahe kommen.

Trotzdem ist zu beachten, dass es sich bei den Terminal Server Sizing Messungen um Auswertungen in einer vereinfachten, idealisierten und standardisierten Umgebung handelt, um vergleichbare Bedingungen für alle Systeme zu schaffen. Zusätzliche Komponenten und Programme sind nicht installiert, und der Terminal Server wird bis an seine Leistungsgrenze belastet. In der Realität wird man Zusatzsoftware wie zum Beispiel einen Virenschanner installiert haben, oder man betreibt weitere Add-Ons von Citrix Presentation Server oder Komponenten aus der Citrix Access Suite, für die Rechenleistung zur Verfügung stehen muss. Der Terminal Server sollte im Normalbetrieb auch nicht bis an seine Leistungsgrenze belastet werden.

Es hat sich auch gezeigt, dass sich allein durch leichte Modifikationen des Benutzerprofils hinsichtlich der Eingabegeschwindigkeit die Benutzeranzahl pro Terminal Server verdoppeln oder halbieren kann. Die Ergebnisse dieser Untersuchungen findet man im Kapitel »[Eingabegeschwindigkeit](#)«.

»Tool for User Simulation«

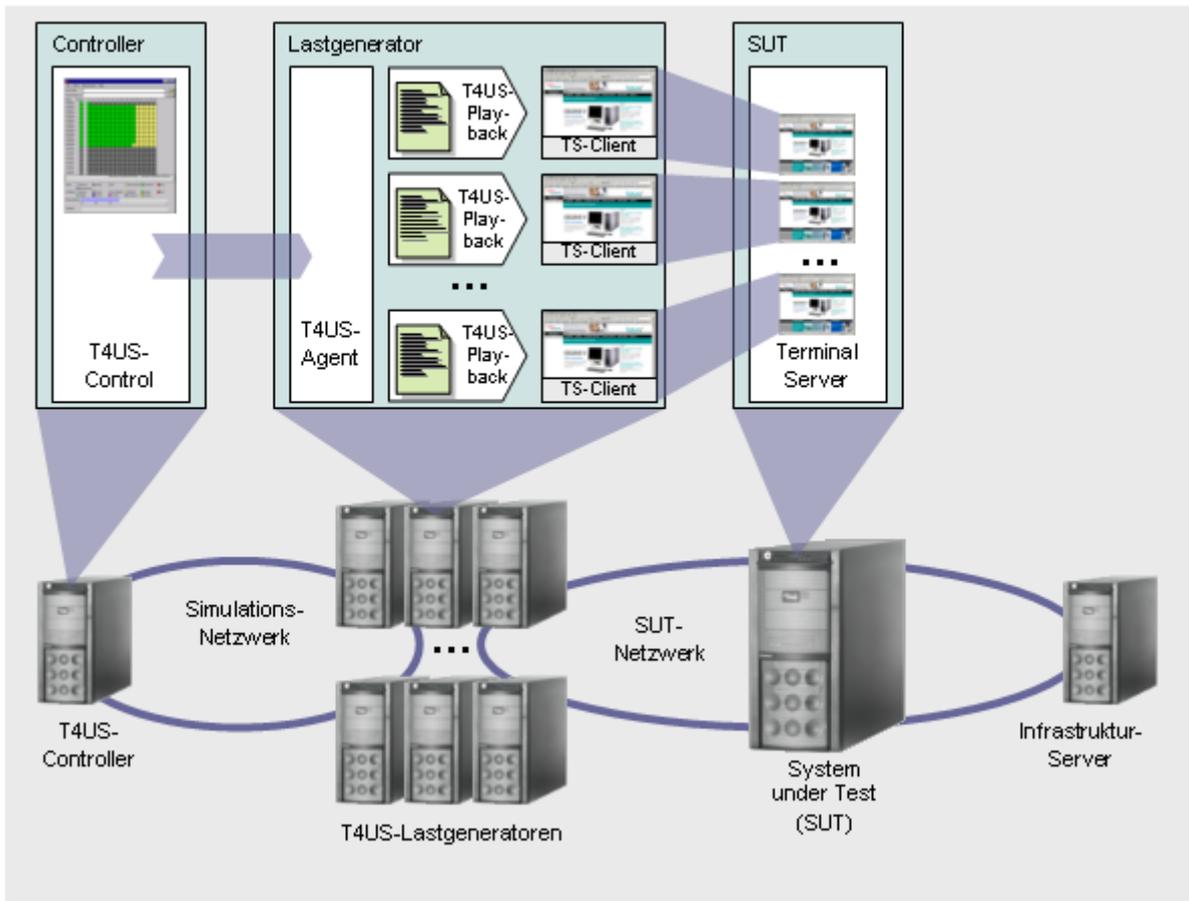
Aus den Gründen, die in den vorhergehenden Abschnitten diskutiert werden, hat Fujitsu Siemens Computers sich entschlossen, einen eigenen Lastsimulator zu entwickeln, der keinen dieser Nachteile besitzt und der unabhängig von dem verwendeten Terminal Server und ohne Einflüsse auf das zu testende System beliebige Benutzerprofile simulieren kann.

T4US, »Tool for User Simulation«, ist ein flexibles Werkzeug, das beliebige Terminal Server-artige Szenarien simulieren kann, unabhängig vom verwendeten Betriebssystem und von der Anwendersoftware, und eine detaillierte Messwerterfassung von Antwortzeiten und Auslastung unterschiedlichster Systemkomponenten vornimmt.

Benutzeraktivitäten können mit Hilfe des Aufzeichnungswerkzeugs **T4US-Record** in Echtzeit aufgezeichnet werden. Dazu gehören die Tastatur- und Mauseingaben, die Zeiten zwischen den einzelnen Eingaben, sowie die Bildschirmausgaben. Alle Aktionen werden in lesbarer Form in einem **T4US-Skript** abgelegt. In der Simulation werden diese aufgezeichneten Eingabedaten mit identischem Zeitverhalten simuliert und die Bildschirmausgaben mit den aufgezeichneten verglichen. Dabei sind die Simulationsläufe jederzeit reproduzierbar. Verschiedene T4US-Skripte können miteinander kombiniert werden, so dass aus den Aufzeichnungen von unterschiedlichen Benutzeraktivitäten beliebige Lastprofile zusammengestellt werden können. Sollen die Lastprofile an verschiedene Umgebungen angepasst werden oder für eine Vielzahl von Benutzern gleichzeitig verwendet werden, so können Teile mit Hilfe von variablen Parametern wie Benutzername, Servername, Domainname usw. parametrisiert werden ohne das Zeitverhalten zu beeinflussen.



Der Lastsimulator von T4US besteht aus drei Komponenten. **T4US-Control** ist die zentrale Steuerkonsole. Über eine grafische Oberfläche wird der gesamte Simulationslauf zentral gesteuert und überwacht. Alle Messwerte werden bereits während der Messung ermittelt und über ein separates LAN, das Simulationsnetzwerk, von den Lastgeneratoren zur Steuerkonsole übermittelt und dort gesammelt. Bereits während der Messung können die Werte automatisch ausgewertet und so zur dynamischen Steuerung des Messlaufs verwendet werden.

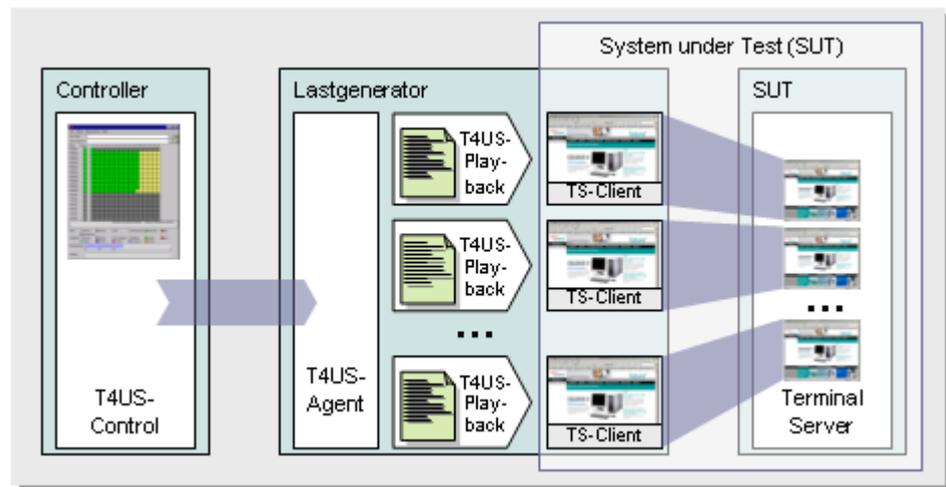


Auf jedem der Lastgeneratoren laufen mehrere Instanzen des **T4US-Playback**. Jedes T4US-Playback »füttert« einen Terminal Server-Client in Echtzeit mit Tastatur- und Mauseingaben anhand der mit T4US-Record aufgezeichneten Skripte und überwacht die Bildschirminhalte des Terminal Server-Clients. Anhand der Bildschirminhalte erfolgt die Synchronisation; das Skript wartet so lange, bis der erwartete Bildschirminhalt vollständig erschienen ist. Durch hoch auflösende Timer wird die Antwortzeit des Terminal Servers ermittelt. Die Synchronisation ist besonders wichtig für ein verlässliches Messwerkzeug, da dadurch einerseits Fehleingaben vermieden werden und andererseits so erst die Reaktionszeit des Terminal Servers deutlich und messbar wird.

Jede Instanz von T4US-Playback kann dabei ein beliebiges Skript ausführen, so dass ein Mix unterschiedlicher Benutzergruppen und asynchrones Benutzerverhalten simuliert werden kann. Auf jedem der Lastgeneratoren läuft ein **T4US-Agent**, der für die Kommunikation mit dem Controller zuständig ist, die Instanzen von T4US-Playback steuert und überwacht und die ermittelten Antwortzeiten zum Controller überträgt.

Mit T4US erfolgt die gesamte Lastsimulation von außen, ohne den Terminal Server-Client zu modifizieren oder zusätzliche Software auf dem Terminal Server zu installieren. Selbst für die Kommunikation zwischen Controller und den Lastgeneratoren wird ein separates Netzwerk verwendet, so dass es keine Einflüsse auf den Datentransport zwischen Terminal Server und Terminal Server-Clients gibt. Somit ist es möglich, nicht nur den Terminal Server zu vermessen, sondern auch die Einflüsse verschiedener Clients oder Client-Optionen, wie z.B. Bildschirmauflösung, -farbtiefen oder Audioausgabe, auf die Netzwerkbandbreite zu ermitteln. Das »System under Test« (SUT), wie man das System, welches vermessen wird allgemein bezeichnet, besteht also nicht nur aus dem Terminal Server selbst, sondern genau genommen aus den Terminal Server-Clients, dem Netzwerk zwischen Clients und dem Terminal Server, sowie dem Terminal Server selbst. T4US wird diesem Sachverhalt gerecht, in dem es keinen Eingriff in diese Client-Netzwerk-Server-Beziehung macht.

Die Terminal Server-Clients und die Komponenten T4US-Agent und T4US-Playback laufen zusammen auf den Lastgeneratoren. Durch



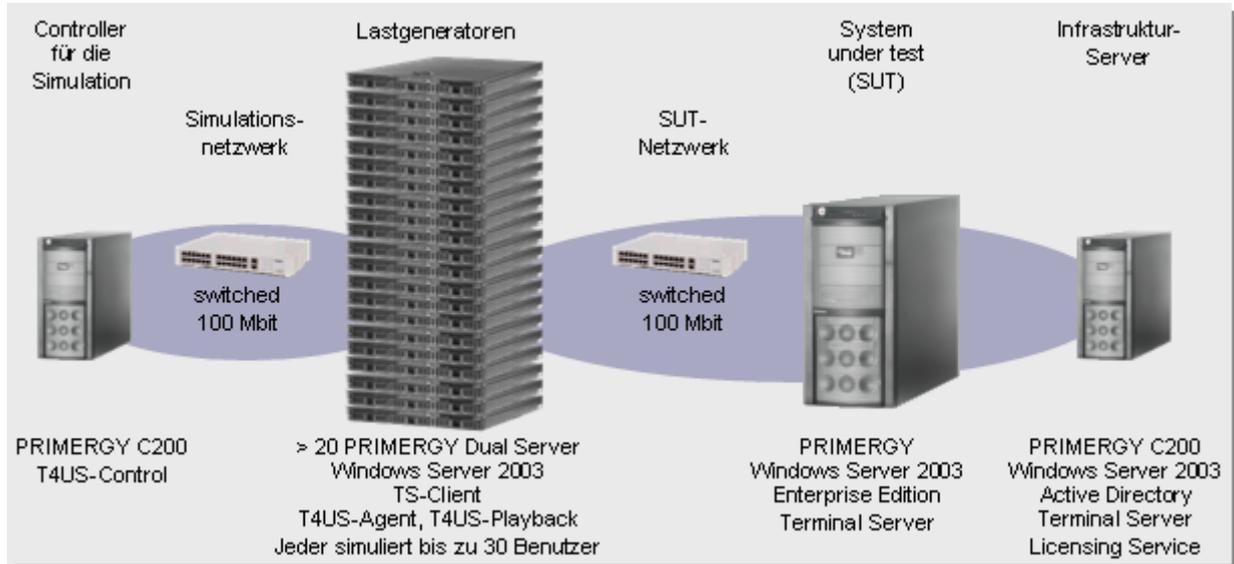
Vergleichsmessungen wird jedoch sichergestellt, dass die Hardware der Lastgeneratoren so dimensioniert ist, dass sie keinen Engpass darstellt und keinen negativen Einfluss auf die Terminal Server-Clients und somit auf das Messergebnis hat. Weiterhin kann man optional einen Lastgenerator als so genannten Referenz-Client betreiben, der nur einen einzigen Benutzer simuliert, während alle anderen Lastgeneratoren eine Vielzahl von Benutzern simulieren. Durch Vergleich der Messwerte des Referenz-Clients mit denen der anderen Clients kann eine Einflussnahme der Lastgeneratoren auf die Messergebnisse ausgeschlossen werden.

Der sich noch im SUT-Netzwerk befindende Infrastruktur-Server stellt dem zu vermessenden Terminal Server Basis-Dienste wie Active Directory, Domain Name Service (DNS) und Terminal Services Licensing zur Verfügung, er wird selbst nicht vermessen.

Messumgebung

Nachdem wir uns im vorangegangenen Kapitel allgemein mit Benutzerklassen, Benutzersimulation und Lastgeneratoren auseinandergesetzt haben, kommen wir nun zu den für die PRIMERGY Server-Familie durchgeführten Performance-Messungen.

Untersucht wurden alle aktuellen PRIMERGY Modelle, die für den Einsatz als typischer Terminal Server geeignet sind, in der folgenden Messumgebung:



Controller (T4US-Control):

- Auf dem Controller kam Windows Server 2003 Standard Edition zum Einsatz.

Lastgeneratoren:

- 20 - 24 Lastgeneratoren mit jeweils zwei Pentium III Prozessoren mit mehr als 1 GHz und 1 GB Arbeitsspeicher wurden eingesetzt.
- Die Lastsimulatoren liefen unter dem Betriebssystem Windows Server 2003 Standard Edition SP1.

Clients:

- Für den Zugriff auf den Terminal Server über das ICA-Protokoll wurde der Citrix Terminal Server-Client (Programm Neighborhood mit 32-bit ICA-Client) verwendet (entweder Version 7.00.17534 aus »Citrix MetaFrame XP Presentation Server« Feature Release 3 oder Version 9.00.32649 aus »Citrix Presentation Server 4.0«).
- Der RDP-Client (»Remote Desktop«) von Microsoft ermöglicht den Zugriff auf einen Terminal Server über das RDP-Protokoll. In Windows Server 2003 Standard Edition ist die Version 5.2.3790.1830 des RDP-Clients enthalten, der das RDP-Protokoll V5.2 unterstützt.

Netzwerk:

- Die Anbindung der Lastsimulatoren an das SUT-Netzwerk erfolgte über ein 100 MBit-Ethernet-Netzwerk, wobei der Terminal Server über den Gigabit-Uplink angeschlossen war. Das Netzwerkprotokoll war TCP/IP.

Terminal Server (System under Test):

- Die vermessenen PRIMERGY Server (System under Test) waren jeweils mit Windows Server 2003 Enterprise Edition ausgestattet. Die Terminal Services waren im Application Server Modus aktiviert.
- Bei den Messungen, bei denen ein Citrix Terminal Server vermessen wurde, war entweder Citrix MetaFrame Enterprise Edition mit Service Pack 3 und Feature Release 3 oder Citrix Presentation Server installiert. Die Terminal Server-Farm bestand nur aus einem Terminal Server. Der Data Store befand sich lokal auf der Systemplatte des zu vermessenden Systems und war als Microsoft Access Datenbank realisiert.
- Auch die Dateien der Benutzer, die während der Messung gelesen und geschrieben wurden, lagen lokal auf dem Terminal Server.
- Die Benutzerprofile wurden standardmäßig auf dem Terminal Server gespeichert.

Infrastruktur-Server:

- Der Infrastruktur-Server stellt dem System unter Test notwendige Dienste zur Verfügung, er selbst wurde nicht vermessen. Der Server wurde so dimensioniert, dass er keinen Engpass darstellt.
- Die Benutzerkonten der simulierten Benutzer wurden auf dem Active Directory Domain Controller angelegt. Ein Login fand immer gegen das Active Directory statt.
- Das Active Directory System dient gleichzeitig als DNS Server und als Terminal Server Licensing Service.
- Der Infrastruktur-Server lief unter dem Betriebssystem Windows Server 2003 Standard Edition SP1.

Diese synthetische Messumgebung vereinfacht eine realistische Kundenumgebung stark, um Einflüsse anderer Systeme auszuschließen und reproduzierbare Ergebnisse zu erhalten. Der Einfluss weiterer Komponenten in einer Terminal Server Umgebung wird im Kapitel »[Infrastruktur](#)« diskutiert.

Lastprofil

Alle Messungen wurden mit einem Medium Lastprofil durchgeführt. Wie im Kapitel »[Benutzer](#)« definiert, arbeitet ein »Medium User« mit nur einer Anwendung und gibt Daten zügig ein. In unserem Medium Lastprofil dient Microsoft Word als Anwendung und der Benutzer schreibt einen bebilderten Text mit einer durchschnittlichen Eingaberate von 230 Anschlägen pro Minute.

Das Lastprofil wurde in einem realen Szenario aufgezeichnet und die Eingabe der Zeichen entspricht der realen Arbeitsgeschwindigkeit eines mit 10 Fingern schreibenden Autors.

Des Weiteren beinhaltet das Medium Lastprofil:

- Jeder Benutzer arbeitet unter einem eigenen Benutzerkonto.
- Die erste Anmeldung (Login) des Benutzers und der erste Start der Anwendung liegen außerhalb der Messstrecke. Jedoch meldet sich der Benutzer nach einmal getaner Arbeit am Terminal Server ab und für einen neuen Durchlauf wieder an. Da die Benutzer versetzt gestartet werden, ergibt sich so während der gesamten Messdauer ein kontinuierliches An- und Abmelden.
- Jeder Terminal Server-Client (Benutzer) startet die Applikation aus seinem Desktop heraus, die Applikation wird bei jedem Skriptdurchlauf gestartet und beendet.
- Jeder Benutzer hat sein eigenes Verzeichnis, in dem die im Text verwendeten Bilder hinterlegt sind. So wird verhindert, dass alle Benutzer die gleichen Bilder-Dateien laden und sich diese nach kurzer Zeit alle im Server File Cache befinden. Jeder Benutzer schreibt bei jedem Skriptdurchlauf ein neues Dokument mit eindeutigem Namen. Nach erfolgreicher Erstellung wird das Dokument mit der Größe von ca. 227 KB auf die Festplatte des Terminal Servers in ein benutzereigenes Verzeichnis gespeichert.
- Die durchschnittliche Eingabegeschwindigkeit liegt bei etwa 4 Zeichen bzw. Cursor-Bewegungen pro Sekunde. Allerdings finden nicht während des gesamten Durchlaufes Eingaben statt, denn es sind diverse, unterschiedlich lange Denkzeiten im Skript eingestreut, wie es einem natürlichen Arbeiten nahe kommt.
- Die Bildschirmauflösung ist 1024x768, die Farbtiefe ist 16-bit.
- Ein Durchlauf des Skripts inklusive Wartezeiten dauert ca. 16 Minuten.

Da es bei diesem Sizing Guide um einen relativen Vergleich der PRIMERGY Modelle untereinander geht, wurde auf Untersuchungen mit weiteren Lastprofilen verzichtet. Dies würde zwar zu einer anderen Anzahl von Benutzern pro PRIMERGY führen, die Relation zwischen den einzelnen PRIMERGY Modellen wäre jedoch die gleiche.

Bei einer Aussage von absoluten Benutzerzahlen auf einem Server muss ohnehin der kundenspezifische Last-Mix analysiert und mit den Leistungsdaten in diesem Papier in Relation gesetzt werden (vgl. Kapitel »[Vergleichbarkeit](#)«).

Messmethode

Zu einer Leistungsmessung gehört neben einem Simulationswerkzeug und einem möglichst realistischen Lastprofil ein Regelwerk, nach dem die einzelnen Messungen durchgeführt und bewertet werden.

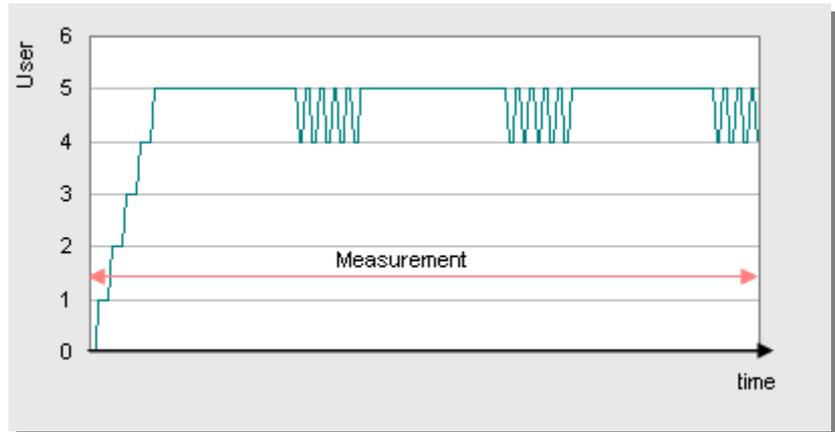
Messarten, Messdauer und Messphasen

Vor der Messung werden grundsätzlich alle Systeme, d.h. Lastgeneratoren und Server under Test, inklusive der T4US-Client Komponenten T4US-Agent und T4US-Playback, neu gestartet. Auf dem Controller-System wird der T4US-Controller jedes Mal neu gestartet.

T4US unterstützt drei verschiedene Messfunktionen:

Referenzmessung mit konstanter Benutzeranzahl

Eine konstante, aber geringe Anzahl von Benutzern lässt ein oder mehrere T4US-Skripte mehrmals durchlaufen. In der Terminal Server Messumgebung wurde das T4US-Skript mindestens dreimal hintereinander von fünf Benutzern ausgeführt. Diese Messdaten werden gesammelt und aus ihnen berechnet der T4US-Controller für jeden einzelnen Messpunkt Vergleichswerte, die als Baseline für die weiteren Messungen dienen.



Messung mit konstanter Benutzeranzahl

Bei der Messung mit konstanter Benutzeranzahl arbeitet eine gleich bleibende Anzahl von Benutzern über einen vorgegebenen Zeitraum mit dem Terminal Server.

Als Resultat der Messung erhält man die Antwortzeiten des Terminal Servers und Performance Counter des Servers.

Die Messung selbst unterteilt sich in drei Phasen:

Startphase (Startup) 15 Minuten	Während der Startphase nehmen nach und nach alle T4US-Playback's auf Befehl des T4US-Controllers ihre Arbeit auf. Hierbei verteilt der T4US-Controller den Start der Skripte gleichmäßig auf die Startphase, die immer 15 Minuten dauert; unabhängig davon, wie viele Benutzer simuliert werden sollen. Dies entspricht der Realität, da leistungsstärkere PRIMERGY Server insgesamt mehr Benutzer bedienen können und diesen daher auch mehr Anmeldungen in der gleichen Zeit zugemutet werden als leistungsschwächeren Systemen. Auf eine ungleichmäßig gestaffelte Verteilung der Anmeldungen, wie sie in anderen Messungen gern gemacht wird, wurde bewusst verzichtet, da in der Realität der Benutzer auch nicht länger mit seiner Anmeldung warten wird, nur weil schon viele Benutzer arbeiten. Die Startphase ist beendet, wenn alle Skripte gestartet sind.
Einschwingphase (Warm-up) 30 Minuten	Während der Einschwingphase laufen alle T4US-Playback's gemäß den eingestellten Skripten ab.
Messphase (Steady State) 60 Minuten	Die jetzt folgenden 60 Minuten dienen der Erhebung von Messdaten. Es werden die Performance Counter des Terminal Servers ausgewertet sowie die von den T4US-Playback's gemeldeten Antwortzeiten des Terminal Servers.

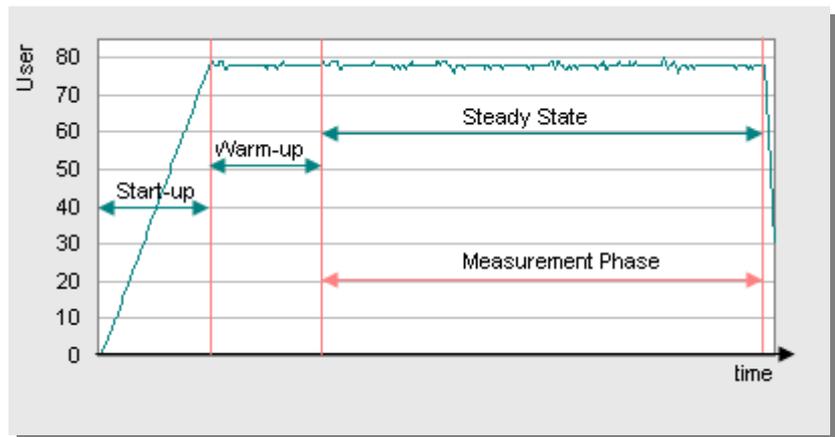
Sollte während einer der Phasen ein Fehler festgestellt werden, führt dieser zum Abbruch und zur Wiederholung der Messung.

Während aller Phasen werden Messdaten erhoben und kontrolliert, aber nur die Messdaten, deren Beginn und Ende vollständig in die Messphase fallen, werden zur Auswertung herangezogen. Die Antwortzeiten des Terminal Servers werden von den T4US-Playback's registriert und an den Controller gemeldet.

Die Performance Counter des Servers werden vom Controller abgefragt. Die Daten des Terminal

Servers im Zeitraum der Messphase werden ausgewertet. Die Daten der anderen beteiligten Systeme wie Lastgeneratoren, Controller und Infrastruktur-Server werden zur Kontrolle überwacht, um sicherzustellen, dass diese nicht überlastet sind oder dass eine Messung durch Seiteneffekte ungültig ist.

Alle Messungen des Terminal Server Sizing Guides V3.0 wurden mit diesem Typ der Messung durchgeführt.



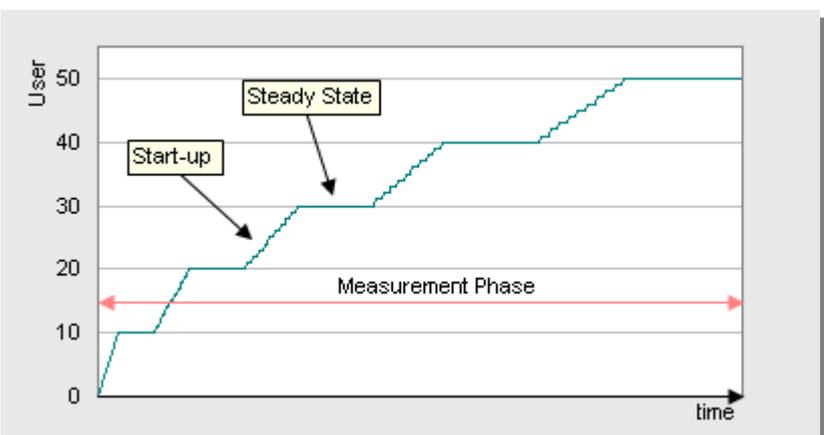
Messung mit variabler Benutzeranzahl

Bei der Messung mit variabler Benutzeranzahl wird die Anzahl der Benutzer, die mit dem Terminal Server arbeiten, nach einer voreingestellten Regel kontinuierlich erhöht, bis der Terminal Server überlastet ist.

Während der gesamten Messung werden die Antwortzeiten des Terminal Servers von dem T4US Controller überwacht. Dieser vergleicht jeden einzelnen Messwert mit einem gespeicherten Referenzwert, der aus einer vorhergehenden [Referenzmessung](#) ermittelt wurde. Als Maßgabe für die Überlastung des Servers werden bestimmte Ende-Kriterien konfiguriert.

Als Resultat der Messung erhält man eine Benutzeranzahl (»Score«).

Bei der Messung mit variabler Benutzeranzahl wechseln sich Phasen, bei denen neue Terminal Server Benutzer hinzukommen, mit Phasen ab, in denen die Benutzeranzahl stabil bleibt. Während der einzelnen Startphasen nimmt nach und nach ein Teil der T4US-Playback's auf Befehl des T4US-Controllers die Arbeit auf. Hierbei verteilt der T4US-Controller den Start der Skripte gleichmäßig auf die Startphase. Eine Startphase ist beendet, wenn alle Skripte der



T4US-Playback's gestartet sind. Während der jetzt folgenden »Steady State« Phase ist die Messung in einem stabilen Zustand, es werden keine neuen Benutzer hinzugefügt. Die Startphasen und stabilen Phasen werden kontinuierlich wiederholt. Dabei können im frühen Zeitraum der Messung viele Benutzer schnell gestartet werden, während die Dauer der Startphasen im Verlauf der Messung immer weiter zunimmt, dadurch wird der Abstand der Benutzeranmeldung vergrößert. Durch diese Dehnung der Benutzeranmeldungen wird das Messergebnis genauer und reproduzierbar.

Sollte während der Messung ein Fehler festgestellt werden, führt dieser zum Abbruch und zur Wiederholung der Messung.

Die Performance Counter des Terminal Servers werden vom Controller abgefragt und über die gesamte Messphase ausgewertet. Die Daten der anderen beteiligten Systeme wie Lastgeneratoren, Controller und Infrastruktur-Server werden zur Kontrolle überwacht, um sicherzustellen, dass diese nicht überlastet sind oder dass eine Messung durch Seiteneffekte ungültig ist.

Alle Messungen des Terminal Server Sizing Guides ab Version 3.1 wurden mit diesem Typ der Messung durchgeführt.

Prozessorauslastung

Es ist insbesondere festzulegen, wann ein Server ausgelastet ist. Denn sicherlich lässt sich auf einem System, auf dem n Applikationen laufen, auch noch eine $n+1$ -te starten. Aber es ist ja nicht zweckmäßig, einen Server beliebig zu überlasten. Dies würde nur ermitteln, wie »dehnbar« die Verwaltungstabellen des Betriebssystems sind. Vielmehr muss man ein Maß für stabiles Arbeiten des Systems finden. Ist dieses überschritten, so wird das System überlastet und wird instabil. Alle Windows Server Betriebssysteme bieten hierfür eine Vielzahl von Performance Countern, die Auskunft über den Systemzustand geben. Ein Indikator für die Aus- bzw. Überlastung des Systems ist die »Processor Queue Length«. Dieser Counter gibt an, wie viele Threads auf ihre Ausführung durch die CPU warten. Steigt dieser Counter signifikant und kontinuierlich an, so ist dies ein Hinweis auf eine Überlastung des Systems. Dabei ist zu bedenken, dass unabhängig von der Anzahl Prozessoren nur ein Counter für die Queue-Länge geführt wird. Auch bei einem hohen Wert für die Queue-Länge muss die prozentuale CPU-Auslastung nicht nahe 100% sein, auch bei niedrigerer CPU-Auslastung von unter 50% kann es zum Ansteigen der Prozessor-Queue kommen. Dies tritt dann auf, wenn sich eine Vielzahl von Prozessen im Idle-Zustand befindet, ein Zustand, der beim Terminal Server-Szenario insbesondere durch eine Vielzahl von Benutzern erreicht wird, die im Prinzip nichts weiter tun, als eine Applikation offen zu halten (vgl. Kapitel »[Anzahl Prozesse](#)«).

Reaktionszeit

Ein zweites Maß für die Stabilität des Servers ist die Antwortzeit, mit der ein Server auf Eingaben des Benutzers reagiert. Sie hängt natürlich unmittelbar mit Prozessor-Auslastung und Prozessor-Queue-Länge zusammen.

Bei den **Messungen mit konstanter Benutzeranzahl** wird der Terminal Server so weit belastet, bis die durchschnittliche CPU-Auslastung über 70% liegt, die Prozessor-Queue signifikant ansteigt oder die Antwortzeit der Applikation sich gegenüber einer Referenzzeit um mehr als 10% verlängert.

Um die Referenzzeit festzulegen, wird auf fünf Lastgeneratoren je eine Instanz des betreffenden Lastprofils gestartet und dreimal erfolgreich durchlaufen. Die Referenzzeiten sind hauptsächlich von den Wartezeiten innerhalb der Skripte begrenzt und unterscheiden sich nur minimal von System zu System. Die Referenzzeiten selbst werden nicht verwendet, um die Leistungsfähigkeit des betreffenden PRIMERGY Systems zu dokumentieren, sondern nur, um die Verlängerung der Antwortzeiten zu berechnen.

Um die Antwortzeiten zu bestimmen, wird aus allen Messdaten der einzelnen Messpunkte, die die Clients während der Messphase ermitteln und an den T4US-Controller senden, der Durchschnitt gebildet und mit den Referenzzeiten verglichen.

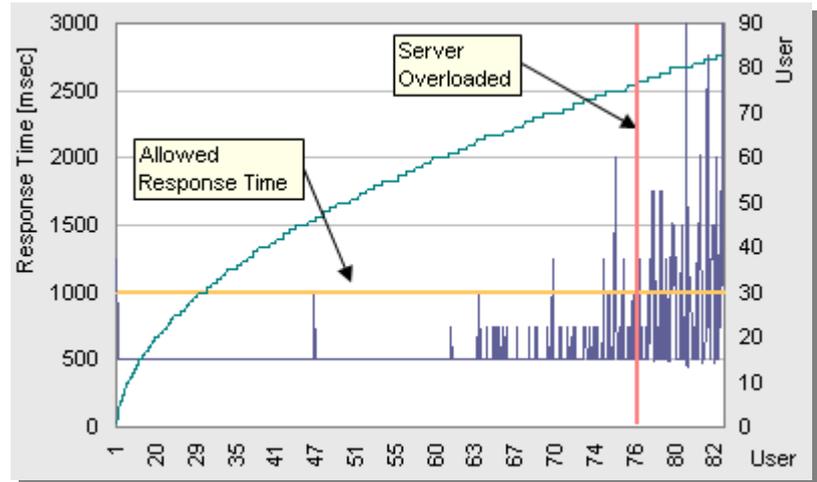
Da von allen Messwerten während der Messphase der Durchschnitt gebildet wird, ist der Prozentsatz, um den sich die Antwortzeiten verschlechtern dürfen, mit 10% nicht sehr hoch.

Bei einer Messung mit konstanter Benutzeranzahl wird die Anzahl Benutzer vorgegeben, mit der der Terminal Server während der Messung arbeitet, und erst nach Durchlauf der Messung kann festgestellt werden, ob der Terminal Server diesen Anforderungen noch gewachsen war oder nicht.

Bei den **Messungen mit variabler Benutzeranzahl** wird der Terminal Server so weit belastet, bis sich die Antwortzeit der Applikation gegenüber einer Referenzzeit so weit verschlechtert, dass sie den vorgegebenen Regeln nicht mehr genügt. Hierbei werden die Performance Counter zwar mitgeschrieben und nach der Messung ausgewertet, jedoch einzig die Antwortzeit des Terminal Servers wird verwendet, um die Anzahl der Benutzer zu ermitteln, die zufriedenstellend mit dem Terminal Server arbeiten können.

Die Maßgabe für die Benutzerzufriedenheit kann im T4US Controller konfiguriert werden. Bei den diesem Terminal Server Sizing Guide zu Grunde liegenden Messungen darf sich jeder einzelne Messwert um 30% (Werte über 1500 ms) bzw. 100% (Werte bis 1500 ms) verschlechtern. Insgesamt müssen 90% der Messwerte innerhalb dieses Rahmens liegen, 10% Ausreißer werden toleriert. Jeder weitere Messwert, der die Grenze von 30% (bzw. 100%) Verschlechterung übersteigt, führt zur Beendigung der Messung. Die so erreichte Benutzeranzahl ist das Ergebnis der Messung und wird als »Score« bezeichnet.

Die Grafik veranschaulicht die Arbeitsweise des T4US Controllers bei der Auswertung der Messwerte. Die waagerechte Linie bei 1000 ms Antwortzeit ist die eingestellte Grenze, die 90% der Messwerte nicht überschreiten dürfen. Einige Ausreißer sind erlaubt, diese werden durch nachfolgende Werte innerhalb des erlaubten Bereichs kompensiert. Werden jedoch zu viele Ausreißer erkannt, gilt der Terminal Server als überlastet. Die Messung wird an dieser Stelle beendet, das Bild zeigt außerdem, wie sich die Messwerte weiter verschlechtern würden, wenn noch weitere Benutzer hinzugefügt würden. In diesem Beispiel ist der Score »76 Benutzer«.



Um die Referenzzeit festzulegen, wurde auf fünf Lastgeneratoren je eine Instanz des betreffenden Lastprofils gestartet und dreimal erfolgreich durchlaufen. Die Referenzzeiten sind hauptsächlich von den Wartezeiten innerhalb der Skripte begrenzt und unterscheiden sich nur minimal von System zu System. Die Referenzzeiten selbst wurden nicht verwendet, um die Leistungsfähigkeit des betreffenden PRIMERGY Systems zu dokumentieren, sondern nur, um die Verlängerung der Antwortzeiten zu berechnen.

Im Vergleich zur Messung mit konstanter Benutzeranzahl, bei der eine Verschlechterung der Antwortzeiten von 10% erlaubt ist, arbeitet die Messung mit variabler Benutzeranzahl mit einem höheren Limit von 30%. Dies resultiert daraus, dass bei der Messung mit variabler Benutzeranzahl jeder einzelne Messwert dieses Limit einhalten muss und nicht nur der Durchschnitt.

Tuning

Auch wenn es von Microsoft und Citrix zahlreiche Artikel zur Optimierung von Betriebssystem- und Terminal Server gibt, so haben wir bei unseren Messungen doch gänzlich auf eine Optimierung verzichtet. Der Grund ist, dass viele dieser Einstellungen nur in bestimmten Umgebungen Sinn machen; in einem anderen Umfeld eingesetzt bewirken sie oft das Gegenteil. Da wir bei dieser Messreihe verschiedene PRIMERGY Systeme mit unterschiedlichen Systemkomponenten untersucht haben, wären die Ergebnisse nicht miteinander vergleichbar gewesen.

Die einzigen Einstellungen, die verändert werden um alle PRIMERGYs den gleichen Testbedingungen zu unterwerfen, sind:

- Das Pagefile des Betriebssystems wurde auf eine feste Größe von 4 GB eingestellt, um eine Fragmentierung zu vermeiden und damit für alle getesteten Server die gleichen Bedingungen vorliegen.
- Für Citrix musste die Grenze von 100 Benutzern pro Server, die durch das eingebaute Load Balancing vorgegeben wird (auch wenn die Terminal Server-Farm aus nur einem Terminal Server besteht) aufgehoben werden.

Die bisher unter Windows 2000 Server notwendige Vergrößerung der Registry auf einem Terminal Server entfällt bei Windows Server 2003.

Ressourcenbedarf

In diesem Dokument wird Terminal Server sowohl unter dem 32-bit Betriebssystem »Windows Server 2003 R2« als auch unter dem 64-bit Betriebssystem »Windows Server 2003 R2 x64« untersucht. Die 32-bit und 64-bit Versionen von Windows Server 2003 R2 basieren dabei auf der gleichen Code-Basis und sind daher direkt vergleichbar. Windows Server 2003 R2 ist bis auf einige zusätzliche Dienste und Tools identisch mit Windows Server 2003 Service Pack 1. Für die 64-bit Messungen wurden die gleichen Randbedingungen verwendet wie bei den 32-bit Messungen. Die simulierten Benutzer arbeiteten in beiden Fällen mit dem Medium Lastprofil unter Verwendung von Microsoft Office 2003.

Für ein Server-System sind die folgenden Performance-relevante Faktoren maßgeblich:

- Rechenleistung
- Arbeitsspeicher
- Disk-Subsystem
- Netzwerk

Je nach Aufgabe des Servers haben die Einzelkomponenten verschiedene Gewichtung in Bezug auf die Gesamtleistung des Servers. Im Folgenden wollen wir diskutieren und anhand von Messergebnissen untermauern, welche Komponenten welchen Einfluss auf die Leistungsfähigkeit eines Terminal Server-Systems haben.

Bei dieser Betrachtung wird zur besseren Vergleichbarkeit gezielt eine einzelne Komponente untersucht und ein Engpass dieser Komponente provoziert, während die anderen Ressourcen ausreichend dimensioniert sind. Daher kann es bei einigen Messreihen auf den ersten Blick so aussehen, als hätte das 32-bit Betriebssystem Vorteile gegenüber dem 64-bit Betriebssystem, da es weniger CPU- und Hauptspeicherressourcen benötigt. Bei moderner Server-Hardware mit Multi-Core-Architektur ist jedoch ausreichend Rechenleistung vorhanden, so dass sich in der Realität die Limitierungen Richtung Arbeitsspeicher und Kernel-Strukturen verschieben werden, von denen das 64-bit Betriebssystem größere Einheiten verwalten kann.

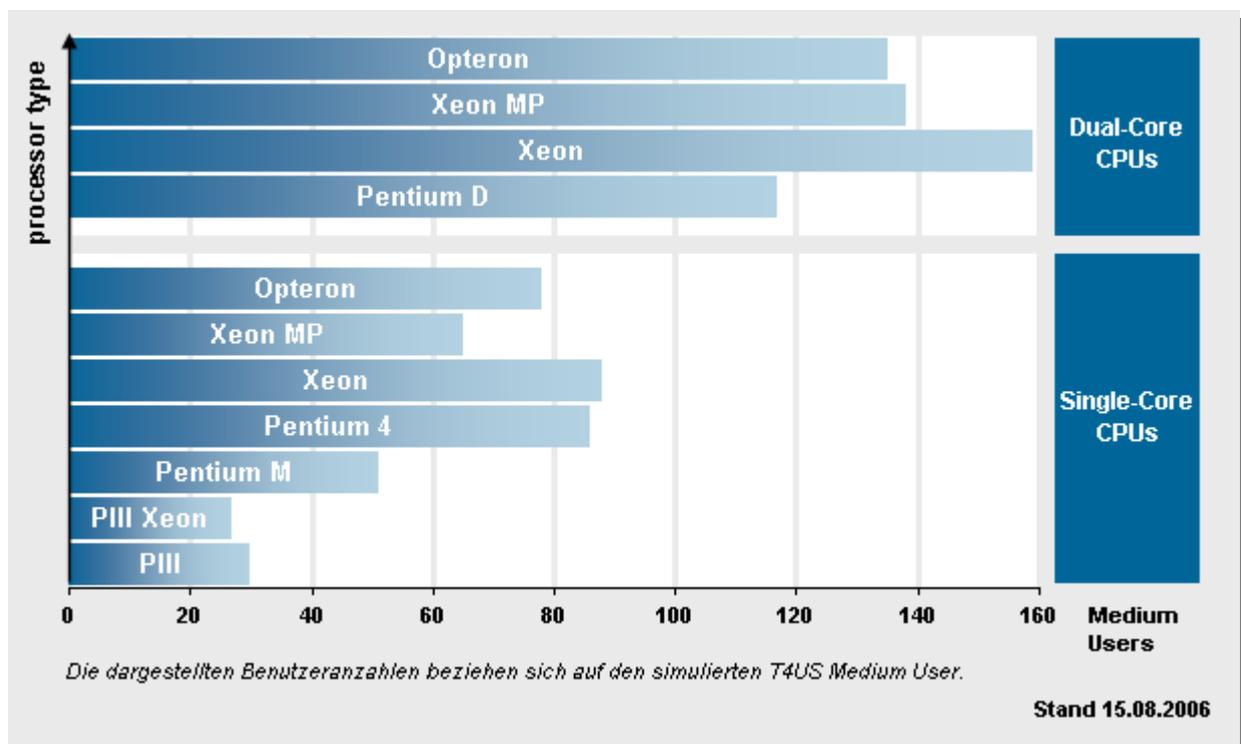
Rechenleistung

Die Rechenleistung eines Systems hängt von den Prozesseigenschaften und der Anzahl Prozessoren ab.

Prozessorotyp

Für Intel-kompatible Server steht heute eine Vielzahl an Prozessorvarianten zur Verfügung, die heute in den PRIMERGY Systemen eingebaut werden: AMD Opteron, Intel Celeron, Intel Pentium 4, Intel Pentium D, Intel Pentium M, Intel Xeon Prozessor und Intel Xeon Prozessor MP. Während der Celeron eher für den Low-End-Bereich designed ist, haben sich heute im Server-Bereich die Prozessoren AMD Opteron, Intel Pentium 4, Pentium D, Xeon und Xeon MP etabliert. Einige der heutigen Prozessorvarianten basieren auf der Dual-Core-Technologie, bei der pro Chip zwei Prozessorkerne enthalten sind.

Folgende Übersicht zeigt die verschiedenen Prozessoren, die heute für Terminal Server im Einsatz sind, mit ihren Leistungsdaten für jeweils eine CPU, wobei es Exemplare mit einem Prozessorkern (Single-Core) und Exemplare mit zwei Prozessorkernen (Dual-Core) gibt. Aus diesen Messergebnissen für eine CPU kann man die Leistung eines Systems mit einer höheren CPU-Anzahl nicht ablesen, da die Steigerung nicht linear ist. Informationen zur Skalierung der Prozessoren findet man in den folgenden Abschnitten zu den Themen »[Taktfrequenz](#)« und »[Anzahl Prozessoren](#)«.



Bei der Betrachtung der Rechenleistung waren alle Systeme mit genügend Arbeitsspeicher ausgebaut, damit bei diesem Vergleich der Arbeitsspeicher keinen Engpass darstellt. Bei der Untersuchung wurde der Arbeitsspeicher nicht variiert. Architekturbedingt kann bei den PRIMERGY Modellen mit AMD Opteron Prozessoren nur die Hälfte der Speicherbänke bestückt werden, wenn lediglich ein CPU-Sockel bestückt ist. Hyper-Threading, sofern vorhanden, war eingeschaltet, da dieses Feature bei Terminal Server-Anwendungen für eine Entlastung des Systems sorgt und dadurch eine höhere Anzahl von Benutzern mit dem Terminal Server arbeiten kann.

Die Prozessoren unterscheiden sich in ihrer Architektur sowie in Taktfrequenz, Caches, ggf. Hyper-Threading-Technologie, Anzahl Prozessorkerne und möglicher CPU-Anzahl, dies wird in den folgenden Abschnitten detaillierter betrachtet.

Taktfrequenz

Die Prozessorlinien gibt es wiederum in verschiedenen Leistungsstufen. Abhängig vom Prozessormodell gibt es bei Intel Prozessoren Taktfrequenzen von 1 GHz bis heute 3.8 GHz und die Geschwindigkeit des Front-Side-Busses reicht von 133 MHz bis derzeit 1333 MHz. Die CPU-Taktfrequenz sagt heutzutage allerdings nicht mehr unbedingt etwas über die Rechenleistung aus. Diese Tatsache ist bereits aus dem Vergleich von AMD mit Intel Prozessoren bekannt, da AMD Prozessoren die gleiche Rechenleistung mit einer niedrigeren Taktfrequenz erreichen. Beim AMD Opteron Prozessor reichen die Taktfrequenzen von 1.8 GHz bis 3.0 GHz. Aber auch Intel Prozessoren der neuesten Technologie sind niedrigerer getaktet als ihre Vorgänger und Strom sparer, bieten aber auch eine hervorragende Rechenleistung durch eine neue interne Architektur.

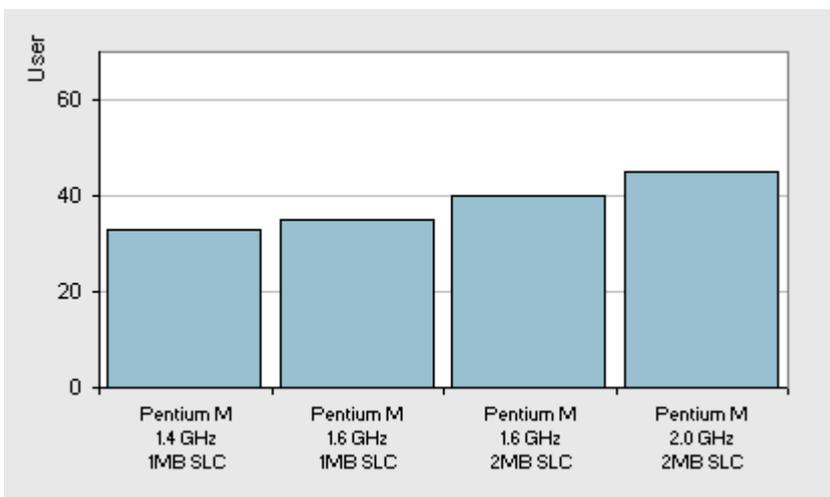
Prozessortyp	CPU Nummer	Taktfrequenzbereich	FSB
Intel Celeron	336, 346, 352	2.8 GHz – 3.2 GHz	533 MHz
Intel Pentium M	725, 755	1.6 GHz – 2.0 GHz	400 MHz
Intel Pentium 4	630 – 650	2.8 GHz – 3.6 GHz 3.0 GHz – 3.4 GHz	800 MHz
Intel Pentium D	820, 930 – 950, 925 – 945	2.8 GHz – 3.4 GHz	800 MHz
Intel Xeon		2.8 GHz – 3.8 GHz	800 MHz
Intel Xeon (Dual-Core)	5050 5060 – 5080 5110, 5120 5130 – 5160	2.8 GHz 3.0 GHz 3.2 GHz – 3.733 GHz 1.6 GHz – 1.867 GHz 2.0 GHz – 3.0 GHz	800 MHz 667 MHz 1067 MHz 1067 MHz 1333 MHz
Intel Xeon MP		2.0 GHz – 3.0 GHz 3.0 GHz – 3.667 GHz	400 MHz 667 MHz
Intel Xeon (Dual-Core)	7020, 7040 7030, 7041 7110M – 7140M	2.667 GHz – 3.0 GHz 2.8 GHz – 3.0 GHz 2.6 GHz – 3.4 GHz	667 MHz 800 MHz 800 MHz
AMD Opteron	244 – 256	1.8 – 3.0 GHz	
AMD Opteron (Dual-Core)	265 – 285 865 – 885	1.8 – 2.6 GHz	

Die Skalierung der CPU-Taktfrequenz bei Terminal Server Anwendungen wurde bei mehreren Prozessorfamilien untersucht. Taktfrequenzen können aus den oben genannten Gründen nur innerhalb einer Prozessorfamilie verglichen werden. Man sieht, dass eine höhere Taktfrequenz auch eine höhere Performance bedeutet, aber die Frequenzsteigerung spiegelt sich nicht 1:1 in der relativen Performancesteigerung wider. Dies erklärt sich aus dem gleich bleibenden Frontside-Bus und der somit unveränderten Geschwindigkeit bei Speicher- und I/O-Zugriffen. Stärkeren Einfluss als die Taktfrequenz hat die Vergrößerung des [Caches](#), der im folgenden Abschnitt im Detail diskutiert wird.

Im Folgenden wird die Skalierung der Intel-Prozessor-Familien Pentium M, Pentium 4, Pentium D, Xeon und Xeon MP sowie der AMD Opteron Prozessoren für die aktuellen PRIMERGY Systeme im Detail dargestellt.

Pentium M

Eine Besonderheit stellen die besonders Strom sparenden Pentium M Prozessoren dar, die in der PRIMERGY BX300 mit der so genannten Blade Server Technologie verwendet werden. Diese erfüllen besondere Anforderungen an Leistungsaufnahme und Wärmeentwicklung, so dass auf dichtestem Raum eine hohe Rechenleistung zur Verfügung gestellt werden kann.

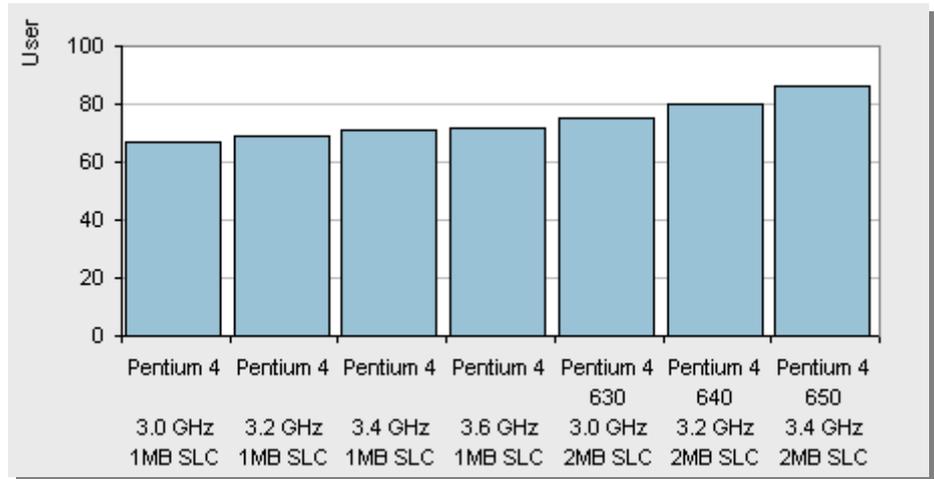


(Medium Lastprofil, Microsoft Office XP, Microsoft Terminal Services)

Die Grafik zeigt die CPU-Skalierung der Pentium M Prozessoren bei einem Mono-CPU-Blade der PRIMERGY BX300. Durch Erhöhung der CPU-Frequenz ist eine Leistungssteigerung zu erkennen, diese resultiert in einer höheren Benutzeranzahl. Die Frequenzsteigerung wird jedoch nicht 1:1 in eine Performancesteigerung umgesetzt, da die Geschwindigkeit der Speicher- und I/O-Zugriffe über den Frontside-Bus gleich bleibt. Durch eine Frequenzsteigerung um 25%, beispielsweise von 1.6 auf 2.0 GHz, erhöht sich die Benutzeranzahl um 12.5%. Die Auswirkungen von [Cachegrößen](#) werden im folgenden Abschnitt im Detail diskutiert.

Pentium 4

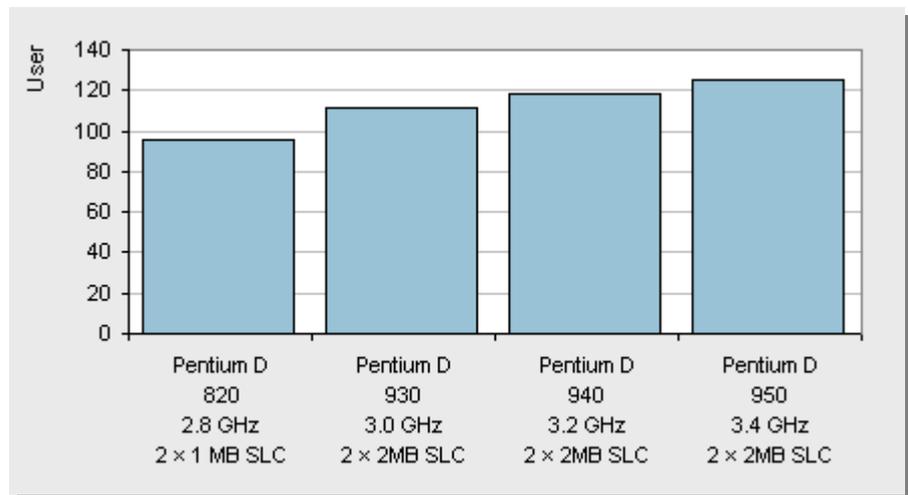
Der Pentium 4 Prozessor findet in den PRIMERGY Servern mit einem CPU-Sockel Verwendung, wie zum Beispiel PRIMERGY TX150 S2, PRIMERGY TX150 S3 oder PRIMERGY RX100 S3, wobei nicht alle Prozessoren der Pentium 4-Klasse in allen PRIMERGYs verwendet werden. Bei dieser CPU-Generation gibt es derzeit Exemplare mit 1 MB und 2 MB Second Level Cache (SLC). Auch bei diesen Messungen ist ein Mehr an Performance durch eine Steigerung der Taktfrequenz zu erkennen. Allerdings kann eine Steigerung der Taktfrequenz nicht 1:1 in eine Erhöhung der Benutzeranzahl umgesetzt werden. Bei einer Steigerung der Taktfrequenz um beispielsweise 7% steigt die Benutzeranzahl nur um 3%. Dies erklärt sich aus dem gleich bleibenden Frontside-Bus und der somit unveränderten Geschwindigkeit bei Speicher- und I/O-Zugriffen. Trotz niedrigerer Taktfrequenz bieten die Pentium 4 Prozessoren mit größerem Cache eine höhere Leistung.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Pentium D

Der Pentium D Prozessor ist eine Dual-Core Variante für 1-Socket-Systeme. In der nebenstehenden Grafik ist die Terminal Server Leistung einiger Pentium D CPUs dargestellt. Auch hier gilt natürlich: eine höhere Taktfrequenz und ein größerer CPU-Cache bieten eine höhere Leistung. Diese Prozessoren kommen in den PRIMERGY Systemen Econel 100, TX150 S4, TX150 S5 und RX100 S3 zum Einsatz, wobei für Terminal Server üblicherweise Rack-Systeme verwendet werden.

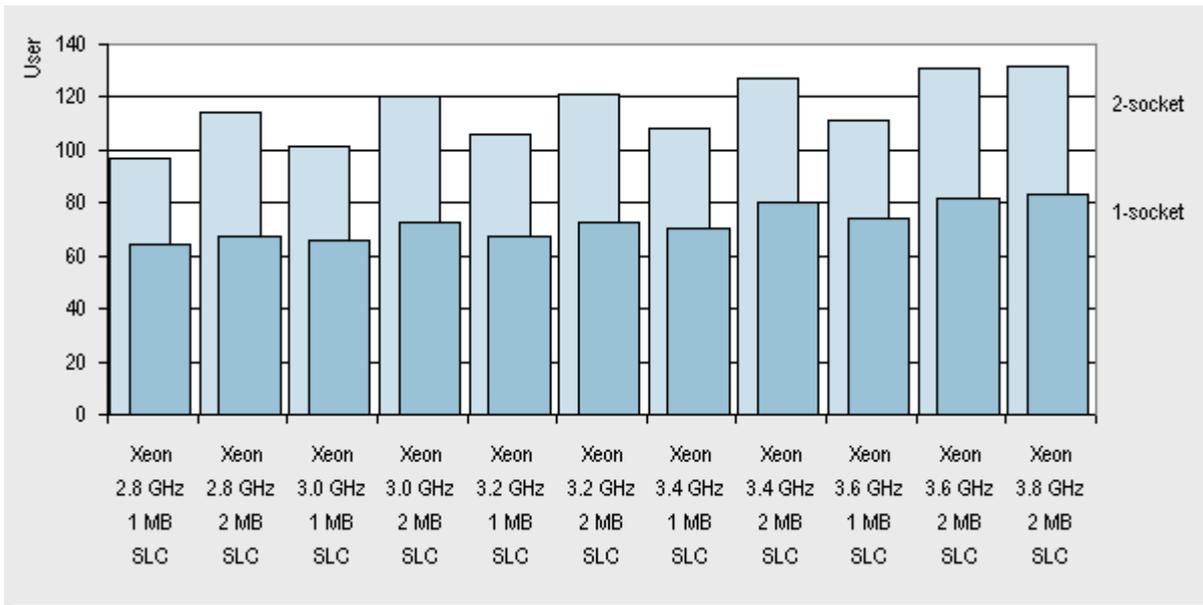


(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Der Pentium D Prozessor unterstützt auch das 64-bit Windows Betriebssystem, allerdings wird man auf einem 1-Socket-System trotz Dual-Core Prozessor nicht so viele Terminal Server Benutzer betreiben können, dass man von den Vorteilen des 64-bit Windows Betriebssystems profitieren könnte.

Xeon (Single-Core Xeon für 2-Socket-Systeme)

Untersucht man einen 2-Socket PRIMERGY Server, hier eine PRIMERGY RX300 S2 mit jeweils einem oder zwei Xeon Single-Core Prozessoren, erkennt man wieder die Leistungssteigerung durch die CPU Frequenz. Beim System mit einer CPU wird in diesem Beispiel eine Frequenzsteigerung um ca. 7% in eine um ca. 2.5% höhere Benutzeranzahl umgesetzt. Dies erklärt sich aus dem gleich bleibenden Frontside-Bus und der somit unveränderten Geschwindigkeit bei Speicher- und I/O-Zugriffen. Beim Dualprozessorsystem kommt hinzu, dass die Performancesteigerung bei den höheren Frequenzen etwas niedriger ausfällt als beim Monoprocessorsystem. Durch weitere Prozessoren wird seitens des Betriebssystems ein höherer Synchronisationsaufwand benötigt. Diese Single-Core Xeon-Generation ist momentan mit 1 MB und 2 MB SLC verfügbar. Die Performancesteigerung durch einen größeren Cache ist höher als die Performancesteigerung, die aus einer Frequenzsteigerung resultiert.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Die Auswirkungen von [Cachegrößen](#) werden später im Detail diskutiert.

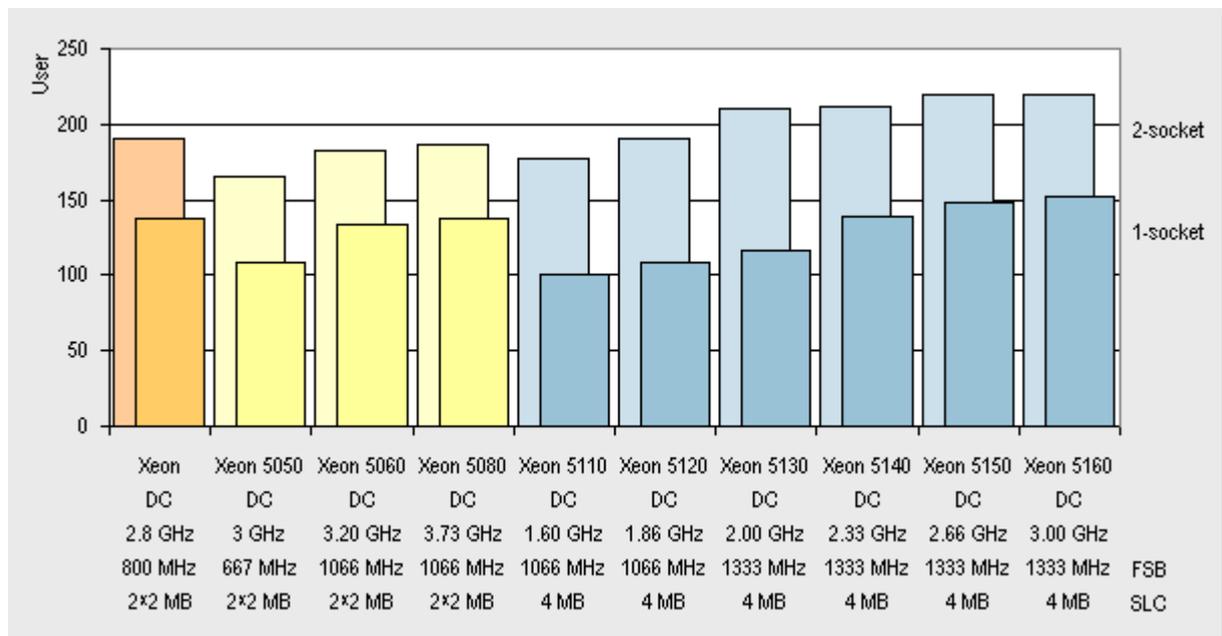
Dual-Core Xeon (für 2-Socket-Systeme)

Die Unterscheidung erfolgt bei neueren Xeon Prozessoren über die Prozessornummer. Als Dual-Core Prozessoren für die 2-Socket-Systeme gibt es folgende Prozessortypen:

- Xeon 2.8 GHz 800 MHz Front-Side-Bus Codename: Paxville
- Xeon 50xx 3 - 3.73 GHz 667 MHz Front-Side-Bus Codename: Dempsey
- Xeon 51xx 1.60 - 3 GHz 1066 MHz Front-Side-Bus Codename: Woodcrest

Alle Prozessoren haben zwei Kerne pro Chip (Dual-Core). Der Second Level Cache (SLC) ist bei den Paxville- und Dempsey-Prozessoren jedem CPU-Kern direkt zugeordnet und hat eine Größe von 2 MB pro Kern. Bei den Woodcrest-Prozessoren beträgt die Größe des Second Level Cache 4 MB und ist beiden Kernen zugeordnet. Hyper-Threading war bei den Paxville- und Dempsey-Prozessoren eingeschaltet, da dieses Feature bei Terminal Server-Anwendungen für eine Entlastung des Systems sorgt und dadurch eine höhere Anzahl von Benutzern mit dem Terminal Server arbeiten kann. Die Woodcrest-Prozessoren bieten dieses Feature nicht an. Die Messungen erfolgten unter 64-bit Windows 2003 (SP1).

Alle Dual-Core Xeon Prozessoren für 2-Socket-Systeme sind zusammen im nachfolgenden Diagramm dargestellt, obwohl sie nicht in den gleichen PRIMERGY Systemen verbaut werden. Die Paxville Prozessoren kommen in der PRIMERGY RX300 S2 zum Einsatz, während die Dempsey und Woodcrest Prozessoren in den 2-Socket-Systemen der PRIMERGY S3 Reihe angeboten werden.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services, 64-bit Windows)

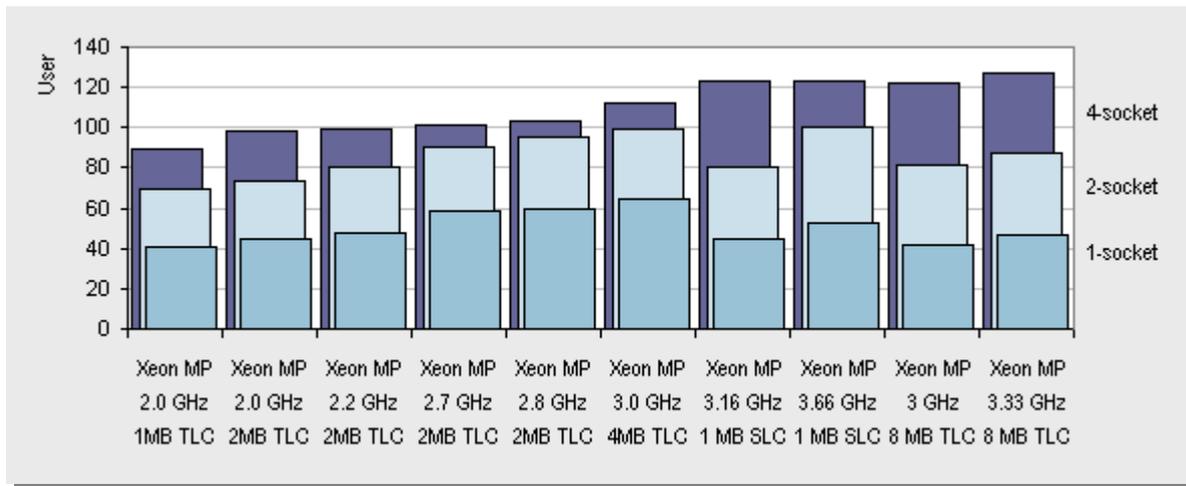
Eine Verdoppelung der CPU von einem auf zwei Prozessoren führte bei den Messungen mit Microsoft Terminal Server zu einer Leistungssteigerung von ca. 35% - 80%. Die Messungen mit den geringer getakteten Prozessoren profitierten mit einem Leistungszuwachs um ca. 80% dabei stärker von dem zweiten Prozessor als die Messungen mit den höher getakteten Prozessoren. Bei dem am höchsten getakteten Prozessor lag der Leistungszuwachs bei ca. 35%. Hier liegt die Anzahl arbeitender Benutzer in einem Größenbereich, in dem nicht allein die CPU-Leistung als begrenzender Faktor wirkt, sondern auch andere unveränderte Faktoren, wie z.B. die Geschwindigkeit bei den Speicherzugriffen.

Am Beispiel der Xeon Prozessoren sieht man den Trend bei der Prozessorentwicklung hin zu weniger Taktfrequenz bei guter Performance. Die Woodcrest-Prozessoren liefern bei wesentlich niedrigerer Taktfrequenz und geringerem Stromverbrauch eine vergleichbare Leistung. Durch die verbesserte Mikroarchitektur basierend auf neuester Technologie wird dies auch ohne Hyper-Threading erreicht.

Xeon MP (Single-Core, ab 4-Socket-Systemen)

Bei der PRIMERGY RX600/TX600-Systemreihe mit jeweils einem, zwei und vier Prozessoren kann man die Skalierung über die Taktfrequenz, Cachegröße und Prozessoranzahl beobachten. Die Performance steigt erwartungsgemäß beim Hinzufügen weiterer Prozessoren. Jedoch ist die Skalierung über die Frequenzen beim Monoprozessorssystem besser als bei den Multiprozessorssystemen. Auch macht sich beim Dualprozessorsystem eine Steigerung der Taktfrequenz deutlicher bemerkbar als beim 4-way System. Dies resultiert aus dem gleich bleibenden Frontside-Bus und der somit unveränderten Geschwindigkeit bei Speicher- und I/O-Zugriffen sowie durch den höheren Synchronisationsaufwand seitens des Betriebssystems, der mit der Anzahl der Prozessoren ansteigt.

Neben der PRIMERGY RX600 kommen Single-Core Prozessoren auch noch in der PRIMERGY RX600 S2 und RX600 S3 zum Einsatz. Im nachfolgenden Diagramm sind alle Single-Core Xeon MP Prozessoren gegenübergestellt, unabhängig davon, in welcher Generation der PRIMERGY RX600 oder TX600 Systeme sie eingesetzt werden können.



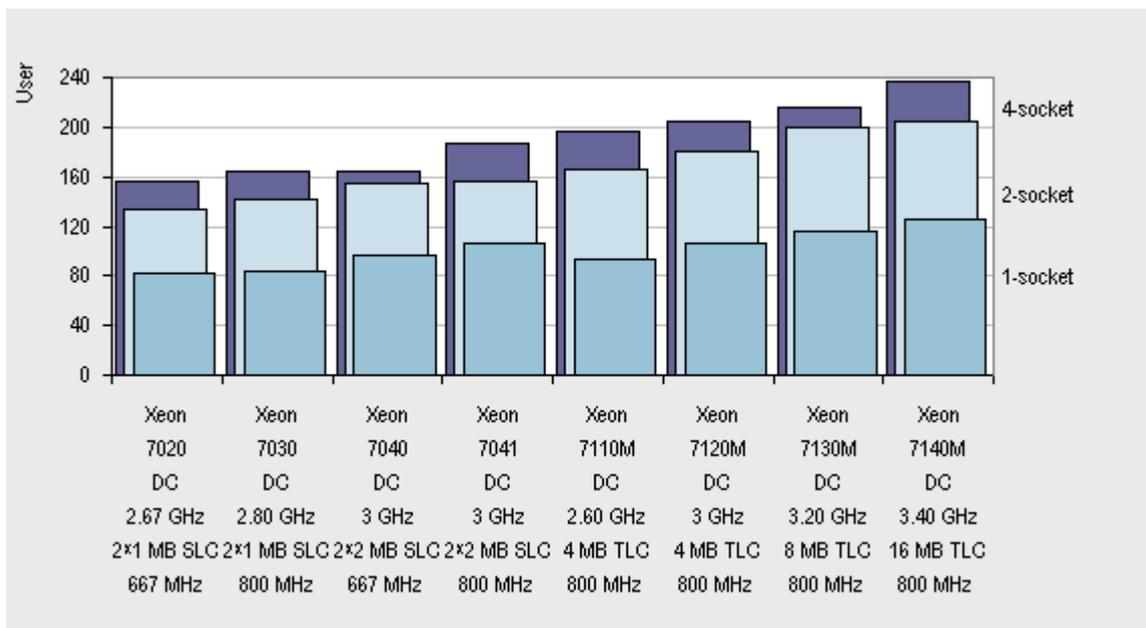
(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services, Windows Server 2003 Enterprise Edition, 32-bit)

Die Auswirkungen von [Cachegrößen](#) und der [CPU-Anzahl](#) werden später im Detail diskutiert.

Dual-Core Xeon (ab 4-Socket-Systemen)

Die multiprozessorfähigen Dual-Core Xeon Prozessoren werden in der PRIMERGY RX600 S2 und PRIMERGY RX600 S3 eingesetzt. Unter einem 64-bit Windows Server 2003 wurden diese Prozessoren in ihrer Skalierung vermessen. Die Prozessoren der Xeon 70xx Reihe mit einem 2 MB Second-Level-Cache (SLC) pro Prozessorkern und FSB-Taktgeschwindigkeiten von 667 MHz und 800 MHz sind unter dem Codenamen »Paxville« bekannt, während die Xeon Prozessoren der Reihe 71xxM mit dem Codenamen »Tulsa« einen von beiden Prozessorkernen gemeinsam benutzten Third-Level-Cache (TLC) haben und immer mit einem 800 MHz FSB betrieben werden.

Im nachfolgenden Diagramm sind alle Dual-Core Xeon Prozessoren gegenübergestellt, unabhängig davon, in welcher Generation der PRIMERGY RX600 oder TX600 Systeme sie eingesetzt werden können. Die Leistung der Prozessoren im Anwendungsszenario eines Terminal Servers ist nicht ausschließlich von der Taktfrequenz abhängig, sondern auch vom Prozessortyp, der Cache-Größe und der Geschwindigkeit des Front-Side-Busses. Vergleicht man innerhalb einer Prozessorgeneration die Leistungssteigerung, die sich aus der Erhöhung der Taktfrequenz ergibt, so erkennt man, dass sich die Frequenzsteigerung nicht 1:1 in der relativen Performance-Steigerung widerspiegelt. Dies erklärt sich aus dem gleich bleibenden Front-Side-Bus (FSB) und der somit unveränderten Geschwindigkeit bei Speicher- und I/O-Zugriffen. Einen deutlich größeren Einfluss auf die Leistung hat jedoch die Größe des CPU-Cache, denn durch diesen Cache wird der Front-Side-Bus entlastet. Die Verdopplung des Caches vom 1 MB auf 2 MB bei den Xeon 7030 und Xeon 7041 Prozessoren in Verbindung mit einer geringen Frequenzerhöhung von 0,2 GHz steigert die Benutzerzahl um ca. 28%. Einen deutlichen Leistungssprung erzielt man durch den Einsatz eines zweiten Prozessors, die Steigerung der Benutzeranzahl liegt zwischen 48% und 78%.



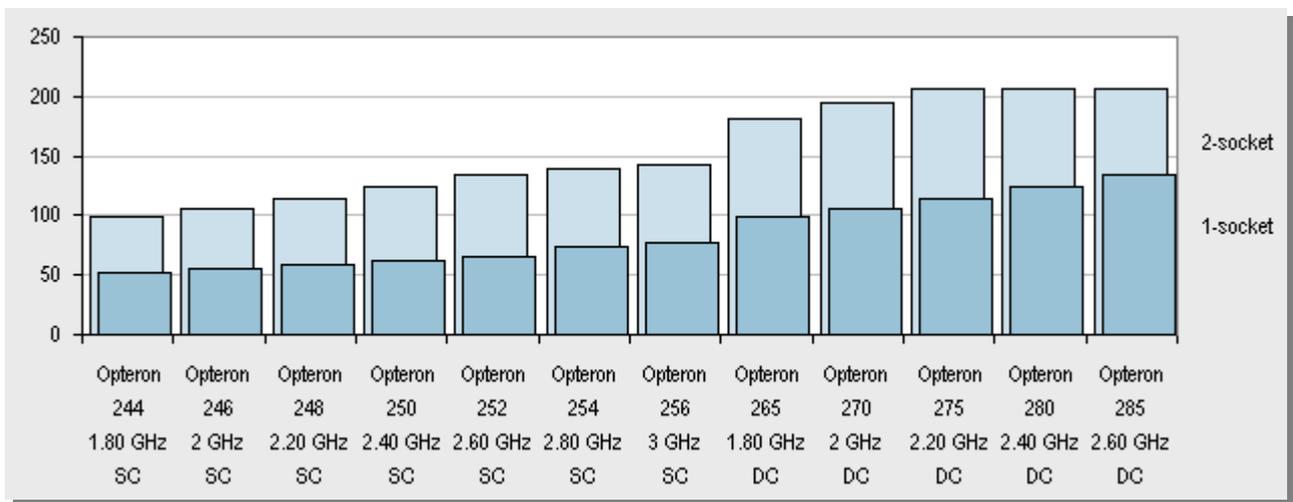
(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

In der Konfiguration mit vier Prozessoren zeigt sich, bei dem hier verwendeten Lastprofil, eine ungünstigere Skalierung verglichen mit einem Ausbau mit zwei Prozessoren. Durch jeden weiteren Prozessor steigt der Synchronisationsaufwand seitens der Hardware und des Betriebssystems an, während die Geschwindigkeit der Zugriffe auf die anderen Ressourcen gleich bleibt. Aus diesem Grund kann die Leistungsfähigkeit des Terminal Servers durch das Hinzufügen weiterer Prozessoren nicht mehr effizient gesteigert werden. Man erkennt, dass die Benutzeranzahl nicht gegen einen Grenzwert läuft, sondern durch Einsatz eines größeren Prozessors weiter erhöht werden kann, aber der Zugewinn an Leistung ist nicht so groß wie bei der Verdopplung der Prozessoren von 1 auf 2. Durch das hier verwendete Medium Lastprofil, das die An- und Abmeldung des Benutzers mit einschließt und bei dem alle Benutzer einen bebilderten Text mit Microsoft Word erstellen, tritt bei der PRIMERGY RX600 S3 im Ausbau mit vier Prozessoren eine Überlastung des Terminal Servers auf, da sich mit steigender Benutzeranzahl viele Benutzer gleichzeitig anmelden und gleichzeitig auf das Dateisystem zugreifen. Die Benutzeranzahl bleibt hinter den Erwartungen zurück, während noch ausreichend CPU- und Hauptspeicherreserven vorhanden sind. Hierbei sind die Prozessoren nur zu 40% - 50% ausgelastet, wobei das Disk-Subsystem keinen Engpass darstellt. Der Bottleneck ist in der Handhabung des Dateisystems seitens Windows zu suchen.

AMD Opteron

Der Rackserver PRIMERGY RX220 und der Blade Server PRIMERGY BX630 sind mit AMD Opteron Prozessoren ausgestattet. Diese gibt es als Single-Core und als Dual-Core Varianten, wobei die Dual-Core CPU die gleiche Charakteristik hat wie die Single-Core CPU. So hat jeder Prozessorkern 1 MB SLC, und auch die Frequenzabstufungen sind die gleichen. Aus diesem Grund werden die Single-Core und Dual-Core Prozessoren in einem Schaubild dargestellt. Dabei muss nicht notwendigerweise jede CPU auch in jedem PRIMERGY Modell verfügbar sein.

Wie nachfolgende Grafik zeigt, steigt die Leistung der Prozessoren mit der Taktfrequenz. Bei den Single-Core Prozessoren skaliert das System von ein auf zwei Prozessoren mit einem exzellenten Faktor von deutlich über 80%. Einen deutlichen Leistungssprung erzielt man durch den Einsatz eines Dual-Core Prozessors. Im oberen Leistungsspektrum bei dem 2-Socket-System mit zwei Dual-Core Prozessoren ist jedoch eine Stagnation zu erkennen, die Benutzeranzahl kann trotz einer höheren CPU-Frequenz nicht mehr gesteigert werden. Dies ist durch das hier bei den Messungen verwendete 32-bit Betriebssystem zu erklären, das bei einer Größenordnung von ca. 200 Benutzern in einen Engpass der Kernel-Ressourcen gerät. Dies wird im Kapitel »[Betriebssystem](#)« im Detail diskutiert.



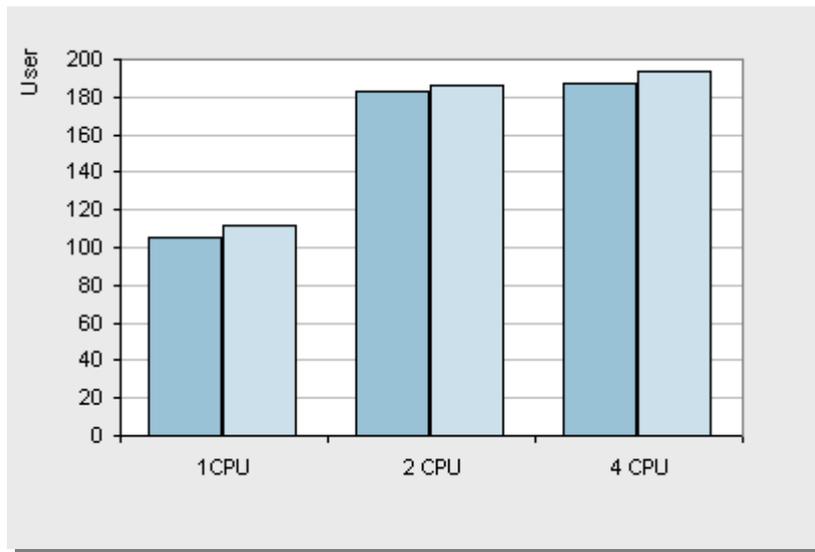
(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Front-Side-Bus

Im Folgenden wird der Einfluss des Front-Side-Bus(FSB)-Taktes bei Intel-basierten Servern untersucht.

Über den Front-Side-Bus wird die Kommunikation zwischen dem Prozessor und der so genannten Northbridge abgewickelt, einem Bestandteil des Chipsatzes, der wiederum mit anderen Komponenten wie zum Beispiel dem Arbeitsspeicher oder dem PCI-Bus verbunden ist. Der FSB wird mit einer bestimmten Taktrate betrieben, welche Einfluss auf die Systemperformance hat.

Im Allgemeinen unterscheiden sich zwei Prozessoren durch mehr Eigenschaften als nur die Taktung des Front-Side-Busses, daher ist ein direkter Vergleich meist nicht möglich. Während vor einiger Zeit pro System meist nur eine Front-Side-Bus-Taktgeschwindigkeit verwendet wurde, werden heute in modernen PRIMERGY Systemen gleich zwei oder drei verschiedene FSB-Geschwindigkeiten angeboten, je nach eingesetztem Prozessortyp.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Unter Laborbedingungen konnte der Einfluss des Front-Side-Bus-Taktes jedoch vermessen werden. Die Grafik zeigt die Performanceverbesserung durch Steigerung der Front-Side-Bus-Taktfrequenz von 667 MHz auf 800 MHz, was einer Steigerung von ca. 20% entspricht. Die Anzahl der Terminal Server Benutzer steigt in allen drei vermessenen Konfigurationen. Ein schnellerer FSB entlastet den Prozessor, die »% Processor Time« ist bei 800 MHz FSB etwas niedriger und gleichermaßen auch die »Processor Queue Length«. Durch diese CPU-Reserve können mehr Benutzer mit dem Terminal Server arbeiten, bevor die Antwortzeiten sich verschlechtern und sich nicht

mehr im zulässigen Rahmen befinden. Die Steigerung der Benutzeranzahl ist jedoch geringer als die Erhöhung der FSB-Taktfrequenz. Terminal Server als Anwendung beansprucht nicht nur Prozessorleistung, sondern zur Gesamtleistung trägt das Zusammenspiel aller Ressourcen wie Prozessor, Speicher, Netzwerk und Disk bei.

Architekturbedingt haben PRIMERGY Systeme, die mit AMD Opteron Prozessoren bestückt sind, keinen Front-Side-Bus. Der Memory Controller ist in die CPU integriert und die Prozessoren untereinander sind durch einen HyperTransport-Link verbunden.

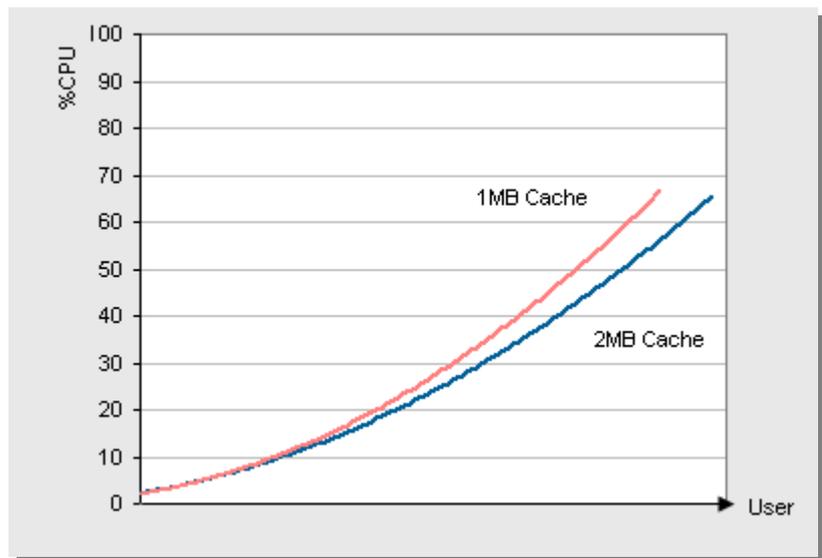
Caches

Ein Cache ist generell ein schneller Zwischenspeicher, der durch die Pufferung von Daten den Zugriff beschleunigt. Bei Intel-CPU's sind die Caches in mehreren Stufen kaskadiert. Man unterscheidet Level 1 Cache, Level 2 Cache (auch Second Level Cache (SLC) genannt) und Level 3 Cache (Third Level Cache (TLC)). Meistens wird bei den Leistungsdaten der CPU's nur der jeweils letzte Cache genannt. Der Cache soll verhindern, dass der Prozessor auf Daten des langsameren Arbeitsspeichers warten muss. Je größer der Cache, umso weniger Speicherzugriffe sind nötig. Aus dieser Zeitersparnis resultiert wiederum eine höhere Rechenleistung.

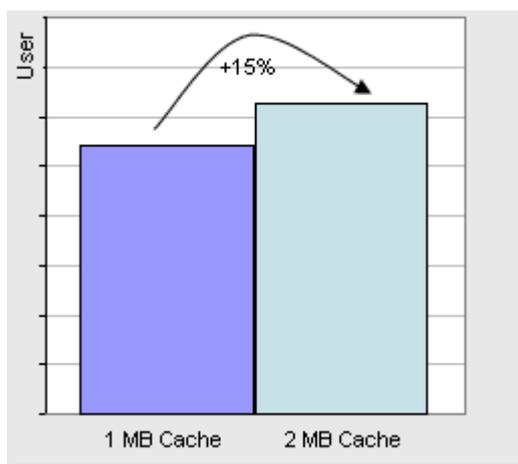
Bei den Intel Prozessoren werden neben unterschiedlichen Taktfrequenzen auch unterschiedliche Cache-Größen angeboten.

Eine Verdoppelung des Caches macht sich erfahrungsgemäß durch eine Performancesteigerung von ca. 10% bemerkbar.

Bei Terminal Server führt eine Verdoppelung des Caches zu einer CPU Entlastung. Mit dem Medium Lastprofil wurden Messungen mit variabler Benutzeranzahl durchgeführt. Bei sonst gleichen Bedingungen wurden jeweils zwei Prozessoren einmal mit 1 MB Cache und einmal mit 2 MB Cache eingesetzt. Hyper-Threading war eingeschaltet. Wie nebenstehende Grafik zeigt, wird der Prozessor des Terminal Server Systems bei dem größeren Cache weniger stark beansprucht. Bei der Messung mit variabler Benutzeranzahl wird die Messung dann beendet, wenn die Antwortzeiten des Terminal Servers den erlaubten Rahmen überschreiten. Durch die CPU-Reserven des Systems mit dem größeren Cache können mehr Benutzer vom Terminal Server bedient werden, bevor diese Grenze erreicht ist.



(Medium Lastprofil, Microsoft Office XP, Citrix MetaFrame)



(Medium Lastprofil, Microsoft Office XP, Citrix MetaFrame)

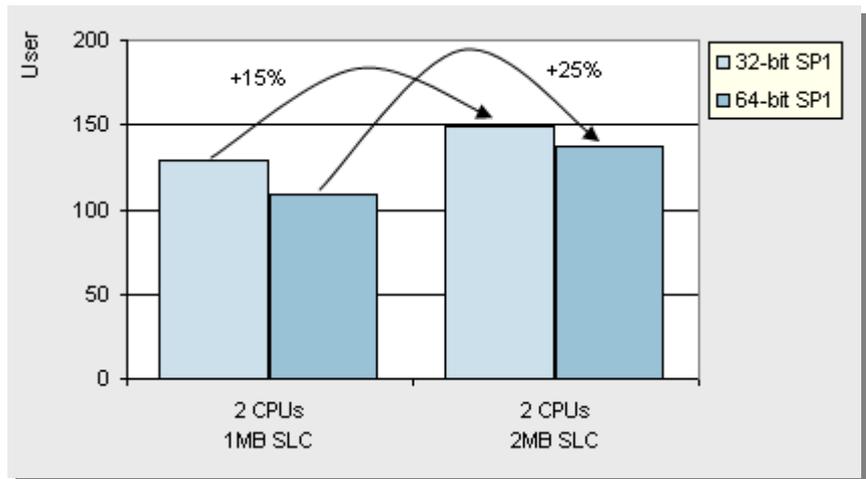
Beim Medium Lastprofil konnte eine um 15% höhere Benutzeranzahl beobachtet werden, wie die nebenstehende Grafik zeigt.

Eine Performance-Steigerung bei größerem Cache wird sich bei einem Heavy User am deutlichsten bemerkbar machen.

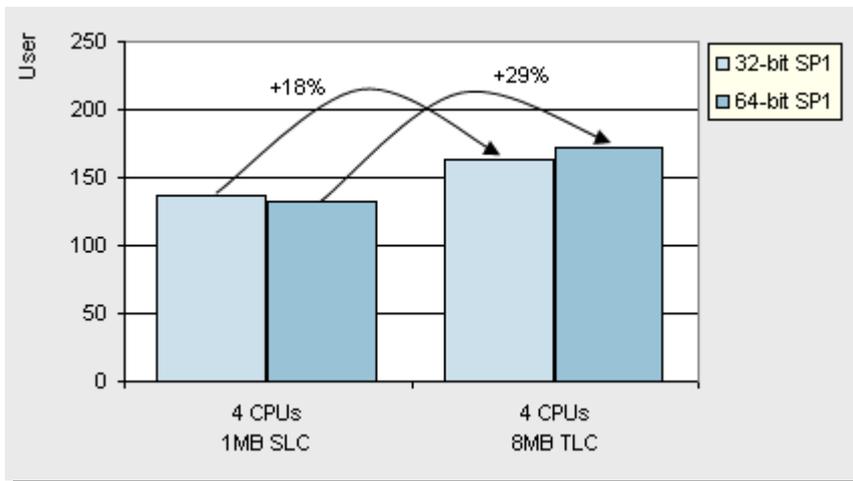
Gerade auch beim Einsatz eines 64-bit Betriebssystems ist auf einen ausreichend großen Cache zu achten, wie im Folgenden diskutiert wird.

Da beim 64-bit Betriebssystem mehr Daten verarbeitet werden müssen, hat die Größe des CPU-Caches hier einen größeren Einfluss.

Das gleiche PRIMERGY System wurde wahlweise mit Xeon Prozessoren mit 1 MB oder mit 2 MB SLC ausgestattet. Wie nebenstehende Grafik zeigt, führt eine Verdoppelung des Caches auf beiden Plattformen zu einer höheren Performance. Das 64-bit System profitiert am meisten von dem doppelt so großen Cache mit einem Performancegewinn von 25%. Das 32-bit System gewinnt bis zu 15% mehr Leistung mit dem größeren Cache. Diese Messungen zeigen, dass ein 64-bit System durch den Adress-Overhead einen größeren Cache benötigt.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Bei einem voll ausgebauten 4-Socket-System mit vier Prozessoren mit großem Cache und ausreichend Arbeitsspeicher kann sich dieser Effekt jedoch umkehren, so dass die insgesamt besten Leistungswerte mit dem 64-bit Betriebssystem erzielt werden. Grund hierfür ist, dass das 32-bit Betriebssystem bei der hier vorhandenen Rechenleistung und erreichten hohen Benutzeranzahl in einen Engpass gerät, der durch die 64-bit Architektur überwunden wird.

Eine generelle Empfehlung leitet sich aus diesen Ergebnissen ab: Für Systeme mit 64-bit Betriebssystem sollten in jedem Fall Prozessorvarianten mit einem großen Cache eingesetzt werden. Dies ist wichtiger als eine geringfügig höhere Taktfrequenz.

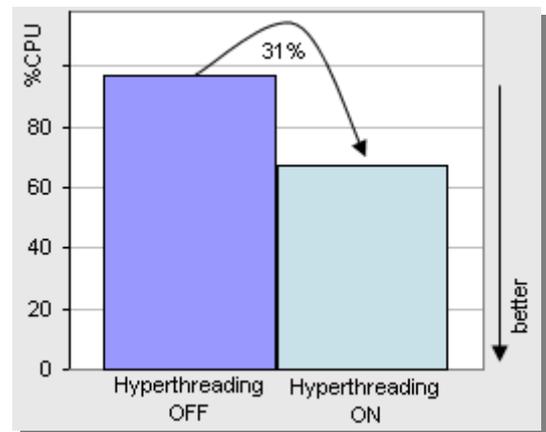
Hyper-Threading

Die meisten Intel Prozessoren der aktuellen Generationen unterstützen Hyper-Threading. Bei den heute angebotenen Prozessortypen wird eine Prozessorfamilie meist durchgängig entweder mit oder ohne Hyper-Threading angeboten. Der Vorteil von Hyper-Threading gegenüber dem klassischen Multiprocessing war ein Mehr an Leistung bei geringeren Kosten. Die neuesten Strom sparenden Intel Xeon Prozessoren für 2-Socket Systeme besitzen keine Hyper-Threading Funktionalität mehr, sondern bieten durch neueste Technologie und Dual-Core Architektur auch ohne dieses Feature hervorragende Rechenleistung. AMD Prozessoren besitzen grundsätzlich kein Hyper-Threading.

Bei Hyper-Threading-Prozessoren sind einige Ressourcen auf dem Chip verdoppelt, so dass diese CPUs nun die Fähigkeit besitzen, zwei Threads parallel ausführen zu können. So werden zwei virtuelle bzw. logische CPUs simuliert. Dem Betriebssystem gegenüber stellt sich eine CPU mit Hyper-Threading als zwei CPUs dar und wird auch so angesteuert. Dies bringt einen Geschwindigkeitsvorteil, wenn Betriebssystem und Anwendungen dafür geeignet sind. Windows als Betriebssystem ist vom Design her Hyper-Threading-fähig und gerade in Terminal Server Umgebungen arbeiten viele einzelne Benutzer mit insgesamt vielen, meist kleineren Anwendungen parallel, so dass von Hyper-Threading eine Performancesteigerung zu erwarten ist.

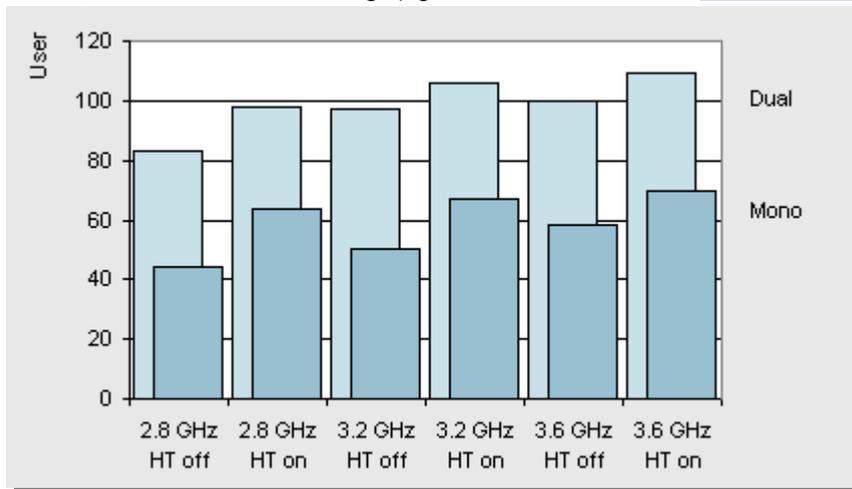
Messungen haben gezeigt, dass der Performancegewinn durch Hyper-Threading bei gering bis mittel belasteten Systemen am größten ist. Bei Systemen, die an ihrer Lastgrenze arbeiten, ist der Performancegewinn geringer. Auch ist der Performancegewinn auf einem Monoprozessorsystem höher als auf Multiprozessorsystemen.

Eine Messreihe wurde auf einer PRIMERGY RX200 mit zwei Prozessoren durchgeführt, einmal mit eingeschaltetem Hyper-Threading und einmal mit abgeschaltetem Hyper-Threading. In beiden Fällen wurde die gleiche Benutzeranzahl von 101 betrieben. Bei Terminal Server-Anwendungen sorgt das Hyper-Threading für eine Entlastung des Systems, die CPU-Last konnte um 31% reduziert werden, wie die nebenstehende Grafik veranschaulicht. Bei der verwendeten PRIMERGY RX200 und 101 Benutzern konnte man ohne Hyper-Threading schon einen CPU Engpass feststellen, die Antwortzeiten des Terminal Servers waren nicht mehr im vorgegebenen Rahmen. Durch die CPU Entlastung durch Hyper-Threading erfüllte das System wieder die vorgegebene Reaktionszeit.



(Medium Lastprofil, Microsoft Office XP, Citrix MetaFrame)

Bei einem Medium Lastprofil können also durch Hyper-Threading ca. 20% mehr Terminal Server Benutzer auf dem gleichen PRIMERGY System betrieben werden. Eine Reduzierung der CPU Belastung kann nicht 1:1 in eine Benutzersteigerung umgerechnet werden, da die CPU-Zeit im Hochlastbereich nicht mehr linear skaliert, sondern stärker ansteigt (vgl. Grafik im Abschnitt »[Verhalten bei hoher CPU-Last](#)«).



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Mit der Messmethode »Messung mit variabler Benutzeranzahl« wurden die absoluten Benutzeranzahlen für Systeme mit und ohne Hyper-Threading ermittelt. Die nebenstehende Grafik zeigt die Messergebnisse. Das vermessene PRIMERGY System war eine PRIMERGY RX300 S2 mit einem oder zwei Prozessoren mit verschiedenen Taktfrequenzen. Für diesen Vergleich wurden Prozessoren mit 1 MB SLC verwendet. Hyper-Threading war alternativ ein- (»HT on«) oder ausgeschaltet (»HT off«). Man erkennt deutlich,

dass mit eingeschaltetem Hyper-Threading eine höhere Anzahl Benutzer mit dem Terminal Server arbeiten können. Bei einem langsameren Monoprozessorsystem ist der Performancegewinn durch Hyper-Threading höher als bei einem schnellen Dualprozessorsystem.

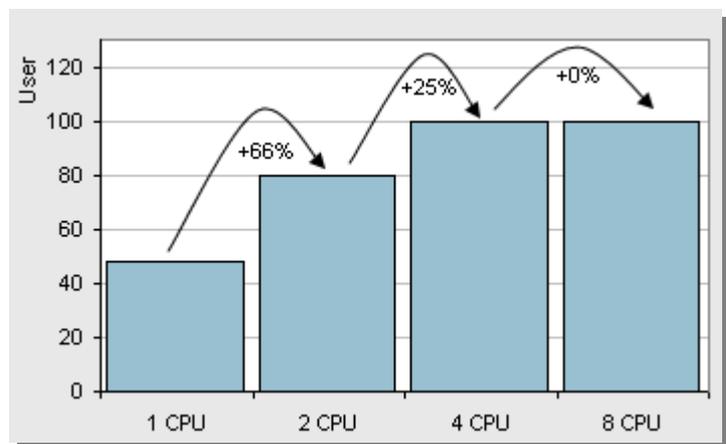
Anzahl Prozessoren

Darüber hinaus lässt sich die Rechenleistung eines Systems durch den Einsatz mehrerer Prozessoren erhöhen. Man bezeichnet den Einsatz mehrerer gleicher physikalischer Prozessoren auch als »symmetric multiprocessing« (SMP). Alle Ressourcen einer CPU sind mehrfach, d.h. auf jeder physikalischen CPU vorhanden, so dass mehrere Prozesse gleichzeitig ausgeführt werden können. Aktuelle Prozessoren der Dual-Core-Technologie besitzen zwei vollwertige Prozessorkerne auf einem Chip. Im Gegensatz dazu tritt beim [Hyper-Threading](#), s.o., eine physikalische CPU wie zwei logische CPUs auf.

Die Skalierung mit wachsender Anzahl Prozessoren ist nur im Idealfall einer optimal parallelisierbaren Anwendung linear. Je mehr Zugriffe jedoch auf gemeinsame Ressourcen wie Arbeitsspeicher, Festplatten oder Netzwerk erfolgen, und somit eine Koordination zwischen den Prozessoren bedingen, umso mehr flacht die Skalierungskurve ab. Im Extremfall kann es bei sehr großen Anzahl Prozessoren und sehr hohem Koordinations-Anteil der Prozessoren untereinander sogar zu einem »Umkippen« der Skalierung kommen. (vgl. auch »[Scale-Up](#)«). Designer von großen Multiprozessorssystemen wirken dem entgegen, indem sie den Prozessoren große Caches beiseite stellen oder Gruppen von Prozessoren bilden und diesen eigenen Arbeitsspeicher und I/O-Komponenten zuordnen. Letzteres bedingt für optimale Performance jedoch speziell angepasste Betriebssysteme und Anwendungen, wie z.B. Windows Server 2003 Enterprise und Datacenter Edition mit »nonuniform memory access« (NUMA) Support.

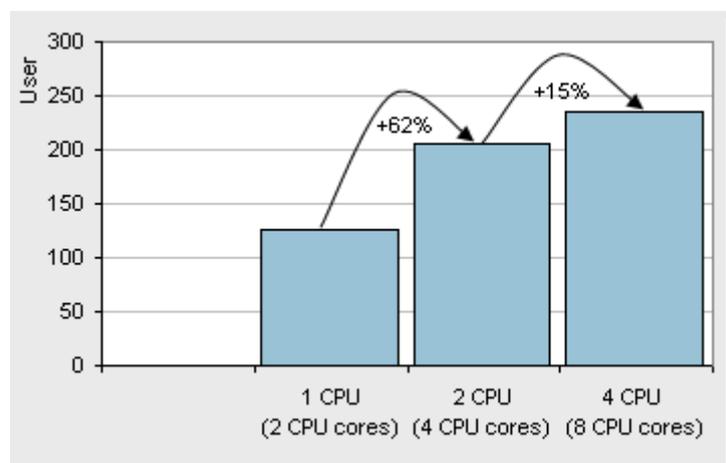
Die Skalierung über die CPU-Anzahl bei Terminal Server wurde an zwei Beispielen untersucht.

Die nebenstehende Grafik zeigt die Skalierung der PRIMERGY RX800 mit einem bis acht Prozessoren Pentium Xeon MP bei gleicher Taktfrequenz von 2.2 GHz. Hyper-Threading war eingeschaltet. Wie beim Dualprozessorsystem PRIMERGY RX300 S2 oben sieht man auch beim dem PRIMERGY RX800 System die Leistungssteigerung durch Hinzufügen von Prozessoren: die Performancesteigerung von einem Prozessor auf zwei Prozessoren beträgt 66% und von zwei auf vier Prozessoren 25%. Natürlich kann sich die Anzahl Benutzer bei einer Verdoppelung der CPU-Anzahl nicht um 100% erhöhen! Durch weitere Prozessoren kommt seitens des Betriebssystems ein höherer Synchronisationsaufwand hinzu; die Prozesse (genauer gesagt, die Threads) müssen verteilt und die Zugriffe auf die anderen Ressourcen, die ja unverändert geblieben sind, koordiniert werden. Bei der Skalierung von vier auf acht Prozessoren sieht man, dass sich »[Amdahl's Gesetz](#)« in der Praxis leider bestätigt. Für Terminal Server Anwendungen bringt eine nochmalige Verdoppelung der Prozessoranzahl keine weitere Erhöhung der Anzahl Benutzer. Der Terminal Server kann keine weiteren Sitzungen mehr verwalten, obwohl noch CPU Leistungsreserven vorhanden sein müssten.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services, Windows Server 2003 Enterprise Edition, 32-bit)

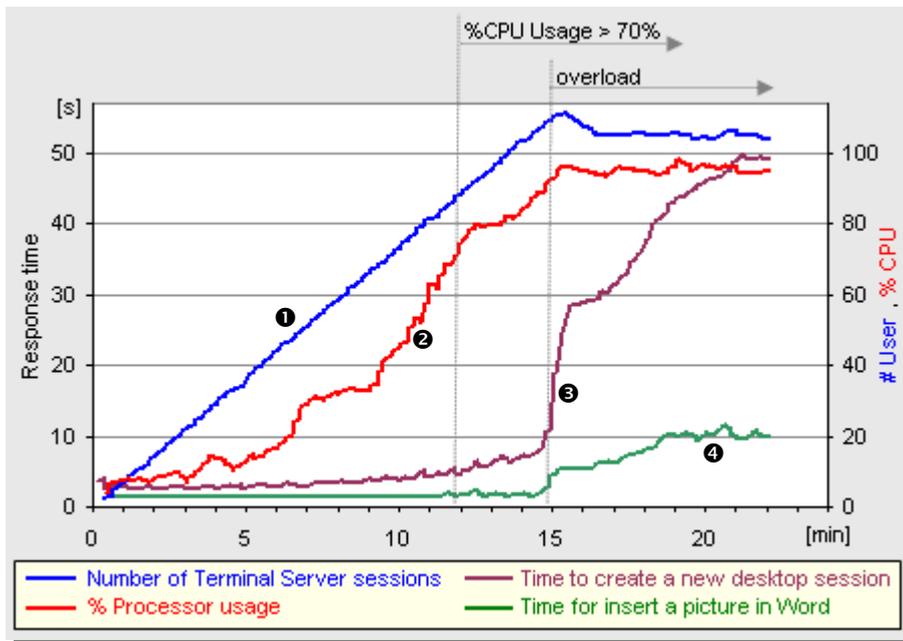
Bei den Dual-Core Prozessoren der aktuellen Generation, die auch Hyper-Threading unterstützen, laufen bereits acht physikalische, d.h. sechzehn logische, Prozessoren parallel. Bei dieser Konfiguration kommt man mit der Benutzeranzahl in Größenordnungen, die ein 32-bit Betriebssystem nicht mehr optimal verwalten kann. Auf einer PRIMERGY RX600 S3 mit einem 64-bit Windows wurde die Skalierung über die Anzahl Prozessoren ebenfalls untersucht. Bei Dual-Core Prozessoren werden mit jedem weiteren Prozessor zwei Prozessorkerne hinzugefügt, d.h. die Skalierung ist mit 62% noch gut, und auch eine weitere Verdoppelung der Anzahl Prozessorkerne bringt noch einmal eine Steigerung der Benutzeranzahl.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services, Windows Server 2003 Enterprise Edition, 64-bit)

Verhalten bei hoher CPU-Last

Setzt man die Anzahl Benutzer, die CPU-Auslastung des Systems und die Antwortzeiten miteinander in Beziehung, so erkennt man, wie sich die Antwortzeiten des Terminal Servers bei steigender Auslastung verhalten. Während einer Messung über 25 Minuten wurde in den ersten 15 Minuten die Benutzeranzahl kontinuierlich erhöht. Jeder Benutzer meldet sich erst an, danach arbeitet er mit Microsoft Word, um sich nach ca. 16 Minuten wieder abzumelden und wieder von vorn zu beginnen.

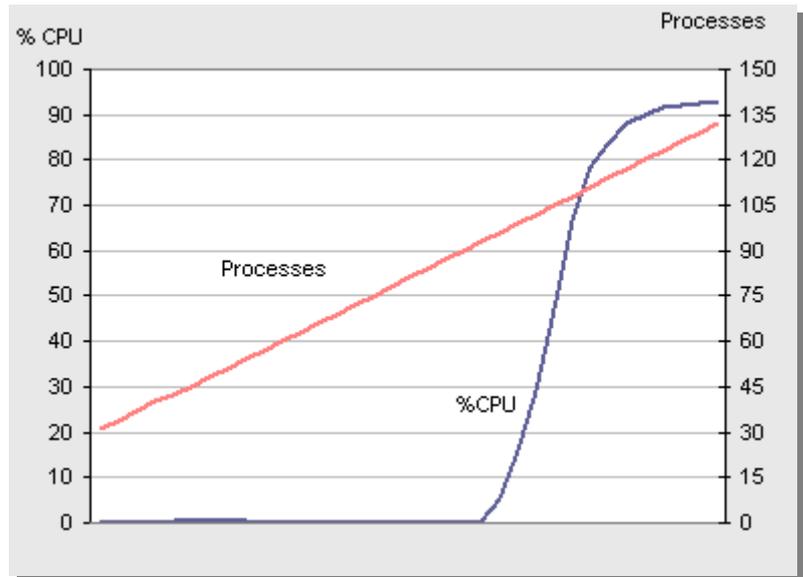


Dadurch verteilen sich die Benutzeranmeldungen (blaue Kurve ❶) auf einen relativ kurzen Zeitraum, was aber der Realität nahe kommt, wenn die Benutzer etwa zur gleichen Zeit ihre Arbeit aufnehmen. Die CPU Auslastung des Terminal Servers stieg ständig an (rote Kurve ❷) bis nah an 100%. Zwei signifikante Messergebnisse wurden beobachtet: Die Zeit, die der Benutzer braucht, um sich am Terminal Server anzumelden (violette Kurve ❸) und die Zeit, um in Microsoft Word ein Bild einzufügen (grüne Kurve ❹). Eine Anmeldung an den Terminal Server beinhaltet nicht nur das Login selbst, sondern auch der Desktop des Benutzers wird gestartet. Dies belastet den Terminal Server mehr als das Einfügen des Bildes in Microsoft Word. Aktionen, die von sich aus den Terminal Server stark belasten, werden unter Hochlast stärker verlangsamt als weniger belastende Aktionen.

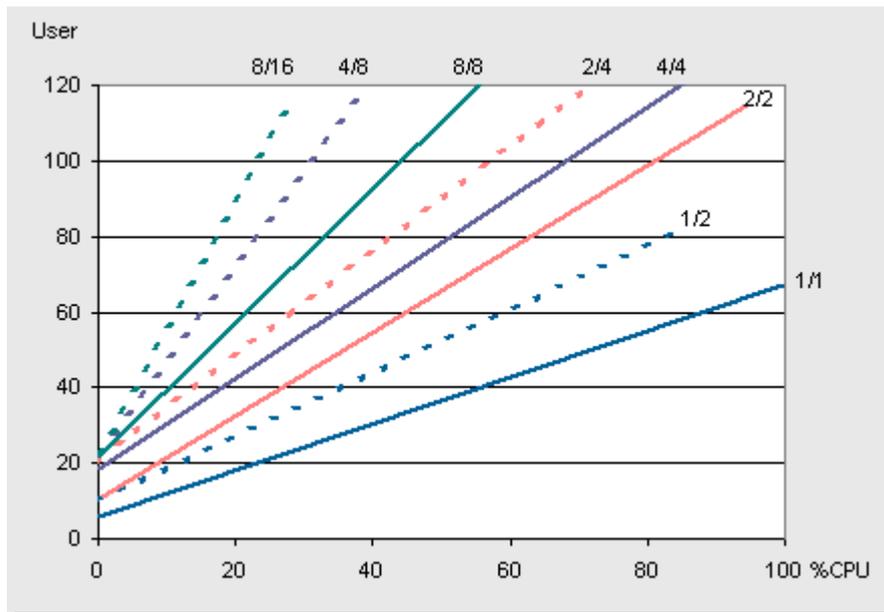
Man erkennt drei Phasen der Terminal Server Auslastung:

- Die CPU-Belastung des Terminal Servers ist unter 70%. Die Antwortzeiten des Servers verlängern sich geringfügig, dies wird der Benutzer aber nicht realisieren.
- Die CPU-Belastung des Terminal Servers ist zwischen 70% und 90%. Aktionen, die den Server stärker belasten, sind verlangsamt, aber die Antwortzeit ist für den Benutzer noch tolerierbar. Aktionen, die den Server nicht so stark belasten, haben im Durchschnitt die gleichen Reaktionszeiten, jedoch können Schwankungen auftreten.
- Wenn die CPU-Belastung des Terminal Servers über 90% ist, ist der Server deutlich überlastet. Je nach Benutzeraktion antwortet der Server wesentlich später, speziell Aktionen wie ein Login brauchen deutlich mehr Zeit. Aber auch einfachere Aktionen sind deutlich verlangsamt. Der Benutzer wird das Antwortzeitverhalten des Servers nicht mehr tolerieren.

Bei allen PRIMERGY Servern, insbesondere aber bei den leistungsfähigeren, kann man ein Verhalten beobachten, bei dem ein scheinbar normal ausgelastetes System durch das Hinzufügen einiger weniger Benutzer überlastet wird. Auf einer PRIMERGY RX600 mit vier Prozessoren, 4 GB Arbeitsspeicher und eingeschaltetem Hyper-Threading wurde untersucht, ob dieses Verhalten nur mit Terminal Server auftritt, oder ob diese Limitierung schon durch das Betriebssystem gegeben ist. Hierzu wurde eine synthetische kleine Applikation verwendet, die abwechselnd das System belastet und wartet, so wie es auch eine Anwendung auf einem Terminal Server System tun würde. Von dieser Applikation wurden kontinuierlich neue Instanzen gestartet. Wie nebenstehende Grafik zeigt, wird das System durch den kontinuierlichen Start der Anwendungen nicht belastet. Doch es gibt einen Zeitpunkt, bei dem durch das Hinzufügen weniger Prozesse die CPU-Lastung sprunghaft bis auf über 90% ansteigt. Bei den Untersuchungen, die mit Terminal Server durchgeführt wurden, ist diese Grenze nicht so deutlich wie in diesem synthetischen Fall, jedoch gibt es auch diesen Punkt, bei dem das »Fass zum Überlaufen« gebracht wird.



Die T4US-Messung mit variabler Benutzeranzahl beendet eine Messung in dem Moment, in dem die Antwortzeiten des Terminal Servers nicht mehr ausreichen, unabhängig davon, wie hoch die CPU Auslastung des Servers zu dem Zeitpunkt ist. Daher ist es interessant, sich die Prozessorauslastung in diesem Moment anzuschauen und diese Ergebnisse zwischen den PRIMERGY Systemen mit unterschiedlichen CPU-Ausbauten in Beziehung zu setzen. Nebenstehende Grafik zeigt auf der Y-Achse die Benutzeranzahl, die das entsprechende PRIMERGY System bedient, während auf der X-Achse die prozentuale CPU-Auslastung aufgetragen ist, die für den Betrieb dieser Benutzeranzahl notwendig ist. Die verschiedenen Geraden zeigen die verschiedenen CPU-Ausbaustufen der Systeme: bei einer durchgezogenen Geraden handelt es sich um ein System ohne Hyper-Threading, eine gestrichelte Gerade zeigt ein System mit eingeschaltetem Hyper-Threading. Die Grafik zeigt nur den Trend der CPU-Auslastung. In der Realität schwankt die tatsächliche Prozessorauslastung um diese Gerade. Die Zahlen bezeichnen die Anzahl CPUs »real/virtuell«, beispielsweise bedeutet »2/4« ein System mit 2 CPUs und eingeschaltetem Hyper-Threading. Man erkennt deutlich, dass, bei gleicher Benutzeranzahl, die Prozessorauslastung bei einem System ohne Hyper-Threading entsprechend höher ist als bei einem System mit Hyper-Threading. Weiterhin ist die CPU-Auslastung umso geringer, je größer ein System ist. Die Benutzerzahl, die ein solches System mit den erwarteten Antwortzeiten noch bedienen kann, wird also nicht immer bei der gleichen CPU-Auslastung erreicht. Während ein Monoprozessorsystem erst bei 90% bis 100% CPU-Last die Antwortzeiten nicht mehr erreichen kann, kann ein Dualprozessorsystem vielleicht schon bei 70% CPU-Last schon keine weiteren Benutzer mehr bedienen. Trotz dieser Unterschiede kann ein größeres System natürlich in Summe mehr Benutzer bedienen, jedoch können seine CPU-Ressourcen nicht so optimal ausgenutzt werden, da andere Komponenten im System begrenzend wirken.



Die T4US-Messung mit variabler Benutzeranzahl beendet eine Messung in dem Moment, in dem die Antwortzeiten des Terminal Servers nicht mehr ausreichen, unabhängig davon, wie hoch die CPU Auslastung des Servers zu dem Zeitpunkt ist. Daher ist es interessant, sich die Prozessorauslastung in diesem Moment anzuschauen und diese Ergebnisse zwischen den PRIMERGY Systemen mit unterschiedlichen CPU-Ausbauten in Beziehung zu setzen. Nebenstehende Grafik zeigt auf der Y-Achse die Benutzeranzahl, die das entsprechende PRIMERGY System bedient, während auf der X-Achse die prozentuale CPU-Auslastung aufgetragen ist, die für den Betrieb dieser Benutzeranzahl notwendig ist. Die verschiedenen Geraden zeigen die verschiedenen CPU-Ausbaustufen der Systeme: bei einer durchgezogenen Geraden handelt es sich um ein System ohne Hyper-Threading, eine gestrichelte Gerade zeigt ein System mit eingeschaltetem Hyper-Threading. Die Grafik zeigt nur den Trend der CPU-Auslastung. In der Realität schwankt die tatsächliche Prozessorauslastung um diese Gerade. Die Zahlen bezeichnen die Anzahl CPUs »real/virtuell«, beispielsweise bedeutet »2/4« ein System mit 2 CPUs und eingeschaltetem Hyper-Threading. Man erkennt deutlich, dass, bei gleicher Benutzeranzahl, die Prozessorauslastung bei einem System ohne Hyper-Threading entsprechend höher ist als bei einem System mit Hyper-Threading. Weiterhin ist die CPU-Auslastung umso geringer, je größer ein System ist. Die Benutzerzahl, die ein solches System mit den erwarteten Antwortzeiten noch bedienen kann, wird also nicht immer bei der gleichen CPU-Auslastung erreicht. Während ein Monoprozessorsystem erst bei 90% bis 100% CPU-Last die Antwortzeiten nicht mehr erreichen kann, kann ein Dualprozessorsystem vielleicht schon bei 70% CPU-Last schon keine weiteren Benutzer mehr bedienen. Trotz dieser Unterschiede kann ein größeres System natürlich in Summe mehr Benutzer bedienen, jedoch können seine CPU-Ressourcen nicht so optimal ausgenutzt werden, da andere Komponenten im System begrenzend wirken.

Die verschiedenen Geraden zeigen die verschiedenen CPU-Ausbaustufen der Systeme: bei einer durchgezogenen Geraden handelt es sich um ein System ohne Hyper-Threading, eine gestrichelte Gerade zeigt ein System mit eingeschaltetem Hyper-Threading. Die Grafik zeigt nur den Trend der CPU-Auslastung. In der Realität schwankt die tatsächliche Prozessorauslastung um diese Gerade. Die Zahlen bezeichnen die Anzahl CPUs »real/virtuell«, beispielsweise bedeutet »2/4« ein System mit 2 CPUs und eingeschaltetem Hyper-Threading. Man erkennt deutlich, dass, bei gleicher Benutzeranzahl, die Prozessorauslastung bei einem System ohne Hyper-Threading entsprechend höher ist als bei einem System mit Hyper-Threading. Weiterhin ist die CPU-Auslastung umso geringer, je größer ein System ist. Die Benutzerzahl, die ein solches System mit den erwarteten Antwortzeiten noch bedienen kann, wird also nicht immer bei der gleichen CPU-Auslastung erreicht. Während ein Monoprozessorsystem erst bei 90% bis 100% CPU-Last die Antwortzeiten nicht mehr erreichen kann, kann ein Dualprozessorsystem vielleicht schon bei 70% CPU-Last schon keine weiteren Benutzer mehr bedienen. Trotz dieser Unterschiede kann ein größeres System natürlich in Summe mehr Benutzer bedienen, jedoch können seine CPU-Ressourcen nicht so optimal ausgenutzt werden, da andere Komponenten im System begrenzend wirken.

Arbeitsspeicher

Den stärksten Einfluss auf die Leistungsfähigkeit des Terminal Servers übt der Arbeitsspeicher aus. Dabei spiegelt sich dies insbesondere in der Antwortzeit wider. Denn Windows verschafft sich bei Bedarf weiteren virtuellen Speicher durch Auslagern (Pagen) von momentan nicht benötigten Daten aus dem Arbeitsspeicher (RAM) in die Auslagerungsdatei (Pagefile) auf Festplatte. Da Plattenzugriffe aber mindestens um die Größenordnung 1000 langsamer sind als Speicherzugriffe, führt dies unmittelbar zum Zusammenbruch der Leistung und zu einem rapiden Anstieg der Antwortzeiten.

In der Praxis setzt heute nicht die Hardware die Grenzen, sondern die Software-Architektur. Die heute zumeist eingesetzte Software im 32-bit Design kann die zur Verfügung stehende Hardware häufig nicht mehr voll nutzen. Im speziellen ergibt sich eine Limitierung bei der Adressierung des Arbeitsspeichers, bei der 32-bit Anwendungen auf 4 GB virtuellen Adressraum begrenzt sind. Ist der Server physikalisch mit mehr als 4 GB Arbeitsspeicher ausgestattet, so kann dieser Speicher zumeist nicht effektiv genutzt werden. Durch die Abhängigkeit zwischen dem Bedarf an Arbeitsspeicher und Rechenleistung können viele Anwendungen auch die Rechenleistung, die 8-way und 16-way Server bereitstellen, nicht ausschöpfen.

Die Prozesse der Benutzer belegen virtuellen Speicherplatz, unabhängig davon, ob die Benutzer gerade aktiv arbeiten oder nicht und sogar auch, wenn sich die Sitzung gerade im Status »disconnected« befindet. Da diese Grenzen den virtuellen Speicherplatz betreffen, würde man durch Hinzufügen von Hauptspeicher keine Verbesserung erreichen. Microsoft hat bereits in der 32-bit Version von Windows Server 2003 Verbesserungen hinsichtlich der Speicherverwaltung durchgeführt: der begrenzte virtuelle 32-bit Adressbereich wird vom Kernel optimaler ausgenutzt, so dass mit Windows Server 2003 mehr Benutzer arbeiten können als mit Windows 2000 Server. Dies gilt natürlich nur für Systeme, bei denen die virtuelle Speicherverwaltung der Engpass ist. Auf Systemen, bei denen die CPU oder der physikalische Speicher der Engpass ist, können bei Windows 2000 Server und bei Windows Server 2003 die gleiche Anzahl von Benutzern betrieben werden.

Erst mit 64-bit Betriebssystemen und 64-bit Anwendungen wird hier Abhilfe geschaffen. Zwar gibt es heute 64-bit Versionen von Windows Server 2003 und 64-bit PRIMERGY Systeme, jedoch mangelt es an 64-bit Anwendungen. Aber auch 32-bit Anwendungen können von einem 64-bit Betriebssystem profitieren.

Folgende Tabelle gibt einen Überblick über die Speicher-relevanten Unterschiede der verschiedenen 32-bit und 64-bit Windows Produkte, die für Terminal Server sinnvoll eingesetzt werden können:

Windows Produkt	Windows Server 2003 R2		Windows Server 2003 R2 x64	
	Standard Edition	Enterprise Edition	Standard Edition	Enterprise Edition
Anzahl Prozessoren	1 – 4	1 – 8	1 – 4	1 – 8
Max. RAM	4 GB	64 GB (mit PAE)	32 GB	1 TB
gesamter virtueller Adressraum	4 GB		16 TB	
Virtueller Adressraum eines 32-bit Prozesses	2 GB		2 GB	
	3 GB wenn mit /3GB Switch gebootet		4 GB wenn mit /LARGEADDRESSAWARE übersetzt wurde	
Virtueller Adressraum eines 64-bit Prozesses	-		8 TB	

Allgemeingültig ist jedoch die Tatsache, dass der Speicherbedarf linear nach der Formel

$$\text{Memory} = \text{Memory}_{OS} + \#_{\text{Client}} \cdot \text{Memory}_{\text{App}}$$

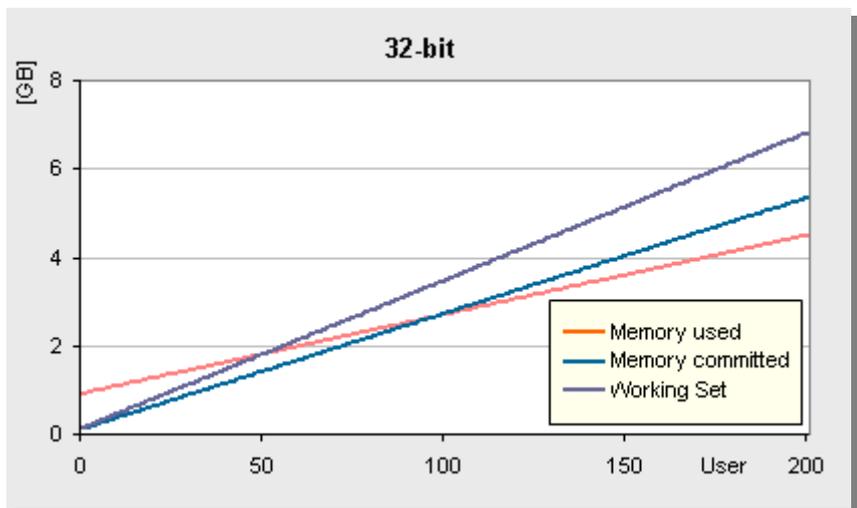
anwächst.

Man beachte dabei nur, dass der Speicherbedarf pro Benutzer stark von der verwendeten Anwendung abhängt. Kennt man jedoch den Speicherbedarf der Anwendung bei einem Benutzer, so lässt sich leicht der Gesamt-Speicherbedarf berechnen.

Nach oben begrenzend wirken hier der maximal mögliche Speicherausbau des Systems sowie die Unterstützung des Betriebssystems (siehe Tabelle oben).

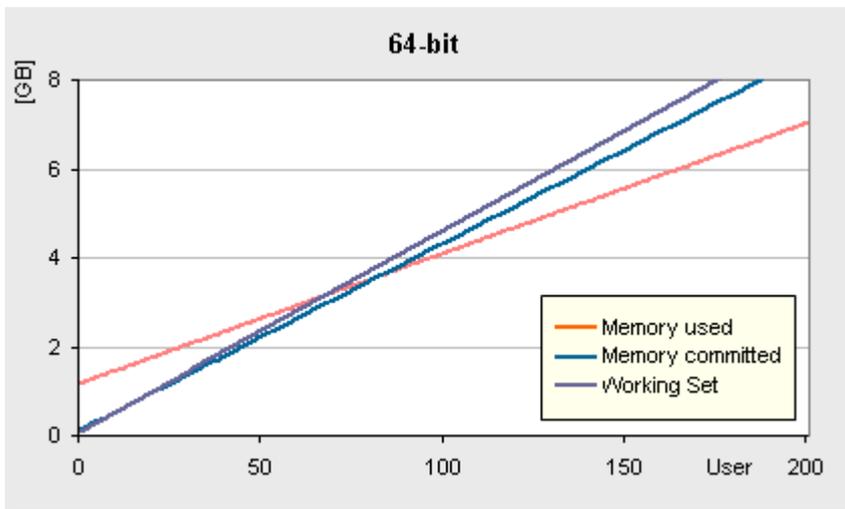
Der Speicherverbrauch einer Windows-Anwendung ist nicht leicht zu bestimmen. Verschiedene Performance Counter wie »Available MBytes«, »Committed Bytes« und der »Working Set« der Anwendungen stehen zur Verfügung. Während »Available MBytes« den aktuell freien physikalischen Hauptspeicher angibt, wird in »Committed Bytes« angegeben, wie viel virtueller Hauptspeicher den laufenden Anwendungen zugesagt wurde. Der »Working Set« einer Anwendung ist der Speicher, den sie bereits verwendet hat. Eine Eigenschaft des Windows Betriebssystems ist dabei zu beachten: wenn genügend Hauptspeicher zur Verfügung steht, wird auch mehr Speicher belegt. Erst wenn der freie Speicher unter einen bestimmten Schwellwert fällt, wird »aufgeräumt« und der Speicher außerhalb des »Working Set« ausgelagert. So kann man an einem System, das nicht im Speichergrenzbereich betrieben wird, den wirklichen Speicherbedarf nicht ablesen.

Trägt man dennoch den aus »Available MBytes« berechneten belegten Speicher, den »Committed« Speicher und das »Working Set« grafisch auf, so erkennt man einen linearen Verlauf, der mit steigender Benutzeranzahl wächst, unabhängig vom PRIMERGY Modell und gleichermaßen bei 32-bit und 64-bit Betriebssystemen. Der Anstieg der Geraden beim 64-bit Betriebssystem ist jedoch steiler. Theoretisch könnte man diese Gerade bis zur maximalen Ausbaustufe der hier betrachteten PRIMERGY Systeme weiterziehen.

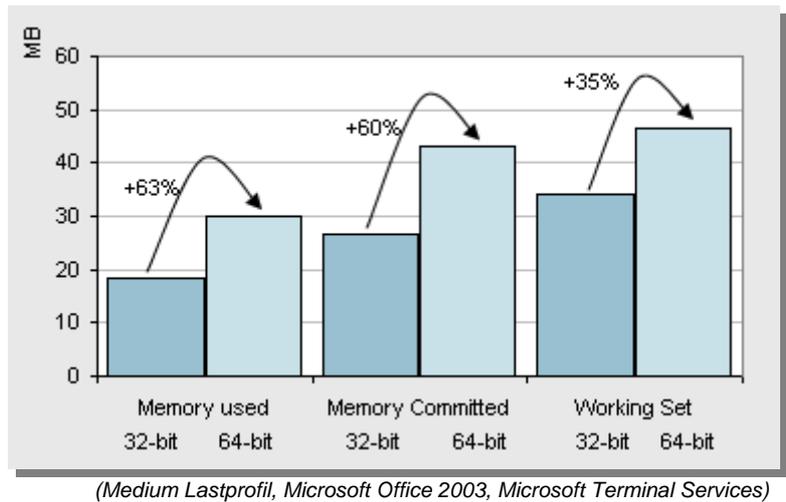


Allerdings sollte man beachten, dass die Gesamtleistung des Systems durch die schwächste Komponente bestimmt wird. Hinzu kommt, dass durch die Architektur des 32-bit Betriebssystems die internen Strukturen

und der virtuelle Adressraum eingeschränkt sind, so dass man den maximal möglichen Speicherausbau der PRIMERGYs für Terminal Server unter 32-bit nicht ausnutzen kann. Dies wird im Kapitel »[Betriebssystem](#)« im Detail diskutiert.



Wie bereits oben erwähnt, belegt ein Terminal Server Benutzer auf einem 64-bit System mehr Arbeitsspeicher als auf einem 32-bit System. Die Anwendung, mit der der Terminal Server Benutzer arbeitet, ist in beiden Fällen Microsoft Word, welches heute nur als 32-bit Version existiert. Die Microsoft Terminal Services liegen als Bestandteil des Betriebssystems in einer 64-bit Version vor. Wie die nebenstehende Grafik zeigt, belegt der gleiche Benutzer, der den Desktop gestartet hat und mit Microsoft Word 2003 arbeitet, auf dem 64-bit System ca. 60% mehr Arbeitsspeicher. Dies ist eine Folge der bereits erwähnten doppelten Breite der 64-bit Adresszeiger.



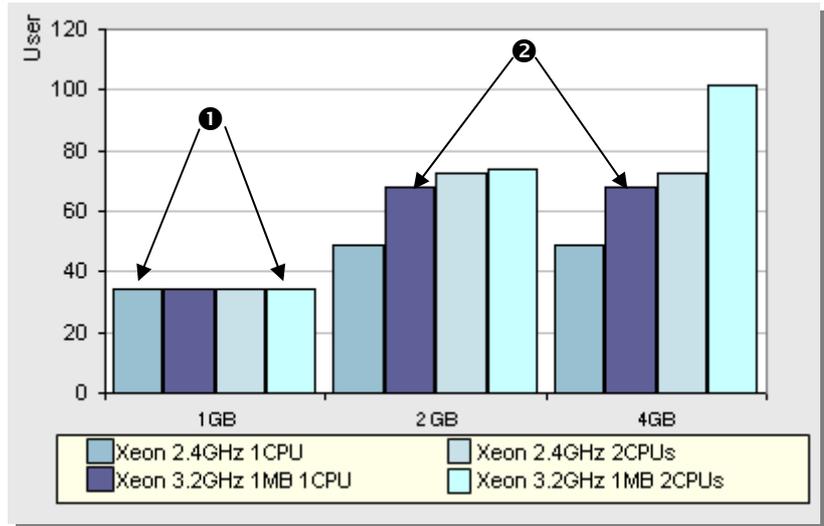
Für unsere Sizing Empfehlungen werden die hier ermittelten Werte verwendet und den Speicherbedarf eher konservativ abgeschätzt, da das System so dimensioniert werden soll, dass der Speicher nicht von vorn herein knapp ist.

Das Betriebssystem (Windows Server 2003 Enterprise Edition) hat einen Grundbedarf von 128 MB, und pro Benutzer bzw. Client werden weitere 20 MB benötigt. Der Grundbedarf des 64-bit Systems erhöht sich auf ca. 150 MB. In dem Messszenario arbeiten allerdings alle Benutzer mit der gleichen Applikation, daher zeigen alle Benutzergruppen den gleichen Speicherbedarf. Der Speicherbedarf ist jedoch von den verwendeten Applikationen abhängig und muss kundenspezifisch ermittelt werden.

Unter Citrix Presentation Server ist der Speicherverbrauch pro Benutzer etwas höher als beim Microsoft Terminal Server, siehe Kapitel »[Microsoft Terminal Server vs. Citrix Presentation Server](#)«.

Wenn die Prozessor-Leistung nicht mehr ausreicht, dann können die angeschlossenen Clients nicht mehr mit akzeptabler Antwortzeit bedient werden, auch wenn noch ausreichend Arbeitsspeicher vorhanden ist. Dies wird durch nebenstehende Grafik am Beispiel der PRIMERGY RX200 verdeutlicht.

Bei einem Speicherausbau von 1 GB kann keiner der eingesetzten Prozessoren seine Leistungsfähigkeit unter Beweis stellen, da der unzureichende Hauptspeicher der begrenzende Faktor ist ❶. Eine Verdoppelung des Hauptspeichers von 1 auf 2 GB führt bei allen vier Varianten zu einer höheren Benutzeranzahl. Verdoppelt man den Hauptspeicher von 2 auf 4 GB, so kann zum Beispiel die PRIMERGY RX200 mit einem 2.4 GHz-Prozessor nicht davon profitieren, da hier die CPU die schwächste Komponente ist ❷. Um für einen bestimmten Einsatzbereich ein optimales System zu finden, gilt es also immer, ein ausgewogenes Verhältnis von CPU und Hauptspeicher herzustellen.

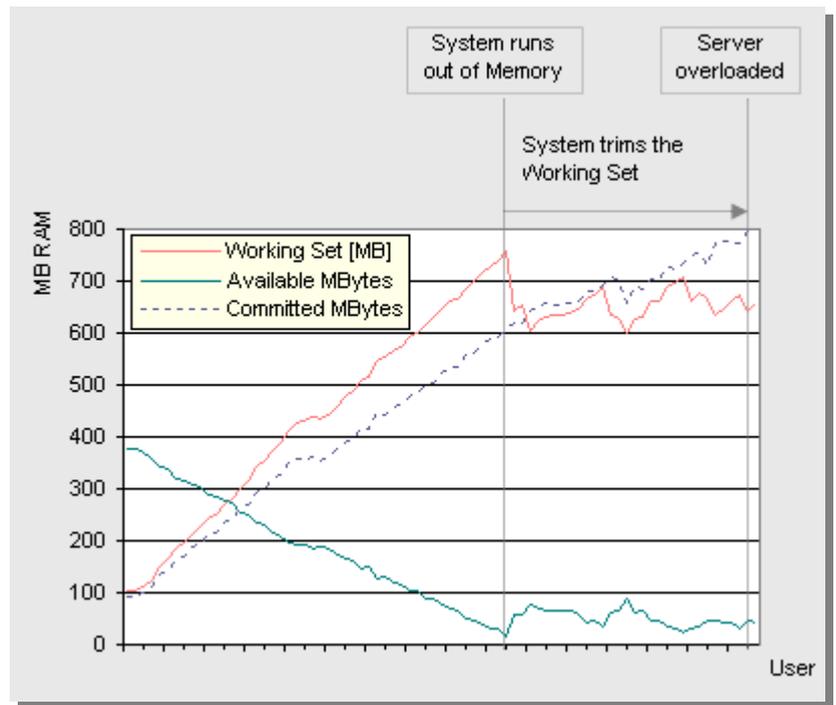


(Medium Lastprofil, Microsoft Office XP, Citrix MetaFrame)

Zu beachten ist, dass die PRIMERGY Modelle mit AMD Opteron Prozessoren die dafür typische direkte Zuordnung des Speichers zu den Prozessoren besitzen. Jeder Prozessor hat »seinen« Speicher direkt angebunden, was in einem schnellen Zugriff resultiert. Speicher, der anderen Prozessoren zugeordnet ist, wird über den Memory Controller der anderen CPU angesprochen, hierdurch ist der Zugriff langsamer. Um die Vorteile dieser Architektur optimal nutzen zu können, ist ein Betriebssystem mit »nonuniform memory access« (NUMA) Unterstützung notwendig, das die einzelnen Gruppen aus Prozessor und Speicher erkennt und wenn möglich den Zugriff eines Prozesses auf den schnell angebundenen Speicher durchführt. Bei dem hier verwendeten Betriebssystem »Windows Server 2003« mit Service Pack 1 in der 32-bit und 64-bit Version wird die Systemarchitektur erkannt und NUMA automatisch verwendet. Aus dieser Systemarchitektur folgt umgekehrt, dass eine PRIMERGY, die lediglich mit einem AMD Opteron Prozessor bestückt ist, auch nur die Hälfte der Speicherbänke nutzen kann, da der zweite CPU-interne Memory Controller zur Ansteuerung fehlt.

Wie bereits erwähnt, hilft sich ein System, das mit zu wenig Hauptspeicher ausgestattet ist, indem es nicht benötigte Daten auf die Festplatte auslagert.

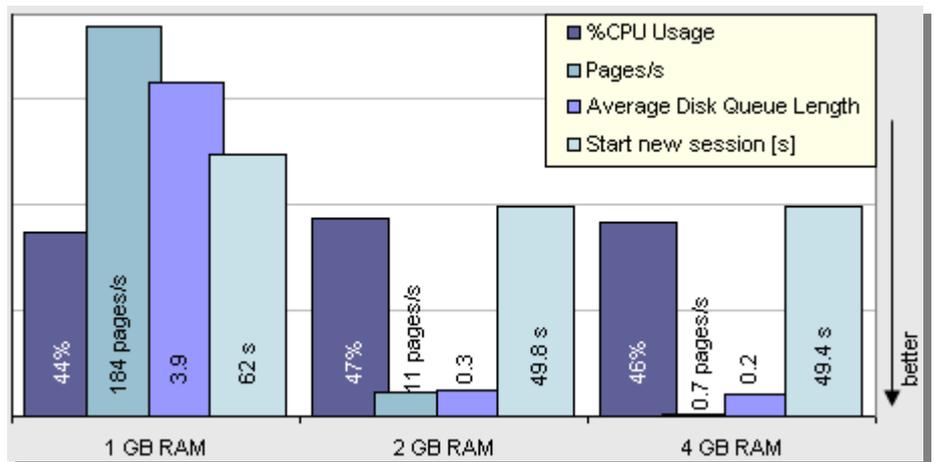
Das Windows Betriebssystem und die laufenden Anwendungen belegen mehr Speicher, wenn noch genügend Hauptspeicher zur Verfügung steht. Erst wenn der freie Speicher unter einen bestimmten Schwellwert fällt, wird »aufgeräumt«. Dies lässt sich durch eine Messung mit einem Microsoft Terminal Server System, das mit nur 512 MB Hauptspeicher ausgestattet war, veranschaulichen. Im nebenstehenden Bild ist der Zeitpunkt, an dem das System einen drohenden Speicherengpass erkennt, an der ersten senkrechten Linie zu erkennen. Der »Working Set« wird verkleinert, dadurch steigt der verfügbare Speicher an. Durch diese Maßnahme des Betriebssystems können auf dem Terminal Server mehr Benutzer arbeiten. Dieses hier beschriebene Auslagern von nicht benötigten Speicherbereichen ist durch die Memory Managing Mechanismen des Betriebssystems bedingt und fällt noch nicht unter das gefürchtete »Paging«. Beim Pagen werden aufgrund des Speicherengpasses Speicherbereiche auf die Festplatte geschrieben, die jedoch noch in Benutzung sind und daher beim nächsten Zugriff wieder von der Platte in den Arbeitsspeicher geladen werden müssen.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Am Beispiel einer PRIMERGY RX200 mit 81 Benutzern wurde untersucht, wie sich die Größe des Hauptspeichers auf die Antwortzeiten des Terminal Servers auswirkt. In diesem Fall erwiesen sich 4 GB RAM als ausreichend, 2 GB RAM als knapp bemessen und 1 GB RAM als zu klein.

Bei nahezu gleicher Prozessorauslastung (linker Balken in jeder Vierergruppe) von 44% bis 46% erkennt man gut den Einfluss der Arbeitsspeichergöße. Bei einem zu geringen Speicherausbau von nur 1 GB, ist der Terminal Server gezwungen, Teile des Betriebssystems auf Festplatte auszulagern (Pagen). Mit 184 Zugriffen pro Sekunde auf das Pagefile (zweiter Balken von links) gegenüber 11 bei 2 GB und 0.7 bei 4 GB Speicherausbau. Die 184 Plattenzugriffe pro Sekunde liegen zudem über dem, was eine einzelne Festplatte leisten kann, wie die Disk-Queue-Länge von 3.9 (dritter Balken von links) zeigt.



(Medium Lastprofil, Microsoft Office XP, Citrix MetaFrame)

Bei einer Ausstattung mit 2 GB Arbeitsspeicher erkennt man, dass der Engpass beseitigt ist. 10 Festplattenzugriffe auf das Pagefile sind völlig normal. Die Zeit für den Start einer neuen Terminal Server Session sinkt von 62 auf 49.8 Sekunden (vierter Balken).

Eine weitere Vergrößerung des Arbeitsspeichers auf 4 GB RAM führt hingegen zu keiner nennenswerten Leistungssteigerung. Zwar sinken die Zugriffe auf das Pagefile weiter, dies wirkt sich jedoch nicht in einem schnelleren Start neuer Sessions aus.

Bei der Kalkulation des Arbeitsspeichers für Terminal Server sollte man zwei Besonderheiten nicht außer Acht lassen:

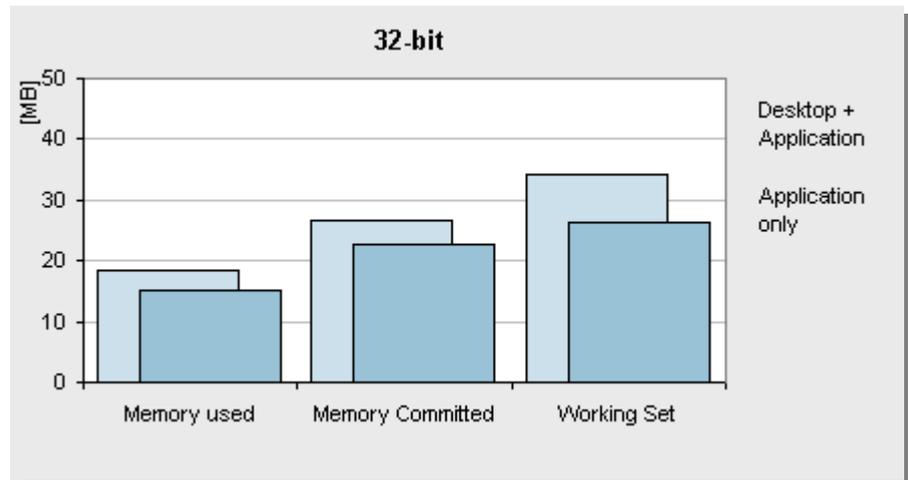
- »Desktop« oder »Published Application«?

Bei Terminal Server Umgebungen muss man dem Benutzer nicht den gesamten Desktop mit allen Anwendungen zur Verfügung stellen, man kann den Zugriff auch beschränken. Microsoft Terminal Server kann man so konfigurieren, dass eine bestimmte Anwendung statt des Desktop gestartet wird. Bei Citrix Presentation Server kann man den Benutzern auch mehrere einzelne Anwendungen (»Published Application«) direkt zur Verfügung stellen, der Start der Applikation innerhalb des Desktop entfällt. Dabei werden pro Benutzer ungefähr 5 bis 10 MB an Hauptspeicher gespart, da der Explorer-Prozess nicht mit gestartet werden muss. Auch wenn es in einer bestimmten Umgebung vielleicht nicht um diesen kleinen Gewinn von Hauptspeicher geht; ein nicht zu vernachlässigender Vorteil dieser Konfiguration ist, dass der Benutzer nur die für ihn vorgesehenen Anwendungen auf dem Server laufen lassen kann. Man kann also die Aktionen der Benutzer besser vorhersagen und auch einschränken.

Der Speicherbedarf eines Benutzers, der eine Anwendung über den Desktop startet, wurde mit dem Speicherbedarf eines Anwenders verglichen, der die gleiche Applikation direkt ohne Desktop startet. Dies wurde auf einer PRIMERGY RX300 S2 mit zwei Prozessoren, 4 GB Arbeitsspeicher und eingeschaltetem Hyper-Threading untersucht. Die Anwendung Microsoft Word aus Microsoft Office 2003 wurde einmal über den Desktop und einmal über den RDP Client direkt gestartet.

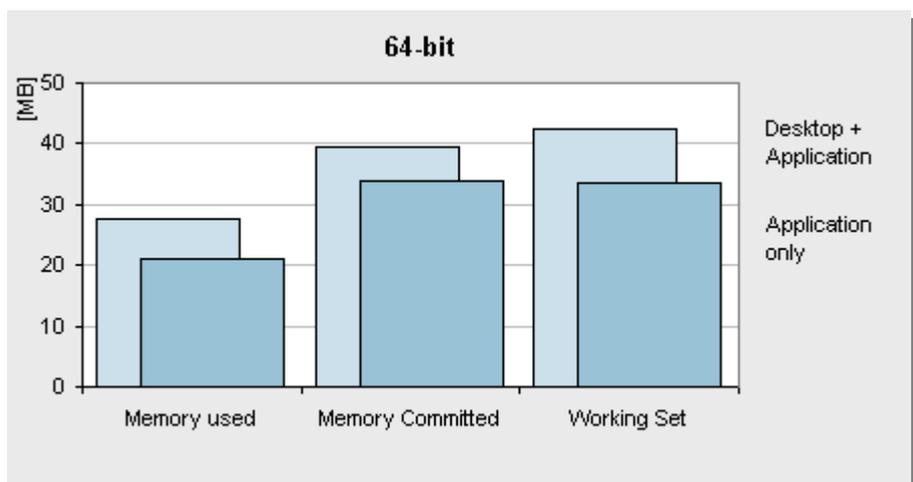
Der Speichermehrverbrauch des Desktop pro Benutzer ist deutlich sichtbar, wie die Grafiken zeigen.

Beim belegten Speicher (»Memory used«) belegt der Desktop beim 32-bit Betriebssystem ca. 3.3 MB mehr RAM, bei »Memory Committed« beträgt der Unterschied ca. 4.1 MB, während der »Working Set« ohne Desktop ca. 8 MB kleiner ist.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Auch beim 64-bit Betriebssystem erhöht sich, wie erwartet, der Speicherplatz pro Benutzer durch die Verwendung der Desktops. »Memory used« wächst um ca. 6.3 MB, »Memory Committed« um ca. 5.8 MB und der »Working Set« erhöht sich um ca. 8.8 MB. Verglichen mit dem 32-bit Betriebssystem verbraucht das 64-bit Betriebssystem generell mehr Speicher pro Benutzer, dies ist durch die Adressierung unter 64-bit bedingt.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

- »Logoff« oder »Disconnect«?

Es ist ein Unterschied, ob der Benutzer die Verbindung zum Terminal Server beendet, indem er sich abmeldet (»Logoff«) oder ob er mit einem »Disconnect« die Verbindung nur unterbricht. Im letzteren Fall läuft die Anwendung auf dem Terminal Server weiter und gibt ihre Ressourcen nicht frei. Der Benutzer kann seine Arbeit dann an dieser Stelle fortsetzen. Wenn der belegte Arbeitsspeicher des Servers analysiert wird, zählen die Verbindungen im Status »Disconnect« natürlich mit. Manche Anwendungen brauchen auch CPU-Ressourcen, während die Verbindung getrennt ist. Sowohl Microsoft als auch Citrix unterstützen »disconnected« Sessions.

Disk-Subsystem

Geht man davon aus, dass ein Server dediziert als Terminal Server eingesetzt wird, und nicht gleichzeitig noch als File-Server oder Datenbank-Server, so werden an das Disk-Subsystem keine großen Anforderungen gestellt. Es muss im Wesentlichen nur das Betriebssystem, den Paging-Bereich und die den Terminal Server-Clients zur Verfügung stehenden Applikationen beherbergen. Die Plattenzugriffe sind dabei gering. Auf das Betriebssystem und die Applikationen wird nur zugegriffen, wenn sie das erste Mal in den Speicher geladen werden. Der Paging-Bereich spielt prinzipiell auch keine Rolle, denn das System muss so konfiguriert sein, dass es nicht in starkes Pagen gerät, anderenfalls ist die Leistungsfähigkeit des Systems in jedem Fall erheblich beeinträchtigt. Weiteres hierzu lesen Sie im Abschnitt »[Arbeitsspeicher](#)«.

Je nach Anzahl Festplatten im Server-System gelten folgende Empfehlungen:

Zwei Festplatten, Konfiguration für Sicherheit: Wenn nur zwei Festplatten zur Verfügung stehen und auf Sicherheit Wert gelegt wird, so sollte aus Sicherheitsgründen eine Spiegelung über RAID 1 eingerichtet werden. Dort befinden sich dann Betriebssystem, Paging-Bereich und Applikationen. Eine Spiegelung kann entweder mit den onboard RAID 1-Funktionalitäten, die fast alle PRIMERGY Server standardmäßig bieten, oder softwaremäßig mittels Windows Server 2003 realisiert werden. Alternativ kann auch ein RAID-Controller eingesetzt werden.

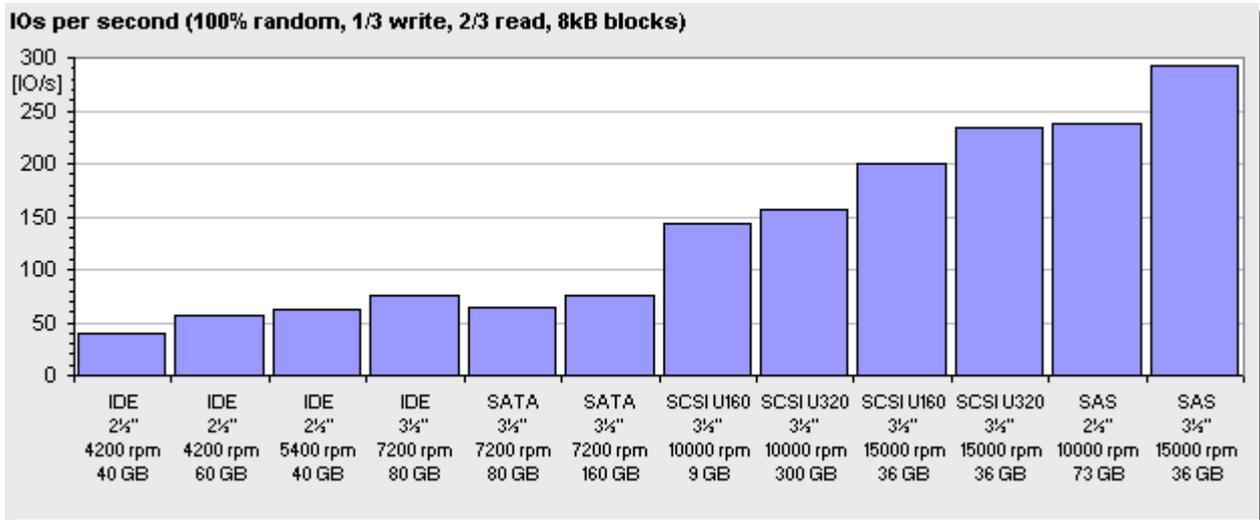
Zwei Festplatten, Konfiguration für Performance: Um eine bessere Performance zu erreichen, sollte das Pagefile nicht auf einer gespiegelten Festplatte konfiguriert werden. Wenn nur zwei Festplatten zur Verfügung stehen, sollten das Betriebssystem und die Applikationen auf der ersten Festplatte gespeichert werden, während das Pagefile auf die zweite Festplatte gelegt werden sollte. Da die Festplatte des Betriebssystems nicht gespiegelt ist, sollte sichergestellt werden, dass dort keine Benutzerdaten abgelegt und dass regelmäßige Backups durchgeführt werden.

Drei oder mehr Festplatten: Stehen insgesamt mindestens drei Festplatten zur Verfügung, sollte man zwei Festplatten mit Hilfe von RAID 1 spiegeln und dort das Betriebssystem und die Applikationen unterbringen. Der Paging-Bereich wird auf die dritte dedizierte Festplatte gelegt, da die Nutzung einer gespiegelten Festplatte durch das Pagefile beträchtliche Auswirkungen auf die Performance haben kann. Da diese Daten alle temporärer Natur sind, ist hier natürlich keine Absicherung durch RAID notwendig. Von einer Konfiguration der Systemfestplatte mit RAID 5 wird aus Performance-Gründen dringend abgeraten.

Die Benutzerdaten und Benutzerprofile werden typischerweise auf entsprechende Disk-Subsysteme oder externe File Server gelegt und nicht auf die lokale Festplatten eines Terminal Servers.

Um einen maximalen Durchsatz zu erreichen, sollten alle Caches, auch die Write-Caches, eingeschaltet werden. Write-Caches der Festplatten tragen erheblich zur Performance-Steigerung bei, und es empfiehlt sich, diese bei allen Festplatten vorhandene Funktionalität auch im produktiven Einsatz zu nutzen. Dabei ist die Verwendung einer USV zum Schutz gegen Stromausfälle und damit verbundenem Datenverlust empfehlenswert.

Prinzipiell sind alle Festplatten für den Einsatz mit Terminal Server geeignet. Selbst relativ langsame Festplatten stellen bei den Terminal Server Anwendungen keinen ernstzunehmenden Engpass da, wenn die verwendete Festplatte auf die Anzahl der mit dem Terminal Server arbeitenden Benutzer abgestimmt ist. Die folgende Abbildung zeigt, dass eine 2½" IDE Festplatte wesentlich langsamer ist als eine high-end 3½" SCSI oder SAS Platte, jedoch sind die Disk-I/O-Aktivitäten von Terminal Server auf der Festplatte, die das Betriebssystem trägt, relativ gering.



Bei unserem Medium Lastprofil erzeugen die simulierten Benutzer, die sich während der Messung kontinuierlich an- und abmelden, eine durchschnittliche IO-Last von 0.45 IOs pro Sekunde auf der Betriebssystemplatte und unter 0.1 IOs pro Sekunde auf der Datenplatte, beides unter einem 32-bit Windows gemessen. Für kleinere Konfigurationen bieten sich daher auch IDE oder SATA Platten an, während bei Konfigurationen mit mehr Benutzern SCSI oder SAS Festplatten zu empfehlen sind.

Wird ein Terminal Server auf moderner Server-Hardware mit Multi-Core-Architektur und ggf. einem 64-bit Betriebssystem aufgebaut, so können damit mehr Benutzer arbeiten und somit muss auch das Disk-Subsystem wachsen, um diesen Anforderungen gerecht werden zu können. Durch mehr Benutzersitzungen finden mehr Zugriffe auf das lokale Pagefile und auf das Disk-Subsystem statt und beim 64-bit Betriebssystem werden auch größere Datenmengen Richtung Pagefile geschrieben. Aus diesen Gründen muss das Disk-Subsystem entsprechend dimensioniert werden. Zum Sizing von Disk Subsystemen wurde ein eigenes White Paper »[Disk Subsystem Sizing – RAID Controller](#)« erstellt.

Ob das Disk-Subsystem leistungsfähig genug ist, kann man mit Hilfe des Windows Performance Monitors überprüfen. Ein signifikanter Indikator für die Belastung des Disk-Subsystems ist der Performance Counter »Avg. Disk Queue Length«. Dieser sollte nicht dauerhaft größer sein als 1 pro Festplatte, die netto zur Verfügung steht. In einem RAID 1 Verband aus zwei Festplatten also nicht größer als 1, in einem RAID 1+0 aus 6 Festplatten nicht größer als 3.

In allen unseren Messungen lag die »Avg. Disk Queue Length« für die Festplatten, die das Betriebssystem trugen, deutlich unter 0.5.

Bei der Messung war der Write-Cache jeweils eingeschaltet. Im produktiven Einsatz ist dabei die Verwendung einer USV zum Schutz gegen Stromausfälle und damit verbundenem Datenverlust empfehlenswert.

Wird das Terminal Server-System hingegen gleichzeitig noch als Datei- oder Datenbank-Server eingesetzt, so gelten für das Disk-Subsystem natürlich zusätzliche Kriterien, wie sie für Datei- oder Datenbank-Server typisch sind. Von einer solchen Konstellation ist jedoch abzuraten, es sei denn, das System wird nur in einem sehr begrenzten Workgroup-Umfeld eingesetzt. Ansonsten sollten dedizierte Systeme für die einzelnen Aufgaben, wie Terminal Server, Datei-Server, Datenbank- oder Applikationsserver, aufgebaut werden. Nur so kann ein Server-System optimal für seinen Aufgabenbereich zugeschnitten werden.

Insbesondere bei einer »load-balanced Terminal Server-Farm« ist es wichtig, die Dateien der Benutzer nicht lokal auf den Terminal Server abzulegen, sondern auf einem zentralen Disk-Subsystem. Hierzu ist ein Network Attached Storage (NAS) oder auch ein klassischer File-Server am besten geeignet.

Man unterscheidet folgende Arten von Disk-Subsystemen:

Direct Attached Storage (DAS) bezeichnet eine Speicher-Technologie, bei der die Festplatten direkt an einen oder mehrere im Server eingebauten Festplatten-Controller angeschlossen werden. Typischerweise wird SCSI in Verbindung mit intelligenten RAID-Controllern eingesetzt. Diese Controller sind relativ preisgünstig und bieten eine gute Performance. Die Festplatten finden entweder im Server-Gehäuse oder in externen Disk-Gehäusen Platz. Ein DAS bietet erstklassige Performance und ist bei kleinen und mittleren Installationen eine gute Wahl. Für große Installationen sind jedoch begrenzte Skalierung, die aufwändige Verkabelung und die eingeschränkte Cluster-Tauglichkeit nachteilig. Da das Disk-Subsystem dediziert einem Server zugeordnet ist, ist es für eine zentrale Ablage der Benutzerdateien in Server-Farmen nicht nutzbar.

Network Attached Storage (NAS) ist im Prinzip ein klassischer File-Server. Ein solcher NAS-Server ist spezialisiert auf die effiziente Verwaltung großer Datenmengen und stellt diesen Speicher über ein LAN anderen Servern zur Verfügung. Intern verwenden NAS-Server typischerweise wieder die Platten- und Controller-Technologie des DAS. Für den Datentransport von und zu den Servern werden klassische LAN-Infrastrukturen genutzt. Dadurch können NAS-Systeme recht kostengünstig aufgebaut werden. Da der Datenspeicher nicht dediziert einem Server zugeordnet ist, lässt es sich ideal von Server-Farmen zur zentralen Ablage der Benutzerdateien nutzen.

Storage Area Network (SAN) lässt sich auf Basis von Fibre-Channel und LAN (iSCSI) aufbauen.

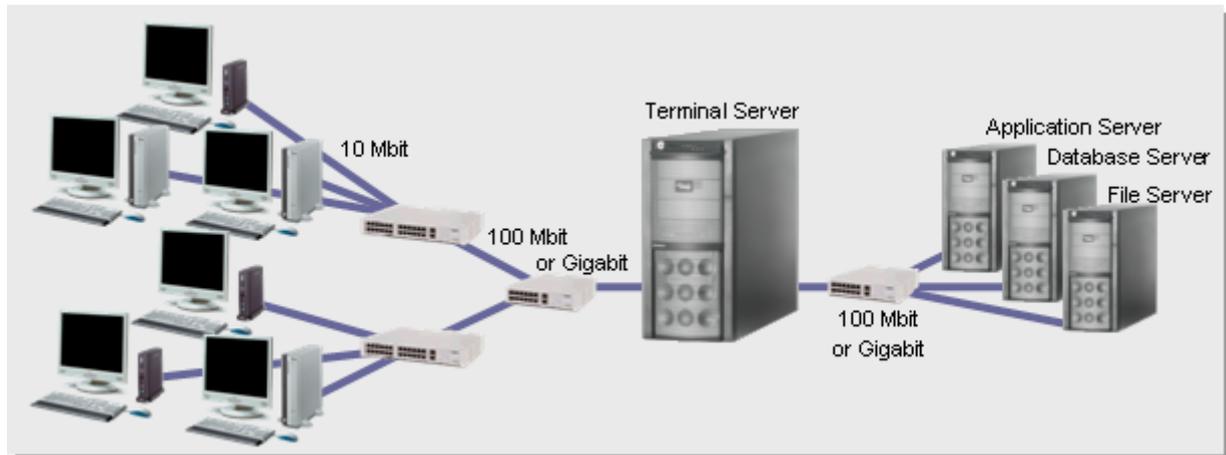
- **SAN auf Fibre-Channel** Basis verwendet im Gegensatz zum NAS nicht das LAN für den Datentransport, sondern ein eigenes Netz hoher Bandbreite auf Basis von Fibre-Channel (FC). Zum Aufbau einer Fibre-Channel Infrastruktur sind spezielle FC-Controller sowie eine spezielle Verkabelung und FC-Switches notwendig.
- **SAN mit iSCSI (IP SAN)** verwendet im Gegensatz zum SAN auf Fibre-Channel Basis das LAN für den Datentransport. Im Laufe des vergangenen Jahres begann »Internet SCSI« (iSCSI) als weitere Form des SAN sich immer weiter durchzusetzen. Dabei ist nicht, wie bei Fibre-Channel, eine komplett eigene Infrastruktur mit speziellen Boards, spezieller Verkabelung und Switches notwendig, sondern es können vorhandene TCP/IP Infrastrukturen verwendet oder gegebenenfalls erweitert werden. Der Aufbau und das Management separater Speichernetze für iSCSI unterscheiden sich nicht vom Aufbau und der Verwaltung von »normalen« TCP/IP-Netzen. Beim Betrieb von iSCSI kann in den meisten Fällen ganz auf spezielle Controller verzichtet werden, da es Softwarelösungen gibt, die auf den in den PRIMERGY Servern eingebauten Netzwerkadapter die notwendige Funktionalität zur Verfügung stellen.

Bei einem SAN werden alle Server und die Storage Systeme miteinander verbunden. Allerdings sind die Daten dediziert den einzelnen Servern zugeordnet, so dass ein SAN, unabhängig davon ob es sich um Fibre-Channel oder um iSCSI handelt, sich nicht zur zentralen Ablage für die Benutzerdateien in einer Server-Farm eignet.

Zum Sizing von Disk Subsystemen wurde ein eigenes White Paper »[Disk Subsystem Sizing – RAID Controller](#)« erstellt.

Netzwerk

Eine wichtige Rolle im Terminal Server-Umfeld kommt dem Netzwerk zu. Die Thin-Clients eines Terminal Servers dürften häufig – aus Ihrer Entstehungsgeschichte als ältere kleine PC-Systeme – über 10-Mbit-Ethernet oder noch langsamere PPP-Anbindungen verfügen. Zur Server-Seite bündeln sich natürlich solche Bandbreiten: wenn pro Client beispielsweise 100 kbit an Bandbreite benötigt wird, so wäre theoretisch mit 100 Clients ein 10-Mbit-Ethernet ausgeschöpft. Bei einer solchen Anzahl an Clients bietet es sich also an, den Server mit einer höheren Bandbreite von 100 Mbit oder Gigabit an den Backbone des Netzwerks anzubinden und die Bandbreite über Switches oder Hubs an die Clients zu verteilen.



Bedingt es das Aufgabenszenario, dass die auf dem Terminal Server ablaufenden Anwendungen auf große Datenmengen, Datenbanken oder gar Host-Anwendungen zugreifen, so empfiehlt es sich, den Terminal Server mit einer weiteren Netzwerkkarte für den dedizierten Zugriff auf diese Serverdienste auszustatten, um so den Datenverkehr in dieser Three-Tier-Umgebung zwischen Server-Server- und Client-Server-Kommunikation zu trennen.

In der Praxis sind dabei jedoch häufig Kompromisse einzugehen, da vorhandene Netzwerk-Topologien zu berücksichtigen sind und der Terminal Server darin zu integrieren ist.

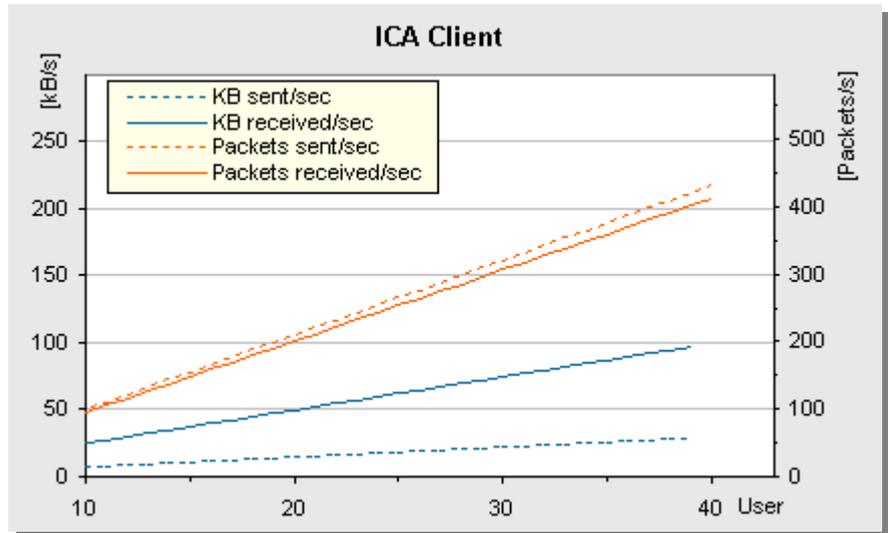
In einer Terminal Server-Umgebung muss folgender Netzwerkverkehr berücksichtigt werden:

- Terminal Server und Active Directory
Hauptsächlich während des Anmeldevorgangs der einzelnen Benutzer in die Domäne werden Informationen aus dem Active Directory benötigt. Dies wird in realen Konfigurationen oft auch über ein Netzwerksegment abgewickelt, das von den Client-Netzwerken getrennt ist.
- Clients und Terminal Server
Tastatur- und Mauseingaben werden zum Terminal Server gesendet, und die Änderungen der Bildschirmdarstellung werden zum Client übertragen.

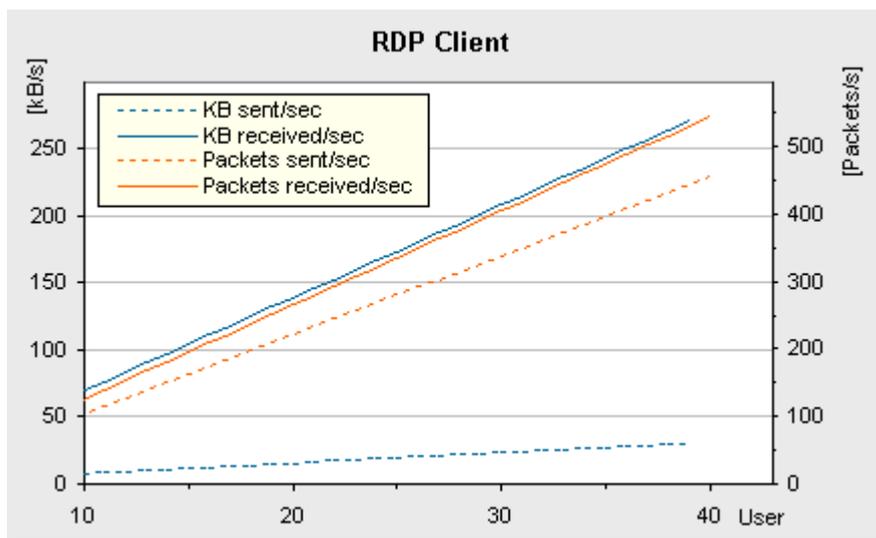
Terminal Server-Clients können über WAN, LAN oder W-LAN angeschlossen werden, wobei die kabellose Anbindung mit W-LAN immer beliebter wird. Gerade das ICA-Protokoll mit seiner vorteilhaften Komprimierung ist optimiert für Verbindungen bis herunter zu 14.4 kbit/s.

Stellt man die während einer Messung aus Sicht des Clients gesendeten und empfangenen Daten des Terminal Servers grafisch dar, so sieht man, dass die Daten, die pro simulierten Benutzer über das Netzwerk übertragen werden, linear skalieren. Diese lineare Skalierung kann bei allen PRIMERGY Systemen beobachtet werden. Erst unter Hochlast wird diese gleich bleibende Steigerung der Datenrate nicht mehr erreicht.

Auf einer PRIMERGY RX200 S2 wurde in diesem Beispiel die Netzlast von bis zu 40 Benutzern untersucht. Die Grafik »ICA Client« veranschaulicht die Messergebnisse. Bei dem Medium Lastprofil werden pro Benutzer durchschnittlich etwa 0.72 KB pro Sekunde in 11.03 Paketen zum Terminal Server gesendet und 2.41 KB pro Sekunde in 10.55 Paketen vom Terminal Server empfangen. Dieser Messung lag das ICA-Protokoll zu Grunde.



(Medium Lastprofil, Microsoft Office 2003, Citrix MetaFrame)

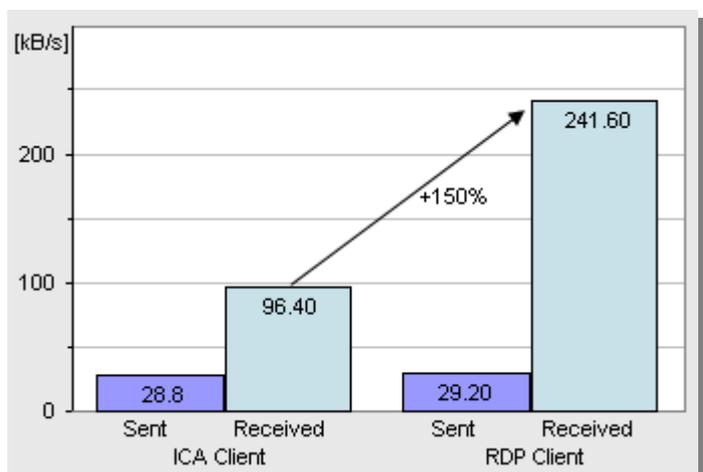


(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Vergleicht man die Netzwerk-Datenrate einer Messung mit dem ICA-Protokoll mit einer vergleichbaren Messung unter Verwendung des RDP-Protokolls, so erkennt man, dass beim RDP-Protokoll in Richtung des Clients (»Received«) mehr Daten übertragen werden müssen. Dies sind bei dem Medium Lastprofil pro Benutzer durchschnittlich etwa 0.73 KB pro Sekunde in 11.59 Paketen in Senderichtung und 6.04 KB pro Sekunde in 13.75 Paketen in Empfangsrichtung, vom Client aus gesehen. Die Grafik »RDP Client« im gleichen Maßstab zeigt das Messergebnis von Microsoft Terminal Server.

Der Unterschied zwischen Microsoft Terminal Services und Citrix Presentation Server wird durch eine andere grafische Darstellung der Messergebnisse noch besser verdeutlicht, wie das nebenstehende Bild zeigt. Bei einer PRIMERGY RX300 S2 mit 40 Benutzern wie in diesem Beispiel werden bei den Microsoft Terminal Services bis zu 150% mehr Daten vom Server zum Client übertragen als bei Citrix Presentation Server.

Im LAN-Umfeld wird das zugrunde liegende Netzwerk normalerweise keinen Engpass darstellen. Im WAN-Umfeld steht jedoch weniger Bandbreite pro Benutzer zur Verfügung.



Benutzerverhalten

Neben der Server Hardware, die in den vorherigen Kapiteln ausführlich diskutiert wurde, gibt es noch weitere Größen, die das Verhalten des Terminal Servers entscheidend beeinflussen. Die Arbeitsweise der Benutzer spielt dabei eine wichtige Rolle. Dies wurde am Beispiel der Eingabegeschwindigkeit untersucht.

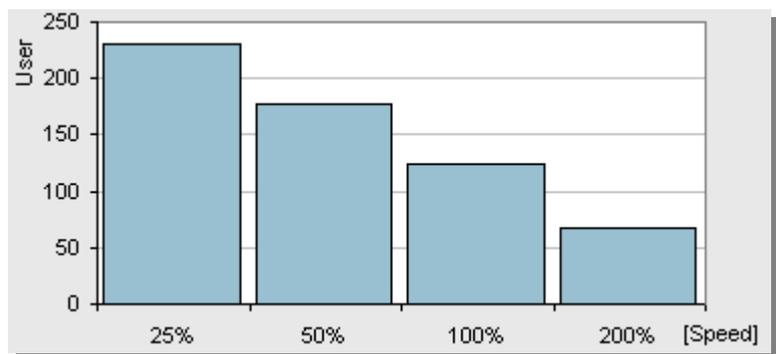
Eingabegeschwindigkeit

Welchen Einfluss hat das Benutzerverhalten auf die Leistungsfähigkeit von Terminal Server?

Die Eingabegeschwindigkeit hat einen maßgeblichen Einfluss auf die Performance eines Terminal Servers. Allein durch die Variation der Eingabegeschwindigkeit kann bei sonst gleichen Bedingungen ein Server unausgelastet sein oder überlastet werden. Unter Laborbedingungen arbeiten alle simulierten Benutzer kontinuierlich und mit der gleichen Eingabegeschwindigkeit. In der Realität wird dies stark variieren.

Bei Terminal Server Anwendungen werden Tastatureingaben und Mausclicks des Benutzers zum Server übertragen und die Änderungen des Bildschirms zurück zum Client. Daher wird jede Aktion des Benutzers mehrere Prozesse auf dem Server System aktivieren und auch Netzwerkverkehr auslösen.

Die Grafik zeigt die Benutzeranzahl, die ein Microsoft Terminal Server mit den vorgegebenen Antwortzeiten bedienen kann, für jeweils 25%, 50%, 100% und 200% Eingabegeschwindigkeit. Die real aufgezeichnete Eingabegeschwindigkeit entspricht dabei 100%. Bei 50% tippt der simulierte Benutzer nur halb so schnell und bei 200% gibt er doppelt so schnell ein. Für die Messung wurde eine PRIMERGY RX300 S2 mit zwei 3.6 GHz-Prozessoren, eingeschaltetem Hyper-Threading und 4 GB Arbeitsspeicher als Terminal Server Testsystem verwendet.



(Medium Lastprofil mit verschiedenen Eingabegeschwindigkeiten, Microsoft Office 2003, Microsoft Terminal Services)

Durch eine Verdoppelung der Eingabegeschwindigkeit von der normalen Geschwindigkeit 100% auf 200% verringert sich die Anzahl Benutzer, die ein Terminal Server mit den vorgegebenen Antwortzeiten bedienen kann um 86.5% von 125 auf 67. Der begrenzende Faktor ist die CPU.

Bei einer Reduzierung der Eingabegeschwindigkeit von 100% auf 50% und 25% kann der Terminal Server mehr Benutzer bedienen. Man erreicht eine Zahl von 177 bzw. 230 Benutzern, jedoch läuft der Server in einen Speicherengpass, während die CPU noch Reserven hat.

Betriebssystem

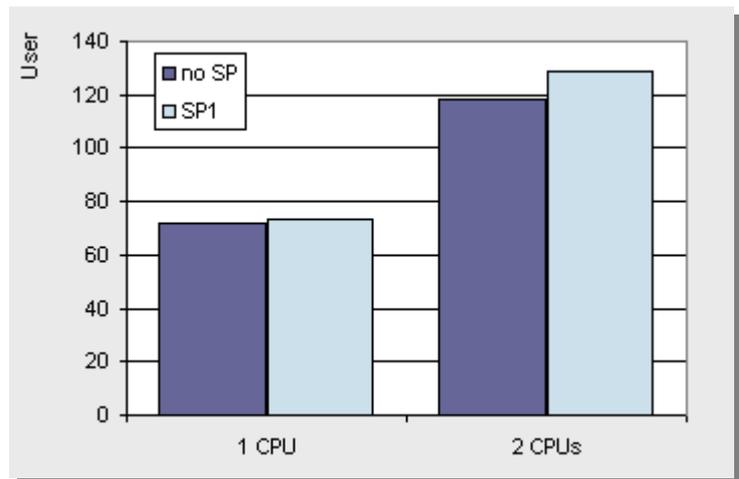
Neben der Server Hardware, die in den vorherigen Kapiteln ausführlich diskutiert wurde, gibt es noch weitere Größen, die das Verhalten des Terminal Servers entscheidend beeinflussen. Eine davon ist das Betriebssystem, das heute entweder in der 32-bit Variante oder in der 64-bit Variante eingesetzt wird. Die Leistung eines Terminal Servers kann von den Einschränkungen des 32-bit Betriebssystems bestimmt werden, und es stellt sich oft die Frage, ob diese Limitierungen durch den Einsatz eines 64-bit Betriebssystems überwunden werden können. Zum Verständnis sind allerdings Grundkenntnisse der Unterschiede der zwei Betriebssystemversionen notwendig, daher wird zuerst ein Überblick über das 64-bit Betriebssystem und die Vor- und Nachteile gegenüber 32-bit gegeben.

Windows Server 2003 R2

Die aktuelle Version von Windows Server 2003 ist das Release 2, kurz R2 genannt. Die 32-bit und 64-bit Versionen von Windows Server 2003 R2 basieren dabei auf der gleichen Code-Basis und sind daher direkt vergleichbar. Windows Server 2003 R2 ist bis auf einige zusätzliche Dienste und Tools identisch mit Windows Server 2003 Service Pack 1. Gegenüber dem Windows Server 2003 enthält Windows Server 2003 R2 einige Performanceverbesserungen für Terminal Server Systeme:

- Verbesserungen beim Memory Management Paging Verhalten.
- Verbesserungen beim Sperren der Registrierung.
- Verbesserungen bei der Kernel Timer Verarbeitung.

Zum Vergleich zwischen Windows Server 2003 und Windows Server 2003 R2 wurden auf einer PRIMERGY RX300 S2 Messungen mit Windows Server 2003 Enterprise Edition mit und ohne installiertem Service Pack 1 durchgeführt. Als Prozessor wurde eine Xeon 3.6 GHz CPU mit 1 MB SLC verwendet. Hyper-Threading war eingeschaltet. Die nebenstehende Grafik zeigt die Performanceverbesserungen mit Windows Server 2003 Service Pack 1. In der gleichen Messumgebung kann das gleiche System mit Service Pack 1 bis zu 9% mehr Benutzer bedienen als ohne Service Pack. Der Performancegewinn durch das Service Pack 1 ist auf einem Monoprozessorsystem geringer als auf einem System mit mehreren Prozessoren.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Alle weiteren Messergebnisse in diesem Dokument wurden mit Windows Server 2003 R2 ermittelt.

64-bit

Heute sind zwei verschiedene 64-bit Plattformen etabliert, die sich grundlegend unterscheiden.

Intel Itanium (IA64)

Bereits seit langem sind 64-bit PRIMERGY Systeme auf Basis von Intel Itanium und Itanium-2 Prozessoren verfügbar. Diese Prozessoren sind jedoch nicht Code-kompatibel zu 32-bit-Anwendungen (x86). x86-Anwendungen werden auf Itanium Systemen lediglich emuliert, was jedoch einen entsprechenden Performance-Verlust durch die Emulationsschichten bedeutet. Obgleich Itanium Systeme seit Jahren verfügbar sind, gibt es bis heute nur sehr wenige Anwendungen, die für Itanium Prozessoren optimiert und kompiliert sind. So ist zwar Windows Server 2003 und somit die Terminal Services für Itanium Systeme verfügbar, jedoch gibt es Citrix Presentation Server nicht für Itanium, und auch Standard-Office Anwendungen, wie Microsoft Office, sind nicht für die Itanium Plattform verfügbar.

Daher wäre es äußerst ineffizient, einen Terminal Server auf einer Itanium-Architektur zu betreiben, denn (fast) alle Benutzer-Anwendungen müssten emuliert werden. Dies ist effektiv langsamer als ein Terminal Server auf Basis der 32-bittigen x86-Architektur.

x64

Einen weicheren Migrationsweg von 32-bit auf 64-bit ermöglichen Prozessoren, die 100% kompatibel zur x86-Architektur sind und Erweiterungen für 64-bit bieten. Hierzu zählen der AMD Opteron mit der AMD64-Architektur und die Intel Pentium und Xeon Prozessoren der neuesten Generation mit EM64T-Architektur. Beide Architekturen werden zumeist kurz mit x64 bezeichnet.

Da bei einem x64-System alle 32-bit Anwendungen ohne Emulation ablaufen, ist der Overhead wesentlich geringer als bei einem Itanium-System. Dennoch ist es natürlich optimal, wenn x64-Versionen der Software verfügbar sind. Im zweiten Quartal 2005 hat Microsoft zeitgleich mit dem Service Pack 1 für Windows Server 2003 auch eine x64-Version von Windows Server 2003 herausgegeben. Beide Versionen liegen heute in der »R2«-Version vor.

Viele systemnahe Anwendungen und alle Anwendungen, die Treiberkomponenten enthalten, müssen speziell für x64-Systeme angepasst werden, da diese im 32-bit-Modus nicht ablauffähig sind. Beispiele für Anwendungen, die speziell für 64-bit angepasst wurden oder werden, sind Citrix Presentation Server, Microsoft SQL Server oder Microsoft Exchange Server. 16-bit-Anwendungen für DOS oder Windows sind generell nicht mehr ablauffähig, dies ist insbesondere ein Problem für Anwendungen mit einem 16-bit-Installer.

Andere 32-bit Anwendungen können direkt unter dem x64-Betriebssystem ausgeführt werden. Hierzu gehört auch Microsoft Office. Das 64-bit Windows stellt das WoW64 (»Windows on Windows64«) Subsystem bereit, um einen reibungslosen Betrieb der 32-bit-Applikationen zu ermöglichen. Im WoW64 werden Zugriffe auf den Speicher, auf die Registry und auf das Dateisystem isoliert. Im Gegensatz zum IA64-System, auf dem die 32-bit Anwendung komplett emuliert werden muss, ist der x64-Prozessor 100% kompatibel zu x86 und daher ist der Overhead einer 32-bit-Anwendung auf einem 64-bit Betriebssystem gering und kann vernachlässigt werden. Allerdings wird eine 32-bit-Anwendung mehr Arbeitsspeicher verbrauchen, wenn sie auf einem 64-bit-Betriebssystem ausgeführt wird, wie viel mehr, hängt von der Anwendung ab.

Nachfolgende Tabelle gibt einen Überblick über die Systemplattformen, die Windows-Produkte und deren Anforderungen an Gerätetreiber und Anwendungen.

Server-Plattform	x86	x64	x64	IA64
Windows Produkt	Windows Server 2003 R2 Standard Edition Enterprise Edition Datacenter Edition	Windows Server 2003 R2 Standard Edition Enterprise Edition Datacenter Edition	Windows Server 2003 R2 Standard x64 Edition Enterprise x64 Edition Datacenter x64 Edition	Windows Server 2003 Enterprise Edition for Itanium-based Systems Datacenter Edition for Itanium-based Systems
Betriebssystem	32-bit	32-bit	64-bit	64-bit
Gerätetreiber	32-bit	32-bit	64-bit	64-bit
Anwendungen	32-bit	32-bit	32-bit und 64-bit	64-bit

Nutzbarer Speicher

Der wesentliche Unterschied zwischen 32-bit und 64-bit liegt in der Menge des nutzbaren Speichers. Einen Überblick über die Unterschiede der verschiedenen 32-bit und 64-bit Windows Produkte, die für Terminal Server sinnvoll eingesetzt werden können, gibt vorab folgende Tabelle, bevor die Auswirkungen dieser Betriebssystemeigenschaften für Terminal Server im Detail betrachtet werden:

Windows Produkt	Windows Server 2003 R2		Windows Server 2003 R2 x64	
	Standard Edition	Enterprise Edition	Standard Edition	Enterprise Edition
Anzahl Prozessoren	1 – 4	1 – 8	1 – 4	1 – 8
Max. RAM	4 GB	64 GB (mit PAE)	32 GB	1 TB
gesamter virtueller Adressraum	4 GB		16 TB	
Virtueller Adressraum eines 32-bit Prozesses	2 GB		2 GB	
	3 GB wenn mit /3GB Switch gebootet		4 GB wenn mit /LARGEADDRESSWARE übersetzt wurde	
Virtueller Adressraum eines 64-bit Prozesses	-		8 TB	
Paged Pool	470 MB		128 GB	
Non-Paged Pool	256 MB		128 GB	
System Page Table Entries	ca. 900 MB		128 GB	
System Cache	1 GB		1 TB	

32-bit

Bei der 32-bit-Version von Windows Server 2003 sind 4 GB eine magische Grenze, mehr lässt sich mit einer 32-bit-Adresse nicht adressieren. Durch die Systemarchitektur von Windows Server 2003 ist der Speicher von 4 GB standardmäßig in 2 GB für das Betriebssystem und 2 GB für die Anwendungen unterteilt. In den 2 GB, die vom Betriebssystem verwendet werden, werden alle Datenstrukturen und Informationen des Kerns gehalten. Hier sind drei besondere Datenbereiche zu nennen: die *Paged Pool Area*, die *System Page Table Entries (PTE) Area* und der *System File Cache*. Speicher aus der *Paged Pool Area* wird von Kernel-Mode Komponenten angefordert, während die *PTE Area* für Kernel Stack Allocations verwendet wird. Im *System File Cache* werden Speicherabbilder von geöffneten Dateien gehalten. Diese Bereiche teilen sich einen Speicherbereich und die Grenze zwischen ihnen wird beim Systemstart fest eingestellt. Wenn dem Betriebssystem während der Laufzeit in einem Bereich der Speicher ausgeht, kann dies nicht aus den anderen Bereichen ausgeglichen werden. Die Folge ist, dass keine neuen Benutzer mehr angemeldet werden können oder dass unerwartete Fehler auftreten.

In Windows Server 2003 wurden Optimierungen vorgenommen um den virtuellen Adressraum nun optimaler auszunutzen. Die voreingestellten Werte wurden angepasst, so dass für die meisten üblichen Konfigurationen eine Verbesserung erzielt wird. Dies hat zur Folge, dass unter Windows Server 2003 mehr Terminal Server-Benutzer auf einem System arbeiten können als unter Windows 2000 Server. Microsoft gibt an, dass bis zu 80% mehr so genannte »Knowledge Worker« und über 100% mehr »Data Entry Worker« betrieben werden können. Dies hängt natürlich stark von dem verwendeten Benutzerprofil und von dem verwendeten System ab und trifft auch nur für die Fälle zu, in denen der virtuelle Adressraum der Engpass gewesen ist und nicht die CPU oder der physikalische Arbeitsspeicher.

Für ganz spezielle Konfigurationen können die Werte natürlich individuell angepasst werden, außerdem kann PAE eingeschaltet werden. Weitere Details entnehmen Sie bitte dem Kapitel »[Physical Address Extension \(PAE\)](#)«.

Nähere Informationen finden Sie bei Microsoft im Knowledge Base Artikel Q247904 sowie im Dokument »Windows Server 2003 Terminal Server Capacity and Scaling«.

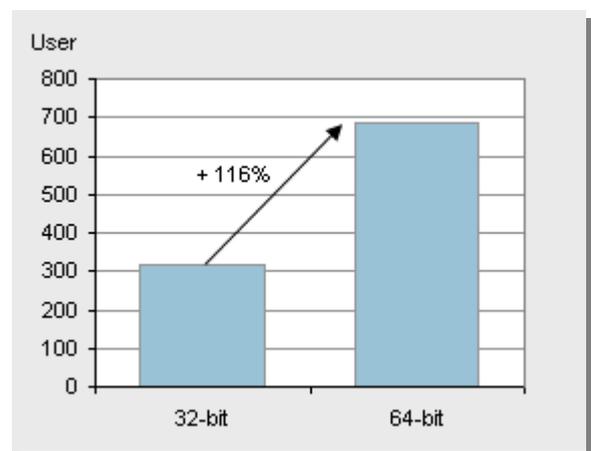
64-bit

Einer der wesentlichen Vorteile von 64-bit ist der erweiterte Adressraum. Heutige Server können problemlos mit mehr als 4 GB Arbeitsspeicher ausgestattet werden. Dieser ist auf 32-bit Systemen nur mit erhöhtem Aufwand adressierbar. Mit 64-bit Windows können theoretisch direkt 2^{64} Bytes = 16 Exabyte adressiert werden. Windows teilt diesen (zumindest aus heutiger Sicht) gigantischen Adressraum in verschiedene Bereiche auf, so dass für den Kernel und jeden 64-bit Prozess jeweils ein Adressraum von 8 TB zur Verfügung steht. Für 32-bit Anwendungen, die im so genannten Kompatibilitätsmodus laufen, steht pro Anwendung ein Adressbereich von 4 GB zur Verfügung, aber auch das ist mehr als bei einem reinen 32-bit Betriebssystem, bei dem es maximal 3 GB sind.

Insbesondere Anwendungen, die Speicher- und nicht CPU-limitiert sind, profitieren von der 64-bit Architektur. Dabei sollte aber auch nicht verschwiegen werden, dass 64-bit Betriebssysteme und 64-bit Anwendungen in der Regel mehr Arbeitsspeicher benötigen als die 32-bit-Versionen, denn alle Adresszeiger sind bei 64-bit doppelt so breit. Im Extremfall führt das bei 64-bit zu einem doppelten Speicherbedarf im Vergleich zu 32-bit.

Vergleich 32-bit und 64-bit

Bei Terminal Server-Systemen mit ihrer Vielzahl von Benutzern kommen normalerweise viele kleinere Anwendungen zum Einsatz. Hier kann es vorkommen, dass die Anzahl der Benutzer nicht durch den Speicherbedarf pro Anwendung begrenzt wird, sondern durch fehlende interne Ressourcen des 32-bit-Betriebssystemkerns. Diese Limitierung wird durch den Einsatz eines 64-bit-Betriebssystems aufgehoben. Wenn ausreichend Hardware-Ressourcen zur Verfügung stehen, also weder die CPU noch der Arbeitsspeicher der begrenzende Faktor ist, dann können auf einem 64-bit System wesentlich mehr Benutzer betrieben werden als auf einem 32-bit System. Für diesen Test wurden 32 GB Hauptspeicher verwendet. Wie nebenstehende Grafik zeigt, ist hier das 64-bit Betriebssystem klar im Vorteil und kann mehr als doppelt so viele Benutzer bedienen als das 32-bit Windows. Der Grund hierfür liegt in der Betriebssystemarchitektur, genauer gesagt, bei den Kernel-Ressourcen. Das 32-bit Betriebssystem hat einen Adressraum von 2 GB zur Verfügung, um seine Daten, unter anderem Kernel-Tabellen und System Cache, zu speichern. Bei hinreichender Rechenleistung erreicht die Benutzeranzahl eine Größenordnung, die für ein 32-bit Betriebssystem zu hoch ist. Die Kernel-Tabellen sind vollständig belegt, was in Paging-Aktivitäten resultiert, obwohl noch über 20 GB Hauptspeicher frei sind. Weiterhin ist auch der System Cache überlastet und kann nicht weiter vergrößert werden, was zu einer schlechten Trefferrate im Cache führt. Hierdurch steigen ebenfalls die Plattenzugriffe sprunghaft an. Sobald auf die Festplatte statt auf den Arbeitsspeicher zugegriffen werden muss, macht sich das sofort in einer Verschlechterung der Antwortzeiten bemerkbar. Würde man den Terminal Server über diese Engpasssituation hinaus weiter belasten, so würden die Applikationen auf Fehler laufen oder Benutzer könnten sich beim Terminal Server gar nicht mehr anmelden.

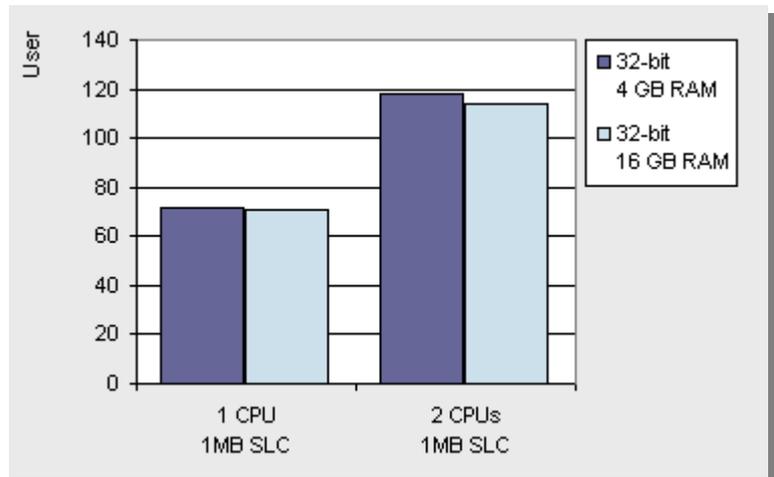


(Light Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

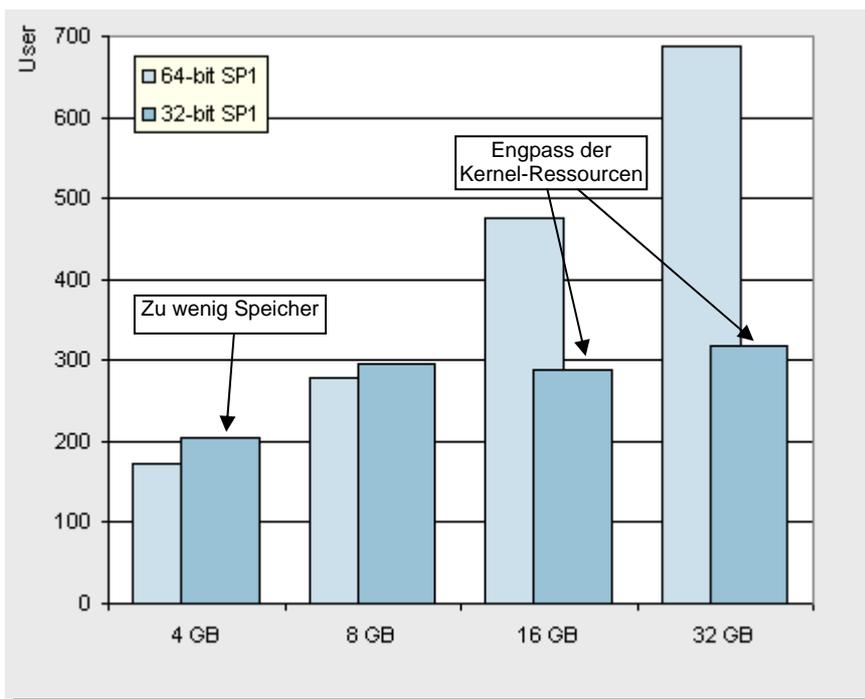
Physical Address Extension (PAE)

Auf einem 32-bit System mit einem 32-bit Adressraum können nur 4 GB Arbeitsspeicher adressiert werden. Um diese Grenze zu überwinden und mehr Speicher zu nutzen, kann auf dieser Plattform PAE eingeschaltet werden.

Wenn ein Terminal Server durch die Rechenleistung und nicht durch den Arbeitsspeicher begrenzt wird, kann PAE einen leicht negativen Einfluss haben. Wie das Diagramm zeigt, führt das Einschalten von PAE auf einem System, bei dem die CPU der Engpass ist, zu einem leichten Performanceabfall von 2% bis 6%. Die heutigen Dual-Core Prozessoren sind jedoch bereits in einem 2-Socket-System so leistungsfähig, dass bei 32-bit Systemen ein Speicherausbau von 4 GB nicht mehr ausreichend sein wird und der geringe Performanceverlust durch PAE durch die höhere Prozessorleistung kompensiert wird.



(Medium Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)



(Light Lastprofil, Microsoft Office 2003, Microsoft Terminal Services)

Sichtbar wird dies auf einem 32-bit System, bei dem nicht die CPU sondern der Arbeitsspeicher der begrenzende Faktor ist. Das gleiche PRIMERGY System wurde mit 4, 8, 16 und 32 GB RAM ausgestattet und die maximale Anzahl Light User ermittelt. Nebenstehende Grafik zeigt die Messergebnisse für das 32-bit Betriebssystem im Vergleich zu 64-bit. Bei 4 GB ist die Größe des physikalischen Hauptspeichers der Engpass, was in Paging-Aktivitäten resultiert. Bei einer Vergrößerung des Speichers auf 8 GB können mehr Benutzer bedient werden, dieser Speicherausbau ist die optimale Konfiguration für das 32-bit System. Darüber hinaus werden die Betriebssystemressourcen zum Engpass, eine weitere Speicheraufrüstung ist daher nicht von Vorteil oder

kann sich sogar negativ auswirken. Dies ist durch die PAE-Adressierung zu erklären, die mit größerem Speicher aufwändiger wird und mehr des ohnehin begrenzten Kernel-Speichers belegt. Im Gegensatz dazu skaliert das 64-bit Betriebssystem auch über 8 GB hinaus, und die Benutzeranzahl kann, genügend Rechenleistung vorausgesetzt, durch eine Vergrößerung des Hauptspeichers weiter gesteigert werden.

Auch unter einem 32-bit Betriebssystem kann ein Terminal Server also mit PAE auf leistungsfähiger Hardware mit einem Speicherausbau über dem 4 GB Limit betrieben werden. Es gibt es jedoch eine Grenze, die je nach Benutzerlast früher oder später erreicht wird, ab der das 32-bit Betriebssystem nicht mehr skaliert. Das 64-bit Betriebssystem ohne die Limitierungen bei den Kernel-Ressourcen kann mit mehr Hauptspeicher ausgestattet werden, was gerade bei den leistungsfähigen Servern mit Multi-Core-Architekturen von Vorteil ist.

Anzahl Prozesse

Auch wenn es paradox klingt: es kann Konstellationen geben, in denen, obgleich ausreichend CPU- und Speicher-Ressourcen zur Verfügung stehen, es zu Leistungsengpässen kommen kann. Auch Disk-I/O oder Netzwerke stellen in dieser Situation keinen Engpass dar, sondern quasi die Systemarchitektur. Diese Situation wird insbesondere von einer großen Anzahl an Benutzern, die viele Prozesse mit einer geringen gleichmäßigen Last induzieren, provoziert. Wie bereits im Abschnitt »[Messmethode](#)« erläutert, spielt die Prozessor-Warteschlange eine nicht zu vernachlässigende Rolle. So gibt es in Abhängigkeit der Prozessor-Leistung und Prozessor-Anzahl einen Punkt, an dem das System keine weiteren Prozesse und somit Clients mehr bedienen kann. Im Prinzip ist hierbei das System nur noch damit beschäftigt, die Prozesse zu verwalten. Oder in »Fachchinesisch« ausgedrückt: In einem Multitasking-Betriebssystem erhält jeder Prozess eine Zeitscheibe, die er nicht schnell genug wieder frei gibt. Diese Situation könnte nur durch kleinere Zeitscheiben und somit größeren Turn-Around-Zeiten behoben werden. Dies hätte jedoch in anderen Lastsituationen negative Auswirkungen auf die Grundlast, die das Betriebssystem erzeugt. Die Anzahl Prozesse, ab der diese Situation eintritt, hängt neben der zur Verfügung stehenden CPU-Leistung und -Anzahl leider aber auch von der verwendeten Applikation ab, sodass keine generelle Formel hierfür angegeben werden kann. Bei Heavy Usern tritt dieser Effekt meist nicht zu Tage, da hier andere Ressourcen, wie Speicher und Rechenleistung die dominierenden Begrenzungsfaktoren sind.

Terminal Server Version

Microsoft Terminal Server vs. Citrix Presentation Server

Welche Unterschiede gibt es zwischen den Microsoft Terminal Services und Citrix Presentation Server?

Die größten Unterschiede zwischen den Microsoft Terminal Services und Citrix Presentation Server liegen im Netzwerkbereich. Dies wurde bereits im Kapitel »[Netzwerk](#)« behandelt. Wenn das Netzwerk keinen Engpass darstellt, liegen die Unterschiede bei der Benutzeranzahl, die ein Terminal Server bedienen kann, im Rahmen der Messungenauigkeit.

Schaut man sich die Performance Counter im Detail an, so kann man, abgesehen vom Netzwerk, weitere Unterschiede zwischen den verschiedenen Terminal Server Implementierungen erkennen. Da sich die Citrix Versionen unterscheiden, sind die Unterschiede in der nachfolgenden Tabelle aufgeführt, wobei die Tendenz bei allen Citrix Produkten gleich ist. Citrix Presentation Server belegt etwas mehr Hauptspeicher pro Benutzer, sowohl der »Working Set« als auch die »Committed Bytes« sind höher, dadurch verringert sich der freie Hauptspeicher »Available MBytes«.

Die Citrix Lösung hat Vorteile bei den Interrupts, es werden bei der gleichen Benutzerlast weniger Interrupts gebraucht. Nachteile zeigen sich jedoch bei den Context Switches, hier treten gegenüber den Microsoft Terminal Services mehr Context Switches auf. Bei den Festplattenzugriffen sind keine nennenswerten Unterschiede zu erkennen. Bei einem Terminal Server im Echteininsatz, der nicht an seiner Leistungsgrenze arbeitet, sind diese Unterschiede jedoch nicht performancerelevant.

Performance Counter	Unterschied zu Microsoft Terminal Server		
	Citrix MetaFrame XP 1.0	Citrix Presentation Server 4.0 (32-bit und 64-bit)	Vergleich
Working Set	+3%	+8% bis +15%	schlechter
Committed Bytes	+5%	+12% bis +20%	schlechter
Available MBytes	-8%	-4% bis -7%	schlechter

Citrix Presentation Server Version

Wie wirkt sich ein Upgrade der Citrix Software aus?

Aus Gründen der Vergleichbarkeit mit älteren Messungen wurden einige Citrix Messungen mit Citrix MetaFrame XP 1.0 FR3 durchgeführt. Das Nachfolgeprodukt war Citrix MetaFrame Presentation Server 3.0 (MPS 3.0). Heute ist Citrix Presentation Server 4.0 aktuell, den es in einer 32-bit Version und in einer 64-bit Version gibt.

Auf einer PRIMERGY RX300 S2 wurde der Unterschied der Citrix Versionen unter Windows Server 2003 Enterprise Edition untersucht. Nur die Software auf dem Terminal Server wurde aktualisiert, der ICA Client auf den Lastgeneratoren wurde unverändert gelassen.

Beim Vergleich von Citrix MetaFrame 1.0 FR3 mit Citrix MetaFrame Presentation Server 3.0 wird die gleiche absolute Benutzeranzahl erreicht, die Prozessorleistung ist hier der begrenzende Faktor. Auch wenn die absolute Benutzeranzahl in beiden Fällen gleich ist, so erkennt man doch Unterschiede zwischen den beiden Versionen, wenn man die während der Messung aufgezeichneten Performance Counter vergleicht. Bei MPS 3.0 ist der Speicherverbrauch insgesamt höher, es werden auch mehr Prozesse pro Benutzer gestartet.

Ob man mit Citrix Presentation Server 4.0 die gleiche Anzahl Benutzer erreicht wie mit den Vorgängerprodukten, hängt vom Hauptspeicherausbau des verwendeten PRIMERGY Systems ab. Citrix Presentation Server 4.0 belegt mehr Hauptspeicher pro Benutzer, auch in der 32-bit Version. Deutlich mehr Hauptspeicher verbraucht die 64-bit Version von Citrix Presentation Server, wie es bei einem 64-bit Programm üblich ist. Die Tabelle zeigt den Speicherverbrauch der verschiedenen Citrix Versionen pro Benutzer, wie er beim Medium Lastprofil in unserer Messumgebung ermittelt wurde.

Keine Unterschiede erkennt man bei den Ressourcen »Disk«, »Netzwerk« und »Prozessor«.

Citrix Version	MB pro User	Committed Bytes pro User	Working Set pro User
MetaFrame XP 1.0	20	33	38
MetaFrame Presentation Server 3.0	20	33	36
Presentation Server 4.0	22	35	38
Presentation Server 4.0 x64	30	45	44

Anwendungen

Auch die Version und bestimmte Einstellungen der Anwendungen, die unter dem Terminal Server zur Verfügung gestellt werden, können die Leistungsfähigkeit des Terminal Servers beeinflussen.

Microsoft Office Version

Inwieweit beeinflusst eine neue Version einer Anwendung die Leistungsfähigkeit von Terminal Server?

Erfahrungsgemäß stellen neuere Versionen einer Anwendung mehr Funktionen zur Verfügung, haben jedoch höhere Anforderungen an das System und benötigen mehr Ressourcen.

Dies wurde am Beispiel von Microsoft Office untersucht. Dabei wurde Microsoft Office XP mit dem neueren Microsoft Office 2003 verglichen. Als Anwendung im Medium Lastprofil diente Microsoft Word, als Terminal Server wurden die Terminal Services von Microsoft mit dem RDP Protokoll verwendet.

Die maximale Benutzeranzahl, die ein Terminal Server bedienen kann, wird durch die Prozessorleistung limitiert und hängt nicht davon ab, welche Office Version verwendet wird.

Jedoch zeigen sich auch hier Unterschiede beim detaillierten Vergleich der Performance Counter. Es wurde festgestellt, dass Office 2003 im Vergleich zu Office XP einen geringfügig höheren CPU-Verbrauch hat, aber die Unterschiede sind nicht signifikant. Bei Office 2003 wird weniger von der Systemplatte gelesen, aber mehr geschrieben. Da das geschriebene Dokument bei den beiden Office-Versionen gleich groß ist, ergeben sich keine Differenzen bei der Datenplatte. Office 2003 belegt mehr Netzwerk-Ressourcen, da mehr Daten auf das Netz gesendet und vom Netz empfangen werden. Dadurch steigt auch die Interrupt-Anzahl. In der Office 2003-Umgebung werden auf dem Terminal Server System weniger Prozesse gestartet. Office 2003 belegt etwas weniger Hauptspeicher-Ressourcen als Office XP. Auf durchschnittlich ausgelasteten Terminal Server Systemen resultieren diese Unterschiede jedoch nicht in einer unterschiedlichen Benutzeranzahl.

Einstellungen für Microsoft Office in einer Terminal Server Umgebung

Obwohl Microsoft Office XP und Microsoft Office 2003 ohne Probleme auf einem Terminal Server installiert werden und ablaufen können, gibt es einige Empfehlungen von Microsoft, die beim Einsatz dieser Anwendungen in einer Terminal Server Umgebung zu einer verbesserten Performance führen können.

Für Office 2003 kann folgender Registrierungsschlüssel gesetzt werden, damit der Verbindungsstatus weniger häufig abgefragt wird:

```
[HKEY_CURRENT_USER\Software\Microsoft\Office\11.0\Outlook\RPC]
"ConnManagerPoll"=dword:0x00000600
```

Falls der RPC-Schlüssel nicht existiert, muss er vorher angelegt werden.

Andere empfohlene Einstellungen der Microsoft Office Konfiguration sind:

- Häufig genutzte Sicherungsoperationen im Hintergrund abschalten:
 - »AutoWiderherstellen-Info speichern« in Word abschalten.
 - »Nicht gesendete Nachrichten automatisch speichern« in Outlook abschalten.
 - Automatische AutoArchivierung von Outlook Nachrichten abschalten.
- Überprüfungen im Hintergrund abschalten:
 - »Grammatik während der Eingabe prüfen« in Word abschalten.
 - »Namen automatisch prüfen« in Outlook abschalten.
- »E-Mail mit Word bearbeiten« abschalten.

Diese benutzerspezifischen Einstellungen müssen für alle Benutzer des Terminal Servers verändert werden.

Dies kann leicht für alle Benutzer konfiguriert werden, wenn die Installationsmöglichkeit genutzt wird, die die Anwendungskompatibilitätsschicht des Terminal Servers bietet. Hierzu muss nach der Installation der Anwendung und bevor sich ein Benutzer angemeldet hat von einem Kommandoprompt aus folgendes Kommando abgesetzt werden:

```
Change user /install
```

Dann wird die anzupassende Anwendung (z.B. Word) gestartet. Im Options-Dialog können die gewünschten Konfigurationsänderungen (beispielsweise das Abschalten der automatischen Grammatikprüfung) durchgeführt und die Applikation geschlossen werden. Danach das folgende Kommando aufrufen:

```
Change user /execute
```

Auf diese Weise werden die Konfigurationsänderungen der Registrierung durch die Anwendungskompatibilitätsschicht gesichert und beim nächsten Logon in die Registrierung jedes Benutzers eingetragen.

Infrastruktur

In unseren Untersuchungen haben wir den Terminal Server immer isoliert betrachtet. In der Messumgebung gab es weitere Komponenten, mit denen der Terminal Server zusammen arbeitet, jedoch waren diese immer konstant und so ausgelegt, dass diese nicht der Engpass waren. In der Realität ist dies jedoch nicht immer der Fall. In diesem Abschnitt soll diskutiert werden, welche weiteren Komponenten der Infrastruktur einen Einfluss auf das Benutzerempfinden in einer Terminal Server-Umgebung haben, was sich in einem negativen Gesamteindruck widerspiegeln könnte.

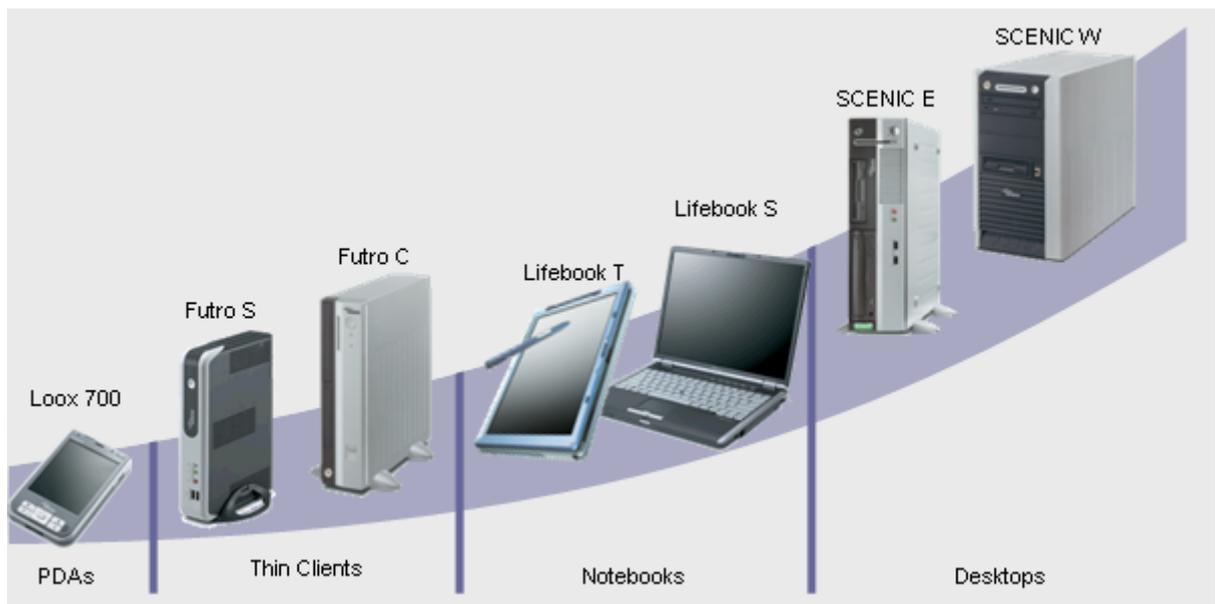
Clients

Neben den Server-Ressourcen und dem Netzwerk gehört natürlich auch der Client zur Gesamtumgebung des Server-based Computing. Also ist auch bzgl. des Clients die Frage berechtigt, welchen Einfluss dessen Leistungsfähigkeit auf die Gesamtkonfiguration hat. Da beim Server-based Computing die eigentliche Applikation auf dem Server abläuft, wird die CPU-Leistung des Clients nur für die Bedienung des Netzwerks und die vom Aufwand her nicht zu vernachlässigende Bildaufbereitung benötigt. Es wurde festgestellt, dass je nach verwendetem Client-System die Gesamt-Ausführungszeit einer Applikation durchaus um bis zu ca. 10% variiert. Diese Unterschiede dürften sich im realen Betrieb aber allenfalls in einem subjektiven Performance-Eindruck des Benutzers widerspiegeln und haben keinen unmittelbaren Einfluss auf die Leistungsfähigkeit des Terminal Server-Systems.

Wie »thin« der Thin-Client sein darf, hängt in erster Linie von den Ansprüchen seitens der eingesetzten Applikationen an die Grafik ab, wie Auflösung, Farbtiefe und Komplexität (Text, Grafik) sowie von ggf. zusätzlichen Anforderungen an lokal auf dem Client ablaufende Anwendungen, die über das reine Server-based Computing hinausgehen.

Es gibt heute mehr und mehr die Anforderung, neben den ortsgebundenen (Thin-)Clients auch mobile Geräte wie Notebooks oder PDAs an Terminal Server anzuschließen. Üblicherweise werden diese Geräte dann nicht mehr über kabelgebundene Netzwerke angeschlossen, sondern über Funk-LANs (W-LAN) oder mobile Funknetze (z.B. »General Packet Radio Services« (GPRS) oder »Universal Mobile Telecommunication System« (UMTS)).

Dies stellt insbesondere beim Ressourcen sparenden ICA-Protokoll kein Problem dar. Den ICA-Client gibt es auch für viele gängige PDAs und das ICA-Protokoll ist durch sein Design besonders auch für langsame Verbindungen geeignet. Für Benutzer, die ständig und ausschließlich mit Terminal Server Anwendungen arbeiten, stellt solch ein Endgerät sicher nicht den idealen Client dar, aber jemand, der mobil arbeiten muss und nur gelegentlich Zugriff auf einen Terminal Server braucht, profitiert von der Flexibilität und Funktionalität dieser Lösung.



Diese Clients sind nur einige mögliche Beispiele für die Palette von Terminal Server-Clients aus dem Hause Fujitsu Siemens Computers.

Active Directory

Terminal Server Benutzer authentifizieren sich im Normalfall in einer Domäne, d.h. der Terminal Server überprüft die eingegebenen Benutzercredentials gegen das Active Directory. Außer in sehr kleinen Workgroup-Umgebungen sollten Active Directory und Terminal Server immer auf verschiedenen Systemen laufen und auf dem Terminal Server selbst sollten keine Benutzer verwaltet werden. An das Active Directory werden die gleichen Anforderungen gestellt wie in einer Umgebung ohne Terminal Server, es sollte aber weder das Active Directory noch das Netzwerk zwischen Active Directory und Terminal Server der Engpass sein.

Benutzerprofile (User Profiles)

In einem Benutzerprofil werden die individuellen Benutzereinstellungen gespeichert. Auch bei einem Login von Terminal Server-Benutzern in einer Domäne in einem Active Directory Umfeld würden deren Benutzerprofile standardmäßig auf dem Terminal Server gespeichert. Insbesondere bei einer load-balanced Terminal Server-Farm wird man die Benutzerprofile allerdings zentral auf einem Server im Netzwerk ablegen wollen, damit der Benutzer immer die gleichen Einstellungen vorfindet, unabhängig davon, auf welchem Terminal Server seine Sitzung ausgeführt wird. Diese Funktionalität ist bereits für so genannte »Wandernde Benutzer« (Roaming User) vorhanden, die sich an verschiedenen Arbeitsplätzen anmelden. Beim Einsatz von Terminal Servern ist zu beachten, dass ggf. verschiedene Benutzerprofile verwaltet werden müssen, wenn sich nämlich das lokale Betriebssystem des Arbeitsplatzes von dem des Terminal Servers unterscheidet bzw. wenn verschiedene Anwendungen vorhanden sind. Aus diesem Grunde kann man ein Terminal Server Benutzerprofil zusätzlich zu einem lokal zu ladenden Benutzerprofil konfigurieren. Eine besondere Variante des serverbasierten Benutzerprofils ist das Mandatory User Profile, ein Benutzerprofil, das der Benutzer nicht ändern kann. Serverbasierte Benutzerprofile sollten generell möglichst klein sein.

DNS

Auch im Terminal Server Umfeld wird DNS verwendet, um die Namensauflösung von Verbindungen zu realisieren. Besonders im Load Balancing Umfeld wird auf diese Weise ein virtueller Name mit einer virtuellen oder realen IP Adresse verknüpft, so dass sich eine Terminal Server-Farm zum Benutzer hin wie ein Server darstellt. Daraus ergibt sich umgekehrt die Forderung, dass DNS immer erreichbar sein muss, damit ein Benutzer eine Verbindung zum Terminal Server aufbauen kann. DNS stellt im Allgemeinen keinen Engpass dar, dieser Dienst sollte nur ausfallsicher und redundant konzipiert sein.

Terminal Services Licensing Server

Bei der Anmeldung eines Benutzers an einen Terminal Server wird dieser einen Lizenzserver suchen und von ihm eine gültige Lizenz für den Zugriff über Terminal Server anfordern. In größeren Konfigurationen wird dieser Lizenzserver ein eigenes System sein.

Backend Server

Gerade in »load-balanced Terminal Server-Farmen« werden die Dateien der Benutzer nicht auf den lokalen Festplatten der Terminal Server Systeme liegen, sondern auf File Servern oder NAS Systemen. Vermutlich werden in größeren Umgebungen auch weitere Dienste wie E-Mail, Datenbanken usw. benötigt, so dass man Server für Anwendungen wie zum Beispiel Exchange, SQL oder SAP R/3 zusammen mit dem Terminal Server vorfinden wird. Wird für die Anbindung der Clients zum Terminal Server nur wenig Netzwerkbandbreite benötigt, so gilt dies nicht für die Anbindung der Terminal Server an die Backend Server. Hier sollte genügend Netzwerk- und Rechenkapazität zur Verfügung stehen. Für die einzelnen Server verweisen wir auf eigene Performance-Untersuchungen und Sizing Guides, da dies den Rahmen dieses Papiers sprengen würde.

Es wird nicht empfohlen, Backend-Dienste auf einem Terminal Server zu betreiben.

Vergleich der Messwerkzeuge

Wie bereits oben erwähnt, können verschiedene Messwerkzeuge als Ergebnis verschiedene Benutzerzahlen liefern, da die Benutzerprofile und Messmethoden unterschiedlich sind. Aber obwohl die absoluten Zahlen verschieden sind, kann man die relative Skalierung des Server-Systems vergleichen.

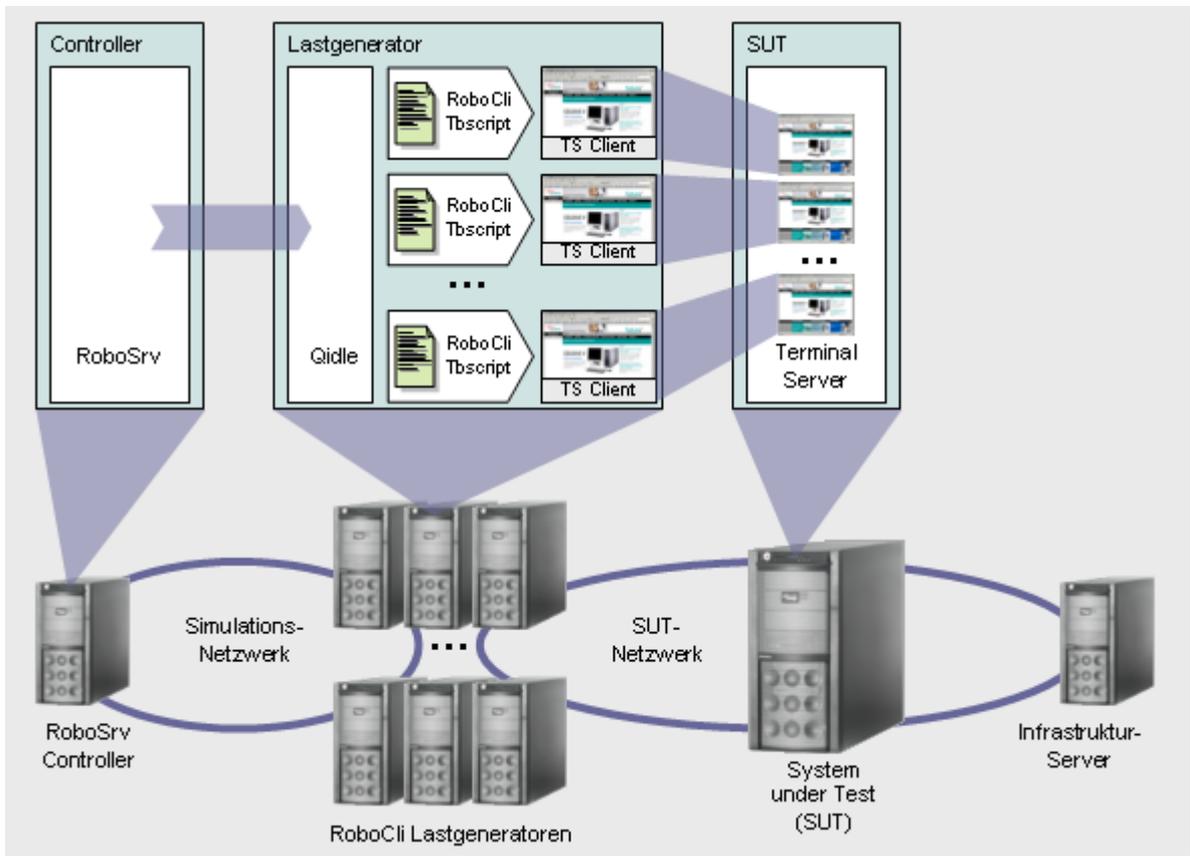
Microsoft Testwerkzeuge und -Skripte

Microsoft hat Testwerkzeuge und –Skripte entwickelt, die auf den Client-Computern eingesetzt werden können, um Benutzersitzungen realitätsnah simulieren zu können. Die Programme zur Durchführung der Tests sind Bestandteil des Windows Server 2003 Resource Kit.

Testwerkzeuge und -Umgebung

»Terminal Services Scalability Planning Tools (TSScaling)« ist eine Sammlung von Werkzeugen, die eine Hilfestellung bei der Planung der Leistungsfähigkeit von Microsoft Windows Server 2003 Terminal Server geben sollen. Sie erlauben es, relativ einfach Last auf einem Server zu simulieren. Hierdurch kann herausgefunden werden, ob eine Umgebung die von ihr erwartete Belastung handhaben kann.

Die Messumgebung besteht aus Hardware und Software-Werkzeugen, wie die nachfolgende Abbildung zeigt:



Das Terminal Server System ist das »**System under Test**«. Als weitere Komponenten befinden sich im Testlabor:

- **Domain Controller und Test Controller:** Dynamic Host Configuration Protocol (DHCP) und DNS Server der Domäne. Verwaltet die 35 Arbeitsstationen inklusive Skriptüberwachung, Softwareverteilung und zentralem Restart der Arbeitsstationen.
- **Arbeitsstationen (35):** Auf jeder der 35 Arbeitsstationen kann eine Vielzahl von Terminal Services Client Sitzungen ablaufen.
- **Mail Server und Web Server:** Dieser Server wird für die Knowledge Worker Tests benötigt.

Die Testprogramme beinhalten die folgenden Automatisierungswerkzeuge:

- **Robosrv.exe** ist das Werkzeug, das die Server-Seite der Lastsimulation steuert. Zusammen sorgen RoboServer und RoboClient für die Automatisierung von Server und Client. RoboServer wird typischerweise auf dem Test Controller installiert und muss laufen, bevor eine Instanz von RoboClient gestartet werden kann. Nachdem je eine Instanz von beiden, RoboServer und RoboClient, läuft, veranlasst RoboServer die RoboClients, die Skripte ablaufen zu lassen, die den Terminal Server in vorgegebenen Intervallen belasten.
- **Robocli.exe** ist das Werkzeug, das die Client-Seite der Lastsimulation steuert. Zusammen sorgen RoboServer und RoboClient für die Automatisierung von Server und Client. RoboClient wird typischerweise auf den Test Client Computern installiert und benötigt einen laufenden RoboServer bevor eine Instanz von RoboClient gestartet werden kann. RoboClient erhält Kommandos von RoboServer, um Skripte ablaufen zu lassen, die den Terminal Server in definierten Intervallen belasten.

Die Testprogramme beinhalten die folgenden Testwerkzeuge:

- **Qidle.exe** wird in einer automatisierten Umgebung eingesetzt und stellt fest, wann eins der momentan laufenden Skripts fehlgeschlagen ist und den Eingriff eines Administrators erfordert. Qidle stellt dieses fest, indem periodisch nachgeschaut wird, ob eine der auf dem Terminal Server angemeldeten Sitzungen für mehr als eine bestimmte Zeit untätig ist. Wenn irgendeine untätige Sitzung existiert, benachrichtigt Qidle den Administrator mit einem Piepton.
- **Tbscript.exe** ist ein Skript-Interpreter der die Client-Seite der Lastsimulation unterstützt. Er führt Visual Basic Scripting Edition Skripte aus und unterstützt bestimmte Erweiterungen, die notwendig sind, um den Terminal Server Client zu steuern. Mit Hilfe dieser Erweiterungen kann der Benutzer Skripte erstellen, die auf dem Client Computer Tastatur- und Mauseingaben erzeugen und, basierend auf den Zeichenfolgen, die die Anwendungen innerhalb der Sitzung ausgeben, die Ausführung synchronisieren.

Test-Skripte

Es wurden zwei Skripte entwickelt, die auf den folgendermaßen definierten Spezifikationen der Gartner Group für den Knowledge Worker und den Data Entry Worker basieren.

Heavy	Knowledge Worker	Sammeln Informationen, werten sie auf und kommunizieren diese in einem Entscheidungshilfeprozess. Die Kosten für Ausfallzeiten sind variabel, aber in hohem Maße erkennbar. Flexible Aufgaben durch Projekte und ad-hoc Bedürfnisse. Treffen eigene Entscheidungen, was gearbeitet wird und wie die Aufgabe bewerkstelligt wird. Beispiele für Arbeitsaufgaben beinhalten Marketing, Projektmanagement, Vertrieb, Desktop-Publishing, Entscheidungshilfe, Data-Mining, Finanzanalyse, ausführendes und überwachendes Management, Design und Authoring.
Light	Data Entry Worker	Geben Daten in Computersysteme ein, beispielsweise Abschriften, Auftragserfassung, Büroarbeit und Fertigung. Darüber hinaus wurde das »Data Entry Worker« Skript im »dedizierten« Modus gemessen, d.h. es wurde kein Windows Explorer für jeden Benutzer gestartet.

Ein detailliertes Ablaufdiagramm, das die Funktionen der Skripte detailliert beschreibt, ist im »[Terminal Server Scaling and Capacity Planning](#)« Dokument enthalten.

Testmethodik

Windows Server 2003, Enterprise Edition und Office 2003 wurden unter Berücksichtigung der Einstellungen, wie im »Appendix B: Terminal Server Settings« des »[Terminal Server Scaling and Capacity Planning](#)« Dokuments beschrieben, installiert. Um zu einem sauberen Zustand aller Komponenten zurückzukommen, wurden Server und Clients automatisch vor jedem Testlauf zurückgesetzt.

Die auf Benutzeraktionen basierenden Antwortzeiten werden verwendet, um zu bestimmen, ob und wann ein Terminal Server überlastet ist. Skripte auf der Client-Seite steuern die Benutzersimulation und zeichnen die Antwortzeiten für eine Reihe von simulierten Benutzeraktionen auf.

Ein Skript enthält mehrere Sequenzen. Eine Sequenz beginnt, wenn das Testskript eine Taste betätigt, die vom Client an eine Anwendung in dessen Session gesendet wird. Dieser Tastendruck resultiert in einer Anzeige der Anwendung. Zum Beispiel öffnet CTRL+F das Menü »File«, das dann den Begriff »Open« anzeigt.

Die Antwortzeit ist die Zeit von dem Tastendruck bis zur Anzeige der Zeichenfolge. Um die Antwortzeit genau messen zu können, wird der Messwert anhand von zwei Ausgangszeitwerten t_1 und t_2 aus einer Referenzzeitquelle berechnet bevor und nachdem ein Tastendruck gesendet wurde. t_1 ist die Zeit wenn der Test Controller die Instruktion zum Client sendet und t_2 repräsentiert die Zeit wenn der Client den Tastendruck zum Server sendet. Eine dritte Messung t_3 wird durchgeführt, wenn die betreffende Zeichenfolge vom Client empfangen worden ist. Die Zeit wird in Millisekunden gemessen. Basierend auf diesen Messwerten wird die Antwortzeit im Intervall ($t_3 - t_2$, $t_3 - t_1$) veranschlagt. In der Praxis ist der Messfehler (die Zeit zwischen t_1 und t_2) geringer als 1 Millisekunde und die Antwortzeiten sind angenähert $t_3 - t_1$.

Für jedes Szenario startet die Test Controller-Arbeitsstation Gruppen von zehn Clientverbindungen auf den Arbeitsstationen mit einem Abstand von 30 Sekunden zwischen den Verbindungen. Nachdem eine Gruppe von zehn Clientverbindungen gestartet wurde, wird ein Stabilisationszeitraum von fünf Minuten abgewartet, in der keine weiteren Sitzungen gestartet werden. Nach diesem Zeitraum startet das Knowledge Worker Skript die vier Anwendungen, die während des Tests genutzt werden, innerhalb von 5 Minuten. Dies verhindert jegliche Beeinflussung zwischen den einzelnen Gruppen der zehn Clientverbindungen.

Mit steigender Benutzeranzahl wird für jede Aktion ein Punkt der Verschlechterung festgelegt, bei dem die Antwortzeiten sich auf einen Wert verschlechtern, der als signifikant erachtet wird.

- Für Aktionen, die anfangs eine Antwortzeit von weniger als 200 ms besitzen, ist der Punkt der Verschlechterung erreicht, wenn die durchschnittliche Antwortzeit über 200 ms liegt und 110% der anfänglichen Antwortzeit überschreitet.
- Für Aktionen, die anfangs eine Antwortzeit von mehr als 200 ms besitzen, ist der Punkt der Verschlechterung erreicht, wenn die durchschnittliche Antwortzeit die anfängliche Antwortzeit um 10% überschreitet.

Dieses Kriterium basiert auf der Annahme, dass ein Benutzer eine Verschlechterung der Antwortzeiten bei Zeiten unter 200 ms nicht bemerkt.

Dieser Test unterstützt das vorherige Testsystem, bei dem der Grenzwert der Systemlast durch ein so genanntes »Canary« Skript ermittelt wurde, das in der stabilen Phase zwischen den Logon-Gruppen ablief. Das Canary Skript wurde zum Ablauf gebracht, bevor sich ein Benutzer am System angemeldet hatte, und die Zeit, die das Skript für einen Durchlauf benötigte (verstrichene Zeit), wird aufgezeichnet. Diese verstrichene Zeit ist die Baseline, sie wird als Grundlinie für die Antwortrate in einer gegebenen Serverkonfiguration erachtet. Mit dieser Methode wurde das Erreichen der maximalen Last festgelegt, wenn die Gesamtzeit für einen Durchlauf des Canary Skripts um 10% höher als der Ausgangswert war. Die Antwortzeit-Methode wird jedoch als genauer angesehen, da sie als Schlüsselparameter die momentane Benutzer Erfahrung misst, auch den Einfluss der Logon-Phasen mit berücksichtigt und mehr Daten für eine Entscheidungsfindung bereitstellt. Die »Canary« Skript-Methode kann aber in Konfigurationen effizienter sein, die nur eine kleine Anzahl Benutzer betreiben und bei denen die Antwortzeit-Methode keine genügende Menge an Antwortzeiten liefert.

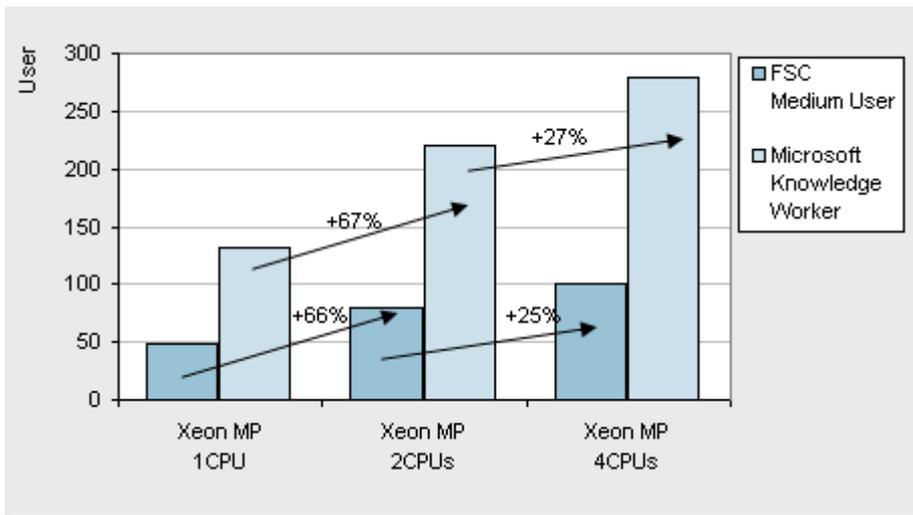
Eine detaillierte Beschreibung findet man im »[Terminal Server Scaling and Capacity Planning](#)« Dokument.

Ergebnisse von Fujitsu Siemens Computers und Microsoft

Um die Werte des Microsoft »Terminal Server Scaling and Capacity Planning« Tools und des Fujitsu Siemens Computers Werkzeugs »T4US« zu vergleichen, wurden sowohl bei Microsoft als auch bei Fujitsu Siemens Computers eine Reihe von Messungen aufgesetzt, bei denen die gleichen PRIMERGY Modelle verwendet wurden:

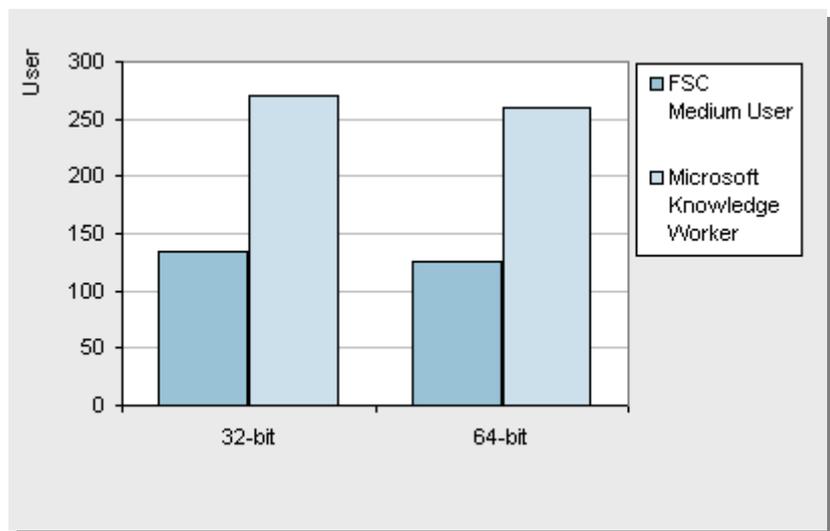
- Ein PRIMERGY RX300 S2 Server mit zwei Prozessoren mit 3.6 GHz, 1 MB SLC und 8 GB RAM. Hyper-Threading war eingeschaltet und Microsoft Office 2003 wurde als Officeanwendung verwendet. Als Betriebssystem wurde Windows 2003 Enterprise Edition in den 32-bit und 64-bit Versionen verwendet.
- Ein PRIMERGY TX600 Server mit bis zu vier Prozessoren mit bis zu 2.8 GHz, 2 MB SLC und 8 GB RAM. Hyper-Threading war eingeschaltet und Microsoft Office 2003 wurde als Officeanwendung verwendet. Als Betriebssystem wurde die 32-bit Version von Windows 2003 Enterprise Edition verwendet.

Microsoft arbeitet mit dem »Terminal Server Scaling and Capacity Planning« Tool, während Fujitsu Siemens Computers das eigene »T4US« Tool einsetzt. Das Medium Lastprofil, das mit T4US verwendet wird, wird im Kapitel »[Lastprofil](#)« im Detail erklärt. Das Benutzerprofil des »Terminal Server Scaling and Capacity Planning« Tools von Microsoft wurde weiter [oben](#) in diesem Abschnitt beschrieben.



Auf dem PRIMERGY TX600 System wurde eine Messreihe aufgesetzt. Wie nebenstehende Grafik zeigt, ist der CPU-Skalierungsfaktor von einer CPU auf zwei CPUs immer 1,6, unabhängig davon, welches Lastprofil benutzt wird. Die Skalierung von zwei auf vier Prozessoren ist für beide Messszenarien 1,25. Die absoluten Benutzerzahlen sind jedoch unterschiedlich.

Eine weitere Messreihe wurde auf einem PRIMERGY RX300 S2 System von Fujitsu Siemens Computers durchgeführt, das sowohl 32-bit als auch 64-bit Betriebssysteme unterstützt. Auf dieser Hardware wurden auch Tests von Microsoft mit 32-bit Windows Server 2003 Betriebssystemen (mit Service Pack 1) und mit der 64-bit Version von Windows Server 2003 durchgeführt. Bei allen Messungen war das System mit genügend Hauptspeicher ausgestattet, um sicherzustellen, dass der Arbeitsspeicher bei diesem Vergleich keinen Engpass darstellt. Wenn diese absoluten Ergebnisse mit den Zahlen des Fujitsu Siemens Computers T4US Medium Lastprofils verglichen werden, wie in dem Diagramm dargestellt, erkennt man, dass die Ergebnisse im gleichen Rahmen skalieren, obwohl sich die absoluten Zahlen unterscheiden.



Mit dem T4US Medium Lastprofil von Fujitsu Siemens Computers können generell weniger Benutzer auf der gleichen Hardware bedient werden als mit dem Lastprofil von Microsoft. Da die vergleichenden Messungen auf der gleichen Hardware durchgeführt wurden, hängt die Benutzeranzahl allein von dem verwendeten Lastprofil ab.

In unseren früheren Sizing Guides wurden ebenfalls Lastprofile verwendet, die höhere Benutzeranzahlen, wie auch Microsofts »Terminal Server Scaling and Capacity Planning« Tool, lieferten, jedoch haben unsere Kunden berichtet, dass diese Werte viel zu hoch waren. Daher wurde das »T4US« Medium Lastprofil mit dem Ziel entwickelt, vergleichsweise niedrigere Benutzerzahlen zu liefern. In unseren neuen Messreihen haben wir dem Rechnung getragen und können daher davon ausgehen, dass die ermittelten Benutzerzahlen denen aus realen Produktionsumgebungen nahe kommen.

Das Microsoft »Terminal Server Scaling and Capacity Planning« Tool liefert eine Benutzeranzahl, die tatsächlich das absolute Maximum der Benutzer darstellt, die in einem bestimmten Szenario unterstützt werden. Diese Zahlen müssen jedoch noch angepasst werden, um einen annehmbaren Belastungswert für reale Produktionsumgebungen zu erhalten.

Wesentlich ist, dass diese verschiedenen Messwerkzeuge vergleichbare relative Zahlen liefern. Der Vergleich von Microsofts »Terminal Server Scaling and Capacity Planning« Tool und Fujitsu Siemens Computers »T4US« Werkzeug zeigt, dass dies gegeben ist und dass beide Werkzeuge dazu benutzt werden können, um die Skalierung von Terminal Server Systemen zu untersuchen.

Resümee

Die Leistungsfähigkeit des Terminal Servers ist durch CPU-Leistung und Hauptspeicher bestimmt.

Den Speicherausbau kann man recht einfach anhand der Formel

$$Memory = Memory_{OS} + \#_{Client} \cdot Memory_{App}$$

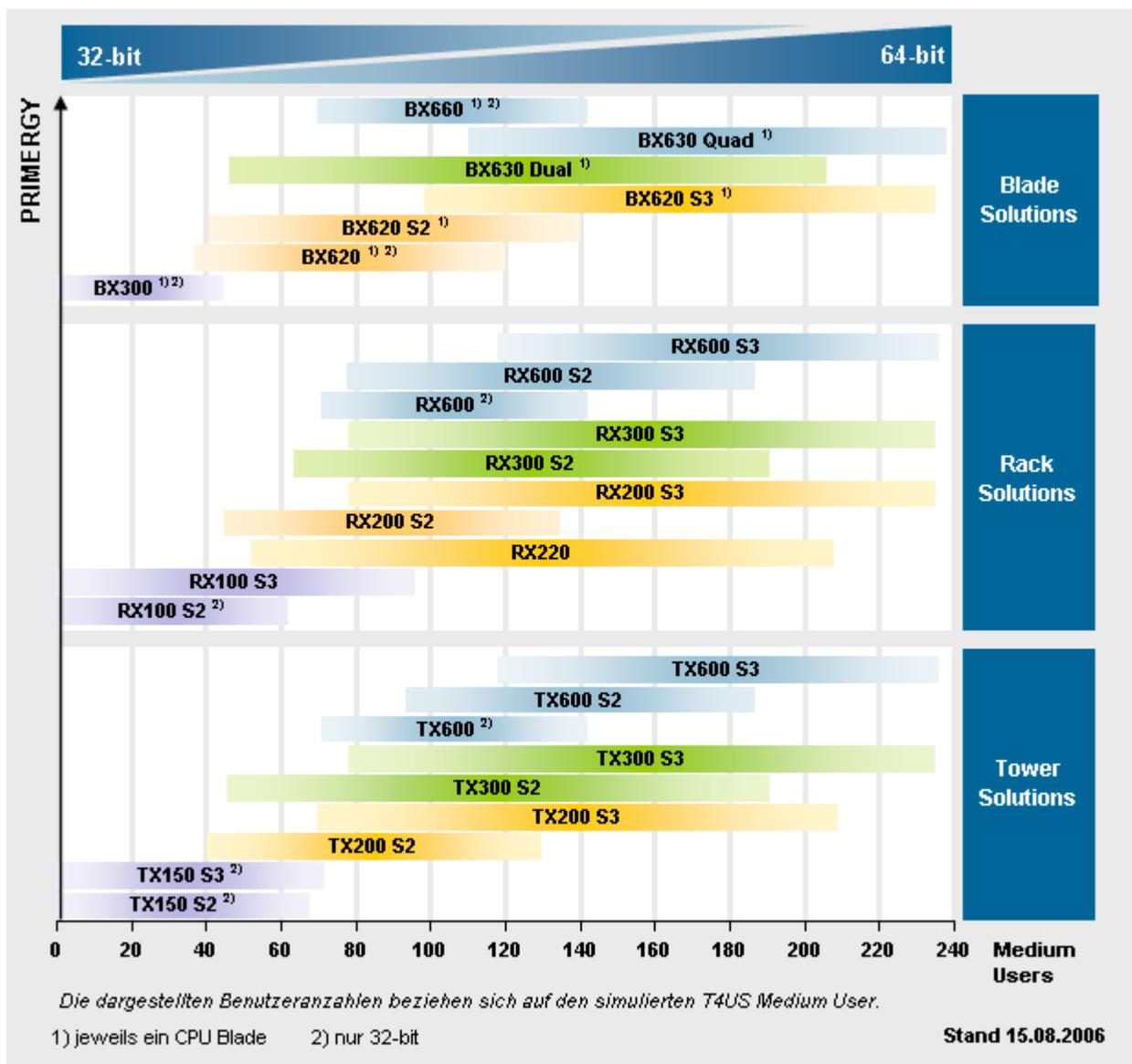
abschätzen. Werden verschiedene Applikationen gleichzeitig verwendet, so ist die Summe des Speicherbedarfs aller gleichzeitig verwendeter Applikationen zu bilden.

Die Anzahl Benutzer bei einer vorgegebenen Speichermenge kann mit folgender Formel berechnet werden:

$$\#_{Client} = \frac{Memory - Memory_{OS}}{Memory_{App}}$$

Für die CPU-Leistung gilt leider keine so einfache Formel. Bei einem Medium oder Heavy User stellt im Allgemeinen die CPU-Leistung den begrenzenden Faktor dar. Für eine Klasse von Light Usern ist die maximale Benutzerzahl eher durch andere Systemressourcen begrenzt.

Die Vielzahl an Prozessoren, mit denen jedes PRIMERGY-Modell ausgestattet werden kann, lässt bereits erahnen, dass es nicht eine bestimmte Anzahl Benutzer gibt, die ein PRIMERGY-Modell bedienen kann, sondern dass jedes PRIMERGY-Modell eine gewisse Bandbreite abdeckt. Auch gibt es keine scharfe Grenze, wo die Leistung eines Modells endet und die des nächst leistungsfähigeren beginnt. Vielmehr gibt es Überlappungen zwischen den Systemen. Die folgenden durch unsere Messreihen gewonnenen Ergebnisse können also nur einen Eindruck für den Leistungsbereich vermitteln. Folgende Grafik zeigt einen Vergleich der PRIMERGY Systeme untereinander.



In dieser Darstellung wird die höchste erreichbare Benutzeranzahl jedes PRIMERGY Systems als Maximalwert verwendet, die mit einer optimalen Hardwarekonfiguration und dem jeweils besten Betriebssystem (32-bit oder 64-bit) erreicht wurde.

Erwartungsgemäß zeigen sich die Vorteile der 64-bit Architektur am oberen Ende des Leistungsspektrums, wo die heutigen Limitierungen der 32-bit Architektur hinsichtlich des Adressraums überwunden werden und so ein Mehr an Speicher und Rechenleistung ausgenutzt werden kann. Bedingt durch die Multi-Core-Architektur moderner Prozessoren stoßen heute bereits Systeme mit 2 CPU-Sockeln in diesen Leistungsbereich vor.

Bei den Terminal Server-Messungen werden weit mehr als 100 Benutzer auf einem Server erreicht. Der erreichte Maximalwert hängt, wie vorher bereits diskutiert, von der Messmethode und von dem gewählten Lastprofil ab. In der Realität ist eine solch hohe Anzahl Benutzer jedoch meist nicht erzielbar, vielmehr liegen hier die Benutzerzahlen deutlich darunter. Diese Ergebnisse müssen an der Realität gespiegelt werden und auf das kundenspezifische Benutzerverhalten relativiert werden. Ob dabei die Ergebnisse durch 10 oder durch 4 geteilt werden müssen oder ob man eine fixe Anzahl, zum Beispiel 20 Benutzer pro Prozessorkern, einplant; diese Empfehlungen können in jedem Fall nur Faustformeln sein.

Wenn Sie einen Terminal Server oder eine Terminal Server-Farm planen, sollten Sie sich die Zeit nehmen, das Benutzerverhalten vorher genau zu analysieren.

- Welche Anwendungen müssen generell über den Terminal Server zur Verfügung gestellt werden?
- Welcher Benutzer nutzt wann und wie oft welche Anwendung?
- Wie intensiv wird die Anwendung verwendet?
- Welche Bildschirmauflösung und Farbtiefe verwendet der Client?
- Welches Netzwerk und Protokoll wird verwendet?
- Welche Antwortzeiten werden erwartet?

Bei größeren Konfigurationen sollte auf eine Pilotphase unter realen Bedingungen nicht verzichtet werden.

Literatur

[L1]	Allgemeine Informationen zu Produkten von Fujitsu Siemens Computers http://www.fujitsu-siemens.de
[L2]	Allgemeine Informationen zur PRIMERGY Produktfamilie http://www.PRIMERGY.de
[L3]	PRIMERGY Benchmarks - Performance Reports und Sizing Guides http://www.fujitsu-siemens.de/products/standard_servers/primergy_bov.html
[L4]	Terminal Server Sizing Guide - 64-bit Technologie http://extranet.fujitsu-siemens.com/vil/pc/vil/primergy/performance/sizing/terminal_server_sizing_guide_-_x64_technologie_de.pdf
[L5]	PRIMERGY Performance Lab Datenbank mit allen Benchmark- und Performance-Ergebnissen der PRIMERGY im Intranet http://ppl.pdb.fsc.net
[L6]	Microsoft Windows 2003 und Terminal Server http://www.microsoft.com/terminalserver
[L7]	Windows Server 2003 Terminal Server Capacity and Scaling http://www.microsoft.com/windowsserver2003/techinfo/overview/tsscaling.msp
[L8]	Microsoft Windows Server 2003 Terminaldienste Bernhard Tritsch, Microsoft Press, ISBN 3-86063-656-1
[L9]	Citrix http://www.citrix.de
[L10]	Gartner Group http://www.gartner.com/

Kontakt

PRIMERGY Performance und Benchmarks
<mailto:primergy.benchmark@fujitsu-siemens.com>

PRIMERGY Produkt Marketing
<mailto:PRIMERGY-PM@fujitsu-siemens.com>