

Lecture 3: Router Design

- Router lay-out in an ISP 10
- Router design in GSR 30
and challenges (McK97)
- Scaling and moving routers degree
why not more routers.
- Research - buffers
 - centralize switch vs. distributed switch.
 - power-aware design and algorithms.
- Takeaway.

- Access trees, Border routers, backbone routers | Show switching capacities.
- PoPs or central offices, stacks based organization
- Observations about link speeds and no degrees.
- Functions of various routers.
- Router design discussions → backbone routers.

Any router:
Control: coarse, but important for overall correctness
Data: fine-grained, fast, simple
Management: coarse, but critical to diagnose faults

Router design: fast data plane implementation in backbone routers.

Data plane functionality: forwarding, packet classification, DSCP marks, filtering, deep packet inspection, route/address lookup

Forwarding line cards and backplane organization, CPU location

Switching capacity needed.

Cross bar → parallelism.

$N \times N$

hierarchical switching systems...

① Scheduler role:

② Packet length: throughput

③ Blocking: HOL → vOB.
input, output ~ delay.

- simple, high throughput, fast and starvation free.
- describe the arbiter \rightarrow implementation.
- throughput is good, but delay is unpredictable
 - priority,
 - speedup. output port memory and link speeds.
- Multicast. - different queue.
 - ↳ usage natural multicast
 - fan-out and no-fanout-split

10

① Router scaling:

traffic volumes grow over time

- \rightarrow roughly doubles every 18 months or so.
- \rightarrow link capacities at edges increase, allowing more traffic to be pumped into core
- \Rightarrow Network switching capacity must also increase correspondingly. \rightarrow scale up core routers
- \rightarrow add a core router hierarchy, more routers and increase switching capacity.
- ~~But do we necessarily need routers which can switch more and more traffic.~~

Second option means we don't necessarily need to build faster and faster high-end routers

But is this feasible

- \Rightarrow central offices are full
- \Rightarrow building central offices expensive
- \Rightarrow increasing switching capacity and replacing older routers is a more viable option.

$1 \text{ Tbps} \rightarrow 10 \text{ Tbps} \rightarrow 100 \dots$

But what is the challenge for the switching fabric.

Racks today can provide 10kW per rack (cooling and placement becomes hard)

at 10kW \rightarrow can't switch more than 2.5 Tbps today

Switching challenge: centralized switch is limited in throughput.

has to scale α^2 where α is scaling factor for link capacity

10

And this is exactly how things have progressed over time so far

The higher and higher capacity switches consume more and more power.

↳ more and more racks to spread power density



one trend: distributed multi-stage switch \Rightarrow unpredictable performance.

A single stage switch ^{fabrics} that is central \Rightarrow good performance

↳ can fit on a single rack.

but scalability is an issue.

Switch hardware has to scale α^2 faster where α is the switch increase in link capacity

Today's racks \leftarrow a max of 10kW per rack (cooling is hard)

\leftarrow 2.5 Tbps due to the limitation of central arbiters

How do you build a single stack switch fabric
that can switch Tbps under 10 kW.

- ② A second challenge is in the design of buffers and
how these are designed, what memories to use.

How much buffer to provision?

→ Rule of thumb to ensure good utilization

40Gbps @ 0.5s → lot of memory.

SRAM . . . fast and cost - - \$ns.

DRAM lot fewer chips, less power-consuming but
slow - - 50ns

Multiple DRAM chips with a wide bus.

- Power and expensive.

Another issue is do we really need such large buffers.

- ③ Power aware design and routing