



RADoN: Storage Network QoS

Andrew Shewmaker¹
Carlos Maltzahn¹
Scott Brandt¹

Tim Kaldewey¹
Richard Golding²
Theodore M Wong²

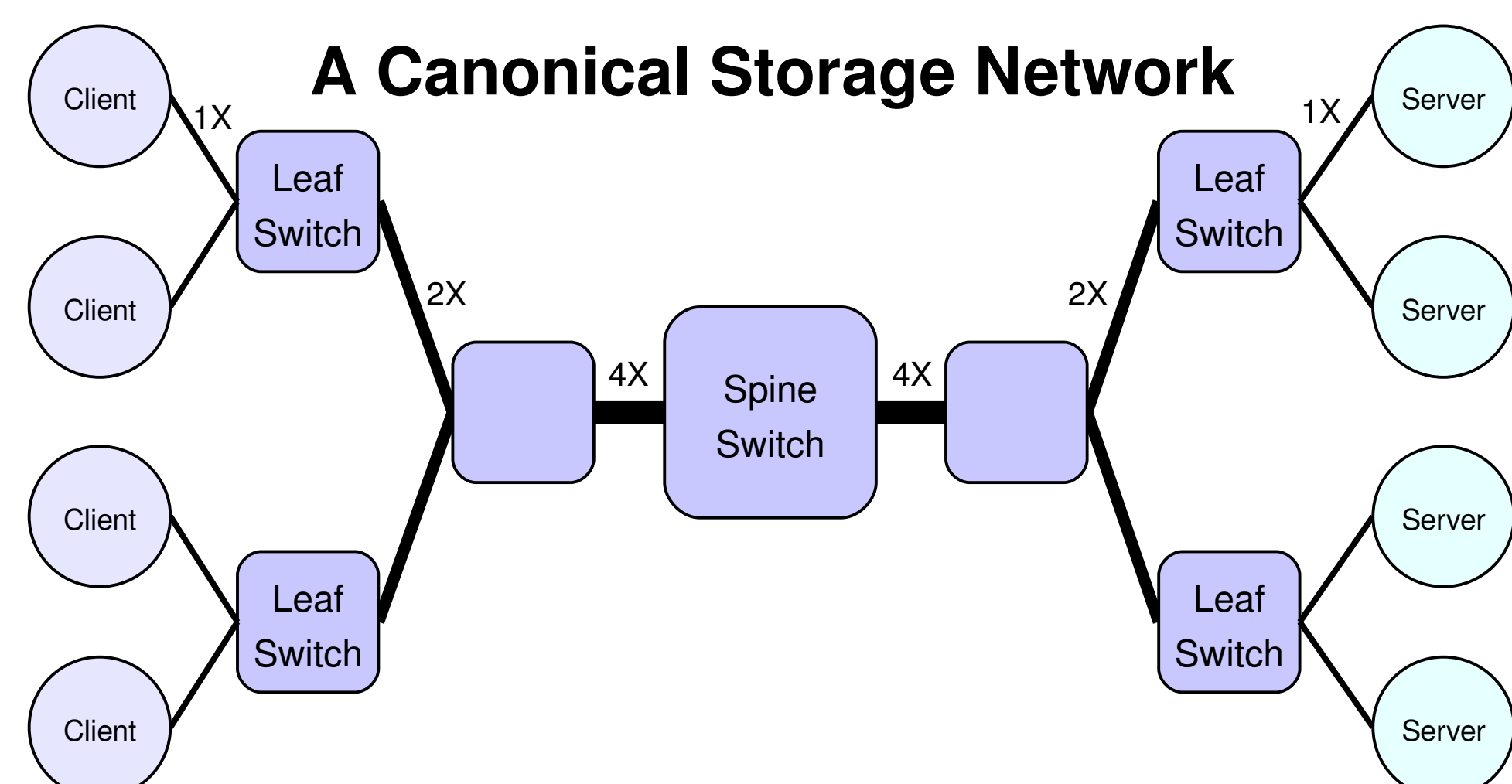
¹University of California Santa Cruz
Computer Science Department
{shewa,kalt,carlosm,scott}@cs.ucsc.edu

¹IBM Almaden Research Center
Storage Systems Department
{rgolding,theowong}@us.ibm.com

The Goal of RADoN

Performance guarantees on standard commodity storage networks

- General
- Flexible
- Fine grained



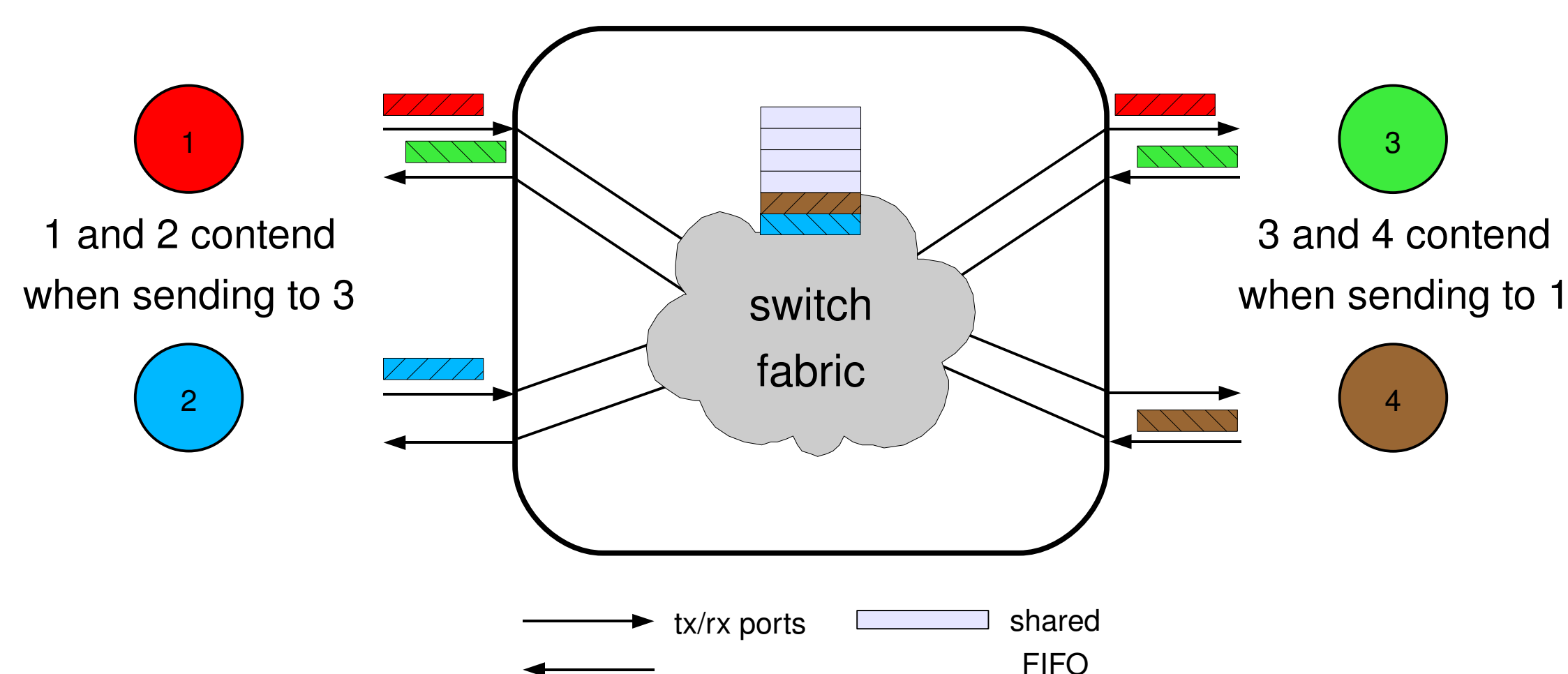
Switch Congestion is the Key Problem

Rate guarantees ensure a feasible volume of network traffic
Congestion in switches due to port contention may still cause

- large variance in delays
- packet loss

Intra and Inter Port Contention in a Simple Switch Model

2 and 4 congest on the queue

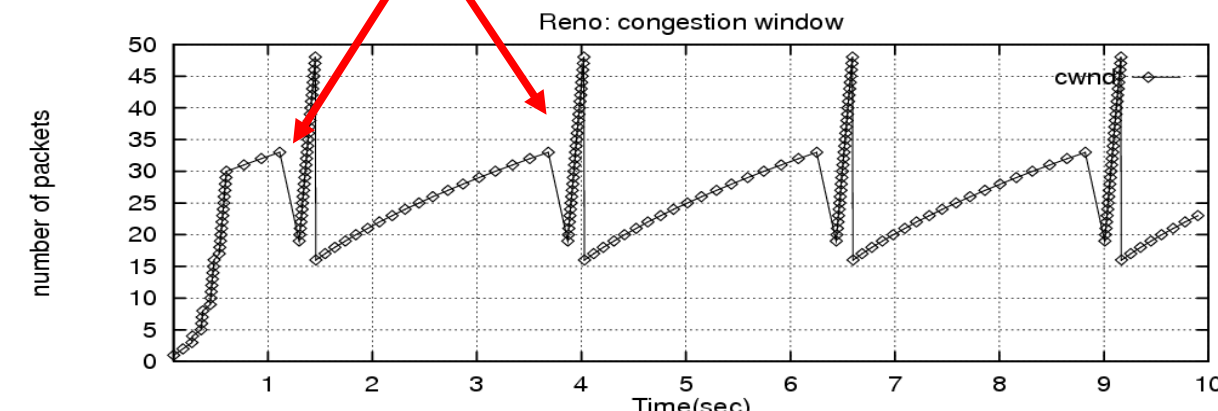


Methods of Congestion Detection

- Traditionally: packet loss, cannot prevent queue overflows
- RADoN uses *Forward Delay* as in TCP Santa Cruz

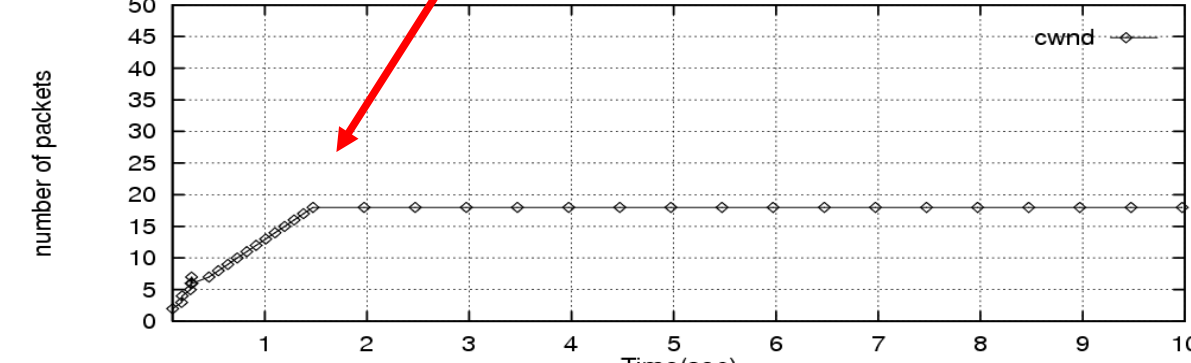
Packet Loss Based Congestion Control

Client Observes Packet Loss

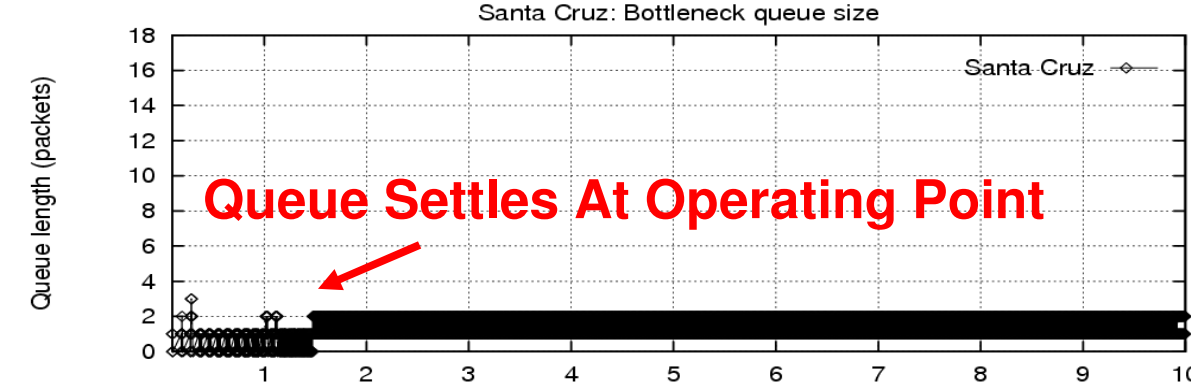
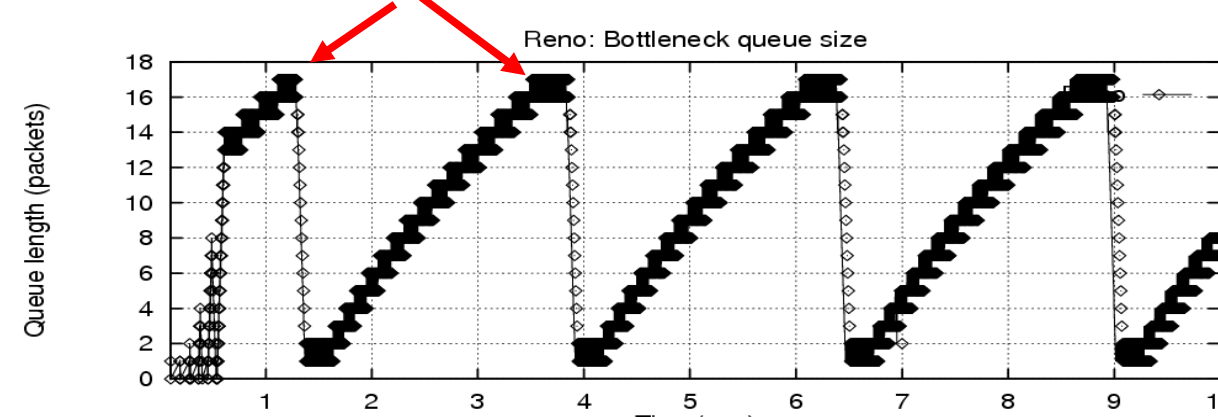


Forward Delay Based Congestion Control

Client Observes Increased Delay



Bottleneck Switch Queue Overflows



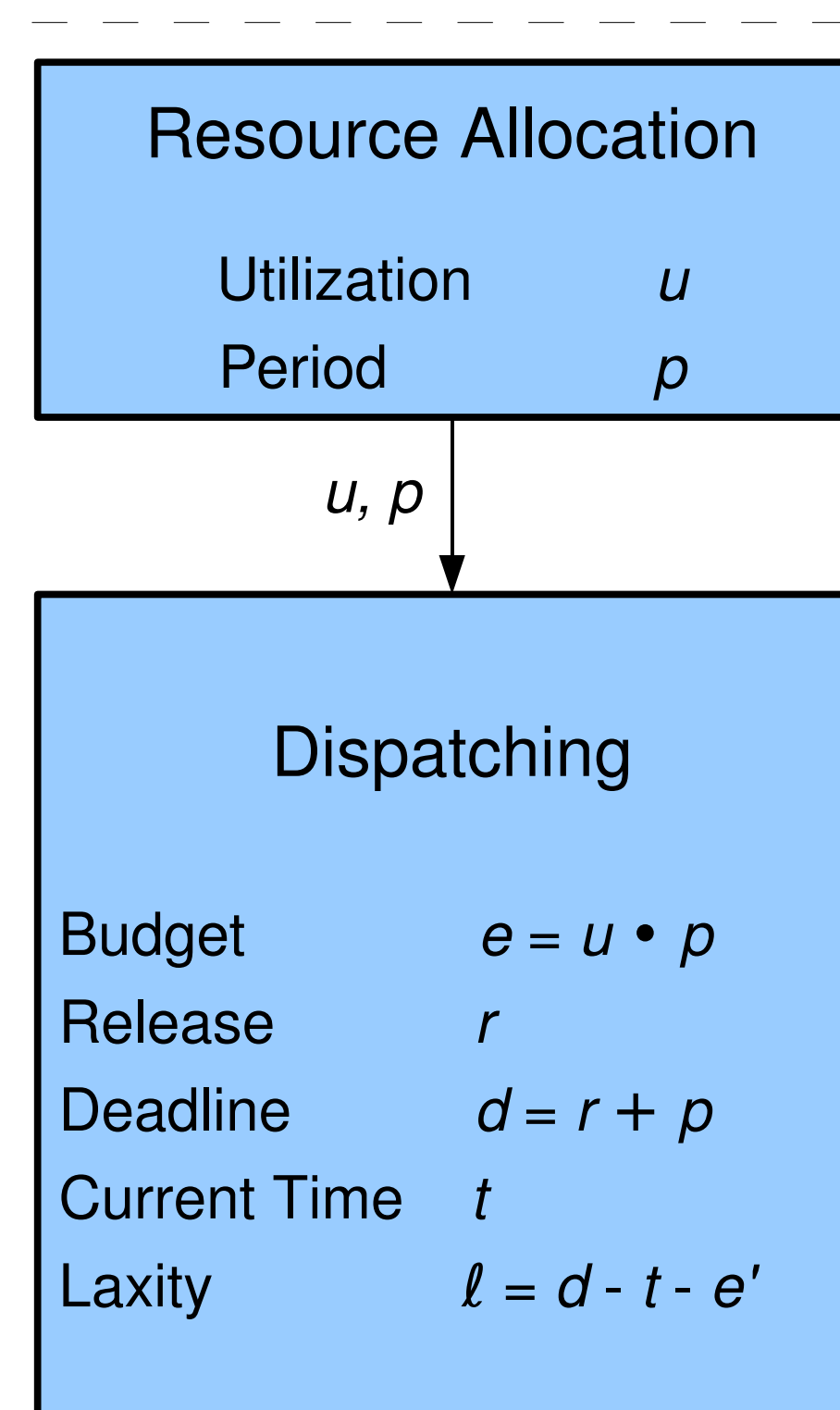
Improving TCP Congestion Control Over Internets with Heterogeneous Transmission Media (1999)
Christina Parsa, J.J. Garcia-Luna-Aceves. Proceedings of the 7th IEEE ICNP

Correct Congestion Detection and Controlled Response

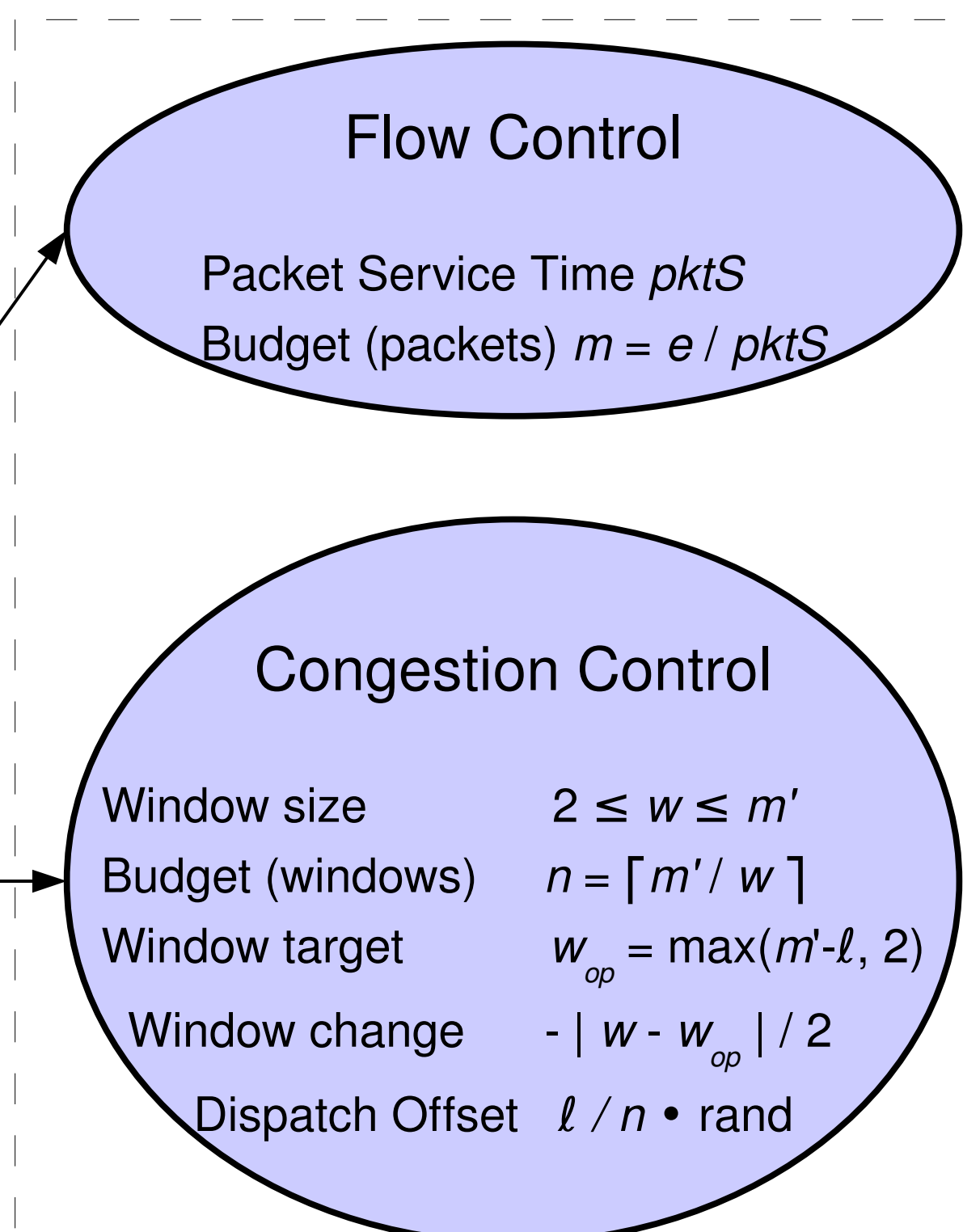
- Corresponding to a share
- Bound by a time limit
- Addressing debilitating synchronization

The RAD model separates scheduling into Resource Allocation and Dispatching.
Proven correct for CPU scheduling, RADoN applies the model to the network resource.

RAD Model



Network Model

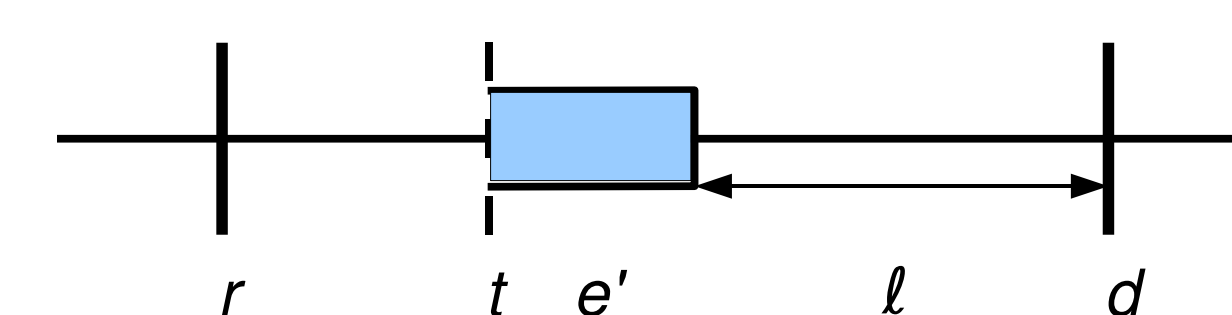


Subscripts omitted for clarity

The key to RADoN is a guaranteed real-time dispatching algorithm

Least Laxity First (LLF) is a guaranteed scheduling algorithm that dispatches jobs in order of criticality.

$$\text{Laxity } l = d - t - e'$$



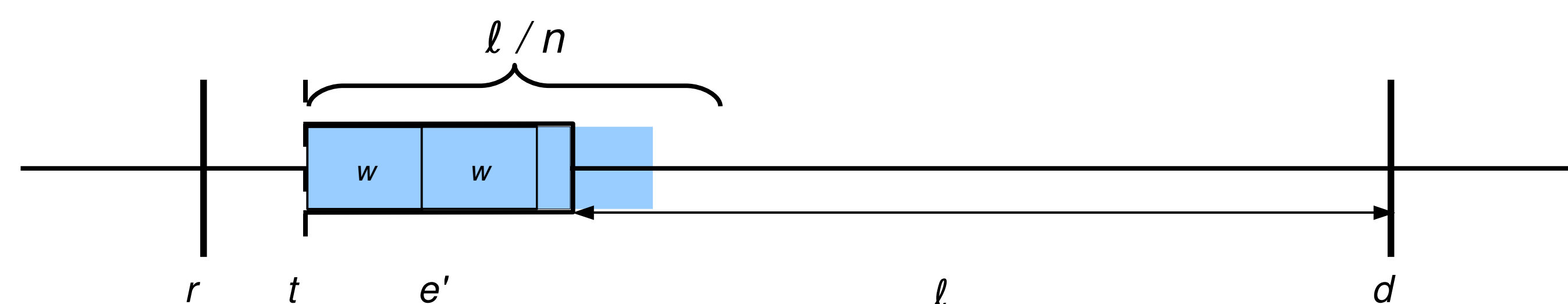
Less Laxity More (LLM) approximates LLF without requiring:

- Global knowledge
- Synchronization

A stream's work for a given period is divided into windows that change in dispatch time and size according to its laxity and remaining work when congestion is detected.

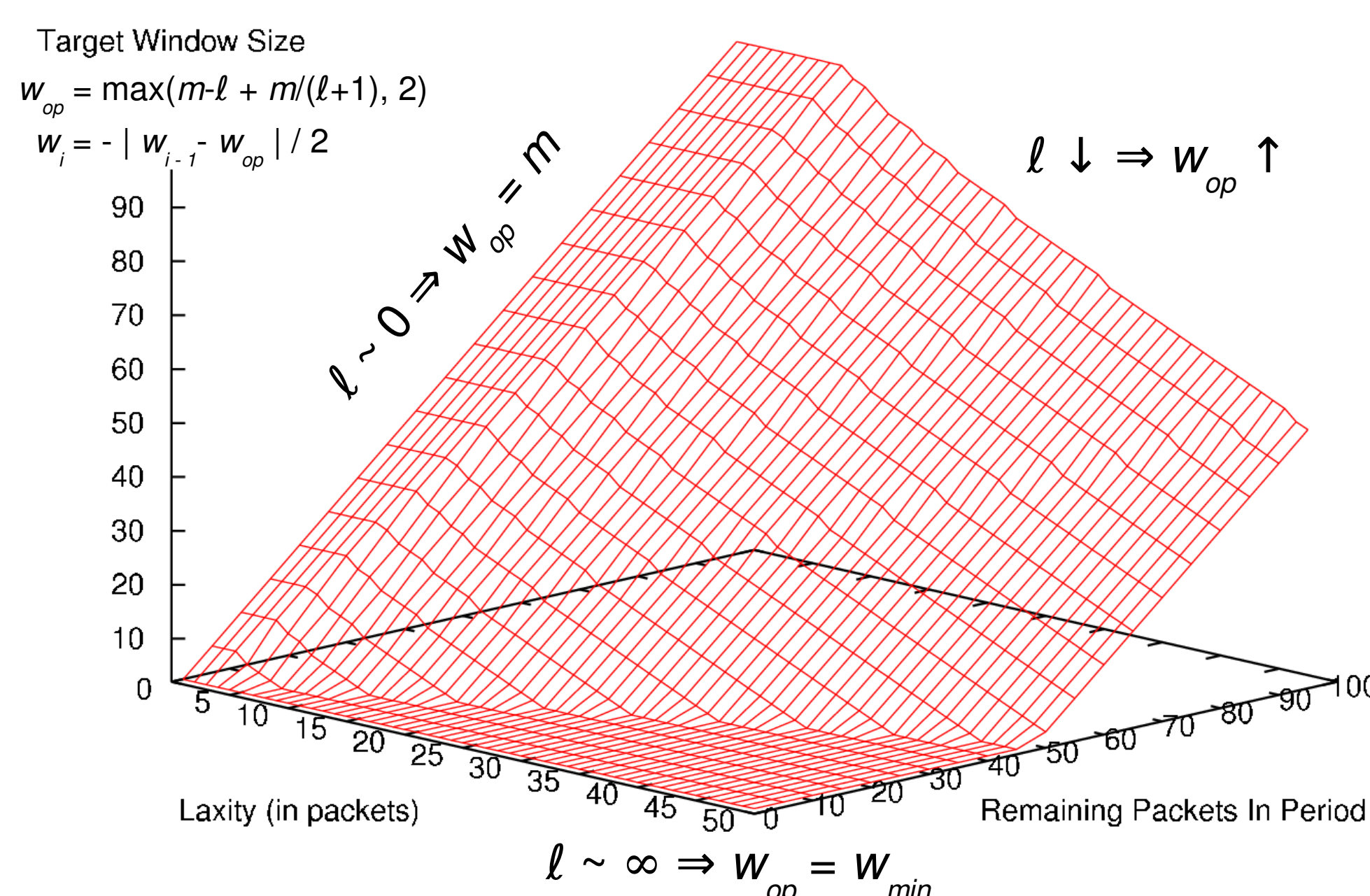
Laxity Based Congestion Response

Dispatch offset range



Window Size Response to Congestion

Target Window Size
 $w_{op} = \max(m-l, m/(l+1), 2)$
 $w_i = \lfloor w_{i-1} - w_{op} \rfloor / 2$



RADoN Summary

Flow = work per period

Congestion detected with *Forward Delay*

Congestion response

- controlled via window size
- bounded by laxity
- proportional to remaining work