



## Priority Control Based on Website Categories in Edge Computing

Noriaki Kamiyama<sup>1,2</sup> Yuusuke Nakano<sup>1,2</sup>  
 Kohei Shiomoto<sup>1</sup> Go Hasegawa<sup>2</sup>  
 Masayuki Murata<sup>2</sup> Hideo Miyahara<sup>2</sup>

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

2016. 4. 11

Copyright©2016 NTT corp. All Rights Reserved.

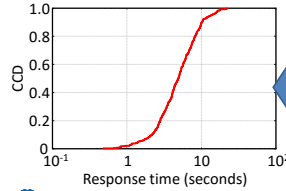
## Increase of Web Response Time



- Web response time\*: longer than 5 seconds in 50% webpages, and longer than 10 seconds in 10% webpages
- Amazon increased revenue 1% for every 0.1 second reduction in web response time.\*\*
- Need to reducing web response time**

\*Web response time: waiting time after clicking hyperlink until entire part of webpage is shown

\*\*R. Kohavi and R. Longbotham, Online Experiments: Lessons Learned, IEEE Computer, Vol.40, No. 9, pp.103-105, Sep. 2007.



Complementary cumulative distribution (CCD) of web response time of most popular 1,000 websites when accessing from Tokyo, Japan, in June 2015



Copyright©2016 NTT corp. All Rights Reserved. 1

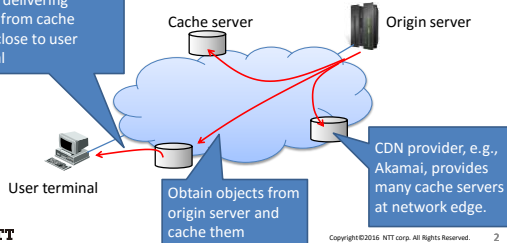
## CDN: Platform Delivering Web Objects



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

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

Reducing HTTP response time by delivering objects from cache server close to user terminal



Copyright©2016 NTT corp. All Rights Reserved. 2

## Edge Computing

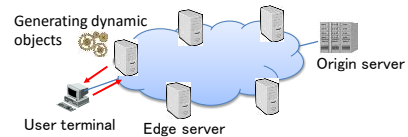


- Many objects are dynamically generated in modern webpages.
- Edge computing\* is effective to deliver dynamic objects efficiently.

\*M. Rabinovich, et al., Computing on the Edge: A Platform for Replicating Internet Applications," WCW 2003. A. Davis, et al., EdgeComputing: Extending Enterprise Applications to the Edge of the Internet, WWW 2004.

Edge servers located at edge nodes

- Caches application codes for generating dynamic objects
- Dynamically generates objects for user requests

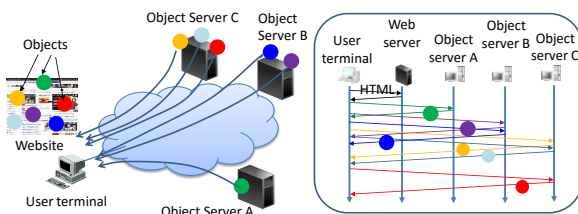


Copyright©2016 NTT corp. All Rights Reserved. 3

## Complexity of Web Traffic Pattern



One website consists of multiple data objects which are delivered from various object servers using HTTP sessions.



Copyright©2016 NTT corp. All Rights Reserved. 4

## Effect of Edge Computing



Geographical deployment pattern may differ among website categories, e.g., Sports and News, and effect of edge computing will depend on website categories.

Yahoo Answers, McAfee SiteAdvisor, ...

Yelp, Groupon, ...



Identical content from North America

Unique content at each region

Contribution of this work

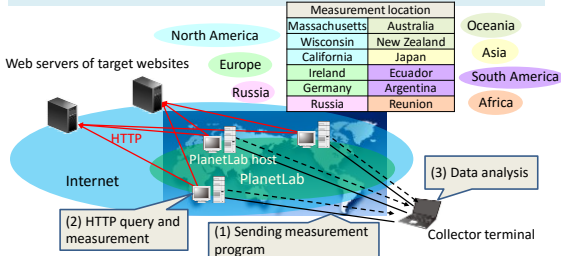
- Propose to differentiate caching priority among website categories
- Roughly analyze effect of category-based priority control in edge computing using active measurement data from 12 locations in world



Copyright©2016 NTT corp. All Rights Reserved. 5

## Measurement Procedure

- Selected 12 PlanetLab hosts as measurement terminals accessing various websites
- Measured various properties, e.g., object count obtained and RTT, by executing program at each PlanetLab host to access various websites sequentially
- Collected measurement results at collector terminal



PlanetLab: overlay network consisting of over 500 hosts worldwide  
Copyright ©2016 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 measurement data 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 ©2016 NTT corp. All Rights Reserved. 7

## 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.msecnd.net	edgesuite.net	cloudfront.net
r.dnscdn.com	edgecastcdn.net	edgekey.net	vo.msecnd.net
s.cdn-care.com	cdnge.net	srtp.net	edgecastcdn.net
cmscdn.staticcache.org	bootstrapcdn.com	akamaithechnologies.com	cdnge.net
g-ex.images-amazon.com	example.com	akamaithechnologies.fr	push-11.cdn.sun.com
max.blurtitcdn.com	akadns.net	akamai.net	ve14.fr3.atl1.linw.net
a.spcdn.com	akam.net	akadns.net	hs-9.cdn77.com
ecx.images-amazon.com	akamaiedge.net	akam.net	myud.net
edgekey.net	akamai.net	akamaistream.net	CloudFlare
edgesuite.net	akamaiedge.net	edgekey.net	Incapsula

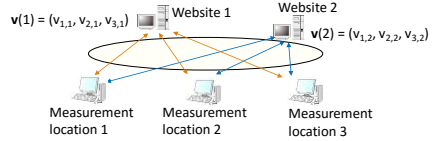
Copyright ©2016 NTT corp. All Rights Reserved. 8

## Clustering Analysis of Webpages 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 webpage 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, 1988

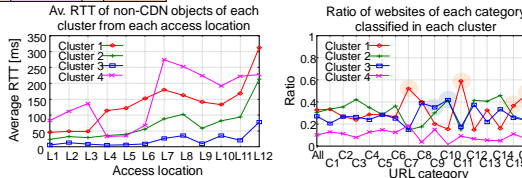
\*\*D. Arthur and S. Vassilvitskii, k-means++: the advantages of careful seeding, ACM SODA 2007



Copyright ©2016 NTT corp. All Rights Reserved. 9

## 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

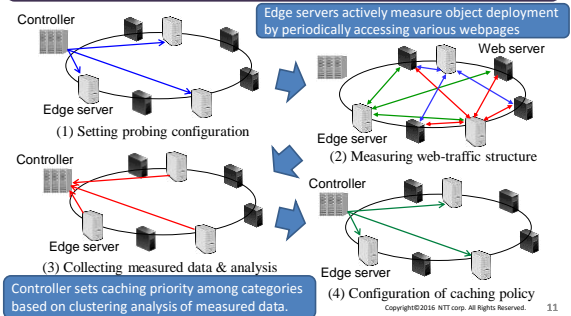


- Cluster 1: RTT was small only in North America. ⇒ **Geographical locality is weak**, and identical content are viewed from various regions.
- Cluster 3: RTT was small in all areas except Africa. ⇒ **Geographical locality is strong**, and unique content are viewed in each region.

Confirmed different tendencies of object deployment among web categories 10

## Platform Measuring Deployment of Web Objects

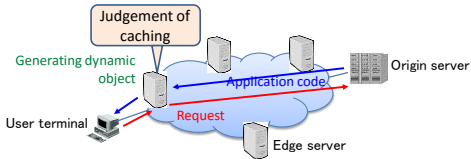
Sets caching priority of each web category at edge servers by continuously measuring geographical deployment of web objects from edge servers in world



Copyright ©2016 NTT corp. All Rights Reserved. 11

### Priority Control among Categories in Edge Computing

Each edge server autonomously makes caching judgement based on web categories according to policy set by controller.



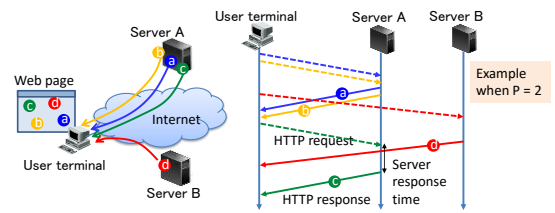
- Controller obtains URL list of website categories using Alexa Web Information Service API\* and sends it to edge servers.
- When cache miss, edge server obtains application code of dynamic object from origin server and judges whether to store code based on URL list.

\*<http://aws.amazon.com/jp/awis/>  
Copyright ©2016 NTT corp. All Rights Reserved. 12



### Roughly Estimating Web Response Time (1)

To investigate potential of differentiating caching policy among web categories in edge computing, roughly evaluate reduction effect of web response time

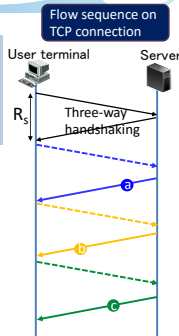


- Number of parallel sessions established with one web server is limited below P
- P = 2 (suggested in the HTTP/1.1 specification)
- P = 4 (Safari 3, Opera 9)
- P = 6 (Explore 8, Firefox 3)

### Roughly Estimating Web Response Time (2)

#### Assumption

- Starts obtaining objects on all TCP co. with all servers
- Fairly obtains objects over all TCP co. with each server
- Continuously receives objects on each TCP connection
- Obtains each object from edge servers with probability H with zero RTT



$D_x$ : estimated time reduced by delivering objects of webpage x from edge servers

$$D_x = \max_{s \in S_x} \left\{ \left\lceil \frac{M_s H}{P} \right\rceil + \lfloor H \rfloor \right\} R_s$$

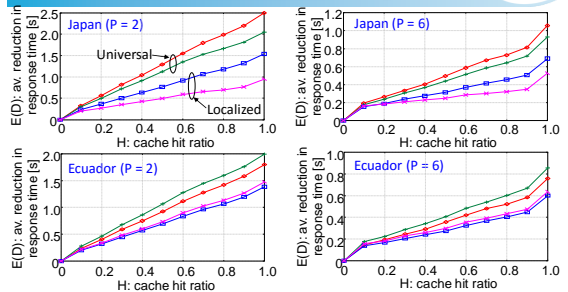
- $S_x$ : set of servers sending objects of page x
- $M_s$ : number of objects obtained from server s
- $R_s$ : average RTT b/w user terminal and sever s

Apply measured value

Copyright ©2016 NTT corp. All Rights Reserved. 14



### Average Reduction in Response Time of Four Categories



Confirm difference of E(D) between Universal websites (adult, society) and Localized websites (home, shopping)

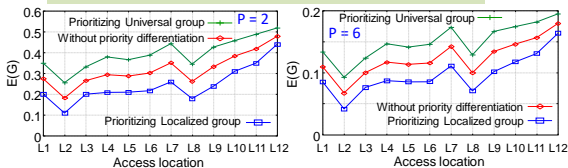
Adult, Society, Home, Shopping



Copyright ©2016 NTT corp. All Rights Reserved. 15

### Effect of Web Category Differentiation in Edge Computing

E(G), average reduction ratio of web response time:  
 $E(G) = E(D) / (\text{average response time without edge computing})$



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

- Compare E(G) among three caching policies:
  - Without priority differentiation: delivering 50% of objects of each category
  - Prioritizing universal group: delivering all objects of Adult and Society webpages
  - Prioritizing localized group: delivering all objects of Home and Shopping webpages

Can improve effect of edge computing by prioritizing universal webpages

16



### Conclusion

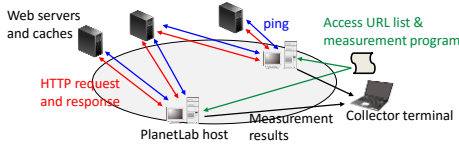
- Actively measured RTT and object count of most popular 1,000 webpages from 12 locations in world using PlanetLab
- Confirmed difference of geographical tendencies of object deployment among website categories
- Universal websites: Adult and Society
- Localized websites: Home and Shopping
- Proposed to differentiate caching priority among web categories in edge computing
- Roughly estimated reduction effect of web response time by edge computing
- Numerical confirmed effect of differentiating caching priority among web categories in edge computing



Copyright ©2016 NTT corp. All Rights Reserved. 17

## Measurement Program

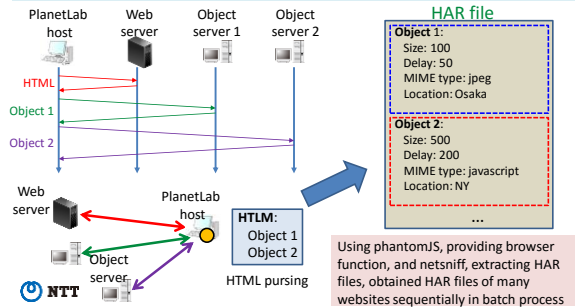
- Generated URL list and sent it to each PlanetLab host
- Starting from 0:00 (midnight) or 12:00 (noon), each PlanetLab host executed following procedures:
  1. Accessed websites according to URL list and obtained HAR (HTTP Archive) files
  2. Extracted information of HTTP response time from obtained 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©2016 NTT corp. All Rights Reserved. 18

## 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

## Example of HAR File

```

{
  "browser": "phantomjs",
  "pages": [
    {
      "url": "http://www.google.com/",
      "totalSize": 1234567,
      "objects": [
        {
          "url": "http://www.google.com/images/logo.png",
          "size": 1024,
          "mimeType": "image/png",
          "location": "Osaka"
        },
        {
          "url": "http://www.google.com/js/jquery.js",
          "size": 5000,
          "mimeType": "text/javascript",
          "location": "NY"
        }
      ]
    }
  ]
}
    
```

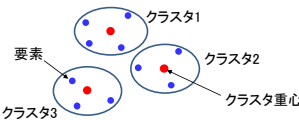
HAR file of www.google.com



Copyright©2016 NTT corp. All Rights Reserved. 20

## 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



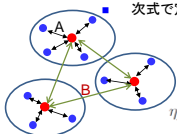
- k-means++法: 距離の離れた要素を初期クラスタの重心に設定することで、分類精度を向上
  - ランダムに一つの要素を選び、クラスタ重心に設定
  - 各要素 x に関して、その最近傍重心との距離 D(x) を計算
  - D(x)<sup>2</sup> に比例する確率に従い、新しいクラスタ重心としてランダムに一つ要素を選択
  - k 個のクラスタ重心が選択されるまで上記処理を反復
  - 以後は k-means法を用いてクラスタを生成



Copyright©2016 NTT corp. All Rights Reserved. 21

## 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



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

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

$x_i^{(j)}$ : 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©2016 NTT corp. All Rights Reserved. 22

## Basic Properties

ID	Category	Website		Object size (kbytes)	Object count	Total size (Mbytes)
		count	size			
C1	Business	59	40	14.70	55.14	0.810
C2	Computers	112	91	16.26	43.63	0.709
C3	News	39	27	13.55	72.45	0.982
C4	Reference	112	109	13.09	43.42	0.568
C5	Regional	80	73	17.77	50.59	0.899
C6	Science	95	86	14.04	52.86	0.742
C7	Society	79	83	15.01	66.86	1.003
C8	Health	86	52	14.27	54.30	0.775
C9	Home	85	47	15.66	55.39	0.867
C10	Shopping	69	68	15.67	70.77	1.109
C11	Adult	112	102	10.49	53.04	0.557
C12	Arts	55	60	15.43	68.18	1.052
C13	Games	87	58	15.28	54.12	0.827
C14	Kids & teens	106	64	13.23	54.59	0.722
C15	Recreation	86	52	13.55	57.30	0.776
C16	Sports	38	53	16.62	86.67	1.440

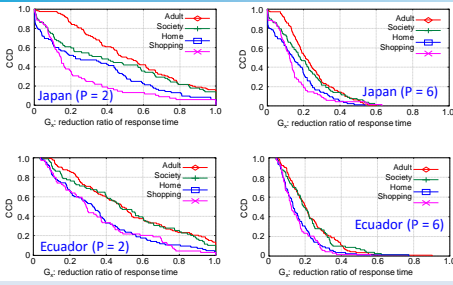
- Entertainment websites, e.g., Arts, Shopping, and Sport, tend to have more objects and larger total data size.
- Information websites, e.g., Business, Computers, Health, and Reference, tend to have fewer objects and smaller total data size.



Copyright©2016 NTT corp. All Rights Reserved. 23

各サイトの応答時間削減率のCCD

Universal: Adult, Society  
Localized: Home, Shopping

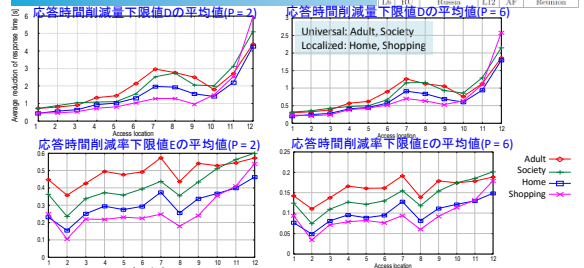


- キャッシュヒット率をH=1とした場合、各Webページの応答時間削減率G<sub>r</sub>のCCDを4つのカテゴリごとにプロット
- UniversalとLocalizedとで、応答時間削減率に明確な差異を確認



アクセス拠点による傾向の差異

ID	Area	Location	ID	Area	Location
L1	NA	Missouri	L7	OA	Australia
L2	NA	Wisconsin	L8	OA	New Zealand
L3	NA	California	L9	AS	Japan
L4	EU	Ireland	L10	SA	Ecuador
L5	EU	Germany	L11	SA	Argentina
L6	EU	Italy	L12	AF	Egypt



- 各地点・4ジャンルのDの平均値(上図)とEの平均値(下図)をプロット
- ジャンルごとのEで見たエッジ配信の効果の順位は、どの地点でもほとんど同一
- 品質が良好な北米も含めて、全地域で、UniversalはLocalizedよりもエッジ配信の効果が高く、効果にジャンルグループ間の差異が見られることを確認

