

# Flowtune

Jonathan Perry

Joint work with Hari Balakrishnan and Devavrat Shah.



# Flowtune is..

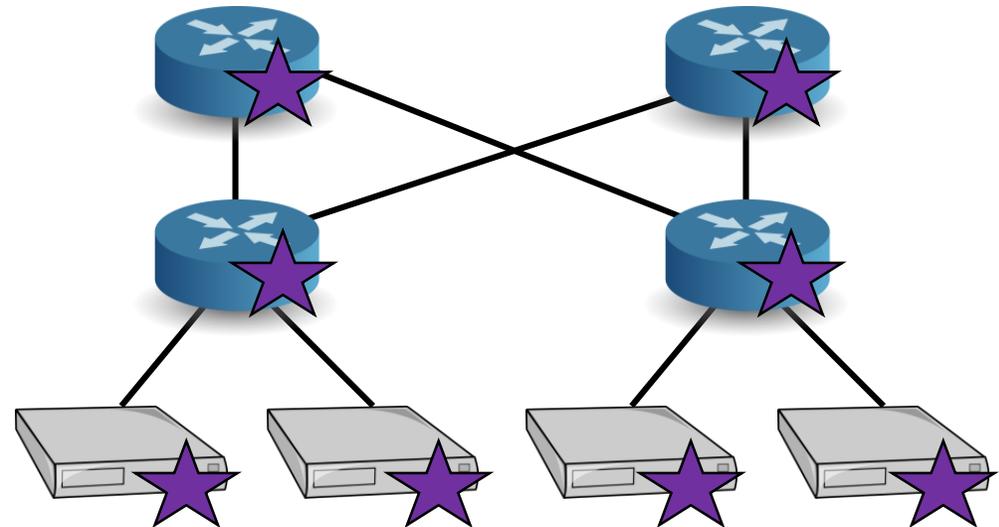
## Allocate network resources

- Quickly
- Explicitly (maximize utility)
- Flexibly (in software)

# Traditional approach is packet-centric

Switch Algorithms

Server Algorithms



Implicit  
Allocation

Several RTT  
to converge

Changes many  
components

# Flowtune's approach

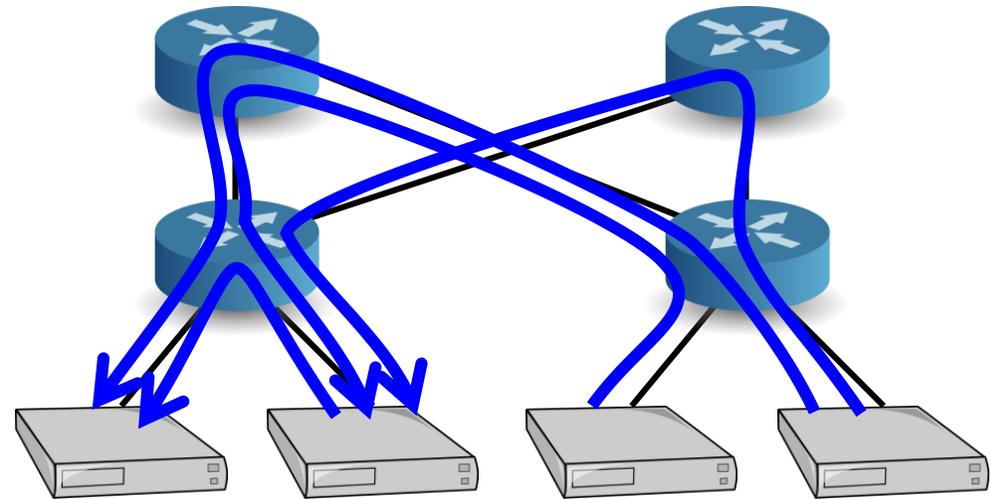
## 1. Flowlet control

Allocation changes *only* when:

- Flowlets arrive
- Flowlets terminate

## 2. Logically centralized

- Reduce RTT dependence



# Example

A → Allocator

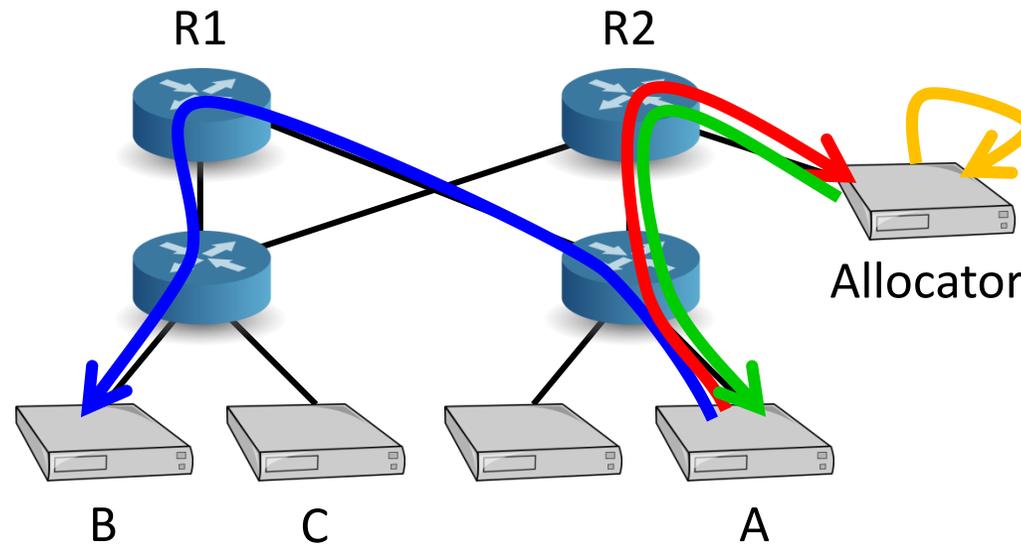
Allocator

Allocator → A

“Hadoop on A has data for B”

Assign rates

“Send at 10Gbps”



# Example

C → Allocator

Allocator

Allocator → A

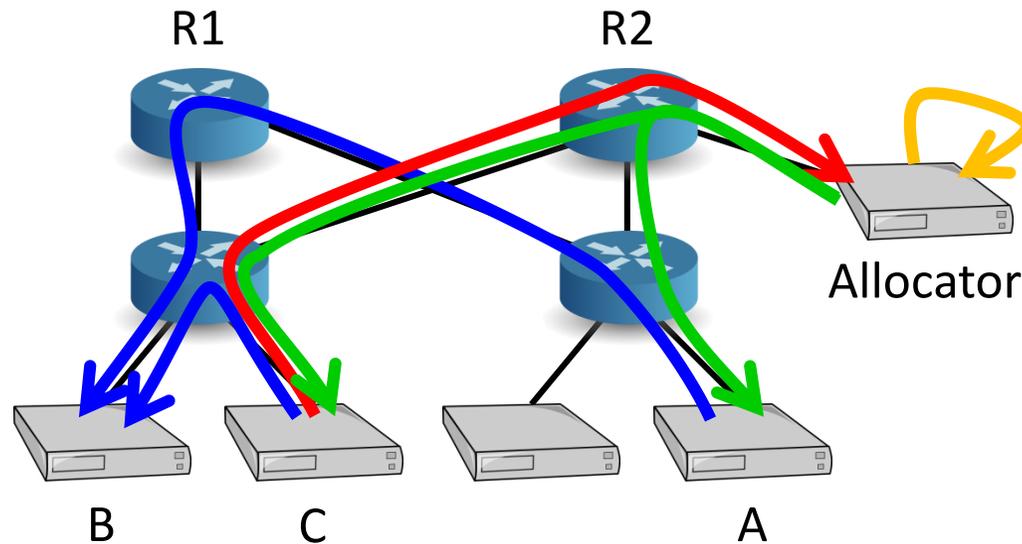
Allocator → C

“ads\_update on C has data for B”

Assign rates

“Send at 1Gbps”

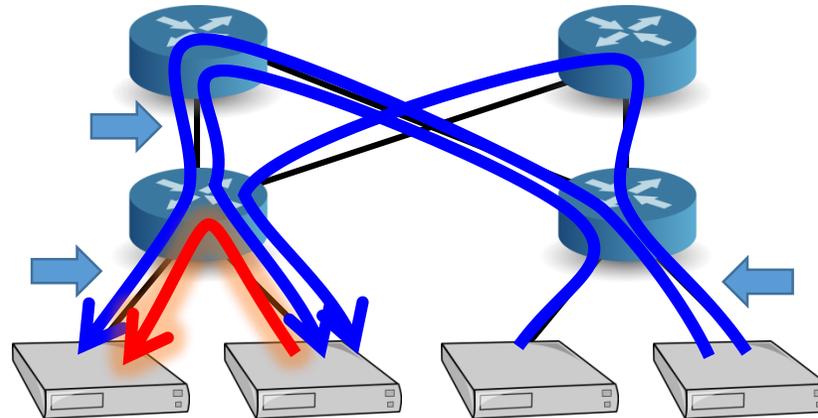
“Send at 9Gbps”



# Why is this hard?

## Need to choose rates given active flowlets

1. Updates cascade!
2. What is the goal? To act like TCP?



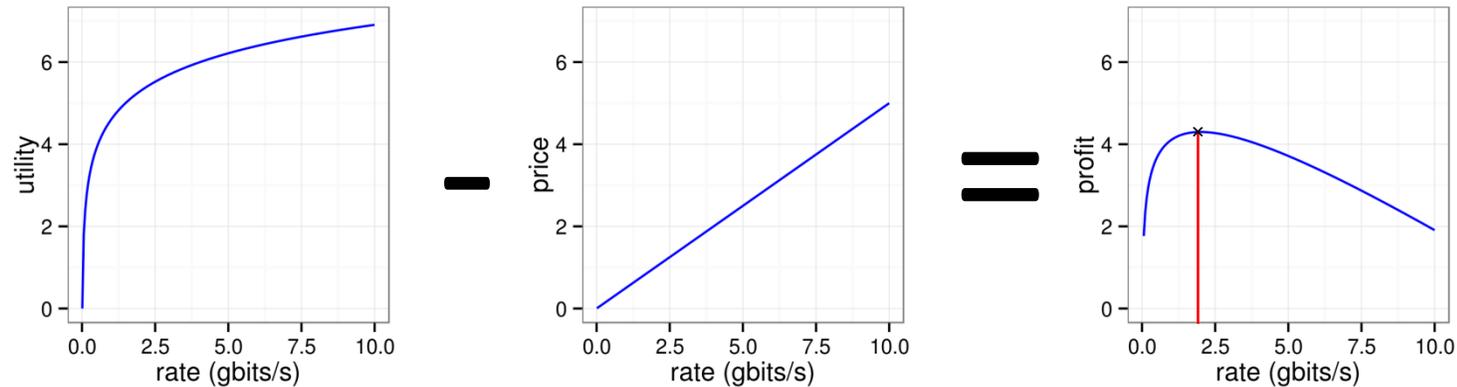
# NUM Iterative Optimizer

1. Each link  $\ell$  chooses price  $p_\ell$

$$\sum_{s \in \mathcal{S}(\ell)} x_s - c_\ell$$

Demand Supply

2. Each flow  $s$  chooses rate  $x_s$



3. Goto 1

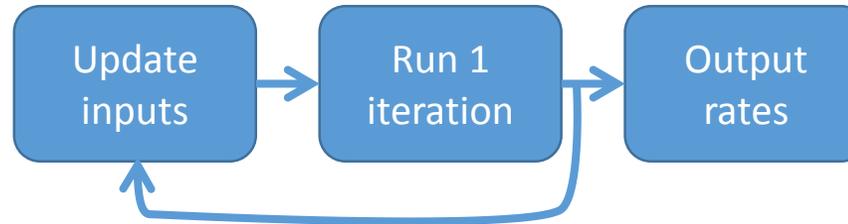
# How to reduce latency?

Solution 1:

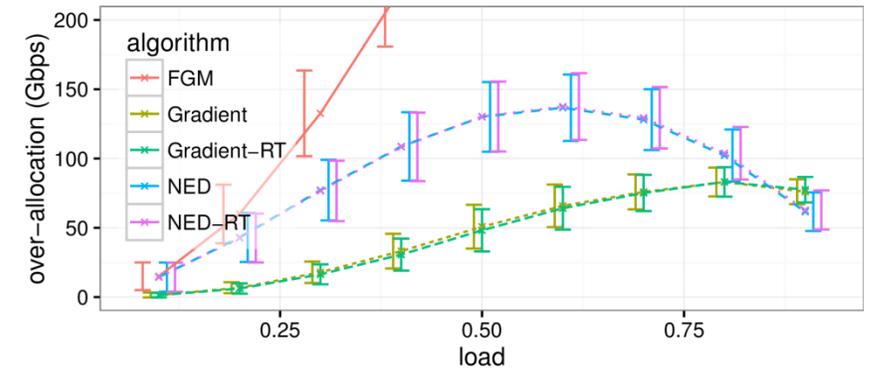


But: too slow!

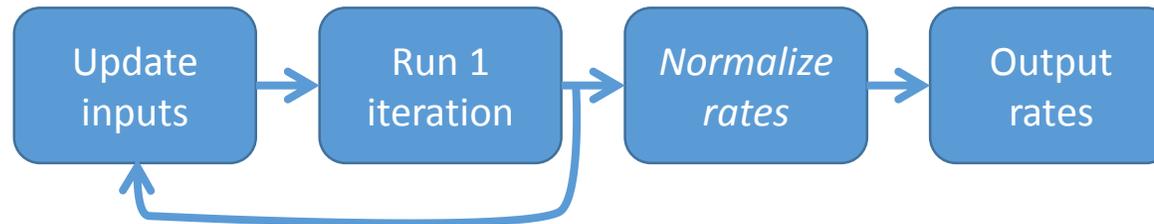
Solution 2:



But: links are over-allocated!



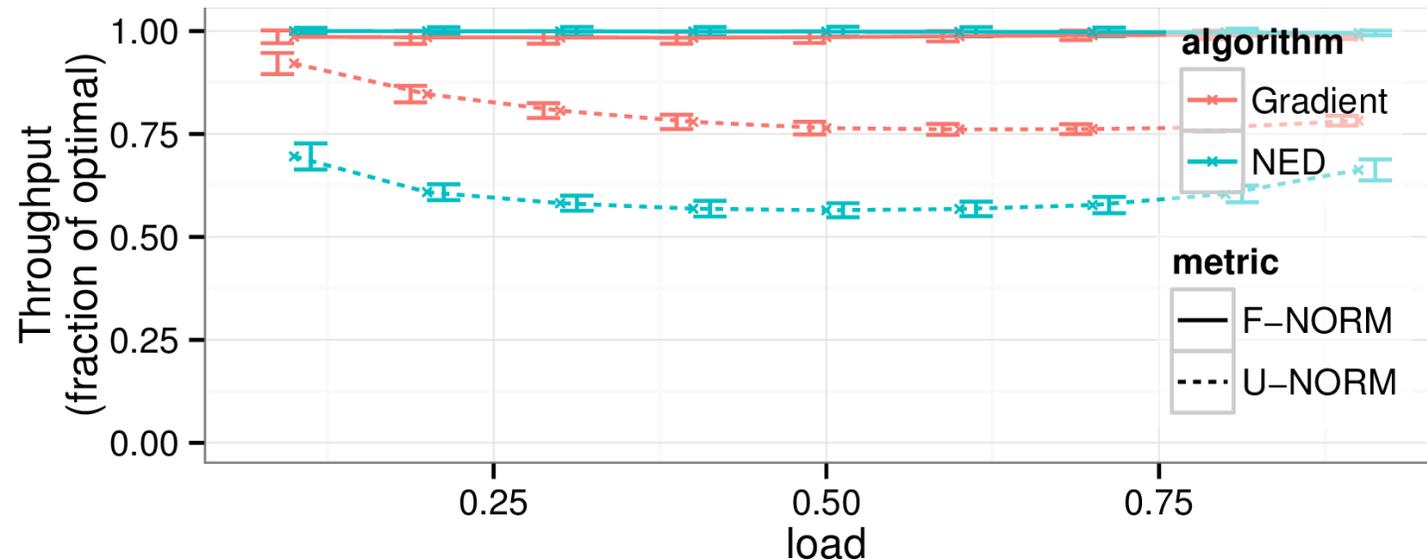
Solution 3:



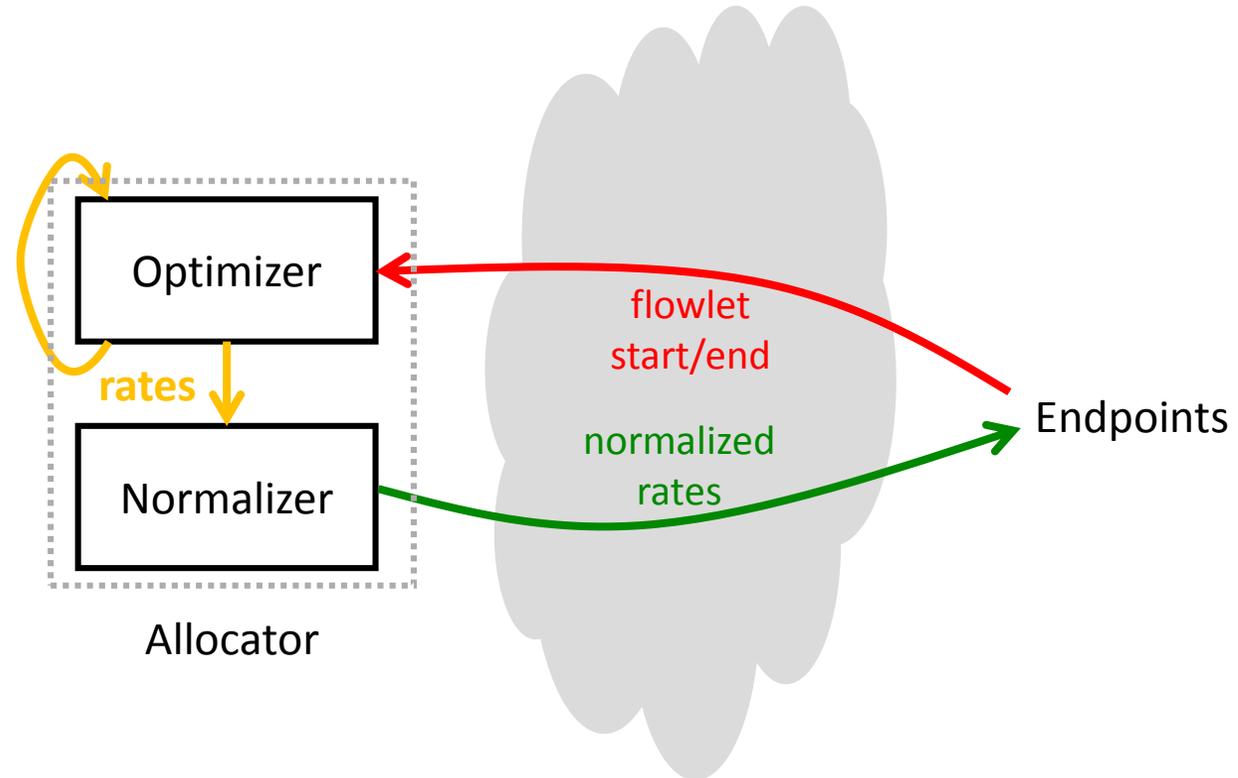
# Flowtune normalizes rates between iterations

- For each flow:

- Find link  $\ell$  on path with largest  $r_\ell = \frac{\sum \text{flow rates}}{\text{link capacity}}$
- Normalize:  $x_s \leftarrow x_s / r_\ell$

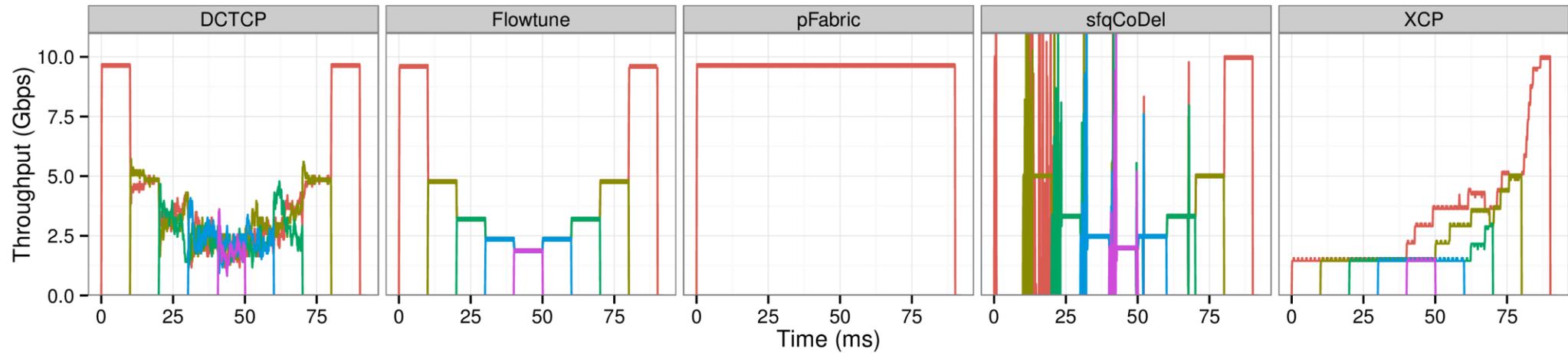


# Architecture



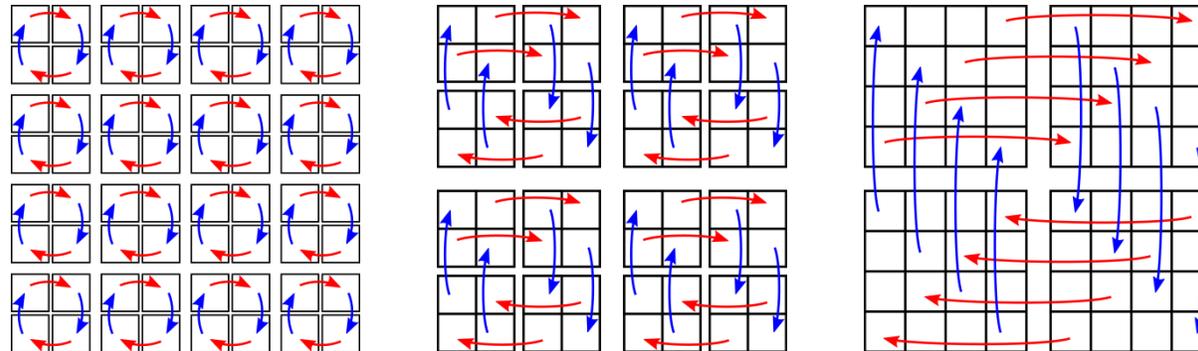
# Flowtune converges quickly to fair allocation

1. Every 10 milliseconds add sender, up to 5 senders
2. Every 10 milliseconds remove sender



# Flowtune scales to 49K flows on 64 cores

Cores	Nodes	Flows	Cycles	Time
4	384	3072	19896.6	8.29 $\mu$ s
16	768	6144	21267.8	8.86 $\mu$ s
64	1536	12288	30317.6	12.63 $\mu$ s
64	1536	24576	33576.2	13.99 $\mu$ s
64	1536	49152	40628.5	16.93 $\mu$ s
64	3072	49152	57035.9	23.76 $\mu$ s
64	4608	49152	73703.2	30.71 $\mu$ s



# Flowtune

## Allocate network resources

- Quickly (centralized)
- Explicitly (maximize utility)
- Flexibly (in software)

