

# High Speed Networks Need Proactive Congestion Control Using Programmable Forwarding Planes!



*Yanya Jose*, Steve Ibanez, Lisa Yan, Nick McKeown,  
Sachin Katti

Stanford University

Mohammad Alizadeh

George Varghese

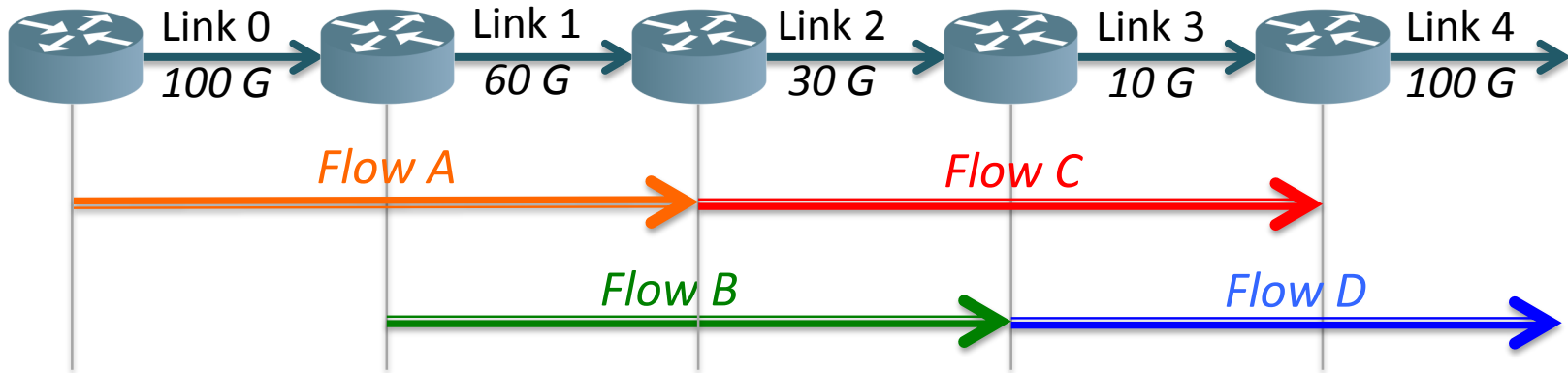
MIT

Microsoft Research

# Outline

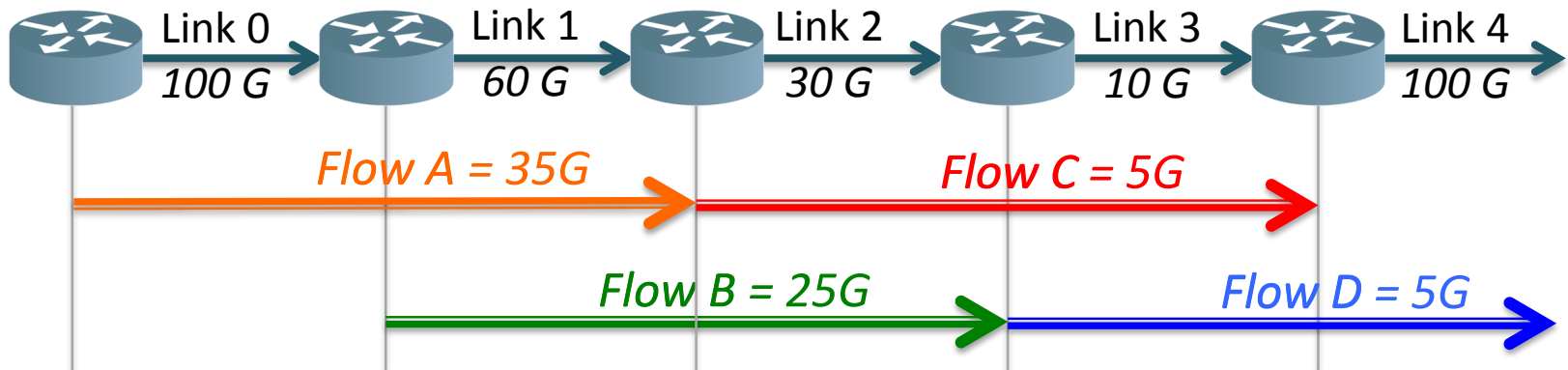
- At 100G speeds, we'll need much faster congestion control schemes
- Letting networks switches directly compute rates is a fast and scalable scheme
- We can realize such a scheme in 100G networks using programmable forwarding planes (stateful data planes)

# The Congestion Control Problem



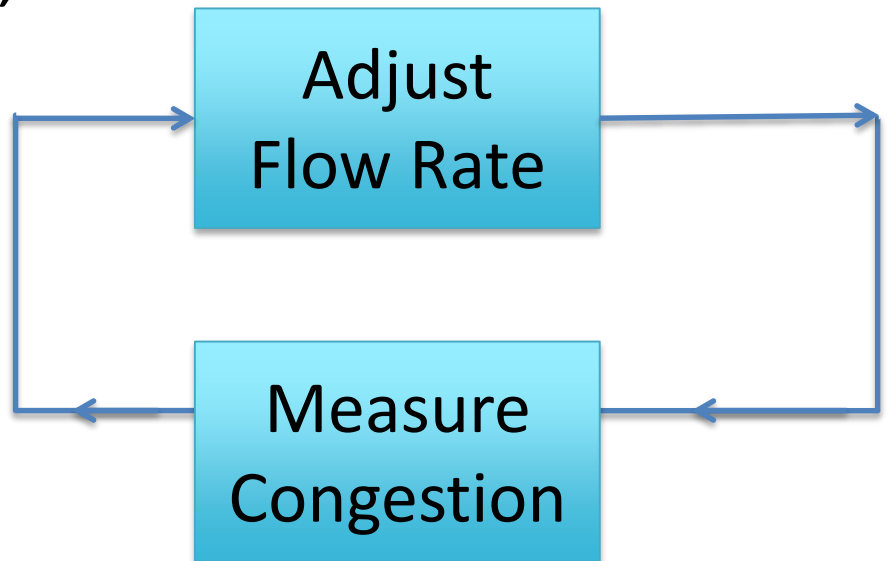
# Ask an oracle.

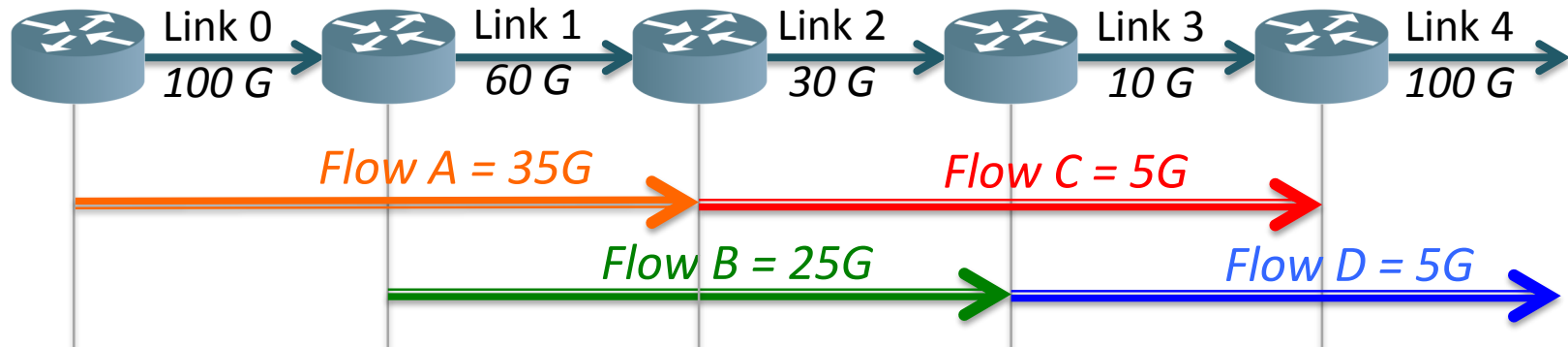
	Link 0	Link	Capacity	Flow	Rate
Flow A	√	0	100	Flow A	35
Flow B		1	60	Flow B	25
Flow C		2	30	Flow C	5
Flow D		3	10	Flow D	5
		4	100		

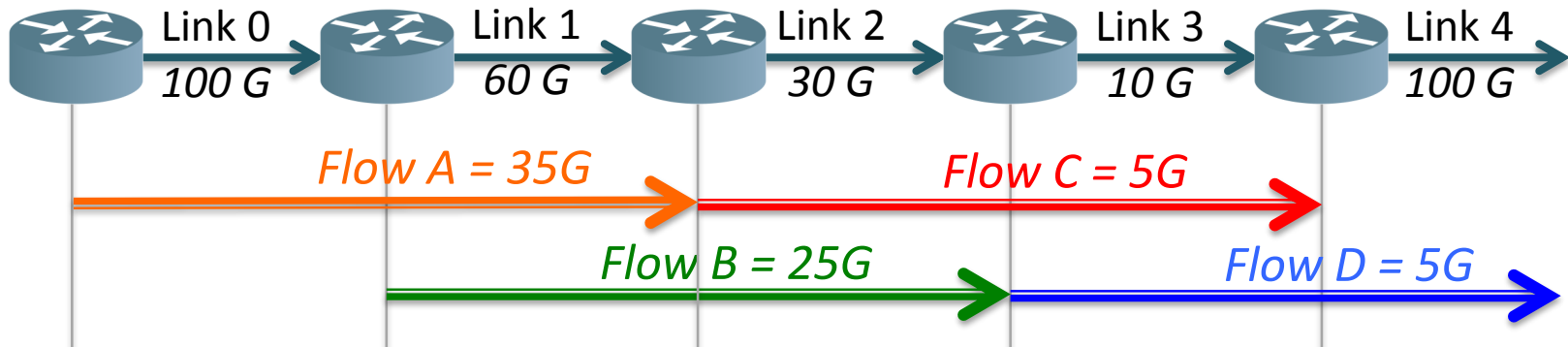
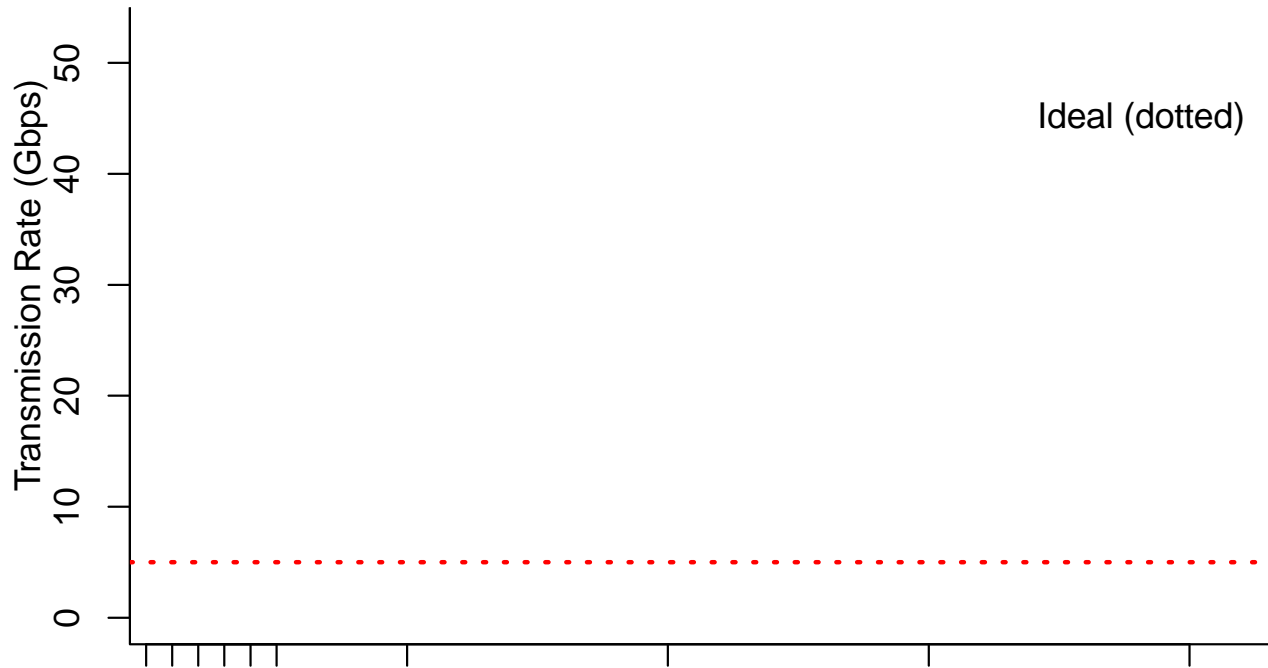


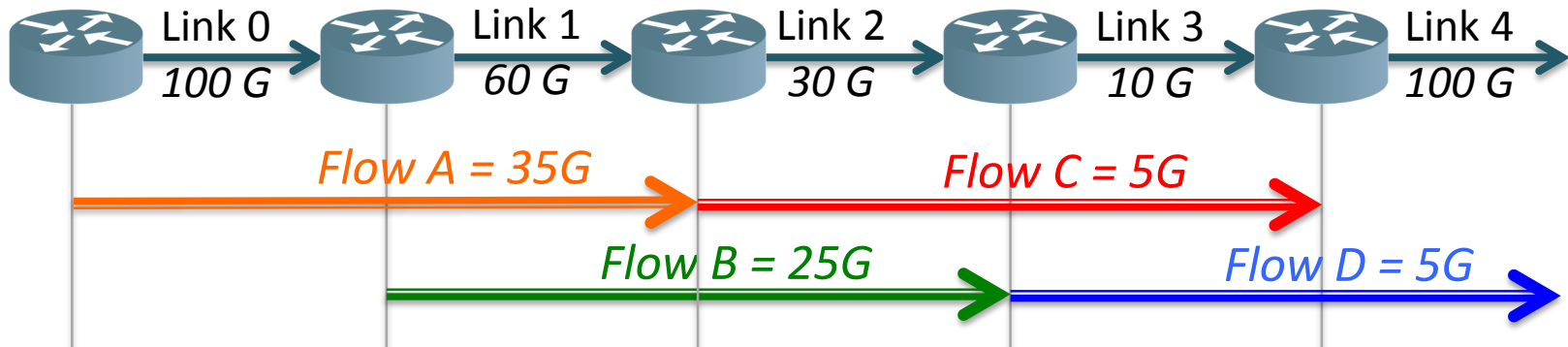
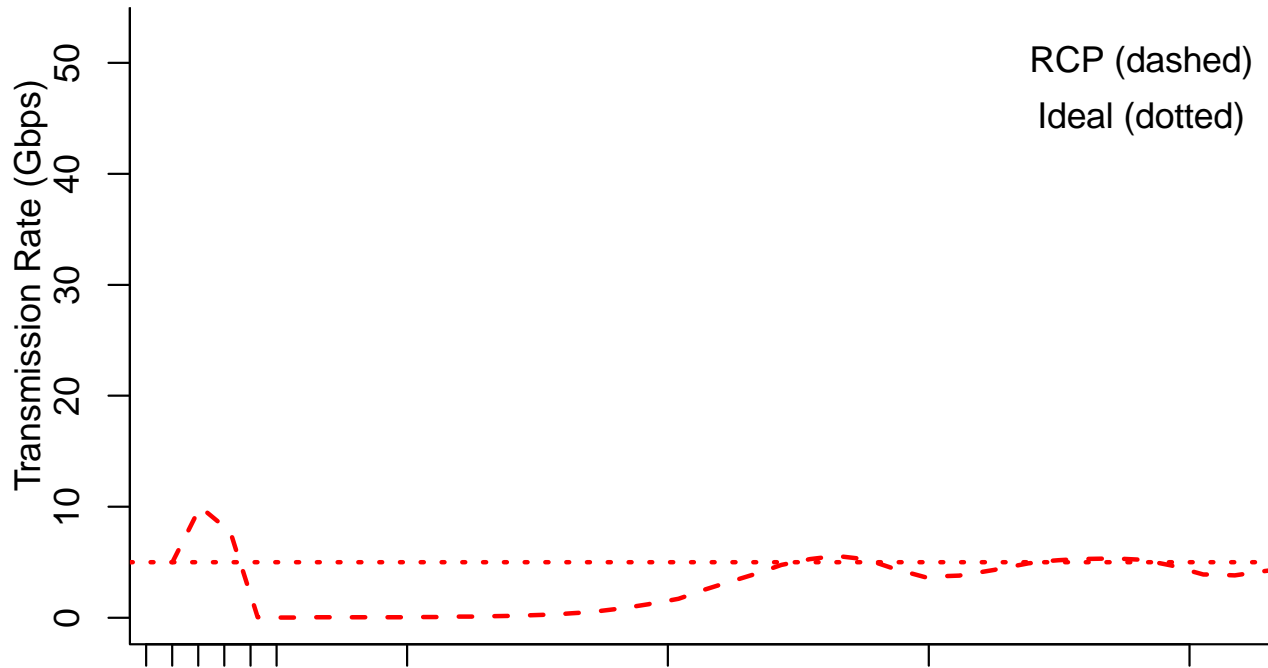
# Traditional congestion control

- No explicit information about traffic matrix
- Measure congestion signals, then react by adjusting rate after measurement delay
- Gradual, can't jump to right rates, know direction
- “Reactive Algorithms”

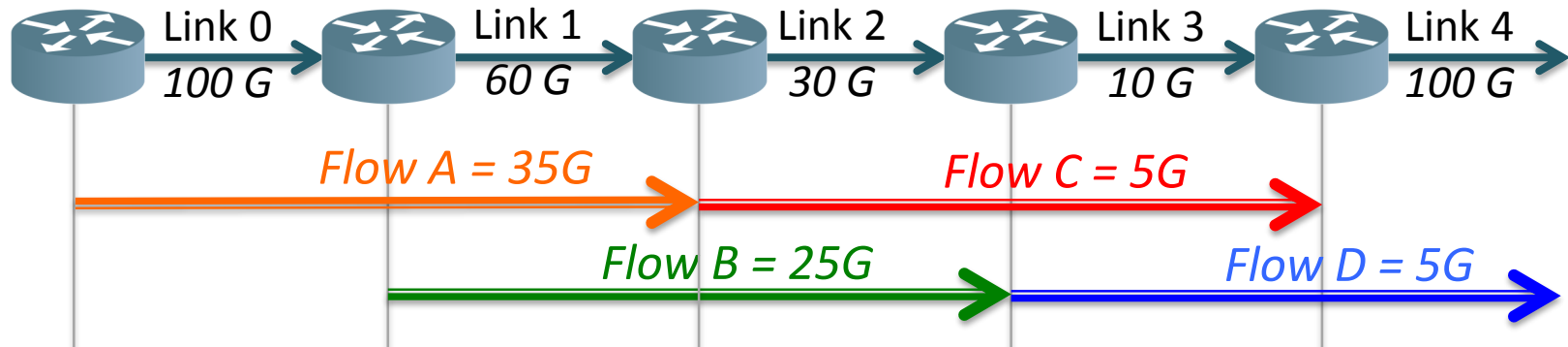
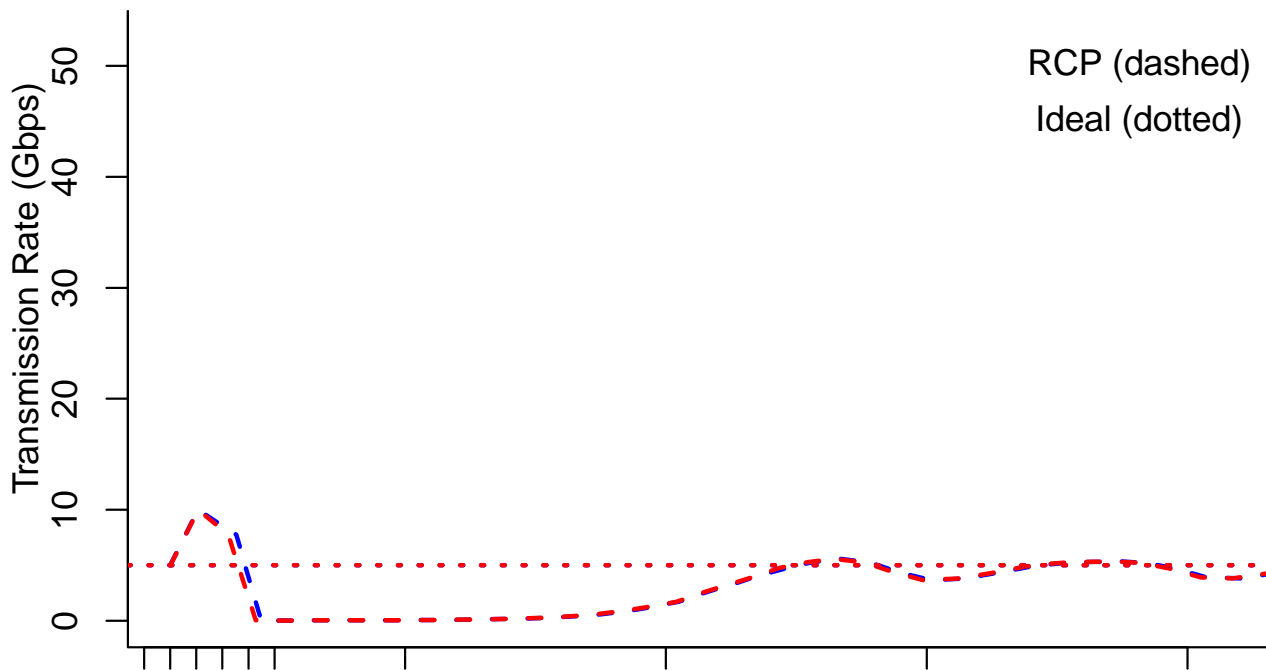


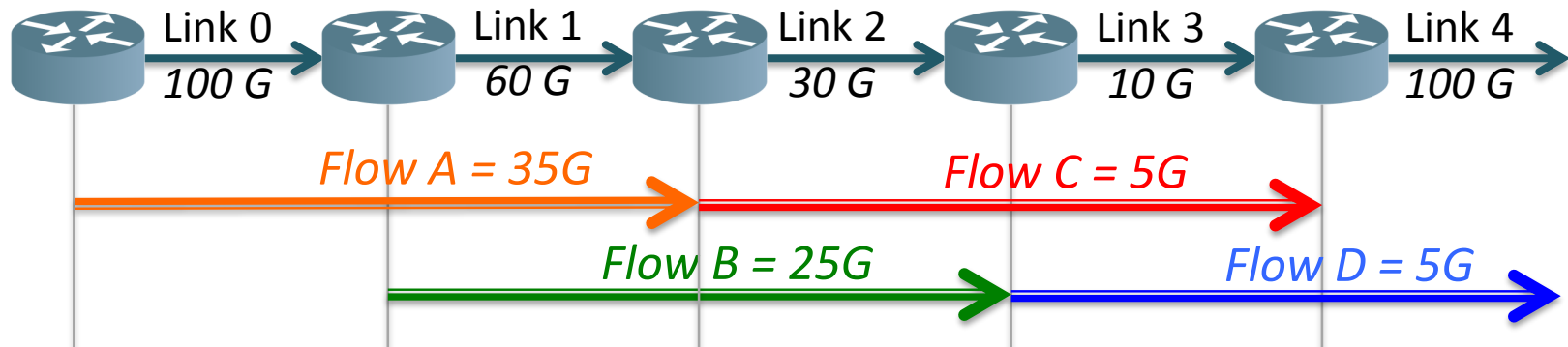
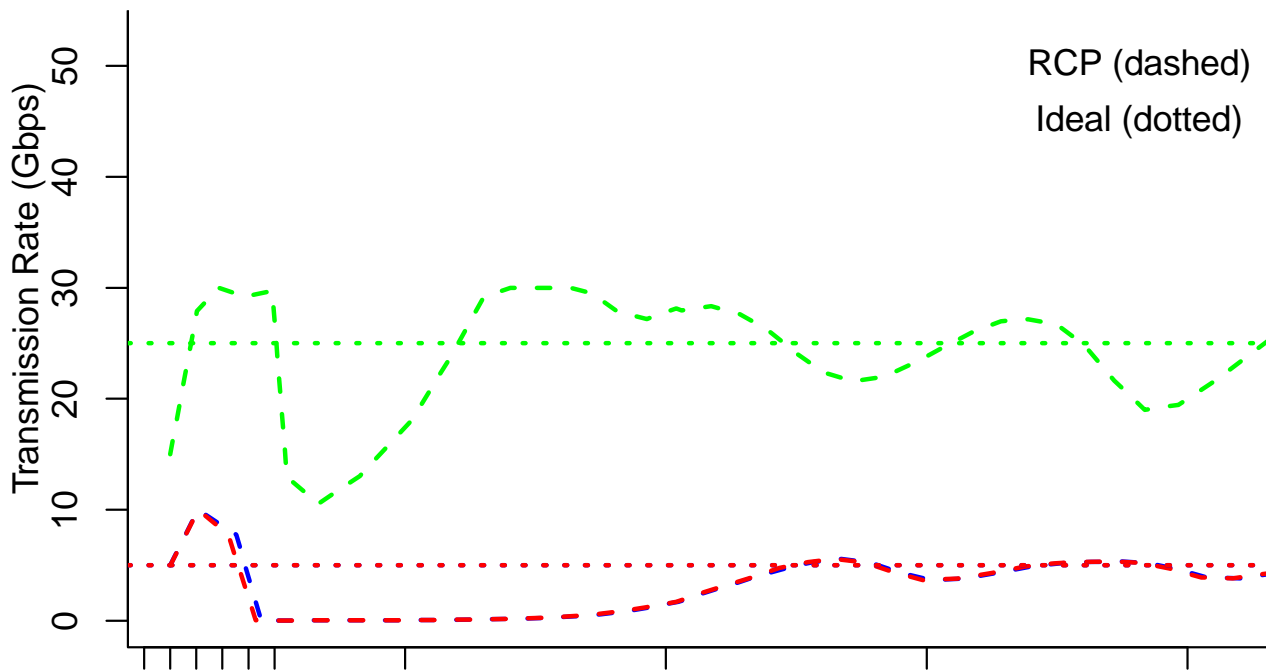


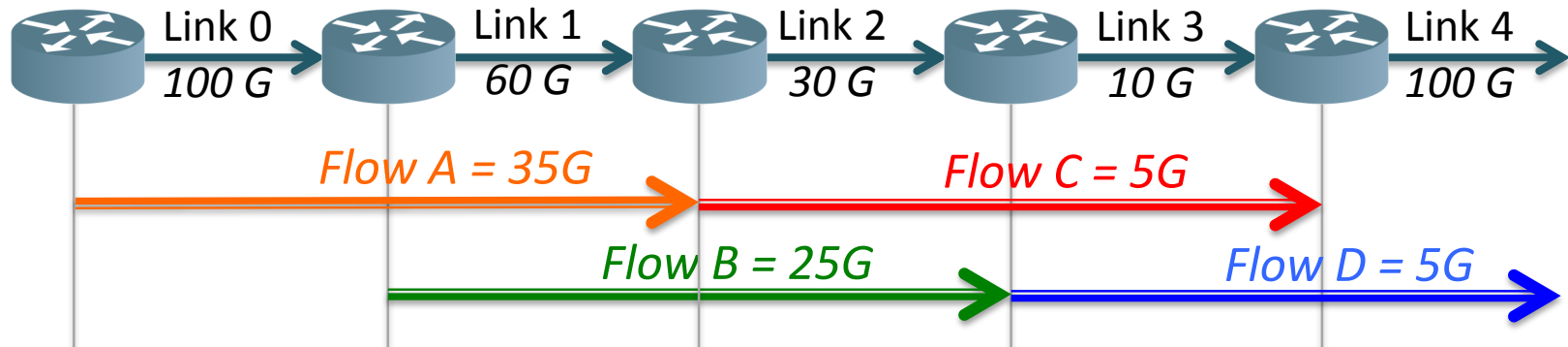
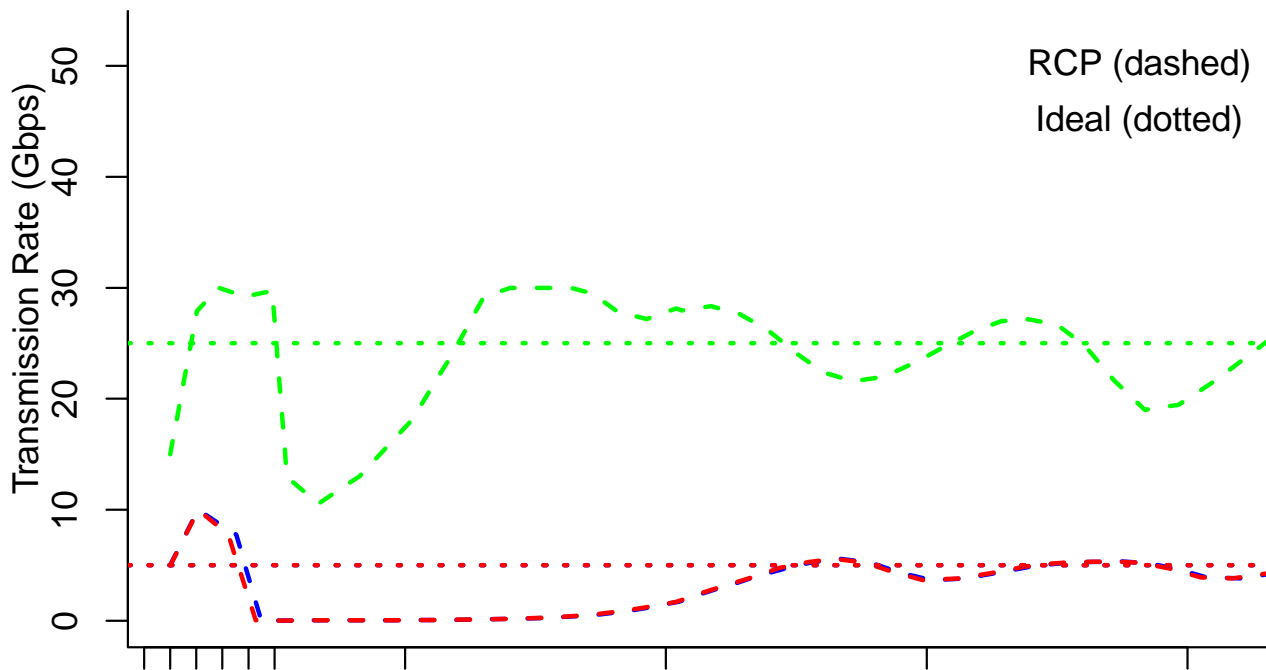


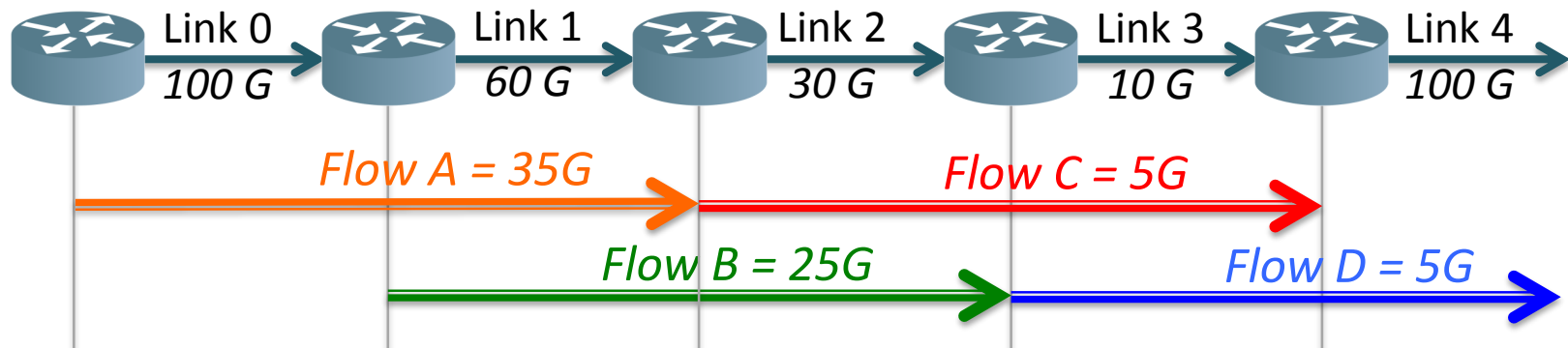
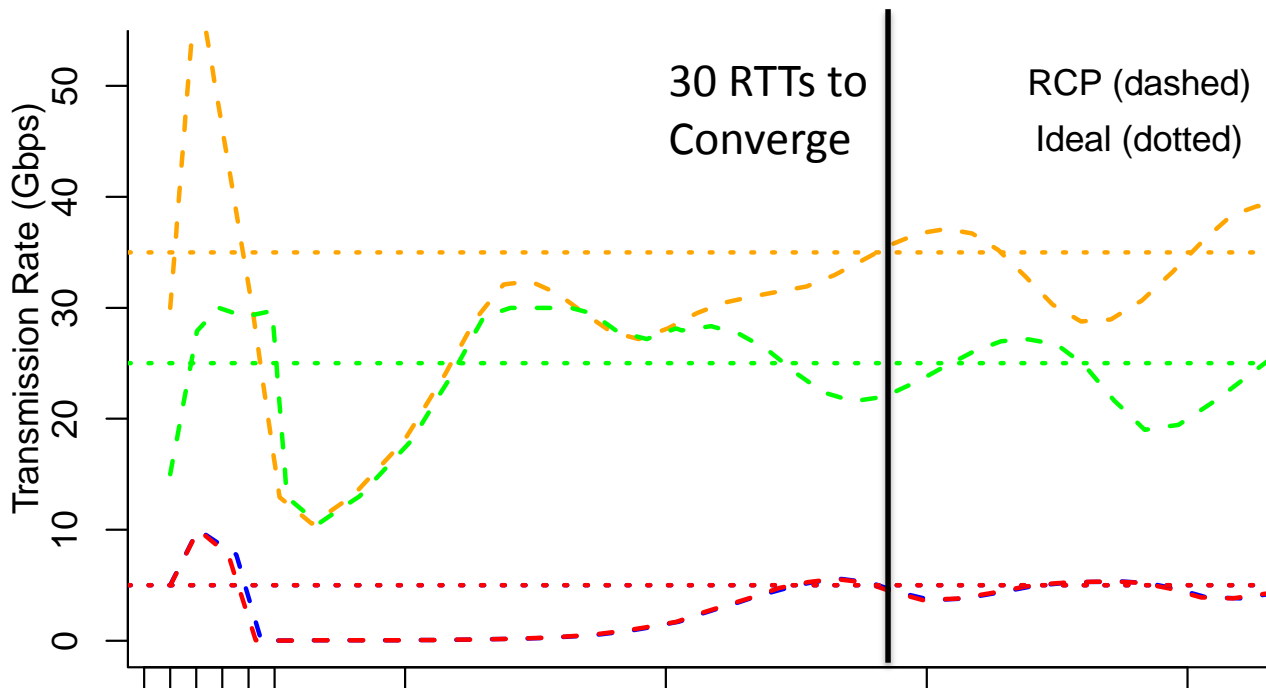








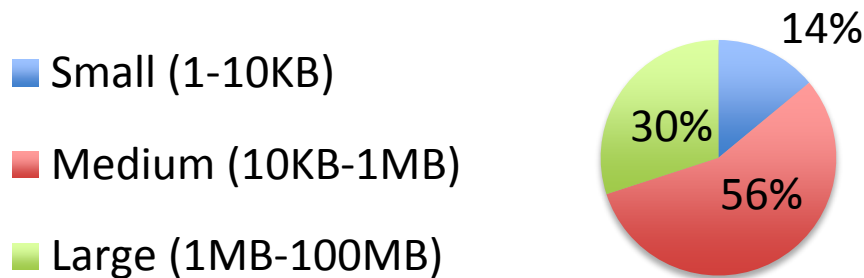




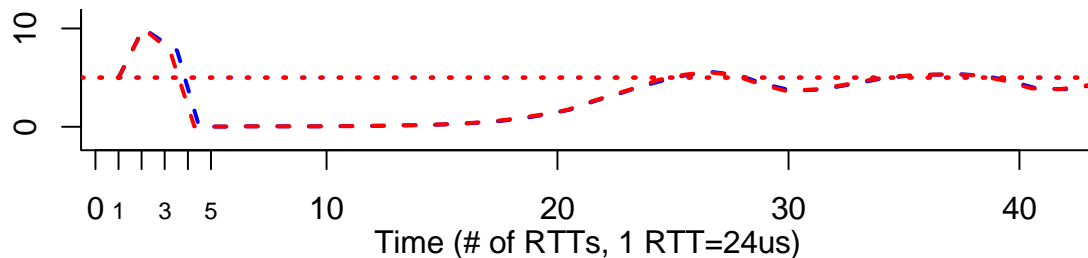
# Relative response times are long

At 100G, a typical flow in a search workload is < 7 RTTs long.

Fraction of Total Flows in Bing Workload



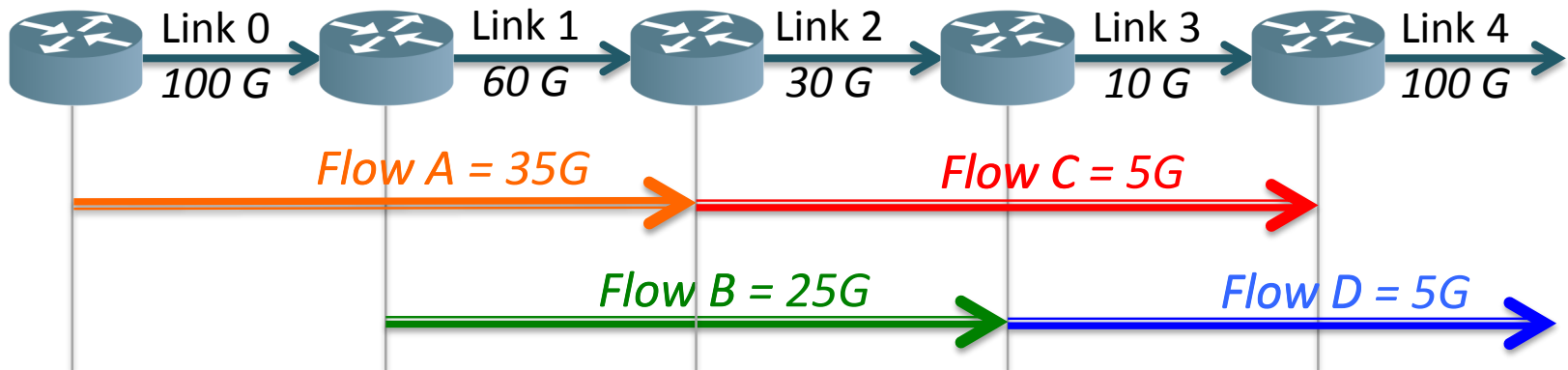
$$1\text{MB} / 100 \text{ Gb/s} = 80 \mu\text{s}$$



Reactive algorithms trade off  
*explicit flow information* for  
*long convergence times*

- Can we use explicit flow information
- and get shorter convergence times?

# Back to the oracle, how did she use traffic matrix to compute rates?



# Waterfilling Algorithm



Link 0 (0/ 100 G)

Link 1 (0/ 60 G)

Link 2 (0/ 30 G)

Link 3 (0/ 10 G)

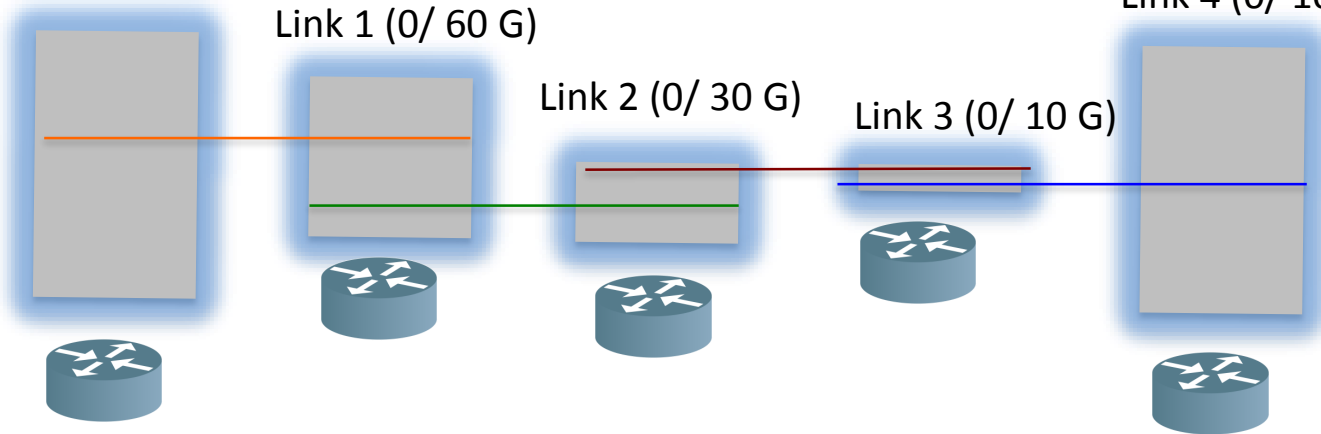
Link 4 (0/ 100 G)

*Flow A (0 G)*

*Flow B (0 G)*

*Flow C (0 G)*

*Flow D (0 G)*





# Waterfilling- 10 G link is fully used



Link 0 (5/ 100 G)

Link 1 (10/ 60 G)

Link 2 (10/ 30 G)

Link 3 (10/ 10 G)

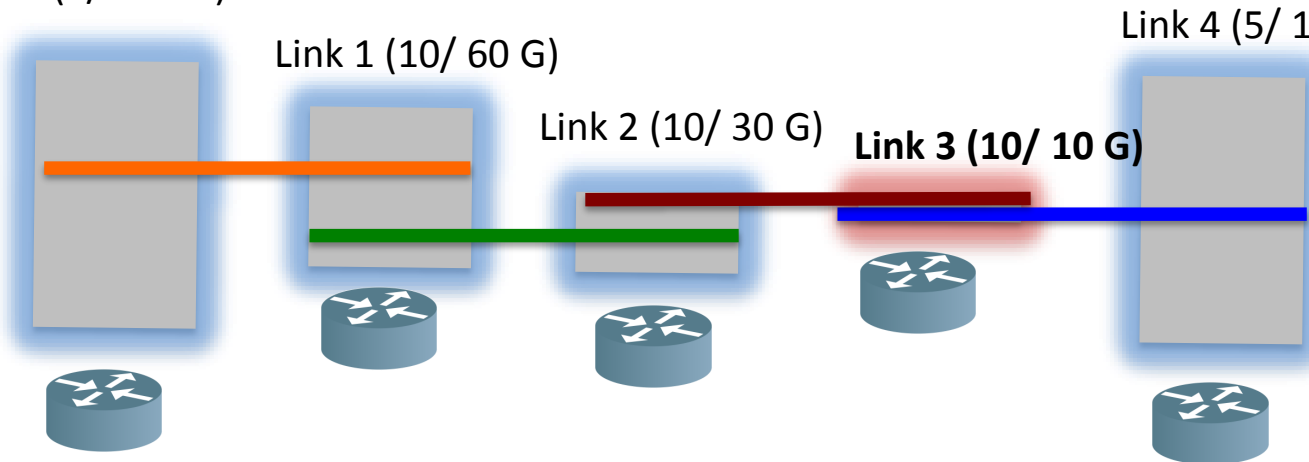
Link 4 (5/ 100 G)

*Flow A (5 G)*

*Flow B (5 G)*

*Flow C (5 G)*

*Flow D (5 G)*



# Waterfilling- 30 G link is fully used



Link 0 (25/ 100 G)

Link 1 (50/ 60 G)

Link 2 (30/ 30 G)

Link 3 (10/ 10 G)

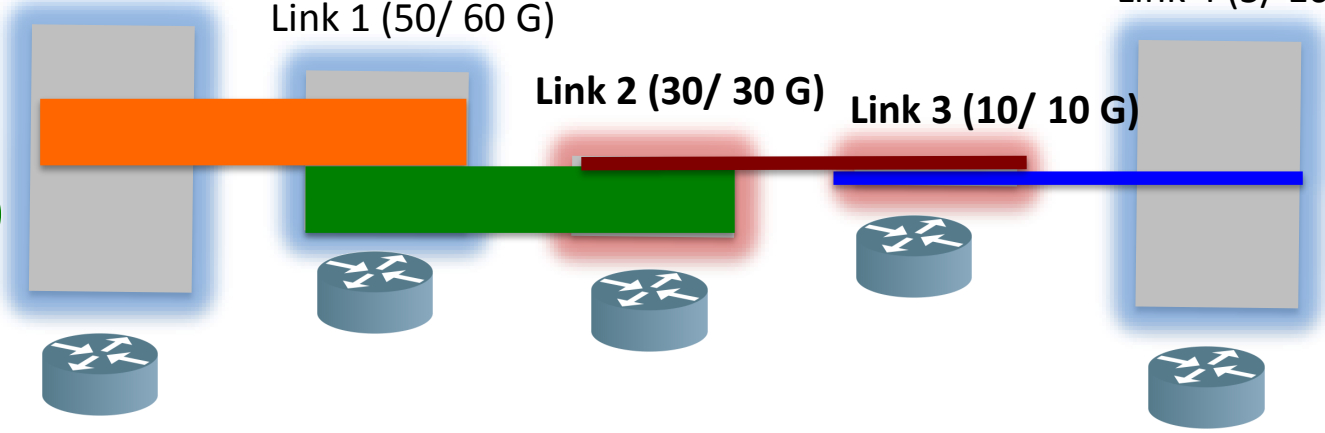
Link 4 (5/ 100 G)

Flow A (25 G)

Flow B (25 G)

Flow C (5 G)

Flow D (5 G)



# Waterfilling- 60 G link is fully used



Link 0 (35/ 100 G)

Link 1 (60/ 60 G)

Link 2 (30/ 30 G)

Link 3 (10/ 10 G)

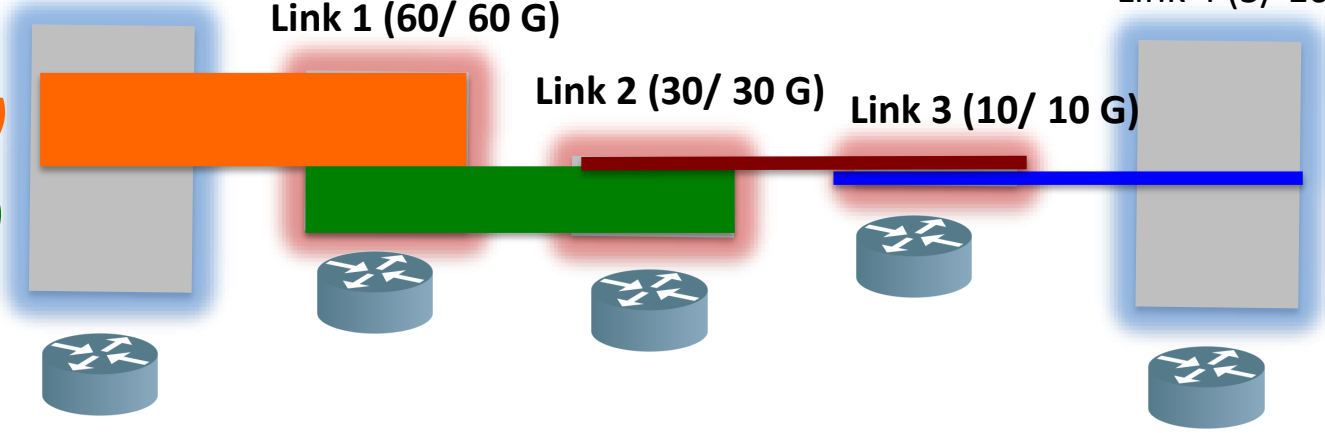
Link 4 (5/ 100 G)

*Flow A (35 G)*

*Flow B (25 G)*

*Flow C (5 G)*

*Flow D (5 G)*



# Fair Share of Bottlenecked Links



Link 0 (35/ 100 G)

Fair Share: 35 G

Link 1 (60 G)

Fair Share: 25 G

Link 2 (30 G)

Fair Share: 5 G

Link 3 (10 G)

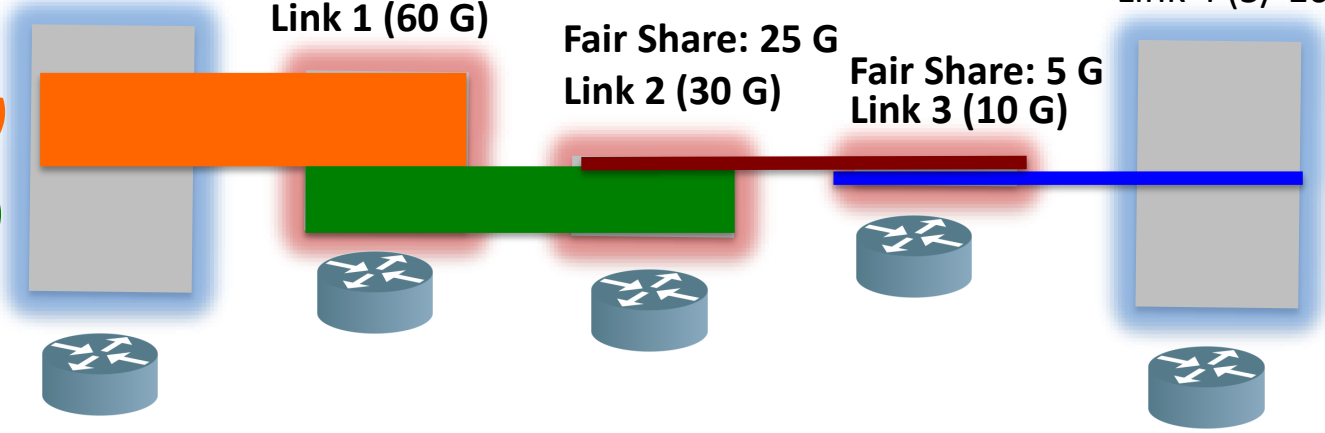
Link 4 (5/ 100 G)

*Flow A (35 G)*

*Flow B (25 G)*

*Flow C (5 G)*

*Flow D (5 G)*



*A centralized water-filling scheme  
may not scale.*

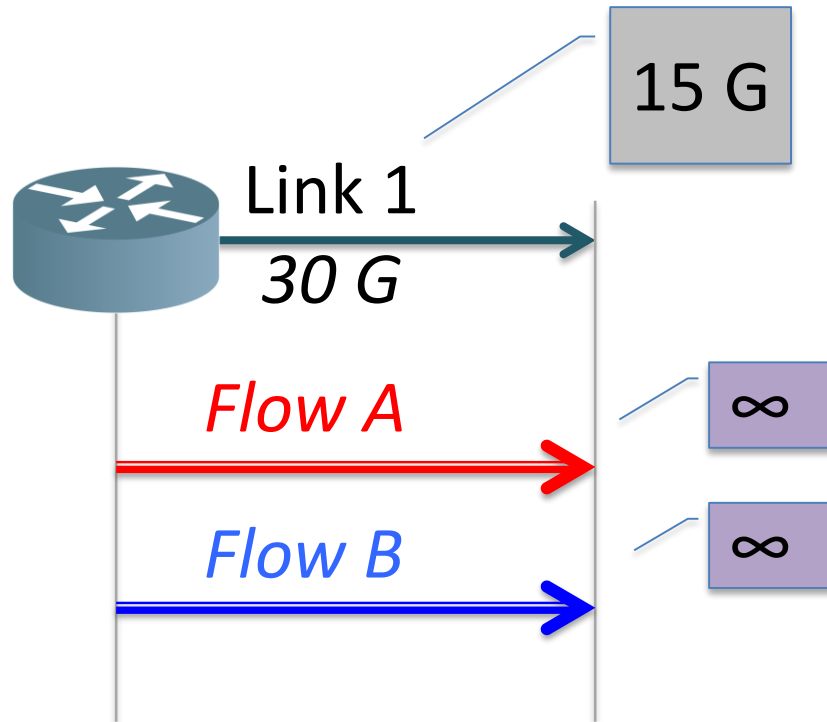
Can we let the network figure out  
rates in a distributed fashion?

# Fair Share for a Single Link

flow	demand
A	$\infty$
B	$\infty$

Capacity at Link 1: 30G

So Fair Share Rate:  $30G/2 = 15G$



# A second link introduces a *dependency*

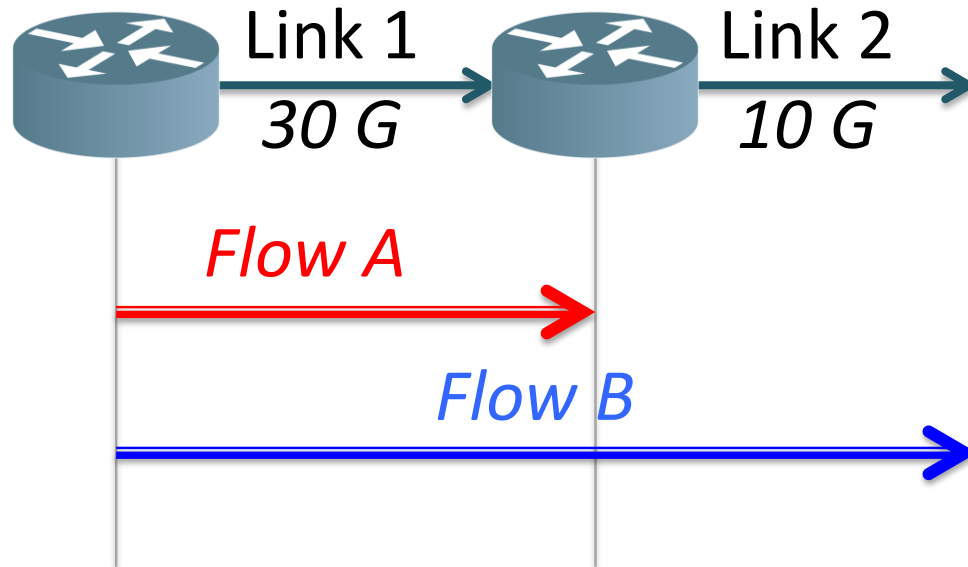
flow	demand
A	$\infty$
B	<del>10 G</del> <b>10 G</b>

Capacity at Link 1: 30G

Demand of Flows restricted at other links: 10G

Number of unrestricted flows: 1

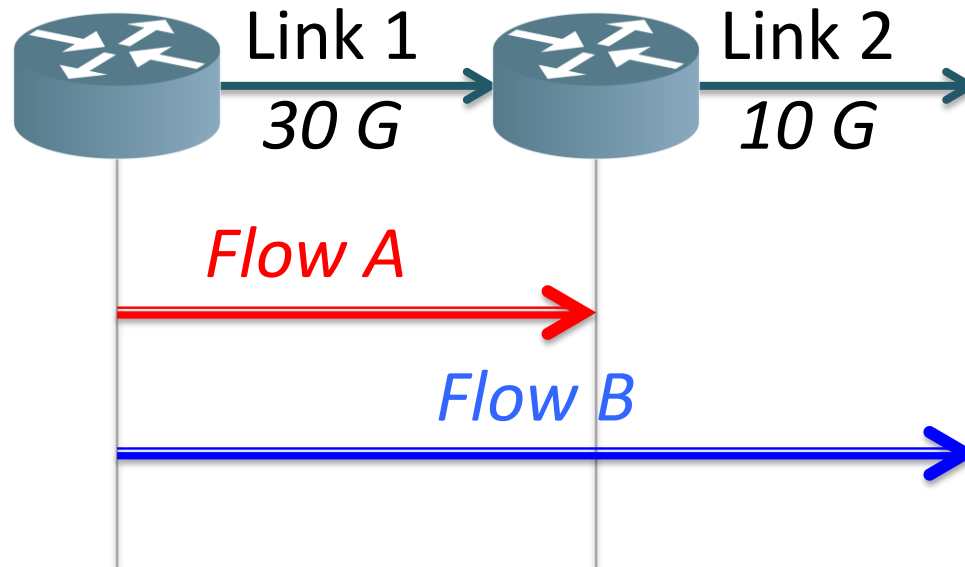
So Fair Share Rate:  $30G - 10G / 1 = 20G$



# Proactive Explicit Rate Control (PERC)

Control Packet For Flow B

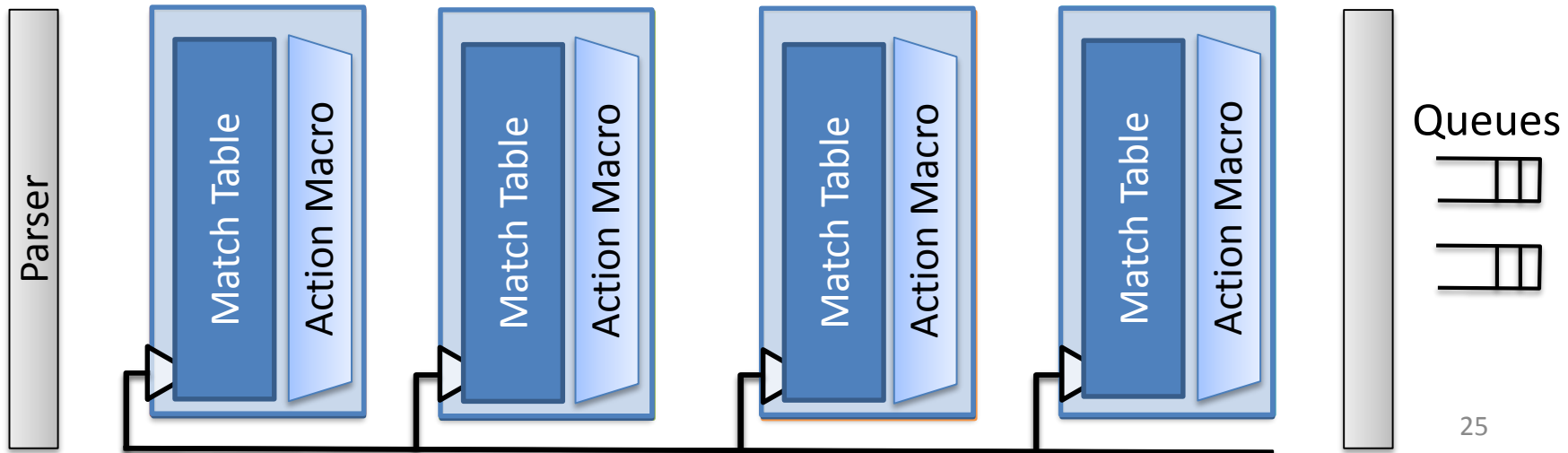
d		$\infty$		$\infty$
f		?		?





# Constraints of Programmable Forwarding Planes at 100 Gb/s

- Limited compute- action  $\sim$  ns, typically primitives like add/ compare etc.
- Limited info. that we can modify per packet.
- Limited area for state and look-up tables  $\sim$  MB, much of which is for L2/L3



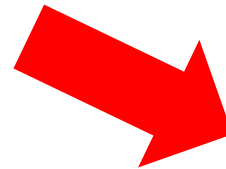
# PERC in P4 $\rightarrow$ NetFPGA



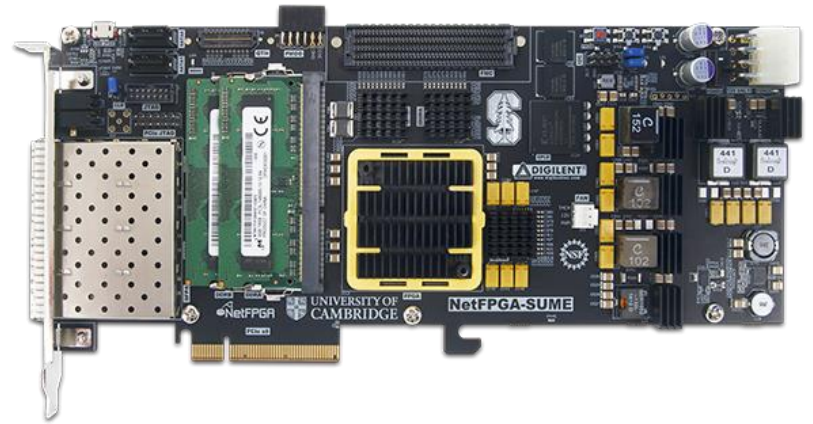
P4 Front end



Xilinx SDNet  
Compilation



NetFPGA SUME Switch



# Division of compute b/n end host & switch

flow	demand
A	$\infty$
B	<b>10 G</b>

Capacity at Link 1: 30G

Demand of Flows restricted at other links: 10G

Number of unrestricted flows: 1

Start Fair Share Rate 30G, 10G, 1  
30G, 10G, 1



# Interesting Questions

- Minimum time for a distributed scheme
- Minimum amount of state for provable convergence
- How many active flows in a max-min fair network?
- Imprecise demands  $\rightarrow$  some reactive component