# Plan

* Feedback summary
* Recitation Qs
* Background on DCs
* Queue game
* DCTCP

## Logistics

* DP prelim report due TODAY 5pm.

* Participation check-in #2 out this weekend.

* Midterm on Tuesday, April 6
  ↳ example exams online

**1.** What is the goal of DCTCP?
- Increase perf for DCs (latency, ...)
- No new hardware

**2.** How does DCTCP differ from TCP?
- ECN — early congestion feedback

**3.** Why does DCTCP differ from TCP?
- Characteristics of DC traffic different than not traffic

# What makes a great research result?    (Spielman)
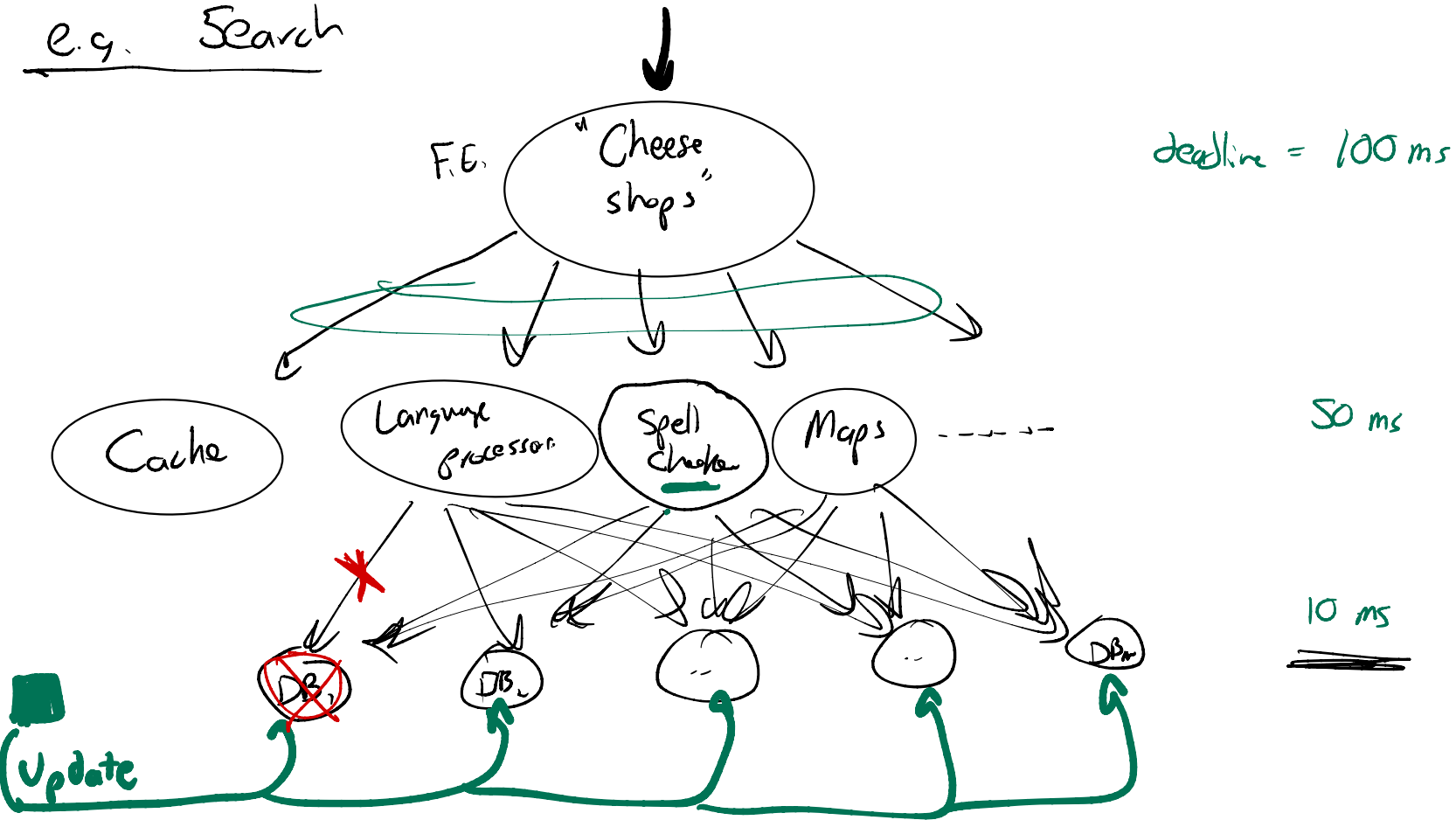
1. Beautiful theory.

2. Works in practice.

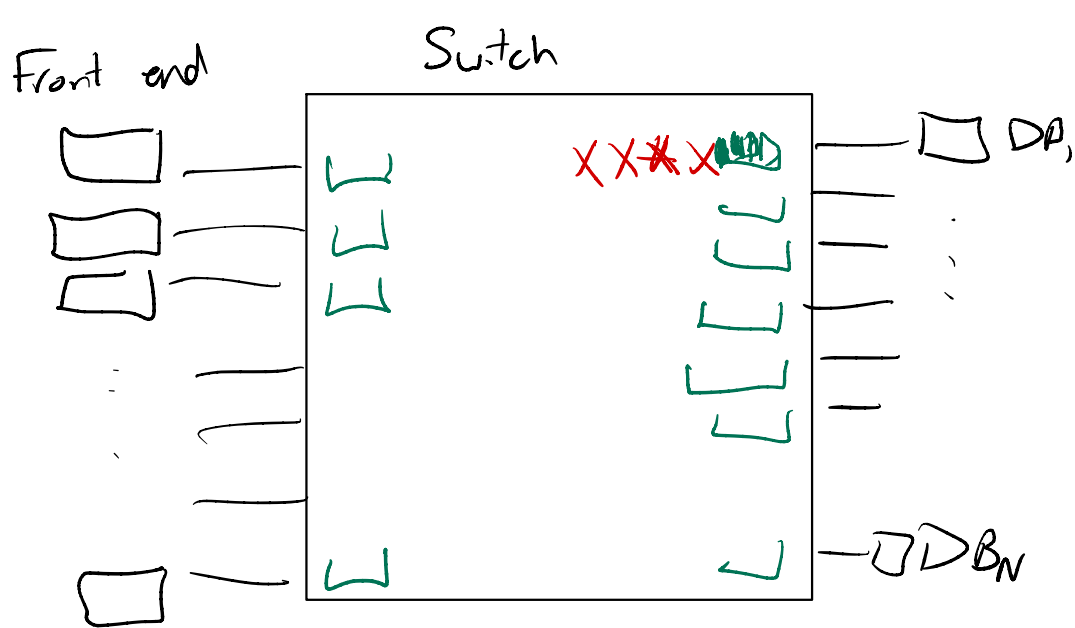3. Solves a problem that people care about.

# Types of flows in DC

1. Query flows      ~2 kB
2. Big traffic      ~50 MB
3. Short message    ~1 MB

throughput

latency sensitive

# e.g.  Search



"Cheese shops"

F.E.

deadline = 100 ms

Cache   Language processor.   Spell Checker   Maps   - - - - -   50 ms

DB   DB   DB   10 ms

Update

# Queues

Front end

Switch

DP₁

DBₙ

Two worries
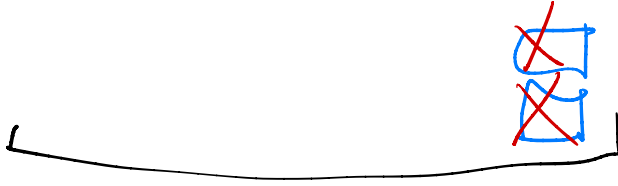
1. Dropped packets

2.

# Queue Game

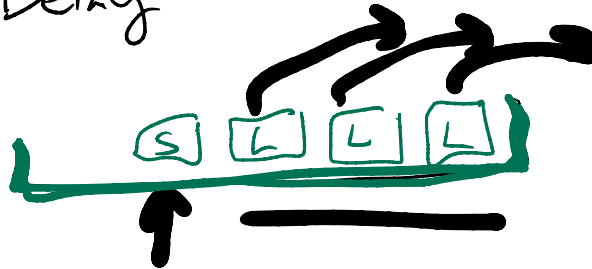**Front-ends**

US

Switch (Amir)

DB Server

1. One short flow  ✓

2. Many short flow same time  "Incast"
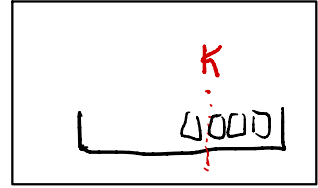
3. Queuing Delay

# What does DCTCP do to fix?

## SWITCH
~ Explicit congestion notification.
  → Set bit on Packet
  → Difference w/ prior art:
     Switch sets bit early

## SENDER
— IF start getting ECN → slow down gradually
  → Normal TCP → really slows down → sharply

Switch



K

4000

# Why doesn't this work on Internet?

- Need to modify both ends and switch ?

- Convergence time depends on RTT

- Feedback is too slow

# Plan

* Feedback summary
* Recitation Qs
* Background on DCs
* Queue game
* DCTCP

Logistics

* DP prelim report due
  TODAY 5pm.

* Participation check-in #2
  out this weekend.

* Midterm on Tuesday,
  April 6
  ↳ example exams online

# DCTCP

1. What is the goal of DCTCP?
   ↳ improve on TCP in DCs    low latency, high throughput

2. How does DCTCP differ from TCP?
   ↳ More clever use of explicit congestion notification
                                                    (ECN)

3. Why does DCTCP differ from TCP?
   ↳ Take advantage of unique properties of DC
   {  0.1 ms  RTT
      ~50 ms  RTT

# What makes a great research result? (Spielman)

1. Beautiful theory

2. Works in practice

3. Solves a problem that people care about.

# Types of Flows

1. Short/low-latency "Query" ~2 kB
2. Long .... background flows ~(00 MB)
   w/ MB
3. Short msg traffic

E.g. Search

Front end

"Cheese Shops"

deadline = 100ms   "Incast"

Cache

Language processing

Spell Checker

Maps   = 50ms

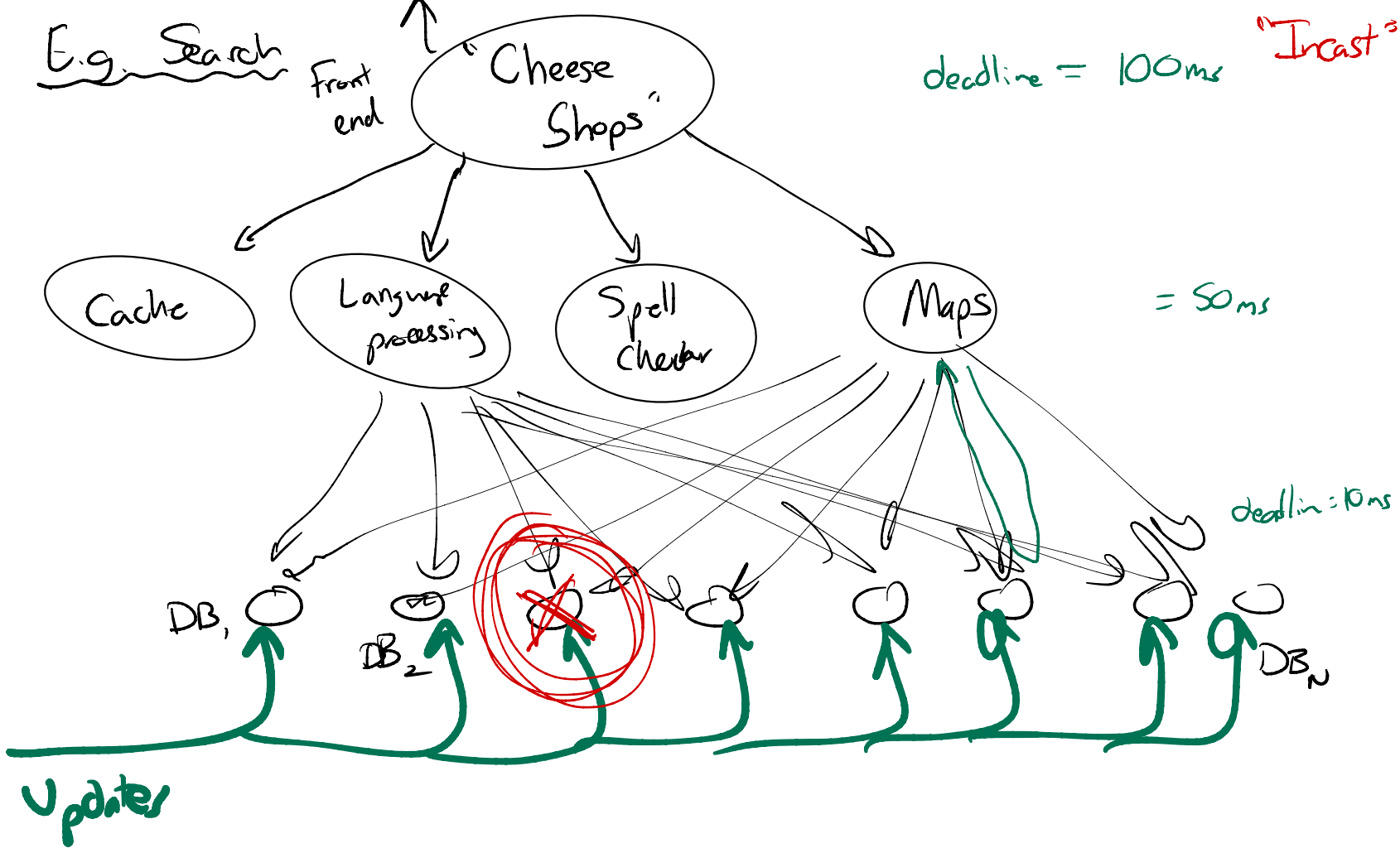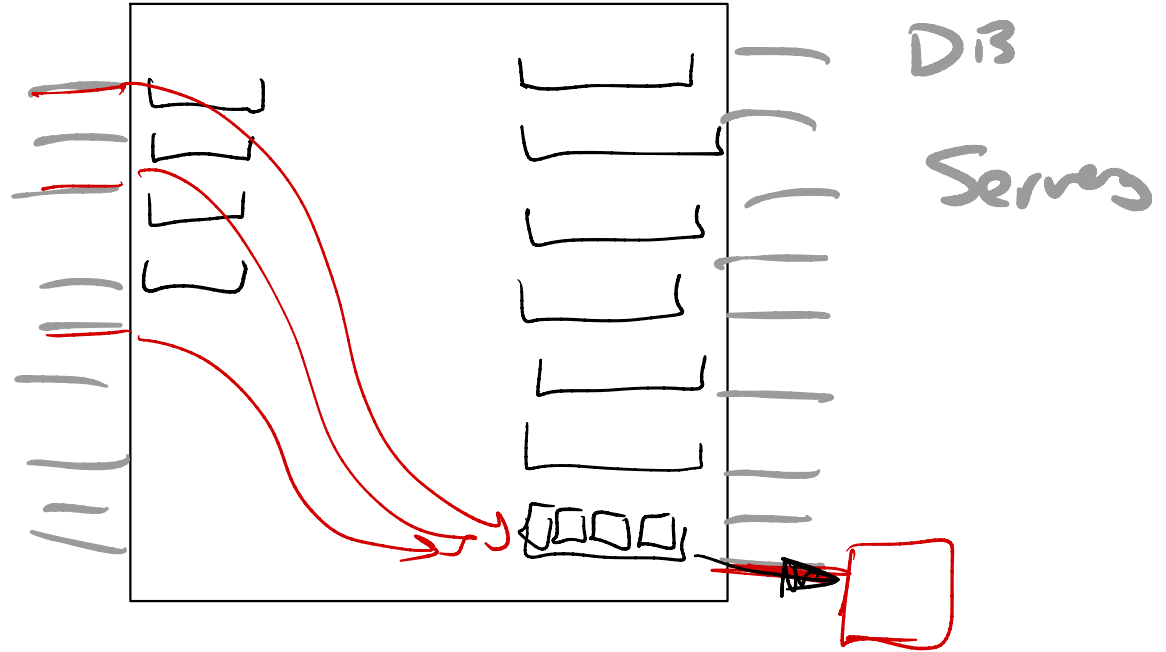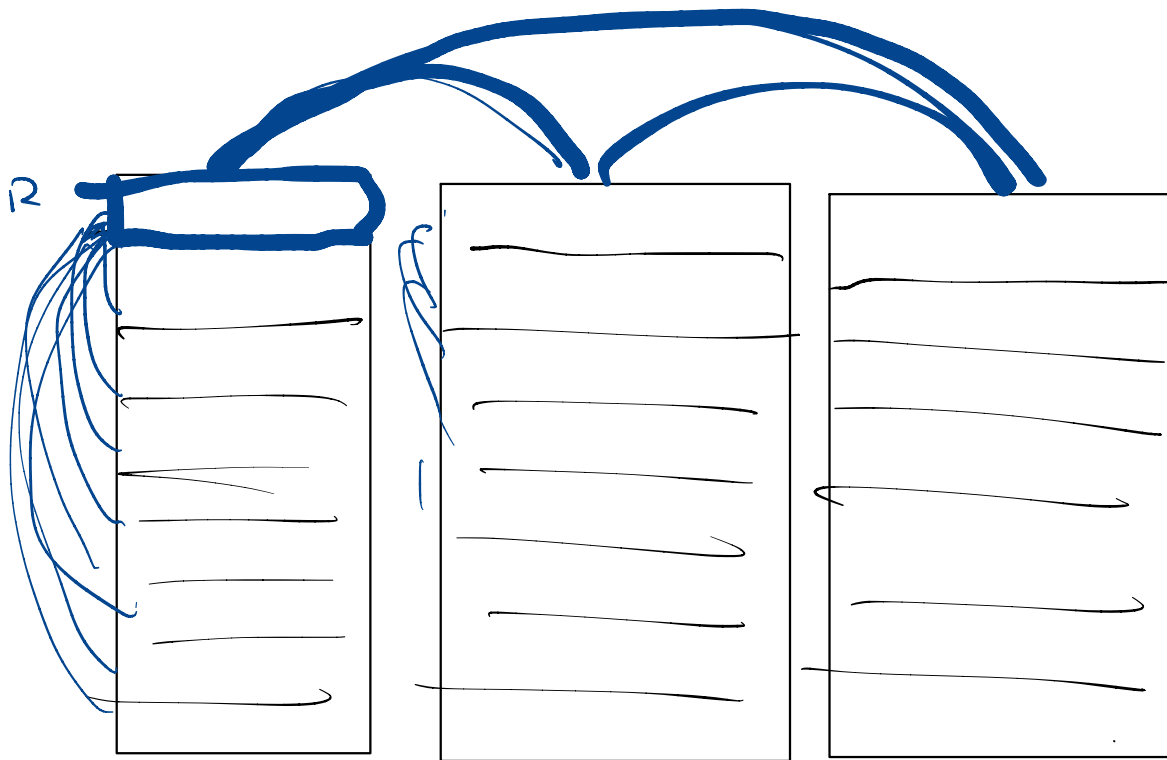deadlin = 10ms

$DB_1$

$DB_2$

$DB_N$

Updates

# Queues



Front-end servers

DB Servers

ToR

# Queue Game
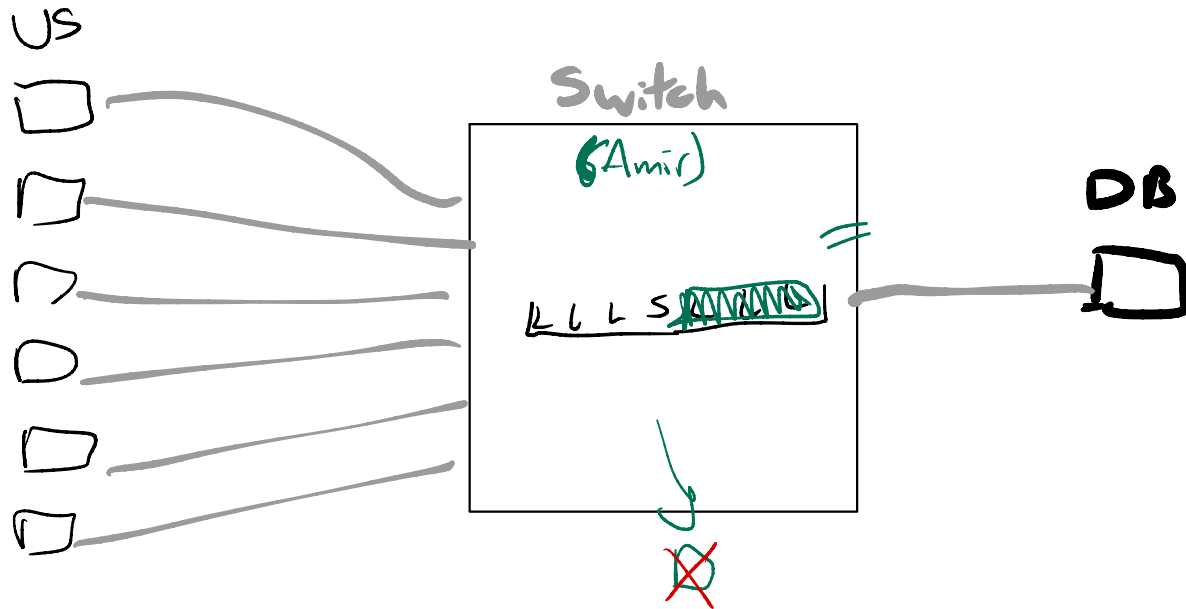
US



Switch
(Amir)

L L L S ~~MMMM~~

DB

1. Drop packet
2. Queuing delay
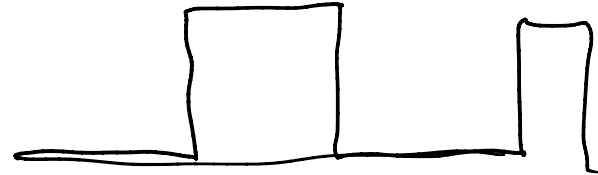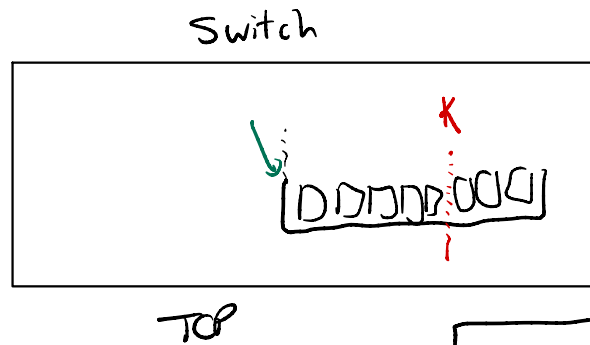
# DCTCP

## Switch



## SWITCH

— Early notification of congestion.

Normal: Wait under buffer full

TCP



DCTCP



## SENDER

— Keeps estimate of queue length
↳ adjusts tx rate accordingly

How is this different from normal TCP?
↳ TCP: Cut window size in half
↳ Quantitative approach

# Why doesn't this work on Internet?

- Scale .... Convergence time = time before every endpoint tx
  at "right" rate

  ↳ Depends on RTT    DC   0.1 ms
                            Net   50 ms

- Deployment

- Less structure ∇₀

-