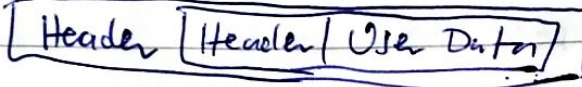


Telematik

2 Internet Congestion Control

2.1 Basics

Packet Switching



forwarded independently

General problem: n users on 1 resource for data transferring

Critical Situations: 2 packets arrive @ same time.

Buffer / drop?

Different speeds, data incoming faster than being able to send

First buffer packets, then drop

" fills up, Latency

End-to-end-Latency: propagation delay + transmission delay + queuing delay

TCP: 3-way handshake for connection establishment

4-way for termination (half-duplex possible)

Full duplex connection

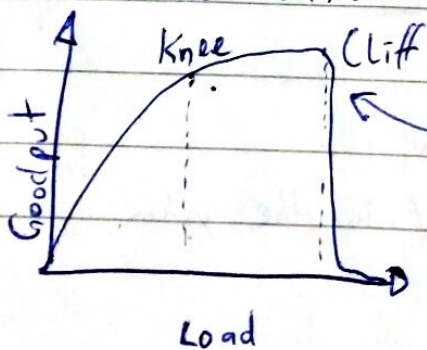
Byte oriented ~~that~~ seq numbers

Go back N

positive cumulative acks

Timeouts

Flow control via sliding window



Network is congested, Buffers full
Goodput = Throughput - Overhead data

Flow control only applies to the capacity of the receiver
⇒ congestion control necessary due to possibility of multiple TCP connections on the same resource

Goal: ~~n~~ rates @ knee

Feedback from the routers necessary

Congestion avoidance and congestion control: avoid traffic @ cliff

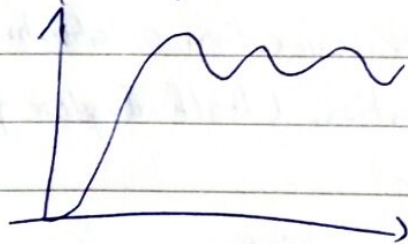
Efficiency: $C \hat{=}$ Capacity of resources

$$\sum_{\text{users}} \text{data rate} = C$$

Fairness: $F(r_1, \dots, r_n) = \frac{(\sum r_j)^2}{N(\sum r_j^2)}$ Jain Fairness Index
between 0 and 1

1 $\hat{=}$ fair, equal data rates $1/N$ unfair

Convergence responsiveness, rate at which speed a connection gets to equilibrium



I Smoothness, oscillations around equilibrium state at steady rate

Distributedness

No central control with knowledge about network state

Window based congestion control

$Cwnd \hat{=}$ Congestion control window

* unacknowledged packets allowed in the system.

Rate based

UDP, controlled by precise timers

Others: hop-by-hop, end-to-end, router based, centralized

Implicit due to timeouts, duplicated acks

Explicit, indicated by the network nodes, ECN

TCP Tahoe

Mechanismen: Slow start (bei Timeout)

Fast retransmit (bei dup. ACKs)

Timeout

Congestion avoidance

Es muss immer gelten

$\text{LastByteSent} - \text{LastByteAcked} \leq \min \{ \text{wnd}, \text{RcvWindow} \}$

RcvWindow $\hat{=}$ Empfangsfenster (Flusskontrollfenster)

SSThres, wann wird von Slow Start zu

Congestion avoidance übergegangen?

Additive Increase, Multiplicative Decrease

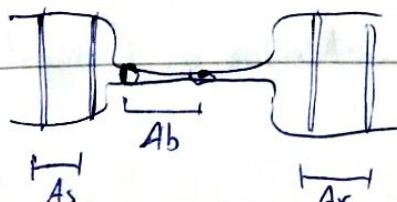
Timeout oder 3 dup ACKs: Slow retransmit

$\text{SSThres} = \max \{ \text{Flight Size} / 2, 2 * \text{MSS} \}$

\neq unacknowledged packets

Fast retransmit bei der 3. Duplizierung (4. ACK)

Self-clocking



P_r : Zwischenankunftszeit zweier Pakete

A_r : Ankunftszeit zweier ACKs

b : Bottleneck

r : receiver

s : sender

Sender sendet Pakete so schnell er kann, der Empfänger sendet ACK so schnell sie ankommen (bei Bottlenecks langsamer), woraufhin der Sender nur in diesem Intervall neue Pakete senden darf \rightarrow Bottleneck bestimmt Geschwindigkeit der Übertragung

3 ways to fail conservation of packages:

1. Connection does not get to equilibrium
2. Sender injects new Package before old has left (\Rightarrow \neg conservation!)
3. Resource along path hinder equilibrium

Slow Start

bei vermuteter Stausituation oder neuer Verbindung erhöhe Cwnd um 1 für jedes erhaltene ACK, bis Ssthresh erreicht ist \rightarrow exp. Wachstum

Maximum segment size (MSS) bezieht sich auf die Nutzdaten (in Bytes)

Go-Back-N: Sendewiederholung ab erstem unbestätigten Segment, Bursts von wiederholten Segmenten

Retransmission Timer

Alternative zum erniedrigen der Flight size (außer per ACK, Paket hat System verlassen)

Basiert auf der Round-Trip-Time (RTT)

Estimated RTT basiert auf mehreren Messungen (Samples)

$$\text{Estimated RTT} = (1 - \alpha) * \text{Estimated RTT} + \alpha * \text{Sample RTT}$$

α typischerweise 0,125

Zeitgeber für Senderwiederholung

$$\text{RetransmissionTimer} = \beta * \text{Estimated RTT}$$

β oft = 2

Schwankungen der RTT durch sich ändernde Last im System

$$\text{Deviation} = (1 - \gamma) * \text{Deviation} + \gamma * |\text{Sample RTT} - \text{Estimated RTT}|$$

$$\text{RetransmissionTimer} = \text{Estimated RTT} + \beta * \text{Deviation}$$

γ oft = 0,25

Mehrfache Senderwiederholungen

Wie groß ist der Zeitraum zwischen zwei Senderwiederholungen des gleichen Segments?

Exponentieller Backoff:

RetransmissionTimer $*=2$, max 60sec.

Kahn's Algorithmus bei Quittierung eines wiederholten Segments:

Estimated RTT & Deviation nicht neu berechnen.

Backoff weiterhin verdoppeln, bis Quittung eines nicht wiederholten Segments eintrifft. Dann alten Algo wieder nutzen.

Unnötige Senderwiederholungen vermeiden \Rightarrow Timer müssen akkurat sein.

Packet loss highly likely due to new TCP connection.

Reduce load of "our" active connection.

Feedback control Algorithm AIMD

applied to congestion control

$$r_j(t+1) = \begin{cases} r_j(t) + a & \text{if no congestion is detected} \\ r_j(t) \times b & \text{if congestion is detected} \end{cases}$$

$$a \in \mathbb{N}, b \in [0, 1] \subset \mathbb{R}$$

Assuming all users receive the same feedback from the bottleneck link \Rightarrow equal resource utilization
 1 MSS equal

Fairness of two connections:

$$F(r_j, r_i) = \frac{(r_i + r_j)^2}{2(r_i^2 + r_j^2)}$$

Optimal operating point: $(c/2, c/2)$

Gieriger Nutzer öffnet mehrere Verbindungen gleichzeitig oder sendet mehr ACKs, bzw. schneller als die Daten wirklich eintreffen

Dies wäre nicht möglich, wenn das Sluickontrollfenster in Bytes zählen würde. Es zählt in ~~Bytes~~ Segmenten

Quittungen fälschen wäre unrentabel, wenn Sender Pakete auslöst. Die dann gelagerten "angekommenen" Daten würden dann ~~fehlen~~.

Appropriate Byte counting

Zähler für quittierte Datenmenge

bytes_acked += ~~Bytes~~ Bytes in Quittung quittiert

if bytes_acked \geq cwnd [Bytes]:

cwnd += 1MSS

bytes_acked -= cwnd [Bytes]

Dies hilft gegen Senden mehrerer Quittungen pro Segment

TCP Reno = TCP Tahoe + Fast ^{Recovery} ~~Retransmit~~
bei leichter Stausituation (Dup. ACKs) kein Slow Start

Fast Recovery

Sendwiederholung des Segments (Fast Retransmit)

$$SSThresh = \max \{ FlightSize / 2, 2 * MSS \}$$

$$Cwnd = SSThresh + 3 MSS$$

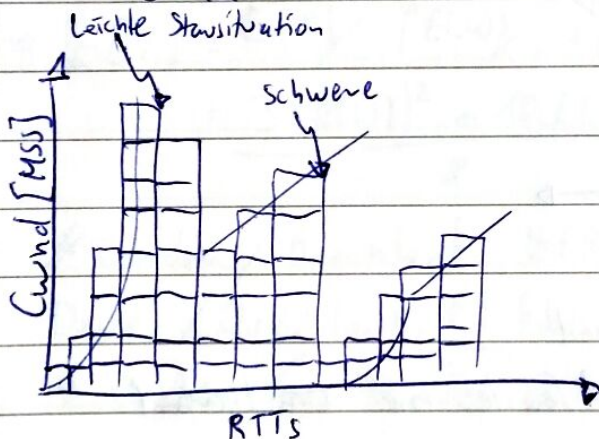
Bei Empfang duplizierte Quittung

$$Cwnd + 1$$

Sende (falls vorhanden) weitere noch nicht gesendete
(neue) Segmente

Bei Empfang neuer Quittung:

$$Cwnd = SSThresh$$



Periodic Model

Zur Analyse von TCP (beispielsweise)

Ist TCP für meine Zwecke geeignet?

Welche Leistung / welches Verhalten kann ich erwarten?

Betrachtung eines langfristigen Verhaltens (SS wird ignoriert)

RTT konstant, Verluste treten periodisch auf.

Kein queuing delay, kein Warten auf timeouts

⇒ Staukontrollfenster bleibt die perfekte periodische Sägezahnkurve.

Variablen: X : Übertragungsrate, Segmente/Zeit

W : Größe des Staukontrollfensters bei Paketverlust

N : # Segmente

A : Dauer eines Durchlaufs

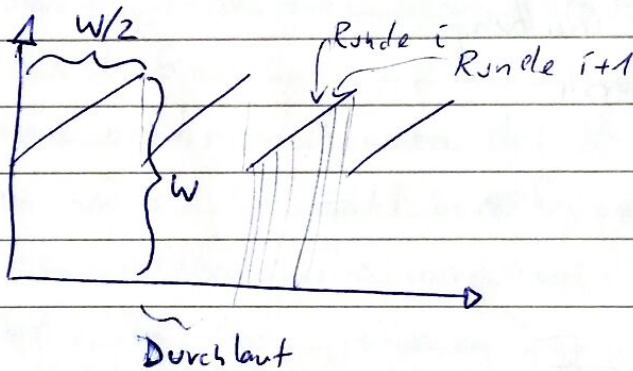
RTT: Umlaufzeit

p : P(Verlust eines Paketes)

MSS: Maximum Segment Size

E : Erwartungswert

O : Größe der zu übertragenden Daten



W_{cur} = aktuelle Größe des cwnd

Annahme: Flusskontrollfenster nie limitierend

Größe minimales Staukontrollfenster: $W/2$

Wachstum um je 1

Länge der Durchläufe: $W/2 * RTT$

Segmente pro Durchlauf: $3/8 W^2$

Es gilt $N = 1/p$

⇒ $W = \sqrt{8/3p}$

Durchschnittliche Übertragungsrates: $1,22 * \frac{1}{RTT \sqrt{p}}$

Bsp. RTT = 200ms, $p = 0,05$ ⇒ $X = 27,4$ Segmente/s

Detailliertes Paketverlustmodell

Verluste nicht mehr periodisch

Durchläufe nicht gleich lang

$\alpha_i \hat{=}$ Nummer des verlorenen Segments in Runde R_j des i -ten Durchlaufs

In $R_j + R$ Runden: $Y_j = \alpha_j + W_j - 1$ Segmente gesendet

$$\Rightarrow E[Y] = E[\alpha] + E[W] - 1$$

Erwartete Position des verlorenen Segments:

$$E[\alpha] = \sum_{k=1}^{\infty} k(1-p)^{k-1} p = 1/p$$

$$\Rightarrow E[Y] = 1/p + E[W]$$

$$E[W] = 2 * E[R] - 2$$

* Segmente im Umlauf: $Y_i = \frac{R_i}{2} \left(\frac{W_{i-1}}{2} + W_i \right) + \beta_i$

$\beta_i =$ * Segmente im Umlauf in der letzten Runde

$$\Rightarrow E[Y] = E[R] \frac{E[W]}{2} + E[\beta]$$

$$= \frac{3(E[W])^2 + 10E[W]}{8}$$

$$\text{* Runden / Durchlauf: } E[R] = \frac{5}{6} + \sqrt{\frac{2}{3}p - \frac{23}{36}}$$

$$\text{Dauer / Durchlauf: } E[A] = RTT * \left(\frac{11}{6} + \sqrt{\frac{2}{3}p - \frac{23}{36}} \right)$$

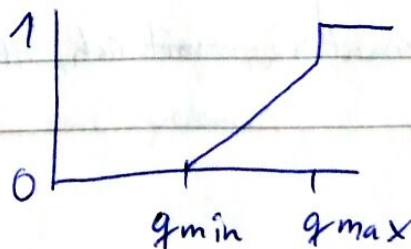
$$\Rightarrow X(p) = \frac{1}{RTT} * \sqrt{\frac{3}{2}p}$$

Im Beispiel: $RTT = 200 \text{ ms}$, $p = 0.05 \Rightarrow X(p) = 23,89 \text{ Segmente/s}$

Active Queue Management (AQM)

RED (Random early detection) $\hat{=}$ verwirft zufällig ein Paket, falls sich die Buffer zunehmend füllen.

Drop probability:



ECN (Explicit Congestion Notification)

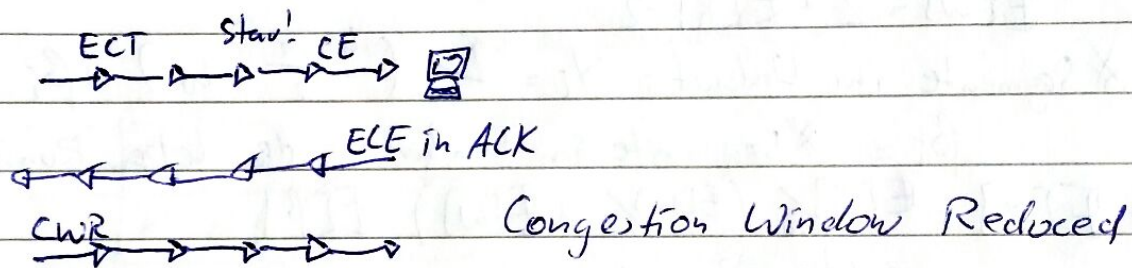
Leite Paket normal weiter, allerdings kann eine ECN gesetzt werden, sollten die Buffer sich zunehmend füllen. Diese ECN lässt den Sender wissen, dass er sein Cwnd zu reduzieren hat.

Dafür werden 2 bits im ToS field genutzt

00 Not-ECT, 01 ECT(1) 10 ECT(0) 11 CE

CE $\hat{=}$ Congestion Experienced ECT $\hat{=}$ ECN Capable Transport

Fähigkeit über ECN wird beim TCP Verbindungsaufbau geklärt.



3. Router

Auf Schicht 3, Vermittlungsschicht.

Dienen zur Wegwahl im Kommunikationssystem, umgesetzt mit Hilfe der Weiterleitungstabelle im Router

leitet Pakete in Richtung Ziel weiter

tauscht Informationen mit anderen Routern aus, nutzt dazu Routingprotokolle

\Rightarrow verteiltes & adaptives Routing

Erstellt Routing- und Weiterleitungstabelleneinträge

Je Paket ein Lookup

Pakete führen Informationen mit sich, IP Adresse, ...

Basiskonkationen:

Überprüfen des IP-Headers

- Versionsnummer

- Länge gültig

- CRC / Prüfsumme

- Lifetime check TimeToLive --;

Prüfsumme neu berechnen und setzen

Lookup, Suche nach dem korrekten Ausgangsport

Fragmentierung

Behandlung von IP-Optionen

Klassifizierung / Priorisierung

Line Speed

Routers need to keep up for every packet size

@ 1000bit/s dauert ein ACK 3,2ns (40 bytes)

$$\frac{320 \text{ bit}}{10^{11} \text{ bit/s}} = 3,2 \text{ ns}$$

~ 50% of all packets transmitted are < 100 bytes long

Core Router

used by service providers

need to handle large amount of data / traffic

fast lookup & forwarding

redundancy for reliability, cost secondary issue.

Enterprise Router

connects end systems in companies, universities, ...

large number of end systems

Support for firewalls, VLAN
 Low cost per port, large number of ports

Edge router / access router

At edge of service provider
 provide connectivity to customer from home, small business
 support for PPTP, IPsec, VPNs

Forwarding

done by prefix matching

largest select the most specific match

can be done in software

efficient data structures required

~~can be done~~ to perform a fast lookup

should be memory efficient. ~~Core~~ ^{Core} routers have

> 500k entries

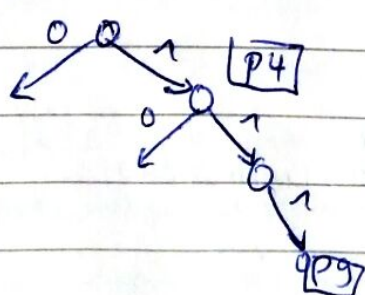
Fast updates >> 100/sec since route changes occur frequently

$N =$ # prefixes

$W =$ length of a prefix

$k =$ length of a stride (only for multibit tries)

Binary Trie



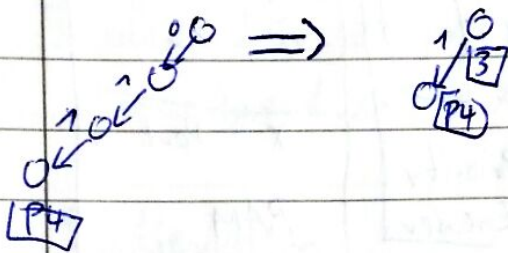
Lookup ~~O(N)~~ $O(W)$

Memory $O(W \cdot N)$

Update $O(W)$

Path compression

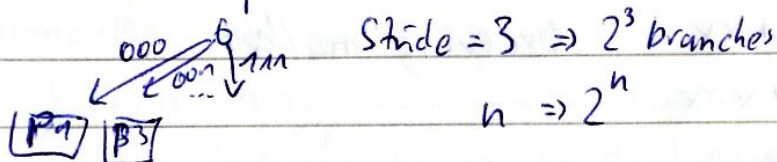
Long sequences of one-child nodes waste memory
store the position of the ~~m~~ bit that next decides which branch should be taken.



lookup $O(w)$
memory $O(N)$
updates $O(w)$

Multibit Trie

match multiple bits at a time instead of a single bit



lookup $O(w/k)$
memory $O(2^k N w / k)$
update $O(w / k + 2^k)$

Hash Tables

can perform lookup in $O(1)$

Used to store the ~~trie~~ trie lookup result

Does an entry for a destination exist in the hash table?

yes \rightarrow no trie lookup ($O(1)$)

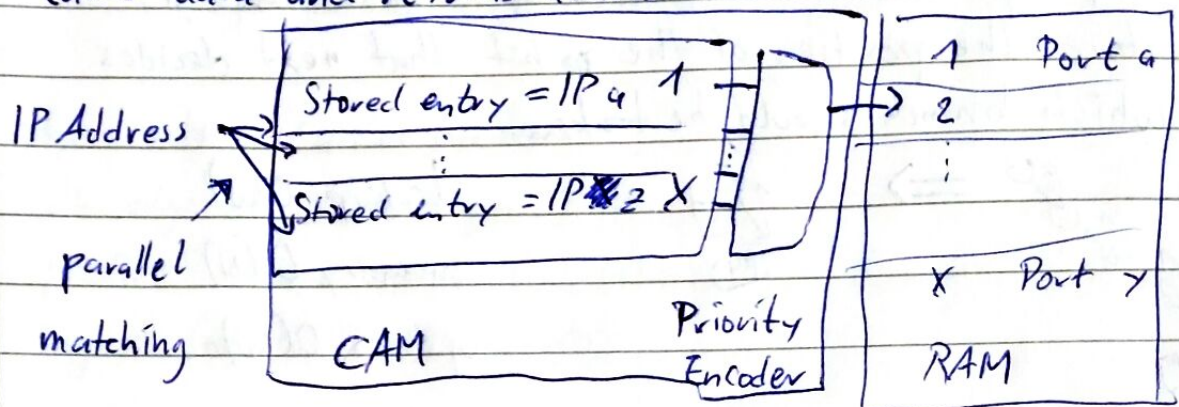
no \rightarrow trie lookup

Longest Prefix matching in Hardware

IP address as RAM address \rightarrow doesn't scale well, unused RAM

Content Addressable Memory (CAM)

takes data and returns address



Ternary CAM (TCAM)

Supports Don't Care state of bit

10011XX Ordering matters

1011XXX

0XXXXXX

High energy consumption

2-3 times transistors like SRAM

strict ordering, updating values requires re-ordering

Router - Architektur

Netzwerkschnittstellen

realisieren den Zugang zum Netz

Schicht 1 & 2

Routing-Prozessor

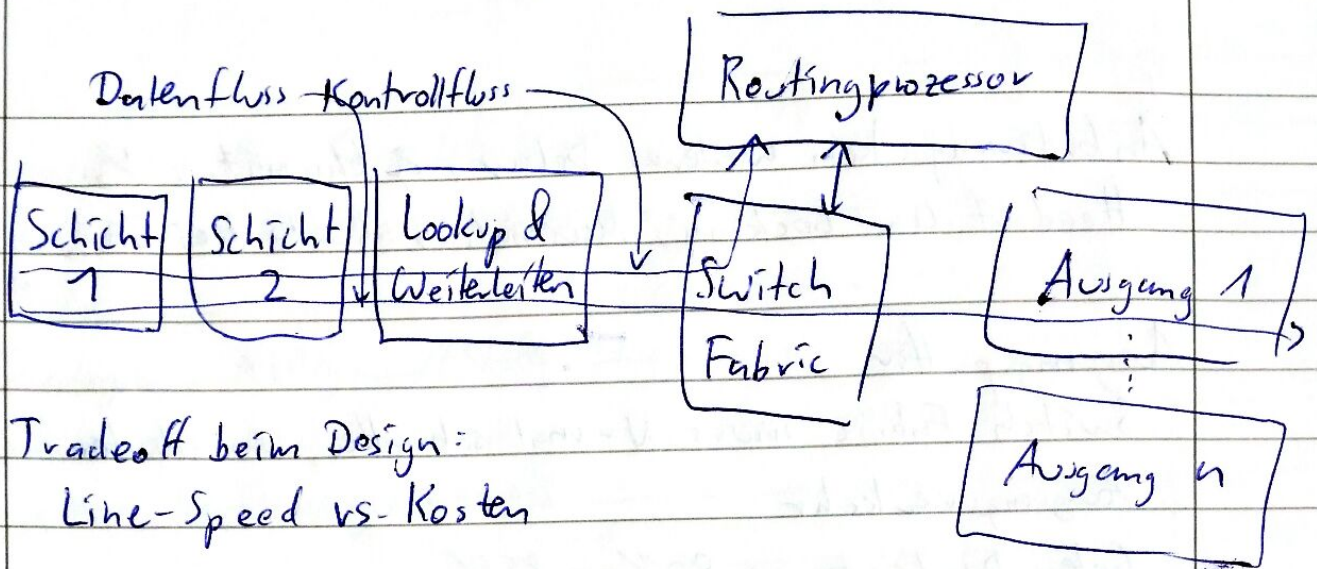
Routing Protokoll

Management Funktionalitäten

Switch Fabric

Backplane

Realisiert die interne Weiterleitung



Tradeoff beim Design:
Line-Speed vs. Kosten

Blockierungen

n Pakete an n Eingängen, die alle den selben Ausgang nehmen wollen.

Vermeidbar durch

Overprovisioning

Interne Verbindungen in Switch-Fabric optimieren mit höherer Geschwindigkeit als die Eingänge

Pufferung

Pakete zwischenspeichern, bis die Resource frei ist.

Backpressure

Signalisierung der Überlastung am Eingangsport
Daraufhin die Last reduzieren

Parallele switch-fabrics

ermöglicht parallelen Transport zu Ausgangsports
benötigt höhere Geschwindigkeit als ~~die~~ an den Ausgängen

Platzierung von Puffern

Eingangspuffer

FIFO oder andere Schedules, Round Robin, Prioritäten

Anforderung: kein weiteres Delay, Zykluszeit = $\frac{1}{2}$
Head of Line blocking, Maximaldurchsatz bei 58%

Ausgangspuffer

Switch Fabric muss N -mal schneller sein als der

Ausgangsverkehr

Guter Durchsatz $\approx 80\% - 85\%$

Verteilte Puffer

Switch Fabric als Matrix

Zykluszeit = $\frac{1}{2}$

Kein Head-of-Line-Blocking

Hoher Speicherbedarf

Zentraler Puffer

Alle Eingänge und Ausgänge mit einem Puffer verbinden.

Organisation: RAM

Zykluszeit = $\frac{1}{2N}$

Adress- und Steuerspeicher für Steuerung von parallelen
Zugriffen

Niedrige Zugriffszeit nötig, Speicher dafür vergleichsweise
klein

Switch Fabric

Gemeinsamen Speicher

Bus/Ring Struktur

Koppelmatrix

Mehrstufige Verbindungsnetze

typische

Grundstrukturen

Bewertung der Grundstrukturen

Das interne Blockierungsverhalten

Puffer

Topologie, # Wege, # Stufen

Steuerungsprinzip

interne Verbindungskonzepte

Bus / Ringstruktur

Konfliktfreier Zugriff durch Zeitduplex

Nötige Kapazität = mindestens Summe aller Eingänge

Einfacher Multi / Broadcast

Geringe # Ports ~ 16



Koppelmatrix

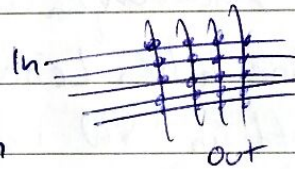
Vollvermaschung der Ein/Ausgänge

Teilweise parallele Vermittlung möglich

Puffer notwendig

Skaliert quadratisch

Eignet sich bei gleichgroßen Paketen



Mehrstufige Verbindungsnetze

Verdrahtungsaufwand geringer als bei Koppelmatrix

Jeder Eingang kann mit jedem Ausgang verbunden werden

Blockierungen möglich

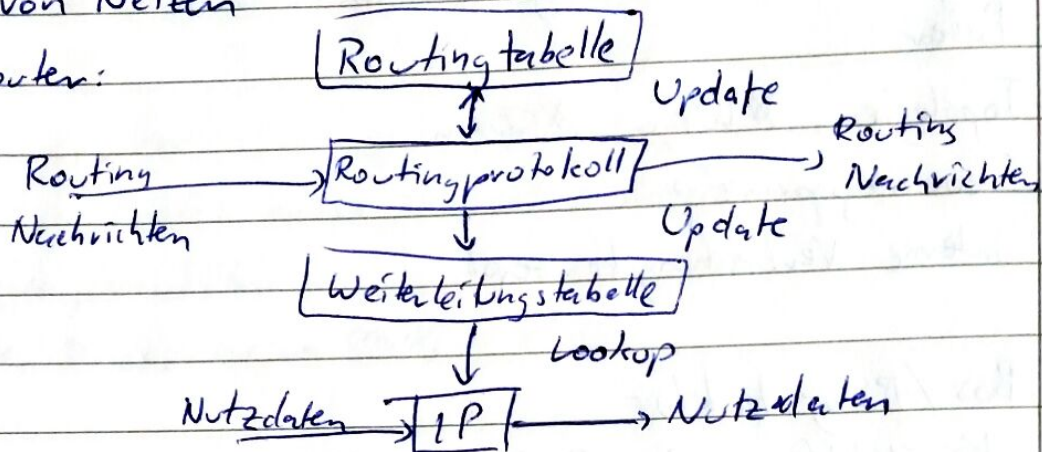


Banyan Netz

4 Internet routing

Netz von Netzen

IP Router:



Routing tabelle:

Von den Routingalgorithmen erzeugt

Einträge: Abbildung von IP Präfix auf next-hop (IP Adresse)

Optimiert für die Anforderungen der Routingalgorithmen
Änderungen der Topologie

Keine hohen Leistungsanforderungen
oft in SW realisiert

Weiterleitungstabelle:

Bei der Weiterleitung von Paketen genutzt

Einträge: Abbildung von IP-Präfix auf Interface (ID & MAC)

Optimiert für Suche nach Ziel-IP-Adresse

Hohe Effizienz erforderlich (Lookup in Line-Speed)

oft durch HW realisiert

Routing Metric

Bewertungskriterium einzelner Übertragungsabschnitte
beeinflusst Bevorzugung / Vermeidung

Ganzzahl, unidirektional, Kosten / Gewicht

In Netzgraphen oft als Zahl an der Kante (Link), dann unidirektional gültig

Routing Policy

Betreibervorgaben zur Routing-Strategie

Üblicherweise im gesamten Netz gültig

Ermöglicht bevorzugen bestimmter Nachbarnetze bei der Wegewahl

Verteiltes adaptives Routing

Router tauschen Routinginformationen per Routingprotokoll aus.

Anpassung der Routen an die aktuelle Situation im Netz

Autonome Systeme

Protokolle innerhalb AS: Interior Gateway Protocol (IGP)

" zwischen AS: Exterior " " (EGP)

Identifizierung durch 32 bit Nummer

Erscheint nach außen als eine Einheit

Besitzt einheitliche Routing Policy

" " IGP, die sich von AS unterscheiden können

Vorteile: Getrennte administrative Domänen

Skalierbarkeit durch 2 logische Ebenen

Routingprotokoll innerhalb eines AS und nicht weltweit

Zwischen AS Austausch von Routinginformationen

Overhead steigt mit der Größe des Netzes

Routinginformationen & auszutauschende Routinginformationen

erhöhen sich

Alle Router brauchen global gleiche Sicht

Distanz-Vektor-Algorithmen: Konvergenz problematisch

Betreiberautonomie

Auswahl des IGP

Verbergen der inneren Struktur

Stub AS

Meist kleine Unternehmen

Anschluss zu genau einem Provider

Kein Transitverkehr

Multihomed AS

~~Transit~~

Große Unternehmen

Anschluss an mehrere Provider

kein Transitverkehr

Transit AS

Tier-1

Provider

oft globale Ausdehnung

Konnektivität herstellen: baue das Netz aus AS, das Internet

Transit: bezahlte Konnektivität

Upstream: Provider (Verkäufer von Transit)

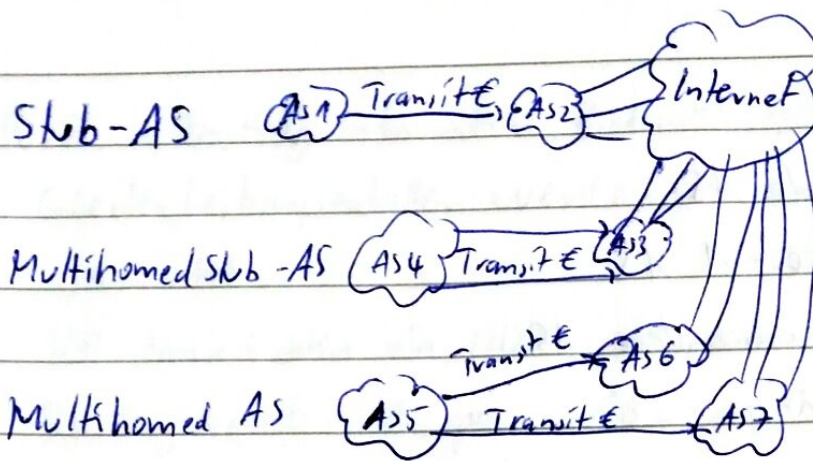
Downstream: Kunde (Käufer " ")

Datenaustausch in beide Richtungen

bezahlt wird i.d.R. nur Downstream

Transit-AS

gehört einem Provider, der Transit bereitstellt



Peering $\hat{=}$ direkte Verbindung zweier AS, aber kein Transit weiter als 1 hop.

Vorteil: Beide AS sparen Transitkosten

Problem: skaliert quadratisch

Public Peering

über IXP, zentrale öffentliche Stelle, an der Vernetzung stattfindet
neutrale Durchleitung auf Schicht 2

monatliche Kosten pro Port

(damit wird Wartung und Betrieb der Switching-Plattform gewährleistet)

Nutzung für Transit, Peering, VPLANs

Verschiedene Public Peering Policies:

open: AS ist offen für Peering mit jedem anderen AS

selective: nur mit gewissen, unter bestimmten Bedingungen

restrictive: AS voll

no peering: AS betreibt kein Peering

Tier-1

verkauften Transit

haben peering zu allen anderen Tier-1 AS

kein Einkauf von Transit

Tier-2

Große, überregionale AS

Downstream von Tier-1 AS

Verkaufen Transit an andere AS

Betreiber von Peering

Tier-3

Kleinere, Regionale AS

Downstream von Tier-2 AS

verkaufen keinen Transit an andere AS, nur an Anwender
betreiben Peering

Content delivery Provider (CDP)

Nähe zu Tier-1 Ports wünschenswert

Aufbau eines Content Delivery Networks (CDN)

Peering mit wichtigen Providern an zentralen Peering-Points

CDN

Weltweites Netz mit eigener ASN (AS Number)

Besieht aus Core Routern und Access Routern

Kunden verbinden sich mit " " " "

Lastbalancierung notwendig, " " sollen nahe am

Kunden sein

IGP

zur Weiterleitung innerhalb eines AS

meist: RIP oder OSPF, IS-IS oder EIGRP

RIP (Routing Information Protocol)

Weiterleitungstabellen werden von einem Anwendungsprozess verwaltet

RIP-Nachrichten via UDP, periodisch

Routing Metrik: ~~#~~ hops, 16 = ∞

alle 30 Sekunden

Timeout bei 180 Sekunden (setze Hop-Count auf 16)

Bei Änderung der Route, sende Triggered Update

Bei Empfang neuer Metrik < 15 : Neuer Eintrag in Routing Tabelle

" " einer Nachricht, dessen Metrik kleiner ist, als ein aktuell gespeicherter Wert = aktualisiere Eintrag

In allen anderen Fällen: ignoriere Nachricht

OSPF (Open shortest path first)

Link state Protokoll (Jeder Router muss den Zustand des gesamten Netzes kennen)

Dijkstra-Algorithmus zur Pfadbestimmung

Router versenden LSA (Link-State-Advertisement)

Inhalt: Besitzer, Links an Ports, Kosten = $\frac{\text{Reference Bandwidth}}{\text{Interface Bandwidth}}$

Jeder Router broadcastet seine LSAs und muss von allen anderen die LSAs empfangen können. Diese speichert er in einer LSD (Link State Database) und berechnet daraus seine Routen

Hello Protocol

contains ID of Router (sender) and expected receiver, if not known 0.0.0.0

If Router sees itself, communication is considered bi-directional

send periodically Hello Messages (10-30 sec)

LSA

header with information to uniquely identify LSA:
advertising router + Sequence number

body contains list of operational links of the router:

Associated cost

Type: router / end-system

Router ID or IP Address

Are associated w/ lifetime

incremented over time and invalid, if Max Age is reached
often @ 1h

⇒ Routers must broadcast periodically ~ 30 min

Initial Synchronization

(re-)start of a router ⇒ empty LSD, may not participate in routing

immediately performs handshake with neighbouring router
hello protocol

link state information is then exchanged
considered adjacent

Coping w/ dynamic changes

if nothing changes, no messages are exchanged
only periodically after 30 minutes

what if interface is down / router on the other site
gets unreachable / configuration changed?

minimum of time between LSA: 5 sec

Reliable Flooding

reception of LSA is acked by adjacent router

Multicast, hop limit = 1

LSA is considered "newer" and then stored, iff
(Seq is higher) OR ((Seq is Equal AND Checksum larger)
OR (smaller age))

If age is equal to maxAge and no LSA from sender is known
⇒ Send ACK & discard

If there is no LSA from sender or received LSA is newer

Store / replace LSA

Send ACK

Update age and flood LSA

If the stored copy is newer

Send stored copy to sender

If LSA and stored copy are the same

Treat it as ACK & discard

Scalability

AS can grow large

LSA flooding and computation overhead doesn't scale well

Idea: divide AS into areas, size ~100 routers

only flood and calculate routes for an area

Areas exchange summary-information w/ each other

⇒ hierarchy level in AS

Area 0 (backbone of AS)

Area Border Routers (ABR) interconnect areas

ABRs create summary-LSAs

it contains a list of reachable destinations within an area
each entry is associated with path cost from the ABR

Addresses are aggregated

after intra area path computation

Inter-Area forwarding goes through the backbone

⇒ end-to-end consists of source area → ABR

ABR → ABR

ABR → destination area

} always
shortest
paths

Broadcast Links

full connectivity between routers → n^2 links

approach: associate LAN with a pseudonode which will
receives broadcast

so called designated Router (DR)

" " Network LSAs

Types of LSAs

Router LSA: includes state of LSP

Network LSA: in multiaccess network, generated by the DR

Network-Summary LSA: generated by ABR, advertising
directions outside areas

Opaque LSA: for Traffic engineering ^(TE) reasons

RIP

vs

OSPF

Metric: hop count

periodic updates, 30sec

easy and memory saving, but problems
with distance vector (count to ∞)

Addresses showcomings at RIP

Large Networks can be divided into areas

Compatible w/ RIP

5 Label Switching

Issues with IP-based routing:

lookup complex but needs to be fast

TE difficult, always one shortest path used, not multiple shortest paths (if any)

no support for data flows, each packet handled independently

Flow $\hat{=}$ a set of packets that belong to the same interest

Microflow $\hat{=}$ a single connection (TCP), many

Macroflow $\hat{=}$ all TCP connections with destination of a given subset, fewer

Packet switching $\hat{=}$ packets are forwarded independent of each other, metadata required, expensive forwarding

circuit switching $\hat{=}$ Connection with resource reservations, no metadata within the stream required (ISDN)

virtual circuits $\hat{=}$ packets are forwarded by labels, connections do not have fixed resource reservation (MPLS)

Label switching is a combination of packet switching & circuit switching

packets are forwarded independently & need metadata (label)

paths established for flows through network

simpler forwarding decision

differentiation possible, TE, load balancing, QoS

Switching @ layer 2 instead of routing @ layer 3

Labels are only locally valid

Label $\hat{=}$ short unstructured ID of fixed length

unique within network, can be swapped

circuit identified through sequence of labels on the path
(\rightarrow virtual circuits)

Label can be placed on Layer 2.5, 1 byte long, 20 bit Label: TTL
or in the IPv6 header → Label Edge Router (LER)
Injected by an edge device on the border of a label switching domain (LSD)
And removed before leaving LSD
Label swapping may occur within LSD, done by Label Switching Router (LSR)

Multi Path Label Switching (MPLS)

based on label switching
operates in data plane
Standardized within the IETF
Increasingly adapted in large AS (faster than IPv6)
fast forwarding due to reduced packet processing
QoS support (for VoIP)
TE support to optimize network utilization
VPN support isolate traffic from other packets on the internet
Multiple Protocol support forwarding regardless of the underlying protocol (IP, ATM, ...)

Routing done by IP

Forwarding done by Label forwarding information base

Weiterleitungs-kategorie

Forwarding Equivalency Class (FEC)

Klasse von Paketen, die gleich behandelt werden sollen

↳ Synonym zu Flow

Gleicher Pfad (TE) ⇒ gleiche QoS

Grundlage für die Label zuzuordnung

Bsp. VoIP-Verkehr, gleiche IPs ~~und~~ mit gleichem Port

Granularität von Weiterleitungsklassen

Grob: wichtig für schnelle Weiterleitung und Skalierbarkeit

Fein: " " differenzierte Behandlung von Paketen

Vergleichbar mit einem Micro/Macroflow

TOS-Typ + Quelladresse + Zieladresse + Quellport + Zielport \Rightarrow Label

\rightarrow Möglichkeit, Datenströme zu differenzieren

Weiterleitungstabelle besteht aus:

(= Label forwarding ~~table~~ information base)

1. Eingehendem Label

2. Ausgehendem Label

3. Ausgehendem Interface

4. Adresse des nächsten hops

Was muss getan werden, damit MPLS läuft?

- definiere FEC

- verteile Labels (von außen, traffic engineer)

vergleichbar mit dem Verteilen von Routinginformationen

Label Bindings

stored in Label-Forwarding-Information-Base

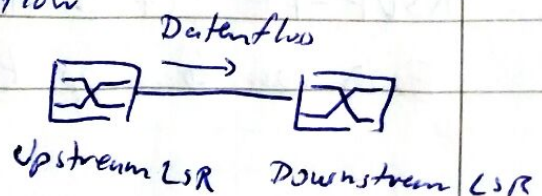
used as incoming label

distributed to neighbouring routers

Roles:

Downstream LSR, in direction of dataflow

Upstream LSR, against " " "



Label Distribution

Unsolicited downstream

LSR generates label as soon as it is ready to forward MPLS packets of the respective FEC

- Upstream Neighbours update forwarding tables
- Label used as outgoing label

Downstream on demand

Downstream LSR generates label binding on demand
Upstream LSR has to request binding for FEC

Usually a MPLS-defined protocol was invented and used de-facto: RSVP-TE used mostly

RSVP (Resource-Reservation-Protocol)

goal: Resource reservation for end-to-end datastreams
establish connection and periodically indicate liveness
they get removed, if no liveness is detected

Path Message

From sender to receiver

Find path to receiver

Every hop is recorded in the message

Resv Message

From receiver to sender

Bandwidth reservation on return path

RSVP-TE (+ traffic engineering)

Extension to support Label distribution

Many additional fields for functionality (fast reroute, ...)

Path Message

from upstream LER to downstream LER

Label request (downstream on demand)

Bandwidth reservation

Resv Message

response to Path Message

Label binding (hop-to-hop)

VPNs implemented by label stacking

Outer label identifies path to LER

inner label " VPN instance / customer

Layer 3 VPN "BGP / MPLS IP VPN"

Customer may use private IP addresses

Different customers can have the same IP address

Provider takes part in routing

use of 64 bit "Route Distinguisher" to identify customer in BGP / MPLS

LER has Virtual Routing and Forwarding Table (VRF) for each customer

Label stacking

Inner MPLS label mapped to VRF

Outer label defines provider, gives path through provider network

→ Different customers can share path as well as different paths for the same customer possible

Layer 3 VPN

Customer Edge Routers announce their subnets to their LER

LERs then signal this information to the other LERs

LERs " create FECs for announced subnet

Inner label also announced over BGP

outer " based on traffic engineering tunnels

Layer 2 VPN

Customer connects Ethernet switch to LER

Inner label then defines output @ LER

No control traffic between customer and provider

In contrast to layer-3 VPNs, where customer announces routes to provider

Label Distribution

Path Message	Backbone
Explicit Route	LER-KA1 , Backbone-Z, Backbone-P, LER-P
Record-Route	LER-KA
Reservation	80Mbit/s
Label	Request new Label

will more down for every Router passed!

Resv Message

Record-Route	LER-KA1, Backbone-KA, Backbone-Z, Backbone-P, LER-P
Label	2, 8, 7, 2 changed @ every Router that was passed by

6. SDN

Basic principles:

1. Separation between data & control plane
2. Control plane has global view
3. different modes of operation applicable
4. processing is based on flow
5. Network is soft-wave programmable

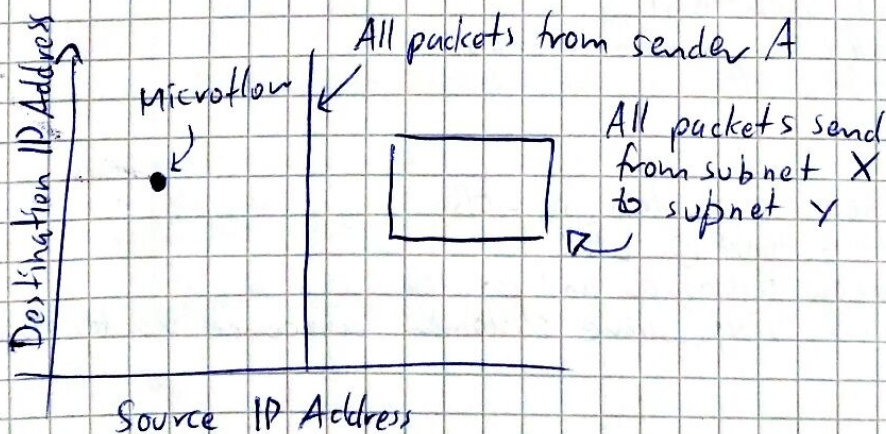
SDN Controller

generates, manages & sends rules for forwarding to SDN routers that they then store and obey.

Rule = \langle Match, Action, Switch \rangle

ID for subset of a flow OP to perform Device on which it should be installed

Flow Space describes flows as objects in a multidimensional space



Global network view

A SDN Controller is logically centralized
assumes that it knows all switches and their configuration
as well as their topology
allows execution of link-state routing without protocol

Static configuration

Fixed settings prescribed by network operator

i.e. static forwarding policies

TCP connections between controller and switches, usually via static IPs

Active Monitoring

SDN controller is actively involved in refining the global view
i.e. periodic polling of information from the switches (for
keepalive & port utilization monitoring)

Topology detection via Link Layer Discovery Protocol

Passive Monitoring

SDN Controller refines the global view by passive observation

i.e. Information learned from handshake w/ switch

" from internal state of switch (e.g. rules programmed etc)

Programmable via Software

via so-called network-apps

SDN controller can execute multiple apps in parallel

programmed via API (application programming interface)

SDN controller then generates rules $\langle M, A, S \rangle$

1. can be proactive: @switch

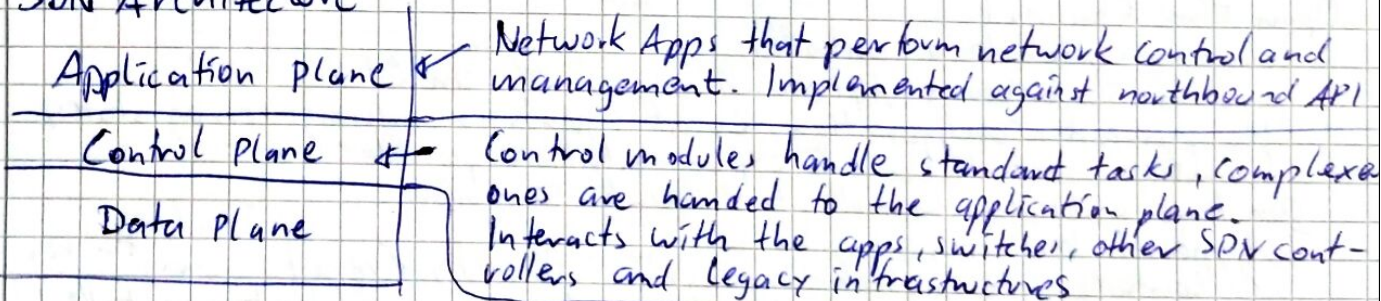
operation is programmed before first affected package arrives

2. reactive: ask controller, if package arrives, for which no rule is known
the former is only feasible w/ flows that are large
don't have any additional delay

controller doesn't necessarily be online

the latter has a higher setup time and controller needs to be reachable
but allows fine-grained flows

SDN Architecture



Infrastructure for the packet forwarding/processing
supports operations like: match, drop, forward

Open Flow (@Data Plane)

Specification for SDN-capable Switches

defined by Open Networking Foundation (ONF)

defines a common logical architecture for SDN-capable switches
rather the functionalities.

implementation stays vendor specific

also defines a standard for secure protocol for controller/switch interaction

OFS $\hat{=}$ Open Flow Switch

Ports: Ingress / Egress = Input / Output

Physically a port / hardware interface

logical abstractions possible for aggregation / tunneling

Port short-names: ALL $\hat{=}$ all egress ports

IN-PORT $\hat{=}$ references in-port of packet

CONTROLLER $\hat{=}$ the interface to its controller

NORMAL $\hat{=}$ back to vendor specific ports

Flow tables

Separate packet streams into flows

flow table entries identify specific flows

Match fields select packets by their header fields

priority determines the relative precedence of an entry

Match fields | Priority | Counters | Instructions | Timers | Cookie | Flag

this enables management functions like:

counting \times processed packets

automatic removal of dead flow (life-time = 0)

set marker values (cookies)

flags indicate how a flow is managed

Matches match packet data to corresponding match field values

Wildcards using bitmasks

empty match fields match all packets / flows

simultaneous matches can occur, highest priority then is taken.

Optional table - miss - flows

Matches all flows \rightarrow enables reactive flow programming

Actions: forward, drop, modify header of a packet, push new tag onto a packet, remove a tag from packet

Open Flow Channel

Connects switch to a controller

provides southbound API functionality of an OFS

Management, configuration, event signals, monitoring (error states, liveliness, statistics...) & experimentation possible

can have multiple channels to different controllers.

3 types of messages:

1. Controller-to-Switch message

(modify / read)

\rightarrow for new flow, port status change

2. Asynchronous

"

(usually switch-to-controller)

3. Synchronous

"

(hello / echo / error)

Pipeline Processing

chaining of multiple flow tables

flow tables are immutable, no recursion possible

Action set

Ingress & Egress processing each have an associated action set

Accumulates actions during pipeline processing

action \rightarrow are deferred until the end of processing unless told otherwise

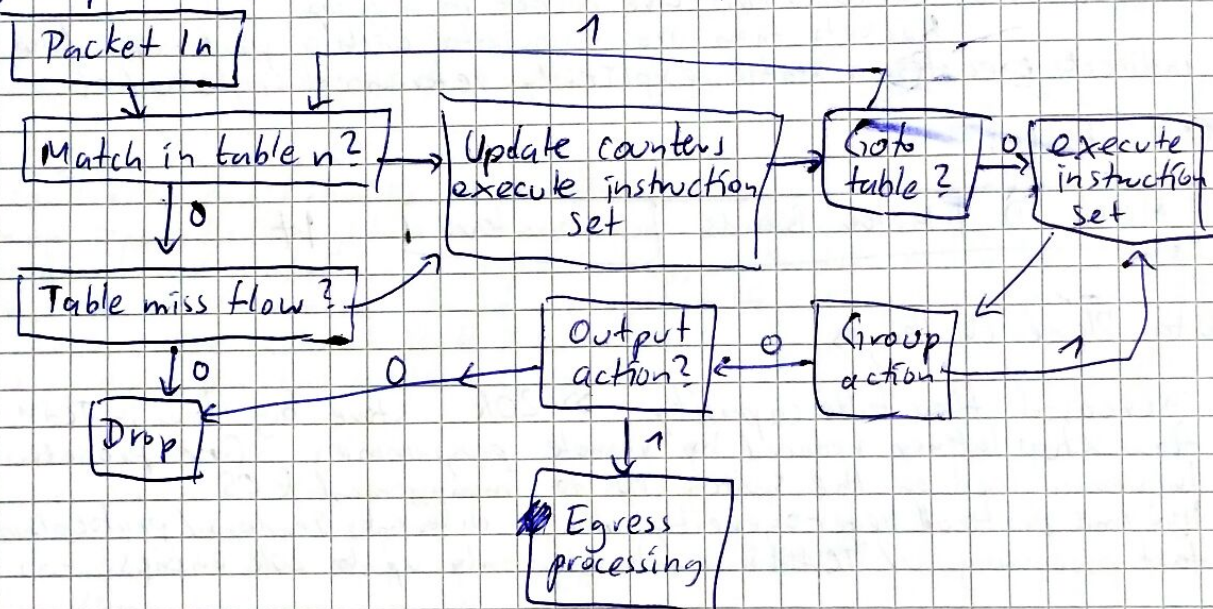
Instructions

control how packets are processed in the pipeline
 instructions are more generic operations than actions
 Flow table traversal can be controlled via go-to statements
 write/clear operations can be applied to the action set
 or executed directly via apply instruction
 Action sets can only have one " of each type. Therefore
 older ones become overwritten
 all instructions are then executed in a static order

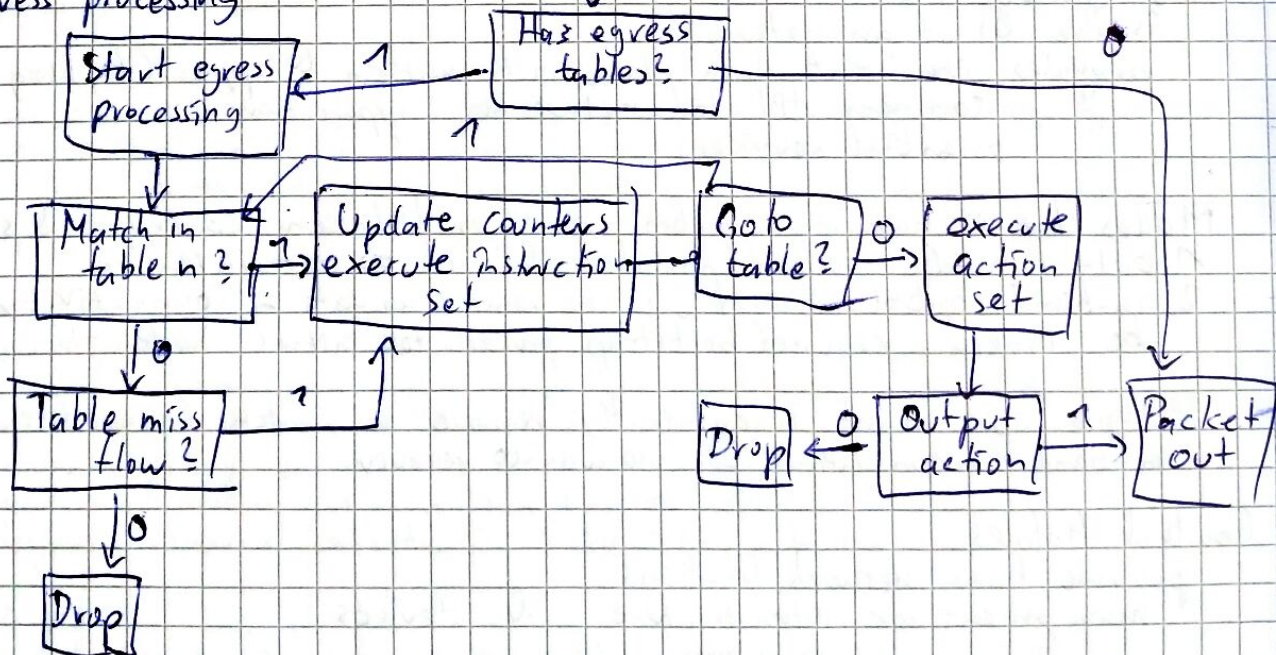
Learning switch

separate tables for learning & forwarding
 Learning table matches source addresses and forward to controller, if
 address was not yet learned
 forwarding table matches destination " and forwards packet if
 yet learned. If not, flood packet
 Only 2-N rules for N hosts
 but HW rarely supports two tables

Ingress processing: table 0, action set empty



Egress processing



Group tables

Group Identifier	Group type	Group Counters	Action Buckets
------------------	------------	----------------	----------------

represent additional forwarding methods

e.g. fast failover, link selection, ...

can be created by ingress processing

optionally groups may perform group actions (group chaining)
effect depends on group ~~priority~~ and their action buckets

Action Buckets

each group has ≥ 0 group buckets

not every action bucket needs to be executed

a group w/o action buckets drops a packet

action buckets contain a set of actions to execute, just like an "set"

Group types have specific semantics

all: execute all buckets in a group, selected by method (hash, round-robin, ...)

select: selects only one of many buckets in a group

fast failover: executes first live bucket in a group

buckets that are associated with a port, up = live

indirect: executes a static, explicitly referenced single bucket

Meter tables

Meter ID	Meter Bands	Counters	it
----------	-------------	----------	----

Data Plane challenges

Increased flow table capacity $\gg 20k$ store big flows in TCAM, others in RAM

flow setup latency reduced by remote programming (use proactivity)

Increased load on the switch CPU for management & OS

HW and protocol heterogeneity due to different vendors / protocols

fast processing w/ TCAMS can handle only up to 20k entries

SDN Control Plane

is the OS of a network

provides consistent data plane information to apps (topology, flows, stats...)

" common API for networking applications

" essential services

Modes for exchanging control messages between controller & switch:

1. out-of-band: dedicated physical channel \Rightarrow cost intensive

2. in-band: control messages use the same " as data \Rightarrow connectivity can be broken, requires bootstrap process for initial switch-to-controller

Multiple apps might independently remove a switch from the system
in-band communication, lat, no way to recover

Control Modules

provide basic network functions

management for infrastructure & devices

topology

statistics

Shortest path forwarding
isolation and security enforced between apps & services
delivery of notifications between SDN switches and apps
App use these control plane functions to realize/implement complex services

Northbound control plane functions to applications
abstraction from underlying hardware

Service Abstraction Layer

Provides basic services to controller functions and apps
Southbound interfaces are treated as plugins
Maps service requests to the most appropriate southbound protocol
⇒ provides uniform access to the underlying network

Plugin Manager

exposes functionality in form of abstract features
i.e. change source or destination fields of IPv4 header
decrease TTL in IPv4 header
configure a queue

Service Manager

maps service requests to feature requests

Distributed Controllers

1 controller can't handle ∞ switches
domain separation might also need more controllers
⇒ problems w/ scalability, responsibility, reliability, incremental deployment
but network view needs to be consistent for all controllers
⇒ synchronize network state information
this is done via the westbound interface

Westbound Interface

enables communication between SDN controllers
for data import & export
coordinated flow setup, exchange of reachability information
algorithms for data consistency
monitoring & notification
can support heterogeneity between controllers

Synchronization

Controller directly apply internal operation & notifies remote controllers of relevant changes in the network
Apps can perform data plane operations on remote switches
operations are then delegated to responsible SDN controller

Abstractions

apps only see an abstracted network
Big switch abstraction, for some apps that doesn't need to know of long router chain, it can be omitted
Network slicing as a form of virtualization
individual users get their own slice of network
requires resource orchestration inside the controller

Challenges

Scalability: logically centralization requires strong HW
control plane can be overloaded by data plane size

Important parameters w/ scalability implications

- number of remotely controlled switches
- number of hosts / flows in the network
- " " messages processed by the controller
- communication delay between switches & "

possible solutions:

Distributed Controllers

parallelism-based optimization such as multi-threading, batch processing, optimized routing schemes that consider flow table size / amount of rules

Challenge	Possible Solution	Tradeoff
Comm. delay, flow setup latency	Proactive mode of Separation, w/ SD cards	Microflows may become invisible
Limited control flow table capacity	Caching mechanisms	additional HW Cache miss increases WCET of flow
High number of Controller events	distributed controller architectures	problems w/ consistency & synchronization

Consistency

multiple controllers inherit the problems of a distributed system

CAP theorem: Consistency: system responds equally, no matter which node responded

Availability: system is always responding

Partition tolerance: system continues to function even if messages are lost or parts of network fail

It is impossible to provide all 3 @ the same time. Maximum of 2

SDN Applications

Network is software programmable, Apps provide functionality

functionality depends on the environment like: data centers, enterprise network, home networks, IXP

Broad classification into 5 wide areas:

TE

Mobility & Wireless

Measurement & Monitoring

Data Center Networking

Security & Dependability

Load balancing in SDN

reactive: load balancing algo triggered for every new connection (microflow)

needed apps: Flow Manager, net manager (monitoring usage)
Host " (monitoring server liveliness)

proactive: preinstall forwarding rules using address space partitioning
using wildcard to reduce rules

each server gets a weight α_i

use binary trees to model address space

normalize weights, so that $\sum \alpha_i = 2^n$, $n \in \mathbb{N}$

Map nodes to leaves according to α_i

optimizations like aggregating sibling nodes

Traffic dynamics

traffic changes over time. No problem for reactive solution.
proactive solution requires transition strategies
temporarily intercept all connections and redirect to controller

Policy Cop for TE

optimize resource utilization

apply (customer specific) policies

QoS with regard to SLAs

Validator monitors network to detect policy violations

Enforcer adapts rules based on network conditions and high level policies

Adaption actions:

packet loss → reconfigure queues or reroute to better path
throughput → " rate limiters to throttle misbehaving flows
latency → reroute to path w/ less congestion & suitable delay
Jitter → " " " " " "
Device failure → " " " that bypasses the failure

7. Network Function Virtualization

Middle box $\hat{=}$ Device on path between client & server that performs other tasks than forwarding / routing
NAT, firewall, proxy, load balancing, intrusion detection

Network Address Translation

allow inside a realm globally non-unique IP-Addresses

can't be used for routing

provides mechanism to connect a realm w/ private addresses to an external realm w/ globally unique IP-Addresses

Exchange globally unique & private addresses when packets traverse network boundaries

Recalculate checksum

Collisions on port numbers need to be resolved

as well as state of every data stream needs to be traced

Firewall

establishes barrier between trusted & untrusted network

drops or forwards packets based on rule set

shallow: only check header of layer 3, 4

deep inspection: look @ data input

stateless: every packet is inspected independently

stateful: remember state of connections (like seq numbers)

Proxy

Intermediate device that acts like server & client

performs requests based on requests of clients

appears to the client as a server and to the server as a client

Transparent: doesn't change data / requests

Non-": modify requests / responses in order to provide additional services

Web caching: Proxy stores responses from server, accessing the same content again is faster, if proxy is closer than the server to the client

Caching used by content delivery providers to ensure audience has low latencies when accessing content. Placement of cached content depends on size of audience in an area.

Middle boxes, often built as proprietary hardware fast, but inflexible & black box for infrastructure operator
Static wiring: hard to setup/teardown/improve/move

Network Function Virtualization (NFV)

Mimic ideas of cloud-computing

Network functions are implemented in software

and therefore decoupled from the hardware they run on

⇒ can be executed on shared HW

⇒ can be easily moved and instantiated everywhere

Servers are often close to network nodes or to end user premises

Network Service

combine multiple network functions

end-to-end behavior of a network service is the combination of the individual network functions passed by on a path

Benefits:

resource sharing

Agility & flexibility

rapid deployment

} reduced costs, no special personnel, x86

NFV Infrastructure (NFVI)

can represent more Points of Presence (PoP)

many enterprises can be in the same NFVI

redirected transparently by SDN

can be done w/ MPLS

each content delivery provider can then easily choose where to open new PoP

Virtualization

abstraction layer between HW and OS on VM, often called hypervisor

acts as a resource broker between HW and VMs

translates between I/O from VM to shared HW

through isolation, multiple coexisting VM instances are possible

allows ~~live~~ live migration of VM to other host

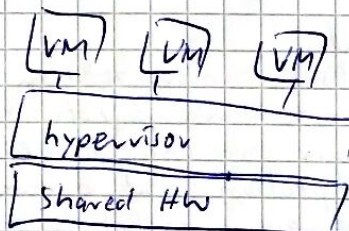
Type-1

runs directly on HW

high-performance

strong isolation

synchronizes access of VM to the HW

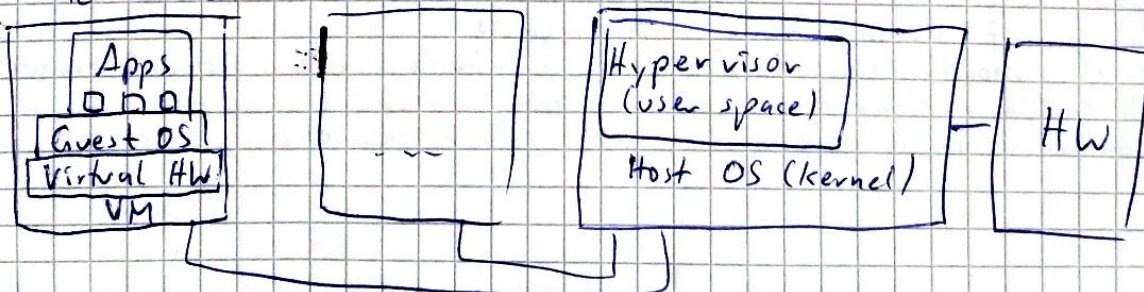


Type-2

runs on an hypervisor OS

executed as an app in user space

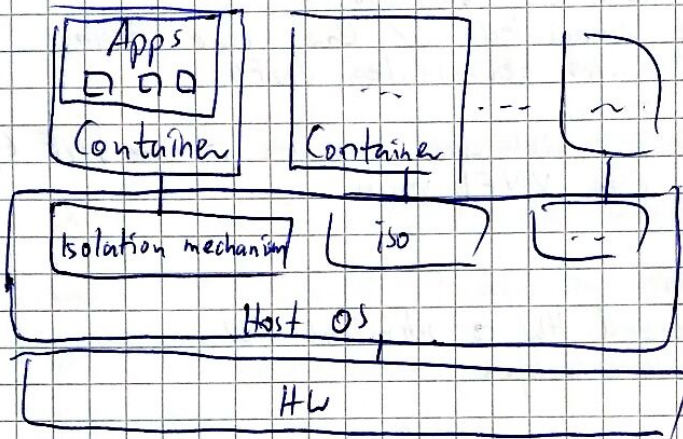
interaction of VM and HW is directed to physical devices through a VM-driver or the host OS



Container based virtualization

single kernel provides multiple instances (containers) of the same host operating system

- no hypervisor involved
- isolation enforced by host OS
 - each container has its own view of OS
- Apps in containers are executed by the host OS
 - apps depend on it
- Kernel synchronizes access of containers to HW



High level view on NFV-Framework

3 main building blocks:

VNF: the actual functions provided by the software

NFVI: HW & SW used to execute VNFs

NFVO: Orchestration of network service management & NFVI resources (MANO)

Functionality (largely) the same as non-VNF networks, but implementation is done in software
Scalable by adding virtualizable HW resources
implementation HW independent

NFVI $\hat{=}$ sum of HW and SW components that build up the environment on which VNF apps run
provides physical platform
manages computing, storage & network resources
usually based on low-cost standard HW

MANO $\hat{=}$ Management & Orchestration
responsible for MANO of NFV resources
includes computing, networking, storage, VMS
all virtualization specific tasks

consists of 3 parts:

1. NFV Orchestrator
2. NFV Managers
3. Virtualized Infrastructure Manager

NFV Orchestrator

performs lifecycle man. of network services
man. and coordinates NFVI resources across multiple virtualized infrastructure man.
may instantiate VNF man.
handles topology between them
policy man.

NFV Manager

responsible for lifecycle man. of ~~all~~ VNFs
multiple VNFs can be associated w/ the same man.
instantiates, configures, scales, terminates VNFs
collects statistics, error logs...
can also provide functions which are specific to a VNF type
requests resources from the VNFI man.

NFVI Manager

man. resources of NFVI
man. mapping of virtualized HW to physical HW
runs hypervisors / software images

Service Function Chain (SFC)

ordered set of network functions
allows creation of composed network services
like advanced caching, authentication → firewall → cache
done by label stack
represented by a forwarding graph (may be cyclic)

Challenges

Security
Performance 100+ GBit/s
placement
reliability
testing & debugging
carrier grade requirements
existence of legacy networks

8. Ethernet Evolution

Aloha

senden, drahtlose Übertragung
wahlfreies senden, Kollisionen möglich
Asynchrone Zugriff

Kollisionserkennung

Explizit: Besetzung durch höhere Schichten
separater Kommunikationskanal erforderlich
Implizit: spiegelte empfangenes Paket, Sender vergleicht dies
Reaktion: Sendewiederholung nach zufälligem Zeitintervall
Abbruch nach zu vielen Sendewiederholungen

Slotted Aloha
 synchroner Zugriff in Zeitschlitzen
 im Mittel weniger Kollisionen
 Eingesetzt bei GSM (Steuerkanal)

Metrik für Kommunikationssysteme $U_{max} \hat{=}$ maximale Auslastung
 wie gut kann das Medium ausgenutzt werden?

hierzu braucht man viele Variablen:

r	Datenrate der Übertragungsstrecke [bit/s]
r_e	" erzielt [bit/s]
N	teilnehmender Systeme
X	zu übertragende Datenmenge [bit]
t_w	Wartezeit, bis das Medium zugeteilt wird [s]
t_s	Sendezeit [s]
U	Auslastung
d	maximale Distanz zwischen zwei Systemen [m]
t_d	Ausbreitungsgeschwindigkeit [m/s]
α	t_a / t_d
p	Wahrscheinlichkeit für Datenübertragungswiederholung

Slotted Aloha W'keit für geglückte Übertragung: $Np(1-p)^{N-1}$

$$\Rightarrow U_{max} = \frac{1}{e} = 0,369 = 36\%$$

Aloha W'keit für geglückte Übertragung: $Np(1-p)^{2(N-1)}$

$$\Rightarrow U_{max} = \frac{1}{2e} = 0,18 = 18\%$$

CSMA basierte Verfahren (Carrier select multiple access)

Medienzuteilung: Zeitmultiplex, variabel, zu fälliger Zugriff
 vor Senden prüfen, ob Medium belegt (Listen before talk)
 Asynchroner Zugriff

CD (Collision Detection) Ethernet
 listen before and while talk

CA (Collision Avoidance) WLAN

Ethernet

Zeitmultiplex, variabel, zufälliger asynchroner Zugriff

Kollisionserkennung durch Mithören an der sendenden Station

Exponentieller Backoff

1-persistent

96 bit keine Aktivität \rightarrow frei

Senden ist 1-persistent

Jamming Signal: 48 bit Präambel

Exponentieller Backoff für Sendewiederholungen

Kollisionserkennung muss vor Beendigung des Sendevorgangs erfolgen,
sonst kein reines CSMA

2 * Ausbreitungszeit \leq Sendezeit

Erreichbar durch das Padding-Field PAD

Ethernet-Rahmen

PR (7)	SD (1)	DA (6)	SA (6)	Type (2)	Data (≥ 1500)	PAD	FCS (4)
--------	--------	--------	--------	----------	----------------------	-----	---------

Präambel für Synchronisierung (101010...)

Start of Frame Delimiter zeigt Beginn an (10101011)

DA Zieladresse

SA Source adresse

Type/Length des darüberliegenden Protokolls oder Länge der Nutzdaten

Datenfeld

PAD Paddingfeld (optional)

FCS Frame check sequence

Zwischen Frames liegt ein Inter Frame Space (IFS) von 96 bit

Beurteilung der Auslastung

$$r_e = \frac{x}{t_s} + t_a$$

$$\Rightarrow U_{\max} = \frac{1}{1+a}$$

$$\lim_{N \rightarrow \infty} U_{\max} = \frac{1}{1+3,349}$$

~~Wichtige~~ ~~Wichtige~~ ~~Wichtige~~

Fast-Ethernet

10 Mbit/s \rightarrow 100 Mbit/s

automatische Erkennung & Umstellung durch Protokolle

Netztopologie: Stern

Halbduplex oder duplex Links

CSMA/CD

Codierung: 4B/5B, NRZI

Stav an Switches

Workarounds mit Backpressure

- Kollision herbeiführen
- ⇒ Sender bricht ab und wiederholt Rahmen (Backoff)
- ⇒ zu viele Wiederholungen → Verbindung beendet
- Medium als belegt vortäuschen (Präambel senden)
- Sender wartet auf Slot
- ⇒ Implizite Flusskontrolle

Duplex Links ohne CSMA/CD brauchen anderes Workaround

⇒ explizite Flusskontrolle:

PAUSE-Funktion

Sende PAUSE-Rahmen mit Zeitangabe

PAUSE mit Zeit = 0 ≙ explizite "GO"

Bremst keine darüber liegenden Protokolle, löst aber auch nicht langfristige Überlastungen im Netz

Control teil von
PAUSE-Funktion im MAC-Layer 2

MAC-Kontrol Frame

PR(7)	SD(1)	DAC(6)	SA(6)	Type(2)	MCC(2)	Parameters(44)	FCSC(4)
-------	-------	--------	-------	---------	--------	----------------	---------

Type = 0x8808

MCC = MAC Control Opcode

0x0001 ≙ PAUSE

Kann nur bei Duplex-Betrieb genutzt werden
wird an vordefinierte Multicast-Adresse geschickt

Parameter: gewünschte Auszeit
fülle mit 0en aus

Automatische Konfiguration

Ziel: Geräte erkennen selber, mit welcher Datenrate & Kodierung eine Kommunikation möglich ist.

Zwischen Endsystemen und Knoten (Hub, Switch)

während der Initialisierung

Austausch von 16 bit langen Nachrichten signalisieren mögliche Konfigurationen, als Folge von Takt- und Datenpulsen gesendet

Gigabit/s-Ethernet

Geänderte Kodierungsverfahren:

8B/10B mit NRZ

Datenrate 1000 Mbit/s (abwärtskompatibel zu 10 Mbit/s und 100 Mbit/s)

Problematisch: Kollisionserkennung bei halb duplex-Links

Carrier Extension

soll Kollisionserkennung sicherstellen

Künstliche Erhöhung der Übertragungszeit ohne die minimale Länge der Rahmen zu erhöhen

Grundlegende Änderung:

Länge des Zeitschlitzes \neq Länge der minimalen Übertragungsdauer
512 B 512 b

\Rightarrow Carrier Extension zwischen 16 und 448 b lang
" " wird physikalisch als non-data-bits gesendet

Frame Bursting

Ziel: effiziente Übertragung kurzer Rahmen

Stationen dürfen einen Burst von Rahmen direkt hintereinander senden

Erster Rahmen mit Carrier Extension, falls notwendig für die CD

Letzter " muss nach mindestens 8192 B beginnen

Schema



P $\hat{=}$ Präambel & Start of Frame delimiter

D $\hat{=}$ Daten

E $\hat{=}$ Carrier Extension

| $\hat{=}$ Inter Frame Space

10 Gbit/s - Ethernet

Glasfaser / twisted pair

Datenrate: 10 Gbit/s

Nur Punkt-zu-Punkt Verbindungen

Ausschließlich Vollduplex-Betrieb

kein halbduplex-Betrieb mehr, keine Hubs

kein CSMA/CD

kein Framebursting

keine Carrier Extension

Weiterhin: 96b IFS

1518 B maximale Größe eines Rahmens

64 B minimale " " "

Zusätzlich: Jumbo-Frames: 9014 B Nutzdaten

Improved coding mechanisms:

Carry multiple bits per transmission step: 16-level pulse amplitude modulation

Interference becomes more problematic:

Near end \rightarrow crosstalk: NEXT

Far " " : FEXT

Crosstalk from neighbouring cables: Alien Crosstalk: AXTLK

Signal-to-noise ratio SNR critical

Not a typical solution for
rather connecting buildings

a network @ workplace / home

10G Base-SR (400m)

10G Base-ER (40km)

10G Base-T (100m)

} ~~64B/66B~~ Code

40 / 100 Gbit/s - Ethernet

Distances up to 40km

multiple parallel channels

twisted pairs / different frequency / wavelengths

Virtual lanes contain data streams

aggregate different data streams into one virtual channel and
enclose them with identifier (64B / 66B → 46B / 66B)

apply bit-level multiplexing, if $\#$ channels $<$ $\#$ virtual channels

- ⇒ Speeds have increased by $\times 10k$
variety of different medias used
extends to WAN
switched Ethernet is dominant
MAC-Protocol often unnecessary
One constant remained:
Ethernet frame format

Spanning Tree

Brücken: Ziel: Kopplung von LANs auf Ebene 2

Filterfunktion: Trennt Intranetverkehr von Internet-Verkehr

Source-Routing Brücken:

Endsystem fügt Wegewahl-Info in Paket ein

Brücken leiten es anhand dessen weiter

Senden über Brücke nicht transparent

Transparente Brücken:

Lokale Weiterleitungsentscheidung von der Brücke

Wegewahlinfo in Tabelle enthalten

Statisch / gelernt

Endsystem muss nichts von der Brücke ~~wissen~~

	Transparente Brücken	Source-Routing-Brücken
Transparenz	✓	✗
Verzögerung	gering	Ggf. Hoch
Routing	nicht optimal	optimal
Alternative Wege	✗	✓
Kapazitätsausnutzung	Schlecht	Gut
Skalierbarkeit	Schlecht	Gut
Reihenfolgeerhaltung	✓	✗

Transparente Brücken $\hat{=}$ Switch

hat für jeden Netzanschluss eine eigene Schicht-1 und MAC-Instanz

Datenpfad über MAC-Relay

implementiert per Weiterleitungsfunktion auf Schicht-2

Kontrollpfad geht durch LLC (Logic Link Control) höher an Schicht-3

Aufgaben: etablieren einer schleifenfreien Topologie

Weiterleiten von Paketen durch Lernen der Lokation von Endsystemen

Aufbau der Filterdatenbasis

Spanning Tree Algorithmus (Brückenprotokoll)

Organisierung der Brücken und Links zu einer Baumtopologie

Weiterleitung nur am Baum entlang möglich

Voraussetzungen:

- MAC-Adresse aller Teilnehmer eindeutig
- Kenntnis " " " "
- Anschlusskennung " " "
- Pfadkosten aller Anschlüsse bekannt

Brücken versenden spezielle Pakete

Bridge Protocol Data Units (BPDU)

beinhalten Kennung der sendenden Brücke

" " angenommen Root-Brücke

Pfadkosten ~~von~~ von sendender - zu Rootbrücke

Algo:

1. Bestimmen der Rootbrücke
2. " " Root-Anschlüsse jeder Brücke (Σ Weg zur Rootbrücke)
3. " " Designated-Brücke für jedes LAN
LAN kann mehrere Brücken konfiguriert haben
Geringste Kosten am Root-Anschluss $\hat{=}$ designated bridge
bei gleichen Kosten regelt die Brückenkennung
4. Schalte alle anderen Brücken aus

Initiell haben alle Brücken keine Topologie-Infos

✓ Brücke: Annahme: ich bin die Rootbrücke
sende periodisch BPDUs

Bei Empfang einer BPDU mit kleinerer Kennung \Rightarrow nimm nicht weiter an, selber Rootbrücke zu sein und sende keine BPDUs mehr

Bei ~~empfang~~ Empfang weiterer BPDUs, ggfs. Update der Konfigurationsdaten.

Kleinere Kennung ^{der Rootbrücke} \rightarrow Update

gleiche " aber kleinere Kosten zur Root \rightarrow Update

gleiche Kennung der Rootbrücke & Kosten gleich aber sendende Brücke hat kleinere Kennung \rightarrow Update

Brücke stellt fest, dass sie nicht designated-bridge ist \rightarrow keine Weiterleitung mehr

Rootbrücke sendet periodisch BPDUs weiter
nur aktive Brücken leiten diese weiter

Ausbleiben von BPDUs \Rightarrow Brücke wieder Root und wiederhole

Nach Stabilisierung werden Pakete über die entsprechenden Ports weitergeleitet
entsprechend der Einträge in der Filterdatenbasis

Filterdatenbasis:

Enthält für die Weiterleitung erforderliche Daten:

Zieladresse, Ausgangsport, Zeitgeber

Statische und dynamische Einträge

↳ vom Sysadmin gestellt

gelernt von durchlaufenden Paketen

verlernt durch Zeitgeber

Pakete mit lokalem Ziel werden nicht über die Brücke weitergeleitet

Empfang eines Pakets

Ziel bekannt?

Identifiziert durch Ziel-MAC-Adresse

Weiterleiten an Port gemäß Tabelle

Unbekannt?

An alle anderen Ausgänge fluten

Quellsystem lernen, ist über Eingangsport erreichbar

Neuen Tabelleneintrag erstellen

Rapid Spanning Tree

Neue Zustände der Ports

Alternate Port: bester alternativer Port zur Root-Brücke

Backup-Port: alternativer Port zu einem Netz

BPDUs werden nun als keep-alive-signal verwendet
=> Erkennung eines Ausfalls schneller möglich

Neue Typen von Ports:

Edge-Port: nur Endsysteme (=> Schleifenfrei)

Point-to-Point: Nachbar ist eine Brücke, Voll duplex

Shared Port: Verbinden zu einem Netz

Echtzeit Ethernet

Für bspw. Windräder, Fahrzeuge

Unterstützt Echtzeitfähigkeit

Kompatibel zu Standard Ethernet

"Time-Triggered-Ethernet" TTE

3 Versandsklassen

Zeitgesteuerte, hochprior. Rahmen (TT), feste, globale Zeitschlitz
Ratenbasierte Datenströme (rate-constrained), Maximale Größe, Mindest-
abstand
Best-Effort-Verkehr (u.a. Standard Ethernet)

TT-Ethernet-Rahmen

in der Empfängeradresse wird die Nachrichtenart kodiert
Critical Traffic (CT) Marker & CT-ID im Header

Entweder TT oder rate constrained

Erkennung von CT über Bitmaskenvergleich

Verwenden für Zeitplan konformes Scheduling und Weiterleitung von
time-triggered Paketen

Ether Cat

Ethernet for Control Automation

Steueraufgaben auf Feld-Bus-Ebene

Effizienter Umgang mit kleinen Datenmengen (8-bit Maswert)

Niedrige Latenz & Jitter für Roboter automation

Audio Video Bridging (AVB)

Für Rundfunk, Veranstaltungstechnik (Konzert, ...)

Zeitsynchronisation zwischen AVB-Stationen < 1ms

Sehr exakte ~~an~~ Echtzeituhren

Jitter < 100ns

Einsatz von Hardware-Zeitstempeln beim Senden & Empfangen

9 Data Center Networking

Typical Services:

IaaS: Infrastructure as a service (for virtual computing, ...)

PaaS: Platform " " " (" web development, ...)

SaaS: Software " " " (" services like Dropbox...)

StaaS: Storage " " " (Storage of a cloud provider)

Data Center 100k+

typically has large # of compute servers and VMs
 switches w/ small buffers
 Should be extensible w/o massive reorganization
 Needs to be reliable (redundancy...)
 " " " highly performant (Load of 100Tbit/s + @ low latency)

Top-of-rack-Switches (ToR)
 connects servers within rack
 48 x GBit/s Ethernet w/ 10Gbit uplink
 small buffers
 42-48 rack units / rack

Challenges:
 Maintain scalability 100k+ - 1m servers
 maximizing throughput while minimizing latency & cost
 guarantee integrity & service availability
 enhance power efficiency

Data Center Network
 Interconnects data servers with each other
 as well as them to the internet (border routers)
 => two types of traffic: external clients <-> servers
 servers <-> servers

Commodity protocols:
 TCP/IP
 Ethernet

Routing / Forwarding within data center
 needed efficient way to communicate between two servers
 => " utilization
 avoid forwarding loops
 detect failures quickly
 provide flexible and efficient migration of VMs between servers
 no change of IP Address to keep TCP connections opened

Layer 3 routing? OSPF?
 routing tables are too slow
 VMs moved require new IP address
 manual configuration

Layer 2 MAC	✓	X	X	✓
Layer 3 IP	X	✓	✓	X
	Plug and play?	Scalability?	Small Switch State?	Seamless VM Migration?

Why not Ethernet?

Internet Tree Topology:
 East-West-Traffic: Between internal servers and server racks
 North-South- " : Result of external request

Fat Tree Networks

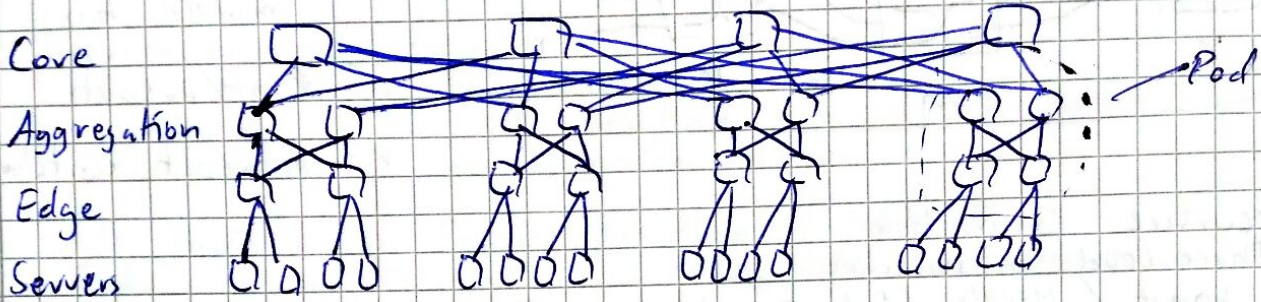
Problem: Connect large number of devices using switches that only have a limited # ports

Characteristics: For any switch: # Links going down to its children is equal to # links going up to its parents.

⇒ the link gets fatter to the top of the tree

But switches w/ high number of ports are expensive

4 Pod Fat Tree



K-Pod fat tree

Each switch: k ports
Up to $k^3/4$ servers using $\frac{5k^2}{4}$ switches and $\frac{k^3}{2}$ links

Edge and aggregation arranged in k -pods
 $k/2$ edge switches and $k/2$ aggregation switches per pod

Edge " " can connect to $k/2$ servers

Core " " : $k/4$, each has non connected to each pod

⇒ Special case of Clos-Network

Private address room: $10.0.0.0/8$

use 10 . pod. switch. \uparrow

Switch where \uparrow is enumerator for edges first, then $(k/2 + 1)$ for aggregation

End hosts are addressed: 10 . pod. switch. ID where ID is enumerator "from left to right" starting w/ 2

Also via Layer 2 forwarding possible

Advantages of k pod fat tree:

All switches are identical, cheap ones can be used
multiple paths w/ equal cost

Disadvantage: high cabling complexity

Load Balancer

Typically on per application
distributes load to servers

that server might delegate that traffic/load to another server

collects info of traffic from servers

can do NAT

hides internal net structure

Hierarchical Topology

Supports scalability

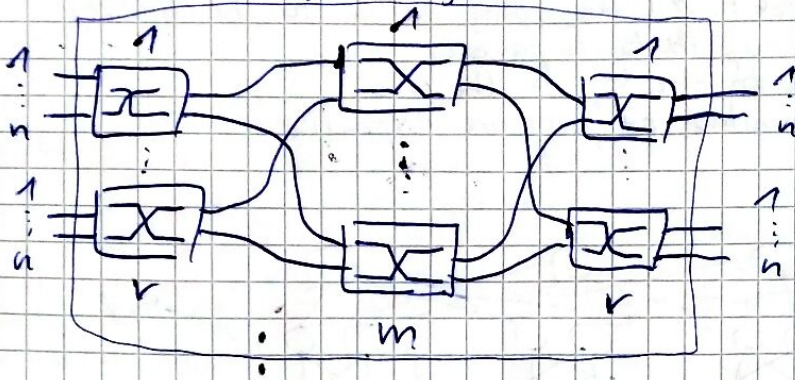
redundancy ⇒ multiple point of failure ⇒ reliability

availability

Clos Network

Multi routed tree

allows forming a large switch w/ smaller ones



Referred to as an

(m, n, v) -switch

$m \hat{=}$ # switches in the middle stage

$n \hat{=}$ # in-/outputs

$v \hat{=}$ # in/out switches

Recursive Construction

Three levels of switches

input / Middle / Output

Clos Network can not be limited to these 3
replace middle switches by 3-level Clos Network

Non-blocking @ $m \geq 2n - 1$

Input switches \rightarrow ToR @ Servers

Middle " \rightarrow Aggregation @ ToR

Output " \rightarrow Intermediate @ Aggregation

Ethernet @ Data Centers

Ziel Vereinheitlichung für SANs (Storage area network)

Aber Ethernet bietet Best-Effort-Dienst

keine Garantien

variablen Verzögerungen

Spanning-Tree zwar schleifenfrei

aber nicht flexibel, keine optimalen Pfade

\Rightarrow Unterauslastung der Netztopologie

Provisionierung mit Q-Header

\rightarrow TPID (Tag protocol identifier), 16bit

\rightarrow TCI (Tag control information)

3b Priorität

1b Drop eligible Indicator

12b VLAN ID

Data Center Bridging

Vereinheitlichte, auf Ethernet basierende Lösung für unterschiedlichste Anwendungen in Datenzentren

Erweiterungen zu Ethernet:

Prioritätsbasierte Flusskontrolle

Enhanced Transmission Selection

Starkontrolle auf Schicht 2

Data Center Bridge exchange

Prioritätsbasierte Flusskontrolle

- ⇒ Datenverluste durch Stausituationen vermeiden
- Erweiterung des PAUSE Flags um 8 Prioritäten
- Differenzierter QoS möglich

Enhanced Transmission Selection

- Reservierung von Bandbreite durch Einführung von Prioritätsgruppen
- Kann mehrere Prioritätsstufen einer Verkehrstyps enthalten
- Garantie einer minimalen Datenrate
- Rest kann von anderen Prioritätsgruppen genutzt werden

Staukontrolle auf Layer-2

Quantized Congestion Notification (QCN)

- Versendet auf Basis der Warteschlangenlänge im Buffer
- an Verursacher per speziellem Rahmen

Data Center Bridge Exchange Protocol (DCBX)

- Erkennung der Fähigkeiten und Konfiguration der Nachbarn
- z.B. prioritätsbasierte Flusskontrolle
- periodische Broadcasts an alle Nachbarn

Alternativen zum Spanning Tree

- Ziele: Mehr Flexibilität & Ressourcennutzung
- bei Topologieänderung, Lastbalancierung, Fast-Recovery
- und zur Nutzung besserer Pfade

Skalierbarkeit für Netze mit > 1000 Brücken

Einsatz des Link-State-Routings, modifiziertes IS-IS

- Basiert im Gegensatz zu OSPF nicht auf IP-Adressen
- keine manuelle Konfiguration
- Unterstützung von Multipath

Nutzung von Tunneling

- Kapselung der Rahmen innerhalb einer Domäne
- Hinzufügen eines neuen Headers für die Weiterleitung

Shortest Path Bridging

- Jede Brücke berechnet kürzeste Wege im LAN
- Pfade müssen symmetrisch sein

Lernen die MAC-Adressen

Unterstützt Equal Cost Multigraph

Gleiche Pfade für Unicast & Multicast

MAC-Adresstabellen zur Berechnung kürzester Pfade genutzt

Ziel unbekannt? → erstmal in kürzesten Pfad weiterleiten.

Pakete werden im LAN gekapselt (MAC in MAC)

Transparent Interconnection of Lots of Links (TRILL)

Neuer Header für Einkapselung

Routing Bridges implementieren TRILL

RB1 Kapselt Rahmen von S

und gibt RB3 als Ziel an, dort befindet sich D

RB3 entkapselt Rahmen



Jede RB berechnet kürzeste Wege zu allen RBs
 Berechnung mehrerer Bäume möglich
 Bäume für Multidestination genutzt
 Weiterleiten, falls Destination unbekannt
 ⇒ Kleine Forwarding Tables, keine Endsysteme
 Header:

Outer Ethernet	TRILL	Innen Ethernet
MAC-Adresse für Punkt-zu-Punkt Weiterleitung. Sieht auf jedem Hop	Nickname der Ingress Bridge und Egress "Multidestination Flag, Hop Count"	Quell- und Ziel MAC-Adresse der Endsysteme

TCP within Datacenters

Wichtiges Hintergrundwissen:
 oft sehr kurze RTTs, Server stehen dicht beieinander μs statt ms
 Incast: oft many-to-one Übertragung
 Multi-Paths, es existieren redundante Pfade im Datacenter
 Lange & kurze Flows gemischt (Suchanfrage vs. Backup)
 Virtualisierungen beeinflusst RTT
 Ethernet bereits für Star- und Flusskontrolle
 Commodity Switches mit kleinen Buffern (sogar geteilt unter Ports)

Ziele

Gleichzeitige Unterstützung von geringen Verzögerungen und hohem Durchsatz
 Hohe Auslastung der Ressourcen / Infrastruktur
 Geringe Kosten

Möglichkeiten

Datacenter hat eine administrative Kontrolle
 Meist homogene Komponenten
 Backward compatibility
 Interner & externer Verkehr gut separierbar

TCP-Reno

benötigt große Buffer
 für große Ende-zu-Ende Verbindungen geeignet
 Verbesserte ECN und RED → DCTCP (Data Center TCP)

TCP Incast-Problem

Häufig, dass mehrere benachbarte Server quasi synchron antworten
 race-condition, Buffer voll
 kann auch Flows betreffen, die gar nicht ~~zu~~ den Incast verursacht haben. Da Buffer geteilt werden (unter Ports)
 ⇒ Retransmission Timeouts
 ⇒ Incast Collapse

Barrier Synchronization

Client hat TCP Verbindungen zu n Servern offen und versendet ab alle n gleichzeitig eine Anfrage. Die langsamste Antwort bleibt wohl aus. Buffer voll. Warte auf Retransmission timer (200ms \gg RTT im Datacenter)
 Zeit, in der TCP-Verbindung keine Daten übertragen kann
 ⇒ Anwendung blockiert

Verbesserung durch Verringerung des Retransmission timers
→ RTT im Data center

Desynchronisierung für retransmissions, sonst passiert das gleiche wieder

⇒ Randomisierung $timeout = (RTO + (rand(0,5) + RTO)) \cdot 2^{backoff}$

Trotz niedriger RTT werden hohe Latenzen beobachtet

Ursachen:

Switches mit gemeinsamen Puffern für mehrere Ports
Incast (kürzere timeouts können helfen)

Abbau der Queue

Langlebige Verbindungen füllen die Puffer
kann zu synchronem Paketverlust führen

TCP-fairness hilft auch nicht...

Data Center TCP (DCTCP)

Ziel: mit commodity switches hohe Bursttoleranz, geringe Latenzen und hohen Durchsatz erreichen

Eigenschaft: DCTCP arbeitet mit geringem Füllstand der Queues ohne den Durchsatz zu reduzieren
Umgesetzt durch Reaktion auf Stärke eines Staus, nicht auf die alleinige Präsenz

ECN-Markierung und RED Schwellwert basieren auf Füllstand der Queue, werden ab Position k gesetzt

DCTCP setzt das ECN-Echo-Flag (E) nur bei markierten Paketen Quittierungen für

Staukontrolle beim Sender

Schätzung für den Anteil α der mit ECN markierten Segmente

$$\alpha = (1-g)\alpha + gF, \quad 0 < g < 1 \quad F = \text{Anteil markiert im letzten Fenster}$$

Für jedes Fenster von Segmenten neu berechnet

α schätzt Wkkeit, dass Länge der Queue $> k$

⇒ reagiert auf Stärke des Staus

Verfahren nach AIMD, falls keine Segmente markiert

$$cwnd = (1 - \alpha/2) cwnd \text{ sonst}$$

Vorteile von DCTCP:

Incast, falls sich die Queues über mehrere RTTs aufbauen, hilft frühes Reagieren von DCTCP

Aufbau einer Warteschlange wird bei k Segmenten instantan behandelt ⇒ mehr Puffer für Bursts vorhanden

Buffer pressure: Flows beeinflussen sich nicht mehr so stark.

10 - TCP Evolution

Selektive Quittungen

Ack ~~like~~ like in RFC 793 have inefficient loss recovery:
 only one lost segment can be recovered per RTT
 multiple losses often lead to retransmission timeouts
 and head-of-line blocking

NACK (negative ACK)

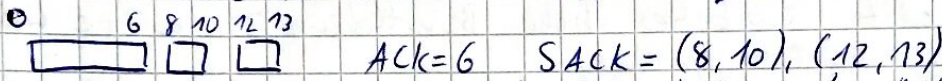
ACKs missing data explicitly

SACK (selective ACK)

explicitly ACK out-of-order data

on incoming in-order data: ACK as usual

" out " " : SACK newly received segment
 " " " " but adjacent to previously SACKed out-
 of order data: SACK whole block



Duplicate SACK (DSACK) Acknowledges segments that arrived multiple times

Forwarding ACK (FACK!) specifies when data may be sent during loss recovery - Replaces Fast Retransmit & Fast recovery

SYN-Cookies

against SYN Flooding $\hat{=}$ spamming SYN-packets w/ fake source address. Server can't handle more ^{TCP} connections

Idea: code timestamp, MSS-size, crypto-hash into seq number
 decode ACK-1 of client and check for correct content
 otherwise drop it \Rightarrow no connection

can't use special options that usually would be coded in seq number,
 MSS can't be arbitrary
 \Rightarrow use only in case of likely DOS-Attack

TCP @ high-speed

RTT = 100ms

Datarate $D = 10 \text{ Gbit/s}$

MSS = 1500 B

Using TCP Reno

assuming 1 packet drop each 5msrd. segments

\Rightarrow time needed to reach 10 Gbit/s (83k segments/s) = 70min

idea: use multiple TCP connections asynchronously.

idea: increase congestion window faster.

TCP CUBIC

growth independent of RTT



Avg. datarates w/ AIMD

$$D = \frac{1}{RTT} \cdot \sqrt{\frac{\alpha}{2} \frac{2-\beta}{\beta} \frac{1}{p}}$$

$\alpha = 1, \beta = 1/2$

$$\Rightarrow D = \frac{1}{RTT} \sqrt{3/2 \cdot 1/p}$$

to reach same increment w/
 CUBIC-TCP

$$W_{TCP} = W_{max} * (1-\beta) + \alpha \frac{1}{RTT}$$

$$\alpha = 3 \frac{\beta}{2-\beta}$$

Compound TCP (CTCP)

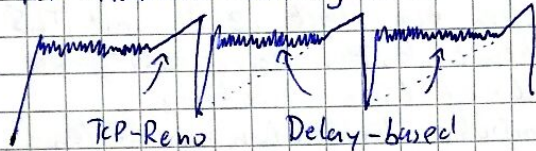
measures RTT growth \Rightarrow implicit queue growth
 \Rightarrow aggressive increase if queues are empty

$$W = W_{REUO} + W_{delay} (RTT)$$

shortest RTT measured = queues empty
needs to adjust RTT continuously

measuring needs to be precise

TCP fairness not guaranteed if initial RTT was luckily low



TCP & Web

high delay in responses \rightarrow less customers on that site \Rightarrow less profit
 \Rightarrow aim to keep delay as low as possible

handshake, slow start, transmission of actual data ... would wide wait

\uparrow 2 RTT

\rightarrow minimum time for transmission: $2RTT + \frac{\text{object size}}{\text{data rate}}$

whereas 2RTT dominates for small objects

If $W = 2MSS$

*K for transmitting object with size o : $k = \left\lceil \frac{o}{W \times MSS} \right\rceil$

\Rightarrow respond time $2RTT + \frac{o}{D} + (K+1) \left\lceil \frac{MSS}{D} \right\rceil + RTT - \frac{W-MSS}{D}$

Server needs to wait for ACK of client \rightarrow

If $\frac{W-MSS}{D} \geq RTT + \frac{MSS}{D} \Rightarrow$ server can send w/o pauses

Elif $\frac{W-MSS}{D} < RTT + \frac{MSS}{D} \Rightarrow$ server needs to wait for ACK

RTT dominates web-delay, because index.html often < 10 kB

Slow Start \Rightarrow small window size @ beginning \Rightarrow RTT dominates delay

Idea: increase initial *MSS in Cwnd (today $\approx 3-4$)

TCP-Fast-Open (TFO)

Idea: Data request in SYN-Packet

\Rightarrow - 1 RTT

but opens room for attacks, DoS via SYN-Data Flooding

Source address spoofing

\Rightarrow TFO Cookie

Client needs to request TFO cookie via regular TCP

Server generates unique - only for this client working cookie
can be used afterwards

Server replies w/ cookie in SYN-ACK Packet

Client can then use this cookie for later use

sends SYN + Cookie + Data request

Server validates source IP & cookie and replies w/ SYN-ACK

and begins immediately w/ data sending

Server can disable cookie in case of DoS Attack assumption

Realitätscheck?

Früher war Engstelle im Protokollstack IP

ist heute ~~HTTP~~ HTTP

Anwendungen benötigen andere Kommunikationsabstraktionen (~~HTTP~~!) HTTP

API

Klassisch: TCP-Bytestrom & Socket Interface

de-facto: HTTP-Interface

um fasst TCP als Teil des Substrates

Transportschicht ist heute eher ein Stack, z.B. mit TLS

Reihenfolge treue bei TCP führt zu hohen Latenzen. Darüberliegende Anwendungen werden gebremst. Head-of-Line-Blocking
Streaming leidet (ohne ein \downarrow), dynamische Anwendungen leiden

SPDY und HTTP2

sollen solche Probleme lösen

höhere Sicherheit

kurze Antwortzeiten

Rückwärtskompatibilität

nur Client und Server müssen angepasst werden

Realisiert durch:

1 TCP Verbindung pro Server

→ keine erneuten 3-Wege-Handshakes

kein Slow Start

Multiplexing bei Anfragen ⇒ 1 Stream via TCP

Priorisierung " "

Server-push $\hat{=}$ senden an Client ohne dessen explizite Nachfrage

Effizientere Header

verpflichtende Nutzung von TLS

Framing

Liefert Einheiten für das Multiplexing

Data Frames, Fragmente der Daten

erfahrungsgemäß bessere Quality of experience

etwas mehr Overhead

Control Frames

um Stream und dessen Frames zu verwalten

Probleme

Anlaufschwierigkeiten

HTTP umging mit vielen Verbindungen dieses Problem

ist nun explizit verboten

höhere initCwnd

Slow-Start-after-fall sollte deaktiviert sein

Verlust eines Segments blockt alle gemultiplexten Streams

⇒ keine Lösung mit TCP

⇒ neues Protokoll

Quick UDP Internet Connections (QUIC)

0-RTT-Aufbau, ähnlich zu TFO

Reihenfolge treue nur innerhalb eines Streams

TLS by default

Starkontrolle angelehnt an TCP

Fairness: eine QUIC mit n Streams verhält sich wie n TCP Verbindungen

QUIC übernimmt Aufgaben von HTTP2, TLS und TCP
optimiert: nutzt seq Num (zuverlässig) und Retransmissions
erhalten neue " " Initialisierungsvektoren verschlüsselt

"0-RTT" durch gemeinsamen Handshake von Transport, Verschlüsselung
und HTTP-Anfrage in der selben RTT
Nur möglich bei Bestehen einer Verbindung von Kurzem (oft bei Google der Fall)

QUIC ~~setzt~~ setzt auf UDP auf.
Könnte direkt IP nehmen, ist aber mit NAT und Firewalls nicht gut
kompatibel

MultiPath TCP (MPTCP)
für Geräte, die drahtlos sind und von Access-Point zu Access Point wechseln
Multi-Homing für Server-Farmen
Datenzentren mit redundanten Pfaden
Identifikation aber statisch über Quell-IP, Ziel-IP, Quell-Port, Ziel-Port
keine Änderung während offener Verbindung möglich
⇒ TCP ist ein Single-Path-Protocol

~~MPTCP~~
⇒ MPTCP
Beispiel: Iphone mit 3G & WiFi
Hat 2 verschiedene IP Adressen, Verbindung kann nicht unigrieren

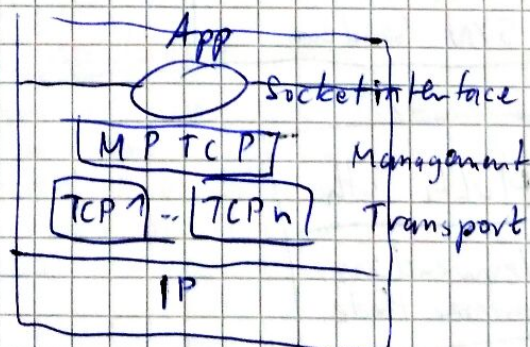
Ziel: Erweiterung für TCP zur Nutzung von mehreren Pfaden
+ Leistung
+ Availability

Muss Kompatibilität zu Anwendungen bereitstellen
sowie " " zu allen TCP Geräten
Middle Boxes ändern z.T. TCP-Optionen und halten den Zustand

MPTCP-Verbindung
ist eine Kommunikationsbeziehung zwischen Sender & Empfänger
besitzt einen oder mehrere MPTCP-Subflows

MPTCP-Subflow
Beziehung zwischen Quell-IP-Adresse und Ziel-IP-Adresse
≙ reguläre TCP-Verbindung
Reihenfolge treu
Endet mit RST oder FIN
3-Wege-Handshake
Können dynamisch zu einer Sitzung hinzugefügt oder gelöscht werden

Struktur



TCP Option MP-Capable im Handshake
 senden Token X, Y zum späteren Hinzufügen/Löschen eines Subflows
 mit

MP Join, Token

MPTCP Segmente haben zwei seqnums. Eine für gesamte MP-Verbindung
 andere für den subflow

⇒ Lückenlose Nummern für subflows

Beachten: Eingangspuffer: einer für die gesamte MPTCP Verbindung
 Scheduling auf höherer Schicht, muss die Reihenfolge und
 Zuordnung erledigen

Steuerkontrolle wäre mit gleichem Algo unfair gegenüber TCP
 neuer Ansatz:

bei Verlust, wie gehabt: $Cwnd * = 0,5$

Balance beim Anstieg: weniger belastete subflows dürfen stärker ansteigen
 Congestion Avoidance:

✗ Quittung, die im subflow r empfangen wurde:

$$Cwnd_r = Cwnd_r + \min\left(\frac{\alpha}{Cwnd_{Gesamt}}, \frac{1}{Cwnd_r}\right)$$

Grundlegende TCP-Verfahren

TCP Tahoe
 Slow Start
 Congestion Avoidance
 Fast Retransmit

TCP Reno
 Wie Tahoe
 + Fast Recovery

Hohe Bandbreiten

Parallel
 Mehrere parallele
 Verbindungen

CUBIC Linux
 Kubische Funktion
 unabhängig von RTT

Compound ^{Windows}
 Proaktiv RTT

Web-Basiert

Initiales Fenster
 $C_{init} \geq 10$

Fast Open
 Daten bereits mit
 SYN Senden

Data Center

Data Center
 Verwendung
 von ECN

Multi Path
 Verwendung
 mehrerer Pfade

Verzögerungsverzerrung

bedingt durch Ausbreitungsverzögerung

auch Frequenzvariierend

⇒ Phasenverschiebung zwischen den Frequenzen

⇒ Intersymbol-Interferenz
beschränkt maximale Datenrate

Rauschen

thermisch, abhängig von der Temperatur

Gleichmäßig verteilt über die Frequenzen (weißes Rauschen)

kann man nichts gegen tun

Modulationsrauschen

durch Signalüberlagerung $f = f_1 + f_2$

Übersprechen

durch Kopplung zweier Kabel, die sich beeinflussen

Impulsrauschen

durch elektromagnetische Störungen, Blitze

analog: Knacken

Digital... fatal

Unregelmäßiger Ausschlag mit hoher Amplitude

Kanal Kapazität

≙ maximale Rate, mit der Daten über einen Kommunikationskanal gesendet werden können

Datenrate ≙ bit/s zur Kommunikation

Bandbreite ≙ höchste - niedrigste Frequenz [Hz]

Rauschen

Fehlerrate = Error/s 0 statt 1 oder 1 statt 0

Ziel: Maximiere Datenrate bei gegebener Bandbreite und minimiere dabei die Fehlerrate

Nach Nyquist gilt:

max. Schrittgeschwindigkeit [baud] = 2 · Bandbreite [Hz]

bei binären Signalen:

Datenrate = $2n$ bit/s bei einer Bandbreite von n Hz

bei mehrwertigen Signalen:

Datenrate $C = 2n \log_2(M)$

M ≙ Signalwerte

Erhöhung von M führt zu einer Erhöhung der Datenrate

bei hohen Datenraten sind mehr bits von einer Störung betroffen, als bei gleicher Störung und niedrigerer Datenrate

Shannon-Hanley-Gesetz

$$C = B \log_2(1 + \text{SNR})$$

Signal to noise ratio (SNR)

C ≙ Datenrate [bit/s]

B ≙ Bandbreite [Hz]

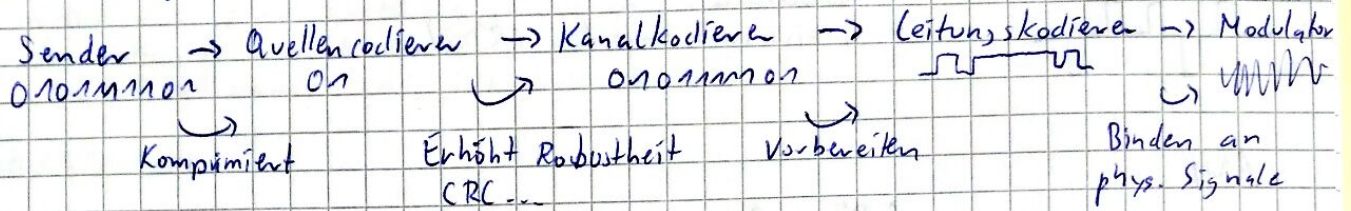
S ≙ Energie des Signals

R ≙ " der Störung

$$\text{SNR [dB]} = 10 \log_{10} \frac{S}{R}$$

Digitale Leitungscode

Kodierung beim Sender



Taktrückgewinnung erforderlich, falls keine separate Taktleitung dabei

Gleichstromanteil

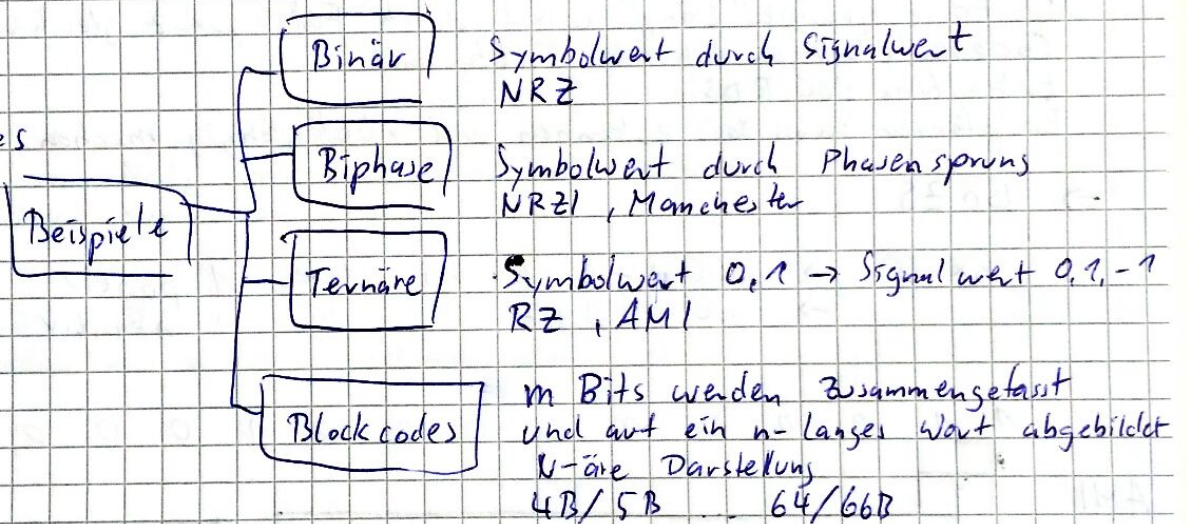
Wegen der angeschlossenen Geräte möglichst vermeiden
Bewertung RDS (Running digital sum)
Gleichstromfrei, wenn statisch bei 0

Baudrate

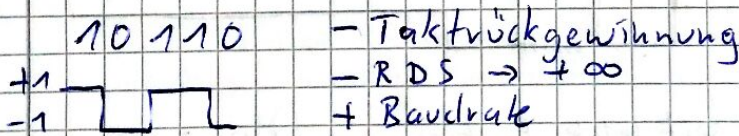
hoch \rightarrow viele Wechsel \rightarrow hohe Anfälligkeit \rightarrow geringere Reichweite

Komplexität

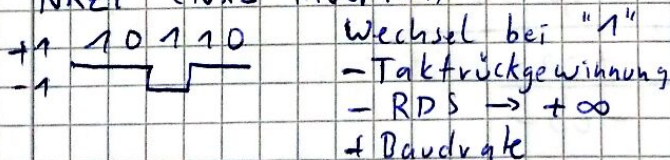
Leitungscode



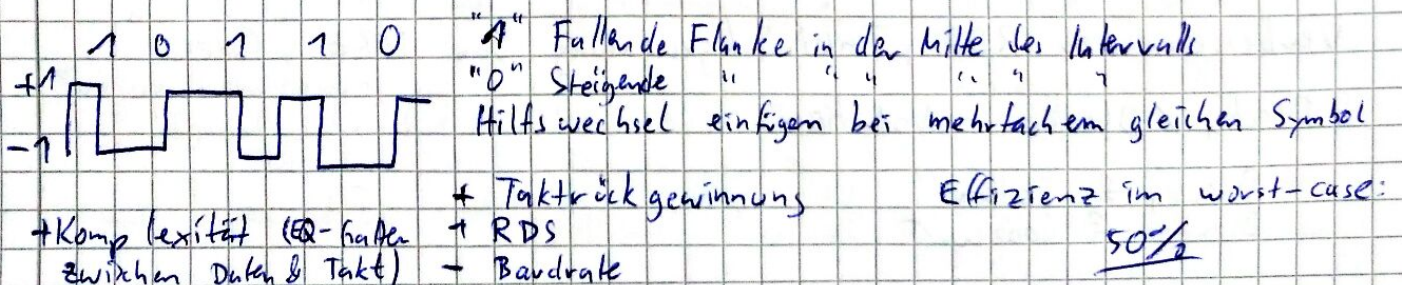
NRZ (non-return to zero)



NRZI (NRZ inverted)



Manchester



NRZ (Return to zero)

1 0 1 1 0

- + Taktrückgewinnung
- Baudrate
- Gleichstromanteil (RDS)



AMI

1 0 1 1 0

- "0" = 0
- "1" = abwechselnd -1 oder 1



- Taktrückgewinnung
- + RDS
- + Baudrate

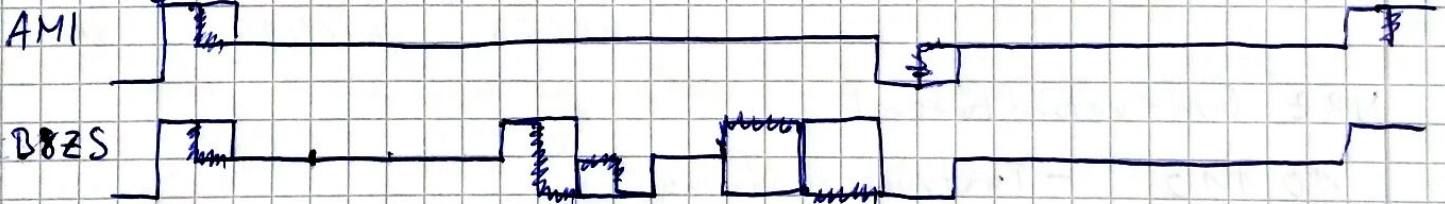
Modifizierte AMI Codes

Verstoße gegen die Regel der abwechselnden Pegeln bei einer 1, um Taktrückgewinnung zu gewährleisten
 0-Folge gewisser Länge wird durch Füllsequenz gleicher Länge ersetzt
 Codeverletzungen erhöhen Taktgehalt
 Erhaltung der RDS
 Empfänger muss das erkennen und rückgängig machen

⇒ B8ZS

8*"0" → 000+-0-f bei letzter 1 positiv
 → 000-+0+- " " " negativ

1 1 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1



HDB3

4*"0" →

* "1" seit letzter Ersetzung

	gerade	ungerade
negativ	000-	+00+
positiv	000+	-00-

+RDS ☺

Block codes

verbessern der RDS, Taktrückgewinnung von NRZ - Codes

Block₁ / Block₂ - Notation

Block₁ ≙ Daten
 Block₂ ≙ Codewort

Block₂ wird so gestaltet, dass er Redundanz ermöglicht
 hat mehr Wörter, als Block₁ (mehr Bits)
 darf ~~aber~~ höherwertig als Block₁ sein
 nutzt nur eine Untermenge seiner möglichen Zustände
 Ausnutzung von Redundanz für
 Taktrückgewinnung und Gleichstromfreiheit

4B/5B Code

Block₁: Länge 4bit, binär
 Block₂: Länge 5bit, binär

Auswahl von Block₂: nicht mehr als eine führende 0
 " " " zwei abschließende 0en

Effizienz: 80% = $\frac{4}{5}$

$\begin{matrix} \hookrightarrow 1000 & 0000 & 0000 & 0001 & 1111 \\ \hookrightarrow 10010 & 11110 & 11110 & 01001 & 11101 \end{matrix}$

Verbleibende Codewörter sind entweder ~~un~~ ungültig oder werden
 anders interpretiert, zum Beispiel als Start-of-Frame-Delimiter

8B/10B

4B/3T

Auswahl von Block₂ mit dem niedrigsten Gleichstromanteil
 Keine Verwendung von 000

64B/66B

Overhead ~3%

Voranstellen von 2 bit Synchronisation

Modulationsverfahren

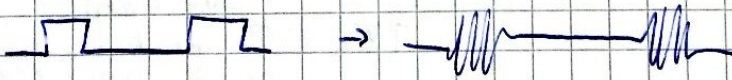
digitales Signal \rightarrow analoges Signal
 Grundlage Fourier-Analyse

$$x(t) = \frac{C_0}{2} + \sum_{n=1}^{\infty} \frac{C_n}{2} \cos(2\pi n f_0 t + \theta_n)$$

Amplitude
Frequenz
Phase

Amplitudenumtastung

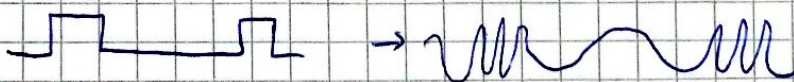
Signal (binär) schaltet zwischen Oszillator und GND



Technisch einfach
 wenig Bandbreite
 störungsanfällig

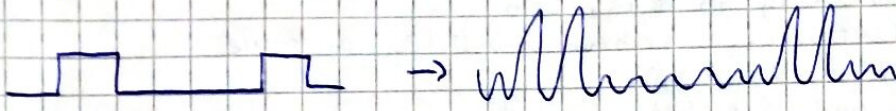
Frequenzumtastung

Signal schaltet zwischen zwei Oszillatoren um



größere Bandbreite
 mit oder ohne Phasen-
 Sprünge

Phasenumtastung,
Signal hält Pegel bei 0, wechselt bei 1



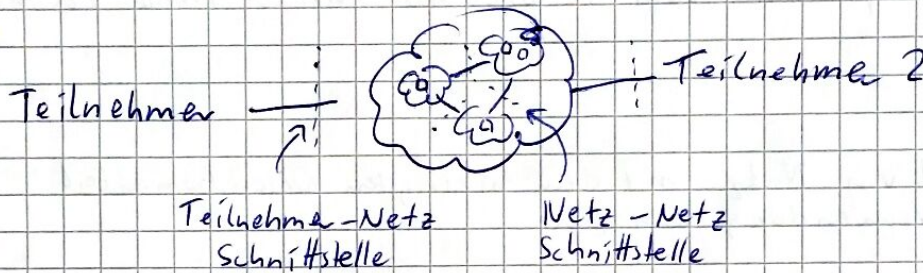
komplexe Demodulation mit Taktrückgewinnung
relativ störungssicher

12. ISDN (Integrated Services Digital Network)

Leitung $\hat{=}$ Reservierung von Ressourcen für die Übertragung von Nutzdaten
keine physikalische Leitung
Ressourcen werden fest einer Verbindung zugeteilt
Keine Metadaten / Pufferüberläufe
Aber ggf. schlechte Ressourcennutzung

Leitungsaufbau durch Signalisierung
Routing und Zuweisung / Rücknahme von Ressourcen
i. d. R. Einsatz von Paketvermittlung

Aufbau



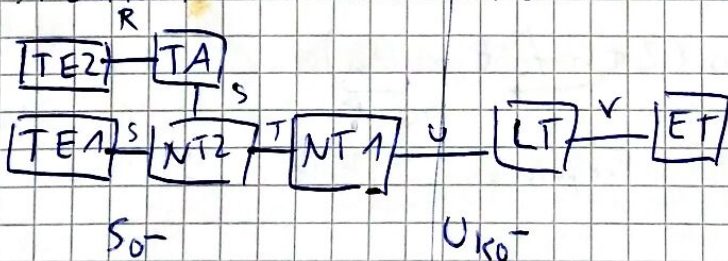
Hierarchie intern

Zentral (vollvermascht)
Hauptvermittlung
Knotenvermittlung
Ortsvermittlung

Ziel von ISDN

Digital bis zum Teilnehmer
Integration unterschiedlicher Dienste (Bild, Ton, Sprache, Daten)

Architektur Teilnehmer Netzseite



- ET $\hat{=}$ Exchange Termination
Vermittlungsabschluss
- LT $\hat{=}$ Line Termination
- NT1 $\hat{=}$ Network Termination 1
Netzabschluss Schicht 1
- NT2 $\hat{=}$ Network Termination 2
Netzabschluss Schicht 1-3
- TE1 $\hat{=}$ Terminal Equipment 1
ISDN-fähig
- TE2 $\hat{=}$ Terminal Equipment 2
nicht-ISDN-fähig
- TA $\hat{=}$ Terminal Adaptor
für TE2

Aufgabe des NT

Übertragungstechnischer Abschluss
der Netzseite (U_{ko}-Schnittstelle)
und der Teilnehmerseite (S_o-Schnittstelle)

Speisung der Teilnehmerinstallation
220V-Netz

Erkennung von Rahmenfehlern

NT1:

Taktückgewinnung
Rahmensynchronisation
Echo-Kanal-Steuerung

NT2:

Vermittlungsfunktionen
Ist optional

ET: (Schicht 1-3)

Medienzugriff auf D-Kanal (Schicht 2)
Signalisierungsfunktionen in Schicht 3
Multiplex- und Demultiplexfunktionen

LT (Schicht 1)

Umsetzen der Übertragungsvorfahren
Ableiten und Generieren von Takten
Fernstromversorgung des Teilnehmerbereiches

Basis-Anschluss

2 B + 1 D Leitungen (2 · 64 kbit/s + 16 kbit/s)

B-Kanal (Schicht 1-7)

zur Nutzdatenübertragung
operiert unabhängig von anderen D-Kanälen
können verschiedene Zielrichtungen haben
" " " Datentypen übertragen
" gleichzeitig operieren

Medienzuteilung fest

wird dediziert zwei Endpunkten zugeteilt
diese Zuteilung ist Aufgabe des D-Kanals

D-Kanal (Schicht 1-3)

16 kbit/s schnell (oder langsam (ol))
Bidirektional zwischen NT und Teilnehmer
Koordiniert bis zu 8 Teilnehmer
CSMA/CD Variante

E(ch)o-Kanal

16 kbit/s

Richtung NT → Endsystem

Für Medienzugriff erforderlich

Carrier Sensing (CS)

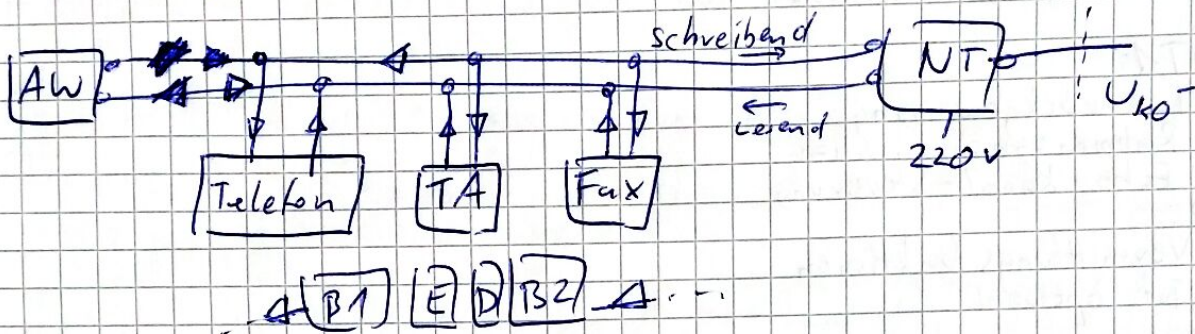
Collision Detection (CD)

Teilnehmer Schnittstelle S₀

Doppelader in jede Richtung, Simplex
verschiedene Topologien möglich

Raummultiplex - Richtungsstrennung

Zeitmultiplex für andere Geräte in der Topologie



Echo-Kanal

Aufgabe: Unterstützung des konkurrierenden Medienzugriffs

Problem: Endsystem kann auf ausgehenden D-Kanal nicht lesen

Lösung: NT reflektiert D-Kanal in E-Kanal

Leitungskodierung: Inverser AMI-Code

1 $\hat{=}$ Nullpegel

0 $\hat{=}$ Wechselnd 1 oder -1

1 0 1 1 1 0



D-Kanal-Medienzugriff

falls 8 bit lang nichts auf dem E-Kanal liest \rightarrow frei

Kollisionserkennung durch sendende Station

Senden ist 1-persistent

Anderes Signal auf E-Kanal als D-Kanal?

(0 überschreibt 1) Aufhören zu senden, weiter prüfen

kein exponentieller Backoff

Anderer sendende Station merkt davon nichts

Wieso 8 bit warten?

Inverser AMI, kann keine 8 bit lang keine Aktivität haben \rightarrow frei

HDLC Protokoll auf Schicht-2

Framebegrenzung durch Flag 01111110

Bitstopfen erhält Datentransparenz

bei 5 '1' wird eine '0' gestopft

Empfänger muss diese wieder entfernen

\Rightarrow Länge des Pakets auch abhängig vom Inhalt

Bei gleichzeitigem Beginn, ist für beide die Teilnehmererkennung im E-Kanal falsch.

D-Kanal & Signalisierung

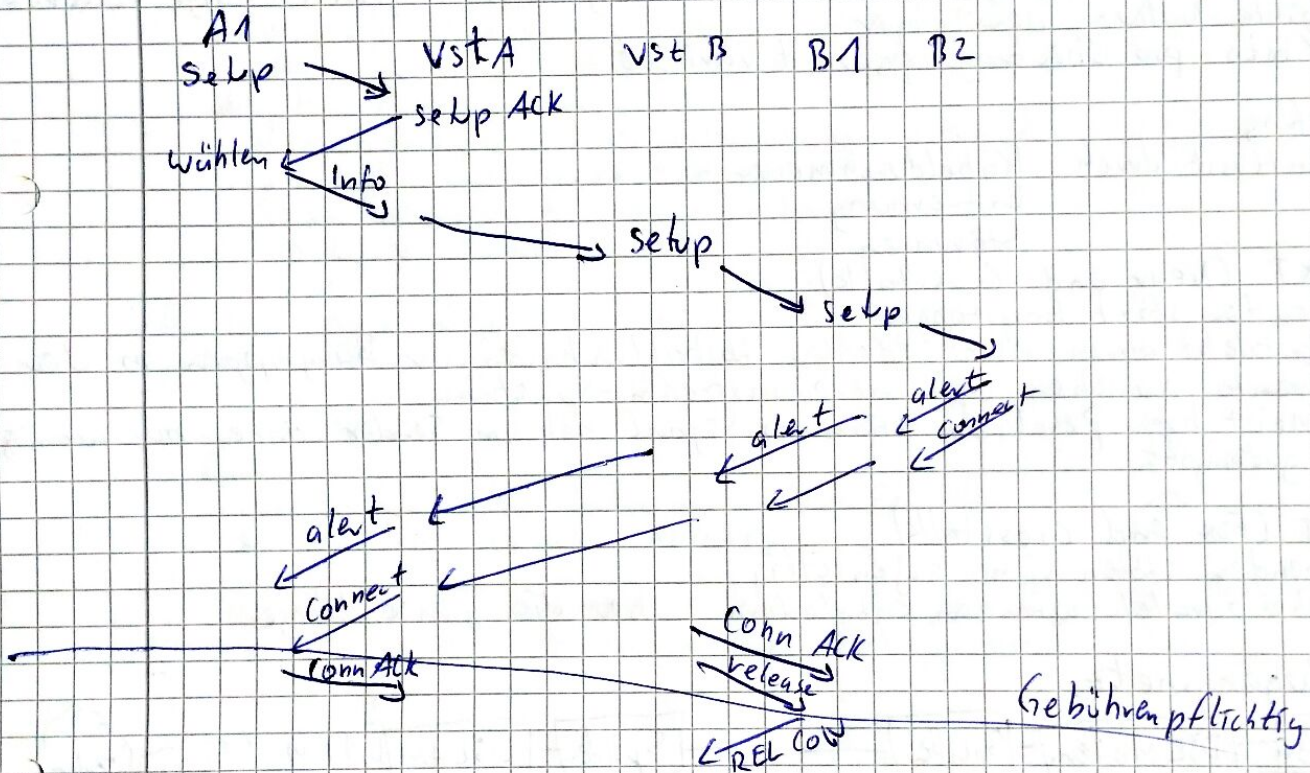
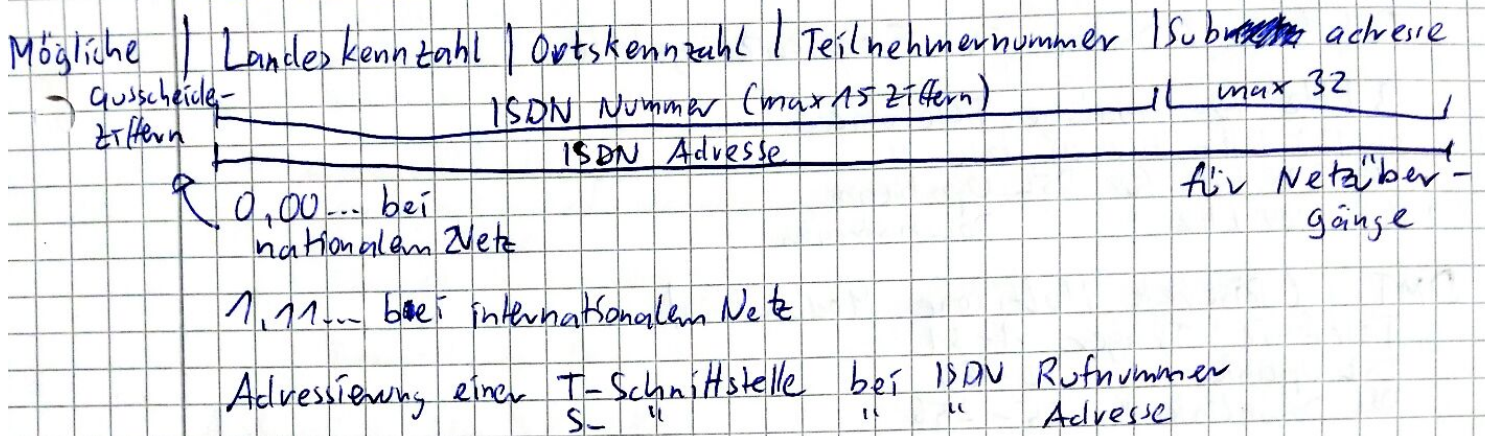
Q.931 Protokoll steuert SS-7

\Rightarrow Aufbau / Abbau von Verbindungen zwischen Endsystemen

und Art des Dienstes: Rückruf / Datenanfrage...

Benötigt ISDN eindeutige Adressen für Schicht 3

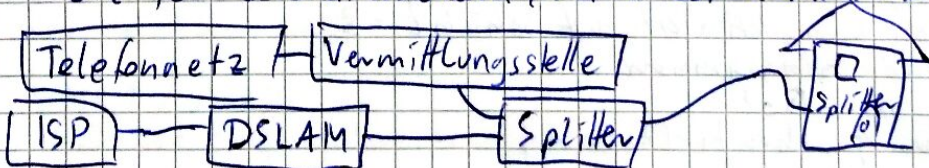
ISDN Adressen



DSL (Digital Subscriber Line)

Ziel: Leistungsfähige Leitung für Teilnehmer
 asymmetrischer Anschluß: 8 Mbit/s Download => ADSL
 576 kbit/s Upload

Bedingt durch das Modell: kleine Anfrage, große Antwort
 Symmetric DSL (SDSL)
 oft für Geschäftskunden, reine Datenanschluß, keine Telefonie



Splitter

Trennt Signal in Telefon / Datensignal
 Sowohl bei der Vermittlungsstelle als auch beim Kunden vorhanden
 Telefonie ist höher prior als Daten - Bei Ausfall bleibt diese erhalten

DSLAM (DSL Access Multiplexer)
 Gegenstück zum DSL-Modem beim Teilnehmer

Datenübertragung

Telefonie / Daten in Frequenz geteilt

3 kHz für analoge Telefonie

0-138 kHz für digitale "

138-276 kHz für DSL Upstream

276-1104 kHz " " Downstream

DMT (Discrete Multitone Modulation)

Telefonie: Träger 1-31

DSL Uplink: " 32-64

DSL Downlink: " 65-256

parallele Übertragung auf mehreren orthogonalen Trägern (niedrige Datenrate)
Flexible Nutzung der Träger
bits pro Übertragungsschritt variiert

Dämpfung

beeinflusst durch Kabeldurchmesser
Entfernung
Störungen

NEXT (Near end Crosstalk)

Sender stört Empfänger

Entsteht durch nicht ideales Abstimmverhalten der Baugruppen an der Grenze zwischen Up- und Downstream-Spektrum

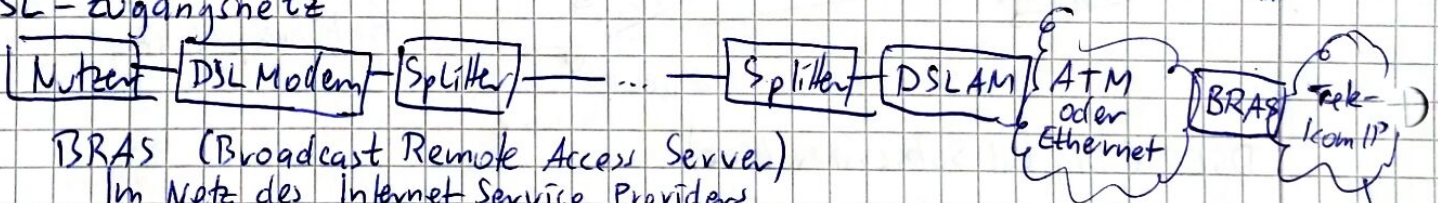
Meist hoher Pegel, da störendes Signal nah am Sender, daher nur wenig gedämpft

FEXT (Far end Crosstalk)

Sender stören sich gegenseitig

Bei parallel laufenden Kabeladern über die ganze Länge

DSL-Zugangsnetz



BRAS (Broadcast Remote Access Server)

Im Netz des Internet Service Providers

... erster Router im Netz Richtung Backbone

Aufgaben

Aggregation mehrerer DSLAM-Verkehre

Bereitstellung von Konnektivität auf Schicht-2

PPP-Sitzung (über Ethernet oder ATM (Asynchronous Transfer Mode))

IP (über ATM)

Bereitstellung von Konnektivität auf Schicht-3

Routing von IP-Datagrammen

Berücksichtigung von QoS

Schnittstelle zu Authentifizierung

Autorisierung

Accounting

RADIUS Protocol (Remote Authentication in user service) auf Schicht-7

Aufbau einer DSL-Verbindung

Erfolgt über PPP in 3 Phasen

1. Establish \rightarrow LCP

Aufbau einer PPP-Verbindung

Aushandlung der Verbindungsparameter (Datenrate / Träger)

" " Authentifizierungsmethode

2. Authentifikationsphase

basierend auf der ausgehandelten Methode (PAP, Challenge Handshake)

3. Network-Phase

Vergabe der IP-Adresse

über IPCP, das PPP Network Control Protocol (NCP) für IP

Bekanntgabe der Adresse des DNS-Servers

Aushandlung der Datenrate

Fixed Rate

Bestimmter Wert, Information aus Dämpfungstabelle
enthält Sicherheitszuschlag zum Erreichen der Datenrate
auch bei einer Störung eines neuen DSL-Anschlusses

Adaptive Rate

Aushandlung der aktuell maximalen Datenrate

ist höher als Fixed Rate

Resynchronisation bei Verschlechterung notwendig

\Rightarrow kurzer Ausfall

Aufbau einer ADSL-Verbindung

Abbruch des vorherigen Ablaufs in Authentifizierungsphase

Erst zu diesem Zeitpunkt ist klar, dass der Nutzer Kunde einer anderen
Provider ist

Danach

Weiterleitung aller Daten an anderen Provider

Neustart des gesamten Ablaufs



Weiterentwicklungen:

ADSL2

Erhöhung der Datenraten

ADSL2+

Nutzung größerer Frequenzen

Verbesserte Kodierung

\Rightarrow weniger Übersprechen

Energieeinsparung, Modem shutdown & wake

Anpassung der Datenrate während Betrieb on LAN

Downstream \rightarrow 25 Mbit/s

Upstream \rightarrow 1 Mbit/s

VDSL1 / VDSL2

Größerer Frequenzbereich \Rightarrow höhere Datenrate \Rightarrow weniger Strecke möglich

Outdoor DSLAM

Glasfaser bis zum Kasten \sim 500 m von Nutzer entfernt

\Rightarrow VDSL2 möglich = 50 Mbit/s kein Problem (bis 200 Mbit/s)

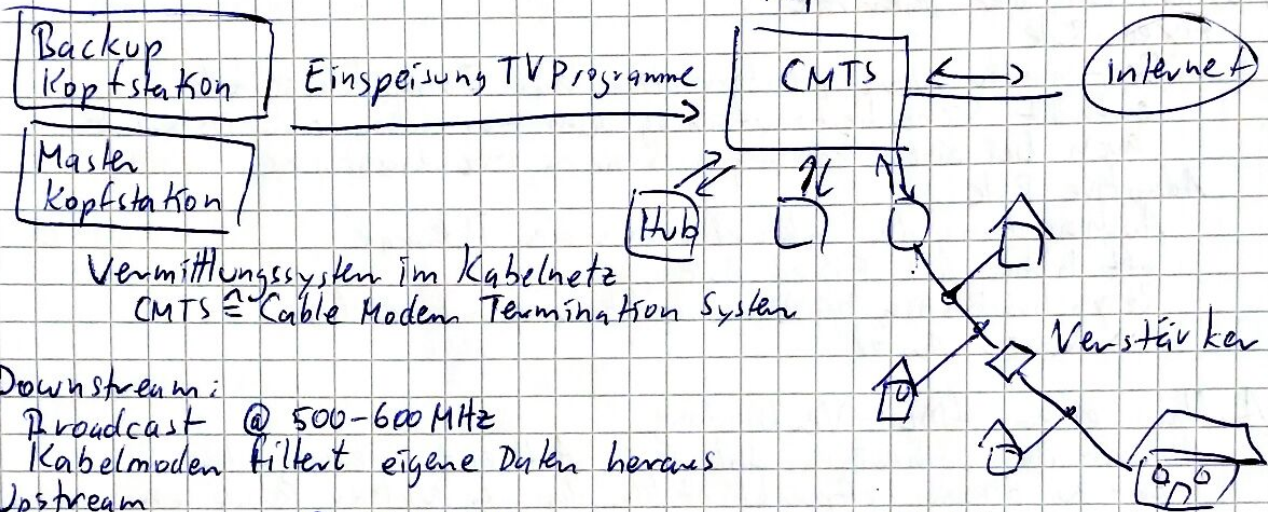
TV-Kabelnetz

Aufgebaut in den 80ern
zuerst ohne Rückleitung
ernster Ausbau
heute auch für Telefonie & Internet

Topologie
Baum nur auf letzter Meile
Ansonsten stark verzweigt

Erreichbarkeit
DE: ~75%

Architektur



Downstream:

Broadcast @ 500-600 MHz
Kabelmodem filtert eigene Daten heraus

Upstream

Zeitmultiplex @ 5-65 MHz bis 1 Gbit/s
Zeitschlitz werden vom CMTS zugewiesen

Shared Medium

⇒ Rate von Nutzung der anderen abhängig

Frequenzbereich:

5-862 MHz (meist auf 570 MHz begrenzt)

DOCSIS (Data over Cable Service Interface Specification)

Erlaubt Nutzung bis 1,7 GHz
⇒ Downstream bis zu 10 Gbit/s
Upstream bis zu 1 Gbit/s

Powerline

Nur für Innenhaus (Babyphone, Sprechanlage)
Daten werden aufmoduliert und die Stromkabel wirken wie Antennen
Störungen vom Funkdienst im gleichen Frequenzbereich
Verschlüsselung notwendig, sonst wireless-Broadcast durchs Land

LTE

Für flaches Land
wie DSL nur drahtlos
geteiltes Medium