# Back-Office Web Traffic on the Internet

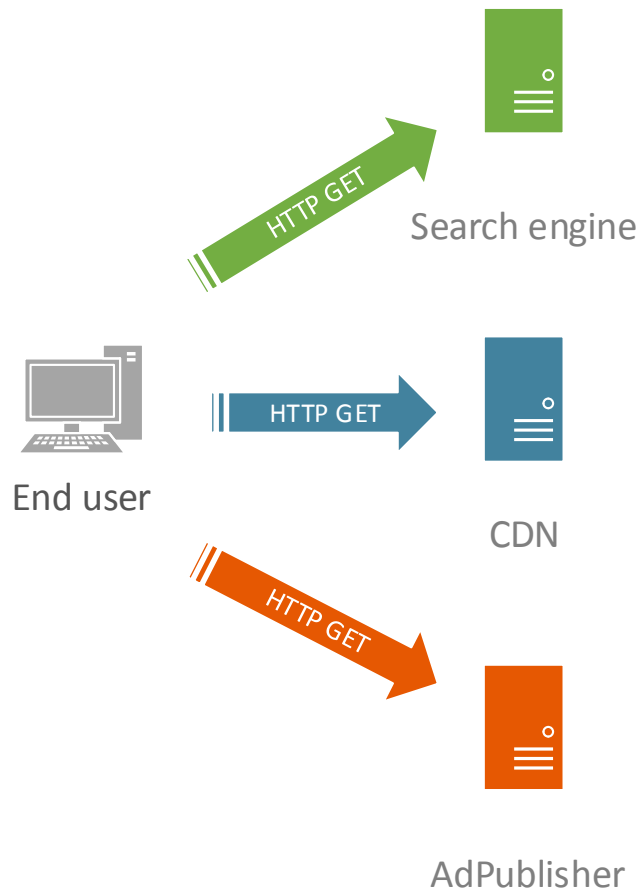Enric Pujol                                   TU-Berlin

Philipp Richter                               TU-Berlin

Balakrishnan Chandrasekaran        Duke University

Georgios Smaragdakis                   MIT / TU-Berlin / Akamai

Anja Feldmann                              TU-Berlin

Bruce Maggs                                Duke University / Akamai
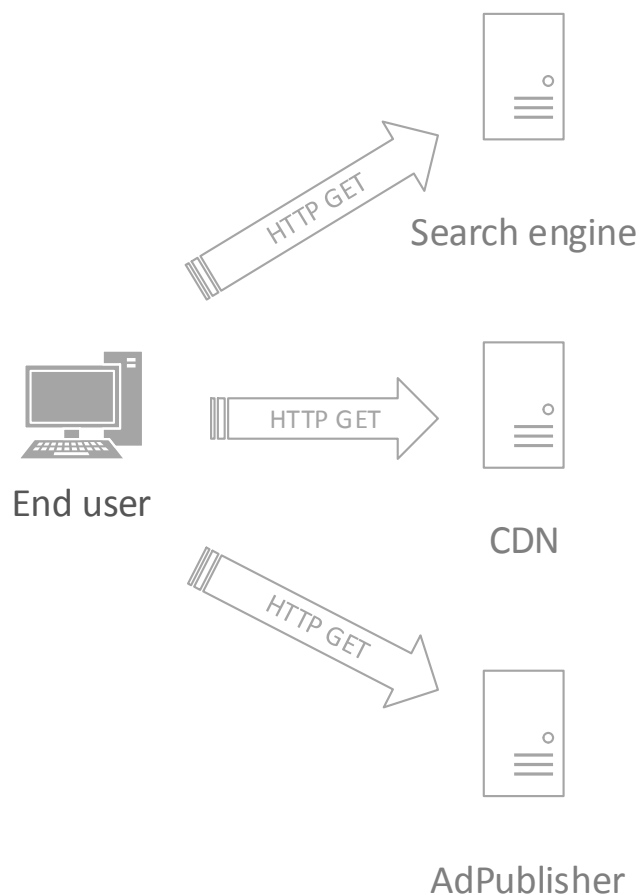
Keung-Chi Ng                               Akamai

# The Web for an end user



**Front-office Web traffic:**
Web traffic between end users and servers

# Behind the scenes...

HTTP GET

Search engine
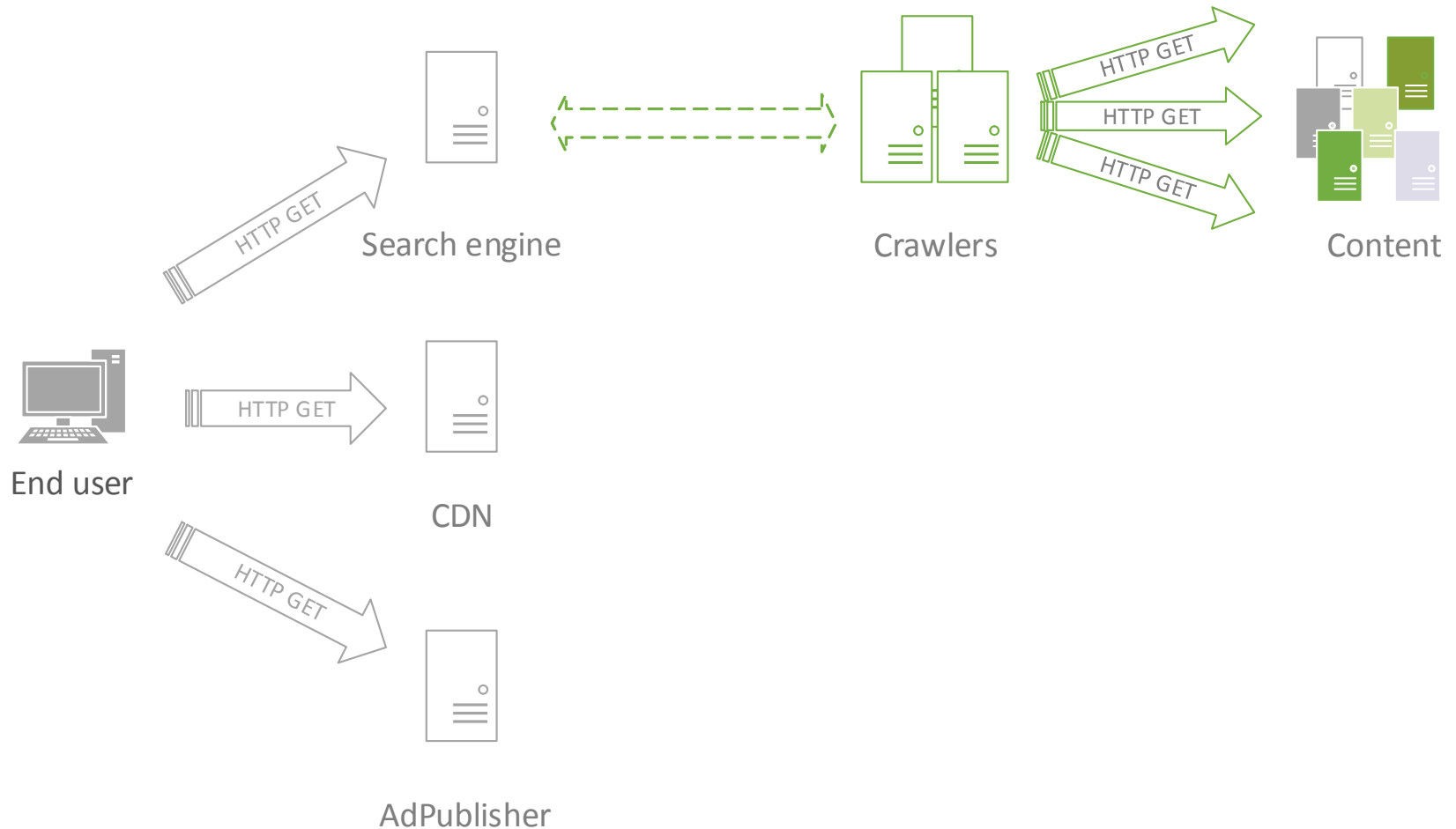
End user

HTTP GET

CDN

HTTP GET

AdPublisher

**Back-office Web traffic:**
Machine-to-machine Web traffic

?

| The front-office | The back-office |

# Search engines: crawlers



HTTP GET

Search engine

Crawlers

HTTP GET
HTTP GET
HTTP GET

Content

End user

HTTP GET

CDN

HTTP GET

AdPublisher

The front-office

The back-office

# Content delivery: proxies



**HTTP GET** — Search engine ← → Crawlers → **HTTP GET** / **HTTP GET** / **HTTP GET** → Content

End user — **HTTP GET** → CDN — **HTTP GET** → Overlay of proxies — **HTTP GET** → Origin

**HTTP GET** → AdPublisher

| The front-office | The back-office |
|---|---|

# AdExchanges: real-time bidding



The front-office

The back-office

Internet Measurement Conference 2014

# Agenda

# Vantage points (VP)

| Type | VP | Daily traffic | Observations |
|---|---|---:|---|
| IXPs | L-IXP | 11,900 TB | SFlow (1/16K) |
| | M-IXP | 1,580 TB | |
| Transit | BBone-1 | 40 TB | Packet sampled (1/1K) |
| | BBone-2 | 70 TB | |
| Content | CDN | 350 TB | 5 locations |
| Eyeballs | RBN | 35 TB | Packet dumps |

## Diverse vantage points: multiple perspectives

# Candidate IPs for the back-office



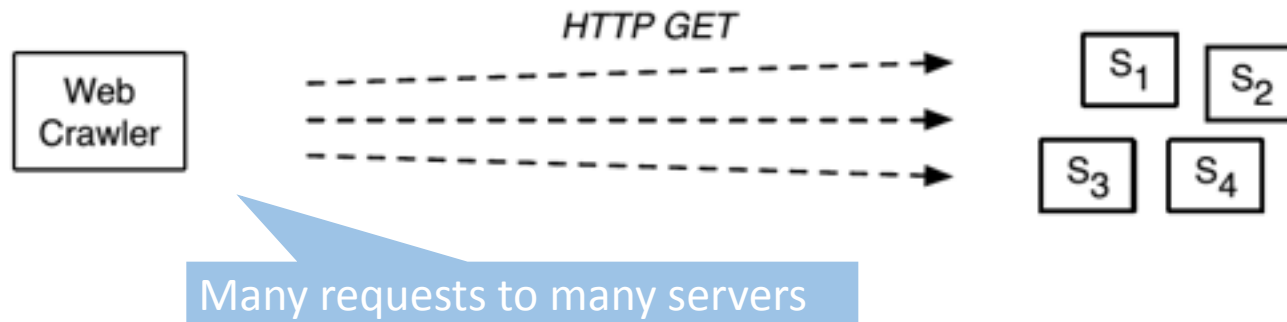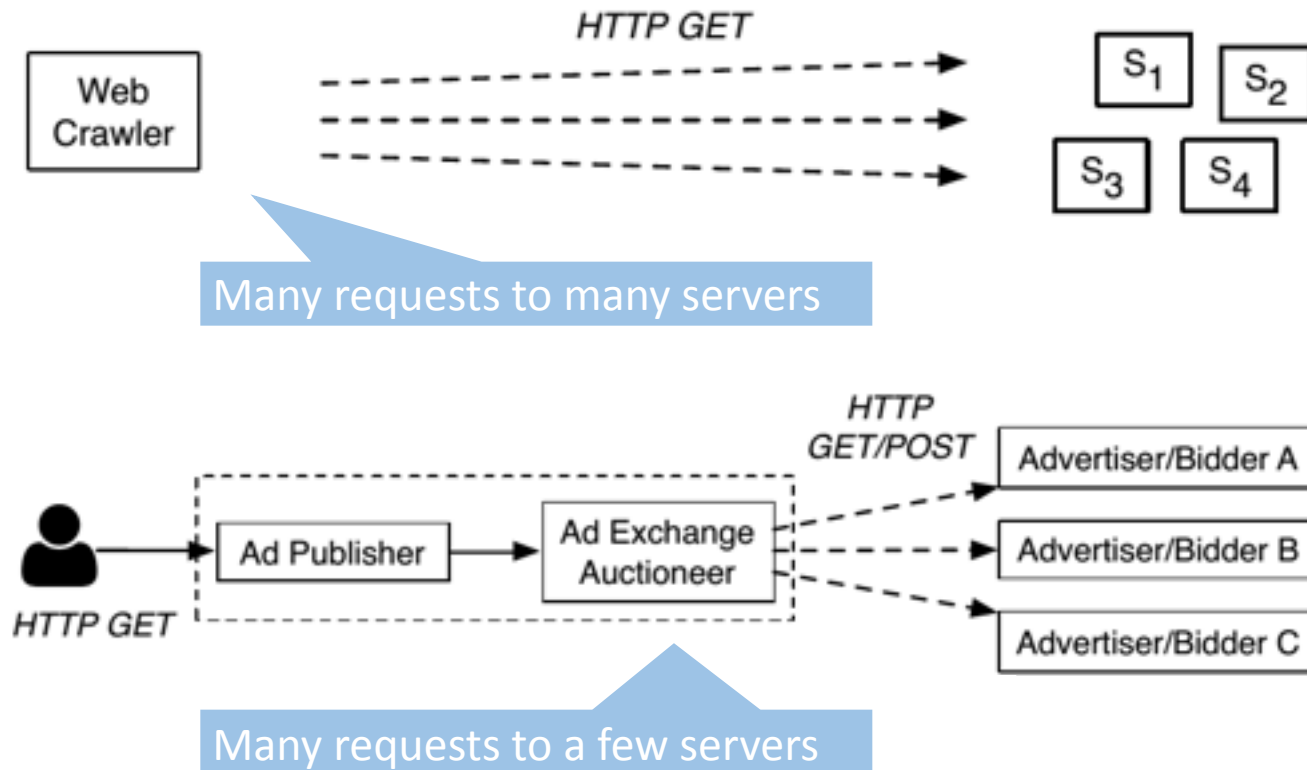HTTP GET/POST          HTTP GET/POST

Send and receive requests

Dual role IPs are prime candidates

# Candidate IPs for the back-office



Dual role IPs are prime candidates

# Candidate IPs for the back-office

HTTP GET

Web Crawler

$S_1$  $S_2$

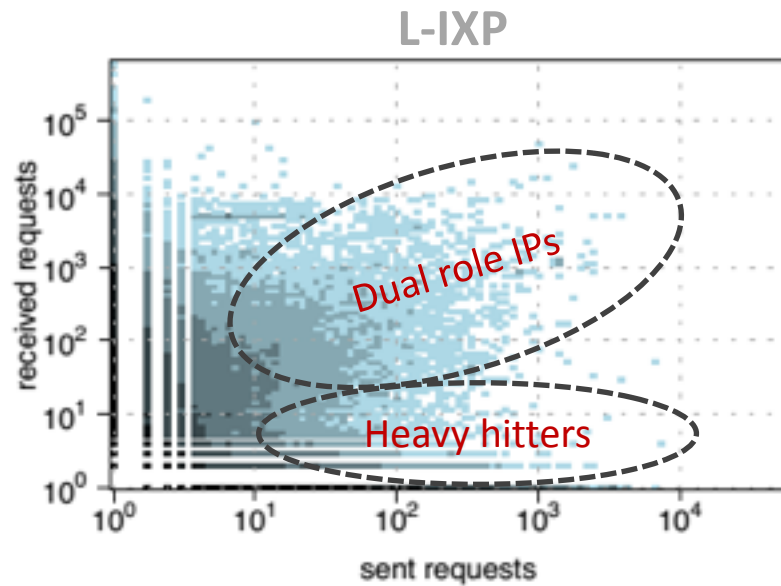$S_3$  $S_4$

Many requests to many servers

Heavy hitter IPs are also prime candidates

# Candidate IPs for the back-office



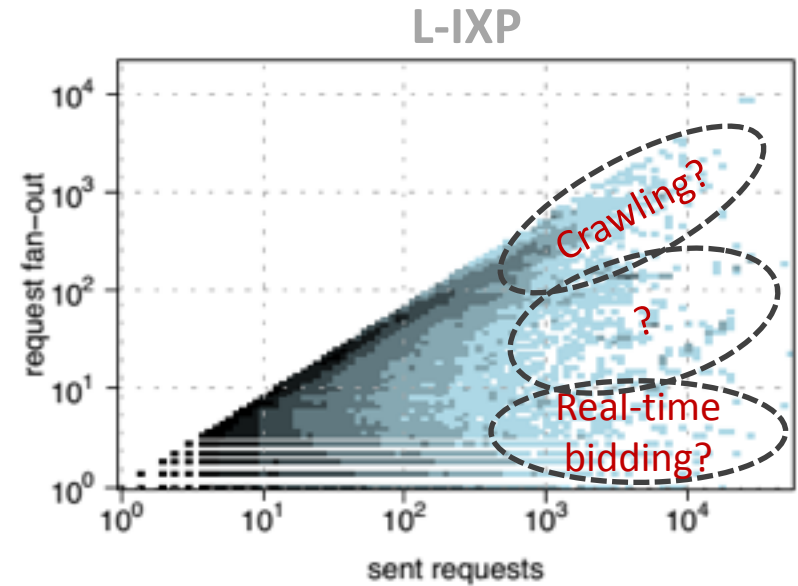Heavy hitter IPs are also prime candidates

# Sources of back-office Web traffic

**L-IXP**

Dual role IPs

Heavy hitters

# Sources of back-office Web traffic



L-IXP

*Dual role IPs*

*Heavy hitters*

L-IXP

*Crawling?*

*?*

*Real-time bidding?*

# Dual-role IPs: active measurements

| | | Client only (%) | Server only (%) | Dual-role (%) |
|---|---|---|---|---|
| **L-IXP** | **Passive** | 96.90 | 2.74 | 0.36 |
| | **Passive+Active** | 93.85 | 2.74 | **3.40** |

ZMap project: Internet-wide scan of Web Servers (scans.io)

## Observations:

1. Most IPs have only client behavior
2. Many servers also show client behavior

## Active measurements augment the number of servers

# Candidates: manual classification

Crawlers:

- Reverse DNS + Origin AS

3.9K IPs, 74% in 2 orgs

L-IXP

# Candidates: manual classification

Crawlers:
- Reverse DNS + Origin AS

Auctioneers:
- URL + Origin AS

> 3.9K IPs, 74% in 2 orgs — L-IXP

> 316 IPs, 4 orgs — L-IXP

# Candidates: manual classification

Crawlers:
- Reverse DNS + Origin AS

<div>3.9K IPs, 74% in <u>2 orgs</u></div>
L-IXP

Auctioneers:
- URL + Origin AS

<div>316 IPs, <u>4 orgs</u></div>
L-IXP

Content Delivery Proxies:
- Origin AS + Reverse DNS (for caches)

<div>36K IPs, <u>8 orgs</u></div>
L-IXP

# Candidates: manual classification

**Crawlers:**
- Reverse DNS + Origin AS

3.9K IPs, 74% in 2 orgs — L-IXP

**Auctioneers:**
- URL + Origin AS

316 IPs, 4 orgs — L-IXP

**Content Delivery Proxies:**
- Origin AS + Reverse DNS (for caches)
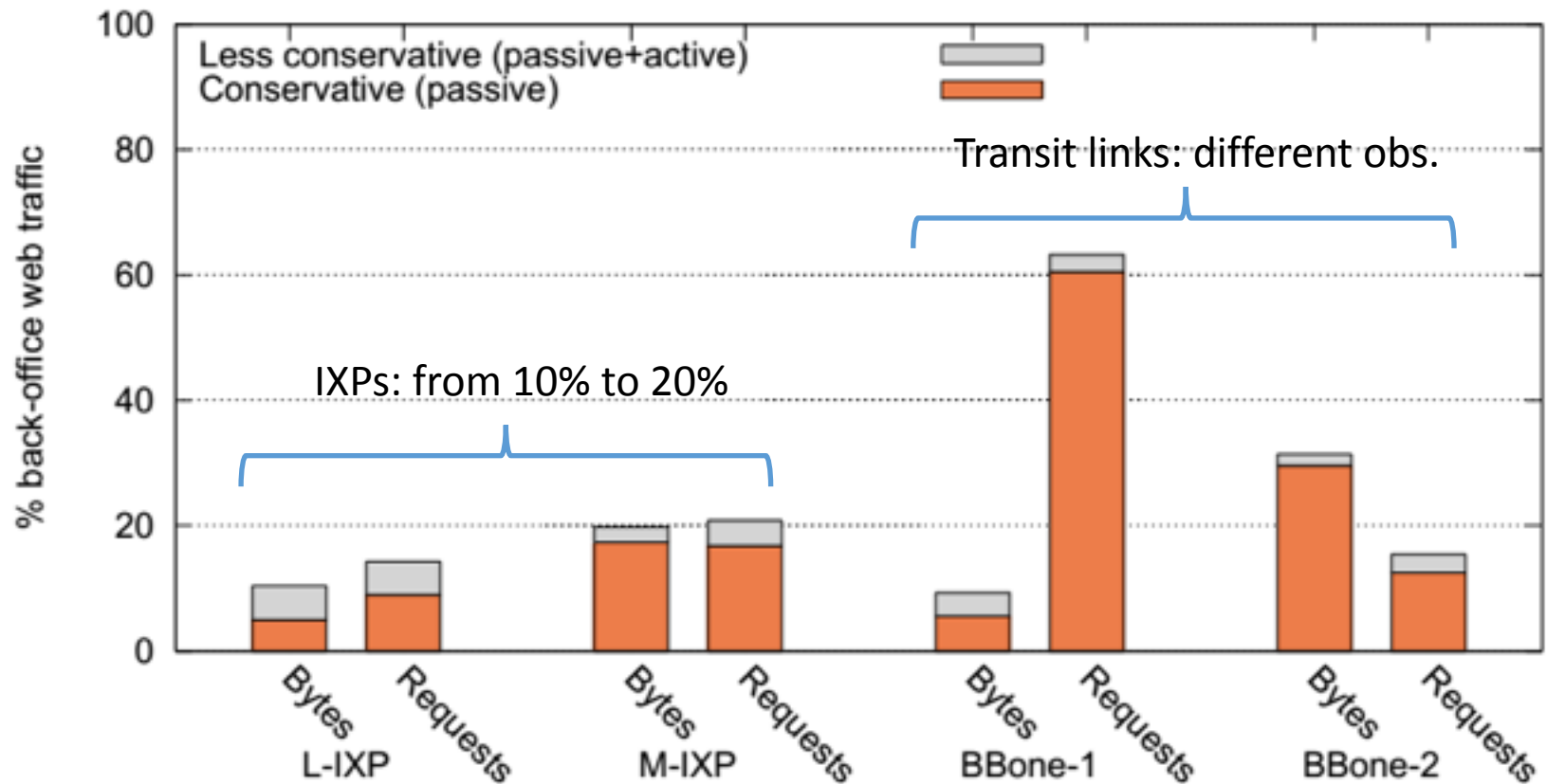
36K IPs, 8 orgs — L-IXP

**Other:**
- Rest of dual-role IPs

151K IPs, mostly in cloud prov. — L-IXP

# Agenda

# Traffic



At least 10% in our VPs

# Traffic: Contribution per class

| L-IXP | | CDPs | Auctioneers | Crawlers | Other |
|---|---|---|---|---|---|
| | **Bytes** | 12.1 % | 1.1 % | 10.3 % | 76.5 % |
| | **Requests** | 11.8 % | 22.5 % | 15.1 % | 50.6 % |

## Observations:

1.  CDPs            big players – significant share
2.  Real-time bidding   many but small transactions
3.  Crawlers        a few orgs – significant share
4.  Other          cloud service providers
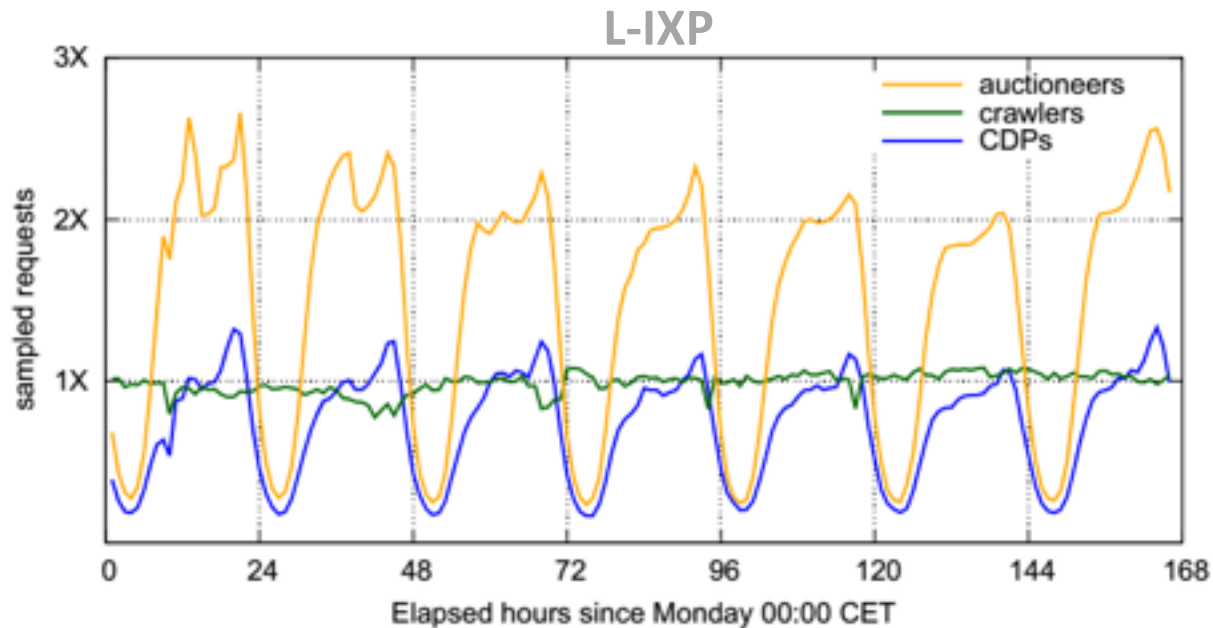
### All classes contribute. More to discover

# Traffic patterns: bytes



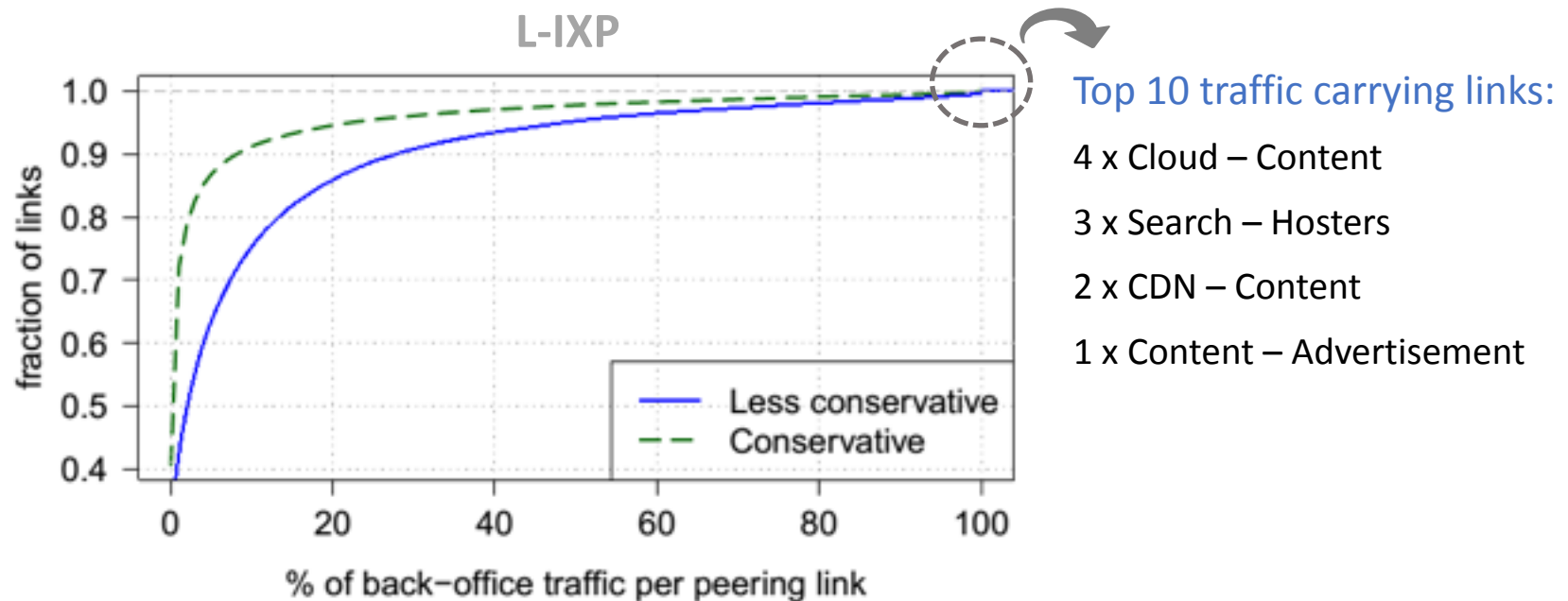% back-office Web traffic increases during off hours in IXPs

# Traffic patterns: requests



**L-IXP**

## Observations:
1. A multiplicative factor of human activity (e.g., RTB)
2. Non-human triggered activity (e.g., crawlers)

# Inter-domain perspective

**L-IXP**



Top 10 traffic carrying links:

4 x Cloud – Content

3 x Search – Hosters

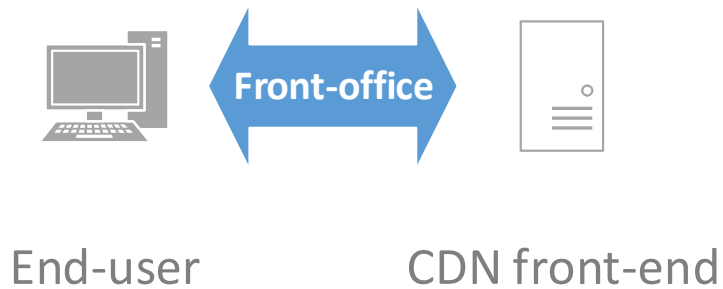2 x CDN – Content

1 x Content – Advertisement

## Back-office traffic appears in many peering links
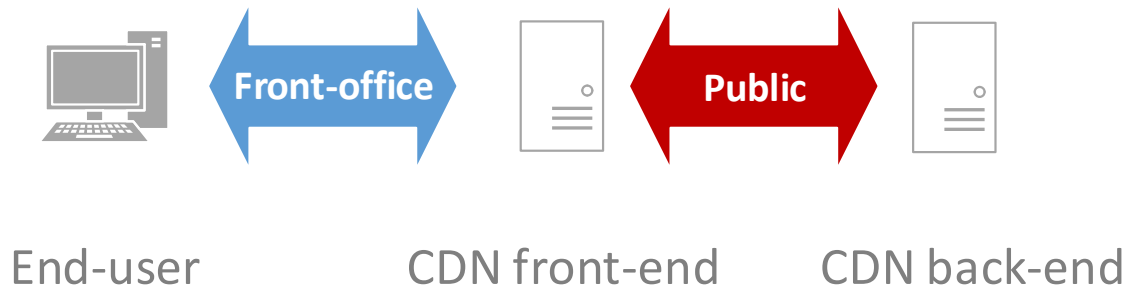
# Agenda

1. Introduction
2. Methodology and datasets
3. Characteristics
   1. Traffic
   2. Patterns
   3. Inter-domain perspective
4. CDN back-office traffic
5. The end-user perspective
6. Summary and implications

# A CDN perspective



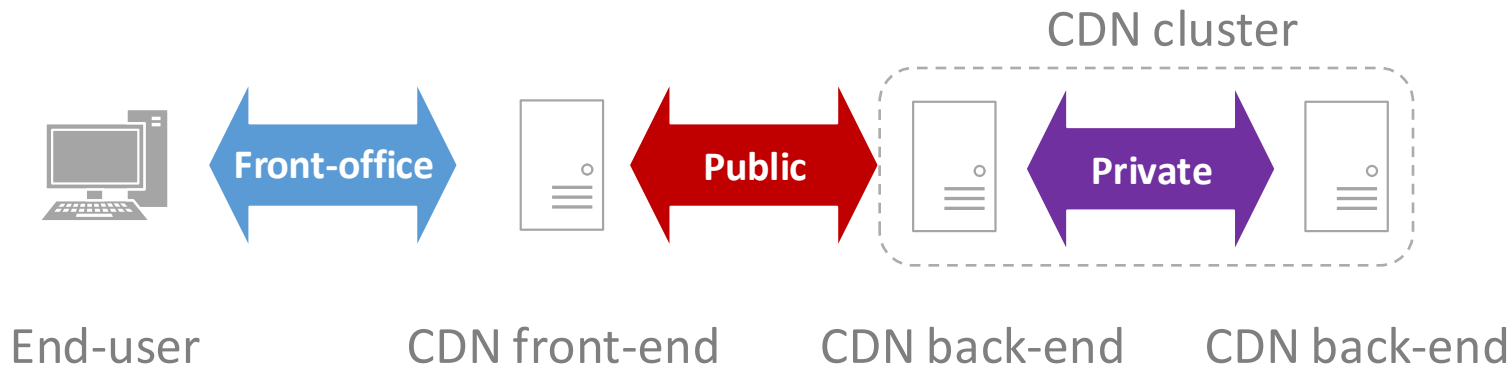End-user          CDN front-end

## Three sub-classes of back-office traffic

# A CDN perspective



End-user       CDN front-end       CDN back-end

Public: front-end back-end over the Internet

# A CDN perspective



Private: within same cluster

# A CDN perspective



CDN cluster

**Front-office** **Public** **Private** **Origin**

End-user      CDN front-end      CDN back-end      CDN back-end      Origin

Origin: inter-organization over the Internet

# Back-office per location



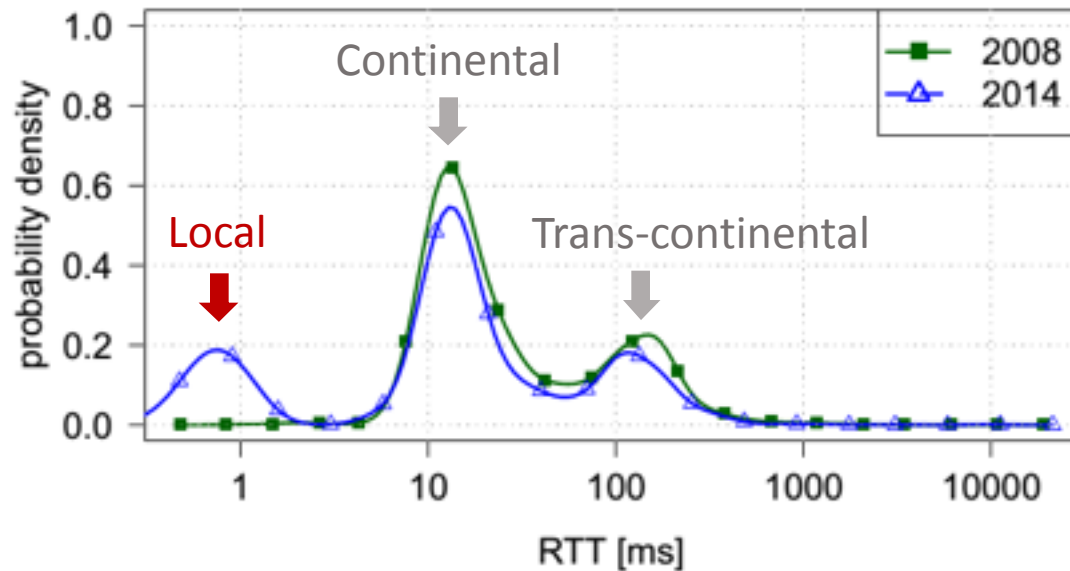CDNs heavily rely on back-office traffic

# Agenda

1. Introduction
2. Methodology and datasets
3. Characteristics
    1. Traffic
    2. Patterns
    3. Inter-domain perspective
4. CDN back-office traffic
5. **The end-user perspective**
6. **Summary and implications**

# The end-user perspective

Residential broadband network: backbone latency (no access)



A smaller front-office: but the back-office may be large

# Summary

1. A back-office to support the Web

2. Significant traffic: bytes and requests

3. Different type of traffic patterns

4. Visible at multiple peering links

   An important yet understudied class of traffic

# Implications

Feasibility to deploy new protocols:

- It is easier to change the back office than the front office

# Implications

Feasibility to deploy new protocols:

- It is easier to change the back office than the front office

Performance evaluation:

- Interactions with the back office
- More users than anticipated

# Implications

Feasibility to deploy new protocols:

- It is easier to change the back office than the front office

Performance evaluation:

- Interactions with the back office
- More users than anticipated

Opportunities:

- ISPs: micro-data centers, virtualized services
- IXPs: co-location strategies
- NSPs: new services e.g., SLAs

# Back-office traffic on the Internet



End user

Search engine

Crawlers

Content

HTTP GET

HTTP GET

CDN

Overlay of proxies

Origin

HTTP GET

AdExchange

AdPublisher

Auctioneer

Advertisers/bidders

HTTP POST

The front-office

The back-office (some examples thereof)