# TCP ex Machina:
# Computer-Generated Congestion Control

Keith Winstein

MIT Computer Science and Artificial Intelligence Laboratory
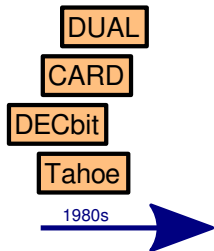
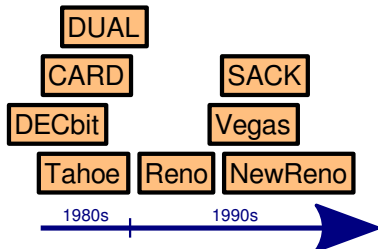http://mit.edu/remy

March 3, 2014



**C S A I L**

Joint work with Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The march of congestion-control protocols



DUAL
CARD
DECbit
Tahoe

1980s

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The march of congestion-control protocols

# The march of congestion-control protocols

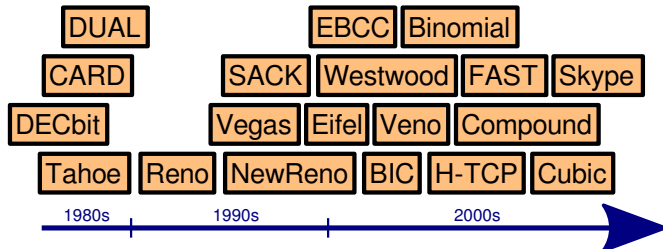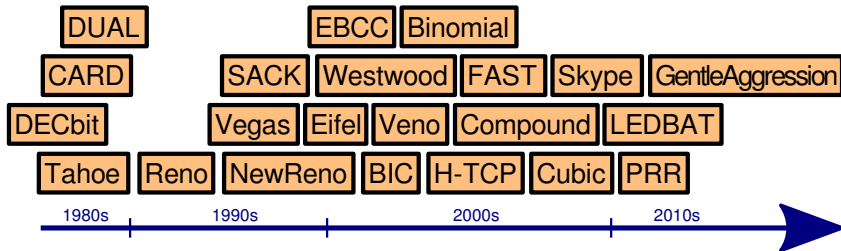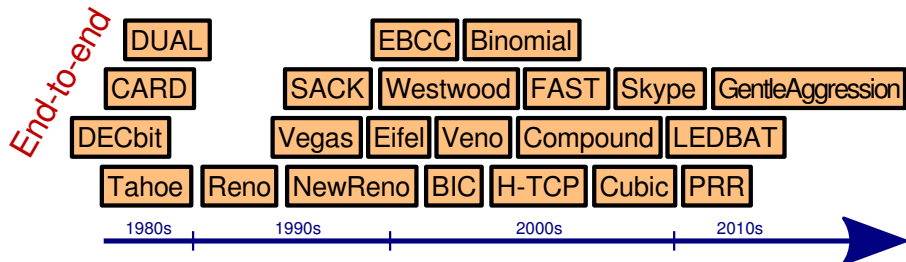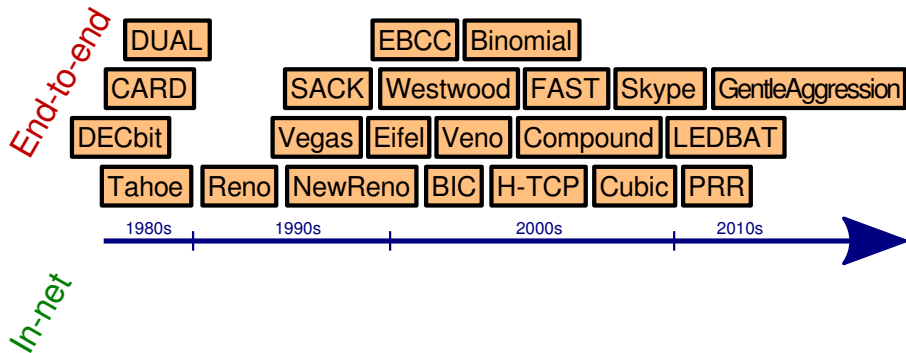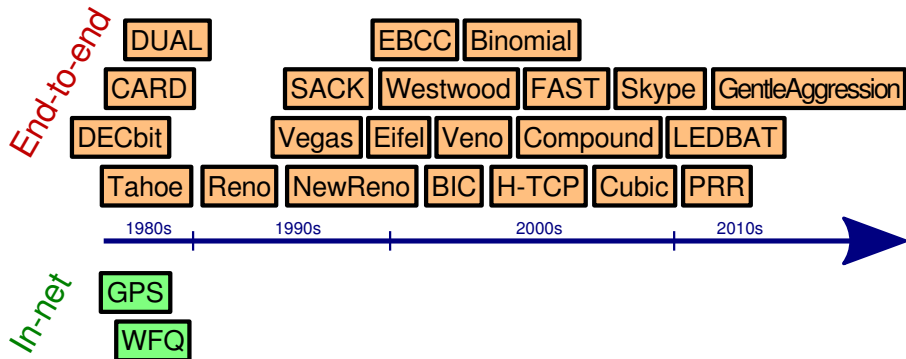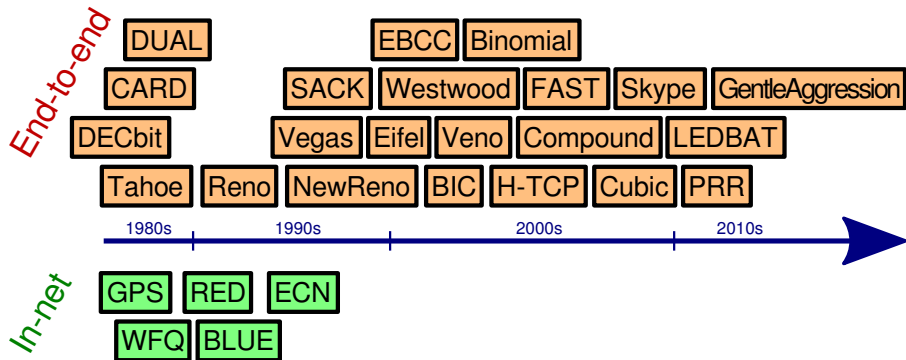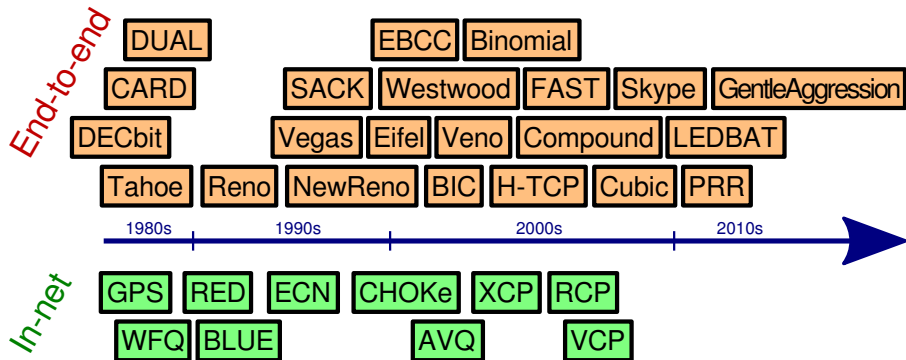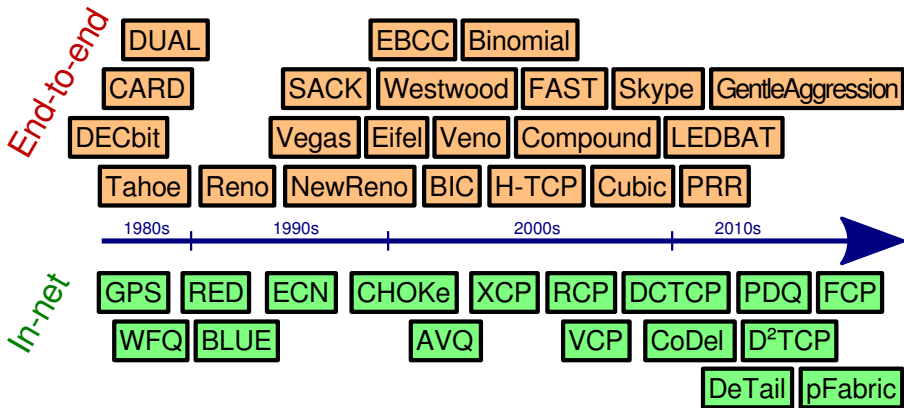Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The march of congestion-control protocols

# The march of congestion-control protocols

# The march of congestion-control protocols

# The march of congestion-control protocols

# The march of congestion-control protocols

# The march of congestion-control protocols



End-to-end

| DUAL | | | | EBCC | Binomial | | | |
| CARD | | SACK | Westwood | | FAST | Skype | GentleAggression |
| DECbit | | Vegas | Eifel | Veno | Compound | LEDBAT | |
| Tahoe | Reno | NewReno | BIC | H-TCP | Cubic | PRR | |

1980s          1990s          2000s          2010s

In-net

| GPS | RED | ECN | CHOKe | XCP | RCP |
| WFQ | BLUE | | AVQ | | VCP |

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The march of congestion-control protocols

# The march of congestion-control protocols

# Rational choice of scheme is challenging

Cubic **vs.** Compound

- ▶ Different missions?
- ▶ Different assumptions about network?
- ▶ One scheme just plain better?

# Networks constrained by a fuzzy idea of TCP's assumptions

- ► Mask stochastic loss
- ► Bufferbloat
- ► Mask out-of-order delivery
- ► No parallel/multipath routing

*Advice for Internet Subnetwork Designers*
(RFC 3819) is 21,000 words!

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

## Apps hack around TCP

- Open lots of flows
- Goose slow start
- Add pacing
- Give up and do it yourself

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Apps hack around TCP

▶ Open lots of flows

▶ Goose slow start

▶ Add pacing

▶ Give up and do it yourself

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Apps hack around TCP

- Open lots of flows
- Goose slow start
- Add pacing
- Give up and do it yourself

# Apps hack around TCP

- ▶ Open lots of flows
- ▶ Goose slow start
- ▶ Add pacing
- ▶ Give up and do it yourself

# Apps hack around TCP

- Open lots of flows
- Goose slow start
- Add pacing
- Give up and do it yourself





**Google**  **MICROSOFT**

**You Tube**

**Chrome** (QUIC)
**BitTorrent** ($\mu$TP)
**Mosh** (SSP)
**IBM Aspera** (fasp)

# Idea: computer-generated protocols

Transport layer should adapt to **whatever**:

- network does
- application wants

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

Transport layer should adapt to **whatever**:

- network does (**model**)
- application wants (**mission**)

# What we built

> **Remy**: a program that generates
> congestion-control schemes offline

**Input:**

- Assumptions about network and workload (**model**)
- Application's objective (**mission**)

**Output:** CC algorithm for a TCP sender  (RemyCC)

**Time:** hours to days

# The basic question of congestion control

> At this moment, do I:
>
> - send a packet
> - not send a packet?

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

## **Missions** of congestion control

**Maximize**

- $\sum_i \log[\text{throughput}_i]$    (proportionally fair throughput)

# Missions of congestion control

**Maximize**

- $\displaystyle\sum_i \log\left[\text{throughput}_i\right]$    (proportionally fair throughput)

- $\displaystyle\sum_i \log\left[\frac{\text{throughput}_i}{\text{delay}_i}\right]$    (proportionally fair throughput/delay)

# **Missions** of congestion control

**Maximize**

- ▶ $\sum_i \log \left[ \text{throughput}_i \right]$   (proportionally fair throughput)

- ▶ $\sum_i \log \left[ \dfrac{\text{throughput}_i}{\left( \text{delay}_i \right)^{\delta}} \right]$   (proportionally fair throughput/delay)

# **Missions** of congestion control

**Maximize**

- $\sum_i \log[\text{throughput}_i]$  (proportionally fair throughput)

- $\sum_i \log\left[\dfrac{\text{throughput}_i}{(\text{delay}_i)^{\boldsymbol{\delta}}}\right]$  (proportionally fair throughput/delay)

- $\min_i \text{throughput}_i$        (max-min throughput)

**Minimize**

- average flow completion time
- page load time
- tail completion time

## Encoding the designer's prior assumptions

- ▶ **Model** of network uncertainty
  - ▶ Link speed distribution
  - ▶ Delay distribution
  - ▶ Topology distribution

- ▶ **Model** of workload
  - ▶ Web browsing
  - ▶ MapReduce
  - ▶ videoconferencing
  - ▶ streaming video (YouTube/Netflix)

# Dumbbell network

# Dumbbell network

# Dumbbell network

# Superrational congestion control

> At this moment,* do I:
>
> - send a packet
>
> - not send a packet?

# Superrational congestion control

> At this moment,* do I:
> - send a packet
> - not send a packet?

**\*** Assuming every node is running the same algorithm.

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

## Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

$T$: simple networks are deterministic, with parameters drawn from a distribution. Flows arrive per a random process.

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

$T$: simple networks are deterministic, with parameters drawn from a distribution. Flows arrive per a random process.

$R$: overall objective awarded at end of long run

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net + whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

$T$: simple networks are deterministic, with parameters drawn from a distribution. Flows arrive per a random process.

$R$: overall objective awarded at end of long run

$\Omega$: acknowledgements from corresponding receiver

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

$T$: simple networks are deterministic, with parameters drawn from a distribution. Flows arrive per a random process.

$R$: overall objective awarded at end of long run

$\Omega$: acknowledgements from corresponding receiver

$O$: deterministic

# Internet congestion control as a Dec-POMDP

$I$: independent endpoint computers

$S$: packets in net $+$ whether each computer has data to send now

$A_i$: { send a packet, don't send a packet }

$T$: simple networks are deterministic, with parameters drawn from a distribution. Flows arrive per a random process.

$R$: overall objective awarded at end of long run

$\Omega$: acknowledgements from corresponding receiver

$O$: deterministic

**Complexity in general:** $\mathcal{O}(2^{2^n})$ ($\mathrm{NEXP}$-hard)

# Remy: tractable search for best policy

- ▶ Best decision given all history: not tractable

- ▶ Instead, remember only **a summary of history**

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# A RemyCC tracks four congestion signals

$r\_ewma_\alpha$: **short-term** moving average of interval between acks
*"How fast are packets arriving (now)?"*

$r\_ewma_\beta$: **long-term** moving average of same
*"How fast are packets arriving (smoothed)?"*

$s\_ewma$: moving average of interval between acked timestamps
*"How fast was I sending?"*

$rtt\_ratio$: ratio of last RTT to smallest RTT so far
*"How long is the queue?"*

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Why these congestion signals?

- Removing any of the four hurts
  - $r\_ewma_\alpha$ hurts the most

- More signals increase search time

- Other signals might help on other networks

# A RemyCC maps each state to an action

$$\textsc{RemyCC}(r\_ewma_{\alpha\beta}, s\_ewma, rtt\_ratio) \rightarrow \langle m, b, \tau \rangle$$

$m$  Multiple to congestion window

$b$  Increment to congestion window

$\tau$  Minimum interval between two outgoing packets

# Runtime for a RemyCC

**On ack:**

- $\langle m, b, \tau \rangle \leftarrow \mathrm{RemyCC}(r\_ewma_{\alpha\beta}, s\_ewma, rtt\_ratio)$
- $\mathrm{cwnd} \leftarrow m \cdot \mathrm{cwnd} + b$

**Send packet if:**

- $\mathrm{cwnd} > \mathrm{FlightSize}$, and
- last packet sent $> \tau$ ago

Find piecewise-continuous $\textsc{RemyCC}()$ that optimizes expected value of objective function

**On ack:**

$$\langle m, b, \tau \rangle \leftarrow \text{RemyCC}(s\_ewma, r\_ewma_\alpha, \ r\_ewma_\beta, rtt\_ratio \ )$$

# Remy example: 2D state space

**On ack:**

$$\langle m, b, \tau \rangle \leftarrow \text{RemyCC}(s\_ewma, r\_ewma_\alpha, \blacksquare)$$

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

## Remy example: **model**

| Quantity | Distribution | Units |
|----------|-------------|-------|
| Link speed | Uniform(10, 20) | Mbps |
| RTT | Uniform(100, 200) | ms |
| $n$ | Uniform(1, 16) | |
| "On" process | $\exp[\mu = 5]$ | seconds |
| "Off" process | same | |

# Remy example: **mission**

$$\sum_i \log \left[ \frac{\text{throughput}_i}{\text{delay}_i} \right]$$

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# One action for all states. Find the best value.

# The best (single) action. Now split it on median.



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize each of the new actions



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Now split the most-used rule

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
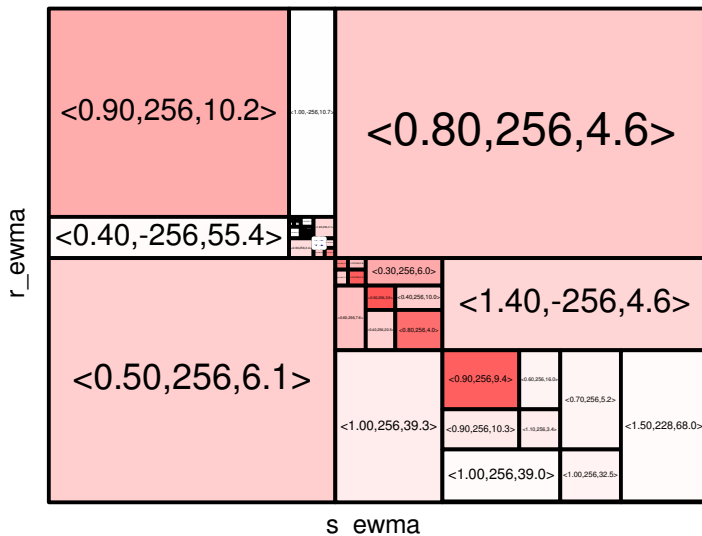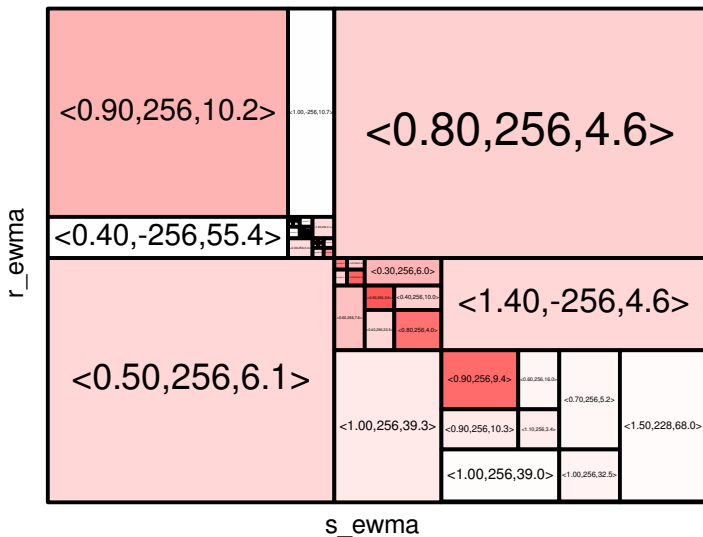
# Optimize
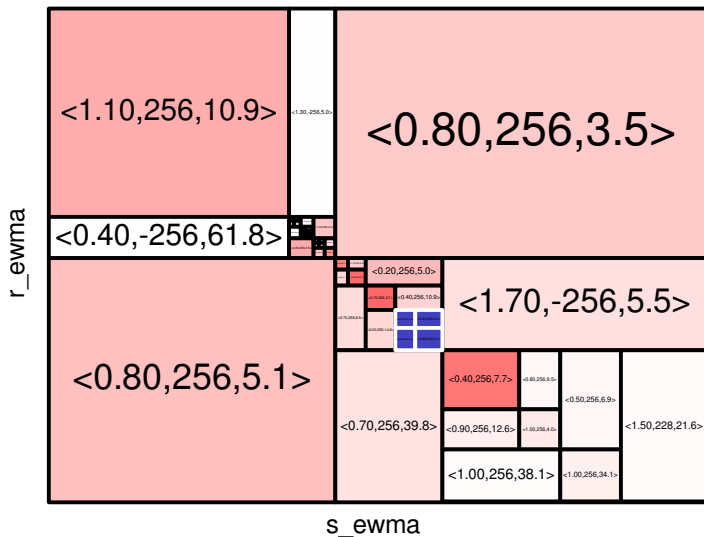


s_ewma

r_ewma

<0.90,5,2.8>  <0.70,6,53.5>

<0.60,19,76.2>  <0.80,5,4.1>  <0.80,5,4.1>

<0.80,5,4.1>  <0.80,5,4.1>

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
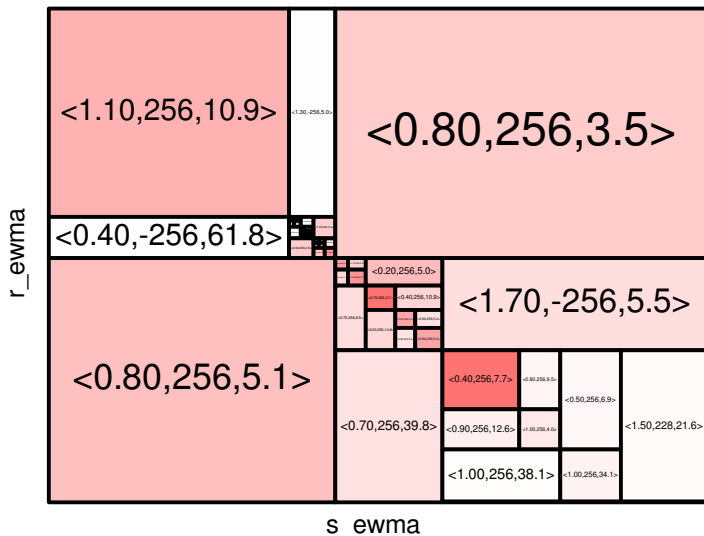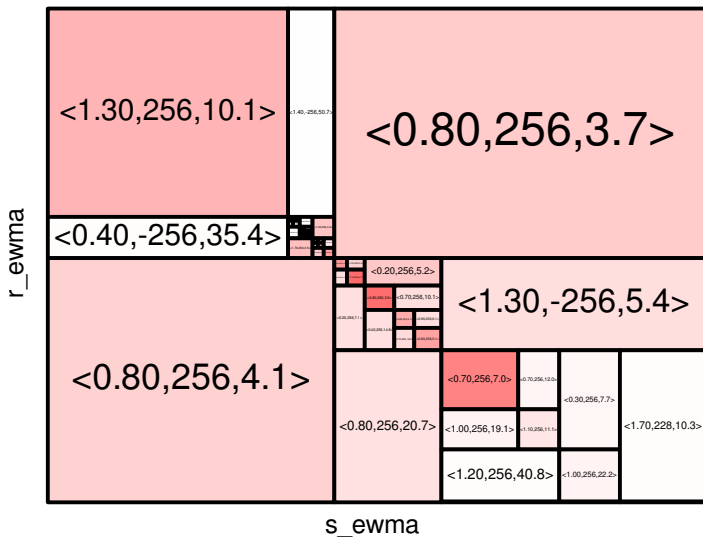
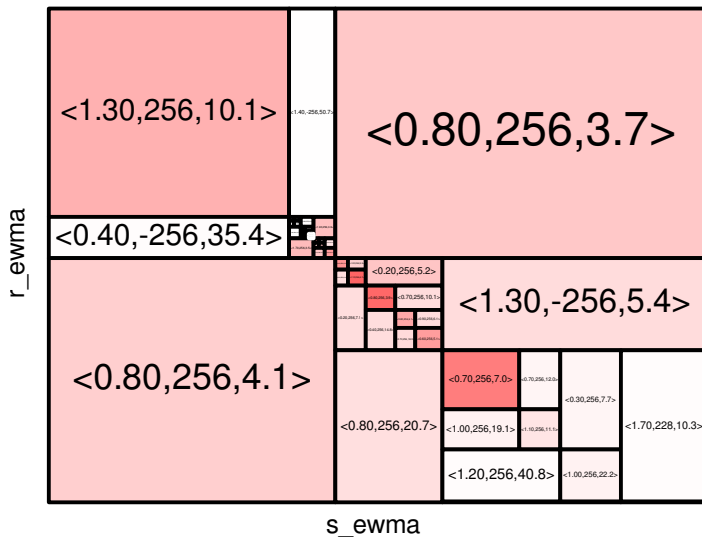TCP ex Machina: Computer-Generated Congestion Control

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

# Split

# Simulate

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate

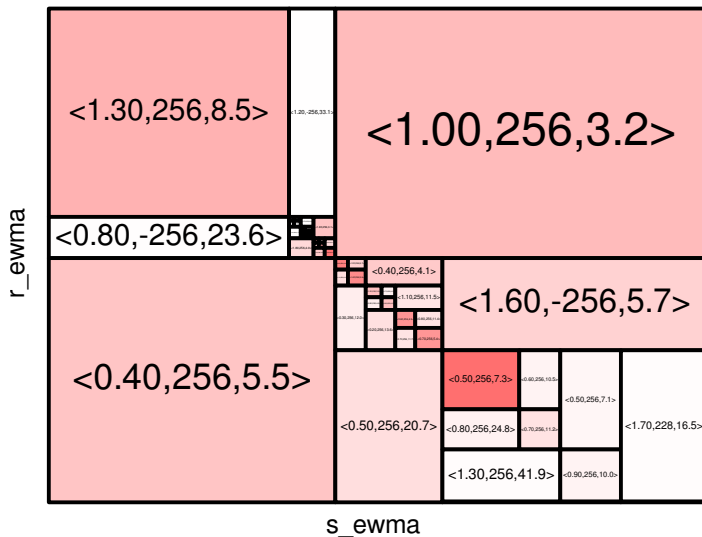

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

# Split

# Simulate



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize



r_ewma

s_ewma

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate



s_ewma

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize



s_ewma

r_ewma
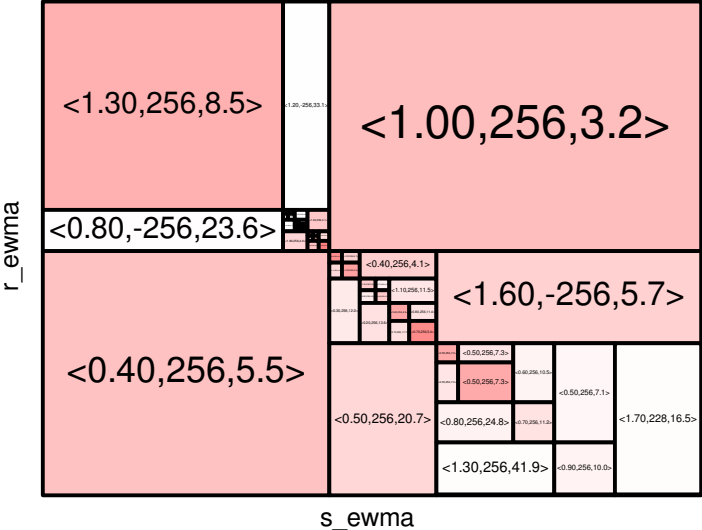
<1.90,256,8.6>
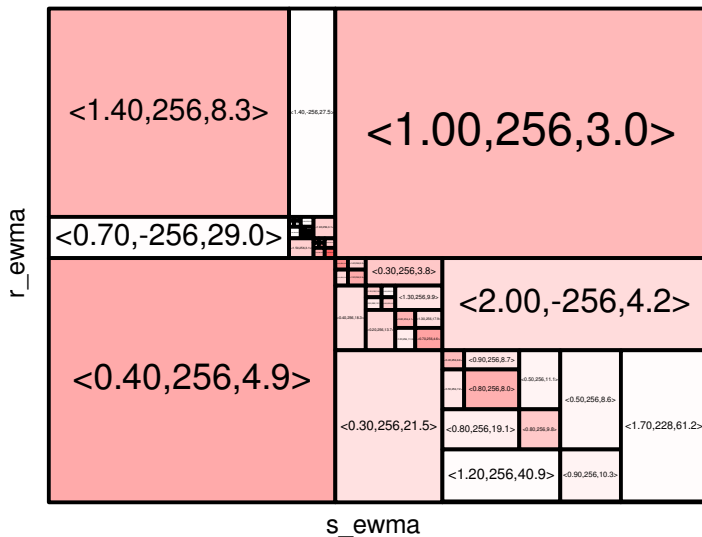<1.10,-256,6.7>
<0.40,256,5.3>
<0.80,-230,24.8>
<0.80,256,17.0>
<1.00,255,3.8>
<0.70,255,3.7>
<1.40,-256,4.2>
<1.80,256,50.6>
<0.90,254,13.1>
<0.80,227,7.1>
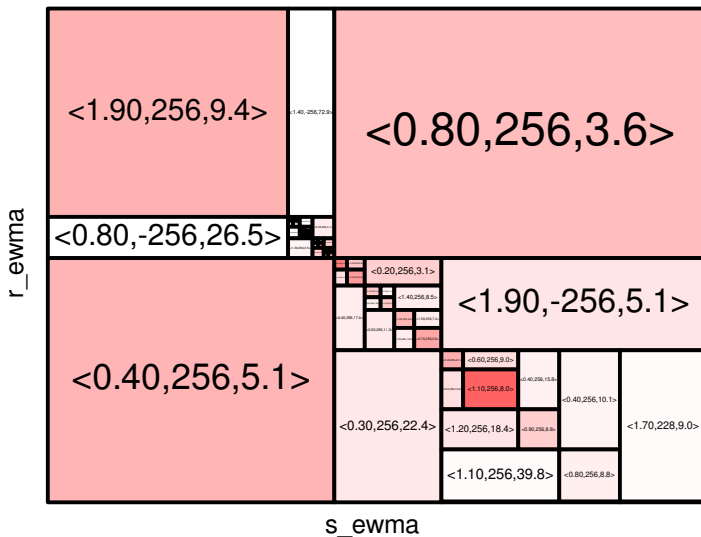
Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
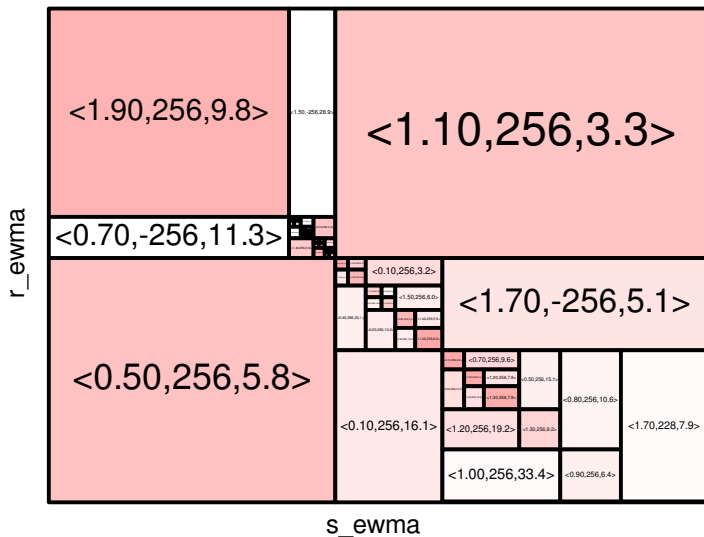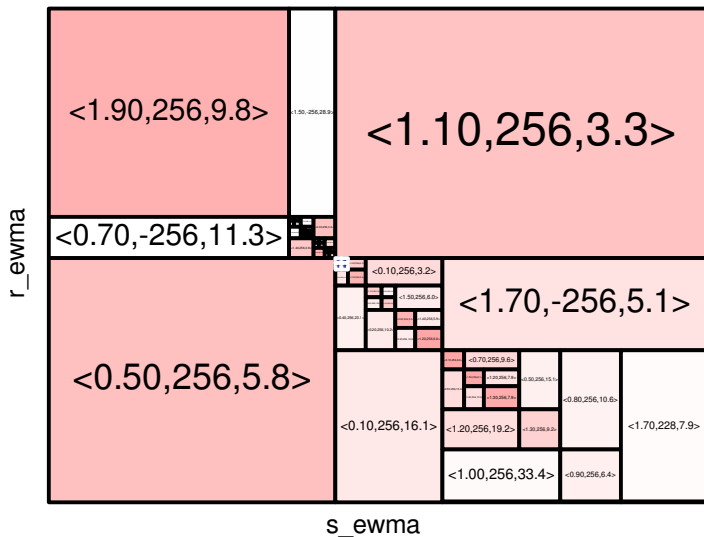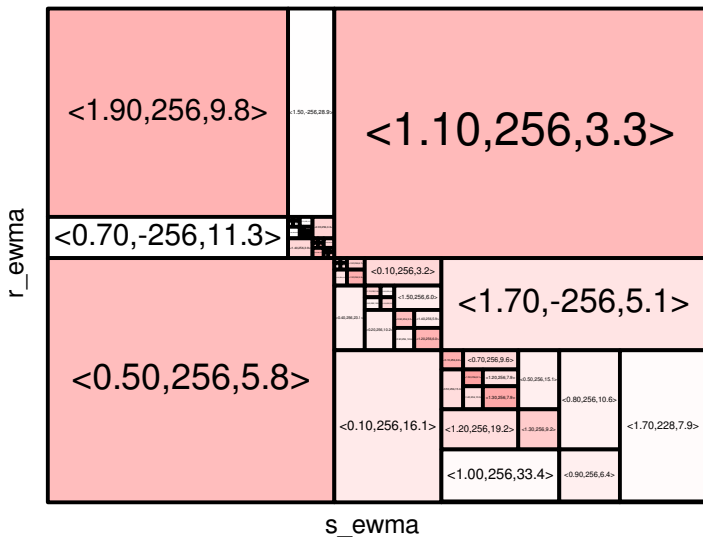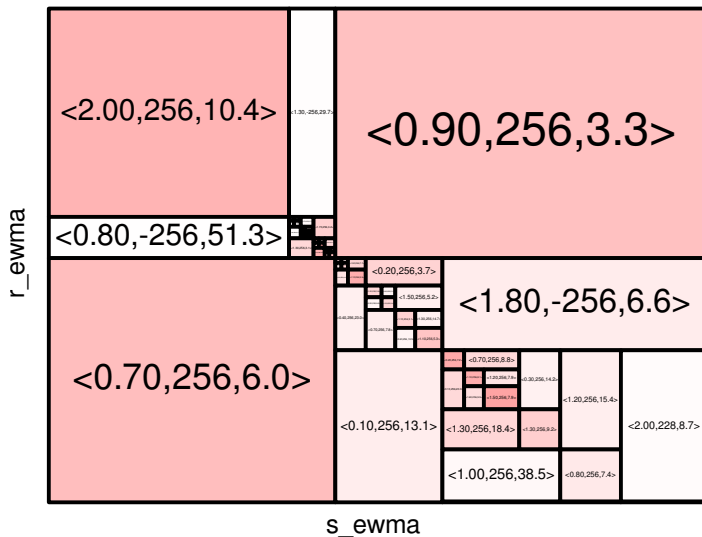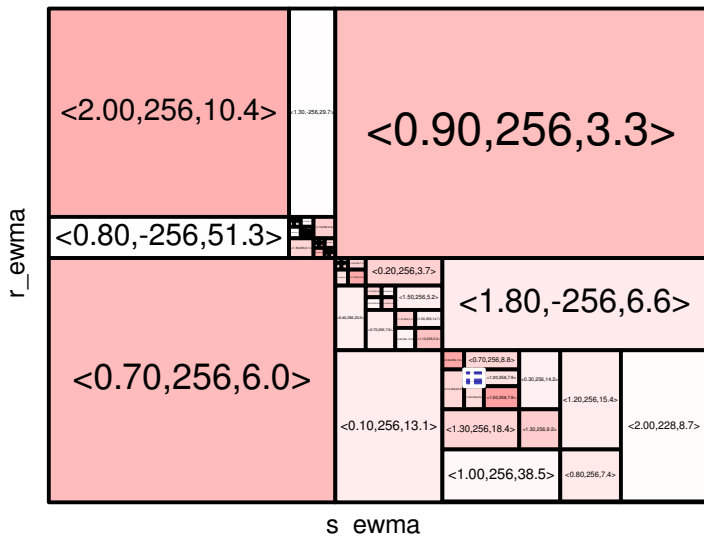
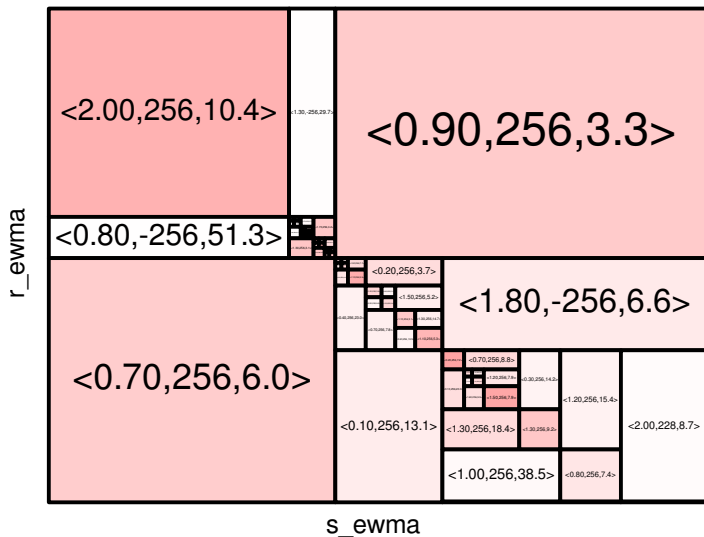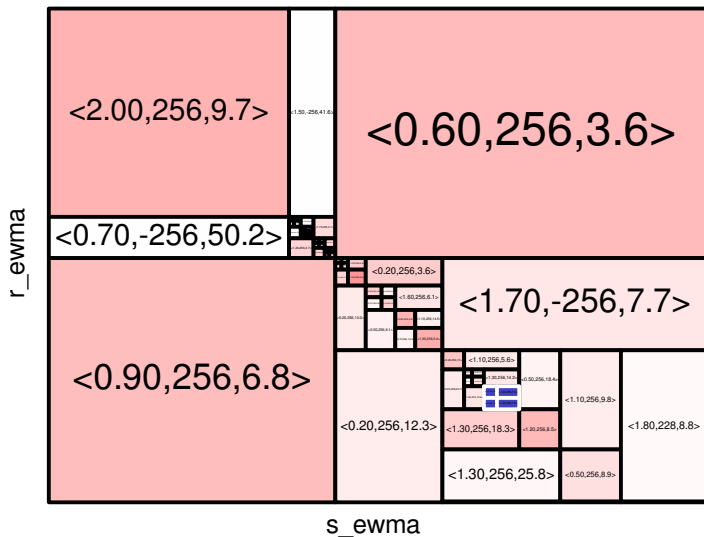TCP ex Machina: Computer-Generated Congestion Control
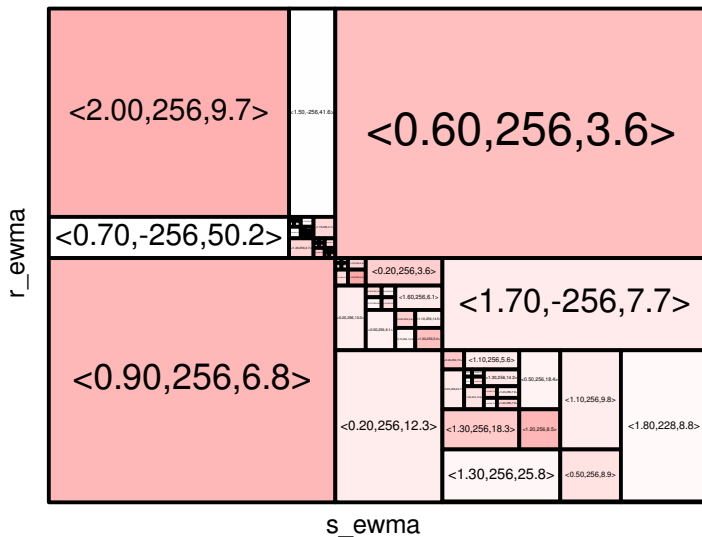
# Split

# Simulate

# Optimize



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate



s_ewma

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate



s_ewma

r_ewma

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split

# Simulate



s_ewma

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Optimize

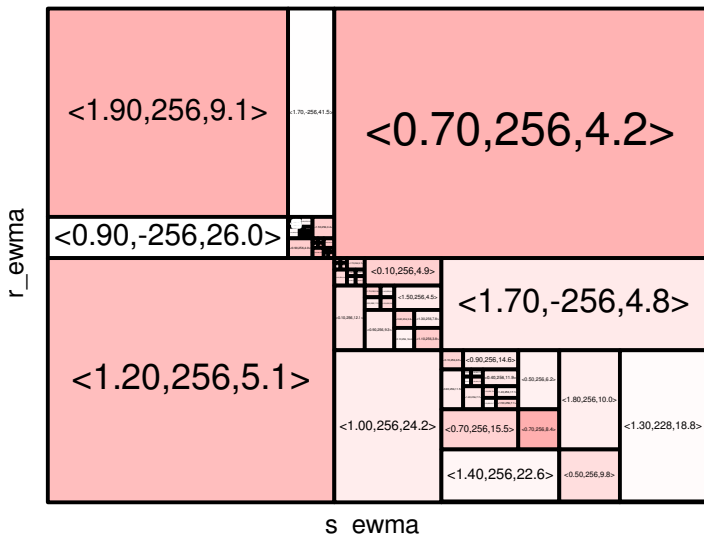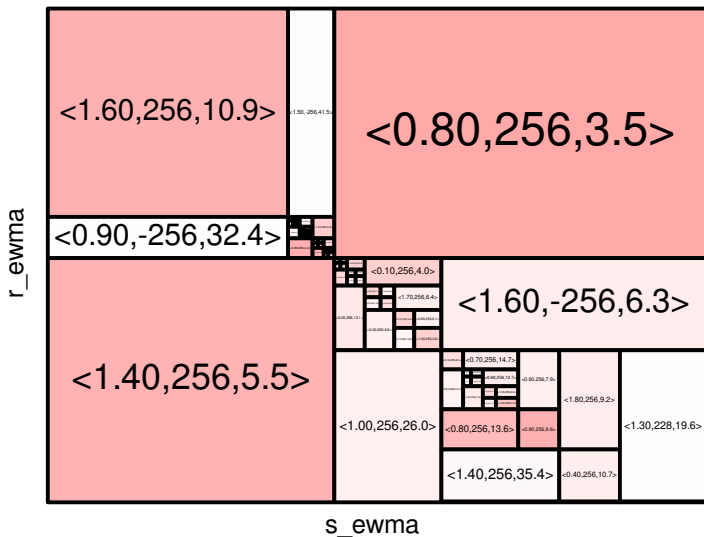Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

TCP ex Machina: Computer-Generated Congestion Control

# Simulate

# Optimize

# Split

# Simulate

# Optimize

# Split

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate

# Optimize

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Split



s_ewma

r_ewma

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Simulate



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
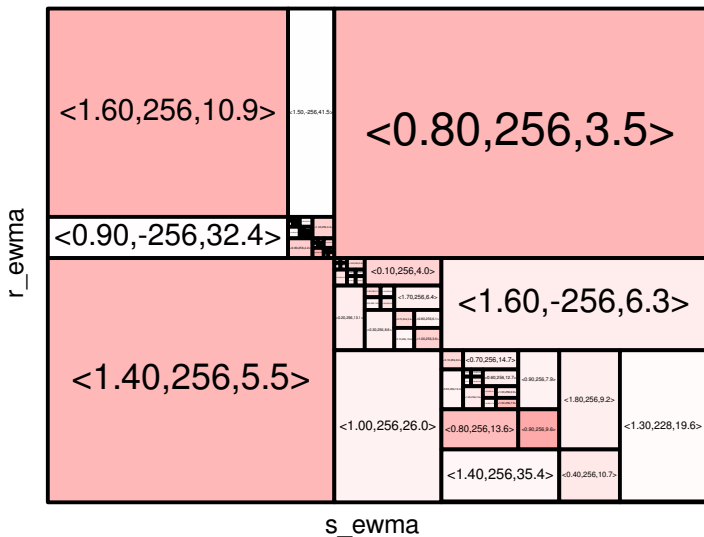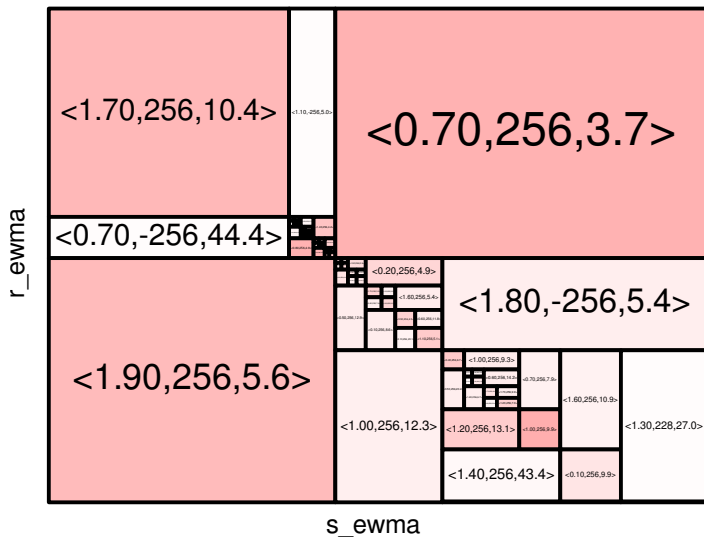
# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC                                                                    .

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

.



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC

.

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Evaluation in ns-2

- ▶ End-to-end comparators: NewReno, Cubic, Compound, Vegas
- ▶ In-net comparators: Cubic-over-sfqCoDel, XCP
- ▶ Simulation setup published for replication

# Scenario 1: fixed-rate network, homogenous senders

## Scenario 1: details

| Quantity | Simulation parameter | Remy assumptions |
|---|---|---|
| Link speed | 15 Mbps | Uniform(10, 20) Mbps |
| RTT | 150 ms | Uniform(100, 200) ms |
| $n$ | 8 | Uniform(1, 16) |
| "On" process | $\exp[\mu = 100]$ **kB** | $\exp[\mu = 5]$ **s** |
| "Off" process | $\exp\left[\mu = \frac{1}{2}\right]$ s | $\exp[\mu = \mathbf{5}]$ s |

**Remy objective:**

$$\sum_i \log\left[\frac{\text{throughput}_i}{(\text{delay}_i)^{\boldsymbol{\delta}}}\right]$$

$$\boldsymbol{\delta} \in \left\{\frac{1}{10}, 1, 10\right\}$$

# Scenario 1: throughput-delay plot

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

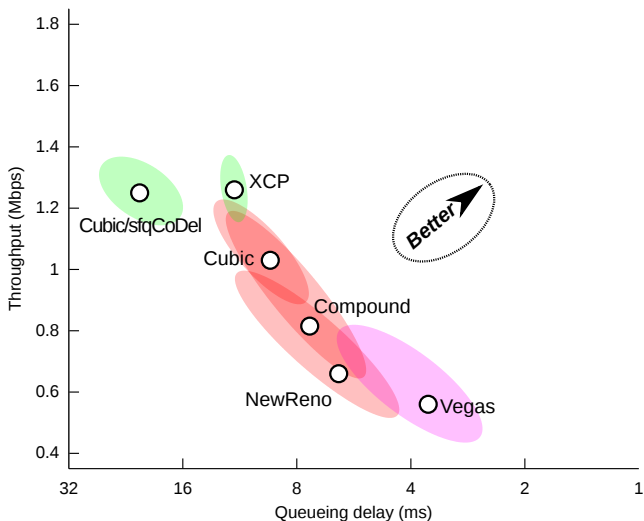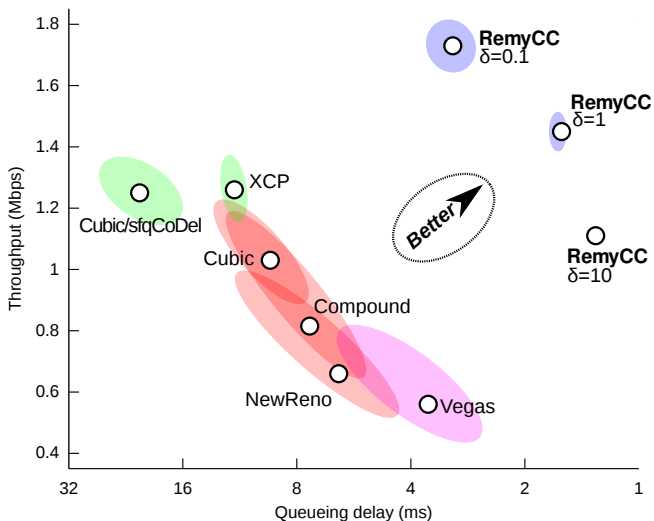TCP ex Machina: Computer-Generated Congestion Control

## Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot

# Scenario 1: throughput-delay plot
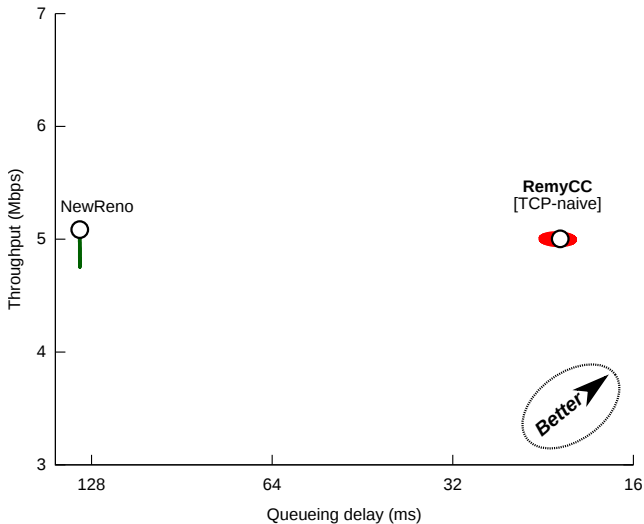
# Scenario 2: Verizon LTE, $n = 8$

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

*From the perspective of a network endpoint,*
**what matters?**

*How difficult is it to learn a good protocol,*
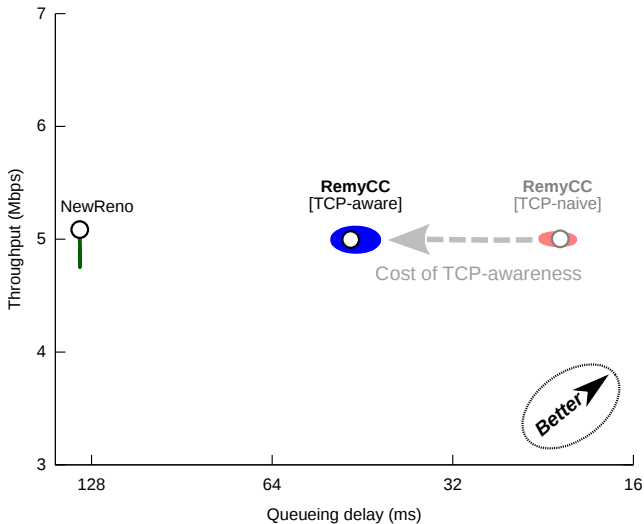*given an* **imperfect** *model of the network?*

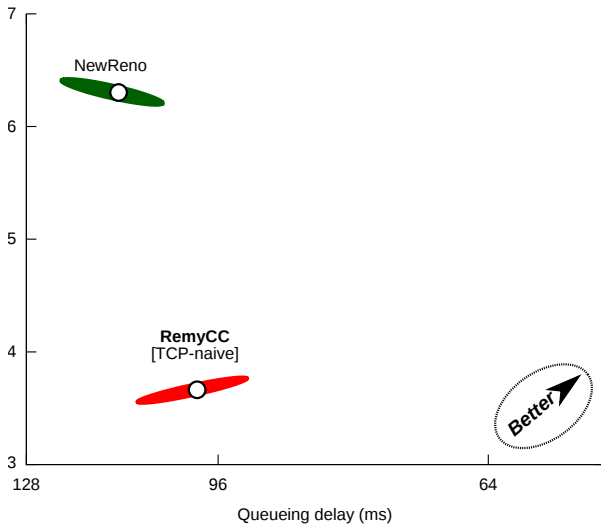Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC competing against itself

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

TCP ex Machina: Computer-Generated Congestion Control

# RemyCC competing against itself

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# RemyCC competing against itself

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

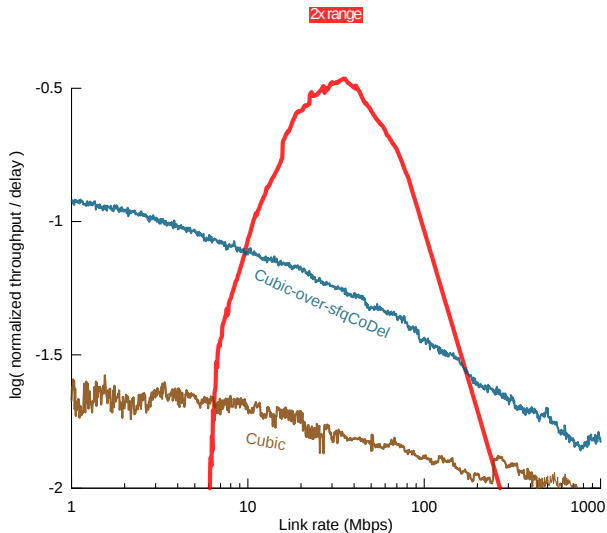# RemyCC competing against TCP NewReno

# RemyCC competing against TCP NewReno

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
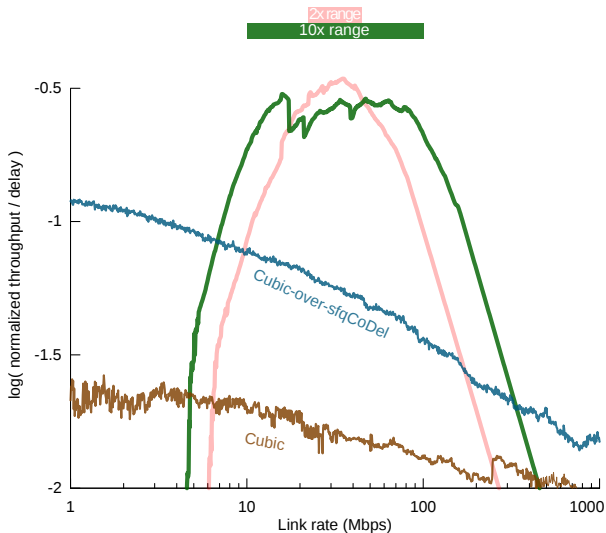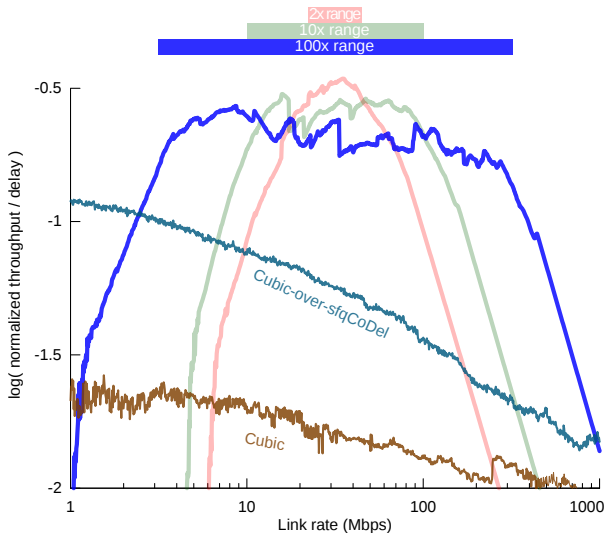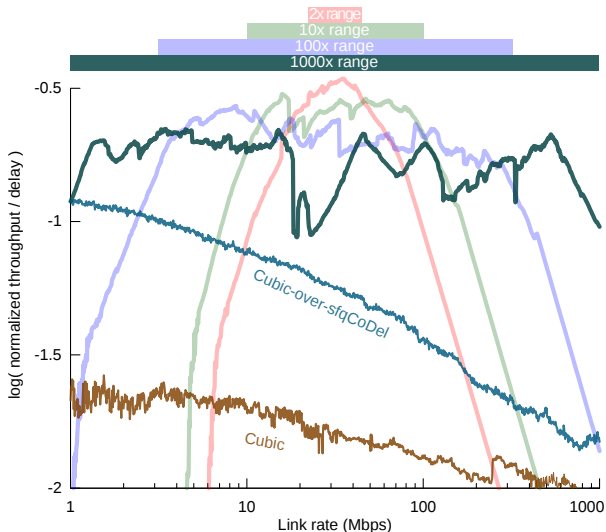
TCP ex Machina: Computer-Generated Congestion Control

# RemyCC competing against TCP NewReno

# The cost of generality

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The cost of generality

# The cost of generality

# The cost of generality



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

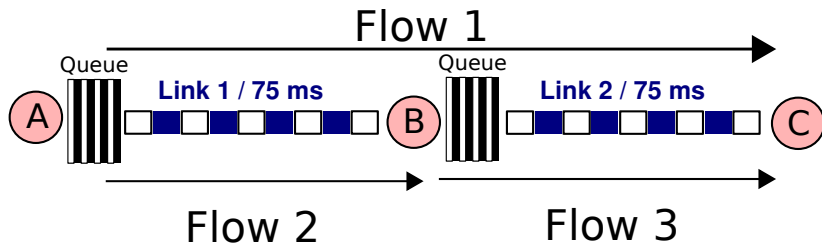# The cost of generality



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The cost of generality



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# The cost of generality



Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

When the model is wrong about the topology

One bottleneck

Flow 1

Queue

A    Link / 150 ms    B

Flow 2

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan
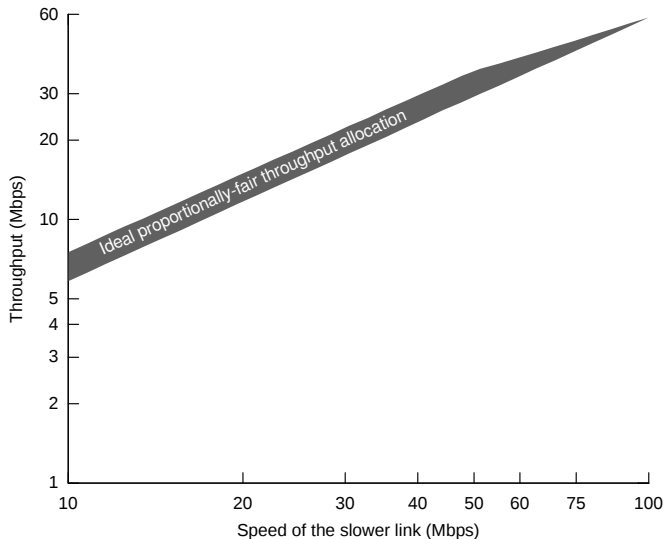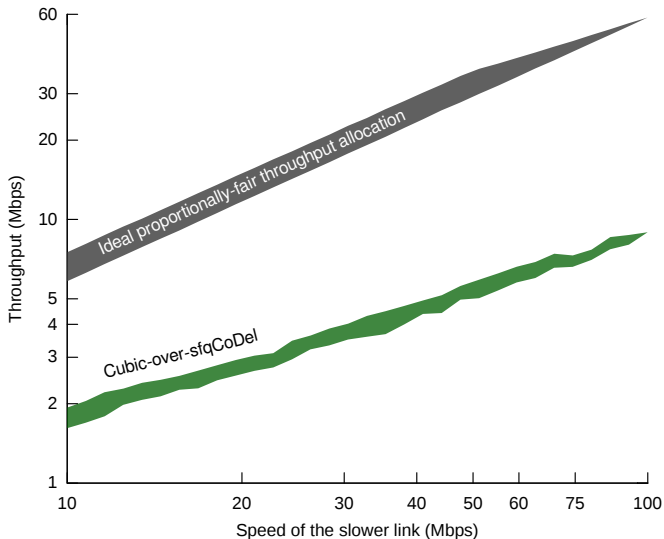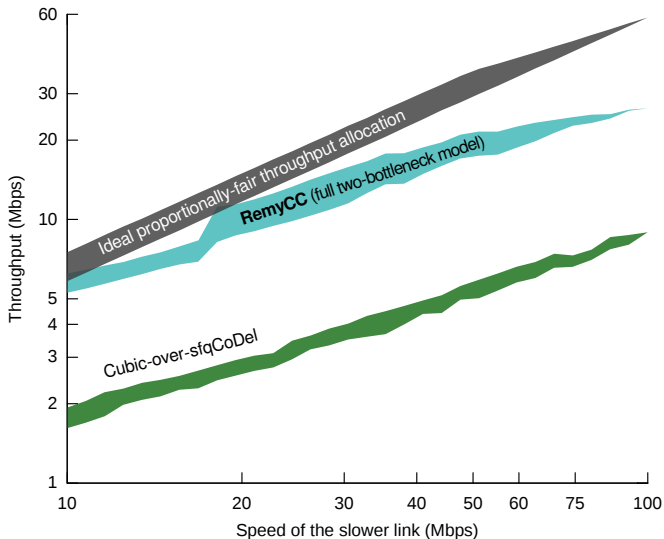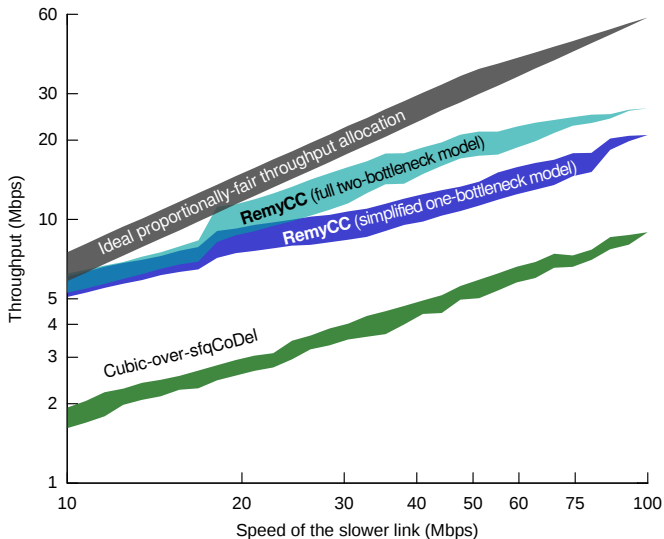
Two bottlenecks

# When the model is wrong about the topology

# When the model is wrong about the topology

# When the model is wrong about the topology

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan

# Systems ex Machina

- Explicit design considerations → **freedom to innovate**
- "If this protocol is the answer, what's the question?"

keithw@mit.edu

http://mit.edu/remy

Keith Winstein, Anirudh Sivaraman, Pratiksha Thaker, and Hari Balakrishnan