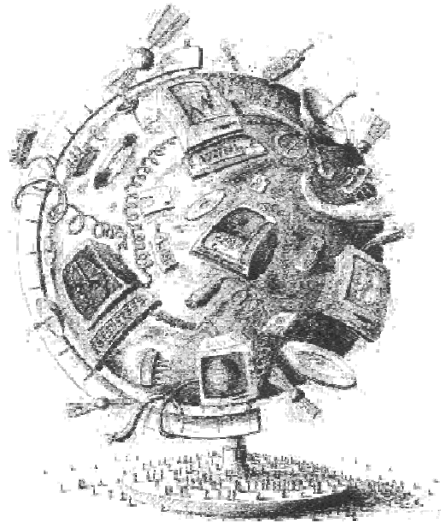


Verteilte Systeme



*Friedemann
Mattern
ETH Zürich*

Prüfungsrelevant ist der Inhalt der Vorlesung, nicht alleine der Text dieser Foliensammlung! Herbst 2011

ETH Eidgenössische
Technische Hochschule
Zürich

© F. Mattern 2011

Wer bin ich? Wer sind wir?



Prof. **Friedemann Mattern**

+ 15 Assistentinnen
und Assistenten



Fachgebiet „Verteilte Systeme“
im Departement Informatik,
Institut für Pervasive Computing



Wer bin ich? Wer sind wir?



Prof. **Friedemann Mattern**

+ 15 Assistentinnen
und Assistenten

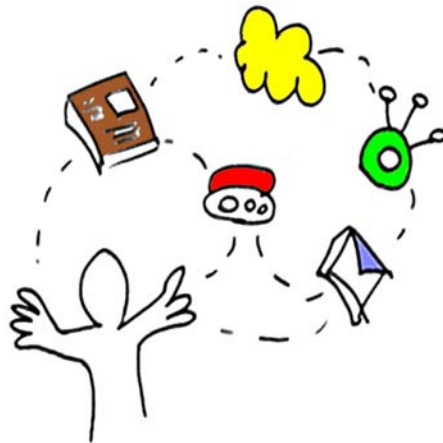
Matthias Kovatsch

Ansprechperson für or-
ganisatorische Aspekte
kovatsch@inf.ethz.ch



Fachgebiet „Verteilte Systeme“
im Departement Informatik,
Institut für Pervasive Computing

Mit was beschäftigen wir uns?



- Infrastruktur für verteilte Systeme
- Internet der Dinge
- Ubiquitous Computing
- Sensornetze
- Verteilte Anwendungen und Algorithmen

Mehr zu uns:
www.vs.inf.ethz.ch

Organisatorisches zur Vorlesung

- Format: 6G+1A: **Vorlesung** und **Praktikum** integriert
 - Mo 9:15 - 12:00, Fr 9:15 - 12:00, jew. NO C 6
 - Praktikum ist inhaltlich *komplementär zur Vorlesung* (mobile Kommunikationsplattformen: Android, HTC Desire)
 - Gelegentliche *Assistentenstunden* (zu den "Vorlesungsterminen") zur Besprechung der Praktikumsaufgaben und Vertiefung des Stoffes
 - Gelegentliche *Denkaufgaben* (ohne Lösung...) in der Vorlesung
- Sinnvolle **Vorkenntnisse** (Grundlagen)
 - 4 Semester der Bachelorstufe Informatik
 - Grundkenntnisse Computernetze und Betriebssysteme (z.B. Prozessbegriff, Synchronisation)
 - UNIX, Java ist hilfreich

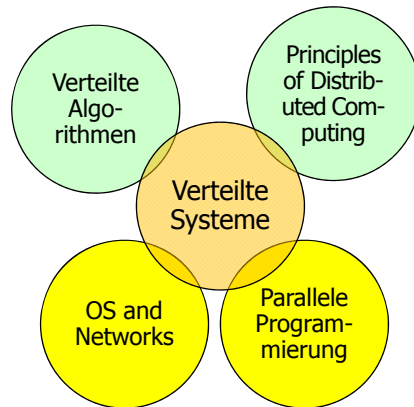
7

Organisatorisches (2)

- **Folienkopien** jeweils einige Tage nach der Vorlesung
 - im pdf-Format bei www.vs.inf.ethz.ch/edu
- Vorlesung ab November: Prof. **Roger Wattenhofer**
- **Prüfung** schriftlich
 - bewertete Praktikumsaufgaben gehen in die Prüfungsnote ein

8

Einordnung der Vorlesung



- „Verteilte Systeme“ ist ein **Querschnittsthema**
- Gewisse Überschneidungen mit anderen Vorlesungen unvermeidlich

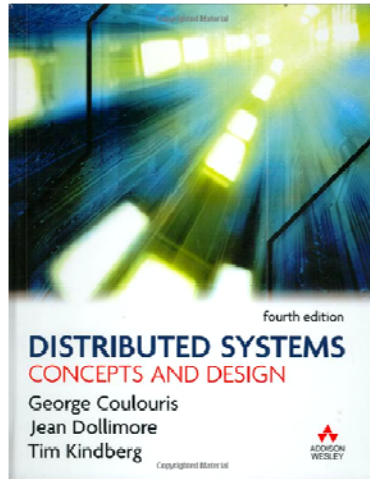
9

Thematisch verwandte Veranstaltungen (Masterstufe)

- Ubiquitous Computing
 - Enterprise Application Integration - Middleware
 - Web Services and Service Oriented Architectures
 - Verteilte Algorithmen
 - Principles of Distributed Computing
-
- Einschlägige Seminare
 - Praktikum („Lab“)
 - Semester- und Masterarbeit

10

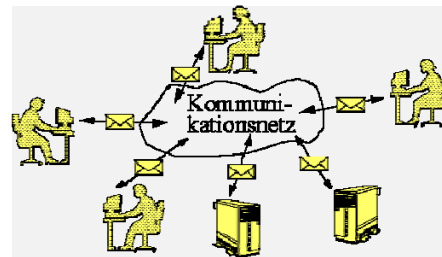
Literatur



- G. Coulouris, J. Dollimore, T. Kindberg: Distributed Systems: Concepts and Design (4th ed.). Addison-Wesley, 2005
- A. Tanenbaum, M. van Steen: Distributed Systems: Principles and Paradigm (2nd ed.). Prentice-Hall, 2007
- Oliver Haase: Kommunikation in verteilten Anwendungen (2. Auflage). R. Oldenbourg Verlag, 2008

“Verteiltes System” – zwei Definitionen

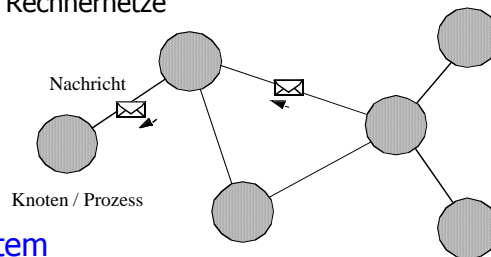
- A distributed computing system consists of multiple autonomous processors that do not share primary memory, but cooperate by sending messages over a communication network. -- *H. Bal*



- A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable. -- *Leslie Lamport*
 - welche Problemaspekte stecken hinter Lamports Charakterisierung?

"Verteiltes System"

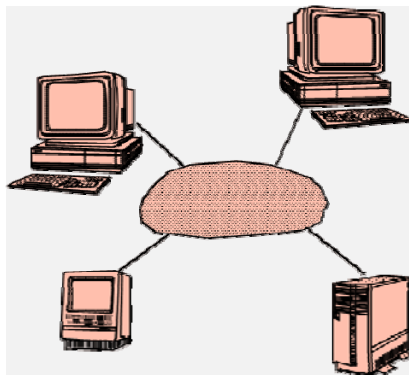
- **Physisch verteiltes System**
 - Mehrrechnersystem ... Rechnernetze



- **Logisch verteiltes System**
 - Prozesse (Objekte, Agenten)
 - Verteilung des Zustandes (keine globale Sicht)
 - Keine gemeinsame Zeit (globale, genaue Uhr)

13

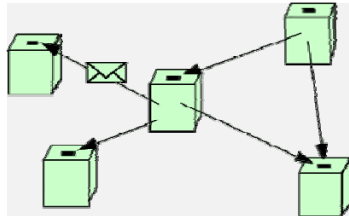
Sichten verteilter Systeme (1)



- **Computernetz mit "Rechenknoten", z.B.**
 - Compute-Cluster
 - Local Area Network
 - Internet
- **Relevante Aspekte:**
 - Routing, Adressierung,...

Zunehmende Abstraktion →

Sichten verteilter Systeme (2)



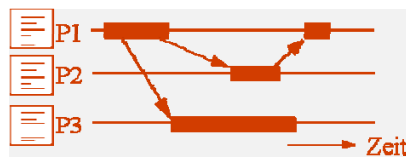
kommunizierende Prozesse,
kooperierende Objekte

- **Objekte** in Betriebssystemen, Middleware, Programmiersprachen
- "Programmierersicht"
 - z.B. Client mit API zu Server

Zunehmende Abstraktion →

15

Sichten verteilter Systeme (3)



- **Algorithmen- und Protokollebene**
 - Aktionen, Ereignisfolgen
 - Konsistenz, Korrektheit

- Man kann verteilte Systeme auf **verschiedenen Abstraktionsstufen** betrachten
- Es sind dabei jeweils **unterschiedliche Aspekte** relevant und interessant

16

Die verteilte Welt



Auch die "reale Welt" ist ein **verteiltes System**:

- viele gleichzeitige ("parallele") Aktivitäten
- exakte globale **Zeit** nicht erfahrbar / vorhanden
- keine konsistente Sicht des **Gesamtzustandes**
- Kooperation durch explizite **Kommunikation**
- **Ursache** und **Wirkung** zeitlich (und räumlich) getrennt

Warum verteilte Systeme?

- **Es gibt inhärent geographisch verteilte Systeme**
 - z.B. Zweigstellennetz einer Bank, Steuerung einer Fabrik (Zusammenführen / Verteilen von Information)
- **Electronic commerce**
 - kooperative Informationsverarbeitung räumlich getrennter Institutionen (z.B. Reisebüros, Kreditkarten,...)
- **Mensch-Mensch-Telekommunikation**
 - E-Mail, Diskussionsforen, Blogs, digitale soz. Netze, IP-Telefonie,...
- **Globalisierung von Diensten**
 - Skaleneffekte, Outsourcing,...

Wirtschaftliche Aspekte

- Outsourcing von Diensten, Verlagerung in eine „Cloud“, kann günstiger als klassische Lösung sein
- Compute-Cluster manchmal besseres Preis-Leistungs-verhältnis als Hochleistungs-computer

Verteilte Systeme als "Verbunde"

- Verteilte Systeme **verbinden** räumlich (oder logisch) getrennte Komponenten zu einem bestimmten **Zweck**
-
- **Systemverbund**
 - gemeinsame Nutzung von Betriebsmitteln, Geräten,...
 - einfache inkrementelle Erweiterbarkeit
 - **Funktionsverbund**
 - Kooperation bzgl. Nutzung jeweils spezifischer Eigenschaften
 - **Lastverbund**
 - Zusammenfassung der Kapazitäten
 - **Datenverbund**
 - allgemeine Bereitstellung von Daten
 - **Überlebensverbund**
 - Redundanz durch Replikation

19

Historische Entwicklung (1)

- **Rechner-zu-Rechner-Kommunikation**
 - Zugriff auf entfernte Daten ("Datenfernübertragung", DFÜ)
 - dezentrale Informationsverarbeitung war zunächst ökonomisch nicht sinnvoll (zu teuer, Fachpersonal nötig)
 - Master-Slave-Beziehung ("Remote Job Entry", Terminals)
- **ARPA-Netz (Prototyp des Internet)**
 - "symmetrische" Kommunikationsbeziehung ("peer to peer")
 - file transfer, remote login, E-Mail
 - Internet-Protokollfamilie (TCP/IP,...)

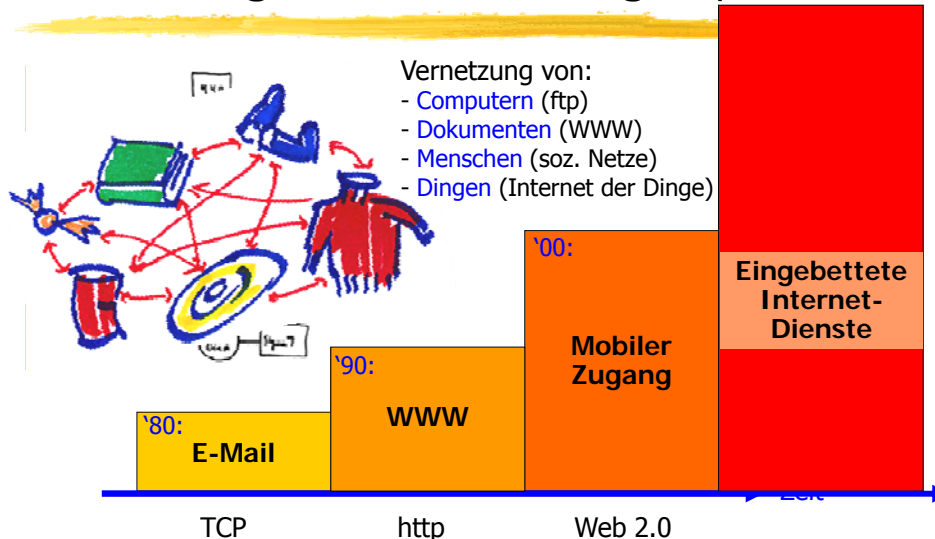
20

Historische Entwicklung (2)

- **Workstation-Netze (LAN)**
 - bahnbrechende, frühe Ideen bei XEROX-PARC (XEROX-Star als erste Workstation, Desktop-Benutzerinterface, Ethernet, RPC, verteilte Dateisysteme,...)
- **Kommerzielle Pionierprojekte als Treiber**
 - z.B. Reservierungssysteme, Banken, Kreditkarten
- **Web / Internet als Plattform**
 - für electronic commerce etc.
 - web services
 - neue, darauf aufbauende Dienste
- **Mobile Geräte (z.B. Smartphones)**
- **Internet der Dinge**

21

Änderung der Vernetzungs"qualität"



Historie von Konzepten

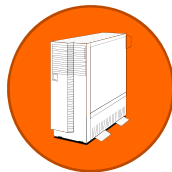
- **Concurrency, Synchronisation**
 - war bereits klassisches Thema bei Datenbanken und Betriebssystemen
- **Programmiersprachen**
 - z.B. „kommunizierende“ Objekte
- **Parallele und verteilte Algorithmen**
- **Semantik von Kooperation / Kommunikation**
 - mathematische Modelle für Verteiltheit (z.B. CCS, Petri-Netze)
- **Abstraktionsprinzipien**
 - Schichten, Dienstprimitive,...
- **Verständnis grundlegender Phänomene der Verteiltheit**
 - Konsistenz, Zeit, Zustand,...

Entwicklung „guter“ Konzepte, Modelle, Abstraktionen etc. zum Verständnis der Phänomene **dauert oft lange** (notwendige Ordnung und Sichtung des verfügbaren Gedankenguts)

Diese sind jedoch für die Lösung praktischer Probleme **hilfreich**, oft sogar **notwendig!**

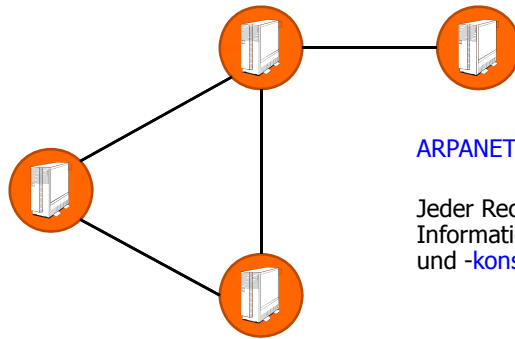
Architekturen verteilter Systeme

- Zu Anfang waren Systeme **monolithisch**



- Nicht verteilt / vernetzt
- **Mainframes**
- **Terminals** als angeschlossene „Datensichtgeräte“ („Datenendgerät“: Fernschreiber, ASCII)

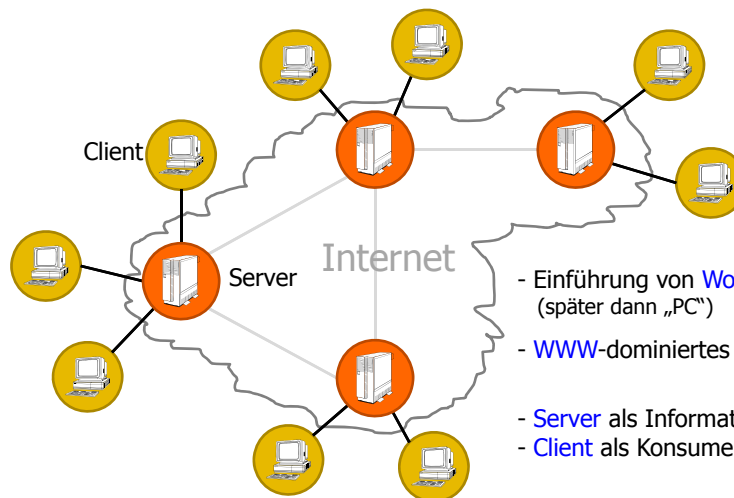
Architekturen verteilter Systeme: Peer-to-Peer



ARPANET 1969

Jeder Rechner **gleichzeitig**
Informationsanbieter
und -konsument

Architekturen verteilter Systeme: Client-Server



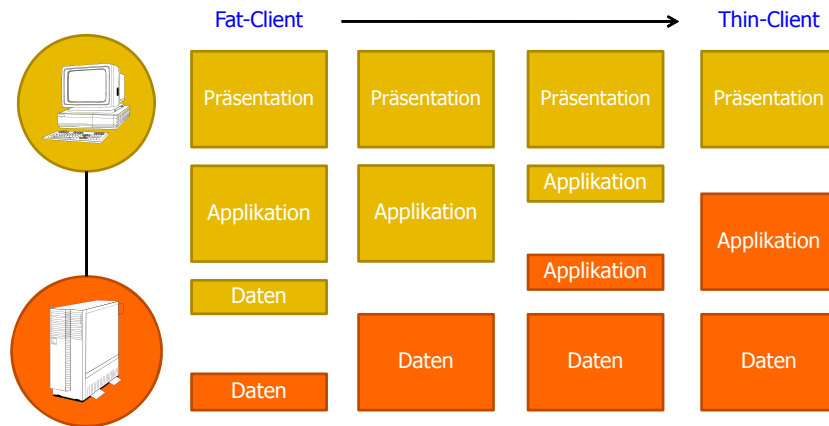
- Einführung von **Workstations**
(später dann „PC“)

- **WWW**-dominiertes Internet

- **Server** als Informationsanbieter

- **Client** als Konsument

Architekturen verteilter Systeme: Fat- und Thin-Client



Architekturen verteilter Systeme: 3-Tier



Architekturen verteilter Systeme: Multi-Tier

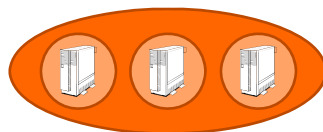


Weitere Schichten sowie mehrere physikalische Einheiten pro Schicht („Compute-Cluster“) erhöhen die **Skalierbarkeit** und **Flexibilität**

Mehrere Webserver ermöglichen z.B. **Lastverteilung**

Verteilte Datenbanken in der Datenhaltungsschicht bietet Sicherheit durch Replikation und hohen Durchsatz

Architekturen verteilter Systeme: Compute-Cluster



Vernetzung kompletter Einzelrechner

Räumlich konzentriert (wenige Meter)

Sehr schnelles Verbindungsnetz

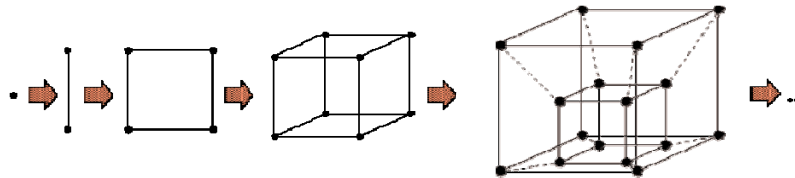
Es gibt diverse **Netztopologien**, um die Einzelrechner (*als Knoten in einem Graphen*) miteinander zu verbinden – diese sind unterschiedlich hinsichtlich

- Skalierbarkeit der Topologie
- Routingkomplexität
- Gesamtzahl der Einzelverbindungen
- maximale bzw. durchschnittliche Entfernung zweier Knoten
- Anzahl der Nachbarn eines Knotens
- Zahl der alternativ bzw. parallel verfügbaren Wege
- ...

Bestimmt die die Kosten und die Leistungsfähigkeit eines Systems

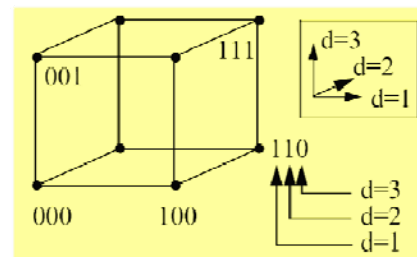
Beispiel: Hypercube-Verbindungstopologie

- **Würfel der Dimension d**
 - Vorteil: einfaches Routing, kurze Weglängen
 - Nachteil: Viele Einzelverbindungen ($O(n \log n)$ bei n Knoten)
- **Rekursives Konstruktionsprinzip**
 - Hypercube der Dimension 0: Einzelrechner
 - Hypercube der Dimension $d+1$: „Nimm zwei Würfel der Dimension d und verbinde korrespondierende Ecken“

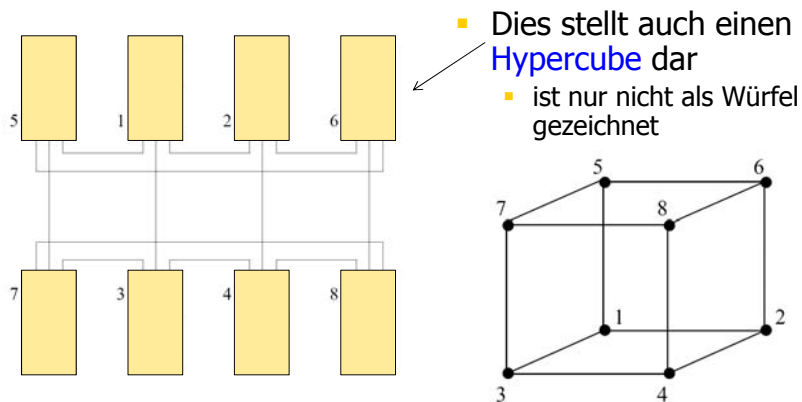


Routing beim Hypercube

- Knoten **systematisch nummerieren** (entspr. rekursivem Aufbau)
- Zieladresse **bitweise xor** mit Absenderadresse
- Wo sich eine "1" findet, in diese Dimension muss gewechselt werden
- Maximale Weglänge: d ; durchschnittliche Weglänge = $d/2$ (Induktionsbeweis als einfache Übung)



Beispiel: Hypercube-Verbindungstopologie (2)

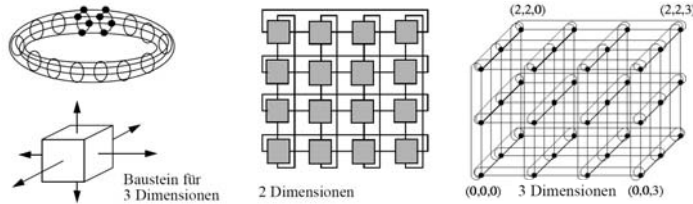


Resümee (1)

- **Verteilte Systeme:** Begriff, Sichtweisen, Eigenschaften,...
- **Warum** verteilte Systeme?
 - Kooperation von a-priori geographisch verteilten Einheiten
 - Verteilte Systeme als "Verbund"
- **Historische Entwicklung** von Systemen und Konzepten
- **Architekturvarianten**
 - Peer-to-Peer
 - Client-Server
 - Fat-Client vs. Thin Client
 - 3-Tier und Multi-Tier
 - Compute-Cluster (Beispiel für Verbindungstopologie: Hypercube und Torus)

Eine andere Verbindungstopologie: Der d-dimensionale Torus

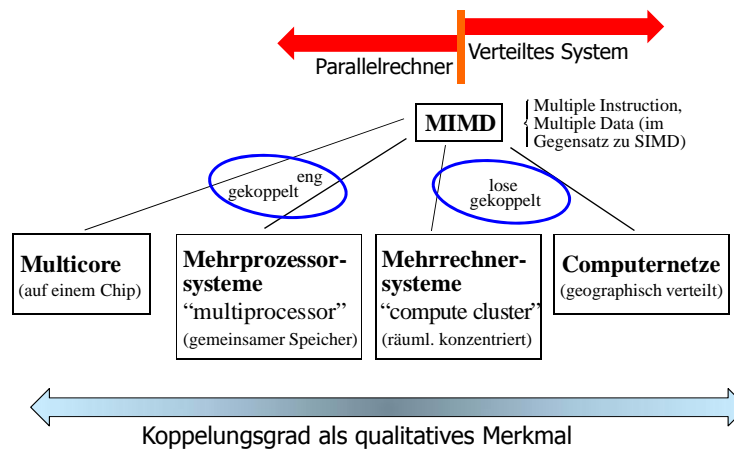
- Gitter in d Dimensionen mit "wrap-around"



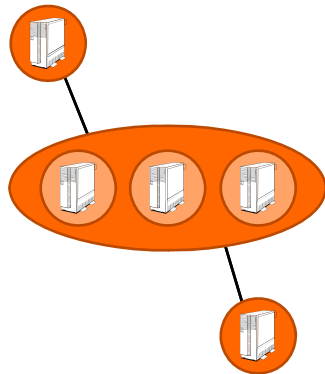
- Rekursives Konstruktionsprinzip:** Nimm w_d gleiche Exemplare der Dimension $d-1$ und verbinde korrespondierende Elemente zu einem Ring
 - Sonderfall **Ring:** $d = 1$
 - Sonderfall **Hypercube:** d -dimensionaler Torus mit $w_i = 2$ für alle Dimensionen i

Es gibt noch einige andere sinnvolle Verbindungstopologien (auf die wir nicht eingehen)

Parallelrechner ↔ verteiltes System



Architekturen verteilter Systeme: Service-Oriented Architecture (SOA)



Eine **Unterteilung** der Applikation in einzelne, unabhängige Abläufe innerhalb eines **Geschäftsprozesses** erhöht die Flexibilität weiter

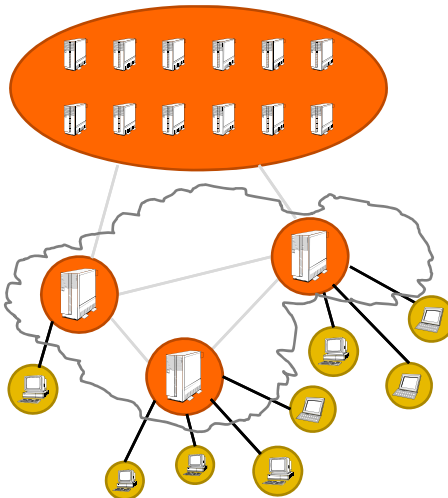
Lose Kopplung zwischen Services über Nachrichten und events (statt RPC)

Services können bei Änderungen der Prozesse **einfach neu zusammengestellt** werden („development by composition“)

Services können auch von **externen Anbietern** bezogen werden

Oft in Zusammenhang mit **Web-Services**

Architekturen verteilter Systeme: **Cloud-Computing**

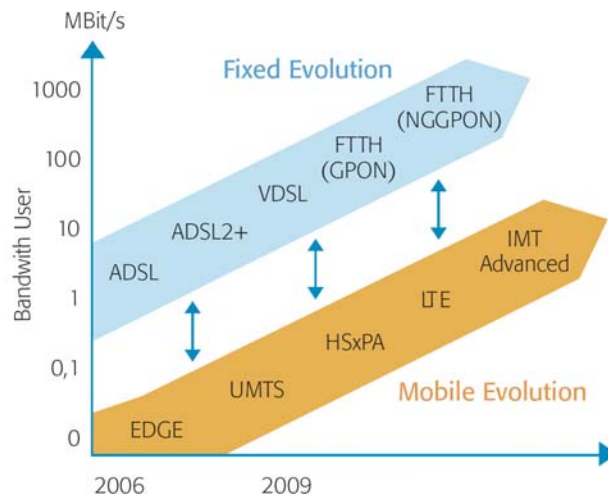


Massive **Bündelung der Rechenleistung** an zentraler Stelle

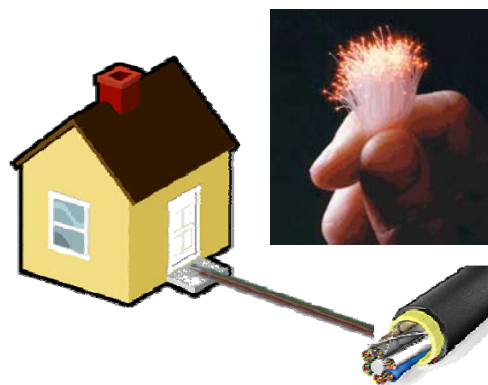
Outsourcen von Applikationen in die Cloud

Internet im Wesentlichen nur noch als **Vermittlungsinstanz**

Motivierender Trend: Stetige Erhöhung der Bandbreite für Endnutzer



Hochgeschwindigkeit ins Haus



- Konvergenz TV, Telekommunikation und Internet
- Technologiewechsel → erhebliche Investitionen
 - wirtschaftliche Faktoren und Bedingungen

- **Telefondraht** → Internet → TV
 - **TV-Kabel** → Telefon → Internet
- } → **Glasfaser**
- „Tripleplay“ (Sprache, Daten, Video)

Cloud-Computing



E-Mail wird beim Provider gespeichert



Cloud-Computing



Fotos werden bei flickr gespeichert

Cloud-Computing



Videos bei **You Tube**

Cloud-Computing



Private Dokumente werden bei einem **Storage Provider** abgelegt

Das **Tagebuch** wird öffentlich „im Netz“ als **Blog** geführt

Cloud-Computing



Informieren tut man sich im Netz

Google



Cloud-Computing



Vernetzen tut man sich bei „social networks“ oder „digital communities“ im Netz



Cloud-Computing



Plattformen im Netz
nutzt man zum

- Kaufen
- Spielen
- Kommunizieren

- ... **Google**
Romance BETA

Cloud-Computing



Vorteile für Nutzer:

- von überall zugreifbar
- keine Datensicherung
- keine Softwarepflege

Kein PC, sondern
billiges Web-Terminal,
Smartphone etc.

Wie Wasserleitungen einst den eigenen Brunnen überflüssig machten...

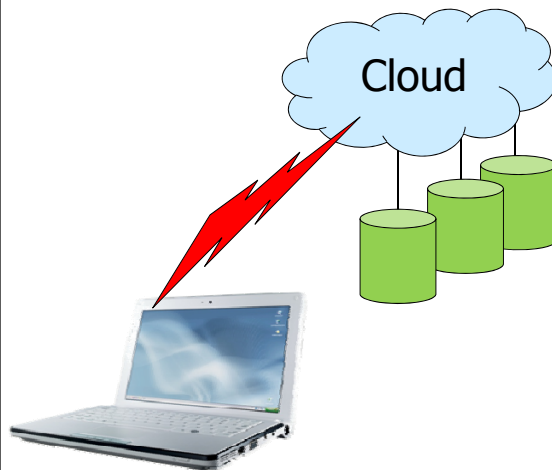
Cloud-Computing



Voraussetzungen?

- Überall Breitband (fest & mobil)
- Netz-Verlässlichkeit (Versorgungssicherheit, Datenschutz,...)
- Wirtschaftlichkeit

Verlässlichkeit?



Voraussetzungen?

- Überall Breitband (fest & mobil)
- Netz-Verlässlichkeit (Versorgungssicherheit, Datenschutz,...)
- Wirtschaftlichkeit

Xdrive®

The Xdrive service is closed.

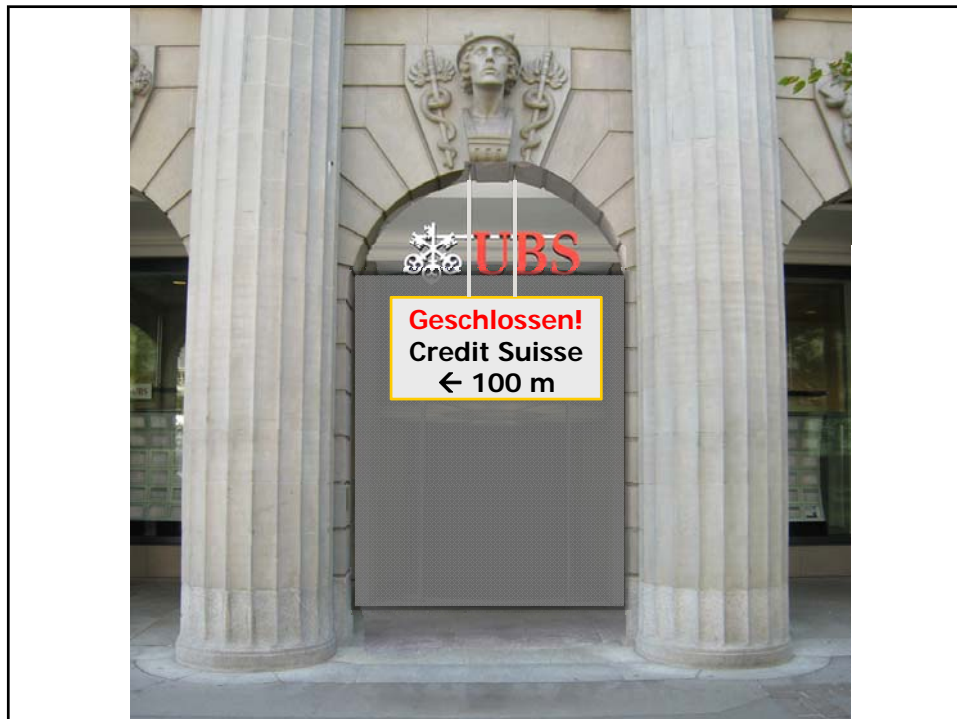
Thank you for having been an Xdrive user.



elephantdrive
ONLINE BACKUP, STORAGE, & SHARING

Start your FREE trial now!

CLICK HERE



Cloud-Computing



Plattformen:

- Wer betreibt sie? Wo?
- Wer verdient daran?
- Wer bestimmt?
- Wer kontrolliert?
- Welche Nationen profitieren davon?

Beispiel: Google-Datenzentren



- Jedes Datenzentrum hat **10 000 – 100 000 Computer**
- Kostet über **500 Mio \$** (Bau, Infrastruktur, Computer)
- Verbraucht **50 – 100 MW Energie** (Strom, Kühlung)
- Neben Google weitere (z.B. Amazon, Microsoft, Ebay,...)

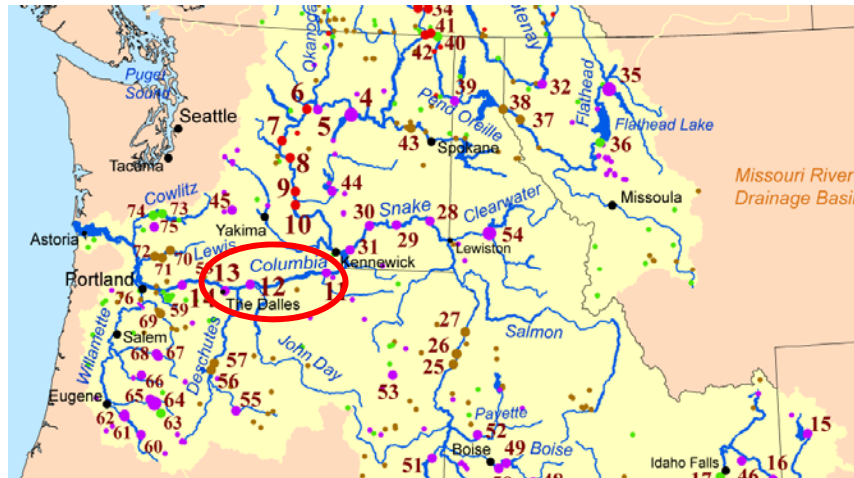
Google Data Center Groningen



Google Data Center Columbia River

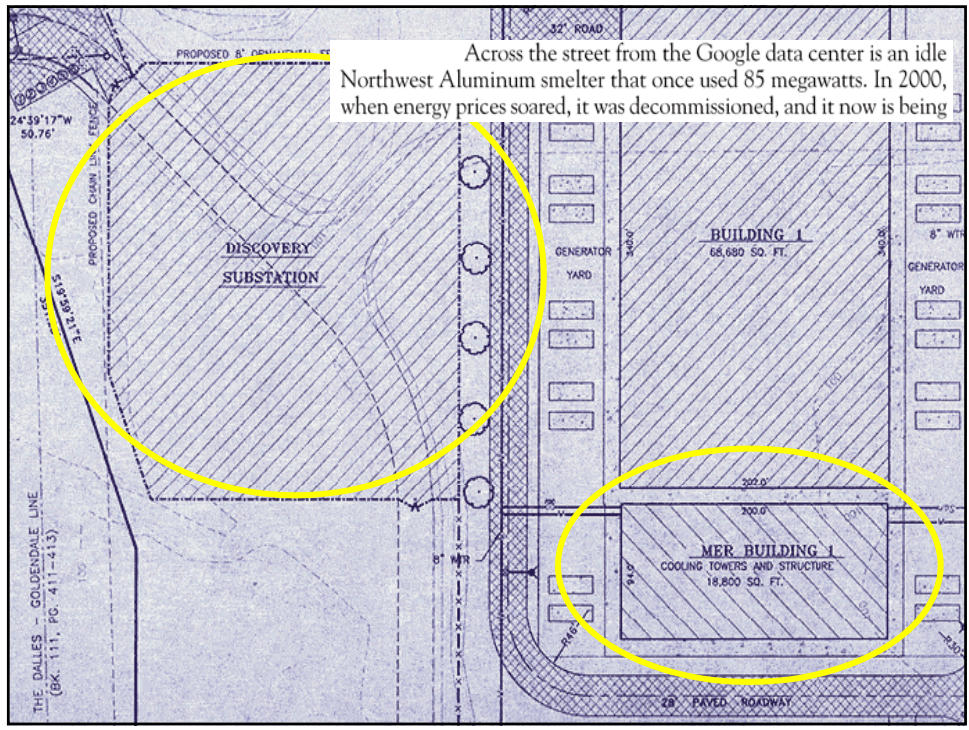
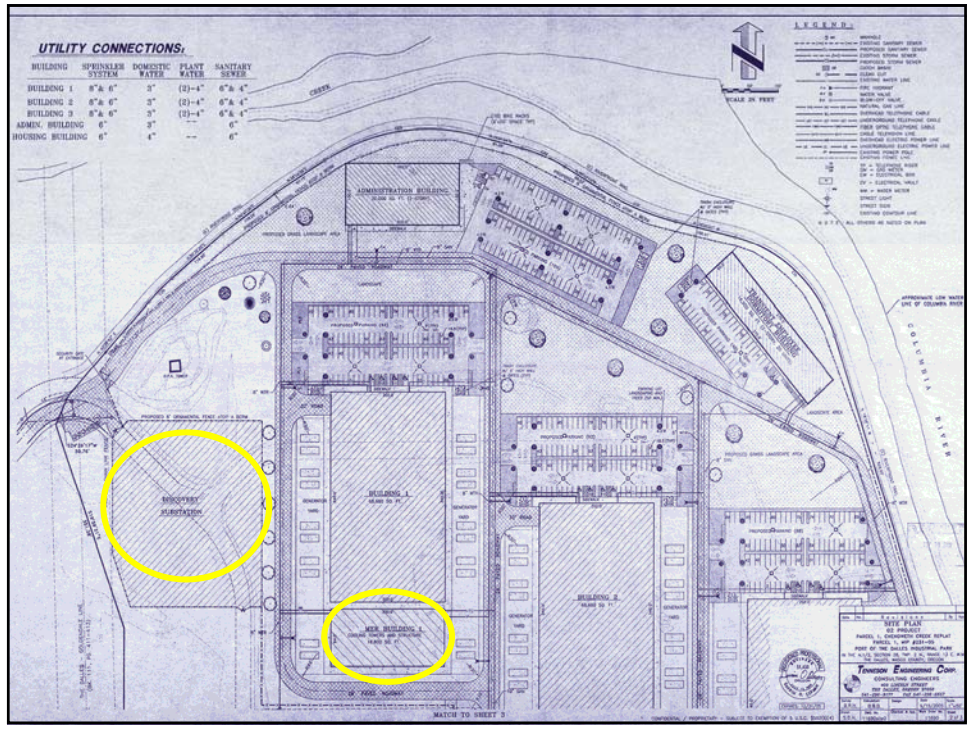


The Dalles, OR, Columbia River



Google Data Center Columbia River





Energiezufuhr „Discovery Substation“: 115 kV / 13.8 kV

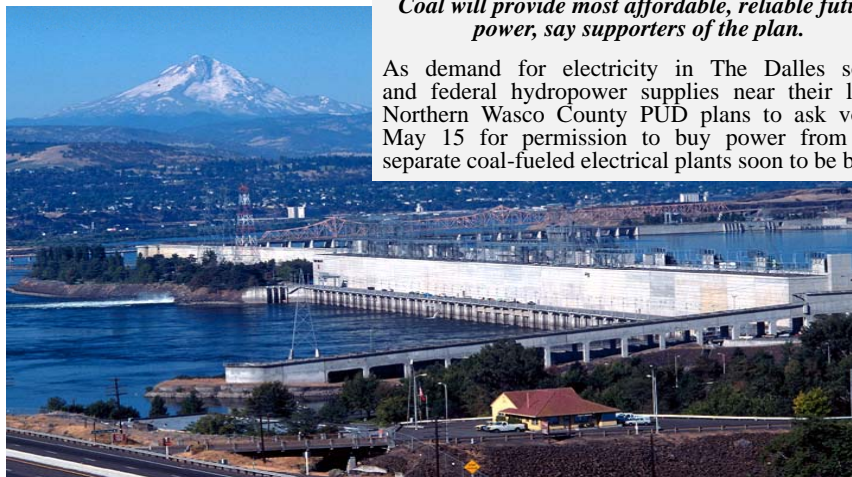


Nahes Kraftwerk: Dalles Dam Power Station, Columbia River

The Chronicle, March 1, 2007 –
PUD to seek vote on coal power

Coal will provide most affordable, reliable future power, say supporters of the plan.

As demand for electricity in The Dalles soars, and federal hydropower supplies near their limit, Northern Wasco County PUD plans to ask voters May 15 for permission to buy power from two separate coal-fueled electrical plants soon to be built.



Innenansicht eines Cloud-Zentrums

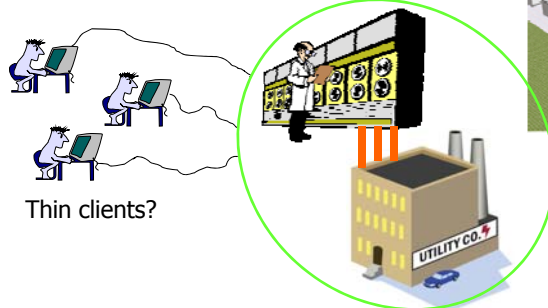


- Effizient wie **Fabriken**
 - Produkt: Internet-Dienste
- **Kostenvorteil** durch Skaleneffekt
 - Faktor 5 – 7 gegenüber traditionellen „kleinen“ Rechenzentren
- Angebot nicht benötigter Leistung auf einem **Spot-Markt**

Das entwickelt sich zum
eigentlichen Geschäft!

Zukünftige Container-Datenzentren

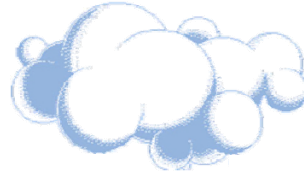
- Hunderte von **Containern** aus je einigen tausend Compute-Servern
 - mit Anschlüssen für Strom und Kühlung
- Nahe an **Kraftwerken**
 - Transport von Daten billiger als Strom



Thin clients?

Cloud-Computing für die Industrie und Wirtschaft

- **Spontanes Outsourcen** von IT inklusive Geschäftsprozesse
 - Datenverarbeitung als Commodity
 - Software und Datenspeicher als Service
- Keine Bindung von Eigenkapital
 - **Kosten nach „Verbrauch“**
- **Elastizität:** Sofortiges Hinzufügen weiterer Ressourcen bei Bedarf
 - virtualisierte Hardware



Markt für „utility computing“
2010: ca. 95 Milliarden EUR

Zusammenfassung Systemarchitekturen

- Peer-to-Peer
- Client-Server (Fat-Client vs. Thin Client)
- 3-Tier
- Multi-Tier
- Service-Oriented Architecture (SOA)
- Compute-Cluster
- Cloud-Computing

Charakteristika und praktische Probleme verteilter Systeme

- Räumliche Separation und Autonomie der Komponenten führen zu **neuen Problemen**:
 - **partielles Fehlverhalten** (statt totaler "Absturz")
 - fehlender globaler **Zustand** / exakt synchronisierte **Zeit**
 - **Inkonsistenzen**, z.B. zwischen Datei und Verzeichnis
- Typw. grosse **Heterogenität** in Hard- und Software
- **Komplexität**
- **Sicherheit** (Vertraulichkeit, Authentizität, Integrität, Verfügbarkeit,...)
 - **notwendiger** als in klassischen Einzelsystemen
 - aber **schwieriger** zu gewährleisten (mehr Angriffspunkte)

Gegenmittel?

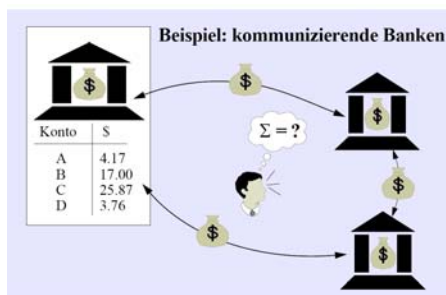
- Gute **Werkzeuge** ("Tools") und **Methoden**
 - z.B. Middleware als Software-Infrastruktur
 - **Abstraktion** als Mittel zur Beherrschung von Komplexität
 - z.B. Schichten (Kapselung, virtuelle Maschinen) oder
 - Modularisierung (z.B. Services)
 - Adäquate **Modelle, Algorithmen, Konzepte**
 - zur Beherrschung der "neuen" Phänomene
-
- **Ziel der Vorlesung**
 - Verständnis der **grundlegenden Phänomene**
 - Kenntnis von geeigneten Konzepten und Methoden

Einige konzeptionelle Probleme und Phänomene verteilter Systeme

- 1) Schnappschussproblem
- 2) Phantom-Deadlocks
- 3) Uhrensynchronisation
- 4) Kausaltreue Beobachtungen
- 5) Geheimnisvereinbarung über unsichere Kanäle

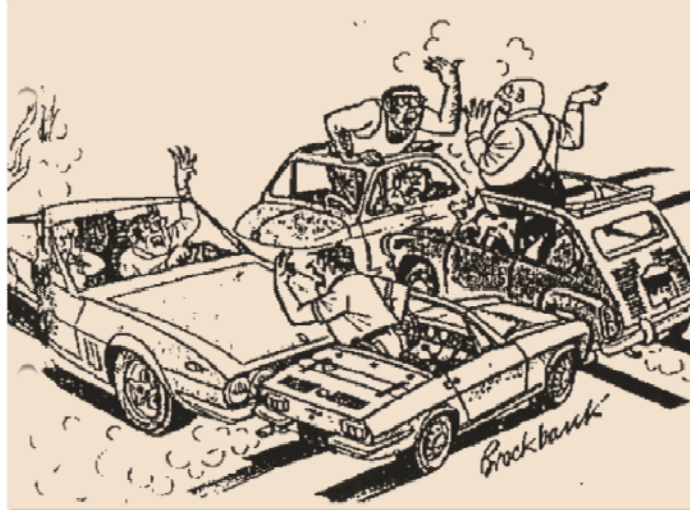
- Dies sind einige einfach zu erläuternde Probleme und Phänomene – es gibt aber noch viel mehr und viel komplexere Probleme konzeptioneller wie praktischer Art
- Achtung: Manches davon wird nicht hier, sondern in der Vorlesung "Verteilte Algorithmen" eingehender behandelt!

Ein erstes Beispiel: Wieviel Geld ist in Umlauf?



- Annahme: konstante Geldmenge
- **Ständige Transfers** zwischen den Banken
- Niemand hat eine **globale Sicht**
- Es gibt keine **gemeinsame Zeit** ("Stichtag")
- Anwendung: z.B. verteilte Datenbank-Sicherungspunkte

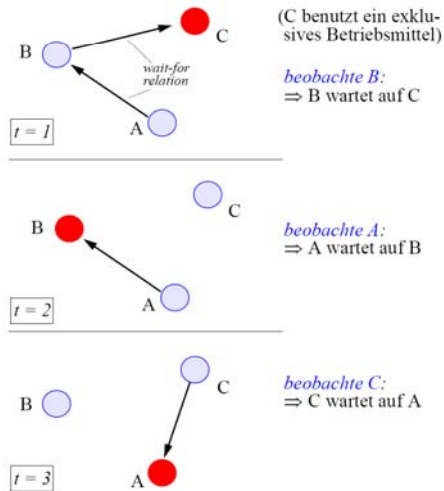
Ein zweites Beispiel: Das Deadlock-Problem



Ein zweites Beispiel: Das Deadlock-Problem

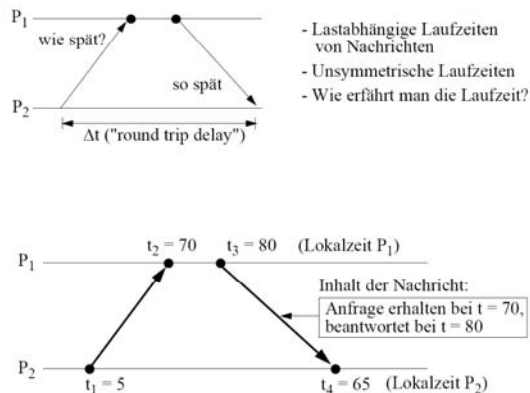


Phantom-Deadlocks



- Aus den Einzelbeobachtungen darf man **nicht** schliessen:
- A wartet auf B und B wartet auf C und C wartet auf A
- Diese **zyklische Wartebedingung** wäre tatsächlich ein Deadlock
- Die Einzelbeobachtungen fanden hier aber zu **unterschiedlichen Zeiten** statt
- **Lösung** ohne Zeitstempel?

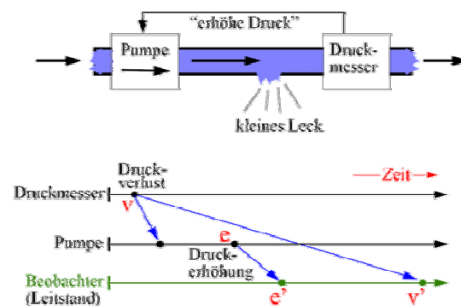
Ein drittes Problem: Uhrensynchronisation



- Uhren gehen nicht unbedingt **gleich schnell!**
- Gilt wenigstens "Beschleunigung ≈ 0 ", d.h. ist konstanter Drift gerechtfertigt?
- Wie kann man den **Offset** der Uhren ermitteln oder zumindest approximieren?

Ein viertes Problem: (nicht) kausaltreue Beobachtungen

- Gewünscht: Eine **Ursache** stets vor ihrer (u.U. indirekten) **Wirkung** beobachten

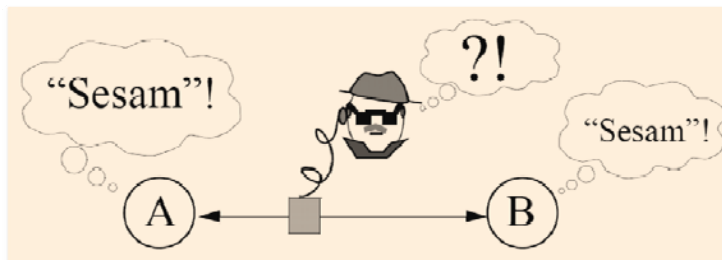


Falsche Schlussfolgerung des Beobachters:

Es erhöhte sich der Druck (aufgrund einer unbegründeten Aktivität der Pumpe), es kam zu einem Leck, was durch den abfallenden Druck angezeigt wird.

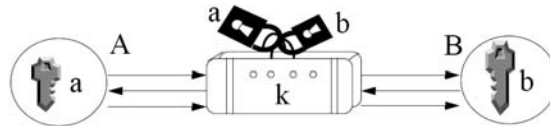
Und noch ein Problem: Verteilte Geheimnisvereinbarung

- Problem: A und B wollen sich über einen unsicheren Kanal auf ein gemeinsames geheimes Passwort einigen



Verteilte Geheimnisvereinbarung (2)

- Idee: Vorhängeschlösser um eine sichere Truhe:



1. A denkt sich Passwort k aus und tut es in die Truhe.
2. A verschliesst die Truhe mit einem Schloss a .
3. A sendet die so verschlossene Truhe an B.
4. B umschliesst das ganze mit seinem Schloss b .
5. B sendet alles doppelt verschlossen an A zurück.
6. A entfernt Schloss a .
7. A sendet die mit b verschlossene Truhe wieder an B.
8. B entfernt sein Schloss b .

- Problem: Lässt sich das so **softwaretechnisch** realisieren?

Kommunikation

Kooperation durch Informationsaustausch

- Prozesse sollen **kooperieren**, daher untereinander **Information austauschen** können
 - mittels gemeinsamer Daten in einem **globalen Speicher** (dieser kann physisch oder evtl. nur logisch vorhanden sein als „virtual shared memory“)
 - oder mittels **Nachrichten**:
Daten an eine entfernte Stelle kopieren

Kommunikation

- Notwendig, damit Kommunikation klappt, ist jedenfalls:
 1. ein dazwischenliegendes **physikalisches Medium**
 - z.B. elektrische Signale in Kupferkabeln
 2. einheitliche **Verhaltensregeln**
 - Kommunikationsprotokolle
 3. gemeinsame **Sprache** und gemeinsame **Semantik**
 - gleiches Verständnis der Bedeutung von Kommunikationskonstrukten und -regeln

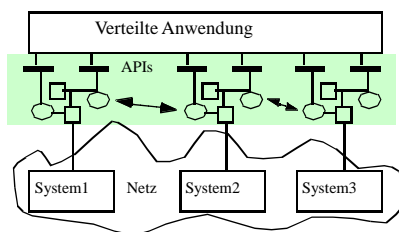
Also trotz Verteiltheit gewisse **gemeinsame Aspekte!**

Nachrichtenbasierte Kommunikation

- **send** → **receive**
- Implizite **Synchronisation**:
Senden vor Empfangen
 - Empfänger erfährt, wie weit der Sender mindestens ist
- Nachrichten sind **dynamische Betriebsmittel**
 - verursachen Aufwand und müssen verwaltet werden

Message Passing System (1)

- Organisiert den Nachrichtentransport
- Bietet **Kommunikationsprimitive** (als **APIs**) an
 - z.B. `send (...)` bzw. `receive (...)`
 - evtl. auch ganze **Bibliothek** verschiedener Kommunikationsdienste
 - verwendbar mit gängigen Programmiersprachen (oft zumindest C)



- Besteht aus Hilfsprozessen, Pufferobjekten, ...
- **Verbirgt Details** des zugrundeliegenden Netzes bzw. Kommunikationssubsystems

Message Passing System (2)

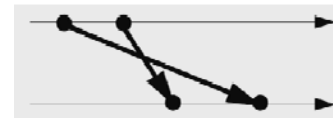
- Verwendet vorhandene Netzprotokolle und implementiert damit neue, „höhere“ Protokolle
- **Garantiert** (je nach „Semantik“) gewisse **Eigenschaften**
 - z.B. Reihenfolgeerhalt oder Prioritäten von Nachrichten
- **Abstrahiert von Implementierungsaspekten**
 - z.B. Geräteadressen oder Längenrestriktionen von Nachrichten etc.
- **Maskiert gewisse Fehler**
 - mit typischen Techniken zur Erhöhung des Zuverlässigkeitsgrades: Timeouts, Quittungen, Sequenznummern, Wiederholungen, Prüfsummen, fehlerkorrigierende Codes,...
- **Verbirgt Heterogenität** unterschiedlicher Systemplattformen
 - erleichtert damit **Portabilität** von Anwendungen

Ordnungserhalt von Nachrichten: FIFO

- Manchmal werden vom Kommunikationssystem Garantien bzgl. **Nachrichtenreihenfolgen** gegeben
- Eine mögliche Garantie stellt **FIFO** (First-In-First-Out) dar: Nachrichten zwischen zwei Prozessen überholen sich nicht: **Empfangsreihenfolge = Sendereihenfolge**



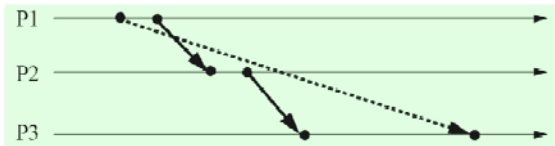
FIFO



kein FIFO

Ordnungserhalt von Nachrichten: kausale Ordnung

- FIFO verbietet allerdings nicht, dass Nachrichten evtl. **indirekt** (über eine Kette anderer Nachrichten) **überholt** werden



Zwar FIFO, aber nicht kausal geordnet

- Möchte man auch dies haben, so muss die Kommunikation **kausal geordnet** sein (Anwendungszweck?)
 - keine Information erreicht Empfänger **auf Umwegen schneller** als auf direktem Wege („Dreiecksungleichung“)
 - entspricht einer „**Globalisierung**“ von FIFO auf mehrere Prozesse
 - **Denkübung**: Wie garantiert (d.h. implementiert) man kausale Ordnung auf einem System ohne Ordnungsgarantie?

Prioritäten von Nachrichten? (1)

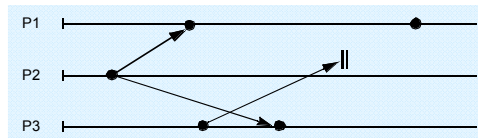
- Achtung: **Semantik** ist a priori nicht ganz klar:
 - Soll (kann?) das Transportsystem Nachrichten höherer Priorität bevorzugt (=?) befördern?
 - Können (z.B. bei fehlender Pufferkapazität) Nachrichten niedrigerer Priorität überschrieben werden?
 - Wie viele Prioritätsstufen gibt es?
 - Sollen auf Empfangsseite zuerst Nachrichten mit höherer Priorität angeboten werden?

Prioritäten von Nachrichten? (2)

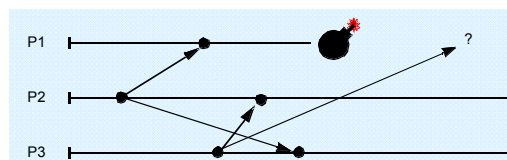
- Mögliche **Anwendungen**:
 - Unterbrechen laufender Aktionen (→ Interrupt)
 - Aufbrechen von Blockaden
 - Out-of-Band-Signalisierung } Durchbrechung der FIFO-Reihenfolge!
- Vgl. auch Service-Klassen in **Computernetzen**: bei Rückstaus bei den Routern soll z.B. interaktiver Verkehr bevorzugt werden vor FTP etc.
- **Vorsicht** bei der Anwendung: Nur bei klarer Semantik verwenden; löst oft ein Problem nicht grundsätzlich!
 - Inwiefern ist denn eine (faule) Implementierung, bei der „eilige“ Nachrichten (insgeheim) wie normale Nachrichten realisiert werden, tatsächlich nicht korrekt?

Fehlermodelle (1)

- Zweck: Klassifikation von Fehlermöglichkeiten; Abstraktion von den konkreten, spezifischen Ursachen
- **Fehler beim Senden / Übertragen / Empfangen:**



- **Crash / Fail-Stop:** Ausfall eines Prozessors:



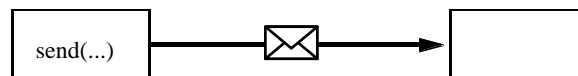
Fehlermodelle (2)

- **Zeitfehler:** Ereignis erscheint zu spät (oder zu früh)
- **„Byzantinische“ Fehler:** Beliebiges Fehlverhalten, z.B.:
 - verfälschte Nachrichteninhalte
 - Prozess, der unsinnige Nachrichten sendet

(solche Fehler lassen sich nur teilweise, z.B. durch **Redundanz**, erkennen)
- **Fehlertolerante Algorithmen** müssen das „richtige“ Fehlermodell berücksichtigen!
 - adäquate Modellierung der realen Situation / des Einsatzgebietes
 - Algorithmus verhält sich **korrekt nur relativ zum Fehlermodell**

Mitteilung vs. Auftrag (1)

Mitteilungsorientiert

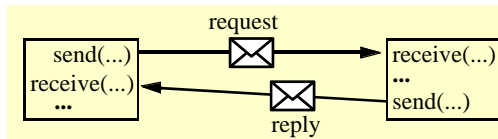


- Unidirektional
- Übermittelte Werte werden der Nachricht typw. als „Ausgabeparameter“ beim **send** übergeben

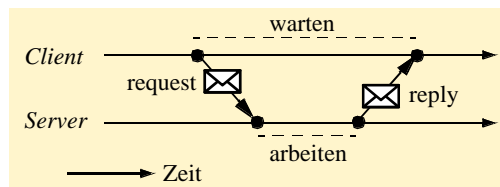
Mitteilung vs. Auftrag (2)

Auftragsorientiert

evtl. zu einem *einzigem* Befehl zusammenfassen



- Bidirektional
- Ergebnis des Auftrags wird als „Antwortnachricht“ zurückgeschickt
- Auftraggeber („Client“) **wartet**, bis Antwort eintrifft



Resümee (2a)

- **Architekturvarianten** verteilter Systeme (Fortsetzung)
 - Service-Oriented Architecture (SOA)
 - Cloud-Computing
- **Cloud-Computing**
 - motiviert durch schnellere / ubiquitäre Netze
 - Trend: „alles“ irgendwo im Netz
 - Beispiele für Datenzentren
 - wirtschaftliche Effekte: Skaleneffekte, Spot-Markt
- **Charakteristika / Problembereiche verteilter Systeme**
 - fehlender globaler Zustand / Zeit; partielles Fehlverhalten
 - Heterogenität; Komplexität

Resümee (2b)

- **Beispielhafte Phänomene und konzeptionelle Probleme**
 - Schnappschussproblem (inkonsistente globale Sicht)
 - Phantom-Deadlocks
 - Uhrensynchronisation
 - kausal inkonsistente Beobachtungen
 - Geheimnisaustausch über unsicheren Kanal
- **Nachrichtenkommunikation**
 - Message-passing-Systeme
 - Ordnungserhalt; Prioritäten von Nachrichten
- **Fehlermodelle**
 - fehlerhaftes Senden / Empfangen / Transportieren von Nachrichten
 - Crash von Prozessen bzw. Prozessoren
- **Kommunikationsmuster**
 - Mitteilung ↔ Auftrag