How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

# *Pathrate* on mobile environments
# An implementation on Android

Antonio Macrì · Francesco Racciatti · Silvia Volpe

July 8, 2013

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

## Contents

**How *pathrate* works**
*pathrate* in mobile environments
**SmartPathrate**

Packet-pair technique
Capacity modes
Packet trains
Capacity estimation methodology

# Section 1

## How *pathrate* works

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
Packet trains
Capacity estimation methodology

## Pathrate

### Pathrate goal
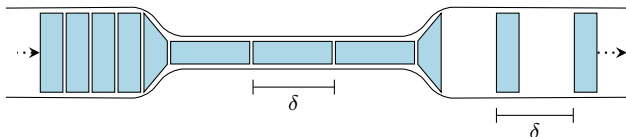
Pathrate is a tool that calculates the capacity of a path.

It works in 2 steps:

phase 1 uses packet pairs to obtain a set of capacity estimations

phase 2 uses packet trains to compute the lower bound of the capacity

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Packet-pair technique**
Capacity modes
Packet trains
Capacity estimation methodology

## Packet-pair technique

### Scenario

Two consecutive *probing packets* leave the sender *back-to-back* and arrive at the receiver with a *dispersion* (spacing) that is determined by the *narrow link* in the path
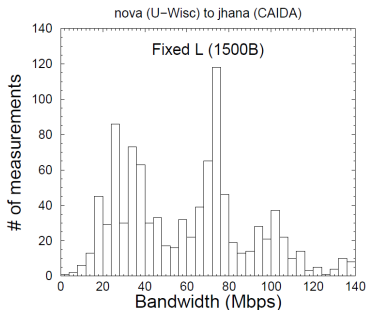


$$b = L/\delta$$

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
**Capacity modes**
Packet trains
Capacity estimation methodology

## Capacity modes

### Packet-pair bandwidth distribution

All capacities are used to compute a probability density function called the *packet-pair bandwidth distribution* $\mathcal{B}$
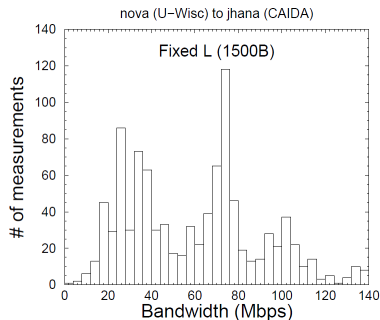
How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
**Capacity modes**
Packet trains
Capacity estimation methodology

# Capacity modes (cont.)

## Problem 1

The distribution of the capacities is *multimodal*

## Problem 2

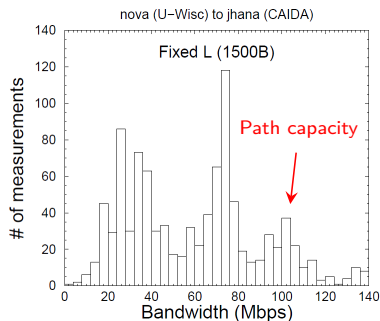The correct path capacity is likely to be the global mode only when the path is *lightly loaded*



nova (U–Wisc) to jhana (CAIDA)

Fixed L (1500B)

# of measurements vs Bandwidth (Mbps)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
Packet trains
Capacity estimation methodology

# Capacity modes (cont.)

**Problem 1**

The distribution of the capacities is *multimodal*

**Problem 2**

The correct path capacity is likely to be the global mode only when the path is *lightly loaded*
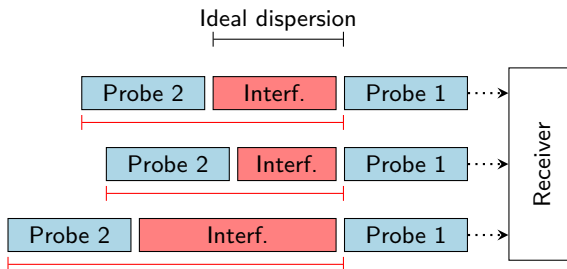
**How *pathrate* works**
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
**Capacity modes**
Packet trains
Capacity estimation methodology

# Capacity modes (cont.)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
**Capacity modes**
Packet trains
Capacity estimation methodology

## Capacity modes (cont.)

**Wrong estimations**

Underestimations due to cross-traffic (SCDR zone)

**How *pathrate* works**
*pathrate* in mobile environments
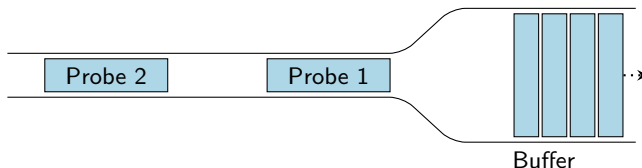SmartPathrate

Packet-pair technique
**Capacity modes**
Packet trains
Capacity estimation methodology

## Capacity modes (cont.)

### Wrong estimations

Overestimations due to non-empty buffers in some post-narrow node (PNCMs)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
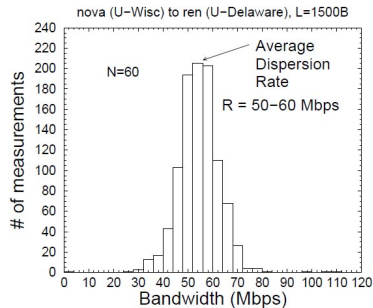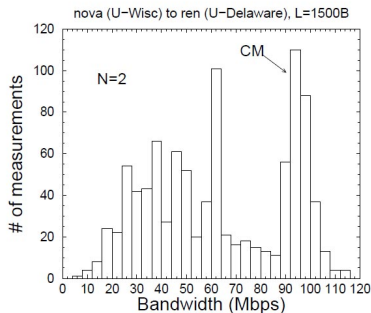Capacity modes
Packet trains
Capacity estimation methodology

## Packet trains

Generalization: using packet trains ($N > 2$ back-to-back packets of the same size $L$) we can calculate the bandwith as:

$$b(N) = \frac{(N-1)L}{\Delta(N)}$$

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
**Packet trains**
Capacity estimation methodology

## Packet trains

Generalization: using packet trains ($N > 2$ back-to-back packets of the same size $L$) we can calculate the bandwith as:

$$b(N) = \frac{(N-1)L}{\Delta(N)}$$

When the train length $N$ is sufficiently large, bandwidth measurements tend toward a single value leading to a unimodal distribution that becomes *independent of N*: this value is called the *Average Dispersion Rate* (ADR)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
**Packet trains**
Capacity estimation methodology

# Average Dispersion Rate

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
Packet trains
**Capacity estimation methodology**

## Capacity estimation methodology

### SCDR and PNCMs countermeasures

Varying probing packet size among different packet pairs makes SCDR modes wider and weaker. Probing packets as large as possible (but not too much) reduce the creation of PNCMs

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Packet-pair technique
Capacity modes
Packet trains
**Capacity estimation methodology**

## Capacity estimation methodology

### SCDR and PNCMs countermeasures

Varying probing packet size among different packet pairs makes
SCDR modes wider and weaker. Probing packets as large as possible
(but not too much) reduce the creation of PNCMs

### ADR as lower bound

The ADR is a *lower bound* of the capacity of the path

**How *pathrate* works**
*pathrate* **in mobile environments**
**SmartPathrate**

Packet-pair technique
Capacity modes
Packet trains
**Capacity estimation methodology**

## Capacity estimation methodology

### SCDR and PNCMs countermeasures

Varying probing packet size among different packet pairs makes SCDR modes wider and weaker. Probing packets as large as possible (but not too much) reduce the creation of PNCMs

### ADR as lower bound

The ADR is a *lower bound* of the capacity of the path

### Final capacity estimate

In the packet pair bandwidth distribution, ignore modes below the ADR mode and choose the one with the maximum *figure of merit*

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Section 2

## *pathrate* in mobile environments

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

## Main problems

Why do not port *pathrate* to Android?

- long running time $(15 \div 30$ mins$)$

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Main problems**
Interrupt Coalescence
IC distortions

## Main problems

Why do not port *pathrate* to Android?

- long running time ($15 \div 30$ mins)
- no care for consumed traffic ($100 \div 180$ MB)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Main problems**
Interrupt Coalescence
IC distortions

## Main problems

Why do not port *pathrate* to Android?

- long running time ($15 \div 30$ mins)
- no care for consumed traffic ($100 \div 180$ MB)
- not suitable for energy-constrained devices

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Main problems**
Interrupt Coalescence
IC distortions

## Not suitable for energy-constrained devices

- CPU frequency is scaled based on battery level or current load

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Main problems**
Interrupt Coalescence
IC distortions

# Not suitable for energy-constrained devices

- CPU frequency is scaled based on battery level or current load
- NIC rate depends on channel conditions

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

**Main problems**
Interrupt Coalescence
IC distortions

# Not suitable for energy-constrained devices

- CPU frequency is scaled based on battery level or current load
- NIC rate depends on channel conditions



- Network device driver may activate *Interrupt Coalescence* (IC)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Interrupt Coalescence

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
**Interrupt Coalescence**
IC distortions

# Interrupt Coalescence in wired environment

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Interrupt Coalescence in wired environment



### Detect IC

*pathrate* detects coalescence by sending a back-to-back packet train and comparing (at the receiver) measured dispersions with the kernel-to-user latency $\delta_{k-u}$

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

## Interrupt Coalescence in wired environment



### Detect IC

*pathrate* detects coalescence by sending a back-to-back packet train and comparing (at the receiver) measured dispersions with the kernel-to-user latency $\delta_{k\text{-}u}$

### Capacity from IC

In case of IC, the capacity of the path can be calculated by dividing the plateau length by the jump height

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

## Interrupt Coalescence in mobile environment

### Problems

- In the NIC buffer there are packets for other applications

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

## Interrupt Coalescence in mobile environment

### Problems

- In the NIC buffer there are packets for other applications
- Kernel can adapt CPU frequency to computational load and battery level

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Interrupt Coalescence in mobile environment

### Problems

- In the NIC buffer there are packets for other applications
- Kernel can adapt CPU frequency to computational load and battery level
- Timer resolution can be insufficient ($30\,\mu s$)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Interrupt Coalescence in mobile environment

### Problems

- In the NIC buffer there are packets for other applications
- Kernel can adapt CPU frequency to computational load and battery level
- Timer resolution can be insufficient ($30\,\mu s$)
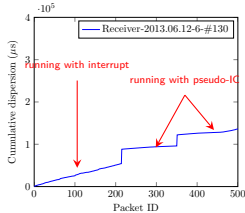- Height of jump depends on many factors, not only packet buffering (e. g. context switches, scheduling mechanisms)

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Interrupt Coalescence in mobile environment

### Problems

- In the NIC buffer there are packets for other applications
- Kernel can adapt CPU frequency to computational load and battery level
- Timer resolution can be insufficient ($30\,\mu s$)
- Height of jump depends on many factors, not only packet buffering (e. g. context switches, scheduling mechanisms)

### Pseudo-IC

Due to these factors the observed behavior is highly variable. Then, in mobile environment it is more correct to talk about pseudo-IC rather than simple IC.

How *pathrate* works
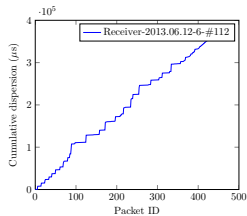***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
**IC distortions**

# Pseudo-IC

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
**IC distortions**

# Pseudo-IC

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
**IC distortions**

# Pseudo-IC

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
**IC distortions**

# Pseudo-IC

How *pathrate* works
***pathrate* in mobile environments**
SmartPathrate

Main problems
Interrupt Coalescence
**IC distortions**

# Pseudo-IC

How *pathrate* works
***pathrate* in mobile environments**
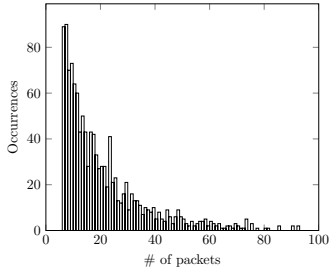SmartPathrate

Main problems
Interrupt Coalescence
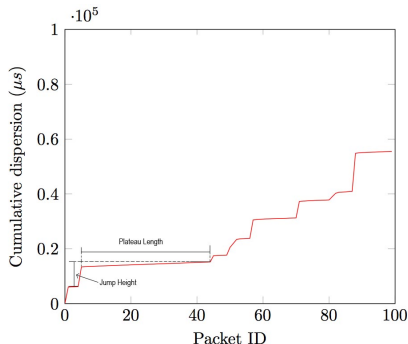**IC distortions**

# Pseudo-IC (cont.)



### Plateaus length & jumps height

If the behavior of pseudo-IC was regular then all the occurrences would be concentrated around a single value. The same speech applies to jumps height.

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

Main problems
Interrupt Coalescence
IC distortions

# Pseudo-IC (cont.)

## Calculation of path capacity

Pseudo-IC patterns are very irregular then it is not possible to calculate capacity with jumps and plateaus.

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
Pseudo-IC

# Section 3

## SmartPathrate

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

**Objectives**
Efficiency
The core procedure
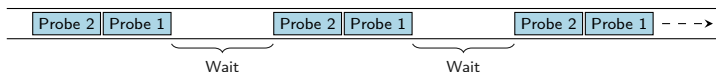Pseudo-IC

## Our objectives

The execution of our application

- should not take a long time
    - user wants to have a result quickly
    - battery life should not be significantly affected

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

**Objectives**
Efficiency
The core procedure
Pseudo-IC

## Our objectives

The execution of our application

- should not take a long time
  - user wants to have a result quickly
  - battery life should not be significantly affected
- should not require too much resources
  - processing complexity of mathematical computations
  - memory usage for buffering partial results

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

**Objectives**
Efficiency
The core procedure
Pseudo-IC

## Our objectives

The execution of our application

- should not take a long time
  - user wants to have a result quickly
  - battery life should not be significantly affected
- should not require too much resources
  - processing complexity of mathematical computations
  - memory usage for buffering partial results
- should use as less data as possible
  - mobile data plan may have traffic limitations
  - more data means more processing

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

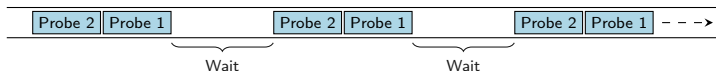Objectives
**Efficiency**
The core procedure
Pseudo-IC

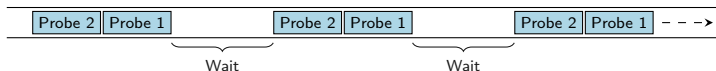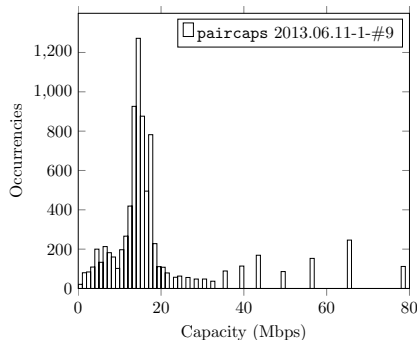## Packet pairs vs packet trains

- Pathrate: spacing between consecutive pairs
  - to avoid late packets interfere with subsequent pairs
  - drop pair in case of interference or packet losses

| Probe 2 | Probe 1 | | Probe 2 | Probe 1 | | Probe 2 | Probe 1 | – – –> |

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

## Packet pairs vs packet trains

- Pathrate: spacing between consecutive pairs
  - to avoid late packets interfere with subsequent pairs
  - drop pair in case of interference or packet losses

| Probe 2 | Probe 1 | | Probe 2 | Probe 1 | | Probe 2 | Probe 1 | – – –> |

Wait                Wait

  - total time in wait state is at least about $12 \div 13$ minutes

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

## Packet pairs vs packet trains

- Pathrate: spacing between consecutive pairs
  - to avoid late packets interfere with subsequent pairs
  - drop pair in case of interference or packet losses

| | Probe 2 | Probe 1 | | Probe 2 | Probe 1 | | Probe 2 | Probe 1 | – – –➤ |

Wait                    Wait

  - total time in wait state is at least about $12 \div 13$ minutes

- SmartPathrate: smaller spacing between trains
  - treat packets from previous trains as cross traffic
  - drop only in case of packet losses

| Probe 4 | Probe 3 | Probe 2 | Probe 1 | | Probe 4 | Probe 3 | Probe 2 | Probe 1 | – – –➤ |

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

# Processing complexity and memory usage



**Mode calculation**

$$3 \cdot O(n\omega)$$
$$\Downarrow$$
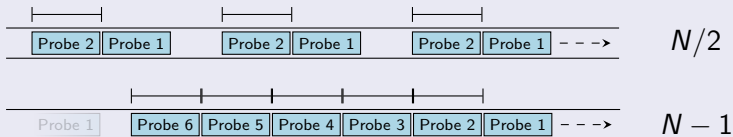$$O(n\omega) + 2 \cdot O(n \log \omega)$$

**Main change**

Linear search
$$\Downarrow$$
Binary search

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

# Data complexity

## Number of capacity estimates



$N/2$

$N-1$

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

## Data complexity



### Number of capacity estimates

| Probe 2 | Probe 1 | | Probe 2 | Probe 1 | | Probe 2 | Probe 1 | – – –> | $N/2$ |

| Probe 1 | | Probe 6 | Probe 5 | Probe 4 | Probe 3 | Probe 2 | Probe 1 | – – –> | $N - 1$ |

ADR

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
**Efficiency**
The core procedure
Pseudo-IC

# Data complexity

## Number of capacity estimates



$N/2$

$N-1$

## Data efficiency

Given a train, try to split it into separate zones, remove coalescence and extract as much information as possible
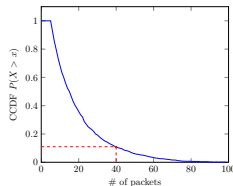
How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

## The core procedure



- Send maximum-size packets
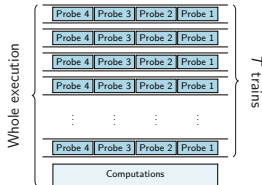- Start with trains of 40 packets
- Gradually increase train length

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

## The core procedure



- Send maximum-size packets
- Start with trains of 40 packets
- Gradually increase train length

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

# The core procedure (cont.)

- Proceed round by round
  - Why rounds?

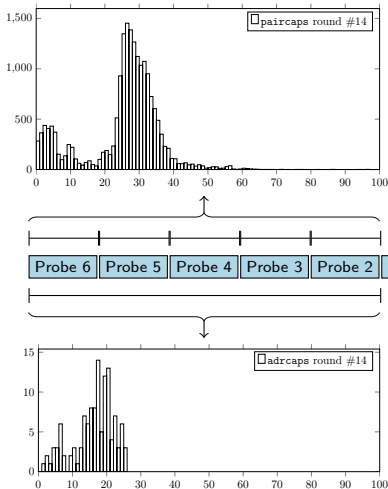How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

# The core procedure (cont.)

- Proceed round by round
  - Why rounds?

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

## The core procedure (cont.)

- Proceed round by round
  - Why rounds?

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

# The core procedure (cont.)



- Proceed round by round
  - Why rounds?

How *pathrate* works
*pathrate* in mobile environments
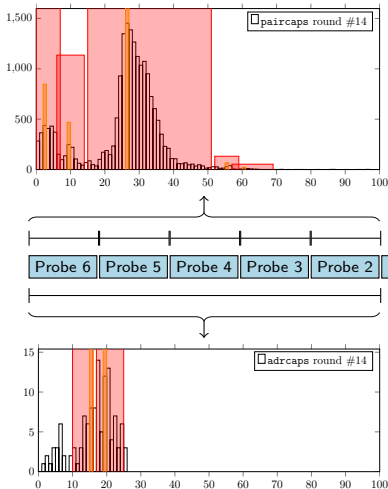**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

## The core procedure (cont.)



### Calculations at each round

- Calculate pair capacities
- Calculate ADR capacities
- Calculate distribution
- Calculate bin width

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
**The core procedure**
Pseudo-IC

## The core procedure (cont.)



### Calculations at each round

- Calculate pair capacities
- Calculate ADR capacities
- Calculate distribution
- Calculate bin width
- Extract modes
- Estimate (temporary) capacity

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

## Detect pseudo-IC

### Filtering capacities from IC

Detecting interrupt coalescence is a critical task

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

## Detect pseudo-IC

### Filtering capacities from IC

Detecting interrupt coalescence is a critical task

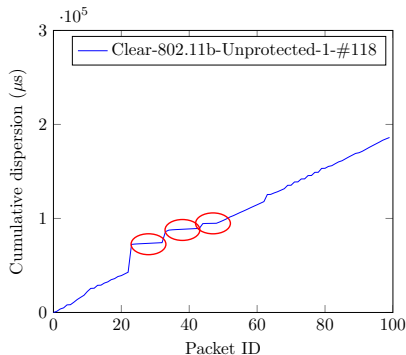### Idea

Why do not reduce the receiver's buffer size by means of
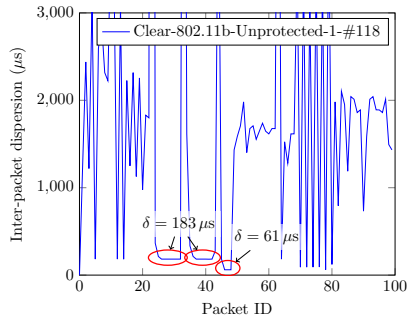setsockopt, to reduce or completely suppress IC?

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

## Detect pseudo-IC

### Filtering capacities from IC

Detecting interrupt coalescence is a critical task

### Idea

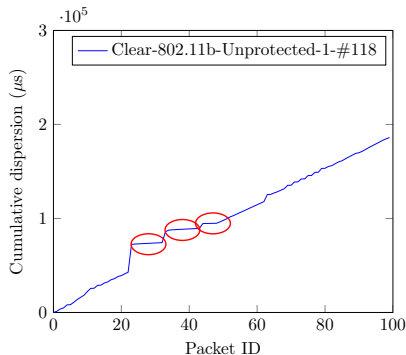Why do not reduce the receiver's buffer size by means of
setsockopt, to reduce or completely suppress IC?

### Wrong

Option SO_RCVBUF affects only a buffer related to the *UDP socket*,
and has nothing to do with the kernel's buffer used by NIC.
In fact, the receiver's buffer should be as large as possible, to reduce
(late) packet drops.

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**
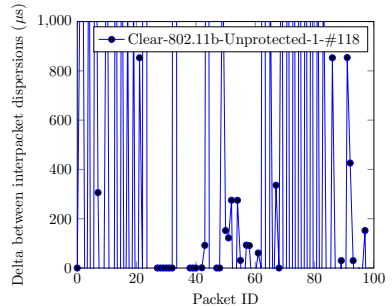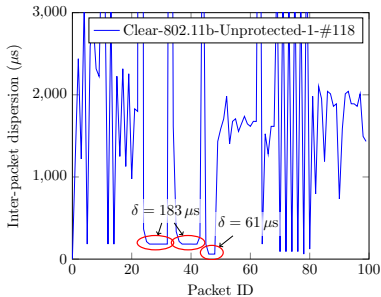
Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)



$$\frac{1500 \times 8 \text{ bit}}{183 \,\mu s} \approx 65.6 \,\text{Mbps}$$

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)



## Two *features*

Packet pair dispersion $\Rightarrow$ Minimum possible dispersion

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)



### Two *features*

| | | |
|---|---|---|
| Packet pair dispersion | $\Rightarrow$ | Minimum possible dispersion |
| Packet pair dispersion | $\Rightarrow$ | Kernel-to-user latency |

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Detect pseudo-IC (cont.)



## Two *features*

| | | |
|---|---|---|
| Packet pair dispersion | $\Rightarrow$ | Minimum possible dispersion |
| Packet pair dispersion | $\Rightarrow$ | Kernel-to-user latency |
| Delta of packet pair dispersions | $\Rightarrow$ | Kernel-to-user latency |

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Final capacity estimate

Procedure:

- Proceed round by round
- Do calculations at the end of each round
- Try to determine the path capacity
- Check for convergence of partial results



Typical results:

- Execution time: $15 \div 30$ mins
- Generated traffic: $100 \div 180$ MB

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

# Final capacity estimate

Procedure:

- Proceed round by round
- Do calculations at the end of each round
- Try to determine the path capacity
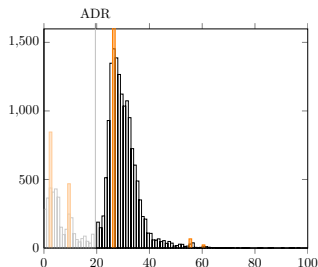- Check for convergence of partial results



Typical results:

- Execution time:     $15 \div 30$ mins   $\Rightarrow$   $1 \div 2$ mins
- Generated traffic:   $100 \div 180$ MB   $\Rightarrow$   $15 \div 30$ MB

How *pathrate* works
*pathrate* in mobile environments
**SmartPathrate**

Objectives
Efficiency
The core procedure
**Pseudo-IC**

## Questions?



Questions

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

## References I

URL: http://vger.kernel.org/~davem/skb_sk.html.

Jonathan Corbet, Alessandro Rubini, and Greg Kroah-Hartman. *Linux Device Drivers*. Third edition. O'Reilly, 2005. URL: http://lwn.net/images/pdf/LDD3/ch17.pdf.

Constantinos Dovrolis, Parameswaran Ramanathan, and David Moore. "Packet-dispersion techniques and a capacity-estimation methodology". In: *IEEE/ACM Transactions on Networking*. Vol. 12. 6. Dec. 2004, pp. 963–977.

Jamal Hadi Salim, Robert Olsson, and Alexey Kuznetsov. "Beyond Softnet". In: *Proceedings of the 5th Annual Linux Showcase & Conference*. Oakland, California, USA, 2001. URL: http://static.usenix.org/publications/library/proceedings/als01/full_papers/jamal/jamal.pdf.

Van Jacobson. "Congestion Avoidance and Control". In: *Proceedings of ACM SIGCOMM*. Sept. 1988, pp. 314–329.

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

## References II

Richard Kelly and Joseph Gasparakis. *Common Functionality in the 2.6 Linux Network Stack*. Tech. rep. 2010. URL: http://www.intel.com/content/dam/www/public/us/en/documents/white-papers/linux-2-6-network-stack-paper.pdf.

Ravi Prasad, Manish Jain, and Constantinos Dovrolis. "Effects of Interrupt Coalescence on Network Measurements". In: *Passive and Active Network Measurement*. Vol. 3015. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 2004, pp. 247–256.

Miguel Rio et al. *DataTAG – A Map of the Networking Code in Linux Kernel 2.4.20*. Tech. rep. URL: http://datatag.web.cern.ch/datatag/papers/tr-datatag-2004-1.pdf.

Rami Rosen. *Linux Kernel Networking – advanced topics*. 2009. URL: http://www.haifux.org/lectures/217/netLec5.pdf.

How *pathrate* works
*pathrate* in mobile environments
SmartPathrate

## References III

📕 Klaus Wehrle et al. *The Linux Networking Architecture: Design and Implementation of Network Protocols in the Linux Kernel*. Prentice Hall, 2004. URL: http://www.6test.edu.cn/~lujx/linux_networking/0131777203_ch25lev1sec3.html.