

## Investigating Structure of Modern Web Traffic

Noriaki Kamiyama<sup>1,2</sup>  
Kohei Shiimoto<sup>1</sup>  
Masayuki Murata<sup>2</sup>

Yuusuke Nakano<sup>1,2</sup>  
Go Hasegawa<sup>2</sup>  
Hideo Miyahara<sup>2</sup>

<sup>1</sup>NTT Network Technology Laboratories  
<sup>2</sup>Osaka University

2015. 7. 3

Copyright©2015. NTT corp. All Rights Reserved.

### Increase of Web Response Time

Complementary cumulative distribution (CCD) of response time of most popular 1,200 websites when accessing from Tokyo, Japan

- Response time of 50% websites: longer than 4 seconds
- Response time of 10% websites: longer than 9 seconds
- **Need to reducing web response time**

Copyright©2015. NTT corp. All Rights Reserved. 1

### Need to Reduce HTTP Response Time

- One website consists of multiple data objects which are delivered from various object servers using HTTP sessions.
- Reducing **HTTP response time** is one of most effective approach for reducing web response time.

HTTP response time:  
RTT plus server response time

Copyright©2015. NTT corp. All Rights Reserved. 2

### CDN: Platform Delivering Web Objects

- 74% of 1,000 most popular websites use CDN\*, and CDN is most common technique for reducing HTTP response time.

\*I. Ott, et al., Content Delivery and the Natural Evolution of DNS, ACM IMC 2012

Obtain objects from origin server and cache them

CDN provider, e.g., Akamai, provides many cache servers at network edge.

Copyright©2015. NTT corp. All Rights Reserved. 3

### Contribution of This Presentation

- Has not clarified actual effect of CDNs on reducing HTTP response time
- Important to clarify structure of modern web traffic to maximize effect of CDN

- Developed and implemented active-based framework for **measuring communication structure of web traffic**
- Measured **distance** and **RTT** to origin or cache servers from client terminals as well as **HTTP response time** of web objects
- Based on experimental results, clarified various tendencies of **geographical distribution of original objects and CDN caches**

Copyright©2015. NTT corp. All Rights Reserved. 4

### Measurement Procedure

1. Selected 12 PlanetLab hosts as measurement terminals accessing various websites
2. Measured various properties, e.g., HTTP response time and RTT, by executing program at each PlanetLab host to access various websites sequentially
3. Collected measurement results at collector terminal

Measurement location		
North America	Massachusetts	Australia
	Wisconsin	New Zealand
Europe	California	Japan
	Ireland	Ecuador
Russia	Germany	Argentina
	Russia	Reunion
		Oceania
		Asia
		South America
		Africa

(1) Sending measurement program

(2) HTTP query and measurement

(3) Data analysis

PlanetLab: overlay network consisting of over 500 hosts worldwide

Copyright©2015. NTT corp. All Rights Reserved. 5

### Measurement Program

- Generated URL list and sent it to each PlanetLab host
- Starting from 0:00 (midnight) or 12:00 (noon), each PlanetLab host executed following procedures:
  1. Accessed websites according to URL list and obtained HAR (HTTP Archive) files
  2. Extracted information of HTTP response time from obtained HAR files
  3. Measured RTT to each object server by sending ping
  4. Obtained domain name of each object server using dig command
  5. Sent measurement results to collector terminal

Copyright © 2015 NTT Corp. All Rights Reserved. 6

### Obtaining HAR Files

- Obtained HTML file initially, and obtained each object embedded in HTML file
- HAR (HTTP Archive) file: outputs various properties of each object in JSON (JavaScript Object Notation) format

Using phantomJS, providing browser function, and netsniff, extracting HAR files, obtained HAR files of many websites sequentially in batch process

Copyright © 2015 NTT Corp. All Rights Reserved. 6

### URL List of Measurement Target

- Selected 300 most popular websites in each of 16 categories based on public information of Alexa\*
- Totally Selected 927 websites from which HAR files were successfully obtained at all 12 measurement locations

\*http://www.alexa.com/topsites

Category	#sites	Category	#sites
Business	40	Home	47
Computer	91	Shopping	68
News	27	Adult	102
Reference	109	Arts	60
Regional	73	Games	58
Science	86	Kids & teens	64
Society	83	Recreation	52
Health	52	Sports	53

Copyright © 2015 NTT Corp. All Rights Reserved. 8

### Classifying Objects Based on CDN Use

- Classified objects into **CDN objects** delivered using CDN or **non-CDN objects** delivered without using CDN
- Listed 44 second-level domains of various CDN providers by manually checking websites of various CDN providers
- Obtained domain names of hosts actually delivering objects, e.g., www.akamai.com/qqq/rrr, by using dig command from URL names, e.g., www.google.com/xxx/yyy.jpg, of objects extracted from HAR files
- Identified CDN objects by comparing second-level domain obtained by dig command with entries of generated list

List of second-level domains of CDN objects

profile.ak.fbcdn.net	cloudfront.net	akamaihd.net	edgesuite.net
static.ak.fbcdn.net	vo.ms.ecnd.net	edgesuite.net	cloudfront.net
r.dnscdn.com	edgecastcdn.net	edgekey.net	vo.ms.ecnd.net
s.cdn-care.com	cdngc.net	snip.net	edgecastcdn.net
cms.cdn.static.cache.org	boots.trapcdn.com	akamitechnologies.com	cdngc.net
g-ex.images-amazon.com	example.com	akamitechnologies.fr	pus-h11.cdn.un.com
max.blurtcdn.com	akadns.net	akamatech.net	ve14.f3.at11.lnw.net
a.es.pncdn.com	akam.net	akadns.net	hs-9.cdn77.com
ex.images-amazon.com	akamaiedge.net	akam.net	nyud.net
edgekey.net	akamai.net	akamaistream.net	CloudFlare
edgesuite.net	akamaiedge.net	edgekey.net	Incapsula

Copyright © 2015 NTT Corp. All Rights Reserved. 8

### Ratio of CDN Objects

ID	Category	ID	Category
C1	Business	C5	Regional
C2	Computers	C6	Science
C3	News	C7	Society
C4	Reference	C8	Health
C9	Home	C13	Games
C10	Shopping	C14	Kids & teens
C11	Adult	C15	Recreation
C12	Arts	C16	Sports

- Having more CDN objects in websites of Computers, News, Society, Shopping, Arts, and Kids & teens
- Having fewer CDN objects in websites of Business, Regional, Science, Adult, and Games

Ratio of CDN objects differed among categories, between 0.2 and 0.5

Copyright © 2015 NTT Corp. All Rights Reserved. 10

### Average HTTP Response Time at Noon

- Non-CDN objects: exceeded 500 ms in 20 - 80% websites
- CDN objects: exceeded 500 ms in just 2 to 40% websites

Confirmed effect of CDN in reducing HTTP response time

Copyright © 2015 NTT Corp. All Rights Reserved. 11

### Clustering Analysis of Websites based on RTT

- Geographical pattern of original objects, i.e., non-CDN objects, and CDN caches delivering CDN objects will differ among access locations even when accessing same website.
- Analyzed geographical tendencies by clustering websites based on average RTT at 12 access locations
  - Applied k-means method based on vectors  $v(y)$  with elements  $v_{xy}$ , average RTT b/w access location  $x$  and objects of website  $y$ .
  - Optimally set cluster count  $k$  using 'JD method'
  - Set initial cluster using k-means++ method\*\*

\*A. L. Jain and R. C. Dubes, Algorithms for Clustering Data, Englewood Cliffs, NJ Prentice-Hall, 1989  
 \*\*D. Arthur and S. Vassilvitskii, kmeans++: the advantages of careful seeding, ACM SODA 2007

NTT Copyright©2015 NTT Corp. All Rights Reserved. 12

### Geographical Distribution of Original Objects

Clustering websites based on RTT of non-CDN objects at midnight

L1	Massachusetts	L7	Australia
L2	Wisconsin	L8	New Zealand
L3	California	L9	Japan
L4	Ireland	L10	Ecuador
L5	Germany	L11	Argentina
L6	Russia	L12	Reunion

C1	Business	C5	Regional	C9	Home	C13	Games
C2	Computers	C6	Science	C10	Shopping	C14	kids & teens
C3	News	C7	Society	C11	Adult	C15	Recreation
C4	Reference	C8	Health	C12	Arts	C16	Sports

- Cluster 1: RTT was small only in North America, and more websites of Society, Adult, Recreation, and Sports were classified. → Geographical locality is weak, and identical content are viewed from various regions.
- Cluster 3: RTT was small in all areas except Africa, and more websites of Home and Shopping were classified. → Geographical locality is strong, and unique content are viewed in each region.

NTT Copyright©2015 NTT Corp. All Rights Reserved. 13

### Geographical Distribution of Cache Servers

Clustering websites based on RTT of CDN objects at midnight

L1	Massachusetts	L7	Australia
L2	Wisconsin	L8	New Zealand
L3	California	L9	Japan
L4	Ireland	L10	Ecuador
L5	Germany	L11	Argentina
L6	Russia	L12	Reunion

C1	Business	C5	Regional	C9	Home	C13	Games
C2	Computers	C6	Science	C10	Shopping	C14	kids & teens
C3	News	C7	Society	C11	Adult	C15	Recreation
C4	Reference	C8	Health	C12	Arts	C16	Sports

- Cluster 1 & 2 & 3: RTT was small in all regions except South America and Africa, and more than 80% of websites of all categories except Adult were classified. → Using CDN deploying cache servers worldwide
- Cluster 4 & 5: RTT was small in North America and Europe, and more Adult websites were classified. → Using CDN deploying cache servers in specific areas

NTT Copyright©2015 NTT Corp. All Rights Reserved. 14

### Conclusion

- Actively measured structure of web traffic from 12 locations in world
  - Classified web objects into non-CDN and CDN objects
  - Measured HTTP response time and RTT of each object
- Ratio of CDN objects differed among categories, between 0.2 and 0.5
- Confirmed effect of CDN in reducing HTTP response time
  - Non-CDN objects: exceeded 500 ms in 20 - 80% websites
  - CDN objects: exceeded 500 ms in just 2 to 40% websites
- Revealed various geographical pattern of original objects
  - Websites with weak locality: identical objects from North America
  - Websites with strong locality: unique objects in each region
- Revealed various geographical pattern of CDN caches
  - Most websites used CDN deploying cache servers worldwide

NTT Copyright©2015 NTT Corp. All Rights Reserved. 15

### Example of HAR File

HAR file of www.google.com

NTT Copyright©2015 NTT Corp. All Rights Reserved. 16

### クラスタリング手法

- k-means法: 非階層型クラスタリング手法の一つで、クラスタの重心を用いて、各要素を k 個のクラスタに分類
  - 各要素を重心の距離が最も近いクラスタに分類する処理をクラスタが収束するまで反復

- k-means++法: 距離の離れた要素を初期クラスタの重心に設定することで、分類精度を向上
  - ランダムに一つの要素を選び、クラスタ重心に設定
  - 各要素 x に関して、その最近傍重心との距離  $D(x)$  を計算
  - $D(x)^2$  に比例する確率に従い、新しいクラスタ重心としてランダムに一つ要素を選択
  - k 個のクラスタ重心が選択されるまで上記処理を反復
  - 以後は k-means法を用いてクラスタを生成

NTT Copyright©2015 NTT Corp. All Rights Reserved. 17

### クラスタ数 k の最適選定

- Jain-Dubes法\*を用いて最適なクラスタ数 k を設定
  - 要素数が n のときに、 $2 \leq k \leq 1 + \log_2 n$  の範囲で各クラスタ数 k のクラスタリングを実施
  - 次式で定義されるコスト p(m) が最小となる k を選択

$$p(k) = \frac{1}{k} \sum_{i=1}^k \max_{1 \leq j \leq k} \left\{ \begin{matrix} \eta_i + \eta_j \\ \xi_{ij} \end{matrix} \right\}$$

$$\eta_j = \frac{1}{n_j} \sum_{i=1}^{n_j} D(x_i^{(j)}, m_j) \quad \xi_{ij} = D(m_i, m_j)$$

$x_i^{(j)}$ : クラスタ j 内の i 番目の要素,  $n_j$ : クラスタ j の要素数  
 $m_j$ : クラスタ j の重心,  $D(a,b)$ : ベクトル a と b 間の距離

- 各クラスタに属する要素のクラスタ重心に対する距離 A の平均値の、二つのクラスタの重心間の距離 B に対する比率を、最小化することに相当

\*A. K. Jain and R. C. Dubes, Algorithms for clustering data, Prentice-Hall, 1988

Copyright © 2015 NTT Corp. All Rights Reserved. 18

### Basic Properties

ID	Category	Website count	Object size	Object count	Total size
		(0.00)	(kbytes)		(Mbytes)
C1	Business	50	40	14.70	55.14
C2	Computers	112	91	16.26	43.63
C3	News	39	27	13.55	72.45
C4	Reference	112	109	13.09	43.42
C5	Regional	80	73	17.77	50.59
C6	Science	95	86	14.04	52.86
C7	Society	79	83	15.01	66.86
C8	Health	86	52	14.27	54.30
C9	Home	85	47	15.66	55.39
C10	Shopping	69	68	15.67	70.77
C11	Adult	112	102	10.49	53.94
C12	Arts	55	60	15.43	68.18
C13	Games	87	58	15.28	54.12
C14	Kids & teens	106	64	13.23	54.59
C15	Recreation	86	52	13.55	57.30
C16	Sports	38	53	16.62	86.67

- Entertainment websites, e.g., Arts, Shopping, and Sport, tend to have more objects and larger total data size.
- Information websites, e.g., Business, Computers, Health, and Reference, tend to have fewer objects and smaller total data size.

Copyright © 2015 NTT Corp. All Rights Reserved. 19

### 平均WaitのCCDを比較(midnight)

- 約20~80%のサイトはNon-CDN-Objの平均HTTP応答時間は500msを超過
- CDN-Objの平均HTTP応答時間が500msを超過するサイトは約5%~20%

NTT 広域アクティブ測定により、CDN-Objとnon-CDN-Objの地理的な配置傾向を分析し、施策&IIIのHTTP応答時間改善のポテンシャルを確認

Copyright © 2015 NTT Corp. All Rights Reserved. 20

### HTTP応答時間改善の施策

- CP(content provider)がCDN使用オブジェクト(CDN-Obj)の比率を増加**  
 現状、CDNを用いずに配信しているObj(non-CDN-Obj)をCDNで配信し、CDN-Objの比率を増加させることでHTTP応答時間を低減
- CPが使用CDNを変更**  
 CDN事業者のキャッシュの地理的な配置、キャッシュサーバ能力、NWスループット、キャッシュ制御ポリシー、等が異なることから、よりHTTP応答時間の低減効果が見込めるCDN事業者にCPが契約先を切替
- CDN事業者がObjのキャッシュ位置や配信サーバ選択を適正化**  
 CDN事業者がキャッシュサーバ配置場所、キャッシュサーバ制御ポリシー(置換方式、キャッシュ判断方式)、配信キャッシュサーバ選択を適正化し、ユーザに近いキャッシュサーバから配信することでHTTP応答時間を低減

本発表: 広域アクティブ測定により、CDN-Objとnon-CDN-Objの地理的な配置傾向を分析し、施策&IIIのHTTP応答時間改善のポテンシャルを評価

Copyright © 2015 NTT Corp. All Rights Reserved. 21

### 取得データ

- 各受信オブジェクトに対して、HARファイルから以下の情報を抽出(GeolPのAPIを用いてホスト名から都市名と座標を取得)

データ項目名	Key
HTTP 送信先ホスト名	"request"."url"
ホストの存在する国名	GeoIP: "country_name"
ホストの存在する都市名	GeoIP: "city"
ホストの緯度	GeoIP: "latitude"
ホストの経度	GeoIP: "longitude"
サイズ (byte)	"response"."content"
総遅延時間 (ms)	"time"
コネクション確立時間 (ms)	"timings"."blocked"
DNS 名前解決時間 (ms)	"timings"."dns"
TCP コネクション確立時間 (ms)	"timings"."connect"
HTTP リクエスト転送時間 (ms)	"timings"."send"
サーバ応答待ち時間 (ms)	"timings"."wait"
レスポンス転送時間 (ms)	"timings"."receive"
SSL/TLS 時間 (ms)	"timings"."ssl"
MIME Type	"response"."content"
	"mimeType"

- digコマンドを用いて、実際に各オブジェクトを配信したサーバのドメイン名を取得し、さらにpingを送付してRTTを計測

Copyright © 2015 NTT Corp. All Rights Reserved. 22

### Overview of Measurement and Analysis

- 世界の12の地点から約1,000のWebサイトにアクセスした際の、配信サーバ距離、RTT、HTTP応答時間等の各種通信特性を測定
- 測定データを、CDN-Objとnon-CDN-Objとに分離して、さらにWebサイトのジャンルごとに、平均値や累積分布を算出
- 12の各測定地点における各特性値に基づきWebサイトをクラスタ分析することで、各WebサイトジャンルのCDN-Objとnon-CDN-Objの地理的な配置傾向を分析

Copyright © 2015 NTT Corp. All Rights Reserved. 23

