# DATA CENTER TO THE HOME

Koen De Schepper, Inton Tsang          Bell Labs
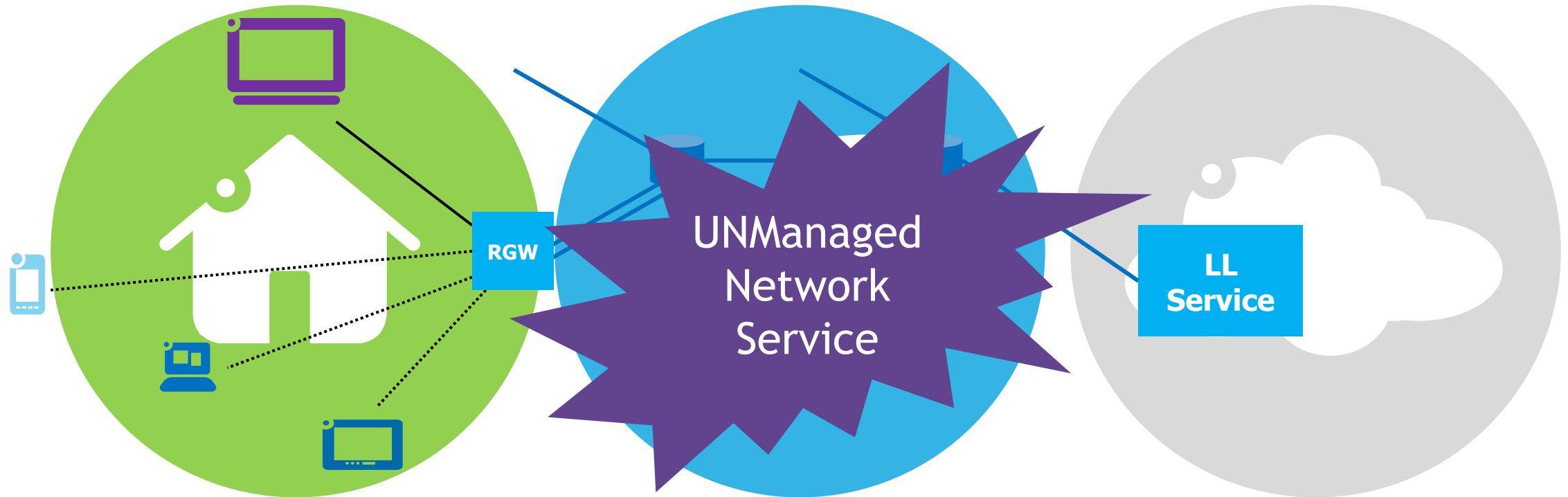
Olga Bondarenko          [ simula . research laboratory ]
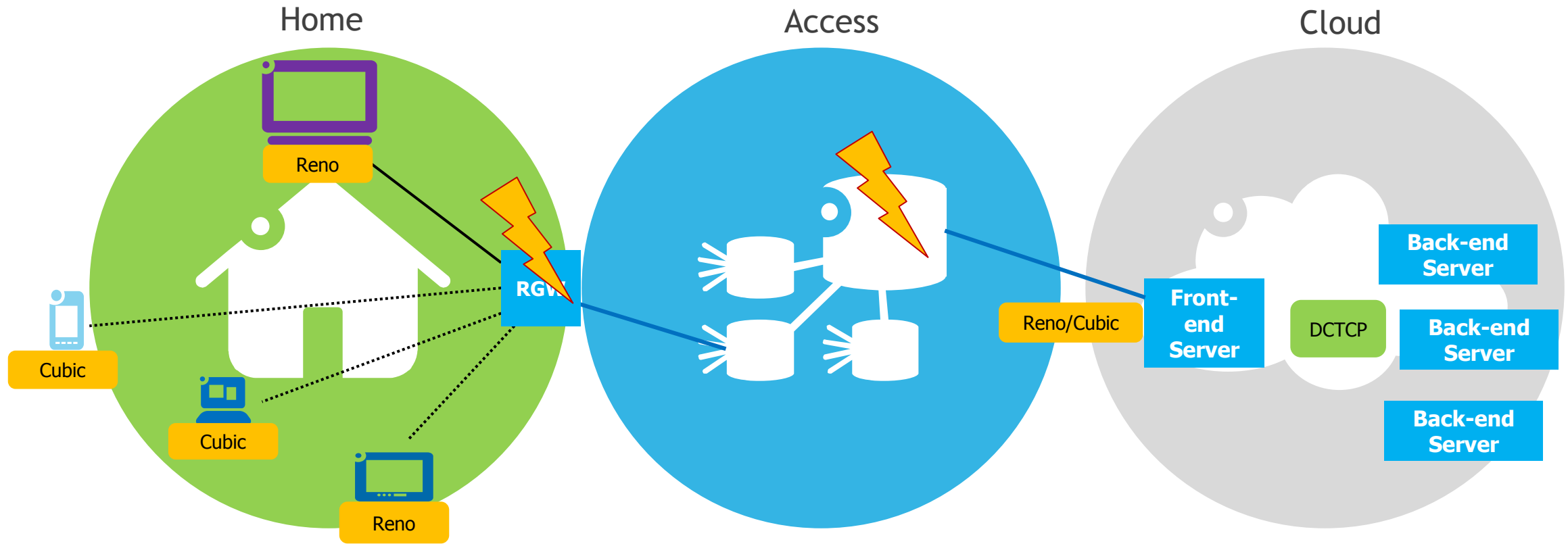
Bob Briscoe          BT

koen.de_schepper@alcatel-lucent.com

March, 2015

R / TE
REDUCING INTERNET TRANSPORT LATENCY

# DCttH OBJECTIVE: UNIVERSAL SUPPORT FOR LOW LATENCY = SUPPORT FOR ADAPTIVE INTERACTIVE APPLICATIONS

RGW

UNManaged Network Service

LL Service

R / TE
REDUCING INTERNET TRANSPORT LATENCY

# INTERACTIVE APPLICATIONS on the INTERNET ?



Home

Access

Cloud

Reno

Cubic

Cubic

Reno

RGW

Reno/Cubic

Front-end Server

DCTCP

Back-end Server

Back-end Server

Back-end Server

**Large queues for high throughput and low drop**
**= Poor Latency**
**= Bad for interactive applications**

**ECN = No drop**
**ECN++ = Small queues**
**= Low latency & High throughput**

R/TE
REDUCING INTERNET TRANSPORT LATENCY

# DATACENTER to the HOME ?



Home

Reno

Cubic

Cubic

Reno

RGW

Access

Cloud

Reno/Cubic

Front-end Server

DCTCP

Back-end Server

Back-end Server

Back-end Server

Windows and Linux 3.18
have DCTCP implementations ready
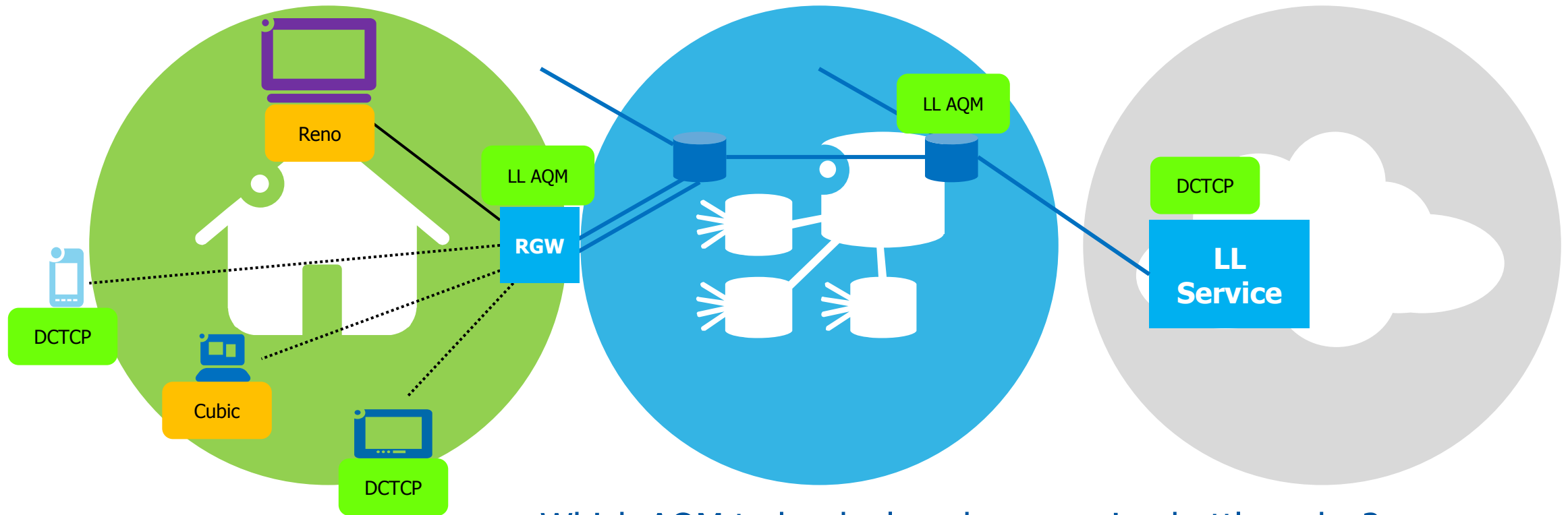
Clients use Reno and Cubic
Can't use DCTCP without causing trouble

Public Internet
does not support DCTCP

DCTCP available on
Windows Server and Linux 3.18
used internally in the data center

R/TE
REDUCING INTERNET TRANSPORT LATENCY

4

# MIGRATION OBJECTIVE: LOW LATENCY ACCESS TO THE CLOUD, EQUAL STEADY STATE THROUGHPUT TO RENO/CUBIC

Can DCTCP be used as Low Latency congestion controller ?



Which AQM to be deployed on queuing bottlenecks ?

Support migration !

# LOWER LATENCY BY SMARTER USE OF ECN
# DATA CENTER TCP

**TCP (Reno)**  ↔  **DCTCP**

### Response to congestion in sender

- Half the congestion window when drop detected in one RTT

- Reduce partially per marked packet; half if all marked in one RTT
  → React according to level of congestion

### ECN feedback in receiver

- Echo Congestion Experienced (CE) until sender acknowledges Congestion Window Reduced (CWR)

- Echo marking state of received packets without acknowledgement → accurate ECN feedback

### ECN marking in network

- Smooth and delay a drop or mark to allow bursts

- Don't smooth or delay queue size
- Shallower marking threshold
  → immediate ECN marking

R / T E
REDUCING INTERNET TRANSPORT LATENCY

# DEMONSTRATED ON A REAL BB RESIDENTIAL TESTBED

# LOWER LATENCY BY SMARTER USE OF ECN
# DATA CENTER TCP

TCP (Reno) ↔ DCTCP

AQM configuration



Average Q size

Instant Q size

Q size variation



Pdf in 250s interval [%]
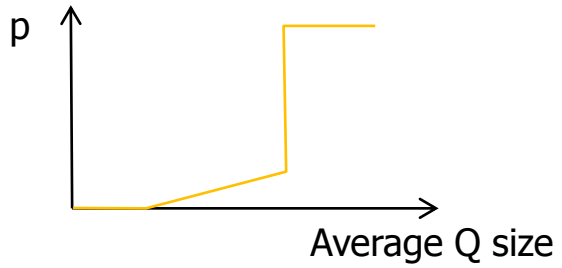
Measured in a BB DSL testbed

RTT = 8 ms (unloaded)

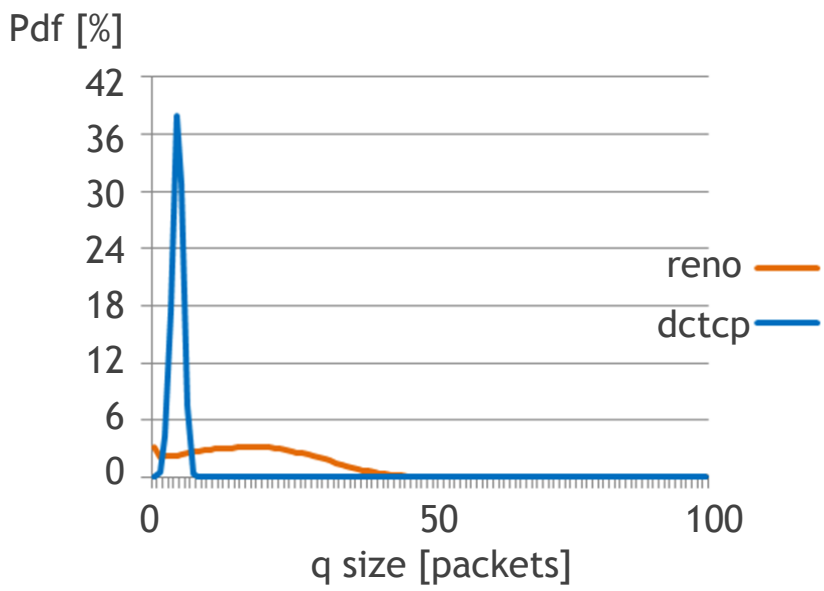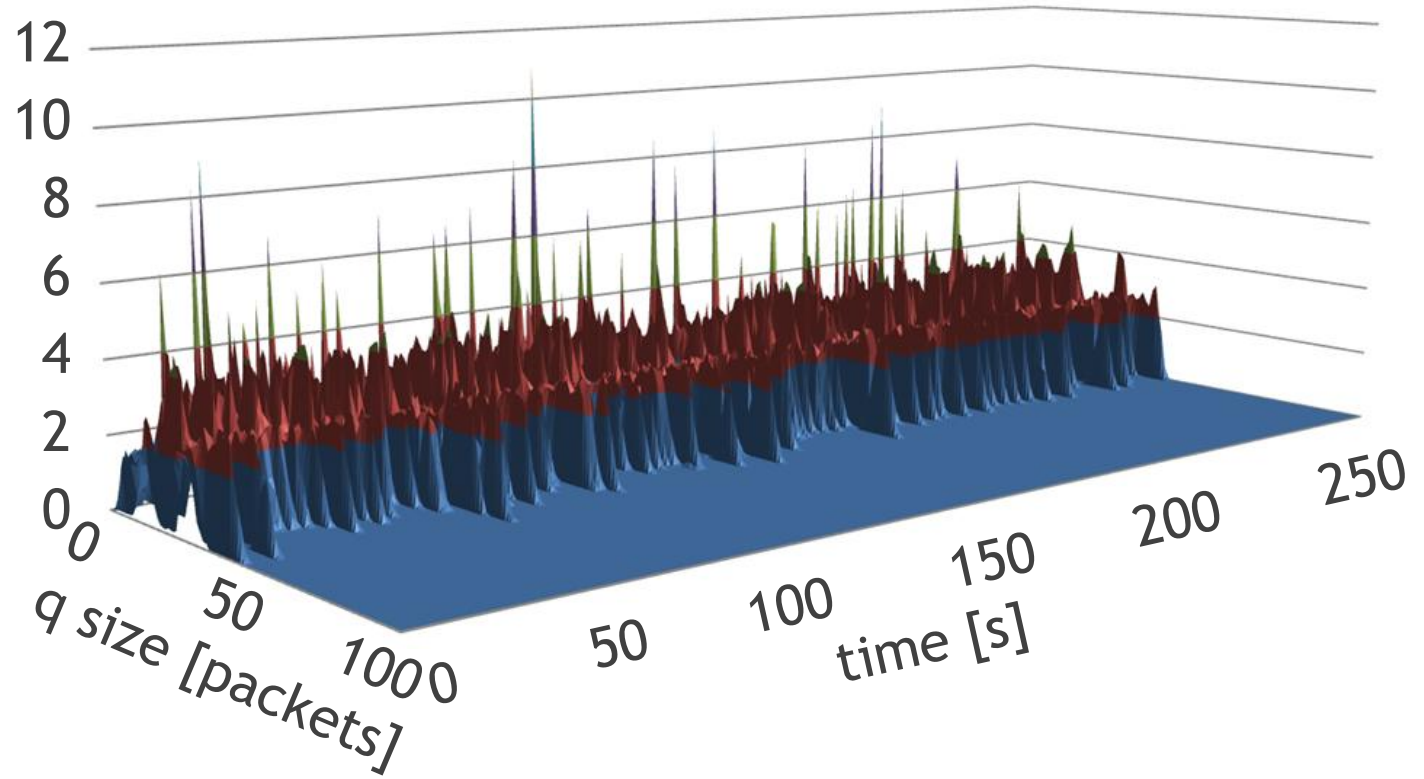BW = 40 Mbps (downstream)

1 steady state flow running alone

Reno/Cubic/DCTCP = Linux kernel 3.18

# QUEUE SIZE AT DEQUEUE
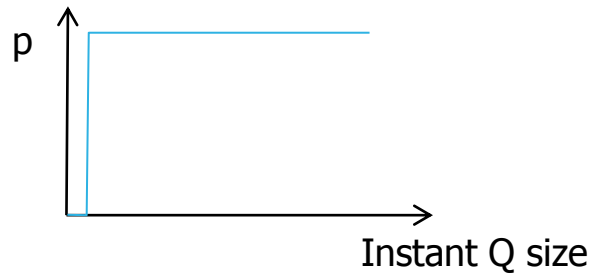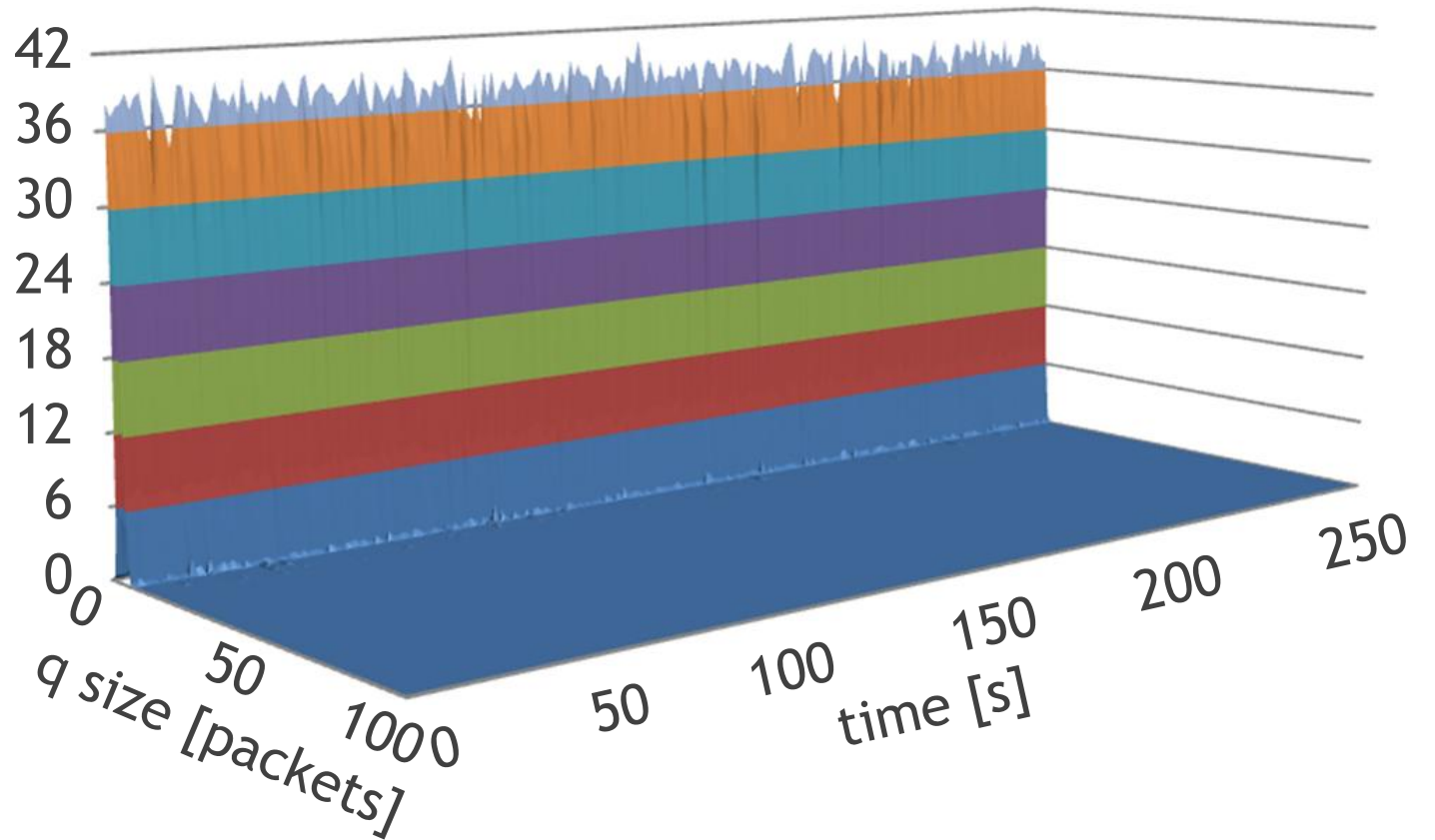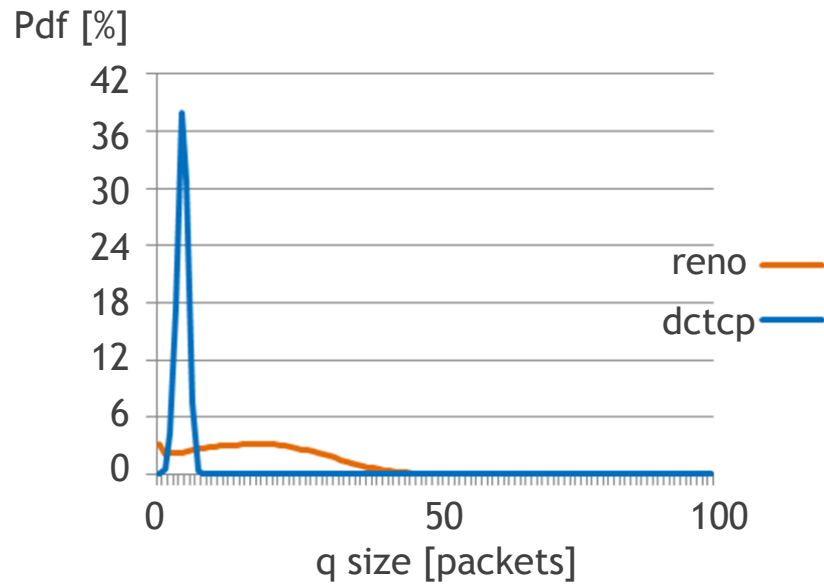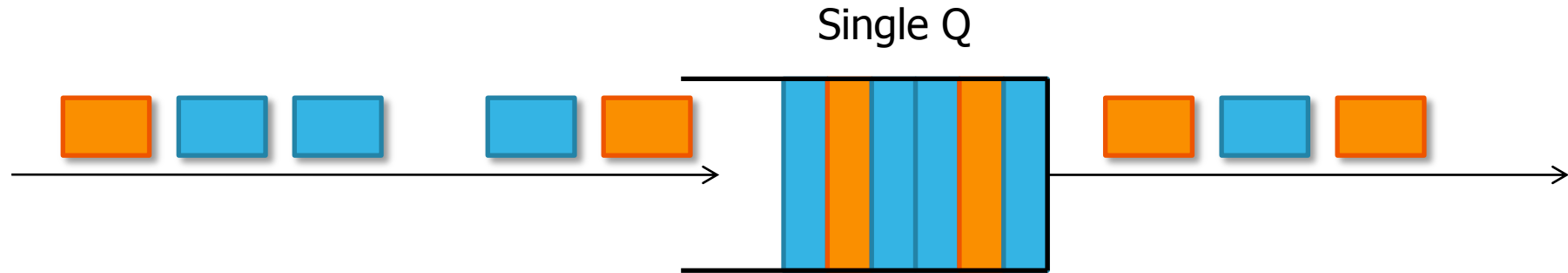# 1 TCP RENO FLOW (STEADY STATE)

# QUEUE SIZE AT DEQUEUE
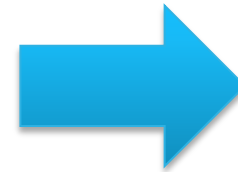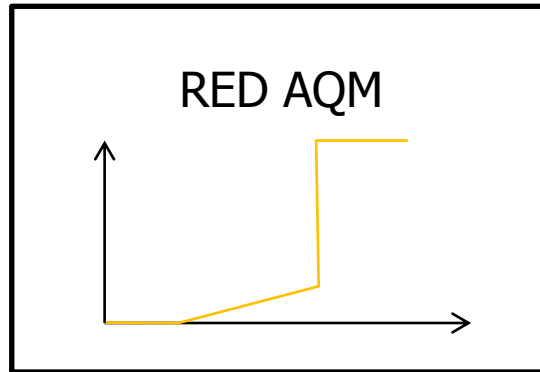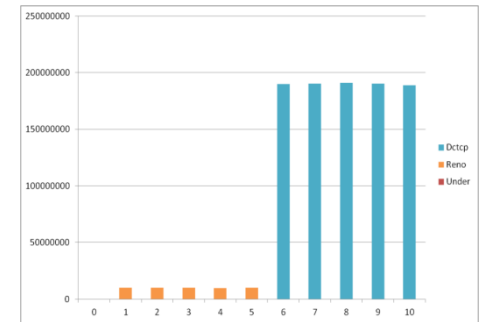# 1 DCTCP FLOW (STEADY STATE)



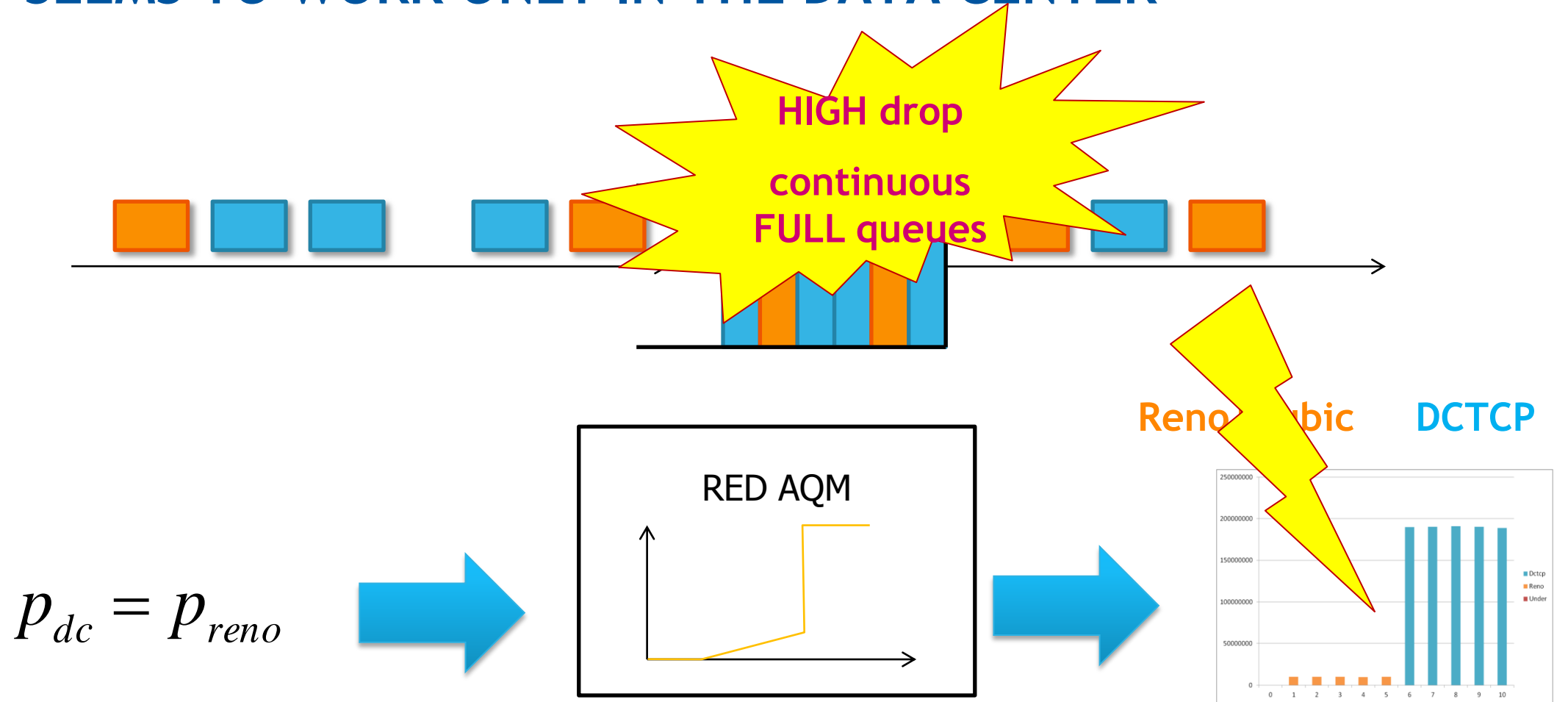Pdf in 1s interval [%]

# DCTCP DOES NOT WORK ON TRADITIONAL RED-ECN



Single Q

RED AQM

Reno|Cubic    DCTCP

$$p_{dc} = p_{reno}$$

# DCTCP SEEMS TO WORK ONLY IN THE DATA CENTER

HIGH drop

continuous
FULL queues

Reno Cubic        DCTCP

$$p_{dc} = p_{reno}$$

RED AQM

# THROUGHPUT:



RTT = 8 ms (unloaded)

BW = 40 Mbps (downstream)

BDP = 27 full sized packets

AQM = RED with recommended
          configuration*

X-axis: 0 – 250 sec

Y-axis: first row:
     0 – (80 / <nbr_flows>) Mbps

Y-axis: other rows
     0 – (80 / <nbr_dctcp>) Mbps

Cubic (= Reno) flows:
0   1   2   3   4   5   6   7   8   9   10

DCTCP flows: 0
1
2
3
4
5
6
7
8
9
10

R / T E
REDUCING INTERNET TRANSPORT LATENCY

* tc qdisc add dev eth2 root red limit 1600000 min 120000 max 360000 avpkt 1000 burst 220 ecn bandwidth 40Mbit

# Q SIZE PDF:

RTT = 8 ms (unloaded)

BW = 40 Mbps (downstream)

BDP = 27 full sized packets

AQM = RED with recommended

configuration*

X-axis: 0 – 300 packets
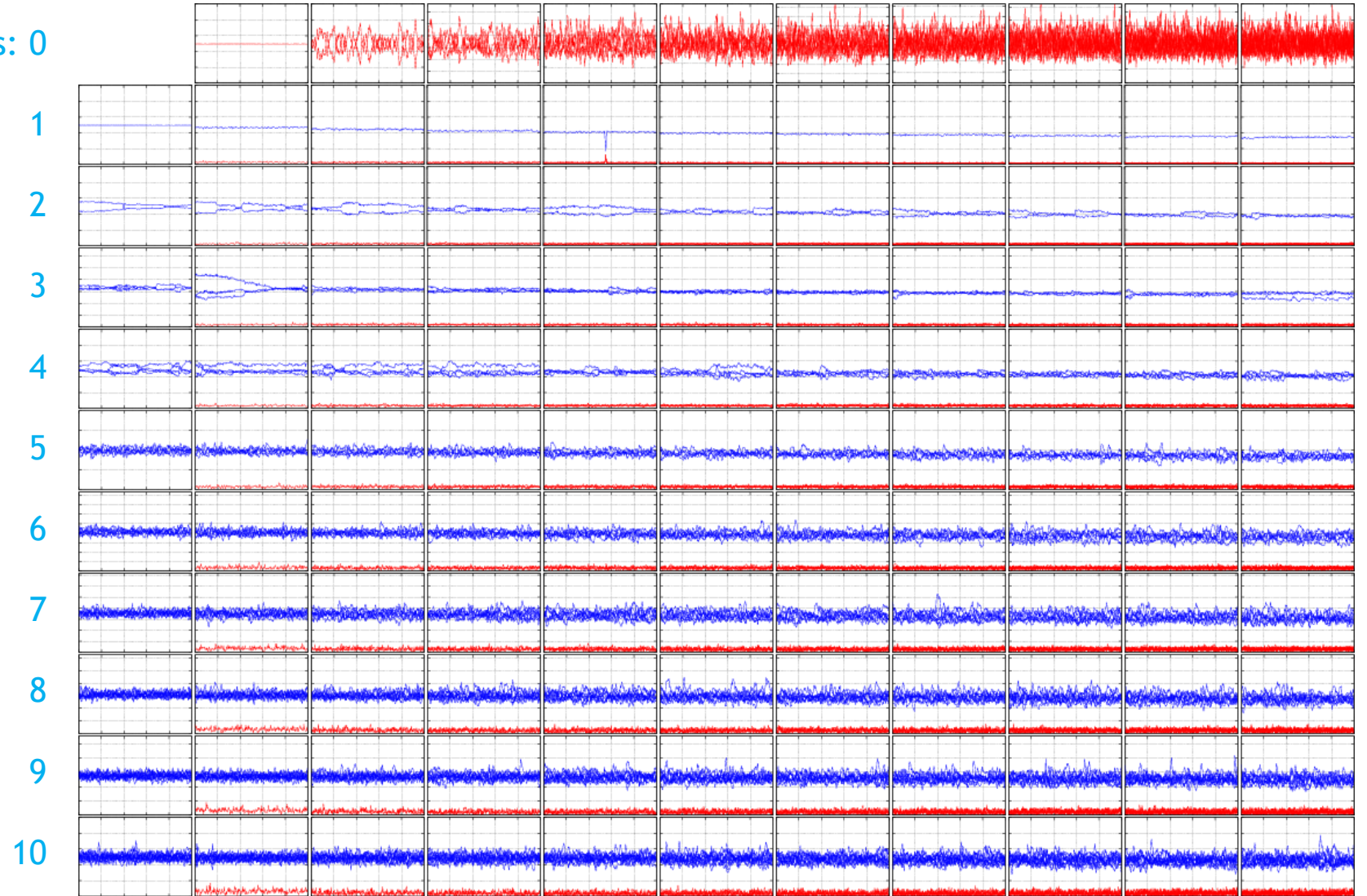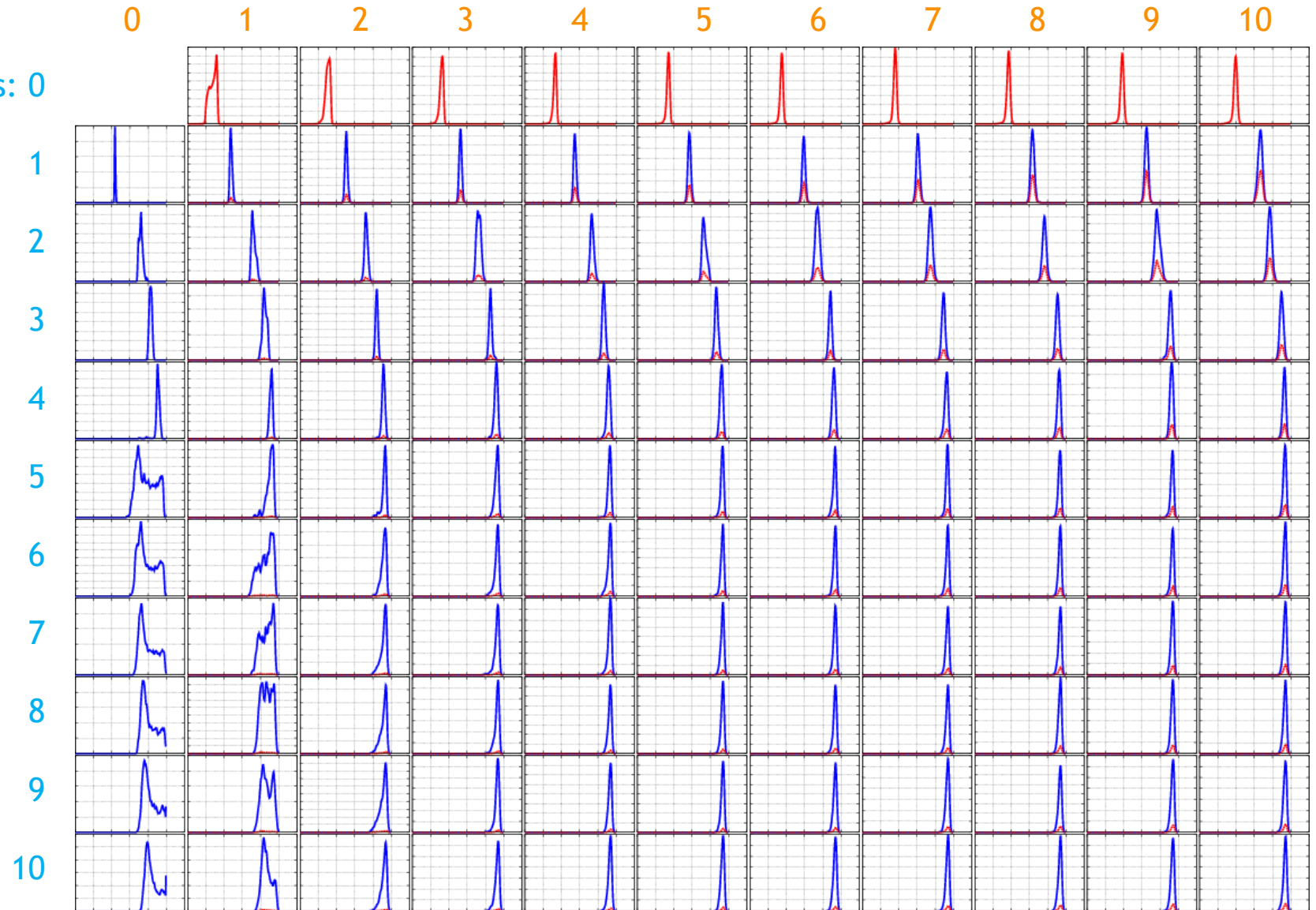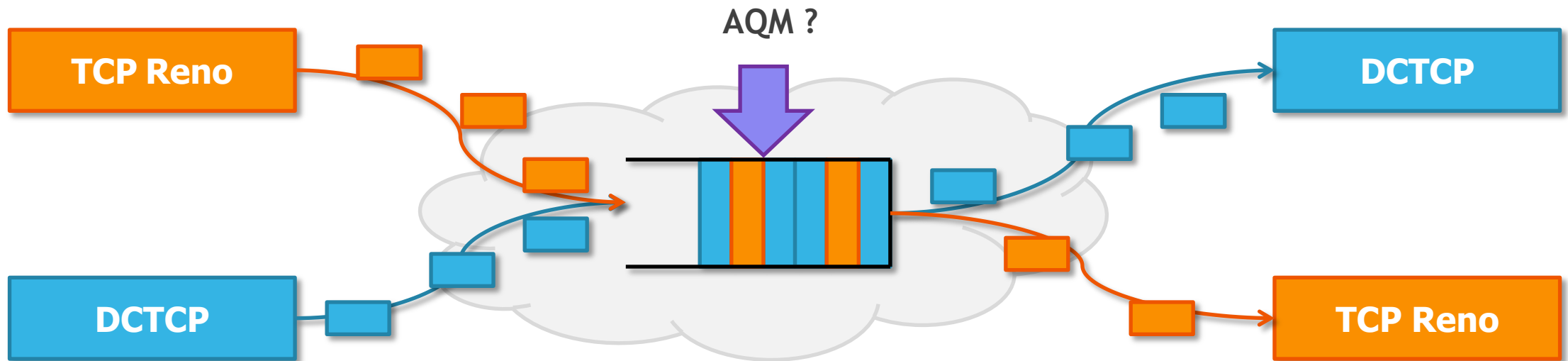
(450 Kbytes, 90 ms)

Y-axis: autoscale count packets

Cubic (= Reno) flows:

DCTCP flows:



* tc qdisc add dev eth2 root red limit 1600000 min 120000 max 360000 avpkt 1000 burst 220 ecn bandwidth 40Mbit

R / TE

REDUCING INTERNET TRANSPORT LATENCY

14

# AQMS FOR EQUAL STEADY STATE RATE MIGRATION PATH FOR NEW CC SCHEMES

- How should an AQM guarantee an equal steady state rate for flows with different congestion control schemes

  - classify packets according to CC schemes
  - align the drop/mark probabilities
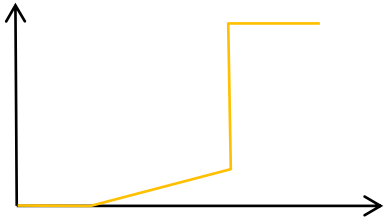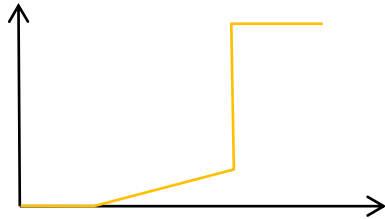
# TCP CONGESTION CONTROL SCHEMES
# STEADY STATE RATE

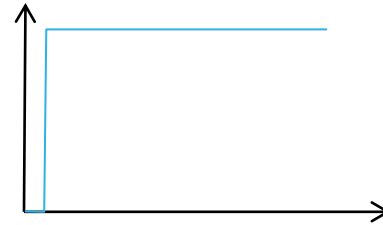- Steady state rate has been calculated for existing CC schemes:

$$r_{reno} = \frac{1.22}{p^{1/2}\mathrm{RTT}}$$

$$r_{cubic} = \frac{1.17}{p^{3/4}\mathrm{RTT}^{1/4}}$$
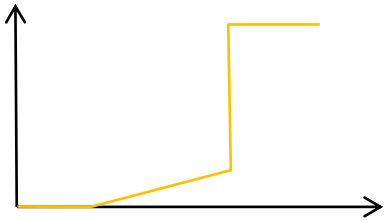
$$r_{dc} = \frac{2}{p^2 \cdot \mathrm{RTT}}$$
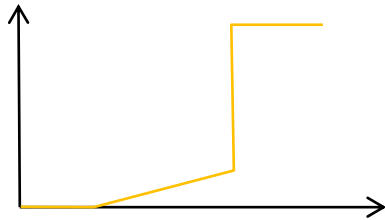
# TCP CONGESTION CONTROL SCHEMES STEADY STATE RATE
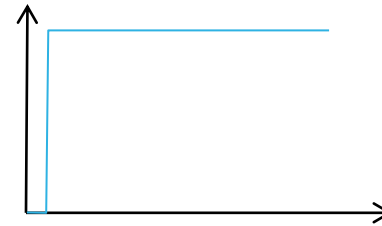
- Steady state rate has been calculated for existing CC schemes:

$$r_{reno} = \frac{1.22}{p^{1/2}\text{RTT}}$$

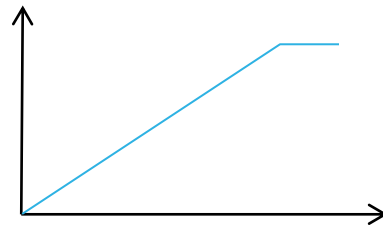$$r_{cubic} = \frac{1.17}{p^{3/4}\text{RTT}^{1/4}}$$

$$r_{dc} = \frac{2}{p^2 \cdot \text{RTT}}$$



- But we calculated that DCTCP running in non-on/off mode behaves as:

$$r_{dc\_p} = \frac{2}{p \cdot \text{RTT}}$$

RATE
REDUCING INTERNET TRANSPORT LATENCY

# TCP CONGESTION CONTROL SCHEMES FAIRNESS BETWEEN DCTCP AND RENO

- Mark/drop probability relation for equal rate and RTT:

$$r_{reno} = r_{dc}$$
$$RTT_{reno} = RTT_{dc}$$

$$\frac{1.22}{p_{reno}^{1/2} RTT_{reno}} = \frac{2}{p_{dc} \cdot RTT_{dc}}$$

$$p_{reno} = \left( \frac{p_{dc}}{1.63} \right)^2$$

# TCP CONGESTION CONTROL SCHEMES
# FAIRNESS BETWEEN DCTCP AND RENO

- Mark/drop probability relation for equal rate and RTT:

$$r_{reno} = r_{dc}$$
$$\text{RTT}_{reno} = \text{RTT}_{dc}$$

$$\frac{1.22}{p_{reno}^{1/2}\text{RTT}_{reno}} = \frac{2}{p_{dc} \cdot \text{RTT}_{dc}}$$

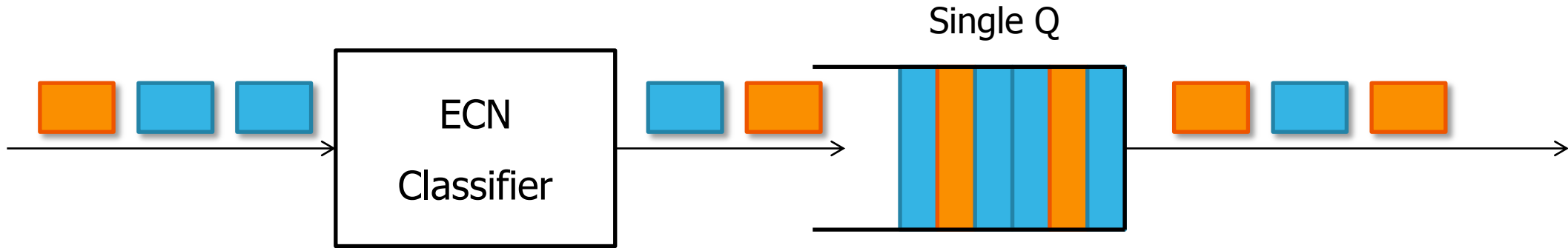$$p_{reno} = \left(\frac{p_{dc}}{1.63}\right)^2$$

**Square is easy!**

**Compare Q size with 2 random variables**

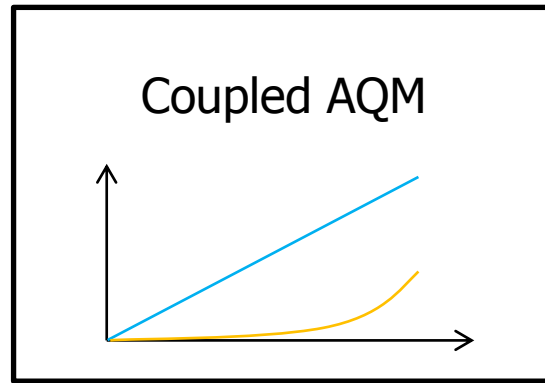$$P = f(Q) \qquad p \Rightarrow \text{Random}() < P$$

$$p^2 \Rightarrow (\text{Random}() < P) \,\&\&(\text{Random}() < P)$$

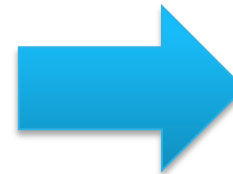$$p^2 \Rightarrow \max(\text{Random}(), \text{Random}()) < P$$

# DCTCP BEHAVES EXACTLY AS RENO
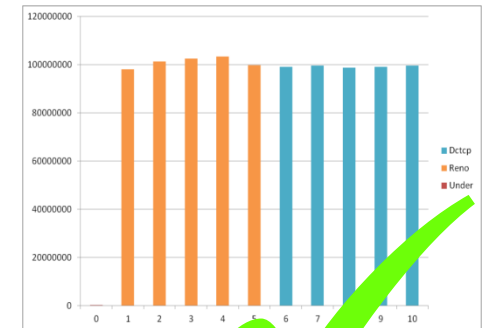# IF WE CORRECTLY CORRELATE MARKING AND DROPPING



Single Q

ECN Classifier

Coupled AQM

Instant Q size

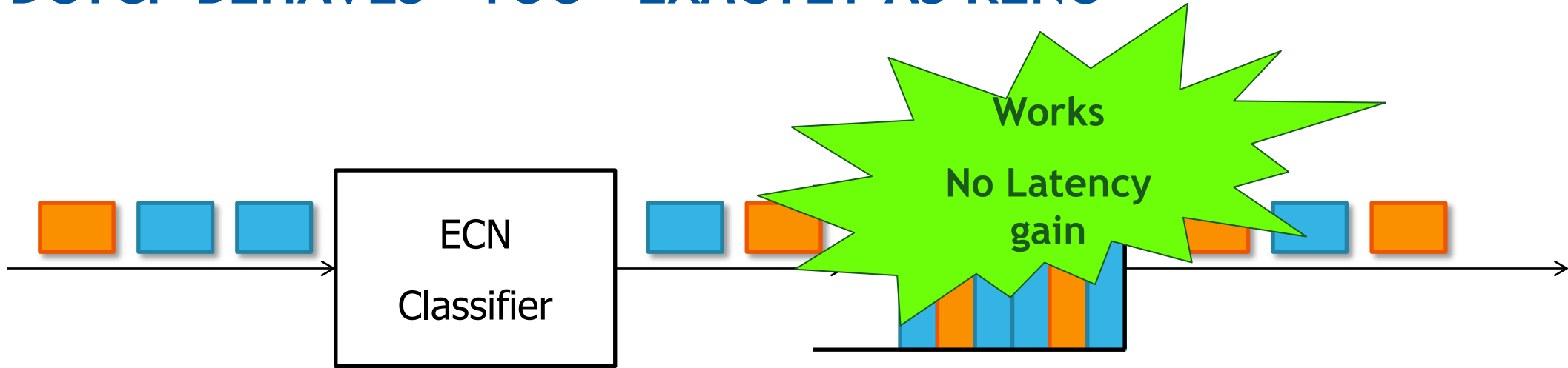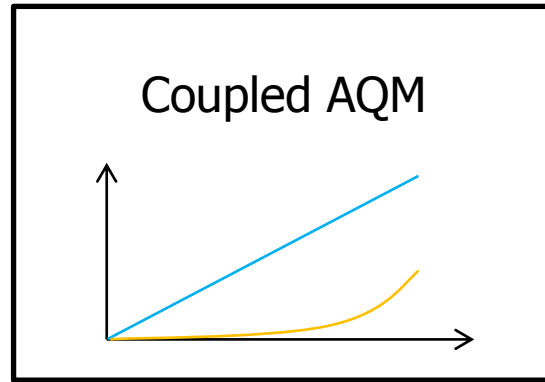$$p_{reno} = \left( \frac{p_{dc}}{1.63} \right)^2$$

**Reno|Cubic***     **DCTCP**

* Under local DC-access conditions (small BDP) Cubic behaves as Reno

Slope starts from the origin to avoid ON/OFF behavior in steady state

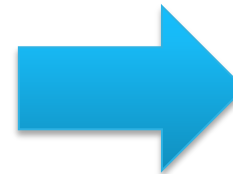R/TE
REDUCING INTERNET TRANSPORT LATENCY

20

# DCTCP BEHAVES "TOO" EXACTLY AS RENO

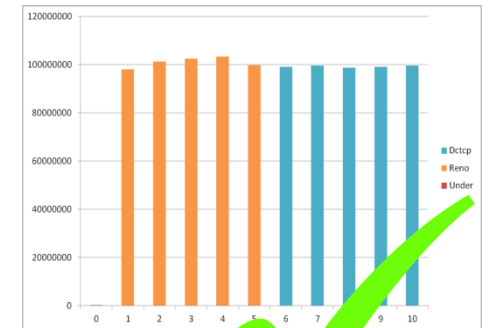ECN Classifier

Works

No Latency gain

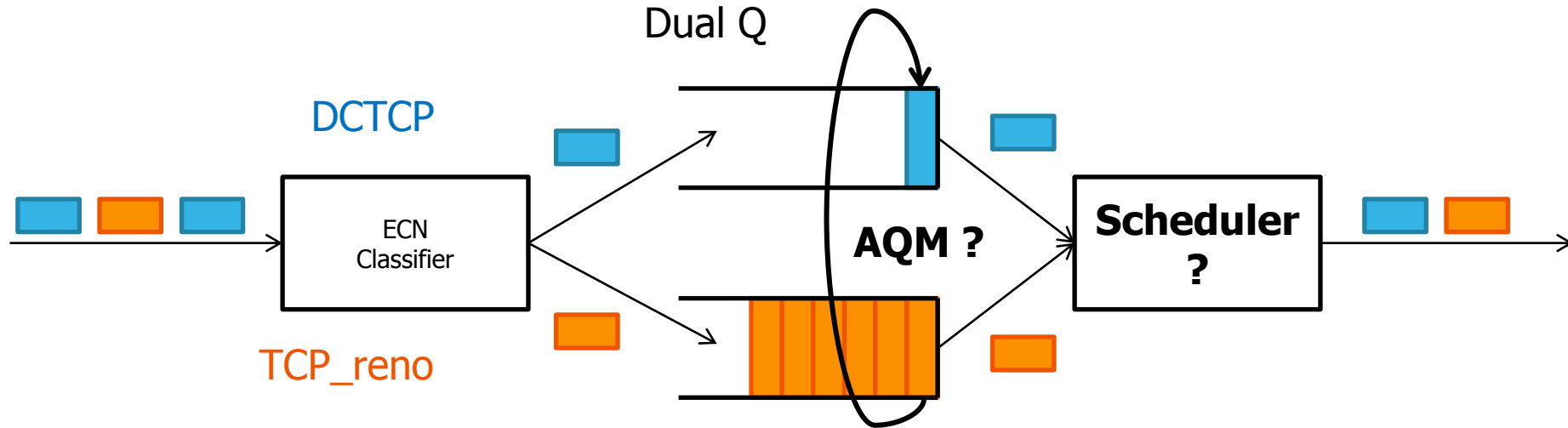$$p_{reno} = \left( \frac{p_{dc}}{1.63} \right)^2$$
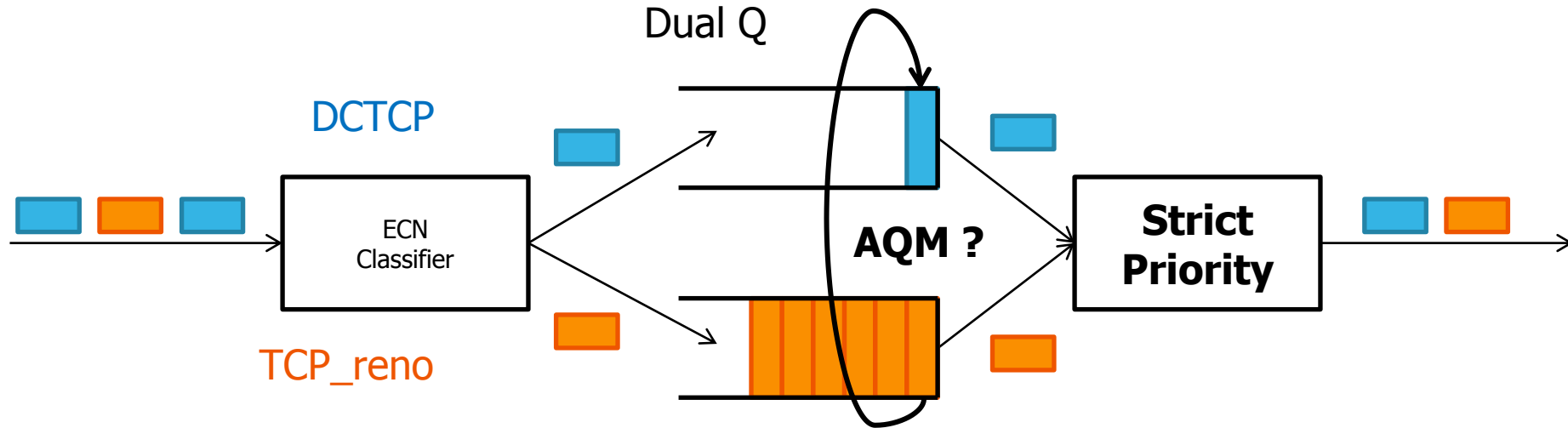
Coupled AQM

Instant Q size

Reno|Cubic    DCTCP

# DUAL QUEUE - LOW LATENCY



$$\frac{r_{reno}}{r_{dc}} = \frac{1.22}{2} \frac{p_{dc}}{\sqrt{p_{reno}}} \frac{\text{RTT}_{dc}}{\text{RTT}_{reno}}$$

$$p_{reno} = \left(\frac{p_{dc}}{8}\right)^2$$

R⁄TE

# DUAL QUEUE - LOW LATENCY



Dual Q

DCTCP

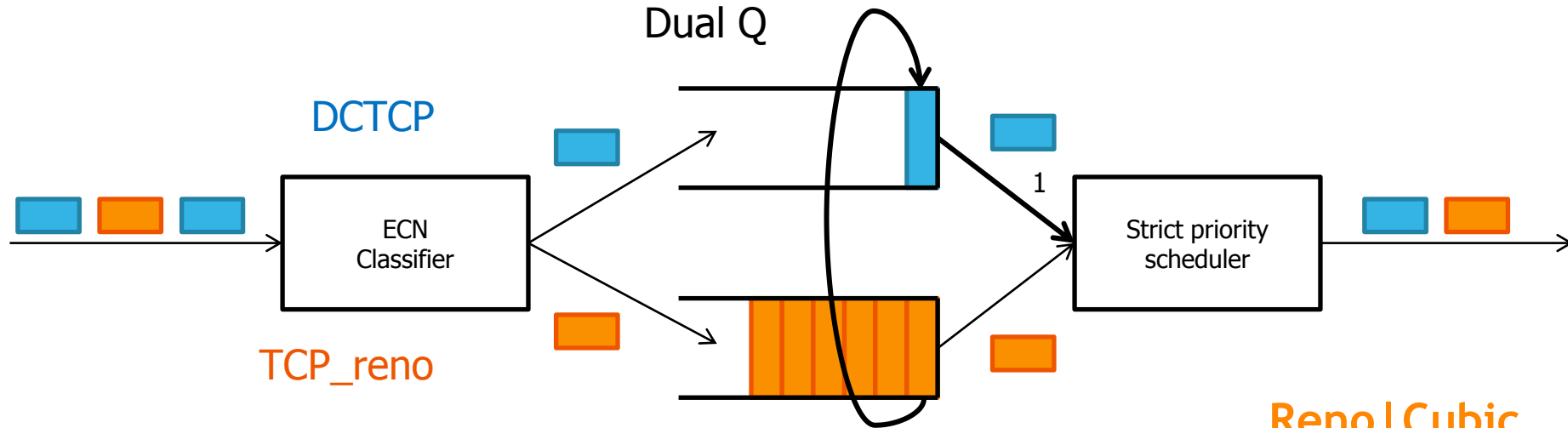ECN Classifier

AQM ?

Strict Priority

TCP_reno

$$1/5 = 8 \text{ ms} /(8 + 32) \text{ ms}$$

$$\frac{r_{reno}}{r_{dc}} = \frac{1.22}{2} \frac{p_{dc}}{\sqrt{p_{reno}}} \cdot \frac{\text{RTT}_{dc}}{\text{RTT}_{reno}} = 5$$

$$p_{reno} = \left(\frac{p_{dc}}{8}\right)^2$$

R/TE
REDUCING INTERNET TRANSPORT LATENCY

# DUAL QUEUE - LOW LATENCY



Dual Q

DCTCP

ECN Classifier

TCP_reno

Strict priority scheduler

1

$$p_{reno} = \left(\frac{p_{dc}}{8}\right)^2$$

Coupled AQM

Instant Q time

Reno | Cubic    DCTCP

REDUCING INTERNET TRANSPORT LATENCY

# DUAL QUEUE - LOW LATENCY

Dual Q

DCTCP

ZERO
Q latency

ECN
Classifier

TCP_reno

1

Strict priority
scheduler

Reno|Cubic    DCTCP

$$p_{reno} = \left( \frac{p_{dc}}{8} \right)^2$$

Coupled AQM

Instant Q time

**Measure Q in time is important for optimal fairness !**

# THROUGHPUT:

RTT = 8 ms (unloaded)

BW = 40 Mbps (downstream)

BDP = 27 full sized packets

AQM = DualQ Coupled

X-axis: 0 – 250 sec

Y-axis: all rows:

0 – (80 / <nbr_flows>) Mbps



DCTCP flows: 0

R / T E

REDUCING INTERNET TRANSPORT LATENCY

# Q SIZE PDF:

RTT = 8 ms (unloaded)

BW = 40 Mbps (downstream)

BDP = 27 full sized packets

AQM = DualQ Coupled

X-axis: 0 – 300 packets

(450 Kbytes, 90/w ms)

Y-axis: autoscale count packets



R / T E
REDUCING INTERNET TRANSPORT LATENCY

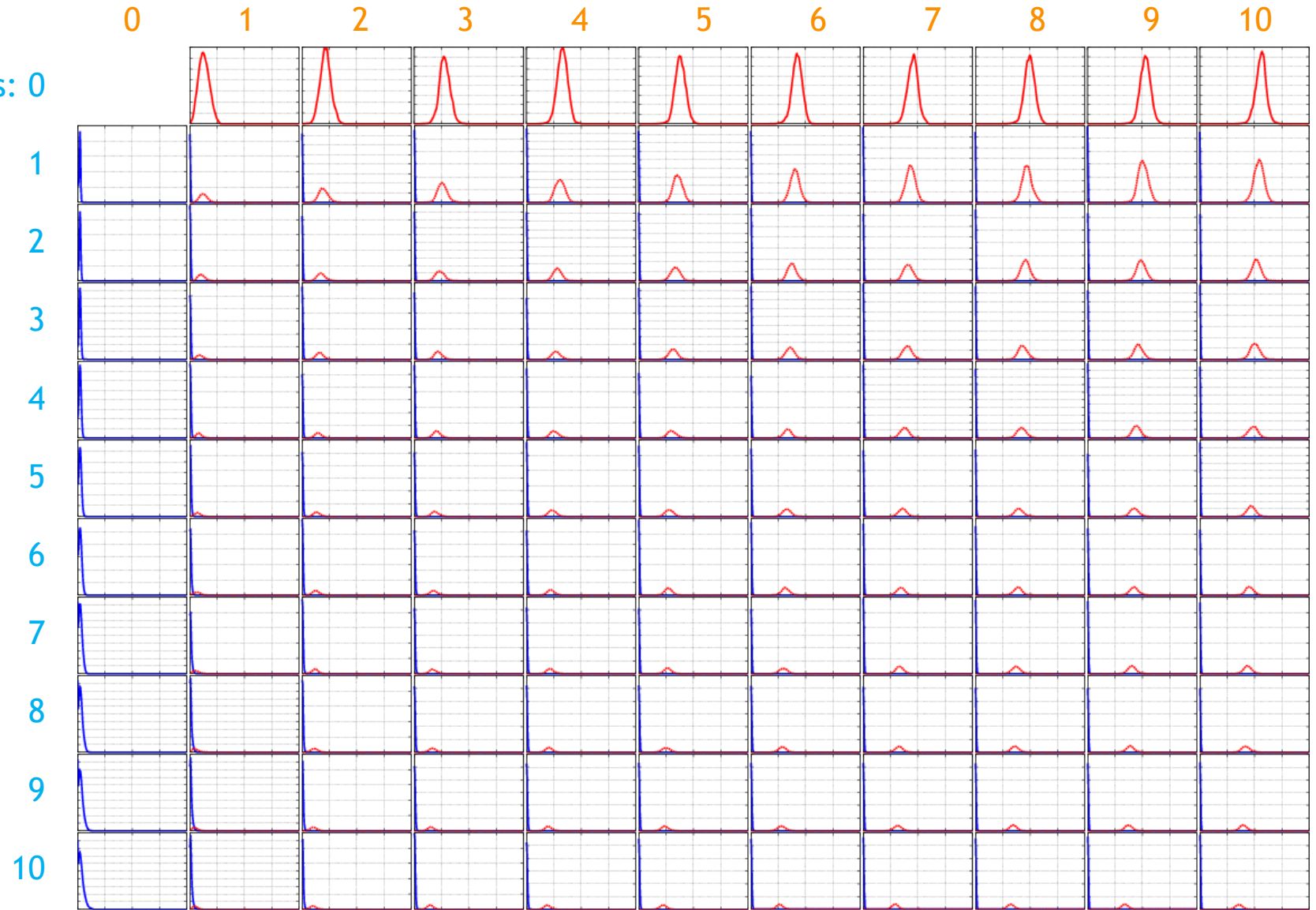# THROUGHPUT RATIO (CUBIC / DCTCP)



RED AQM

DualQ Coupled AQM

# DETAILED IMPLEMENTATION

3 parameters:

- Reno slope (bits)
- DCTCP slope (bits)
- DCTCP threshold (Q size)

$$p_{reno} = \left(\frac{p_{dc}}{8}\right)^2$$

$$S_r = S_d - 3$$

# ADAPTIVE INTERACTIVE APPLICATIONS

- Panoramic interactive video



- Video/Voice conferencing



- Remote control, ....

# FUTURE WORK & CONCLUSIONS ?

- Dynamic behaviour to be investigated (expected to be 5x better due to 5x latency reduction)
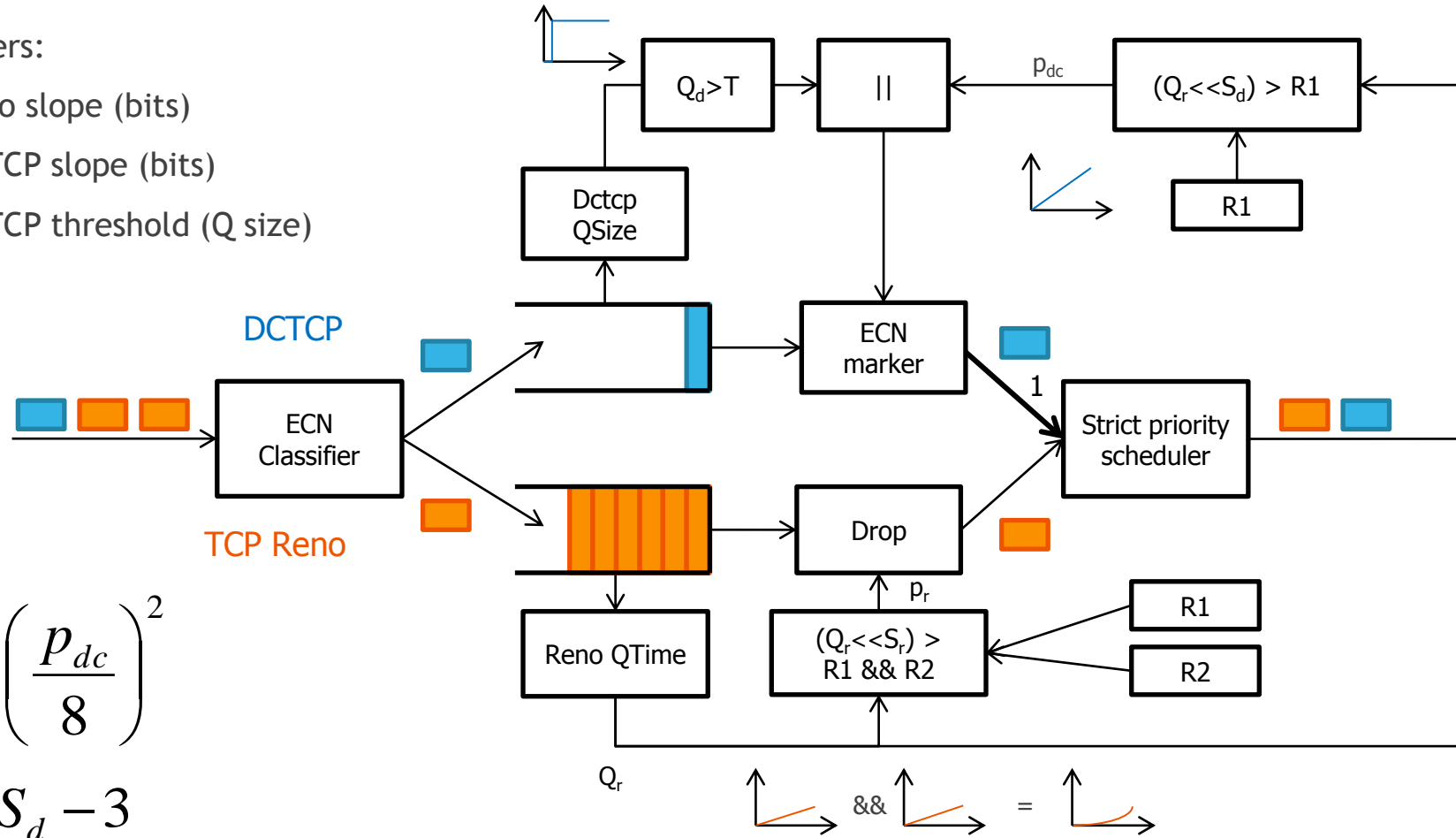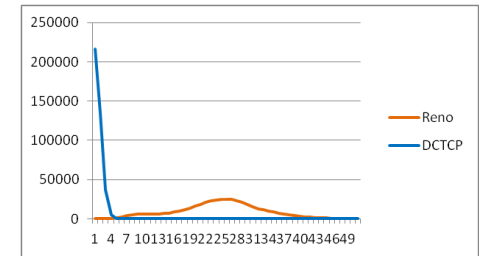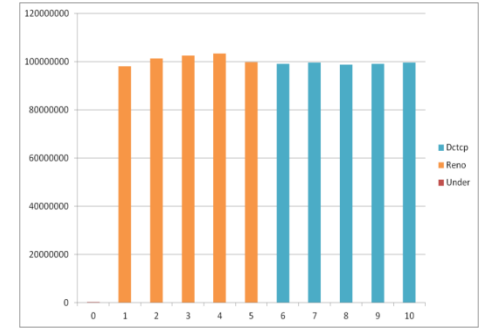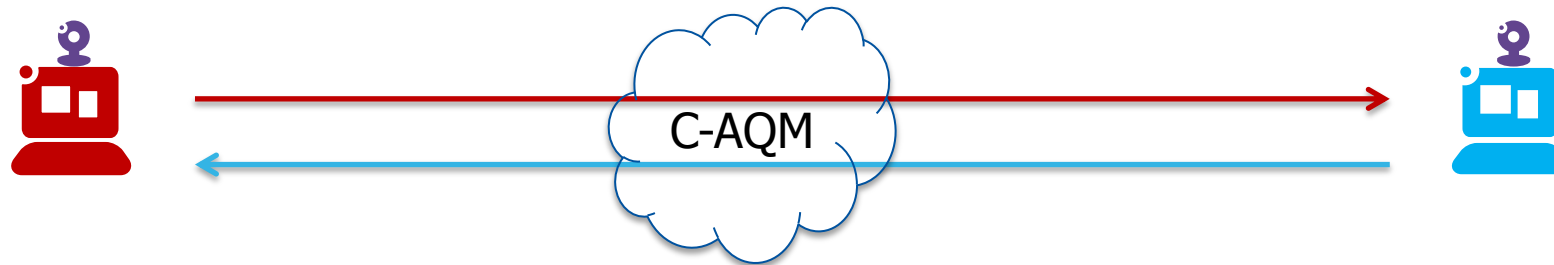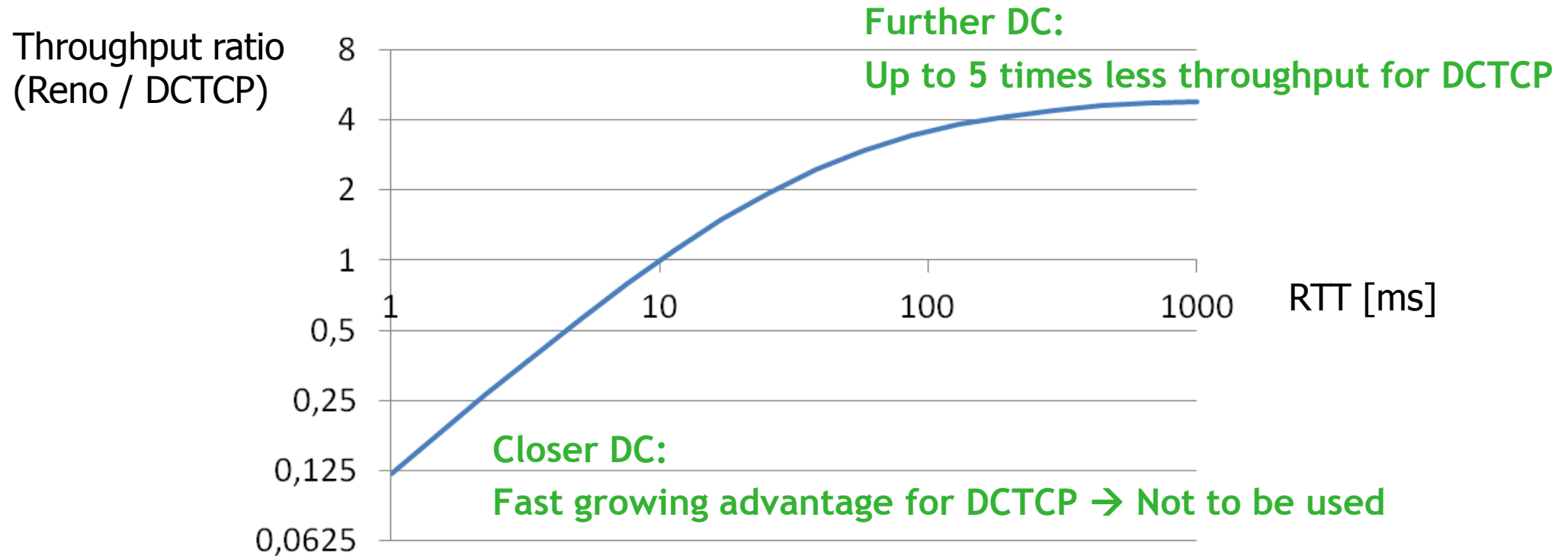
- Unmanaged Low Level network Service ➔ Native support for Adaptive Interactive Applications

- Better usage for ECN (marking can be more often than dropping)
  - ➔ x/p relation for ECN based congestion controller (x determining the marking rate)
    - ➔ $p^2$ relation between mark and drop in AQM

- Backwards compatible: DCTCP should respond to drop as Reno
  - Currently 3.18 Linux implementation fails on this aspect

- Steady state throughput fairness between DCTCP and Reno|Cubic
  - Only if DCTCP flows are terminated to nearby (local) datacenter
  - If longer RTT, DCTCP flows are getting lower throughput than Reno|Cubic
  - ➢ Reno|Cubic fallback if throughput is too low and base RTT is too long ?
  - ➢ Define a TCP congestion controller which is less (/not) dependent on the RTT ?

R / T E
REDUCING INTERNET TRANSPORT LATENCY

# Questions

koen.de_schepper@alcatel-lucent.com

R / TE
Reducing Internet Transport Latency

# BASE RTT FAIRNESS
# WHAT IF THE DATACENTER IS FURTHER OR CLOSER

Throughput ratio
(Reno / DCTCP)

**Further DC:**
**Up to 5 times less throughput for DCTCP**

**Closer DC:**
**Fast growing advantage for DCTCP → Not to be used**

RTT [ms]

Coupled AQM configured for 10ms base RTT and 40ms Reno queue time (1/5 RTT ratio)

R / TE
REDUCING INTERNET TRANSPORT LATENCY

# DCTCP STEADY STATE THROUGHPUT WITH SLOPE-RED

Per "long" RTT: $\quad W \leftarrow W + 1 \qquad$ (1)

And also per RTT: $\quad W \leftarrow W\left(1 - \dfrac{\alpha}{2}\right) \qquad$ (2)

In steady state if (1) is compensated $\quad W \leftarrow W - 1 = W\left(1 - \dfrac{1}{W}\right) \quad$ so from (2) if $\quad \dfrac{\alpha}{2} = \dfrac{1}{W} \quad$ (3)
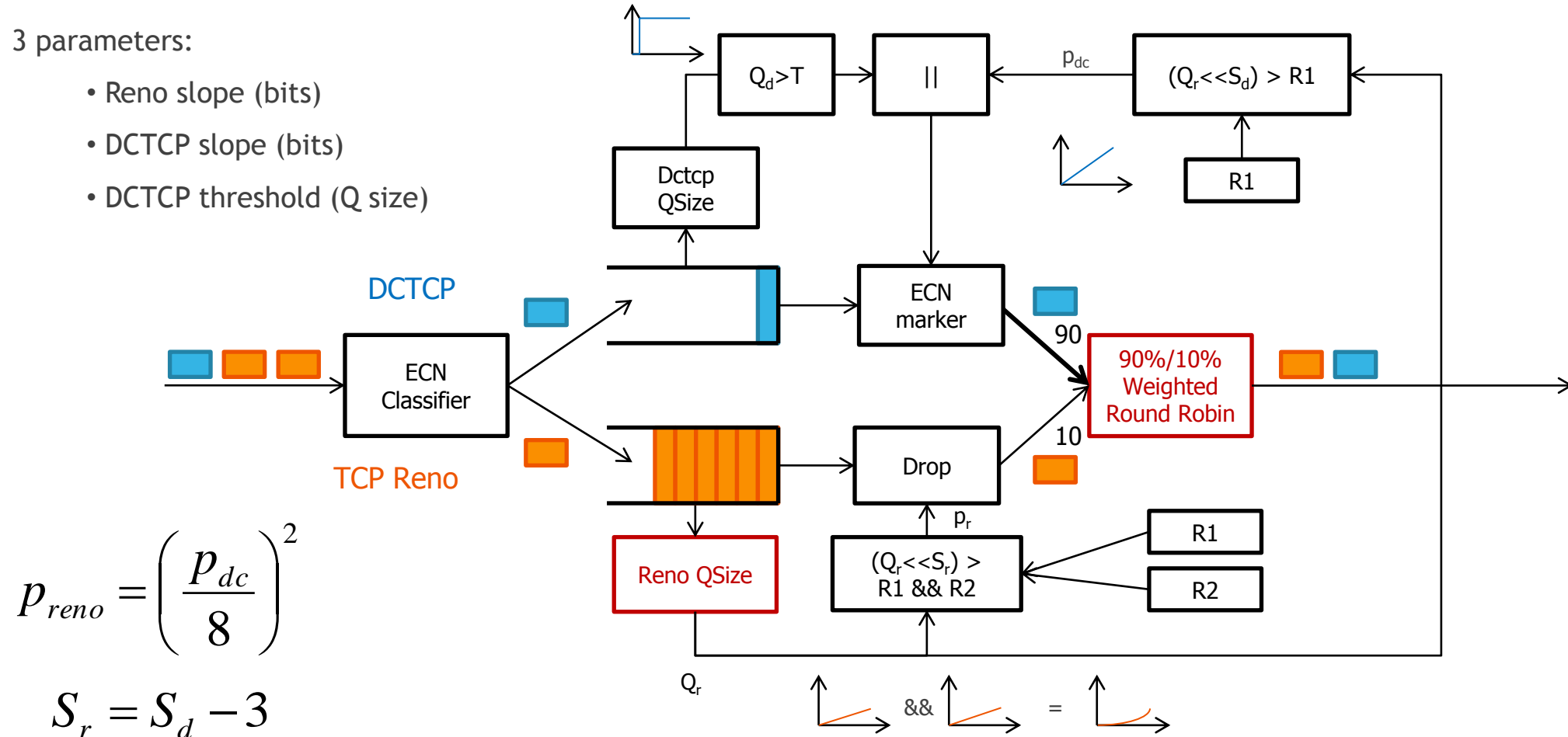
As $\quad \alpha \leftarrow (1 - g)\alpha + gp \quad$, $\quad$ if $p$ is stable in steady state $\qquad \alpha = p \qquad$ (4)

The instantaneous rate $\quad r = \dfrac{W}{rtt} \qquad$ (5) $\quad$ thus (3,4,5) $\quad r = \dfrac{2}{p \cdot rtt}$

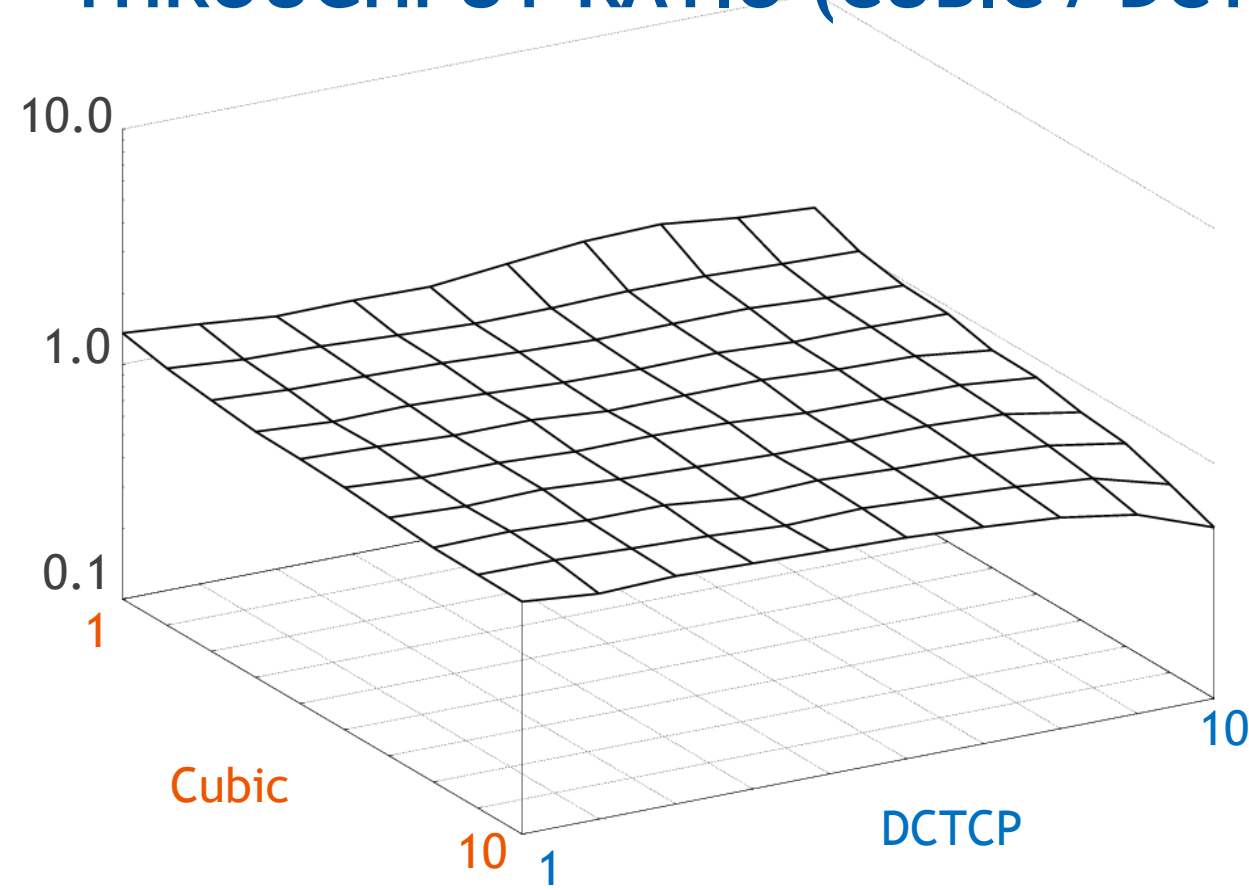# QUEUE SIZE BASED COUPLED AQM

3 parameters:

- Reno slope (bits)
- DCTCP slope (bits)
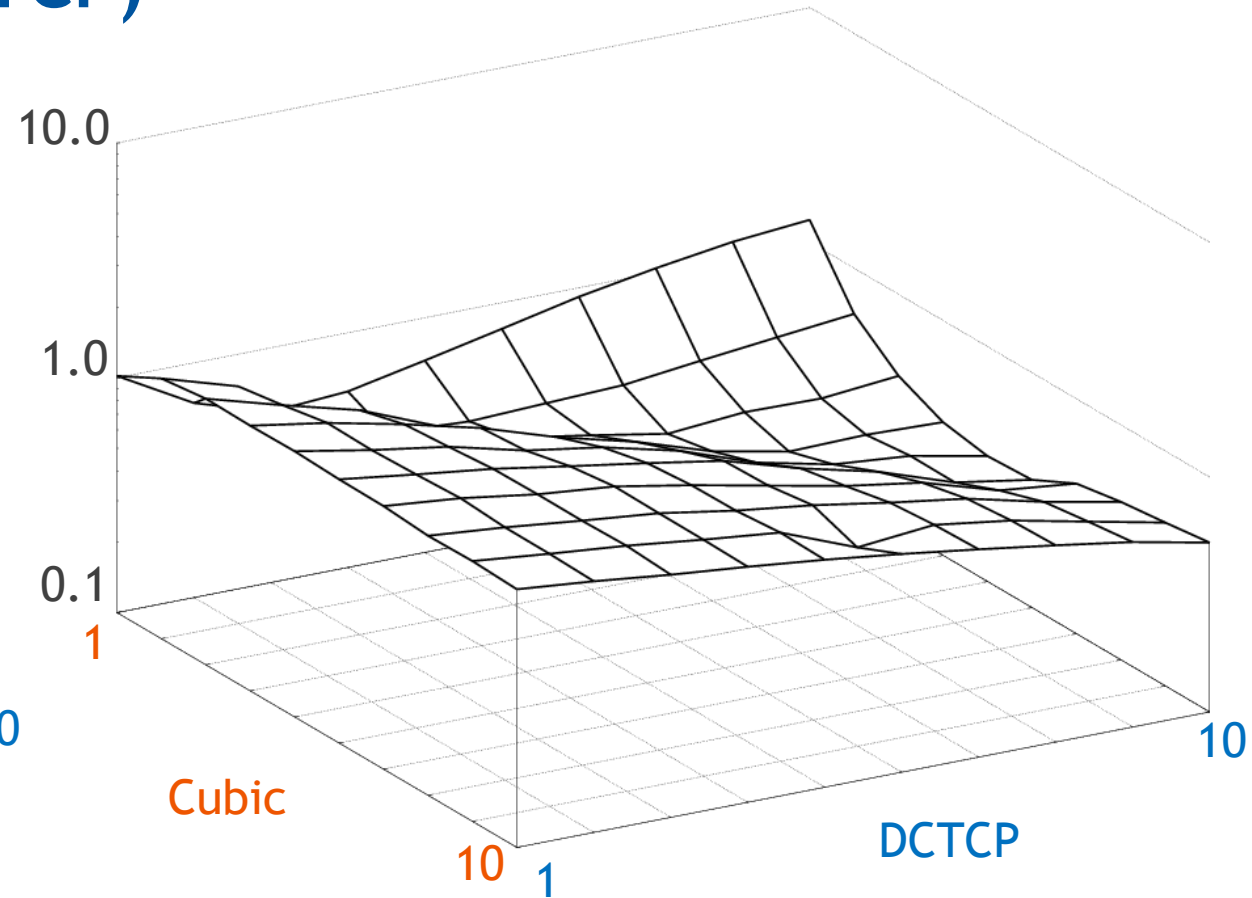- DCTCP threshold (Q size)

$$p_{reno} = \left(\frac{p_{dc}}{8}\right)^2$$

$$S_r = S_d - 3$$

# QUEUE SIZE BASED COUPLED AQM THROUGHPUT RATIO (CUBIC / DCTCP)



Qtime - DualQ Coupled AQM

Qsize - DualQ Coupled AQM