

Zusammenfassung Computernetzwerke und verteilte Systeme



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Alex Praus, Tobi Kratz
13. Februar 2021

Inhaltsverzeichnis

1 Quick Tour	2
1.1 Making two devices communicate	2
1.2 Connecting many computers	3
1.3 Organizing the mess - and connecting 'Alien' computers	3
1.3.1 OSI	3
1.3.2 5 Layer of the Internet	5
1.3.3 network types	5
2 Routing	5
2.1 introduction	5
2.1.1 Forwarding	5
2.1.2 Routing	6
2.2 Routing Algorithms	6
2.2.1 Examples	6
2.3 Distance Vector Routing	7
2.3.1 Count to Infinity Problem	7
2.4 Link-State Routing	8
2.5 Hierarchical Routing	8
2.5.1 About BGP and Internet Routing	9
2.6 Mobile Routing	9
2.6.1 DSDV	10
2.6.2 DSR	10
2.7 Overlay Routing	11
3 IP & Internetworking	11
3.1 IP Header	11
3.2 Addressing	12
3.2.1 Naming and DNS	12
3.2.2 IPv4 Addresses	12
3.2.3 NAT	13
3.2.4 ARP	13
3.3 IPv6	13
4 Transport Layer	14
4.1 Segmentation(+Blocking). Multiplexing, Addressing	15
4.2 Connection Control	15
4.2.1 Phases	15
4.2.2 Error Control	16
4.2.3 Two Army Problem	16
4.3 Flow Control	16
4.3.1 Rate and Credit based Flow Control	17
4.4 Error Control	17
4.5 Congestion Control	18
4.5.1 solving congestion Control	18

4.6	(Internet) UDP	19
4.7	(Internet) TCP	19
4.7.1	Overview	19
4.7.2	Error & Connection Control in TCP	20
4.7.3	Flow & congestion Control in TCP	21
4.8	Other Protocols	22
5	Queueing Theory	22
6	Referenzen	22
6.1	Quicktour	22
6.2	Routing	23
6.3	IP & Addressing	23
6.4	Transport	23

1 Quick Tour

1.1 Making two devices communicate

Verschiedene Möglichkeiten, um Kommunikation herzustellen: direkte physische Verbindung. Hier werden Daten als Bits übertragen. Eine 1 könnte als steigende Taktflanke und eine 0 als fallende dargestellt werden, es gibt jedoch viele Probleme bei der Implementation (00, 11, wechselnde Taktfrequenz...).

low level properties of communication Es gibt verschiedene Werte, die bei low-level communication wichtig sind:

- Delay (Latenz) $d = \text{distance} / \text{Propagation speed } v$. (v ist im Vakuum $= c$, in Kupfer etwa $\frac{2}{3} c$)
- Data rate r (Datenrate) = Data size / data rate (Bsp. bits/second).
Wichtig ist hier, dass nicht die Geschwindigkeit gemeint ist, mit der die Daten transportiert werden (siehe v), sondern in welcher Rate die Bits auf die Leitung gelegt werden.

Wenn der Sender eine Daten sendet, werden diese nicht auf Sender Seite gespeichert, sondern lediglich beim Empfänger. Allernhöchstens sind die Daten während der Übertragung im Kabel gespeichert.

Beispiel Latenzberechnung (delay & data rate): Wir senden 1250 Bytes ($10^4 b$) über 6 Meter mit einer Geschwindigkeit von 10 Mbps ($10^7 \frac{b}{s}$). Der Delay d beträgt dabei $d = \frac{6m}{c} = \frac{6m}{3 \cdot 10^8 \frac{m}{s}} = 2 \cdot 10^{-8} s = 20ns$. Die Data rate beträgt $r = \frac{10^4 b}{10^{-3} s} = 10^7 \frac{b}{s} = 10^7 \frac{b}{s}$.

Types of physical communication Die Typen lassen sich in 3 Fälle aufteilen:

- Simplex: Eine Seite kann nur senden, die andere nur empfangen (one-way). Empfänger lassen sich jedoch beliebig skalieren.
- Half Duplex: Beide Seiten wechseln sich ab mit senden und empfangen (vgl. Telefonat/Gespräch).
- Duplex: Beide Seiten können senden wie und wann sie lustig sind (vgl. Streitgespräch).

Während Simplex und Half Duplex einfach realisierbar sind, ist (Full) Duplex eine technische Herausforderung.

Realizing Half Duplex Denkbar wären 2 Kabel, eins je Host, das wäre jedoch eine Verschwendung, da diese Kabel nie gleichzeitig genutzt werden würden. Es gibt zwei Ansätze: Time division duplex (TDD) und on-demand duplex. Beim TDD hat jeder Host eine feste Zeit T zum senden, und es wird zwangsläufig abgewechselt. Beim on-demand duplex gibt jeder Host die Länge der nächsten Bit Sequenz am Anfang direkt bekannt (pre-announce).

Realizing Full Duplex Hier wären 2 Kabel eher anwendbar, bedeuten aber den doppelten Aufwand. Auf kurzer Distanz kann jedoch auch Full Duplex mit einem Kabel realisiert werden, in dem jeder Host eine leicht andere Frequenz benutzt (z.B. bei WiFi). Auch TDD ist umsetzbar. Während A sendet, speichert B die zu sendenden Daten. Nach T sendet B dann den Stack während A speichert. Das ist jedoch nicht analog umsetzbar, da z.B. Sprachdaten sich nicht in Portionen aufteilen lassen.

1.2 Connecting many computers

Jeden Computer direkt mit jedem anderen zu verbinden wäre zwar möglich, skaliert jedoch eher so semi (<https://www.reddit.com/r/cablegore/top>). Switches wäre eine Lösung, jedoch laufen dann mehrere Connections über ein Kabel, es kommt also zu Einbußung in der Data Rate. Außerdem besteht das Problem, wie man eine Connection über einen switch herstellt (vgl. Vermittlung beim Telefon). Das führt jedoch dazu, dass eine Connection nur einfach genutzt werden kann, ist ein Host also mit einem anderen verbunden, kann ein Host nicht mehr erreicht werden.

Packet switching Statt also ein Circuit für eine Verbindung zu blocken, teilt der switch die Daten in packets und sendet diese. So wird die Leitung nur für die Länge eines Packets geblockt und es kann schneller gewechselt werden.

Probleme: Anfang und Ende bestimmen? Wie bestimmt man wo das Packet hin soll?

Beispiel Ablauf: »store-and-forward«switching:

1. receive a complete packet
2. store the packet in a Buffer
3. Find out the packet's destination
4. decide where the packet should be sent next (benötigt Kenntnis über Netzwerk Topologie)
5. forward the packet to his next hop of its journey

Multiplexing

»Organizing the forwarding of packets over such a single, shared connection is called multiplexing.«

Auch beim Multiplexing ist TDM (Time Division Multiplexing) (nur ein Packet gleichzeitig) und FDM (Frequenz Division Multiplexing) (mehrere Packets auf gleichzeitig auf unterschiedlichen Frequenzen) möglich. Bei optischen Verbindungen WDM (wavelength division multiplexing) statt FDM. Es gibt jedoch noch weitere Formen, die hauptsächlich für Wireless transmission geeignet sind: CDM (Code Division Multiplexing) und SDM (Space Division Multiplexing).

Multiplexing lässt sich auch abstrahieren, wenn zum Beispiel mehrere connections eines höherer Layers eine lower-level connection nutzen sollen (s.c. upward multiplexing). Allgemein lässt sich in dem Kontext von shared Resources sprechen.

Forwarding and next hop selection Bekanntes Problem: Wie weiß eine Host/router/switch, wo er Packets hinschicken soll? Wie kennt er die beste/schnellste Verbindung? Es gibt einfache Ansätze:

- Flooding: alle Packets an alle Nachbarn senden
- Hot-potato routing: So schnell wie möglich an einen/mehrere zufällige Nachbarn senden

Sinnvoller wäre jedoch, wenn sich der router die besten wege merkt, e.g. durch Routing Tabellen. Diese können durch 2 Arten gesammelt werden: aktiv und passiv.

Passiv: Aktiv traffic beobachten und daraus Schlüsse ziehen.

Aktiv: Informationen aktiv senden und empfangen unter den verschiedenen Routern (routing protocols).

Jedoch bei Netzen der größe unsere Internets können Routing Tabellen schnell sehr groß werden. Hier kann man das Netz in mehrere Teilnetze splitten (divide et impera).

1.3 Organizing the mess - and connecting 'Alien' computers

Der Schlüssel ist »Simplification by abstraction«.

Z.B. das Modell des DS (distributed Systems) hilft beim verstehen von CN (computer networks). ein DS besteht aus AS (autonomous systems) und CSS (communication subsystems). Oft referiert wird auch auf die Abstraktion des OSI (Open Systems Interconnection) Schichten Modell.

1.3.1 OSI

Part 1: Concepts and Terms

»A protocol defines the format and the order of messages exchanged between two or more communicating entities, as well as the actions taken on transmission and/or reception of a message or other event«

Analog zum 'Computer' Protokoll kann man sich ein menschliches Protokoll vorstellen, z.B. der definierte Ablauf beim Telefonieren (Hallo, Hallo..... Tschüß, Tschau). Weiter ist das OSI Modell in Schichten aufgeteilt. Analog wäre z.B. Layer 2 ein Manager, der einen Brief in Auftrag gibt (Layer 1), der dann von der Post (Layer 0) transportiert wird. Ein Layer N hat bietet also Services für den Layer n+1 an. Zwei Layer n kommunizieren nur mit den für diesen Layer vorgesehenen Protokollen, diese können jedoch auf Services von Layer n-1 zugreifen (der unterste Layer kann natürlich nur eigene Services nutzen).

Services im OSI Modell lassen sich in 2 Gruppen aufteilen:

- **connection-orientated services:** Diese haben meist 3 Phasen:
 - CON (Connection Establishment)
 - DAT (Data Exchange)
 - DIS (Disconnect)
- **connectionless services:** Bei diesen fallen CON und DIS weg und es werden nur Daten ausgetauscht.

Im OSI Modell werden Nachrichten höherer Layer als Daten Einheiten (Data Units) tieferer Layer transportiert. Hierfür gibt es einige common Notations:

- packet: Ist die Einheit, die transportiert wird (kann aus Fragmenten bestehen)
- datagram: wird bei connectionless services als Ersatz für Pakete verwendet
- frame: 'fertig und verpackt' um versendet zu werden.
- cell: kleinere packet mit einer definierten Größe
- PDU (Protocol Data Unit): eine (N)-PDU ist definiert durch: (N)-PCI + (N)-SDU
- PCI (Protocol Control Information): wird nur von peers genutzt.
- SDU (Service Data Unit): Ist die zu versendete payload eines höheren Layers. (N)SDU = (N+1)-PCI + (N+1)-SDU

Part 2: 7 Layer Model

Layer 1: physical Layer (PH) Senden von Bits durch (de-)aktivieren von Signalen auf Leitungen.

Beispiel: 1000 Base-T

Layer 2: data link layer (D) Sendet Pakete als Frames um Fehler zu erkennen und zu beheben, um fehleranfällige Hosts zu schützen. Kann auch Flow Control verwenden, um langsame Hosts zu schützen.

Beispiel: Ethernet

Layer 3: network layer (N) Ziel ist es den Packet-Stream zwischen zwei Hosts zu ermöglichen. Koordinierung der Pfade von Host zu Host. Konkret: Routing Wege finden und Pakete weiterleiten und Fehler beheben, z.B. durch 'flow control'. Up- und downward multiplexing ist möglich. Außerdem kann congestion control auf diesem Layer angewendet werden.

Beispiel: IP

Layer 4: transport layer (T) Logische Verbindungen zwischen zwei Prozessen (nicht nur zwischen zwei Computern), Fehlerkorrektur und Paketzusammensetzung für den N-Layer.

Beispiel: TCP

Layer 5: session layer (S) Koordiniert Session z.B. bei HTTP mit dem Session-Cookies und hilft Nutzer und Anwendung bei der Konstruktion und dem spannen von aufeinanderfolgenden Verbindungen.

Beispiel: HTTPS

Layer 6: presentation layer (P) Setzt die Daten in eine unabhängige Form um, um unabhängig davon die Übermittlung - gg. mit Kompression und Verschlüsselung - zu ermöglichen.

Beispiel: LDAP

Layer 7: application layer (A) Stellt Funktionen wie Daten Ein- und Ausgabe zur Verfügung.

Beispiel: XMPP

1.3.2 5 Layer of the Internet

Im Internet verschmelzen Layer 5,6,7 oft und oft sind die Übergänge nicht klar definiert.

Layer 1 Übermittlung von Frames als Stream von Bits

Layer 2 Daten von Layer 3 in Frames verpacken und an direkte Nachbarn weiterleiten

Layer 3 Daten vom Client zum Web-Server weiterleiten, Router-to-Router communication, außerdem können Peer-Verbindungen durch hop-to-hop realisiert werden

Layer 4 Verlässliche Verbindung zum Web-Server herstellen und sicherstellen, dass die Daten auch in der richtigen Reihenfolge ankommen, jedoch keine congestion control.

Layer 5,6,7 HTTP Anfragen erstellen, Ebene 4 aufrufen (TCP).

1.3.3 network types

Wie schon erwähnt gibt es CO und CL networks. Beispiel für CO ist z.B. Das Telefon, CL die Post. CO haben durch die durch Handshakes vor allem bei kurzen Verbindungen eine hohe 'extra Last' durch die zusätzlichen Daten des Handshakes, lassen sich jedoch besser skalieren, da sie nicht einen Status gebunden sind (stateless). CO sind jedoch durch den Status der connection zuverlässiger. Bei einem 'verstopften' Network können mit CL immernoch Daten versendet werden, es kann nur sein, dass diese verspätet ankommen, CO haben jedoch Schwierigkeiten. Es ist möglich CO auf CL aufzubauen.

connection-oriented Networks Im CO network ist der erste Schritt mit einem handshake eine connection herzustellen. Nach dem handshake wissen beide Seiten von der connection und der Datenaustausch kann stattfinden. Die connection ist dabei nur ein loser status, auf dem Basis jedoch andere Eigenschaften (flow control, congestion control...) aufgebaut werden können. Bei CO networks implementiert nicht direkt andere Eigenschaften der Verbindung. Dinge wie reliability, flow control und congestion control sind für CO networks nicht notwendig. Diese können z.B. mit TCP ermöglicht werden.

connectionless networks keine handshakes, direkter Fahren (Daten Austausch). Wenn der Datenaustausch vorbei ist, kommt auch kein DIS mehr. CL networks brauchen wenig Aufwand, da keine connection gepflegt werden muss, es kann jedoch sein, dass der receiver nicht bereit zum Empfangen ist. CL implementieren keine reliability, flow control und congestion control.

2 Routing

2.1 introduction

Im Routing werden dafür gesorgt, dass Pakete vom Empfänger an das richtige Ziel kommen. Bei direkten Verbindungen besteht das Problem nicht, jedoch ist das bei großen Netzwerken keine Option, wenn jeder Host mit jedem anderen verbunden werden muss. Wenn Switches genutzt werden, muss diesen jedoch gesagt werden, wie das Netzwerk aufgebaut ist (Netzwerk Topologie). In diesem Kapitel geht es um das Problem, wie ein Host den besten Weg zu einem Ziel findet.

Building a large network Bei großen Netzwerken ist flooding und Hot-potato routing keine Option, da mit jedem Host die Anzahl an Paketen steigt und so mit den beiden Routing ineffizienten Methoden das Netzwerk schnell an sein Limit kommt. Ziel ist es eine effiziente Methode zu finden, die Pakete möglichst schnell ans Ziel bringt, ohne dabei unnötig viel Traffic zu erzeugen. Im folgenden werden zwei Begriffe genutzt:

- **Routing:** determine route taken by packets from source to destination. (Basis: Routing algorithms).
- **Forwarding:** move packets from router's input to appropriate router output.

2.1.1 Forwarding

Wenn Pakete von einem Netzwerk in ein anderes Netzwerk geleitet werden sollen, wird ein Router eingesetzt. (Heute haben uns bekannte Router mehrere Aufgaben, die früher von verschiedenen Geräten übernommen wurden: hub, bridge, switch, gateway.) Wenn also Pakete von einem Netzwerk in ein anderes gesendet werden sollen, übernimmt der Router die Koordination und leitet das Paket (forwarded) in das entsprechende Ziel Netzwerk. Hängt das andere Netzwerk direkt am selben Router, handelt es sich um ein single hop. Wenn mindestens 2 Router zwischen den Netzwerken sind, handelt es sich um ein Multi-Hop.

2.1.2 Routing

Routing findet für gewöhnlich auf Layer 3 statt, dort ist das Ziel Pakete von Host A zu Host B möglichst effizient zu transportieren bzw. erstmal einen Pfad zu finden, auf dem das möglich ist. Dies wird i.d.R. von Routing Algorithmen durchgeführt. Das Internet besteht aus mehreren AS, die alle wieder aus Teilnetzen bestehen. Jedes AS führt dabei selbst routing Algorithmen aus, um die besten Wege zu finden.

- **CONS (CO+NS)** In CO Networks Routing Algorithmen werden meist in der CON Phase durchgeführt. Im COTS (CO+TS) wissen nur die Endsysteme, dass sie verbunden sind, in CONS hingegen wissen alle Systeme auf der Route, dass die Systeme verbunden sind.
- **CLNS (CL+NS)** In CL Networks wird nicht bei jedem Packet der Algorithmus durchgeführt, das würde einen zu großen Overload bedeuten. Manchmal wird beim ersten Packet einer Verbidnung der Algorithmus durchgeführt, das ist allerdings für sich schnell ändernde Netzwerke keine Option. Im Inetnet z.B. dies in regelmäßigen Abständen, oder wenn sich große Teile ändern.

optimizing Routing Algorithmen Routing Algorithmen haben oft unterschiedliche Kriterien, nach denen sie arbeiten:

- Average packet delay
- Total throughput
- individual delay (kann jedoch mit anderen Kriterien im Widerspruch stehen)

Am meisten jedoch wird nach dem Kriterium, des minimal-hop-count gearbeitet, da dieser oft einen Kompromiss aus allen Kriterien bedeutet, es gibt jedoch keine Garantie dafür.

2.2 Routing Algorithms

Routing Algorithmen werden meist in zwei Arten aufgeteilt:

- **Non-adaptive Routing Algorithms**

Diese agieren unabhängig vom State des Netzwerks. Beispiele sind flooding oder preconfiguration.

- **Adaptive Routing Algorithms**

Nehmen den aktuellen Status des Netzwerks mit in Betracht, wenn sie Routing Entscheidung treffen. Beispiel Hierfür wären distance-vector-routing oder link state routing. Das Problem hierbei ist, dass bei sich ändernden Netzwerken die Routen häufig neu entschieden werden. Algorithmen diesen Types sind trotzdem sinnvoll in eignen, fest bekannten Netzen. Bekommt ein Link z.B. so viel Traffic, das Pakete verloren gehen, wird der Link als Broken markiert und es kommt zu noch größeren Ausfällen. Es gibt dort auch 3 Unterarten:

- Centralized adaptive routing
- Isolated (aka. local) adaptive routing
- Distributed adaptive routing

2.2.1 Examples

Flooding (non-adaptive) Hier wird jedes einkommende Packet an alle bekannten Nachbarn weiterleitet. Das Problem dabei ist, dass so Netze schnell überlastet werden. Gibt es zum Beispiel Schleifen, kann es schnell zu einer Flut an nicht aufhörenden Paketen kommen. Eine Lösung für dieses Problem wäre z.B. das implementieren von TTL (Time to live) oder Sequence Number. TTL werden meist in Hops angegeben und werden bei jedem Hop um 1 dekrementiert. Hat ein Packet ein TTL von 0, wird es weggeworfen. Eine Sequence Number wird beim ersten Router initialisiert. Jeder Router führt eine Tabelle mit Sequence Numbers, die er schon einmal geroutet hat. Kommt ein Packet mit einer Sequence Number, die er schon kennt, wird dieses Packet weggeworfen.

Flooding macht jedoch durchaus Sinn in sich schnell ändernden Netzen, z.B. bei WLAN oder Mobilfunk oder wenn alle Pakete Multicast sind und so wie so mehrere Ziele haben.

Static Routes (non-adaptive) Static Routes sind großartig für statische, vorhersehbare Umgebungen. Das Problem ist, dass sich das Internet regelmäßig ändert und statische routen dann viel Wartungsaufwand bedeuten.

Centralized Adaptive Routing (adaptive) Es gibt einen Zentralen Control Center (RCC), der regelmäßig Informationen über die Topologie von allen Routern bekommt und dann einen Idealen Routing Graph erzeugt (z.B. Dijkstra). Das Problem hierbei ist, dass das Netz zusammenbricht, wenn nur der RCC ausfällt. Außerdem werden Routen 'in der Nähe' des RCC bevorzugt, was dort zu einer hohen Last führt während 'abgelegene' Router meist wenig Aufgaben haben. Ebenso bekommen Router die näher am RCC sind schneller die neuen Routing Informationen, was zu unterschiedlichen States führen kann.

Isolated (aka. local) adaptive Routing (non-adaptive) Es werden Entscheidungen über Routen nur lokal getroffen. Beispiele sind Hot potato Routing und Backward learning.

- **Hot Potato Routing**
Idee ist es, die Pakete so schnell wie möglich loszuwerden, wobei nicht beachtet werden muss, zu welchen Host die Pakete geschickt werden. Dieser Algorithmus ist nicht sehr effektiv, es gibt jedoch einige use-Cases in denen diese Art noch genutzt wird (peering/discovering).
- **Backward Learning Routing**
Bei diesem Algorithmus werden im Packet Header Source Adresse und Hop Counter hinzugefügt, Router lernen also im laufenden Betrieb über die Topologie und passen die Routen im Betrieb an. Jedoch müssen in jungen Netzwerken andere Algorithmen genutzt werden (z.B. hot potato / flooding). Wenn der Hop-Count == 1 ist, kommt das Packet von einem direkten Nachbarn. Bei einem Hop-Count $n > 1$ ist die source n hops away.

Distributed Adaptive Routing Durch Graph Abstraction die besten Routen finden. Knoten sind dabei Router und Kanten die physikalischen Links a.k.a. hops. Die Kosten eines links sind dabei z.B. delay, \$, oder der congestion Layer. Die Kosten eines Pfades sind dann alle link Kosten vereint. Ein guter Pfad wird meist als der, mit den geringsten Kosten bezeichnet, es kann aber auch nach anderen Kriterien gesucht werden (e.g. min-hop-count).

Algorithmen hier lassen sich weiter klassifizieren:

- **Decentralized** Jeder Router kennt die Kosten zu seinen Nachbarn. Auch Distance Vector Routing gehört hierzu (z.B. BGP oder RIP)
- **Global** Alle Router kennen die komplette Topologie und alle link kosten. Hierzu gehören Link state Algorithmen z.B. Dijkstras oder OSPF.
- **Static** (nicht adaptiv) Routen ändern sich sehr selten
- **Dynamic** (adaptiv) Routen können sich oft ändern, Hier werden also regelmäßig updates in den Routen gemacht.

2.3 Distance Vecotr Routing

Beim Distance Vector Routing tauschen direkte Nachbarn Informationen über Routen mit ihren Nachbarn aus. Jeder Host pflegt eine Tabelle, in der jede mögliche Ziehladresse eine Reihe und jeder Nachbar eine Spalte hat. In der Tabelle werden dann die "Kosten" der Route eingetragen und mit jeder Iteration verbessert. Konkret schreibt man dann für Route von X to Y via Z als nächsten Hop:

$$D^X(Y, Z) = c(X, Z) + \min_w \{D^Z(Y, w)\}$$

Mit einem Routing Algorithmus wird dann eine "Distance Table/Matrix"gebaut, mit der dann Routing Tablelen aufgestellt werden, aus denen dann der Distance Vector an die Nachbarn announced werden kann. DVR hat jedoch einige Probleme (count to infinity), jedoch handelt es sich um einen sehr simplen Algorithms.

DVR Protokolle sind iterativ und Distributed:

- **Iterativ** Das heißt sie laufen nicht unendlich, sondern stoppen sobald keine weiteren Verbesserungen möglich sind. Außerdem sind sie β self-terminating" d.h. es gibt kein Stop Signal o.ä. Eine Iteration wird dabei ausgelöst indem entweder ein loakler link sich ändern z.B. in den Kosten oder wenn es eine Nachricht eines Nachbarn gibt, dass der Link dort sich geändert hat.
- **Distributed** Des weiteren tauschen sie Informationen nur mit direkten Nachbarn aus und kennen auch nur den State dieser. Eine Node informiert einen Nachbarn dann über neue Routen, wenn sich die Kosten zu einer Destination verringert haben.

2.3.1 Count to Infinity Problem

Gegebene Situation: Wir haben 3 Host A,B,C. A ist mit B mit einem Cost von 1 verbunden, und B ist mit C mit einem Cost von 2 verbunden. Daraus ergeben sich folgende 3 Routing Tabellen: Durch einen Ausfall verschwindet jetzt die Verbindung zwischen B und

A			B			C		
TO	COST	VIA	TO	COST	VIA	TO	COST	VIA
B	1	B	A	1	B	A	3	B
C	3	B	C	2	C	B	2	B

C. A announced an B jedoch, dass es eine Route zu C mit dem Cost von 3 gibt. B versucht nun also C via A zu erreichen: Nachdem B dann diese Information an A sendet, aktualisiert A dann seine Route zu C, da diese über B geht und sich die Kosten erhöht haben.

A			B			C		
TO	COST	VIA	TO	COST	VIA	TO	COST	VIA
B	1	B	A	1	B	A	-	-
C	3	B	C	4	A	B	-	-

die Route von A nach C via B ist dann wie gewohnt die Route von B nach C + die Kosten von A nach B. A announced das dann wieder an B, der ja nach C über A routet. Er addiert darauf also die Kosten von B nach A. Das läuft dann ungebremst so weiter, bis ins unendliche,
Möglichkeiten dieses Problem zu lösen:

Poisend Reverse Methode Wenn die Route von A nach C über B geht, sagt A dem Host B, dass seine Kosten nach C unendlich sind. In einem kleinen Netzwerk wird dann innerhalb weniger Iterationen ein stabiler State erreicht, in größeren Netzwerken besteht das Problem jedoch immernoch, z.B. in einem Netzwerk in dem A,B,C jeweils direkt verbunden sind, und C dann noch eine Verbindung zu D hat. Alle Kosten sind gleich. Fällt dann die Verbindung CD aus, bekommt A immernoch Falsche Routen zu D von B und umgekehrt.

Split Horizon Wenn Host B seine Routen updatet und das an A sendet und A daraufhin einige Änderungen übernimmt, sendet A diese Änderungen nicht wieder an B, sondern nur an seine anderen Nachbarn.

2.4 Link-State Routing

Beim Link State Routing sammelt für gewöhnlich ein Zentraler Knoten (RCC) Informationen und gibt diese dann an alle anderen Router im Netzwerk weiter. Die Netzwerktopologie ist somit dann allen Router im Netzwerk bekannt. Der RCC baut aus allen gesammelten Informationen einen Grapgen (V,E) , wobei V ein set an vertices (nodes) ist und E für die Edges (links) steht. $c(v,w)$ sind dann die Kosten der Kanten. Wenn eine Kante nicht in E ist, ist c unendlich. Das Ziel ist es dann den günstigsten Pfad von node s (source) zu node v zu finden. Hierfür wird meistens Dijkstras verwendet. Jeder Router versendet regelmäßig per flooding "Link state packages" mit Informationen, die er zu seinen Nachbarn gesammelt hat (delay, hop count...) und verseht diese mit einer sequence Number und einem "age flag". Wenn ein Router so ein Packet bekommt, das er jedoch schon kennt (sequence Number) oder es abgelaufen ist, wirft er es weg.

Vergleich LSR und DVR Link State Routing und Distance Vecotr Routing im direkten Vergleich:

	LSR	DVR
Message complexity	mit n Knoten und E Kanten werden jedes mal $O(n \cdot E)$ nachrichten versendet	Austausch findet nur zwischen nur zwischen Nachbarn statt
Speed of Covergence	ein $O(n^2)$ Alogorithmus braucht $O(n \cdot E)$ Nachrichten	Variiert stark. Es kann zu Routing Schleifen kommen. Count-to-infinity Problem
Robustness	Es können verfälschte Link Kosten announced werden. Jede Router nutzt nur die eigenen Tabellen.	Es können verfälschte Path Kosten announced werden. Die eigenen Routen werden von anderen Routern genutzt.

Algorithmen wie LSR und DVR sind für Netze konzipiert, die sich selten ändern und physikalisch verbunden sind. Sie haben vorallem Schwächen bei Mobilen Netzen, bei z.B. folgenden Punkten:

- **High dynamics** z.B. durch stendig wechselnde Nachbarn und Links
- **Power conservation** regelmäßiges senden von Routing Packeten verbraucht Strom und Leistung
- **Low bandwidth links** wenn z.B. Routing Informationen nicht versendet werden können
- **Asymmetry** Links können in der Geschwindigkeit variieren
- **Interference** Störsignale
- **High redundancy** ein Gerät ist mit vielen anderen Verbunden / "meshed"

2.5 Hierarchical Routing

In der Realität sind große Netze nicht so ideal wie bisher beschrieben, sondern sind meist nicht flach wie wir sie beschrieben haben und Router unterscheiden sich oft fundamental. Im Internet heute gibt es über 1 Milliarde links, die alle zu erreichen würde Routing Tabellen explodieren lassen und der Austausch dieser Tabellen würde jeden Link sprengen. Deswegen ist das Internet in Teilnetze aufgeteilt: Autonomous Systems.

Autonomous Systems Jedes AS hat eine eigene Nummer (z.B. AS421220) und kennt eine Route zu jedem anderen AS. Es gibt im heutigen Internet etwa 60000 solcher Teilnetze und alle sind unterschiedlich groß. Verschiedene AS sind mit physikalischen Links verbunden (peers), über die Daten ausgetauscht werden. Jeder Router innerhalb eines AS muss also nur die Routen zu anderen Routern im AS kennen, und die Route zu seinem Gateway. Innerhalb eines AS (intra-AS) hat der Administrator freie Wahl für die Nutzung von Routing Algorithmen, es kann z.B. RIP, OSPF oder IGRP verwendet werden. Für die Kommunikation zwischen AS (Inter-AS) gibt es jedoch einen festen Standard: BGP (Border Gateway Protocol).

Für die inter-AS Kommunikation wird ein Gateway Router benötigt, der den Transfer von Daten zu anderen AS koordiniert. Dieser Router pflegt Routing Tabellen mit anderen Gatewayroutern. Der Vorteil dieses Modells ist, dass es für das reale Internet besser skaliert. Durch die Reduzierung von Peers, die in Routing Tabellen gepflegt werden müssen, kommt es zu selteneren Updates der Tabellen, was es leichter macht, schnelle Routen zu finden. Bei Inter-AS Routing gibt es jedoch Regulierungen, welcher Host über wen peers darf, innerhalb eines AS wird es sowas nicht geben, da ein AS nur eine begrenzte Anzahl Admins hat. Inter-AS Communication kann deswegen eher Performance orientiert sein.

2.5.1 About BGP and Internet Routing

Zum Vergleich die Routing Protokolle RIP & OSPF:

- **RIP**
 - DVR
 - Erstmal 1983 aufgetaucht. RFC 1058 von 1988
 - Am minimalen Hop-Count orientiert
 - Poison Reverse
- **OSPF**
 - LSR
 - RFC 1131 (inzwischen Version 2 & 3)
 - am meisten verwendete IGP
 - verwendet TCP zum Übertragen von Routing Informationen
 - Multicast support

Border Gateway Protocol BGP ist der heutige Standard für EGP. Durch regelmäßige "hello" Pakete erfahren andere Netzteilnehmer von der Existenz des AS. BGP peers bauen dann eine Session auf und tauschen via TCP Routing Informationen aus. Wenn AS1 z.B. AS2 sagt, dass es einen gewissen Prefix routen kann, garantiert AS1 AS2 dann allen Traffic weiterzuleiten. BGP kann auch innerhalb eines AS verwendet werden, man spricht dann von iBGP. Um den Unterschied dann zu inter-AS Communication zu spezifizieren, wird in dem Kontext dann von eBGP gesprochen.

Eine advertised Route in BGP besteht immer aus einem Prefix und einem Attribut set. Die zwei wichtigsten Attribute sind

- **AS-PATH:** Beinhaltet den Pfad, der für die Route genutzt wird (z.B. AS5 AS8 AS10)
- **NEXT-HOP:** Der nächste AS Router für den nächsten HOP

Über die Jahre wurden immer mehr AS registriert und die Zahl steigt weiter (>6000). Das führt dazu, dass es wieder mehr Hosts gibt, für die alle Routing Einträge gepflegt werden müssen, was die Länge der Tabellen explodieren lässt. Durch die hohe AS Zahl, kommt es auch regelmäßiger zu Änderungen im System, was Änderungen der Routing Tabellen bedeutet. Große Netze haben zudem eine höhere Fehleranfälligkeit.

BGP Security In BGP stellt jeder AS selbst ein, welches Subnetz er verwaltet. Kommt es bei dem Einstellen zu Fehlern ein AS announced ein Subnetz, welches er nicht wirklich organisiert, kommt es zu falschen Routen. So geschehen z.B. 24.02.2008 der Pakistan Telekom (AS17557). Eine generelle Einschätzung zu Routing Security siehe RFC4593. Mit BGP können aber auch Security Operations implementiert werden, z.B. durch Setzen einer IP TTL von 255 und es werden nur Routing Informationen mit einer TTL >= 254 verarbeitet. Außerdem können BGP Sessions mit MD5 Signaturen versehen werden.

2.6 Mobile Routing

Mobile Networking ist gut geeignet, um ein schnelles Netzwerk an Orten ohne gute Infrastruktur aufzubauen. Anders als im Internet sind beim Mobile Routing ganz andere Probleme zu bewältigen. Da zwischen Hosts keine physikalische Verbindung besteht, wird Routing zwischen Hosts schwierig. Doch durch sich wechselnde Positionen o.ä. verändern sich Links zwischen Hosts sehr schnell. Das und andere sind Faktoren, mit denen DVR und LSR nicht umgehen können, diese sind für statische Netze konzipiert. Routing Algorithmen können in zwei Kategorien aufgeteilt werden:

- **Proaktiv**

- Routing Informationen werden unabhängig vom aktuellen Traffic regelmäßig und unabhängig generiert
- Alle Internet Routing Algorithmen sind Proaktiv, auch die oben kennengelernten Beispiele
-
- Beispiel für ein Reactive Mobile Routing Algorithmus: DSDV

- **Reactive**

- Routen werden erst dann berechnet, wenn Daten übertragen werden sollen
- Bei CO wird eine Route (wenn noch keine Vorhanden ist) beim connection Setup berechnet
- Bei CL beim ersten Packet
- Beispiel für ein Reactive Mobile Routing Algorithmus: DSR

Hierarchical Mobile Routing In Mobilien Netzen können Ideen des Mobile Routings verwendet werden. Teilnetze können in Cluster aufgeteilt werden. Innerhalb eines Clusters kann proaktives Routing implementiert werden, für die Kommunikation zwischen Clustern können reaktive Algorithmen verwendet werden. Bei kleinen Clustern können alle Nodes sich gegenseitig kennen, es ist also ideal für reaktive Routing. Zwischen Clustern findet seltener Kommunikation statt, und es können mit vielen kleinen Clustern hier am besten reaktive Algorithmen verwendet werden.

2.6.1 DSDV

DSDV ist eine Erweiterung des DVR spezialisiert für Mobile Netze. Es gibt 2 große Extensions:

1. Sequence Numbers for all Routing updates
 - Verbesserungen gegen Schleifen und Inkonsistenzen
2. Decrease update frequency
 - Zeit zwischen ersten und bestem announcement eines Pfades wird gespeichert

Trotz dessen ist DSDV noch ein Proaktiver Algorithmus und ähnelt sehr stark dem DVR.

2.6.2 DSR

DSR (RFC 4728) geht einen Schritt weiter Richtung Reaktivem Routing. Er baut auf zwei simplen Ideen auf:

1. Routing wird in "Path discovery" und "path maintenance" aufgeteilt
2. regelmäßiges Updaten wird vermieden

Er ist für statische und dynamische Netze geeignet und kann mit bis zu 200 Knoten arbeiten. "Source Routing" bedeutet dabei, dass der Sender für die Bestimmung des Pfades verantwortlich ist.

Path Discovery Wenn ein Sender ein Packet versenden will, jedoch noch keine Route für das Ziel hat, startet er die Path Discovery. Es wird dabei ein Packet per flooding und broadcast mit der Ziel Adresse und einer einmaligen ID versendet. Wenn ein Host ein solches Packet erhält und er ist das Ziel, sendet er das Packet mit dem Pfad über das es gekommen ist zurück. Das erste Packet dass diesen Host erreicht kam über dem schnellsten Pfad, jedes weitere Packet kann dann verworfen werden. Es sind noch weitere Optimierungen möglich: Wenn die Topologie (bzw. die maximale Länge) des Netzes bekannt ist, kann statt einer ID eine counter/TTL hinzugefügt werden. Der Host kann dann am Counter sehen, wie viele Hops das Packet hinter sich hat und es wegwerfen, wenn der Counter größer als der maximale Diameter ist. Eine zweite Verbesserung wäre, wenn Hosts discovery Packete speichern, die sie weiterleiten sollen. Die Informationen dieser Packete können dann für die Suche einer route genutzt werden, wenn das Gerät selbst Packete versenden will. Ist ein Path gefunden, muss jedoch sichergestellt werden, dass dieser auch über die Länge der Verbindung offen bleibt.

Path Maintenance Nicht genutzte Paths werden nach einer Zeit aus der Routing Tabelle gelöscht. Es ist möglichkeit, ist z.B. als Sender auf eine Bestätigung des Empfängers auf layer 2 zu warten oder so eine explizit zu erfragen. Auch kann eine Station schauen, ob Nachbarn das Packet weiterleiten, wenn diese Technologie unterstützt ist. Falls es ein Problem gibt, können Pfade neu gesucht werden oder es kann versucht werden, den Empfänger zu erreichen und ihm mitzuteilen, dass es ein Problem gab.

2.7 Overlay Routing

Im Overlay Routing sitzt ein Virtuelles Netzwerk auf einem existierenden. Hier können die gleichen Algorithmen verwendet werden. Wird ein Packet in einem Overlay network an den Nachbarn versendet, wird das Packet über das darunter liegende Netz transportiert und kann dort auch über mehrere Hops geleitet werden, während Sender um Empfänger im Overlay Network denken, sie sind direkte Nachbarn.

3 IP & Internetworking

Das Internet ist anders implementiert als das ideale Netzwerk. Statt das OSI Modell wird das Internet mit dem TCP/IP Modell in 4 Schichten aufgeteilt:

- **http, ftp,...** Application Layer (OSI 5-7)
- **TCP/UDP** Transport Layer (OSI 4)
- **IP** Internet Layer (OSI 3)
- **Physical** Network Interface (OSI 1-2)

Layer 1 bis 2 (Physical & IP) sind Hop by Hop orientiert, während 2 bis 4 (IP, TCP/UDP & http,ftp) e2e orientiert sind. Im weiteren Verlauf wird es konkret um Layer 2 IP (Network Layer, OSI 3) gehen. Beispiel für Protokolle dieses Layers sind IGMP, ICMP, ARP, IP, RARP.

Sendet ein Host A ein Datenpaket via IP an Host B, versieht A dieses Packet mit einer Destination Adresse. Dieses Packet wird dann von Routern weitergeleitet, bis es an Ziel kommt. Ein Beispiel um Routing mit IP zu erklären:

Host A (192.0.0.1) sendet ein Packet an 192.0.1.1. Host A hat jedoch keine direkte Verbindung sondern ist mit einem Router verbunden, der die Adresse 192.0.0.1 hat. In einer Routing Tabelle dieses Routers gibt es z.B. folgende Einträge: 192.0.0.0/24 über Link Interface 1 und 192.0.1.0/24 über Link Interface 2. Bekommt der Router also ein Packet mit der Destination Adresse 192.0.1.1, so ist diese in der zweiten Range drin, und er sendet das Packet über Link 2 weiter.

Die Information über die Destination IP steht im IP Header.

3.1 IP Header

Tabelle 1: IP Header, length: <- 32 bits ->

VER ¹	IHL ²	type of service	length ³	
16-bit identifier ⁴			flgs ⁴	fragment offset ⁴
TTL ⁵		protocol ⁶	Internet Checksum ⁷	
Source IP (32b format)				
Destination IP (32b format)				
Options (if any)				
Data (variable length, typically a TCP or UDP segment)				

Der Type of Service referiert auf priority Informationen, dieser Teil des Headers wird eher selten benutzt. Die Total Length ist die Summe inklusive Header und TCP/UDP Segment. Der identifier kann genutzt werden, wenn Pakete in mehrere Fragmente gespalten wurden, diese Pakete wieder zusammenzusetzen. Wird ein Packet "fragmentiert" bekommen alle bis auf das letzte Fragment den MF (more Fragments) Flag. Ist ein Packet nicht fragmentiert, bekommt es das DF (Don't Fragment) Flag. Im Fragment Offset wird die Position des Fragments typischerweise in 8er Oktets angegeben.

¹IP protocol Version Nummer (z.B. 4 oder 6)

²Header länge 32b Word

³Totale Länge in Bytes

⁴für fragmentation/Wiederherstellen

⁵Max Hop Anzahl. Wird bei jedem Hop um eins dekrementiert

⁶Protokoll des höheren Layers, dem das Packet zugestellt werden soll

⁷Checksum des Headers

3.2 Addressing

In den bisher kennengelernten Routing Algorithmen speicher der Router Adressen in einer Tabelle, mit dem entsprechenden Pfad, den er dorthin routet. In großen Netzen kann dann so eine Tabelle sehr schnell sehr groß werden. Außerdem ist es nicht unüblich, dass eine Gerät mehrere (MAC) Adressen hat, da eine MAC-Adresse nicht Host spezifisch ist, sondern Interface spezifisch.

	Bispiel	Verteilung
MAC Adresse	70:12:a5:42:93:10	Flach und Permanent (Hersteller)
IP Adresse	172.20.42.10	Topologisch (meistens)
Hostname	tu-darmstadt.de	hierarchisch

3.2.1 Naming and DNS

Bei einem Flat namespacing bestehen Namen nur aus einfachen Strings, die von einer Zentralen Stelle verwaltet werden. Hierbei kann es schnell zu Doppelungen kommen, da in einer langen Liste schnell der Überblick verloren werden kann.

Beim Hierarchical Name Space hingegen wird die Vergabe von Namen dezentralisiert. Es gibt die ersten Namen TLD (Top Level Domain), z.B. de, die dann die Unteradresse tu-darmstadt an einen anderen Host delegiert. Das kann rekursiv weiterlaufen, tu-darmstadt kann dann z.B. Informatik (informatik.tu-darmstadt.de) an einen weiteren Host delegieren, ohne de darüber zu informieren. Wie beim Routing auch zeigt sich, dass flache Strukturen nicht funktionieren, in der Praxis wird deswegen beim namespacing auf das hierarchische Modell gesetzt.

3.2.2 IPv4 Addresses

IPv4 Adressen sind in Blocks aufgeteilt, die für verschiedene Zwecke gedacht sind.

- CLASS A: Organisation ≤ 16 Millionen Hosts. First Bit 0, first 8 for Net, last 24 Bit Host.
- CLASS B: Organisation ≤ 65 Tausend Hosts. First Bits 10, first 16 for Net, last 16 Bit Host.
- CLASS C: Organisation ≤ 255 Hosts. First Bits 110, first 24 for Net, last 8 Bit Host.
- CLASS D: Multicast Adressen. First Bits 1110, last 28 Bits for Multicast Adresses
- CLASS E: Reserviert (Privat, Dokumentation...) Firsts Bits 1111, last 28 reserved.

Eine Netzwerkadresse wird angegeben durch die festen Werte (First Bits & Net) und die Host bits, wobei die alle 0 sind. Die letzte Adresse (Host bits sind alle 1) ist für Broadcast reserviert. 127.0.0.0/8 sind für loopback Anwendungen reserviert, Pakete an diese Adresse werden am lokalen Host bearbeitet. Das /8 eben nennt man eine Netzmaske. Diese gibt an, wie groß das Teilnetz ist, also bei einem CLASS A z.B. /8, weil die ersten 8 Bits fest sind, CLASS B /16 und so weiter. Alternativ schreibt man auch: 11111111.00000000.00000000.00000000 (255.0.0.0, /8).

Einem Host Interface wird in der Regel ein /32, also eine genaue IP Adresse zugewiesen. Ein Host kann mehrere Interfaces, also auch mehrere Adressen haben. Einem Teilnetz wird dann ein Subnetz zugewiesen, in dem sich alle Hosts dieses Netzes befinden. Für lokale Netze werden in der Regel /24 verwendet. Ein Router, der mehrere Netze verwaltet, mehr sich dann nicht jede Einzelnde Adresse, sondern nur die Teilnetze je Interface. Ein Beispielrouter verwaltet 3 Netze auf je einem Interface. Netz 1 läuft über Interface 1 mit 192.0.1.0/24, Netz 2 Interface 2 mit 192.0.2.0/24 und Netz 3 Interface 3 192.0.3.0/24.

Subnetting Ein Größerer IP Bereich kann auch in mehrere Subnetze geteilt werden, so kann von einem /16 z.B. ein /24 abgespalten werden. von den 16 Hosts bits werden die ersten 8 dann zu einer Subnet ID, die das Subnetz identifiziert. Die Netzmaske wird dann für Hosts innerhalb dieses Subnetzes ebenfalls angepasst. Aus einem Netz mit $256 \cdot 256 - 2$ (Broadcast & Network Address) Adressen kann so ein Netz mit 254 Subnetzen je 254 Hosts gebaut werden. Ein Subnetz ist jedoch nicht darauf beschränkt ein /24 zu sein, es können Variable Längen für Subnetze verwendet werden zwischen /17 und /32. In diesem Kontext spricht man von Variable Length Subnet Mask (VLSM). Damit lassen sich Wartungsarbeiten für ISPs z.B. verringern. ein ISP announced z.B., dass er sich um ein /20 kümmert, und leitet dann Teile davon an unterschiedliche Organisationen weiter (z.B. /23 Subnets). Das erlaubt es, IP Adressen mehr effizient zu nutzen. Wenn eine Organisation z.B. 2k Adressen braucht, rein ein CLASS C Netz nicht mehr aus. Es wird ein Class B Netz zugewiesen, auch wenn dann 63k Adressen ungenutzt bleiben. VLSM erlauben es, Adressen mehr effizient zu nutzen. So kann der Organisation z.B. aus einem CLASS B Netz eines ISPs ein /21 zugewiesen werden, dass dann $2^{32-21} = 2048$ Adressen ermöglicht zugewiesen werden.

Advertising Doch wie kommt ein Host am Ende an seine IP Adresse? Koordinierung durch Menschen kann sehr umständlich werden und schnell aus dem Ruder laufen. Es kann ein DHCP eingesetzt werden, der ein zugeteiltes Subnetz an verbundene Geräte verteilt. Ein solcher DHCP Server wird meistens auf dem default Gateway betrieben, und weist IP Adressen anhand der MAC Adresse eines Interfaces zu.

3.2.3 NAT

In der Summe gibt es $2^{32} = 4$ Milliarden IPv4 Adressen. Was zu den 70er noch unfassbar viel war, wird heute knapp. IPv4 Adressen gehen heute für über 20\$ das Stück über den Tisch, bei großen zusammenhängenden Netzen teils noch mehr. Eine Lösung wäre IPv6 (2^{128} Adressen), dafür sind einige Menschen aber zu faul. Eine¹ Lösung für das Problem ist NAT.

Bei einem NAT hat ein Gateway eine öffentliche Adresse einem Gateway zugewiesen. Hinter diesem Gateway können sich mehrere Hosts befinden, die mit Adressen aus einem local Network Bereich versehen werden. Pakete, die aus dem lokalen Netzwerk versendet werden, bekommen dann die Source Adresse des Gateways. Für Hosts außerhalb des Netzes sieht es dann so aus, als käme der Gesamte Traffic nur von einem Host (Gateway).

NAT hat jedoch den Vorteil, dass nicht direkt auf Geräte innerhalb eines lokalen Netzes zugegriffen werden kann, außerdem kann sich die öffentliche IP Adresse ändern (z.B. durch wechseln des ISP), ohne dass das lokale Netzwerk neu konfiguriert werden muss. Jedoch greift der Router durch die Manipulation in des IP Headers in den Traffic ein. Ein Router sollte jedoch maximal 3 Layer bearbeiten. Außerdem funktionieren e2e oder p2p Anwendungen nicht mehr richtig, wenn ein Host innerhalb eines NATs ist. Außerdem ist jeder Host auf maximal 65k connections beschränkt (max anzahl Ports). Bei einem NAT sind alle Geräte jedoch auf insgesamt 65k Ports beschränkt, da das gateway nur soviel Ports hat.

NAT Translation Wenn ein Computer innerhalb eines lokalen Netzes hinter einem NAT Pakete aus dem Netz verschicken will, versteht er jedoch den IP Header mit seiner lokalen Source Adresse. Damit z.B. ein Antwort Packet am Router wieder ankommt muss dieser jedoch seine öffentliche IP als Source IP in den Header eintragen. Außerdem wird der Port geändert, damit die connection am Router mit dem Port des Hosts aufrecht erhalten werden kann. Ankommende Pakete werden dann wieder vom Router statt mit der NAT-Adresse mit der lokalen Adresse des Hosts versehen und innerhalb des Netzes geroutet.

3.2.4 ARP

Wie können Host innerhalb eines lokalen Netzes wissen, welche MAC Adresse hinter einer IP steht? ARP (Address Resolution Protocol) Anfragen! Bei einer ARP Anfrage sendet ein Host via Broadcast an alle Geräte im Netzwerk eine Anfrage: Wer hat IP X? Tell IP Y! Der Host mit IP X fühlt sich angesprochen und antwortet mit seiner MAC Adresse. Host X speichert dann die IP/MAC Kombination in einer Tabelle, bis diese Information abläuft. Es sind jedoch Probleme wie ARP cache poisoning denkbar.

3.3 IPv6

Da IPv4 Adressen inzwischen eng werden, wurde ende der 90er IPv6 (RFC 2460) introduced. Neben dem Ziel, genug Adressen für alle Geräte zu haben und um z.B. NAT los zu werden, gibt es noch weitere Vorteile.

- 'traffic class' und 'flow labels' ermöglichen QoS
- Flexible Anzahl an Routing hierarchie Leveln
- Serverless Plug-and-Play
- Breite e2e und p2p, IP layer Authentication & encryption möglich
- Besserer Support für Mobile Devices

IPv6 Header Im Vergleich zum IPv4 Header fallen einige Felder weg.

Der v6 Header hat insgesamt die doppelte Länge (40 bytes), im Vergleich zum v4 Header. Im Vergleich zum v4 Header ist die gröÙe hier fix, es können jedoch im Next Header viel TCP/UDP oder IPv6 extension Header genutzt werden. Traffic Class ist das Äquivalent zu Type of Service. Hop Limit ist Äquivalent zum alten TTL.

IPv6 Addresses v6 Adressen werden als 8x16 Bit Blöcke in hex Nummern geschrieben, wobei mit :: alles mit 00 gefüllt wird. Beispiel: fd42:4242:f31a:: = fd42:4242:f31a:0000:0000:0000:0000:0000. Es kann auch genutzt werden um Blöcke aufzufüllen, z.B fd42:4242:f31a::abba = fd42:4242:f31a:0000:0000:0000:0000:abba. Die 128 Bits einer IPv6 Adresse sind wie folgt aufgeteilt: Wobei:

¹ Nicht schöne!!!!

Tabelle 2: IPv6 Header, length: <- 32 bits ->

VER		Traffic Class		Flow Label	
Payload Length		Next Header		Hop Limit	
Source IP (32b format, 16 Bytes)					
Destination IP (32b format, 16 Bytes)					

3	13	8	24	16	64
001	TLA ID	Res	NLA ID	SLA ID	Interface ID

- TLA: Top Level Aggregation (IANA teilt diese ISPs zu)
- Res: Reserviert für vergrößerung des TLA oder NLA
- NLA: next-level (Kann zur Organisation von Routing Hierarchie verwendet werden)
- SLA: Site-level (Organisationen können den Bereich für Subnetting nutzen)
- Interface ID: Generierter Suffix für z.B. Endgeräte. Hier kann auch ein v6 Suffix aus der MAC Adresse generiert werden:
 - in 2x24 Bits aufteilen, FFFE Einfügen und 7tes Bit invertieren. z.B. 9042:fe34:3413 → 9242:feff:fe34:3413
 - Diese Methode ist jedoch sehr umstritten, da Geräte so weltweit getrackt werden könne. stattdessen lassen sich auch zufällig generierte Suffixe verwenden.

Bisher ist jedoch kein fester Stichtag für die Umstellung von IPv4 auf IPv6 vorgesehen, und der Übergang scheint sehr schleppend zu sein. IPv6 kann jedoch auch zwischen Netzen übertragen werden, die eigentlich nur IPv4 betreiben. IPv6 Pakete können dann als Payload in IPv4 Paketen versendet werden. IPv6 Routen sind i.d.R. auch in der Lage IPv4 Routen zu handeln.

Further Improvements Die Umstellung auf IPv6 betrifft jedoch auch Protokolle wie DHCP. Das Ziel hierbei ist es jedoch, Adressing und Routing zu vereinfachen. Neighbor Discovery durch ARP Anfragen fällt auch weg und wird durch ICMPv6 ersetzt. Auch hilfreich sind Methoden wie Link-local, wo mit fe80::/10 direkt übertragen werden können, bevor sie global geroutet werden. Es gibt zwar eine v6 Alternative zu DHCP: DHCPv6 es wird jedoch empfohlen Stateless Address Autoconfiguration (SLAAC) zu verwenden.

4 Transport Layer

In Chapter 2 beschrieben, funktioniert Routing von Paketen global am besten abstrakt. Jedoch gibt es Anwendungen und Fälle, in denen man eine direkte Verbindung zwischen 2 Applications auf je einem Host herstellen will. Der Sender teilt Nachrichten der Application in Segmente und schickt diese weiter an den Network Layer, um zum Ziel geroutet zu werden. Im Network Layer des Empfängers werden die Segmente dann zu Nachrichten zusammengebaut und an die App weitergereicht.

Im Transport Layer wird die logische Verbindung zwischen Prozessen hergestellt, während im Vergleich der Network Layer die logische Verbindung zwischen Hosts herstellt. Im Internet sind Beispiele für Protokolle des Transport Layers TCP & UDP.

- TCP
 - Congestion Control
 - Flow Control
 - connection Setup & teardown (COTS)
- UDP
 - connectionless
 - Best-effort Implementierung von IP auf dem Transport Layer

Keines dieser Protokolle bietet jedoch eine Lösung für Delay & Bandwidth guarantees.

4.1 Segmentation(+Blocking). Multiplexing, Addressing

Die theoretische Beschränkung von IP Paketen sind 64k mit Header. TCP Streams sind zum Teil jedoch länger. Um Applications ein Verständnis zu geben, wie Pakete verschickt und SDUs zusammengebaut werden, muss dieser Vorgang transparent geschehen. Pakete können über Multiplexing auf dem Host an alle Applications gesendet werden, und Applications holen sich nur die Pakete, die sie brauchen. Das kommt jedoch mit einem großen Sicherheitsproblem, da nun Anwendungen wie Trojaner einfach den Traffic anderer Programme mithören.

Sockets and Ports Das OSI Modell sieht für die Kommunikation zwischen Programmen und zur Identifizierung dieser die Nutzung von lokalen CEPs vor. In Systemen könnte der System PID zur Identifizierung benutzt werden, der müsste jedoch für alle Hosts erreichbar sein, außerdem ändert sich dieser bei jedem Neustart des Prozesses oder des Hosts. Im Internet Modell hat man sich jedoch für die Einführung von Ports entschieden. Ein Prozess kann auf einen festen Port announce oder nach einem zufälligen Port fragen. Eine Application ist dann durch einen socket (IP:Port Kombination) erreichbar, wobei die IP die des Empfänger Hosts und die Port Nummer die der Empfänger Application ist. Es gibt gewisse Ports, die für Anwendungen reserviert sind a.k.a. "well-known" Ports (IANA: 0-1023). Beispiele sind 22/tcp für ssh, 25/tcp für smtp oder 443/tcp für https. Prozesse die solche Ports öffnen wollen brauchen in der Regel root Privilegien.

Um demultiplexing im Internet zu nutzen, empfängt der Host ein IP Datagram, bestehend aus source und destination Socket. Diese Datagrams transportieren je ein Transport-Layer segment als Application Data. Das Problem hierbei: Wenn IP eines Tages ersetzt wird, funktionieren TCP/UDP nicht mehr, da die Nutzung von IP:Port hard in die Implementierung gecoded sind.

4.2 Connection Control

Connection Control funktioniert natürlich nur bei CO Protokollen, in unserem Beispiel also nur TCP. Wie schon in Kapitel 1 angeschnitten, gibt es 3 Phasen der einer connection:

- CONNECT (herstellen der Verbindung)
- DATA (Transfer von Daten)
- DISCONNECT (Schließen der Verbindung)

Um zwischen Network und Transport Layer Phasen zu unterscheiden, schreibt man dort jeweils ein N- bzw. T- Prefix zu den Phasen also z.B. T-Connect oder N-Data.

4.2.1 Phases

Die verschiedenen Schritte können confirmed oder unconfirmed ablaufen. T-Connect wird immer confirmed, T-Data wird unconfirmed übertragen (in seltenen Fällen kann das auch confirmed laufen) und T-Disconnect können confirmed und unconfirmed ablaufen. Ein Beispiel für ein Ablauf ist auf 4:20 zu finden. Auch ein State-Diagramm, das die Übergänge zu verschiedenen State beschreibt ist auf 4:24 enthalten.

Establishment Beim T-Connect gibt es mehrere Methoden, die aufgerufen werden (können)

- T-Connect.Request(Destination Address, Source Address)
- T-Connect.Indication(Destination Address, Source Address)
- T-Connect.Response(Responding Address)
- T-Connect.Confirmation(Responding Address)

Wobei die Destination Address die Adresse des Transport Service Users ist, also die Adresse die zum Transport gerufen wird. Die Source Address ist die Adresse des aufrufenden Service Users und die Responding Address die Adresse des Antwortenden Service Users.

Data Transfer

- T-Data.req(userdata)
- T-Data.ind(userdata)

Die userdata ist dabei die Payload, die von der Application transportiert werden soll.

Connection Release Der Grund für einen Disconnect kann verschieden sein. Es kann zu einem Release kommen durch abruptes verlieren der Verbdng oder als Folge des Connects. Das verlieren von TSDUs (T + SDU) ist möglich.

- T-Disconnect.req(userdata)
- T-Disconnect.ind(cause, userdata)

Der cause für einen Disconnect kann z.B. ein request vom remote User, Quality of service below minimum, error oder auch unknowm sein. Beispiele für Abläufe verschiedener teardowns sind auf 4:23 zu finden.

4.2.2 Error Control

Error Handling in CONNECT Bleibt ein T-connect.req (CR) unbeantwortet, es kommt also zu einem Timeout, sendet der Host einfach erneut ein CR. Kommt das erste CR jedoch doch nur verspätet beim Empfänger an, und dieser Antwortet auf das erste CR mit einem CC (Connection confirmation) und bekommt dann das zweite CR, sollte durch einbeziehung des Application Layers der zweite CR unbeantwortet bleiben. Was passiert jedoch, wenn das CC verloren geht und wie geht man damit um? Denn der Empfänger (sender of CC) erwartet jetzt einen Datenaustausch. Eine Lösung wäre hier der Three-Way-Handshake

Three-Way-Handshake Beim Three-Way-Handshake Wird das durch ein T-Connect.rsp ausgelöste CC nochmals bestätigt. Das kann entweder direkt ber Data oder per ACK passieren. Das schützt jedoch nicht davor, dass wenn CC und ACK verloren gehen, der Empfänger nicht sagen kann, ob es sich um eine neue oder alte Anfrage handelt. Eine Lösung wäre die Einführung von Sequence Number ins CR, ACK und CC. Host A sendet eine CR (seq=x) and Host B. B antwortet mit einem ACK(seq=y,ack=x) und A sendet dann z.B. DATA(seq=x,ACK=y). Nur wenn die richtigen Sequence Numbers angehängt sind, findet dann der Datenverkehr statt. Anderer Fall. Es kommt ein alter CR bei B an, dieser Antwortet mit der seq Number, da er nicht weiß wie alt er ist. A empfängt ein ACK zu einem alten CR und sendet jedoch ein REJECT, da die Anfrage nicht mehr relevant ist.

Connection Rejection & Release Eine angefragte Verbindung kann jedoch auch Rejected werden. Antwortet ein Empfänger z.B. auf ein CR mit einem DR (Disconnect Request), antwortet der Sender dann mit einem DC (Disconnect Confirm). Normalerweise sendet A ein DR an B, hört jedoch auf eventuelle Daten, die noch unterwegs sein könnten. Er wartet also auf das ACK (DC) von seinem DR und kann währenddessen noch Daten Empfangen. Empfängt er ein DC, kann er sich sicher sein, dass B keine Daten danach mehr sendet. So ein Teardown nennt sich explicit. Ein Implicit Teardown wäre ein Teardown der Network Layer Connection. Das Ziel eines Release ist es, dass beide Seiten alle Daten Empfangen haben und sich nichts mehr mitzuteilen haben. Doch was passiert wenn ein DR/DC/DT verloren geht und nue gesendet werden muss?

Wie kann man sicher gehen, dass A weiß was B weiß und B das weiß und A weiß dass B das weiß...

4.2.3 Two Army Problem

Wie koordinieren sich zwei Armeen, die sich nur Boten durch Tal voller Feinde schicken können und wie stimmen sie ab, dass sie gemeinsam angreifen? Man kann sich nie sicher sein, dass die eigene Nachricht angekommen ist. Das ist das gleiche Problem wie beim connection Release, bei dem ist das Risiko jedoch nicht so hoch, und es können Risiken in Kauf genommen werden. Die Lösung für den Connection Release ist es, einenen Timer parallel zum DR zu starten. falls ein DC/ACK verloren geht, timeoutet die connection automatisch und beide Hosts wissen von der geschlossenen Verbindung.

4.3 Flow Control

Bisher wurden Packete einfach per best Effort versendet. Problematisch wird es jedoch, wenn der Sender eine bessere Verbindung hat als der Empfänger. Wie kann der Sender sicher gehen ob und wenn ja welche Packete angekommen sind? Eine Bestätigung für jedes angekommene Packet vom Empfänger würde diesen eventuell noch mehr einschränken. Wird dieses Verfahren dennoch genutzt, kann der Sender einfach die Packete, die nicht geACKt wurden nach einer gewissen Zeit nochmal versenden, was den Empfänger jedoch nochmal mehr belastet. Eine andere Lösung ist Flow Control. Das Ziel ist es, langsame Empfänger vor schnellen Sendern zu schützen. Flow Control kann auf zwei Layern implementiert werden: Im Link Layer, dort wird einem Überfluss an "forwading segments"vorgegriffen. Auf höheren Layern, z.B. auf dem Network oder Transport Layer wird vor einem Überfluss an Conections geschützt. Auf dem Link Layer ist die implementierung simpel, da Segmente eine feste Größe haben, anders auf dem Transport Layer, da können PDUs stark in der Größe variieren.

Buffer Allocation Ein mit dem Flow Control starj zusammenhängendes Problem ist die Buffer Allocation. Ein Empfänger kann z.B. zwar nicht durch seinen Link eingeschränkt sein, sondern durch das Verarbeiten der Packete. Es wird dann ein Buffer geeigneter Größe erfordert, der eingehende Packete bis zu Bearbeitung/Weiterleitung zwischenspeichern kann. Um als Empfänger die Rate, in der Daten eingehen zu kontrollieren, gibt es verschiene Möglichkeiten. Der Empfänger verlangsamt den Sender, wenn es keinen freien Buffer Space mehr gibt (das kann explizit oder implizit geschehen). Dem kann man vorgeifen, indem der Sender initial Buffer Space anfragt und dieser dann allociert wird, oder der Empfänger announced Ich habe noch X Buffer Space Verfügbar zur Zeit"

Alternating-Bit-Protocol Es wird angenommen, dass es genug freien Buffer Space nach dem Empfangen eines Packets gibt. Ein Empfänger sendet eine Bestätigung für jedes eingehende Packet. Der Sender wartet nach jedem gesendeten Packet auf ein ACK des Empfängers. Der Empfänger sendet das ACK erst, wenn das Packet verarbeitet ist, er kann so also kontrollieren, wann und wie oft Pakete eintreffen. Kommt ein ACK beim Sender nicht an, gibt es vier Fälle"

1. Packet loss
2. ACK-loss
3. ACK late (heavy Traffic)
4. Ack late (Flow control)

Der Sender kann jedoch nicht zwischen den vier Fällen unterscheiden. Nach einer gewissen Zeit sendet er das DT PDU also nochmal. Das Alternating-Bit-Protocol ist in der Theorie eine gute Methode für Flow Control, aber eignet es sich auch in realen Umgebungen? In Einem Beispiel senden wir 8KBit von Frankfurt direkt nach NYC über einen 1Gbit/s Link. Wir gehen davon aus, dass durch die Hosts kein Delay entsteht. Die Propagation durch die Verbindung ist

$$T_{prop} \approx \frac{6200km}{300.000 \frac{km}{s}} \approx 20ms$$

Die Zeit, die das 8KBit Packet braucht, ist etwa

$$T_{trans} \approx \frac{8KBit}{10^9 \frac{Bit}{s}} \approx 8\mu s$$

Weitere nehmen wir an, dass ein Ack etwa $1\mu s$ Transmission benötigt. Bis der Sender also 8Kbit übertragen hat und der Sender ein ACK PDU versendet hat und dieses beim Sender ankommt vergehen also $\approx 40.009ms$. Die Utilization des Links ist also etwa:

$$U_{sender} = \frac{0,008}{40,009} \approx 0.0002$$

Stop and Continue Eine alternative einfache Möglichkeit ist Stop-and-Continue Methode. Der Sender sendet Daten, bis er vom Empfänger ein Stop bekommt. Er pausiert das senden bis zu einem continue Signal. Das Problem ist jedoch, dass es bei einem überlasteten Host schwer werden kann, die Stop nachricht zu versenden, wenn es schon zu einem Overflow kommt. Auch kann es sein, dass der Sender Pakete versendet, obwohl der Empfänger schon eine Stop Nachricht versendet hat, da die Nachricht noch nicht angekommen ist (Delay). Der Empfänger muss das Stop also schon früh genug versenden, um einen Overflow proaktiv zu vermeiden. Diese Methode funktioniert nur bei Full-Duplex Links. Ein Beispiel für die implementierung ist das XON/XOFF Protokoll.

4.3.1 Rate and Credit based Flow Control

Beim Rate based Flow Control können im Vergleich zum Stop and Continue Änderungen in der Rate angekündigt werden. Wenn sich also der maximal verfügbare Buffer ändert kann der Empfänger das mitteilen. Jedoch auch hier gibt es das Problem, dass der Sender erst zu spät davon mitbekommt. Der Übergang zwischen Phasen kann anders als bei Stop and Continue flüssiger ablaufen, da der Sender nicht direkt den Empfang stoppen muss, sondern ihn auch nur einschränken kann. Es gibt jedoch keine Möglichkeit für den Sender sicherzustellen, dass Pakete ankommen.

Idee des Credit Based Flow Control ist es, dem Sender regelmäßig Credits zu geben, um Pakete zu senden. Hat der Sender keine Credits mehr, pausiert er das senden, bis er wieder neue credits vom Empfänger bekommt. Alternativ kann der Sender auch einmalig eine absolute Anzahl bekanntgeben.

Auch lassen sich die Konzepte mit einer Error correction kombinieren. Der Sender kann dem Empfänger "permitter", daten zu senden und er kann das empfangen von Daten acknowledge. Eine ACK Nachricht bedeutet dabei, die Nachricht ist korrekt angekommen. Der Sender kann also seinen timeout-timer sowie die lokale Kopie des Packets löschen. Problem ist jedoch, wenn durch Flow Control das ACK Packet ausbleibt. Wenn der Sender timeoutet, sendet er das Packet dann nochmal.

Eine Protokoll des Credit Based Flow Control ist das Sliding Window Protokoll. Ein Beispiel ist auf 4:46.

4.4 Error Control

Es kann nie garantiert werden, dass alle Pakete ankommen und noch weniger, dass das immer erkannt wird. Wird jedoch erkannt, dass ein Packet verloren gegangen ist, gibt es bei den meisten Protokollen keine Reaktion. Das Packet wird einfach gedroppt. Meistens wird das durch eine fehlerhafte Checksumme im Header erkannt. Da die Checksumme für den IP Header jedoch die gleiche wie für die Payload ist, entsteht die meiste Fehlererkennung an der Stelle, kann jedoch falsch interpretiert werden. Ein Empfänger kann auch beim feststellen dass ein Packet fehlerhaft war ein negative-ACK versenden und die neusendung anfordern. Das hat den Vorteil, dass dann nicht erst der timeout beim Sender auslaufen muss, bis das Packet neu verschickt wird. Wird ein Fehler z.B. durch eine falsche Sequence Number erkannt und ein NACK versendet, kann es jedoch sein, dass ein Packet über einen anderen Weg geroutet wurde, weil ein schnellerer gefunden wurde, und das vorherige Packet noch z.B. im Router hängt. Einzig verlässlich in dem Punkt ist das System: timeout @sender → resend.

Automatic Repeat reQuest (ARQ) Der erste Typ ARQ war das Alternating Bit Protocol. Das lässt sich jedoch weiter verbessern: Go-back-N. Geht ein Packet z.B. mit der Sequence Number 2 verloren oder löst ein Error aus, sendet der Empfänger kein ACK für diesen Frame. Er sendet jedoch weiter die Folge Frames (z.B. 3-8), dessen Frames vom Empfänger Link jedoch weggeworfen werden, da dieser noch auf den Frame mit der Seq. Number 2 warten. Timeoutet dann der Frame 2, wird dieser erneut vom Sender gesendet. der Empfänger empfängt diesen Frame korrekt und sendet dann ein ACK. währenddessen sendet der Sender die Folge Frames von 2 (3-8) ebenfalls nochmals da diese auch kein ACK bekommen haben. (Beispiel siehe 4:49).

Eine Adaption dieses Protokolls wäre, wenn der Empfänger nicht auf den Timeout wartet, sondern für den fehlerhaften Frame ein NACK sendet. Dieser Frame wird dann erneut gesendet, in der Zeit buffert der data link des Empfängers Frames mit der folgende seq. Number.

4.5 Congestion Control

Bei Kommunikation zwischen zwei Hosts gibt es immer ein Bottleneck, dass die Geschwindigkeit beschränkt. Ein Netz ist dann auf genau dieses Bottleneck beschränkt. Werden mehr Daten gesendet, als das Bottleneck übertragen kann, kommt es zu einem congestion collapse. Bei solch einem collapse gehen viele Pakete verloren und die Effizienz sinkt unter das Bottleneck. Je nachdem wo das Bottleneck ist, kann sich das ganze auch in sich selbst steigern, ist z.B. der Router durch den Buffer Space beschränkt und es kommt zu einem collapse wobei der Sender ein Protokoll nutzt, dass nicht angekommene Pakete nochmal sendet, wird der Router mit einer noch größeren Flut an Paketen überschwemmt. Im folgenden sind λ_{in} die Daten die Host A versendet und λ_{out} die, die Host B empfängt. Im Idealfall ohne Packetverlust ist $\lambda_{in} = \lambda_{out}$

Scenario 1 Wir gehen von einem Router mit unbegrenztem Buffer aus und Pakete werden nicht neu versendet. Die Throughput steigt bis zum Maximum an, bleibt dann stabil, der Delay steigt exponentiell an, wenn es zu einer congestion kommt.

Scenario 2 Hat der Router jedoch einen beschränkten Buffer und Host A sendet nicht bestätigte Pakete nach einiger Zeit nochmal. Host A sendet jetzt $\lambda'_{in} = \lambda_{in} + \text{retransmitted packages}$. Perfekt wäre nur wirklich verloren gegangene Pakete neu gesendet werden, λ'_{in} wäre $> \lambda_{out}$. Werden nicht wirklich verlorene Pakete neu gesendet, wird λ'_{in} höher als nötig.

Wie bei den Szenarien gemerkt, wird congestion control essentiell, um einen Schneeball Effekt zu verhindern.

4.5.1 solving congestion Control

Eine Globale Lösung wäre, die Send Rate an den Bottleneck anzupassen. Das kann jedoch nur global passieren und jeder Router müsste seine Maximal Kapazität bekanntgeben. Diese Lösung macht das Problem nur komplizierter und kleine Konfigurationsfehler können das gesamte System betreffen. Deswegen wird Congestion Control zu einem lokalen Problem gemacht. Meistens wird es zwischen Sender und Empfänger gelöst. Das Ziel ist es, so viele Pakete wie möglich zu verschicken, ohne jedoch die Leitung zu verstopfen. Außerdem sollte Congestion Control Fair ablaufen, jedoch unter welchen Bedingungen? Ein 3-Hop Weg ist genauso viel Wert wie 3 1-Hop ways? Priorisierung in Anwendungen (Video Call > Downloads)?

Design Options for Congestion Control Es gibt 2 fundamentale Ansätze, Congestion Control zu gewährleisten:

- Closed-Loop - Der Sender bekommt Feedback über das aktuelle Verhalten und kann die Kapazität anpassen.
 - implementierung am Sender
 - implementierung am Empfänger
- Open-Loop - Das System von vorne herein funktioniert und keine Korrekturen im Betrieb notwendig sind.
 - explicit Feedback - Im Fall wenn es passiert, den Sender informieren
 - implizit Feedback - Der Sender passt die Rate selbst an ohne direktes Feedback vom Empfänger sondern durch Daten aus dem Betrieb (z.B. ausbleibende ACKs). Beispiel: TCP

Possible Actions Um Netze vor Congestion zu schützen, gibt es andere Alternativen als einfach den Traffic zu pausieren. Um häufige Congestions entgegen zu wirken, kann die Netz Kapazität z.B. durch mehr Router erhöht werden. Auch kann ein Router das Allokieren von zusätzlichem Traffic einschränken, wenn das Netz schon fast ausgelastet ist, jedoch sind Informationen über den Netzwerk State selten verfügbar. Ebenso können alle Session einen Teil der verfügbaren Load reduzieren, um allgemein mehr Kapazität zu schaffen. Dafür wird jedoch Netzwerk Feedback benötigt, das geht also nur in closed-loops. Die Klassifizierung der Pakete kann Zentral am Router oder am Host stattfinden. Im Internet geschieht das (außer beim Drop von Paketen) am End-Host. Die einzelnen Klassen bekommen entweder einen gewissen Byte Betrag pro Sekunde (rate based) oder eine gewisse Anzahl an Sequence Numbers/Bytes im Netz verarbeitet werden dürfen (Es gibt auch andere, jedoch unpopulärere Methoden z.B. Credit Based Congestion Control).

Package Dropping Entsteht ein collapse durch einen vollen Router Buffer, müssen Pakete aus diesem gedroppt werden, doch welche? Eine Methode ist die drop-tail queue, bei der das neue Paket weggeworfen wird. Das spielt z.B. gut mit dem go-back-n Protokoll zusammen. Für andere Anwendungen, z.B. Telefon sind neue Pakete wichtiger als alte, hier würden am besten Pakete gedroppt, die schon lange in der Queue sind. Eine solche Aktion ist auch immer, ein implizites Feedback. Der Sender erkennt das nicht ankommen des Pakets. Entsteht eine Congestion durch einen vollen Buffer, muss der Traffic verringert werden.

proactiv Actions Prinzipiell ist ein Router nicht mehr in der Lage, richtig zu funktionieren, wenn er einmal eine volle Queue hat, es sollte also um jeden Preis vermieden werden, z.B. indem man bereits bei z.B. 90% einen warning state aufruft. Eine andere Methode ist das Choke Packet. Ein Router mit vollen Buffer sendet so ein solches Packet an den Sender, dass diesem sagt, die Rate zu verringern. Das Problem ist jedoch, in einem congested Netz gehen mehr Pakete verloren als sonst, es kann also sein, dass der Sender ein solches Packet nie erhält. Auch kann ein Router statt ein extra Packet zu senden, ein warning Bit in allen ausgehenden Paketen setzen, dass allen Hosts sagt, dass der Router überlastet ist, oder das kurz bevor steht. Eine vielleicht etwas radikalere Lösung ist die Random Early Detection (RED), die mit einer Rate zufällig Pakete droppt, auch wenn der Router noch nicht überlastet ist, um diesen genau davor zu schützen. Diese Rate je nach Anzahl Pakete im Buffer variieren.

Feedback Actions Sobald ein Sender Informationen über eine congestion erreicht hat, muss der Traffic reduziert werden. Doch wie das geschieht, hängt von den Protokollen des Transport Layers ab.

4.6 (Internet) UDP

RFC 768 UDP ist eine Best Effort Lösung für Paketübertragung im Internet. UDP ist connectionless, es gibt also keine Handshakes zwischen sender & Empfänger. UDP Pakete können verloren gehen und ohne congestion Control einfach "in die Welt" geblasen werden. UDP wird für z.B. Streaming benutzt, wo es nicht so schlimm ist, wenn ein Packet mal wegfällt, aber auch DNS Anfragen oder SNMP setzen auf UDP. In-Order delivery oder reliability können auf dem Application Layer implementiert werden. Ein UDP Segment ist wie folgt aufgebaut:

source Port #	dest Port #
length	checksum
Application data (message)	

UDP Checksum Das Ziel ist es, Fehler in übertragenen Daten zu erkennen, z.B. geflippte Bits. Der Sender teilt für die UDP Checksum in 16 Bit Sequenzen. Er addiert dann jeweils die Sequenzen zusammen und nimmt davon das 1er Komplement und schreibt das in das Checksum Feld. Der Empfänger muss dann nur die Segmente wieder zusammen addieren und die Checksumme draufrechnen. Ist das Ergebnis 1 (0xFFFF), ist das Packet wahrscheinlich fehlerfrei, ansonsten wird das Packet gedroppt. Ein Beispiel für eine solche Berechnung ist auf Folie 4:79.

4.7 (Internet) TCP

4.7.1 Overview

Im Gegensatz zu UDP hat TCP einige Eigenschaften, die es geeigneter machen für den Gebrauch in den meisten Netzen.

- **Point-to-Point** Es gibt einen Sender und einen Empfänger
- **Reliable, in-order byte stream** Es gibt keine direkten Nachrichten Grenzen
- **Pipelined** TCP Congestion und Flow Control setzt auf window-sized Version
- **Full duplex data** Bi-directional communication in der selben Verbindung möglich. Es gibt eine maximum segment size (MSS)
- **send & receiver Buffer** Es gibt im T-Layer (unter dem Socket zum A-Layer) einen TCP Buffer, der auch die zu sendenden Pakete puffert
- **Connection-Oriented Communication** Vor dem Übertragen von Daten werden Kontroll Nachrichten ausgetauscht (Handshake), die den State initiieren
- **Flow & congestion Control** In TCP sind Implementierungen für diese beiden Ansätze vorhanden

source Port #								dest Port #
sequence Number ¹								
acknowledgement number ²								
head len	not used	U ³	A ⁴	P ⁵	R ⁶	S ⁷	F ⁸	receive window ⁹
checksum ¹⁰								Urg data pointer
Options (variable length)								
application data variable length								

Ein TCP Segment hat dabei den folgenden Aufbau:

Der Socket eines Host, z.B. ein Webserver auf Port :443 kann dabei mehrere connections gleichzeitig haben. Beim TCP (CO) De-multiplexing wird dann zwischen dest und source Socket unterschieden. Dieser TCP hat dann genau 4-Tupel:

- Source IP Address
- Source Port Number
- Dest IP Address
- Dest Port Number

Ein Empfänger nutzt diese 4 Tupel um ein Packet/Segment dem richtigen Socket zuzuweisen. Webserver haben z.B. für jede Verbindung zu einem client einen eigenen Socket. Gewisse Webserver können auch für jede HTTP Anfrage einen eigenen Socket brauchen.

4.7.2 Error & Connection Control in TCP

Sequence Numbers werden bei TCP genauso wie ACKs genutzt. Ein ACK gibt dabei die Sequence Number des nächsten erwarteten Packets an. Ein Beispiel für einen Ablauf ist auf 4:87.

Um TCP timeouts richtig einzustellen, also die Zeit bis zu einem timeout, gibt es einige Faktoren zu beachten. Der Timeout darf nicht kleiner als die Round Trip Time (RTT) sein, sonst kommt es zu unnötigen retransmission, sie darf jedoch auch nicht zu lange sein, sonst kommt es zu unnötig langen verzögerungen beim neusenden von Packeten. Die RTT kann man z.B. berechnen, indem man die Zeit zwischen senden eines Packets und ankommen des ACKs misst. Dieses verfahren nennt man SampleRTT. Da das von Packet zu Packet variieren kann, kann man das Ergebnis durch runden mehrerer Rechnungen glätten.

TCP Connections können aktiv oder passiv hergestellt werden:

- **Active Mode:** Ein Host fragt eine Connection bei einem anderem Host an.
- **Passive Mode:** Eine Anwendung informiert TCP, dass sie auf einem Socket connections annimmt. Es kann jedoch auch Port unspezifisch sagen, dass alle eingehenden Verbindungen akzeptiert werden. Sobald eine Verbindung eingeht, wird ein neuer Socket aufgemacht, der als Endpunkt für die Verbindung gilt.

Der Ablauf einer TCP Connection ist auf 4:90 beispielhaft aufgeführt, außerdem sind auf 4:91 & 4:92 Szenarien aufgeführt, wie sich TCP in einigen Fällen verhält, wenn ACKS ausbleiben, nach timeouts ankommen oder anderes.

Fast Retransmit Geht ein Packet verloren, kann es lange dauern, bis der timeout ausläuft, und das Packet neu gesendet wird. Eine Lösung für dieses Problem ist Fast Retransmit. Erhält ein Sender aufeinanderfolgend 3 mal das gleiche ACK, da zwar nachfolgende Packete angekommen sind, der Empfänger jedoch noch das verlorene Packet erwartet sendet TCP schon vor einem timeout das Packet neu. Das ist guter Kompromiss zwischen nicht unnötig lange warten und nicht unnötig Packete versenden.

¹ Es wird nach anzahl der Bytes gezählt, nicht der Segment

² Es wird nach anzahl der Bytes gezählt, nicht der Segmente

³ URG: urgend data. In der Regel nicht genutzt

⁴ ACK: ACK # ist valid

⁵ PSH: Push Data now. In der Regel nicht genutzt

⁶ RST, connection Establishment/teardown commands

⁷ SYN, connection Establishment/teardown commands

⁸ FIN, connection Establishment/teardown commands

⁹ # Bytes, die der Receiver bereit ist, zu akzeptieren

¹⁰ Internet Checksumme (vgl. UDP)

Send and Receive Buffers in TCP TCP pflegt jeweils auf Sender und Empfänger Seite einen Buffer um beim Sender Packete zu speicher, die z.B. noch nicht geack't sind, um diese ggf. neu zu senden oder Daten der Application zu speichern, bevor diese versendet werden. Ältere Versionen von TCP haben nach dem Go-Back-N Prinzip gearbeitet und out-of-order Packete weggeworfen. Inzwischen werden auf Empfänger Seite Packete, die neuer als erwartet sind zwischengespeichert, um nicht alle neu senden zu müssen.

4.7.3 Flow & congestion Control in TCP

In TCP kann der Empfänger den advertise size von seinem Buffer announce. Der Buffer size ist dabei:

$$\text{Advertised window} = \text{MaxRcvdBuffer} - ((\text{NextByteExpected} - 1) - \text{LastByteRead})$$

der advertiste Speicher beschränkt dabei den Sender in der Anzahl an Daten, die dieser verschickt. Der Sender achtet dabei darauf, dass

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{AdvertisedWindow}$$

bzw.

$$\text{EffectiveWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

Self Clocking Winwo Bisher gingen wir davon aus, dass neue Packete direkt verschickt werden, wenn ein ACK eingetroffen ist. Problem jedoch, wenn ACKs schneller einkommen, als die Application Daten in den send Buffer schreibt, werden viele kleine Packete verschickt. Man spricht dabei vom »silly window syndrom«. Eine Lösung hierfür ist Nagle's Algorithmus. Wenn die Verfügbaren Daten und das advertiste Window \geq MSS sind, wird ein volles Segment verschickt. Ansonsten wird geschaut, ob bereits an Packet verschickt wurde, das noch ncith geACKt wurde, dann wird gewartet bis MSS voll ist. Falls kein Packet verschickt ist, werden die neuen Daten sofort verschickt.

Congestion Control TCP kontrolliert arbeitet mit implicit feedback um congestion zu verhindern, konkret mit der Information über dropped Packages. Es kann weiter nicht unterschieden werden, warum Packete gedroppt werden, deswegen wird davon ausgegangen, dass diese wegen congestion gedroppt wurden. TCP ermittelt mit diesen ein congestion window, dass im laufenden Betrieb reflektiert und aktualisiert wird. Zusätzlich zu den Limits der Flow Control beschränkt der Sender die Anzahl an Packeten durch

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{CongWin}$$

Um die darstellung von congestion Control zu vereinfachen, wird in Zukunft das flow Control Window ignoriert.

Hat TCP jetzt genau so viele Packete in das Netzwerk gesendet, wie das congestion Window erlaubt, wird gewartet, bis diese das Netz wieder verlassen haben. Durch das erreichen von ACKs ist sichergestellt, dass die Packete nicht mehr im Netz sind und es können neue gesendet werden.

Congestion Window (AIMD) Bekommt TCP ein ACK, geht es nach Greedy von der Annahme aus, dass das Netz noch Kapazität hat und erhöht das congestion window. Bleiben jedoch ACKs aus und es kommt zu einem timeout, wird das congestion window verkleinert. Doch um wieviel? Um das Netz vor einem congestion collapse zu schützen, müssen drastische Maßnahmen getroffen werden, und das congestion window wird auf 50% gekürzt (halbiert)! Das Window kann jedoch nie kleiner als 1 sein.

Das kürzt jedoch die mehr als nötig, wenn Packete nicht wegen congestion, sondern z.B. durch Übertragungsfehler gedroppt werden. Das passiert in physischen Netzen selten, ist jedoch grade bei Mobilien Netzen ein Problem. Das ganze nennt sich multiplicative decrease.

Wird jedoch das Congestion Window durch ACKs erhöht, wird es nicht so stark erhöht, da es sonst zu einem collapse kommt. TCP testet sich langsam an das Maximum heran und erhöht es immer nur um einen kleinen Teil. Wenn alle gesendeten Packete innerhalb der RTT (Round Trip time) ein ACK erhalten haben, sendet TCP in Zukunft ein Packet mehr pro RTT.

Das ganze kann auch per additive increase passieren, dabei wird das Congestion Window pro erreichtes ACK um folgenden Wert erhöht:

$$\text{Congestion Window} += \text{MSS} \times (\text{MSS} / \text{Congestion Window})$$

TCP setzt insgesamt auf ein AIMD (adaptive increase multiplicative decrease) "Sägezahn" Pattern. (Beispiel siehe 4:105). AIMD start jedoch mit einem congestion Window von 1 oder 2 und brauch länger, um anzulaufen, was gerade bei schnellen Netzen unnötigt Zeit kostet. Es kann also anfangs das congestion Window z.B. je ACK verdoppelt werden, bis es einmal durch loss halbiert wird. Ab diesem Punkt wird dann auf linearen anstieg umgestellt.

Packet Burst Ein Problem jedoch: Ein Packet geht verloren. Nach einem timeout wird das Packet neu gesendet und es kommt ein ACK zurück. Pakete wurden seit dem timeout zurückgehalten und das congestion window zwar halbiert. jedoch wird jetzt auf einmal das gesamte (halbierte) Congestion Window auf einmal ausgeschöpft und ein "Packet Sturm" losgeschickt. Dadurch wird es wahrscheinlich zu einem hohen Packetverlust kommen. Die Lösung ist, in dem Fall wieder slow zu starten mit einem congestion window von 1. Jedoch gibt es jetzt eine Idee von der maximalen Größe des Congestion Windows. Dieser wird kann also congestion threshold genutzt werden. Es wird also mit dem exponentiellen slow start bis zum threshold increased und dann auf additive increase gewechselt.

Ein Beispiel Ablauf ist auf 4:111.

4.8 Other Protocols

5 Queueing Theory

Beim Queueing gibt es 5 wesentliche Punkte:

1. Arrival Process: How do customers/requests arrive?
 - Typically described as probability distribution of interarrival times
 - $A(t) = P[\text{time between arrivals} \leq t]$

6 Referenzen

6.1 Quicktour

- TDD (Time division duplex)
- TDM (Time Division Multiplexing)
- FDM (Frequenz Division Multiplexing)
- WDM (wavelength divisionmultiplexing)
- CDM (Code Divison Multiplexing)
- SDM (Code Divison Multiplexing)
- DS (Distributed Systems)
- AS (autonomous System)
- CSS (communication subsystem)
- CN (computer networks)
- OSI (Open Systems Interconnection)
- CON (Connection Establishment)
- DAT (Data Exchange)
- DIS (Disconnect)
- PDU (Protocol Data Unit)
- PCI (Protocol Control Information)
- SDU (Service Data Unit)
- LAN (Local area network)
- LLC (Logical Link Control)
- MAC (Medium Access Control)
- CO Network (connection-orientated network)
- CL Network (connectionless network)
- e2e (End-to-end)

6.2 Routing

- NS (network layer services)
- TS (transport service)
- TTL (Time to live)
- RCC (Routing Control Center)
- DVR (Distance Vector Routing)
- LSR (Link-State Routing)
- AS (autonomous System)
- RIP (Routing Information Protocol)
- OSPF (Open Shortest Path First)
- IGRP (Interior Gateway Routing Protocol)
- BGP (Border Gateway Protocol)
- IGP (Interior Gateway Protocol)
- EGP (Exterior Gateway Protocol)
- DSDV (Destination Sequenced Distance Vector)
- DSR (Dynamic Source Routing)

6.3 IP & Addressing

- VLSM (Variable Length Subnet Mask)
- DHCP (Dynamic Host Configuration Protocol)
- NAT (Network Address Translation)
- SLAAC (StatelessAddress Autoconfiguration)

6.4 Transport

- CEP (Sockets)
- COTS (connection-oriented Transport services)
- ARQ (Automatic Repeat reQuest)
- MSS (Maximum Segment Size)
- RTT (Round Trip time)
- AIMD (adaptive increase multiplicative decrease)