# Understanding Network Failures in Data Centers: Measurement, Analysis and Implications

**Phillipa Gill**
University of Toronto

Navendu Jain  & Nachiappan Nagappan
Microsoft Research

# Motivation

## Amazon: Networking Error Caused Cloud Outage

April 29th, 2011 : Rich Miller

Last week's lengthy outage for the **Amazon Web Services** cloud computing platform was caused by a network configuration error as Amazon was attempting to upgrade capacity on its network. That error triggered a sequence of events that culminated in a "re-mirroring storm" in which automated replication of storage volumes maxed out the capacity of Amazon's servers in a portion of their platform.

# Motivation



**$5,600 per minute**

We need to understand failures to prevent and mitigate them!

# Overview

**Our goal:** Improve reliability by understanding network failures

1. Failure **characterization**
   - Most failure prone components
   - Understanding root cause
2. What is the **impact** of failure?
3. Is **redundancy** effective?

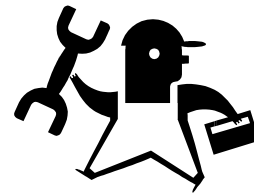**Our contribution:** First large-scale empirical study of network failures across multiple DCs

- Methodology to extract failures from noisy data sources.
- Correlate events with network traffic to estimate **impact**
- Analyzing implications for future data center networks

# Road Map

Motivation

**Background & Methodology**

Results

1. Characterizing failures

2. Do current network redundancy strategies help?

Conclusions

# Data center networks overview



Internet

Access routers/network "core" fabric

Load balancers

Aggregation "Agg" switch

Top of Rack (ToR) switch

Servers

# Data center networks overview



Which components are most failure prone?

Internet

How effective is redundancy?
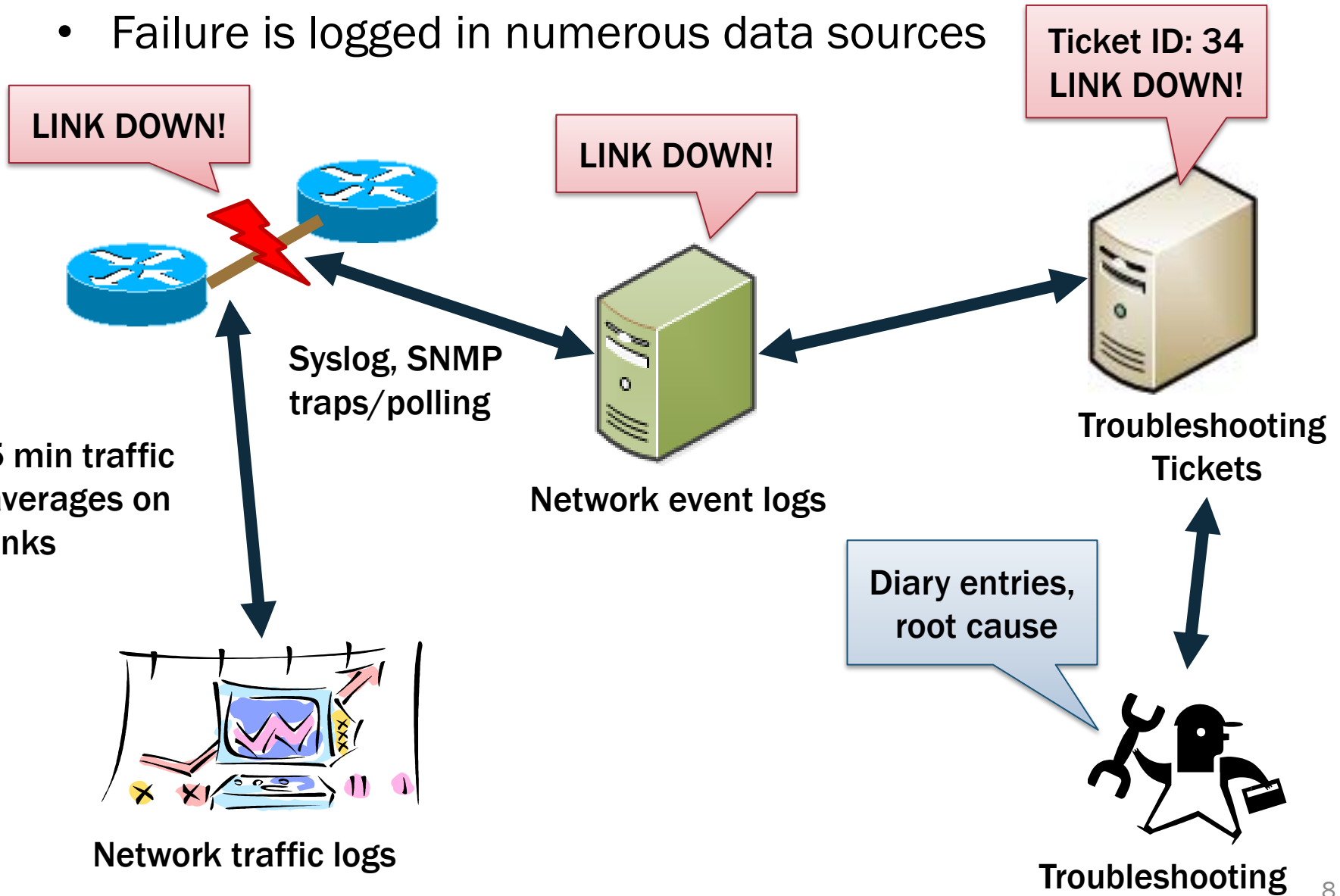
What is the impact of failure?

What causes failures?

# Failure event information flow

- Failure is logged in numerous data sources

**Ticket ID: 34 LINK DOWN!**

**LINK DOWN!**

**LINK DOWN!**

**Syslog, SNMP traps/polling**

**5 min traffic averages on links**

**Network event logs**

**Troubleshooting Tickets**

**Diary entries, root cause**

**Network traffic logs**

**Troubleshooting**

# Data summary

- One year of event logs from Oct. 2009-Sept. 2010
  - Network event logs and troubleshooting tickets

- Network event logs are a combination of Syslog,  SNMP traps and polling
  - Caveat: may miss some events e.g., UDP, correlated faults

- Filtered by operators to *actionable* events
  - … still many warnings from various software daemons running

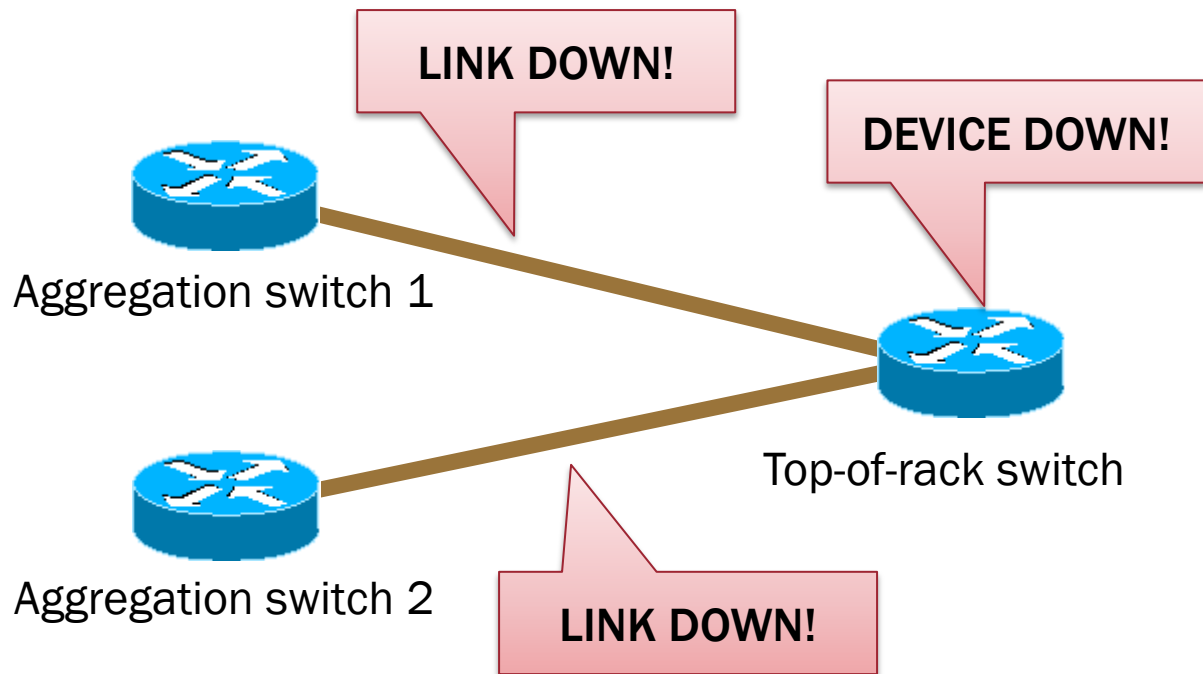**Key challenge: How to extract failures of interest?**

# Extracting failures from event logs

- **Defining failures**
  - **Device failure:** device is no longer forwarding traffic.
  - **Link failure:** connection between two interfaces is down. Detected by monitoring interface state.

- **Dealing with inconsistent data:**

  - **Devices:**
    - Correlate with link failures
  - **Links:**
    - Reconstruct state from logged messages
    - Correlate with network traffic to determine impact
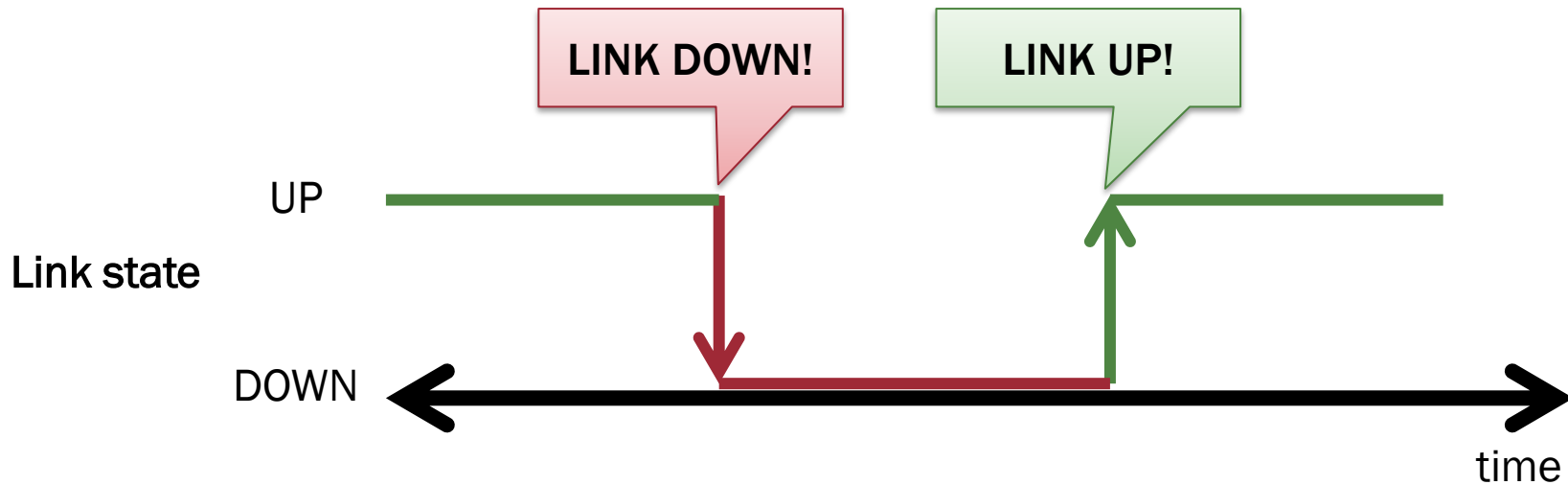
# Reconstructing device state

- Devices may send spurious DOWN messages
- Verify **at least one** link on device fails within five minutes
  - Conservative to account for message loss (correlated failures)

LINK DOWN!

DEVICE DOWN!

Aggregation switch 1

Top-of-rack switch

Aggregation switch 2

LINK DOWN!

**This sanity check reduces device failures by 10x**

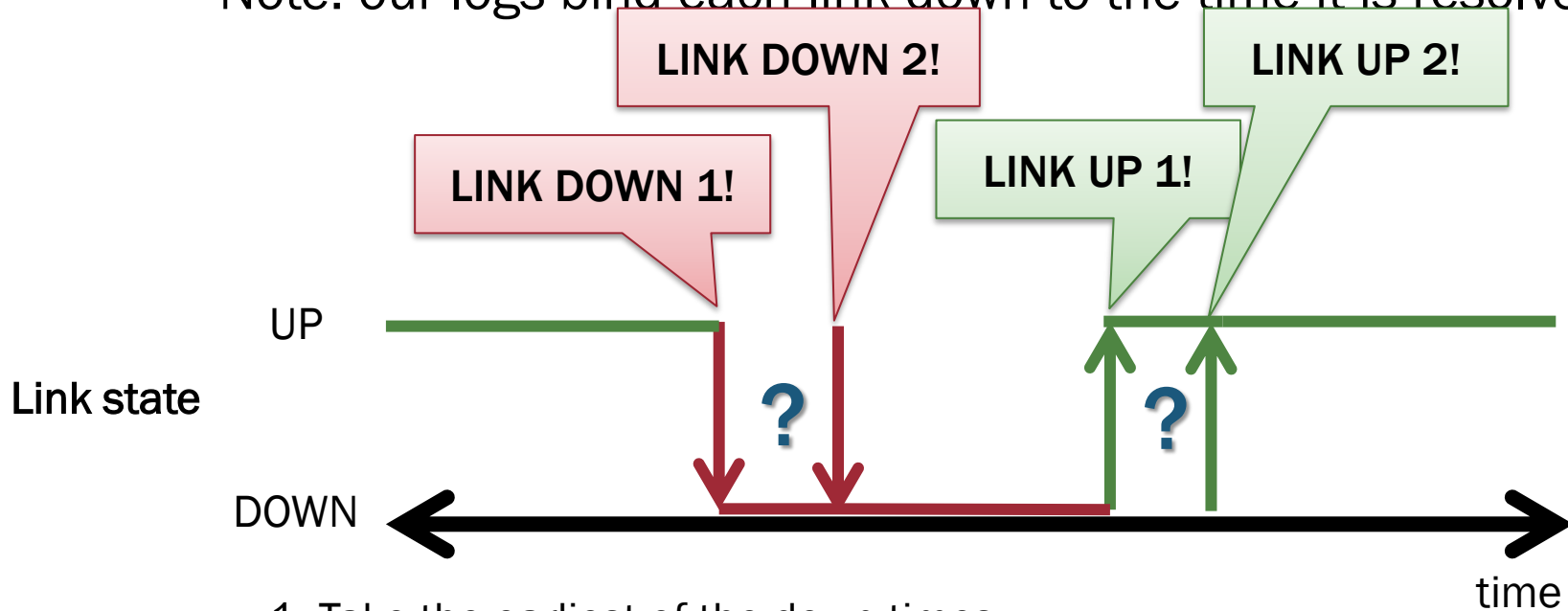# Reconstructing link state

- Inconsistencies in link failure events
  - Note: our logs bind each link down to the time it is resolved



**What we expect**

# Reconstructing link state

- Inconsistencies in link failure events
  - Note: our logs bind each link down to the time it is resolved
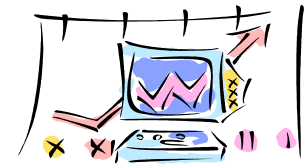


LINK DOWN 2!

LINK UP 2!

LINK DOWN 1!

LINK UP 1!

UP

Link state

?                    ?

DOWN

time

1. Take the earliest of the down times

2. Take the earliest of the up times

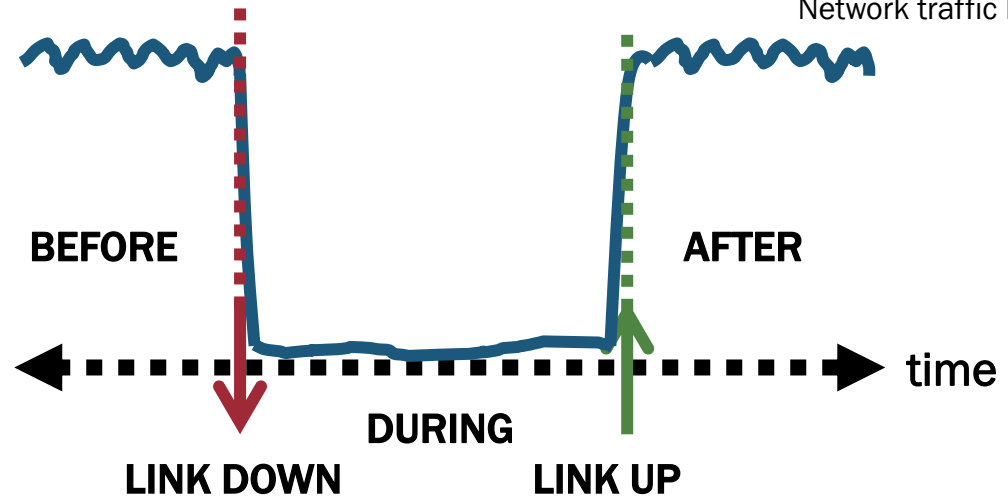**How to deal with discrepancies?**

# Identifying failures with impact

Network traffic logs

**Correlate link failures with network traffic**

**Only consider events where traffic decreases**

$$\frac{traffic\ during}{traffic\ before} < 1$$

BEFORE          AFTER

time

DURING

LINK DOWN          LINK UP

- **Summary of impact:**
  - 28.6% of failures impact network traffic
  - 41.2% of failures were on links carrying **no traffic**
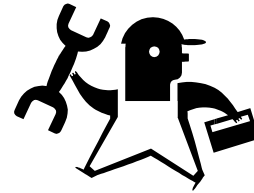    - E.g., scheduled maintenance activities
- **Caveat:** Impact is only on network traffic **not necessarily applications!**
  - Redundancy: Network, compute, storage mask outages

# Road Map

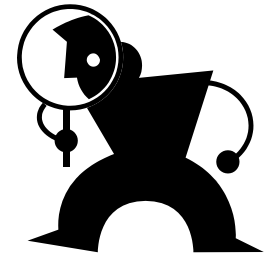Motivation

**Background & Methodology**
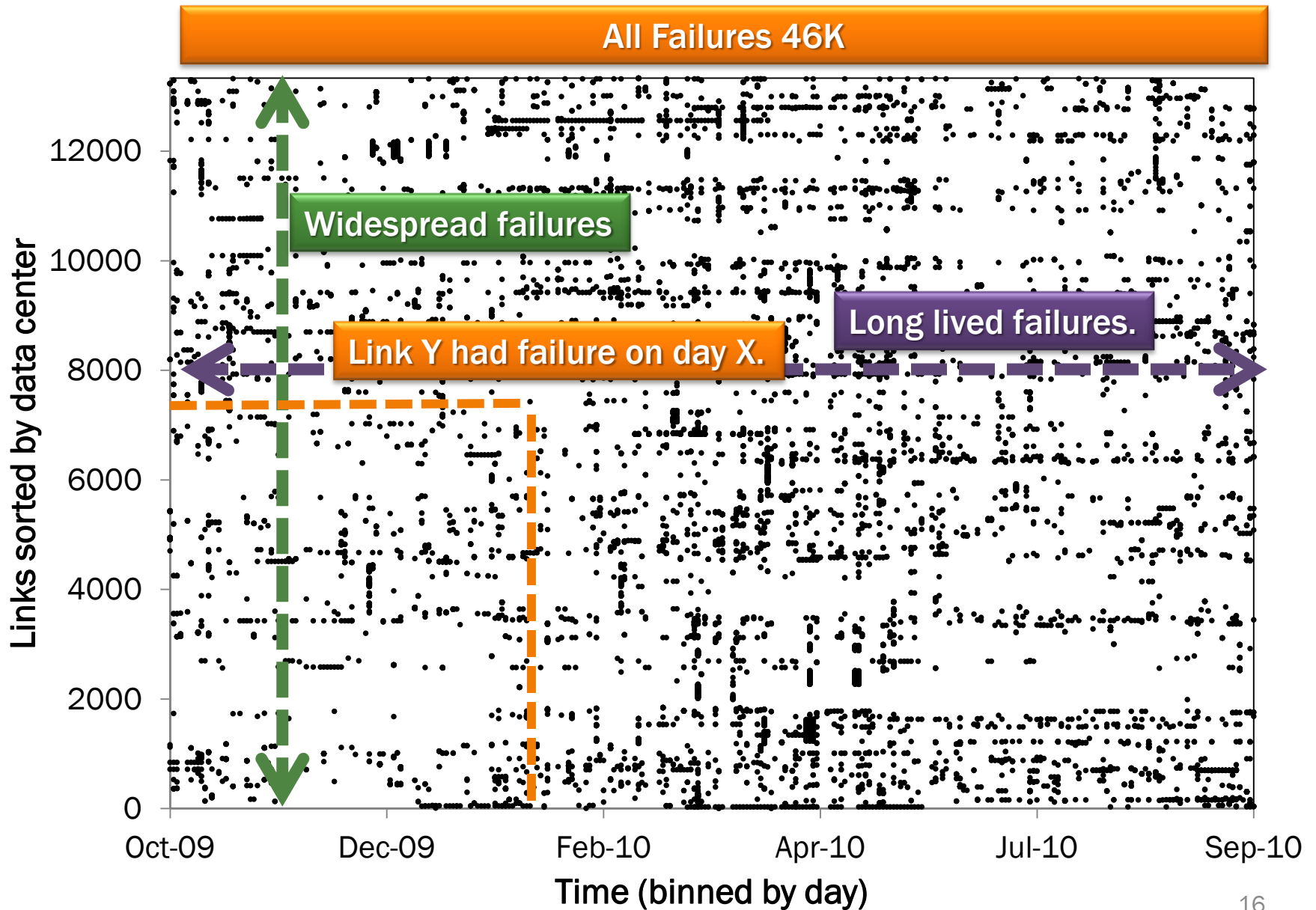
## Results

1. **Characterizing failures**
   – Distribution of failures over measurement period.
   – Which components fail most?
   – How long do failures take to mitigate?

2. Do current network redundancy strategies help?

Conclusions

# Visualization of failure panorama: Sep'09 to Sep'10



All Failures 46K

Widespread failures

Long lived failures.

Link Y had failure on day X.

Links sorted by data center

Time (binned by day)

# Visualization of failure panorama: Sep'09 to Sep'10



Failures with Impact 28%

Component failure: link failures on multiple ports

Load balancer update (multiple data centers)

Links sorted by data center

Time (binned by day)

# Which devices cause most failures?

# Which devices cause most failures?



Top of rack switches have few failures...
(annual prob. of failure <5%)

...but a lot of downtime!

Load balancer 1: very little downtime relative to number of failures.

Legend:
- failures
- downtime

Data:
- Load Balancer 1: failures 38%, downtime 2%
- Load Balancer 2: failures 28%, downtime 18%
- Top of Rack 1: failures 15%, downtime 66%
- Load Balancer 3: failures 9%, downtime 5%
- Top of Rack 2: failures 4%, downtime 8%
- Aggregation Switch: failures 4%, downtime 0.4%

Y-axis: Percentage (0% to 100%)
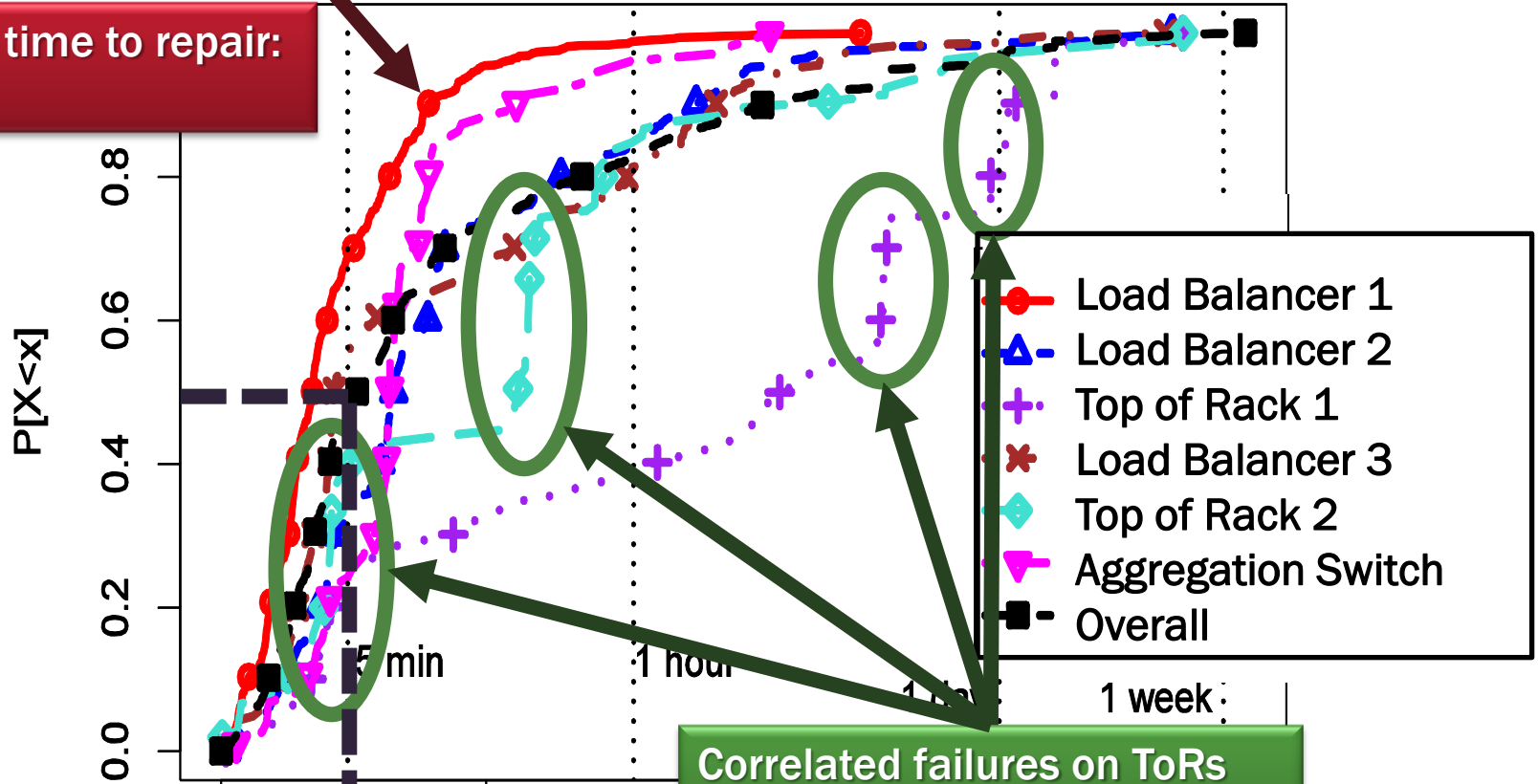X-axis: Device Type

# How long do failures take to resolve?

# How long do failures take to resolve?



**Load balancer 1:** short-lived *transient* faults

**Median time to repair: 4 mins**

**Median time to repair: 5 minutes**
**Mean: 2.7 hours**

**Correlated failures on ToRs connected to the same Aggs.**

**Median time to repair:**
**ToR-1: 3.6 hrs**
**ToR-2: 22 min**

Legend:
- Load Balancer 1
- Load Balancer 2
- Top of Rack 1
- Load Balancer 3
- Top of Rack 2
- Aggregation Switch
- Overall

Axis labels:
- P[X<x]
- Time to repair (s)
- 5 min
- 1 hour
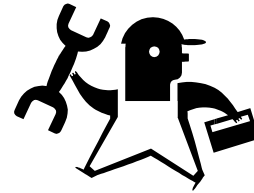- 1 week
- 1e+03
- 1e+04
- 1e+05
- 1e+06

# Summary

- Data center networks are highly reliable
  - Majority of components have four 9's of reliability


- Low-cost top of rack switches have highest reliability
  - <5% probability of failure
- …but most downtime
  - Because they are lower priority component


- Load balancers experience many short lived faults
  - Root cause: software bugs, configuration errors and hardware faults


- Software and hardware faults dominate failures
  - …but hardware faults contribute most downtime

# Road Map
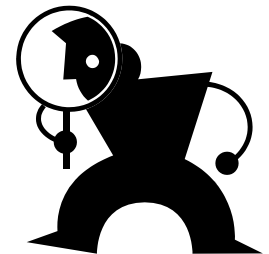
Motivation

Background & Methodology

Results

1. Characterizing failures
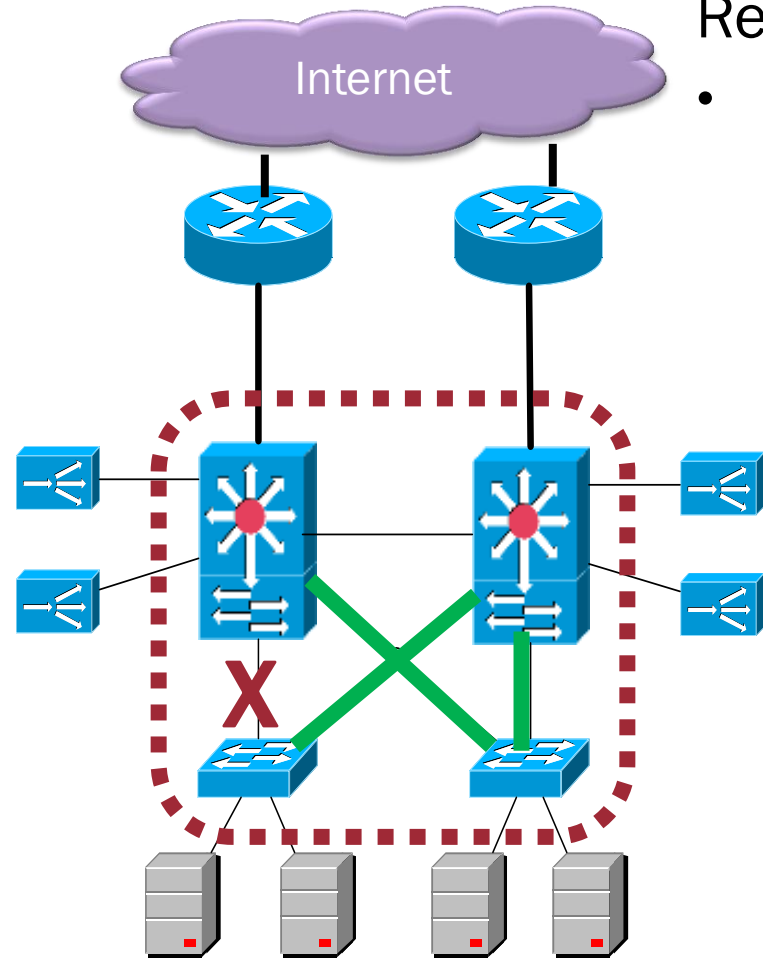
**2. Do current network redundancy strategies help?**

Conclusions

# Is redundancy effective in reducing impact?



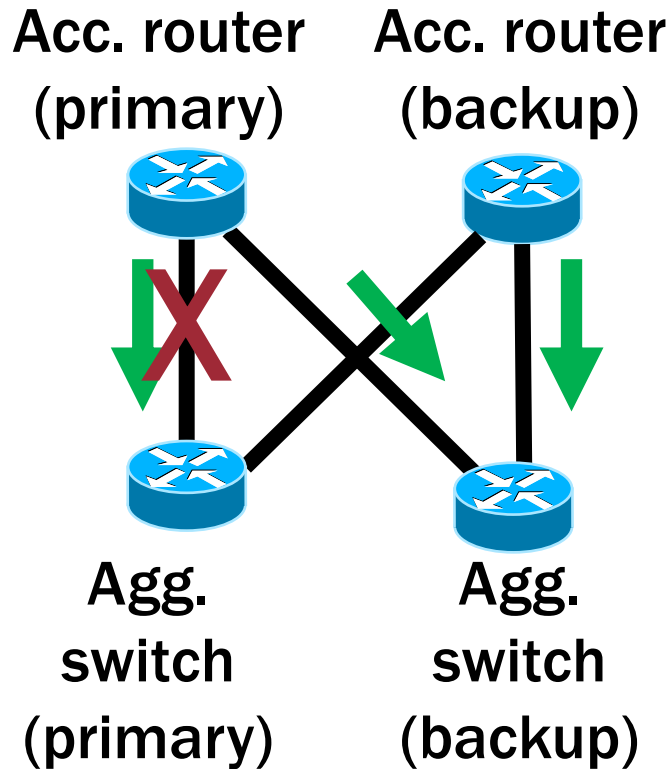Internet

**Redundant devices/links to mask failures**

- This is expensive! (management overhead + $$$)

**Goal:** Reroute traffic along available paths

**How effective is this in practice?**

# Measuring the effectiveness of redundancy

Acc. router (primary)    Acc. router (backup)

Agg. switch (primary)    Agg. switch (backup)

**Idea:** compare traffic before and during failure
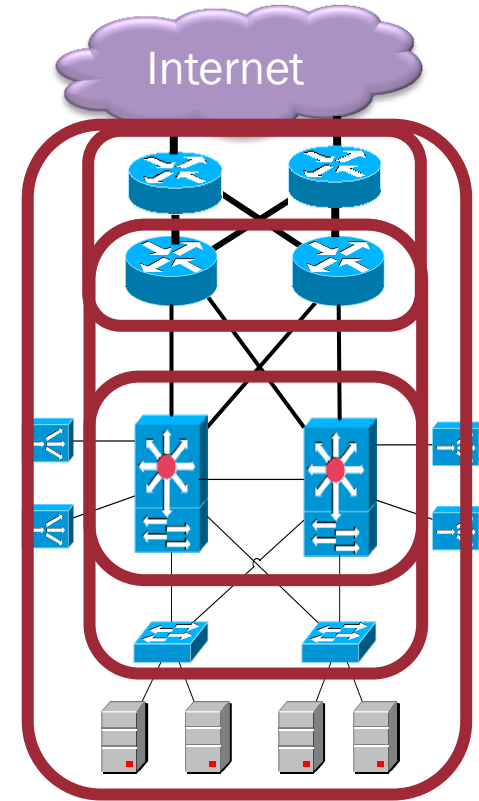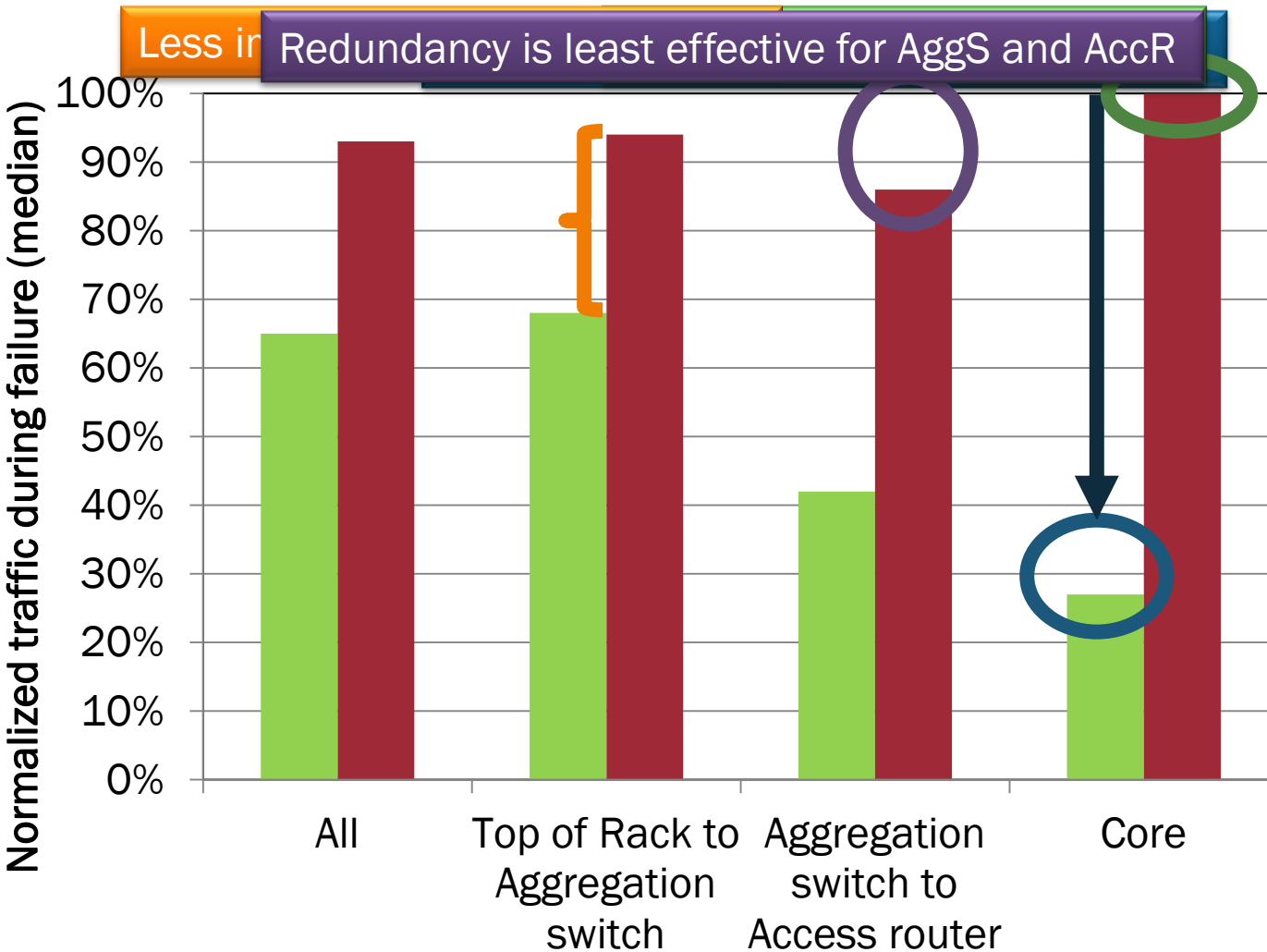
Measure traffic on links:

1. Before failure

2. During failure

3. Compute "normalized traffic" ratio:

$$\frac{traffic\ during}{traffic\ before} \sim 1$$

Compare normalized traffic over redundancy groups to normalized traffic on the link that failed
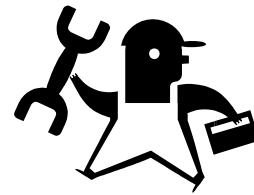
# Is redundancy effective in reducing impact?



Less in...

Redundancy is least effective for AggS and AccR

Y-axis: Normalized traffic during failure (median) — 0% to 100%

X-axis categories: All, Top of Rack to Aggregation switch, Aggregation switch to Access router, Core

Internet

**Overall increase of 40% in terms of traffic due to redundancy**
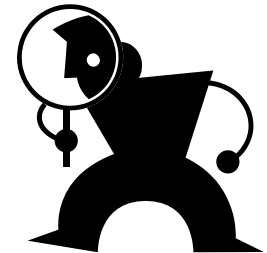
# Road Map

Motivation

Background & Methodology

Results

1. Characterizing failures

2. Do current network redundancy strategies help?
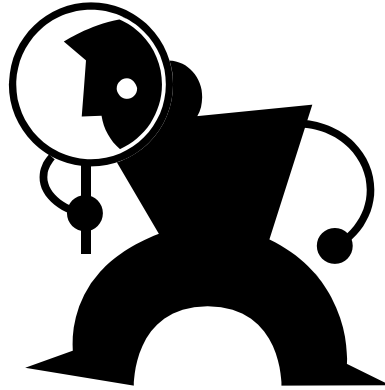
Conclusions

# Conclusions

- **Goal: Understand failures in data center networks**
    - Empirical study of data center failures

- **Key observations:**
    - Data center networks have high reliability
    - Low-cost switches exhibit high reliability
    - Load balancers are subject to transient faults
    - Failures may lead to loss of small packets

- **Future directions:**
    - Study application level failures and their causes
    - Further study of redundancy effectiveness

# Thanks!

Contact: phillipa@cs.toronto.edu

**Project page:**

**http://research.microsoft.com/~navendu/netwiser**