# Cellular Networks and Mobile Computing
# COMS 6998-8, Spring 2012

Instructor: Li Erran Li
(lierranli@cs.columbia.edu)
http://www.cs.columbia.edu/~coms6998-8/
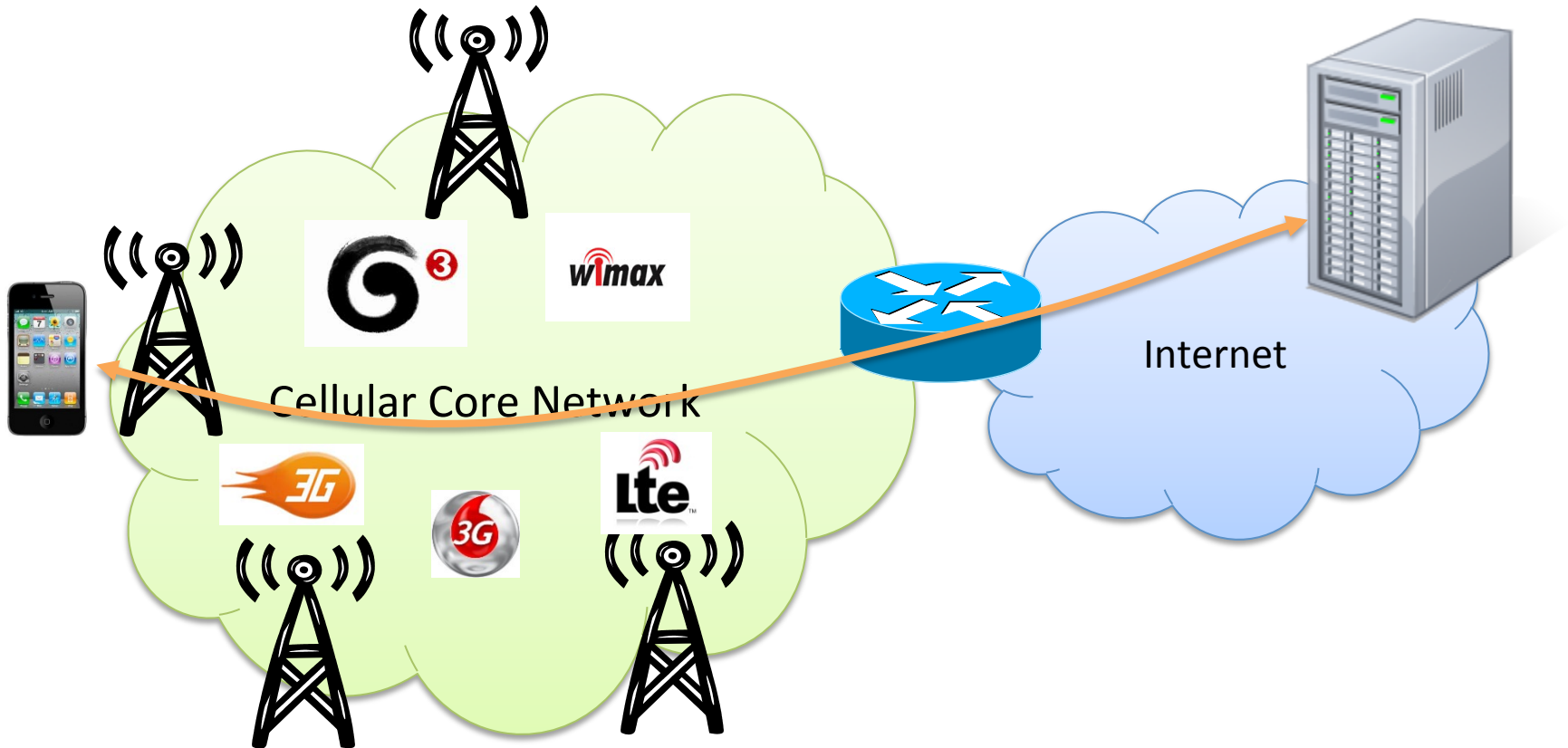3/26/2012: Cellular Network and Traffic Characterization

# An Untold Story of Middleboxes in Cellular Networks
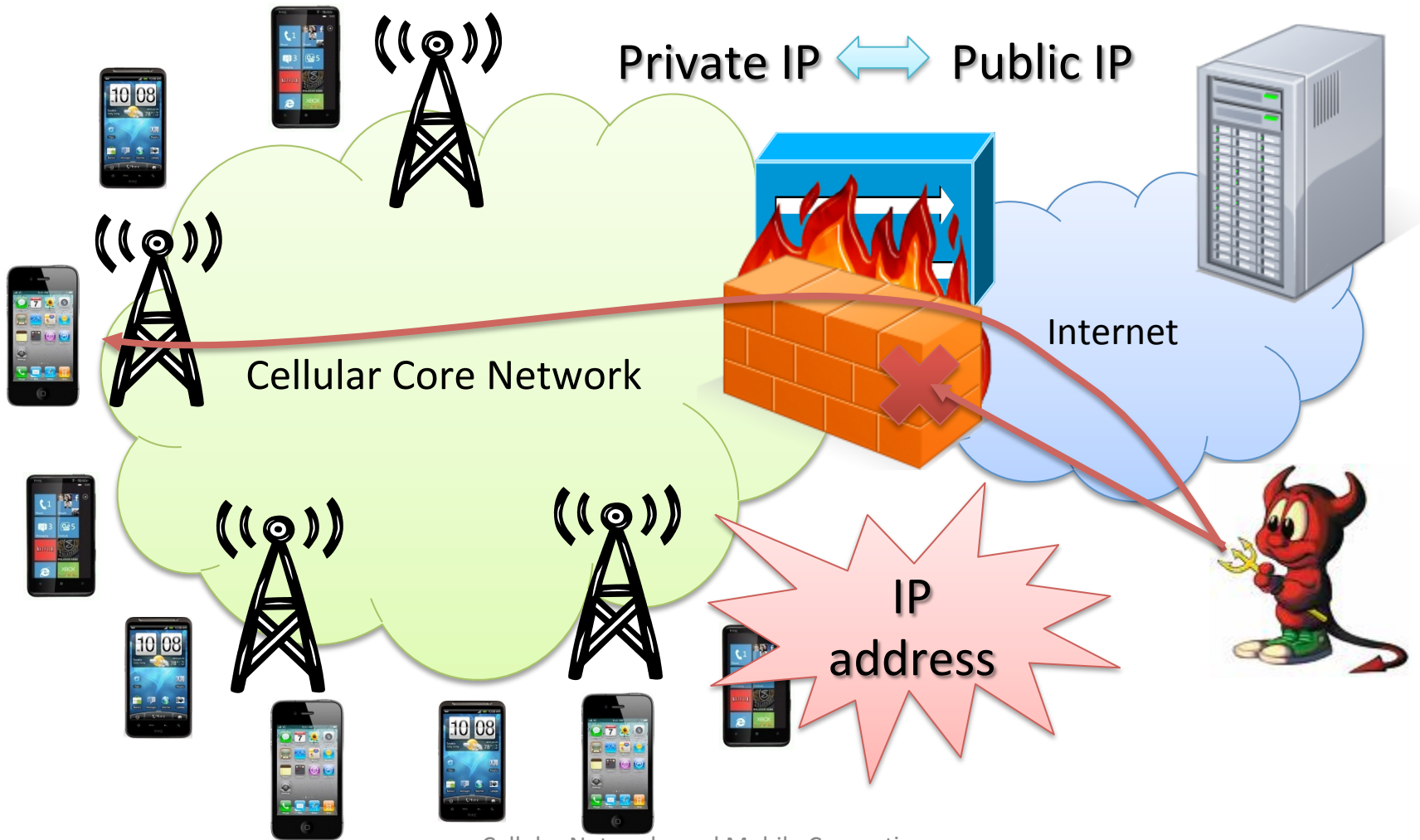
Zhaoguang Wang[1]

Zhiyun Qian[1], Qiang Xu[1], Z. Morley Mao[1], Ming Zhang[2]

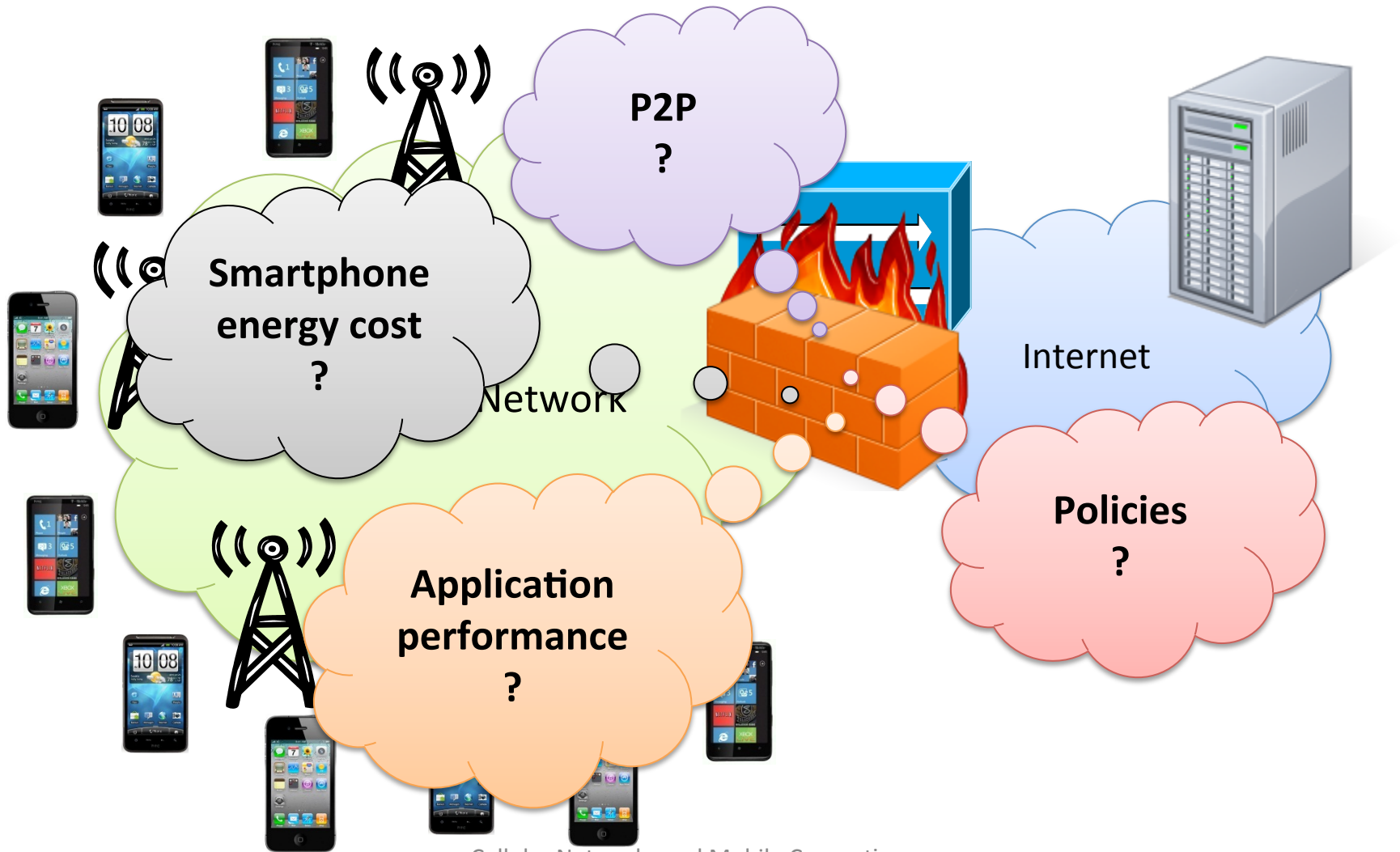[1]University of Michigan    [2]Microsoft Research

# Background on cellular network

# Why carriers deploy middleboxes?

Private IP ⟷ Public IP

Cellular Core Network

Internet

IP address

# Problems with middleboxes



P2P
?

Smartphone
energy cost
?

Network

Internet

Application
performance
?

Policies
?

Courtesy: Z. Wang et al.

# Challenges and solutions

- Policies can be complex and proprietary
  - √ Design a suite of end-to-end probes

- Cellular carriers are diverse
  - √ Publicly available client Android app

- Implications of policies are not obvious
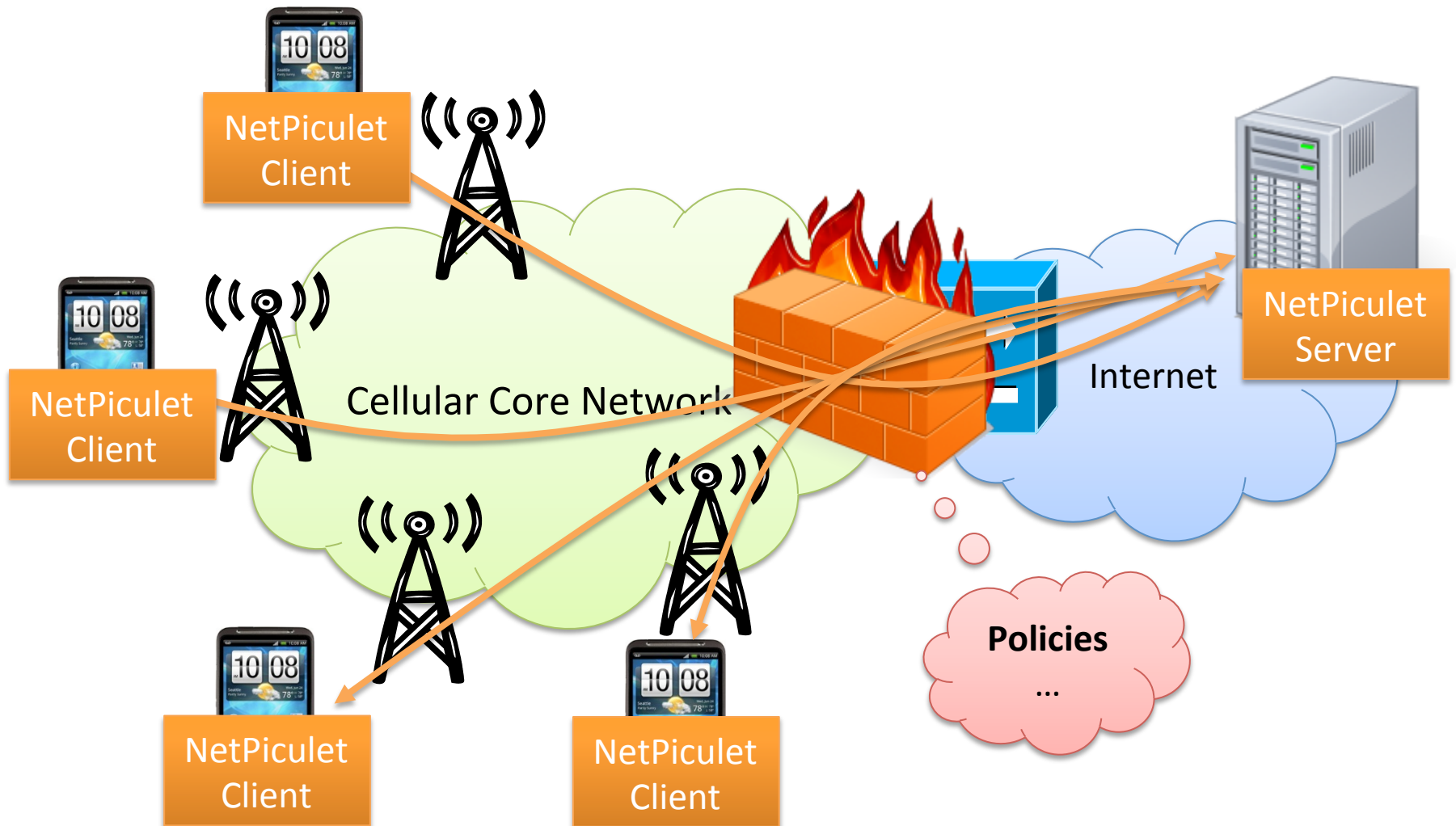  - √ Conduct controlled experiments

# Related work

- Internet middleboxes study
  - [Allman, IMC 03], [Medina, IMC 04]
- NAT characterization and traversal
  - STUN[MacDonald et al.], [Guha and Francis, IMC 05]
- Cellular network security
  - [Serror et al., WiSe 06], [Traynor et al., Usenix Security 07]
- Cellular data network measurement
  - WindRider, [Huang et al., MobiSys 10]

# Goals

- Develop a tool that accurately infers the NAT and firewall policies in cellular networks

- Understand the impact and implications
  - Application performance
  - Energy consumption
  - Network security

Courtesy: Z. Wang et al.

# The NetPiculet measurement system

# Target policies in NetPiculet

| Firewall | IP spoofing |
| | TCP connection timeout |
| | Out-of-order packet buffering |
| NAT | NAT mapping type |
| | Endpoint filtering |
| | TCP state tracking |
| | Filtering response |
| | Packet mangling |

# Target policies in NetPiculet

| Firewall | IP spoofing |
| --- | --- |
| | TCP connection timeout |
| | Out-of-order packet buffering |
| NAT | NAT mapping type |
| | Endpoint filtering |
| | TCP state tracking |
| | Filtering response |
| | Packet mangling |

Courtesy: Z. Wang et al.

# Key findings

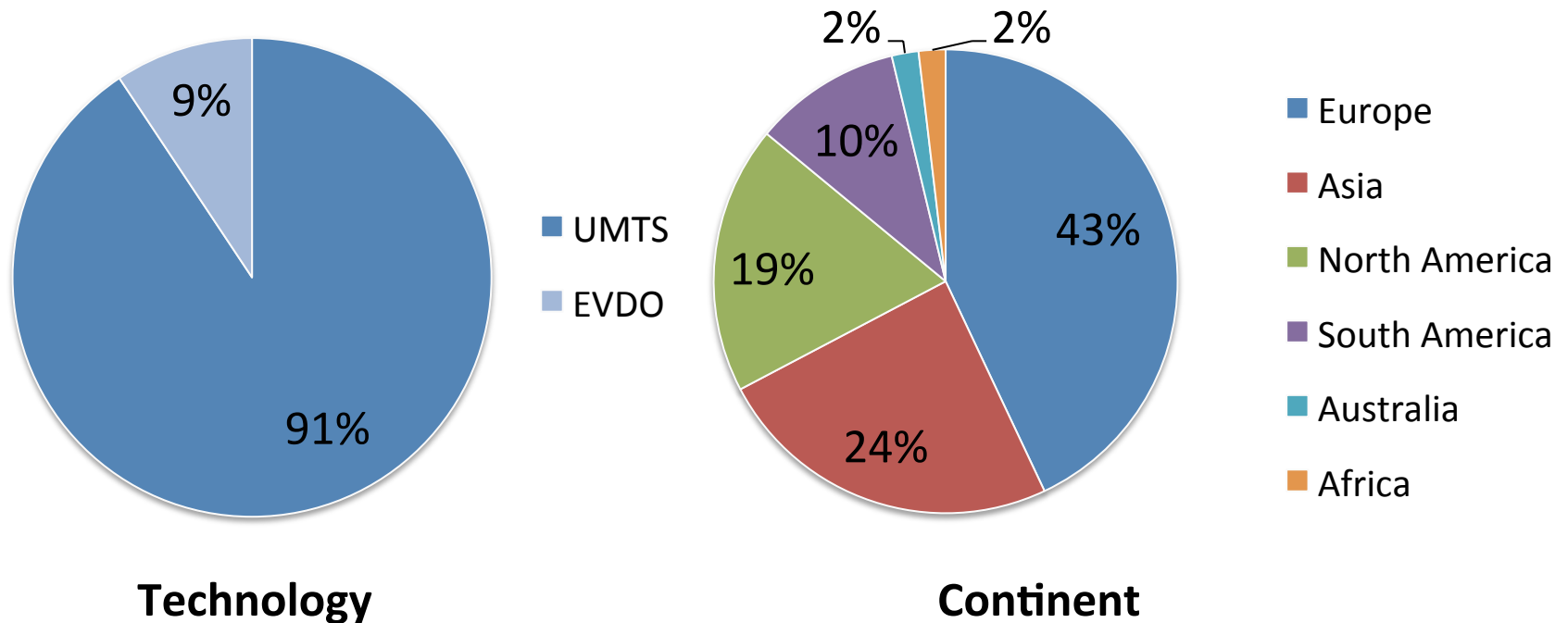| Firewall | Some carriers allow IP spoofing<br>**Create network vulnerability** |
|----------|--------------------------------------------------------------------|
|          | Some carriers time out idle connections aggressively<br>**Drain batteries of smartphones** |
|          | Some firewalls buffer out-of-order packet<br>**Degrade TCP performance** |
| NAT      | One NAT mapping linearly increases port # with time<br>**Classified as random in previous work** |

Courtesy: Z. Wang et al.

# Diverse carriers studied

- NetPiculet released in Jan. 2011
  - 393 users from 107 cellular carriers in two weeks



**Technology**

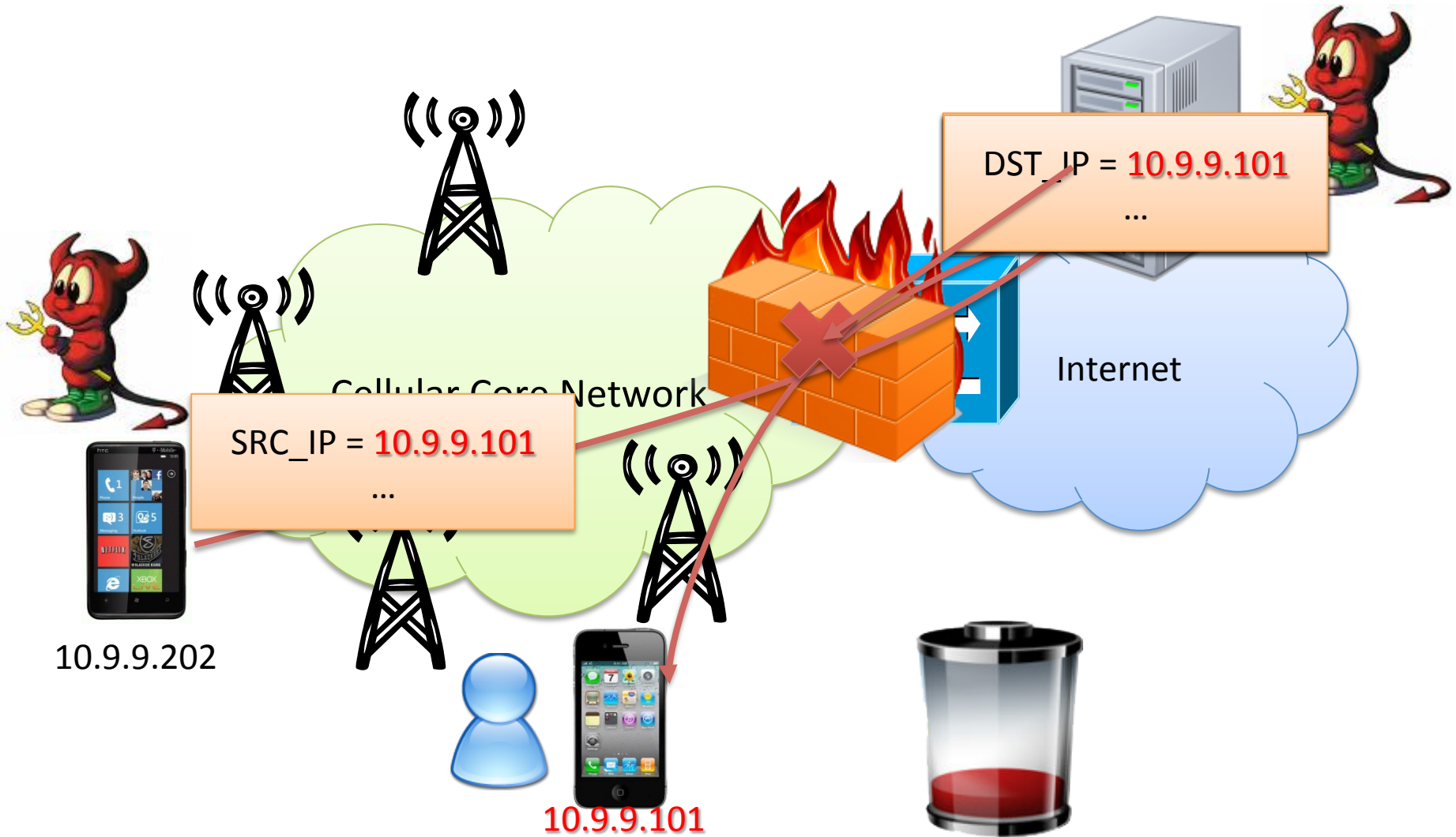| | |
|---|---|
| ■ | UMTS |
| ■ | EVDO |

9%

91%

**Continent**

2%     2%

10%

19%     43%

24%

| | |
|---|---|
| ■ | Europe |
| ■ | Asia |
| ■ | North America |
| ■ | South America |
| ■ | Australia |
| ■ | Africa |

# Outline

1. • IP spoofing
2. • TCP connection timeout
3. • TCP out-of-order buffering
4. • NAT mapping

Courtesy: Z. Wang et al.

# Outline

1. • IP spoofing

2. • TCP connection timeout

3. • TCP out-of-order buffering
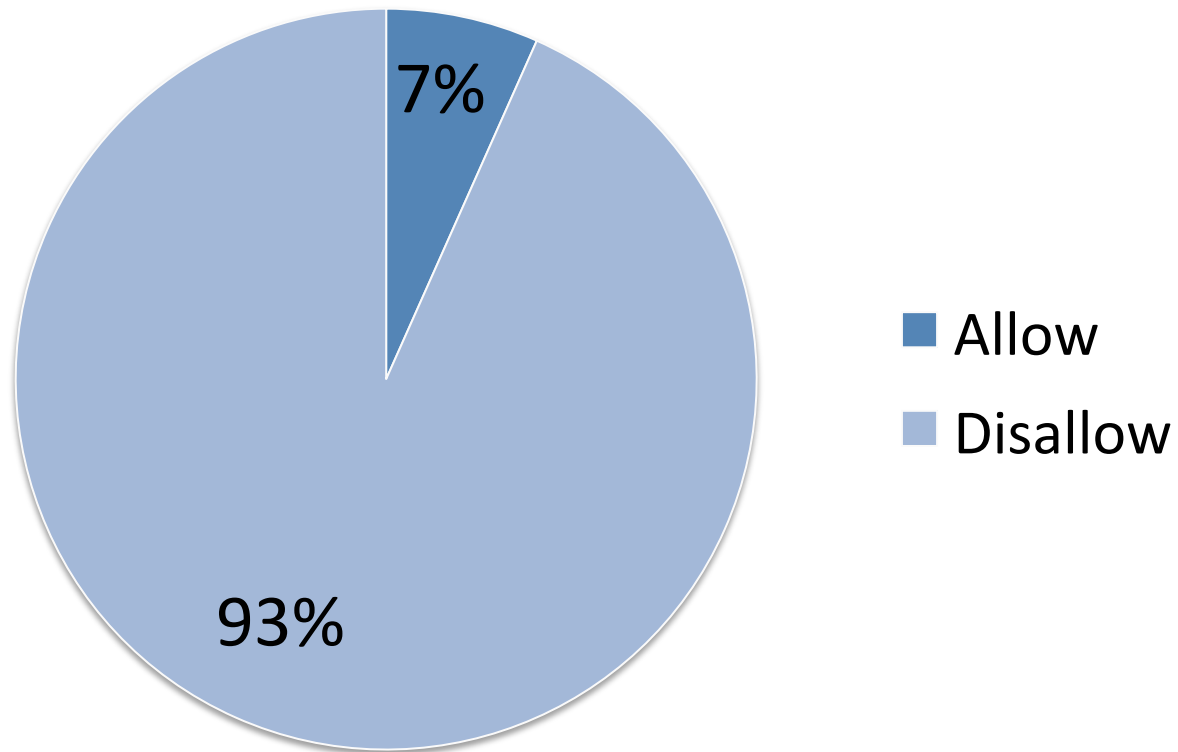
4. • NAT mapping

# Why allowing IP spoofing is bad?



DST_IP = 10.9.9.101
...

SRC_IP = 10.9.9.101
...

Cellular Core Network

Internet

10.9.9.202

10.9.9.101

Courtesy: Z. Wang et al.

# Test whether IP spoofing is allowed



SRC_IP = 10.9.9.202
PAYLOAD = 10.9.9.101

NetPiculet Client

10.9.9.101

Internet

NetPiculet Server

Allow IP spoofing!

# 4 out of 60 carriers allow IP spoofing

IP spoofing should be disabled



7%

93%

- Allow
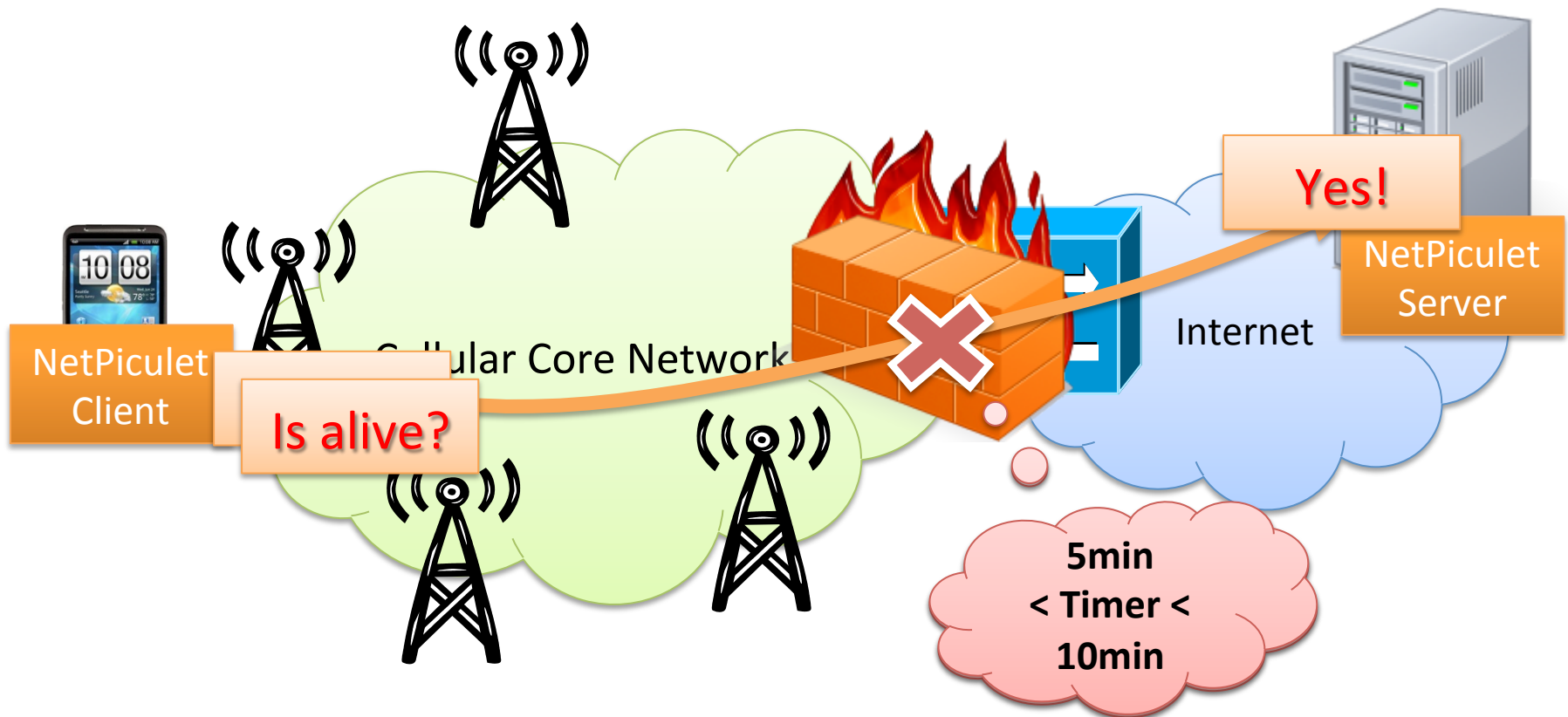- Disallow

# Outline

1. • IP spoofing

2. • TCP connection timeout

3. • TCP out-of-order buffering
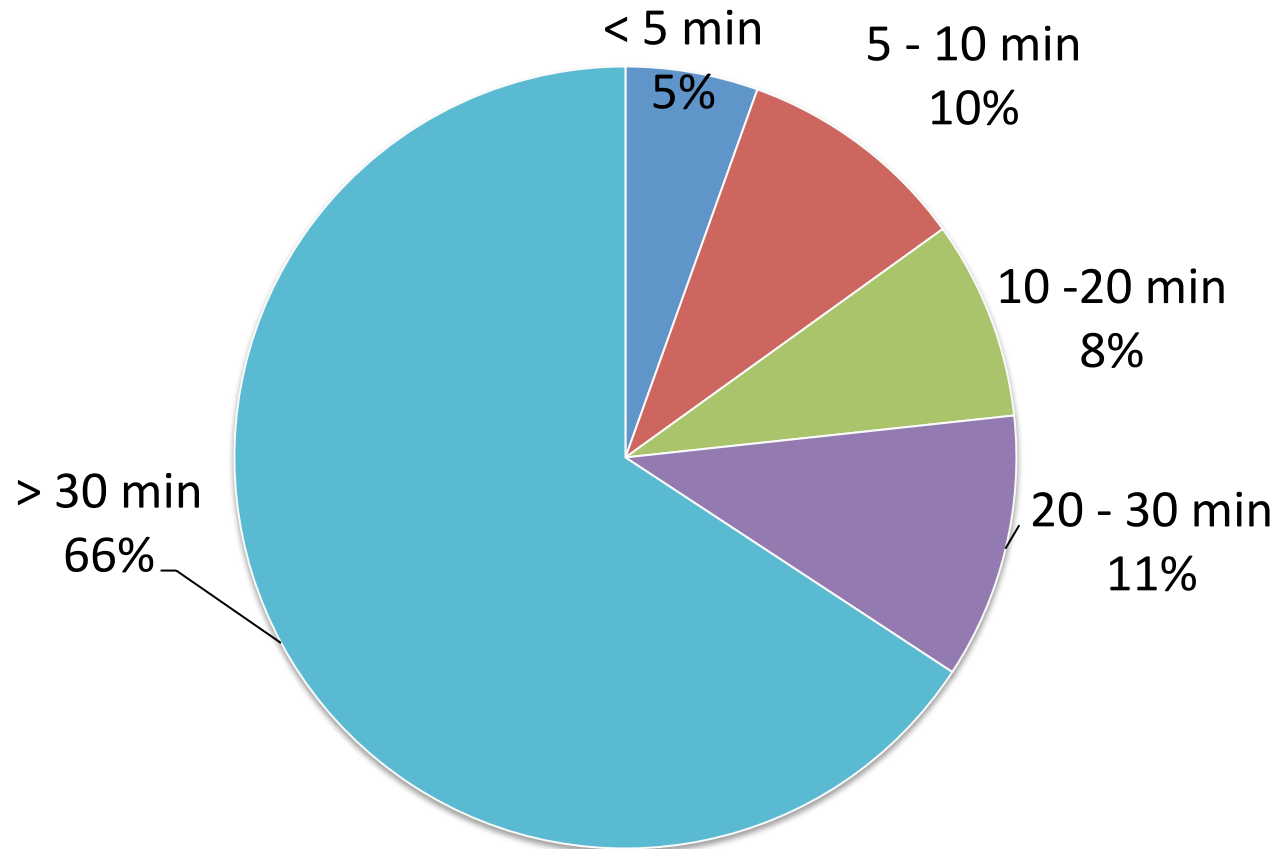
4. • NAT mapping

# Why short TCP timeout timers are bad?



Cellular Core Network

Internet

KEEP-ALIVE

Terminate Idle TCP Connection

# Measure the TCP timeout timer



Time = 10min

NetPiculet Client

Is alive?

Cellular Core Network

Yes!

NetPiculet Server

Internet

5min < Timer < 10min

# Short timers identified in a few carriers

**4 carriers set timers less than 5 minutes**



< 5 min
5%

5 - 10 min
10%

10 -20 min
8%

20 - 30 min
11%

> 30 min
66%

Cellular Networks and Mobile Computing (COMS 6998-8)
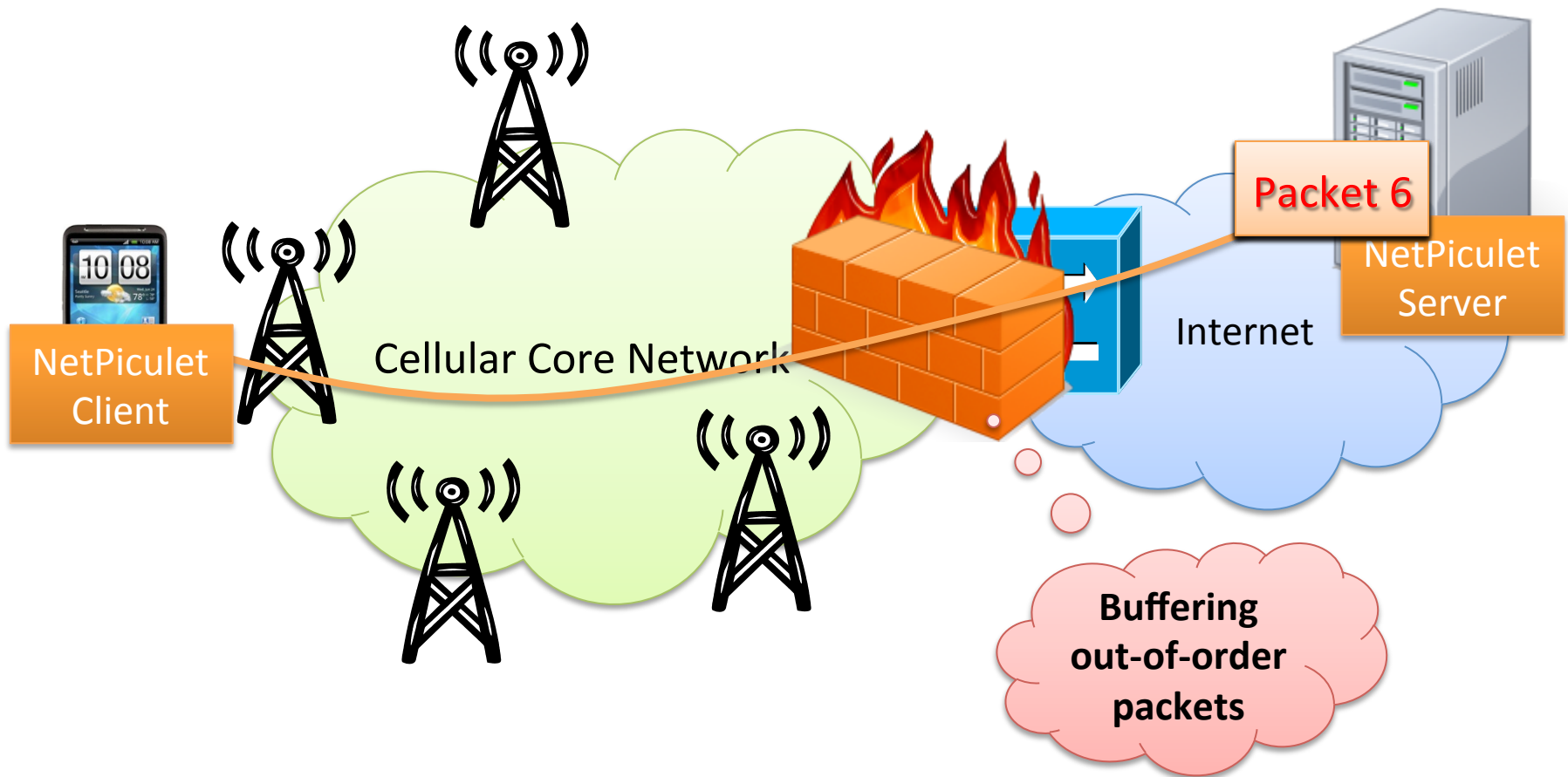
Courtesy: Z. Wang et al.

22

# Short timers drain your batteries

- Assume a long-lived TCP connection, a battery of 1350mAh
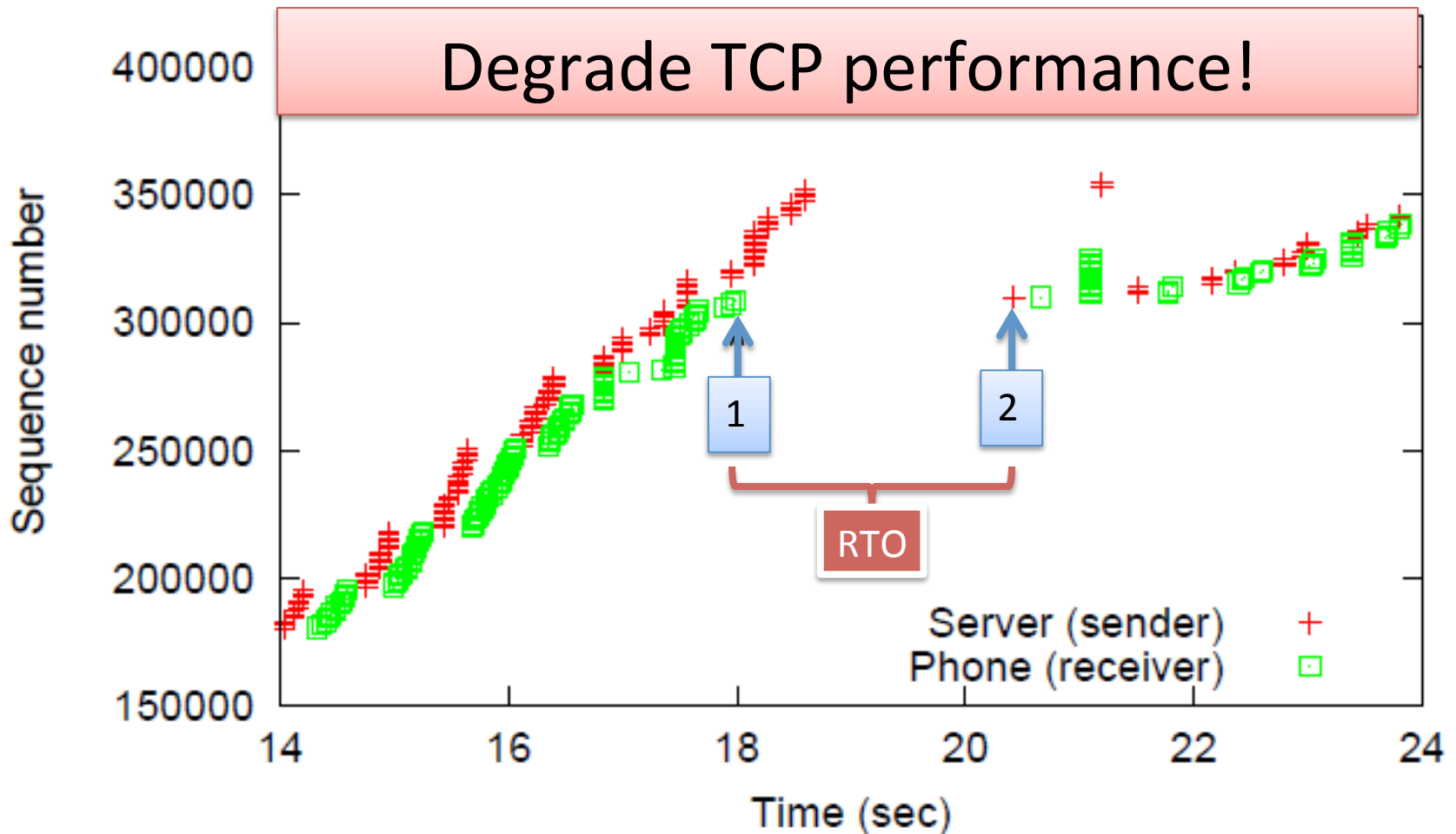- How much battery on keep-alive messages in one day?

Courtesy: Z. Wang et al.

# Outline

1. • IP spoofing

2. • TCP connection timeout

3. • TCP out-of-order buffering
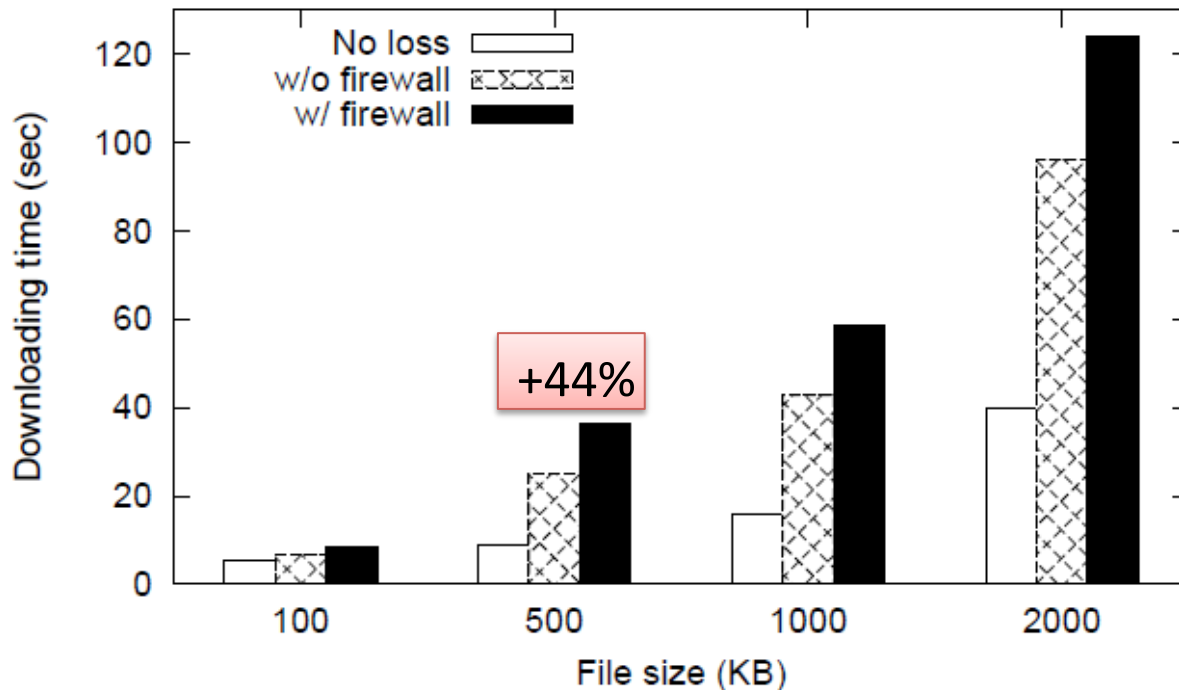
4. • NAT mapping

# TCP out-of-order packet buffering



Packet 6

NetPiculet Server

NetPiculet Client

Cellular Core Network

Internet

**Buffering out-of-order packets**

# Fast Retransmit cannot be triggered



Degrade TCP performance!

# TCP performance degradation

- ## Evaluation methodology
  - Emulate 3G environment using WiFi
  - 400 ms RTT, loss rate 1%



Longer downloading time

↓

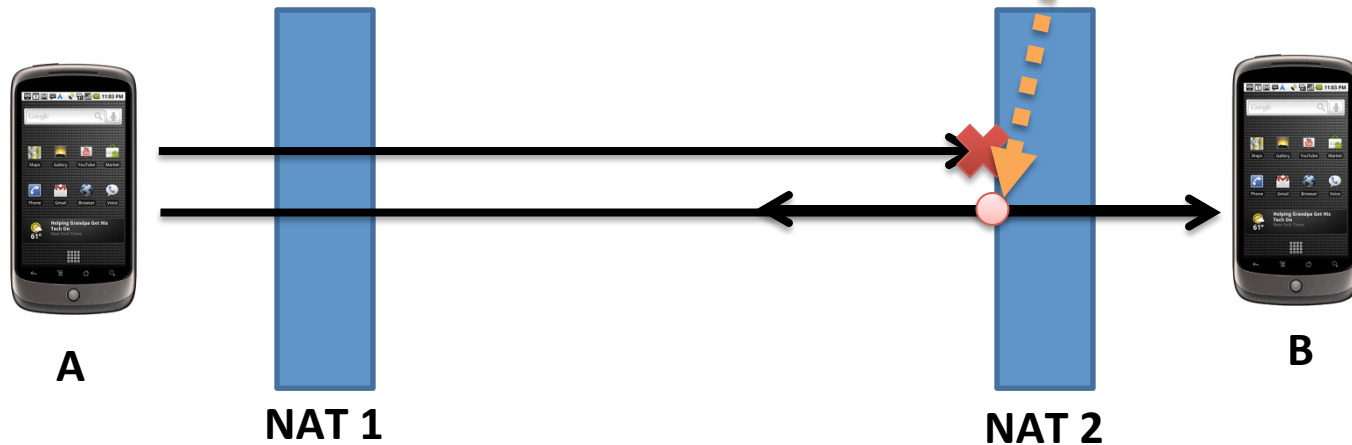More energy consumption

Courtesy: Z. Wang et al.

# Outline

1. • IP spoofing

2. • TCP connection timeout

3. • TCP out-of-order buffering

4. • NAT mapping

# NAT mapping is critical for NAT traversal
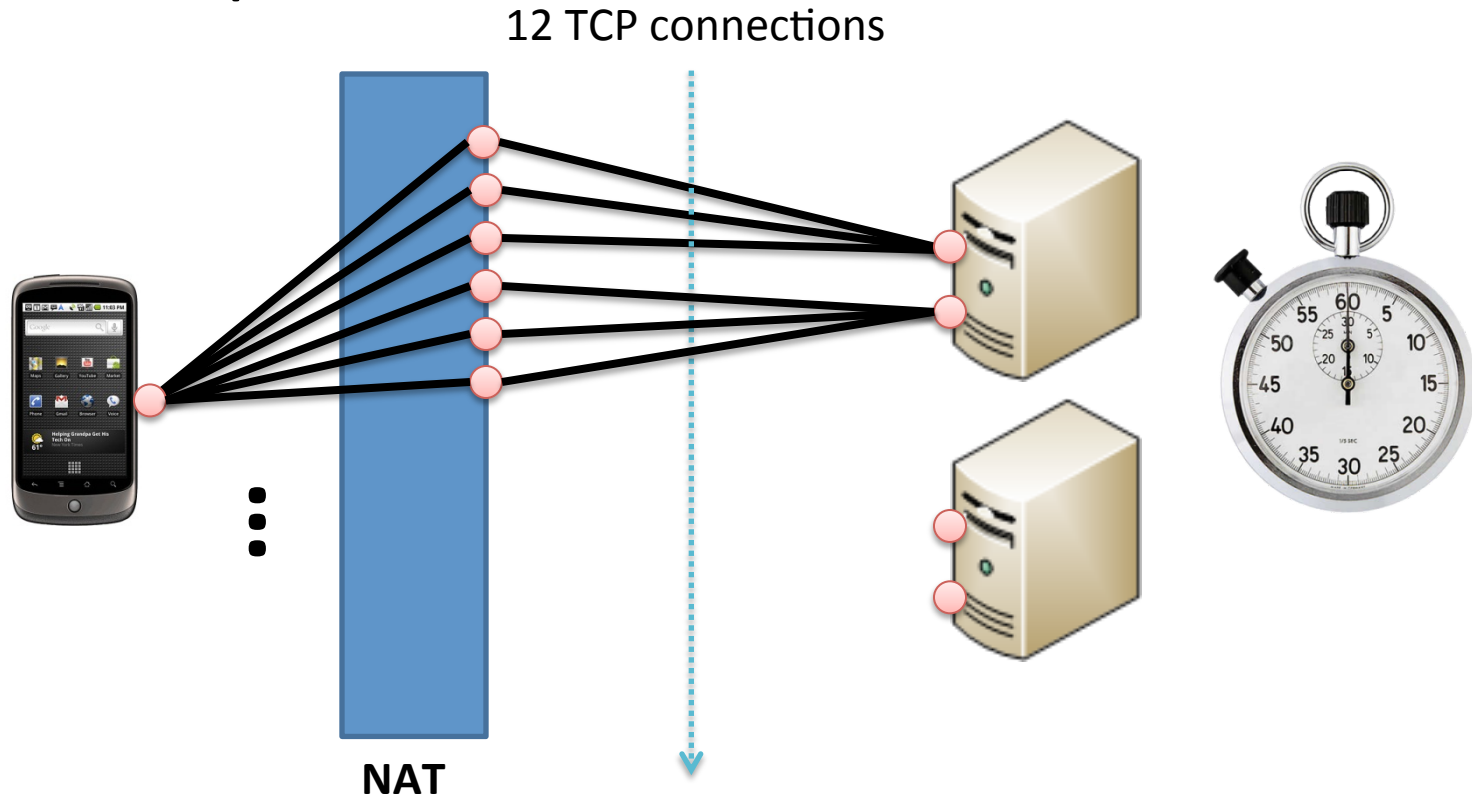


Use NAT mapping type for port prediction

A

NAT 1

NAT 2

B

# What is NAT mapping type?

- NAT mapping type defines how the NAT assign external port to each connection

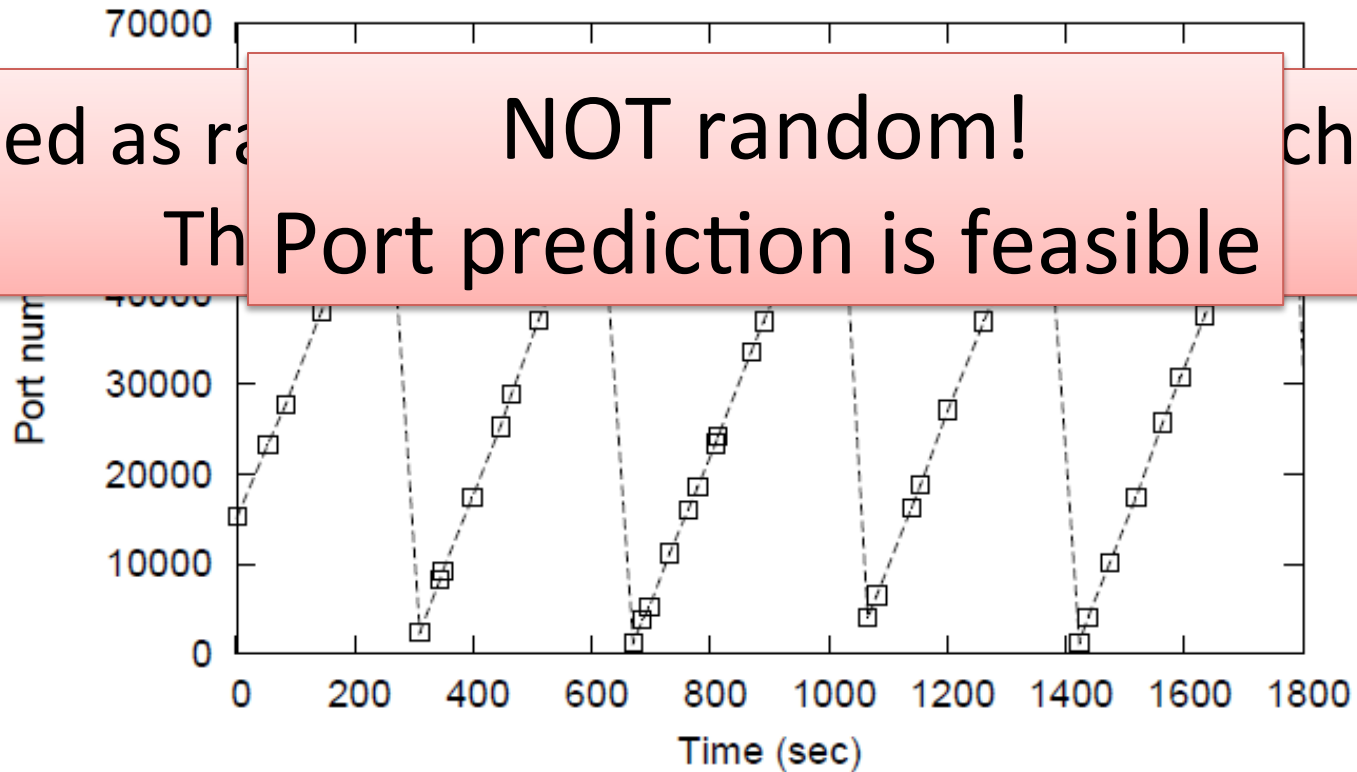12 TCP connections



**NAT**

Courtesy: Z. Wang et al.

# Behavior of a new NAT mapping type

- Creates TCP connections to the server with random intervals

- Record the observed source port on server



Treated as ra ... chniques
Th

NOT random!
Port prediction is feasible

Courtesy: Z. Wang et al.

# Lessons learned

| | |
|---|---|
| **Firewall** | IP spoofing creates security vulnerability<br>**IP spoofing should be disabled** |
| | Small TCP timeout timers waste user device energy<br>**Timer should be longer than 30 minutes** |
| | Out-of-order packet buffering hurts TCP performance<br>**Consider interaction with application carefully** |
| **NAT** | One NAT mapping linearly increases port # with time<br>**Port prediction is feasible** |

# Conclusion

- NetPiculet is a tool that can accurately infer NAT and firewall policies in the cellular networks

- NetPiculet has been wildly deployed in hundreds of carriers around the world

- The paper demonstrated the negative impact of the network policies and make improvement suggestions

# Cellular Data Network Infrastructure Characterization &
# Implication on Mobile Content Placement

Qiang Xu*, Junxian Huang*, Zhaoguang Wang*
Feng Qian*, Alexandre Gerber++, Z. Morley Mao*

*University of Michigan at Ann Arbor
++AT&T Labs Research

# Applications Depending on IP Address

- IP-based identification is popular

    – Server selection

    – Content customization

    – Fraud detection

- Why? -- IP address has strong correlation with individual user behavior
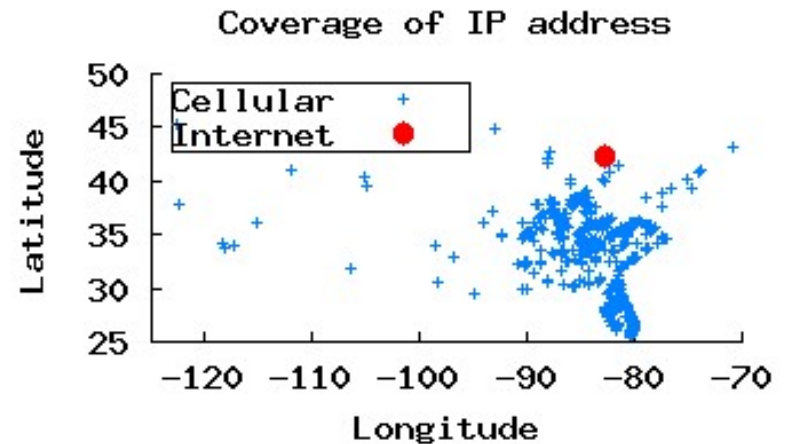
Courtesy: Q. Xu et al.

# Cellular IP Address is Dynamic

- Cellular devices are hard to geo-locate based on IP addresses

  – One Michigan's cellular device's IP is located to

| IP Address | Country | Region | City |
|------------|---------|--------|------|
| | UNITED STATES | PENNSYLVANIA | DOYLESTOWN |
| | Net Speed | | ISP |

IP2Location™ Live Product Demo

MaxMind GeoIP City/ISP/Organization Edition Results

| Hostname | Country Code | Country Name | Region | Region Name | City |
|----------|--------------|--------------|--------|-------------|------|
| 166.137.136.51 | US | United States | NY | New York | New York |

- /24 cellular IP addresses are shared across disjoint regions

Coverage of IP address

Courtesy: Q. Xu et al.

# Problem Statement

- Discover the cellular infrastructure to explain the diverse geographic distribution of cellular IP addresses and investigate the implications accordingly



* The first several IP hops are in GGSN data center
* Cellular IP addresses are allocated by GGSN data center
* GGSN data centers could be far away due to wireless hops

- – The number of GGSN data centers
- – The placement of GGSN data centers
- – The prefixes of individual GGSN data centers

Courtesy: Q. Xu et al.

# Challenges

- Cellular networks have limited visibility
  - The first IP hop (i.e., GGSN) is far away -- lower aggregation levels of base station/RNC/SGSN are transparent in *TRACEROUT*
  - Outbound *TRACEROUTE* -- private IPs, no DNS information
  - Inbound *TRACEROUTE* -- silent to ICMP probing

- Cellular IP addresses are more dynamic [BALAKRISHNAN *et al.*, IMC 2009]
  - One cellular IP address can appear at distant locations
  - Cellular devices change IP address rapidly

# Solutions

- Collect data in a new way to get geographic coverage of cellular IP prefixes
  - Build Long-term and nation-wide data set to cover major carriers and the majority of cellular prefixes
  - Combine the data from both client side and server side

- Analyze geographic coverage of cellular IP addresses to infer the placement of GGSN data centers
  - Discover the similarity across prefixes in geographic coverage
  - Cluster prefixes according to their geographic coverage

# Previous Studies

- Cellular IP dynamics
  - Measured cellular IP dynamics at two locations [Balakrishnan *et al.*, IMC 2009]

- Network infrastructure
  - Measured ISP topologies using active probing via TRACEROUTE [Spring et al., SIGCOMM 2002]

- Infrastructure's impact on  applications
  - Estimated geo-location of Internet hosts using network latency [Padmanabhan et al., SIGMETRICS 2002]
  - On the Effectiveness of DNS-based Server Selection [Shaikh et al., INFOCOM 2001]

# Outline

- Motivation

- Problem statement

- Previous Studies

- **Data Sets**

- Clustering Prefixes

- Validating the Clustering Results

- Implication on mobile content placement

# Data Sets

▸ DataSource1 (server logs): a location search server

- ▸ millions of records
- ▸ IP address, GPS, and timestamp

```
...
timestamp      lat.   long.    address
1251781217   36.75  -119.75
166.205.130.244
1251782220   33.68  -117.17  208.54.4.78
```

▸ DataSource2 (mobile app logs): an application deployed on iPhone OS, Android OS, and Windows Mobile OS

- ▸ 140k records
- ▸ IP address and carrier

```
device:
    <ID:C7F6D4E78020B14FE46897E9908F83B>
    <Carrier: AT&T>
address:
    <GlobalIP: 166.205.130.51>
...
```

▸ RouteViews: BGP update announcements

- ▸ BGP prefixes and AS number

```
...|95.140.80.254|31500|166.205.128.0/17|31500 3267 3356 7018 20057|...
...|95.140.80.254|31500|208.54.4.0/24|31500 3267 3356 21928|...
```

# Map Prefixes to Carriers & Geographic Coverage

- Correlate these data sets to resolve each one's limitations to

| DataSource1 | | | |
|---|---|---|---|
| **address** | **lat.** | **long.** | |
| 166.205.130.244 | 36.75 | -119.75 | |
| 208.54.4.11 | 33.68 | -117.17 | |

| RouteViews |
|---|
| **prefix** |
| 166.205.128.0/17 |
| 208.54.4.0/24 |

| DataSource2 | |
|---|---|
| **address** | **carrier** |
| 166.205.130.51 | **AT&T** |
| 208.54.4.11 | **T-Mobile** |

| prefix | lat. | long. |
|---|---|---|
| 166.205.128.0/17 | 36.75 | -119.75 |
| 208.54.4.0/24 | 33.68 | -117.17 |

| prefix | carrier |
|---|---|
| 166.205.128.0/17 | AT&T |
| 208.54.4.0/24 | T-Mobile |

| prefix | carrier | lat. | long. |
|---|---|---|---|
| 166.205.128.0/17 | AT&T | 36.75 | -119.75 |
| 208.54.4.0/24 | T-Mobile | 33.68 | -117.17 |

# Outline

- Motivation

- Problem statement

- Previous Studies

- Data Sets

- **Clustering Prefixes**

- Validating the Clustering Results

- Implication on mobile content placement

Courtesy: Q. Xu et al.

# Motivation for Clustering -- Limited Types of Geographic



- Prefixes with the same geographic coverage should have the same allocation policy (under the same GGSN)

Courtesy: Q. Xu et al.

# Cluster Cellular Prefixes

- 1. Pre-filter out those prefixes with very few records (<u>todo</u>)

- 2. Split the U.S. into N square grids (<u>todo</u>)

- 3. Assign a feature vector for each prefix to keep # records in each grid

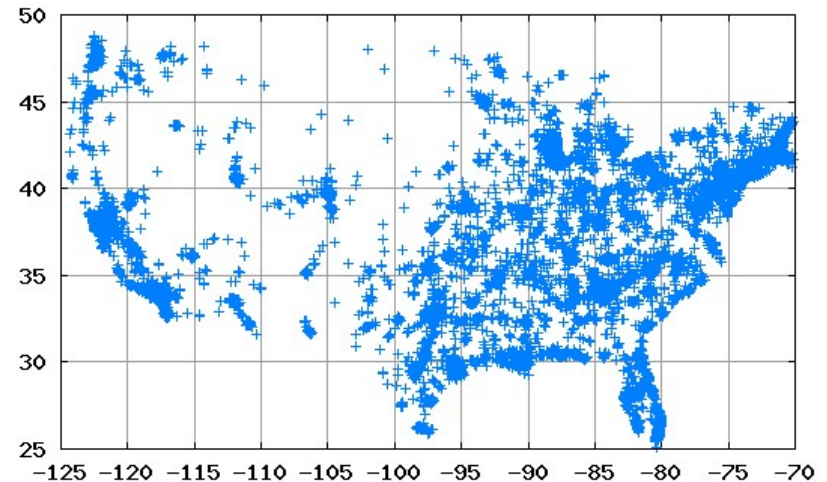- 4. Use bisect k-means to cluster prefixes by their feature vectors (<u>todo</u>)

▸ How to avoid aggressive filtering?
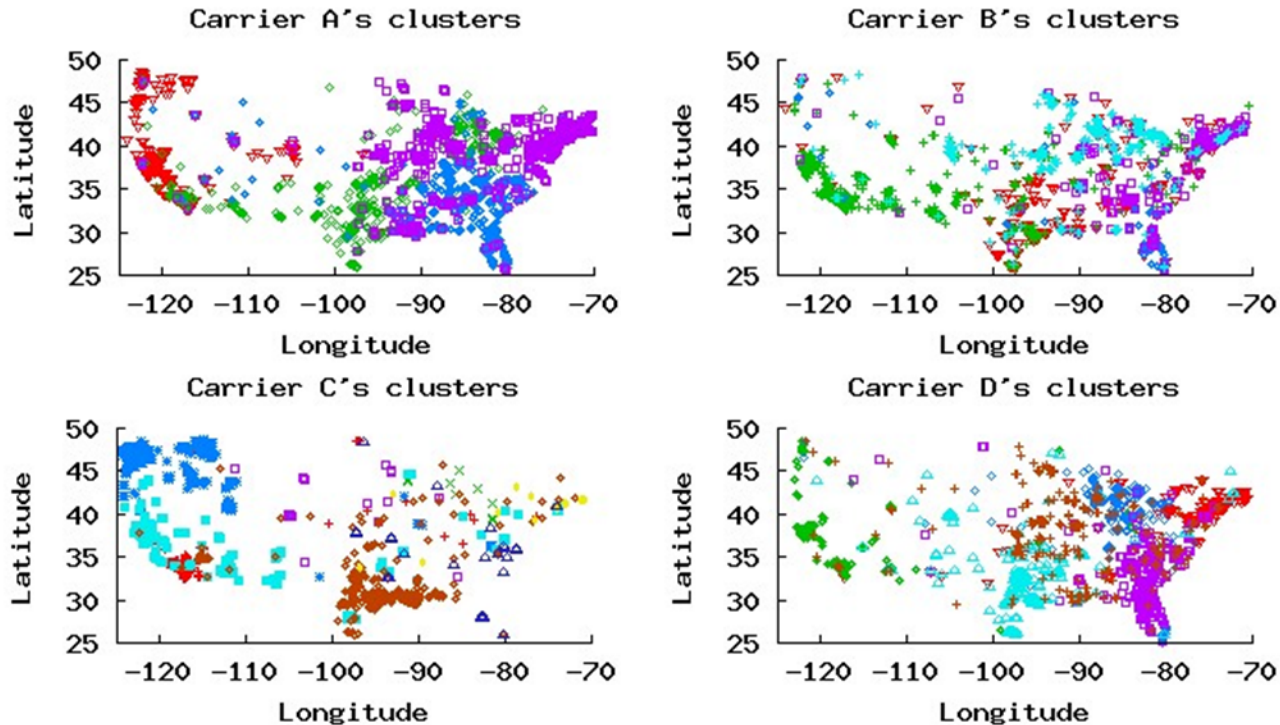  ▸ keep at least 99% records

▸ How to choose N?
  ▸ # clusters is not affected by N while N > 15 && N < 150
    ▸ The geographic coverage of each cluster is coarse-grained

▸ How to control the maximum tolerable SSE?

Courtesy: Q. Xu et al.

# Clusters of the Major Carriers



All 4 carriers cover the U.S. with only a handful clusters (4-8)
- All clusters have a large geographic coverage
- Clusters have overlap areas
  - Users commute across the boundary of adjacent clusters
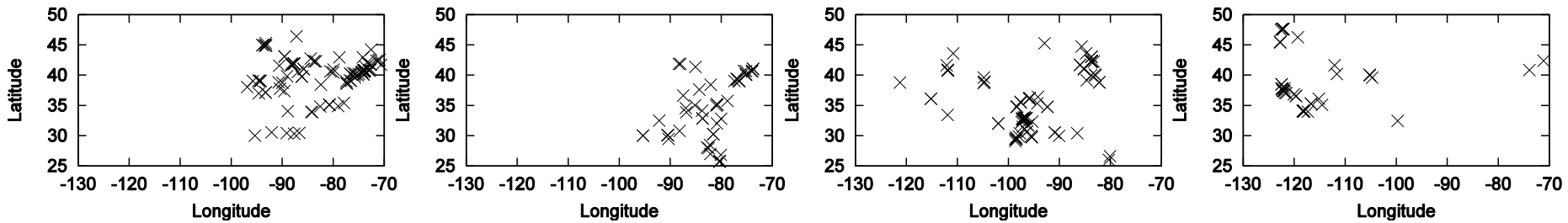  - Load balancing

# Outline

- Motivation

- Problem statement

- Previous Studies

- Data Sets

- Clustering Prefixes

- **Validating the Clustering Results**

- Implication on mobile content placement

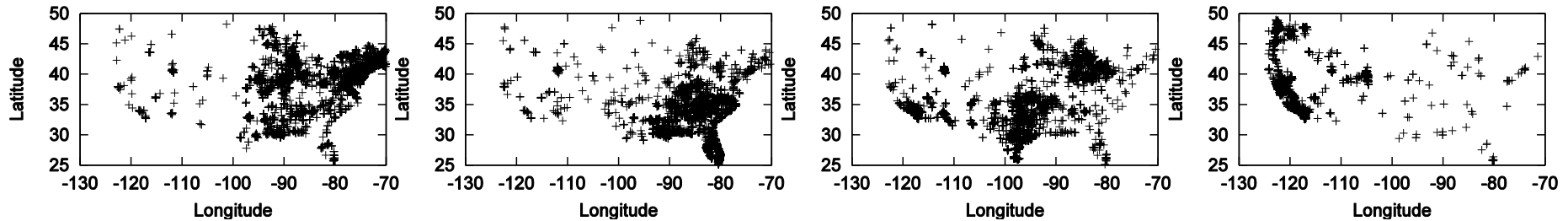# Validate via local DNS Resolver (DataSource2)

- Identify the local DNS resolvers

  - Server side: log the incoming DNS requests on the authoritative DNS resolver of **eecs.umich.edu** and record (id_timestamp, local DNS resolver)

- Profile the geographic coverage of local DNS resolvers

  - Device side: request **id_timestamp.eecs.umich.edu** and record the (id_timestamp, GPS)

Courtesy: Q. Xu et al.

# Validate via Cellular DNS Resolver (Cont.)

- Clusters of Carrier A's local DNS resolvers



- Clusters of Carrier A's prefixes



Courtesy: Q. Xu et al.

# Clustering Results

- Goal -- "...discover the cellular infrastructure to explain the diverse geographic distribution of cellular IP addresses..."
  - All 4 major carriers have only a handful (4-8) GGSN data centers
  - Individual GGSN data centers all have very large geographic coverage
- Goal -- "...investigate the Implications accordingly..."
  - Latency sensitive applications may be affected
    - CDN servers may not be able close enough to end users
    - Applications based on local DNS may not achieve higher resolution than GGSN data centers

# Outline

- Motivation

- Problem statement

- Previous Studies

- Data Sets

- Clustering Prefixes

- Validating the Clustering Results

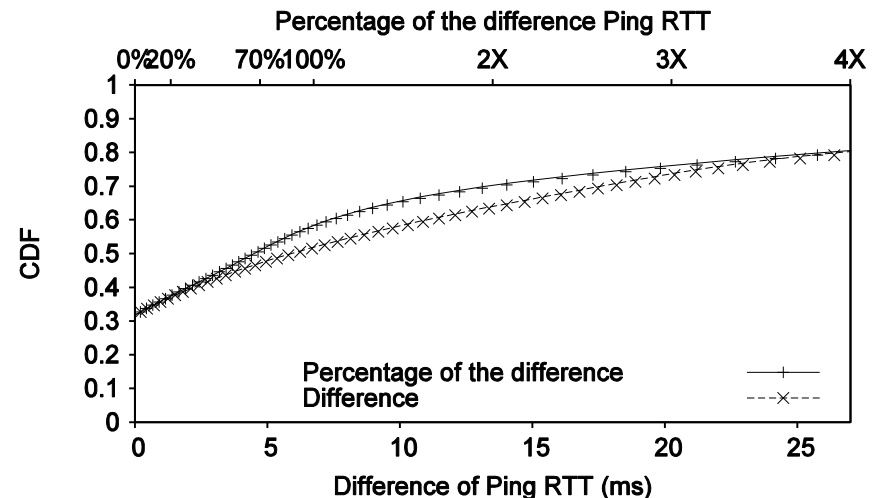- **Implication on mobile content placement**

# Routing Restriction:
## How to Adapt Existing CDN service to Cellular?

- Where to place content?
  - Along the wireless hops: require infrastructure support
  - Inside the cellular backhaul: require support from cellular providers
  - On the Internet: limited benefit, but how much is the benefit?
- Which content server to select?
  - Based on geo-location: finer-grained location may not available
  - Based on GGSN: location of GGSN

# Server Selection (DataSource2)

- Approximately locate the server with the shortest latency

  – Based on IP address

  – Based on application level information, e.g., GPS, ZIP code, etc.

- Compare the latency to the Landmark server (1) <span style="color:orange">closest to device</span> with the latency to the Landmark server (2) <span style="color:orange">closest to the GGSN</span>

  – Estimate the location of GGSN based on *TRACEROUT*

▸ Select the content server based on GGSN!



Courtesy: Q. Xu et al.

# Contributions

- ## Methodology
  - Combine routing, client-side, server-side data to improve cellular geo-location inference
  - Infer the placement of GGSN by clustering prefixes with similar geographic coverage
  - Validate the results via *TRACEROUTE* and cellular DNS server.

- ## Observation
  - All 4 major carriers cover the U.S. with only 4-8 clusters
  - Cellular DNS resolvers are placed at the same level as GGSN data centers

- ## Implication
  - Mobile content providers should place their content close to GGSNs
  - Mobile content providers should select the content server closest to the GGSN

# Questions?