# A Light-weight Content Distribution Scheme for Cooperative Caching in Telco-CDNs

Takuma Nakajima, Masato Yoshimi, Celimuge Wu, Tsutomu Yoshinaga

The University of Electro-Communications

# Summary

- **<u>Proposal</u>**: A light-weight scheme to utilize cache servers to reduce growing internet traffic

- **<u>Major Contributions</u>**:
  - **<u><span style="color:red">Enhanced traffic reduction by distributing contents</span></u>** with a simple grouping scheme of cache servers and contents
  - **<u><span style="color:red">Following a rapid change in access patterns</span></u>** by utilizing a hybrid caching scheme of LFU and LRU algorithms

- **<u>Evaluation</u>**:
  - Case study using a backbone network in Japan and YouTube access patterns
  - Comparison of traffic reduction and computational overhead with a sub-optimal result calculated by Genetic Algorithm

# Outline

- **Introduction**
  - Rapid growth of video traffic
  - Efficient utilization of cache servers

- **Proposal**
  - Adjusting cache distribution by a simple scheme
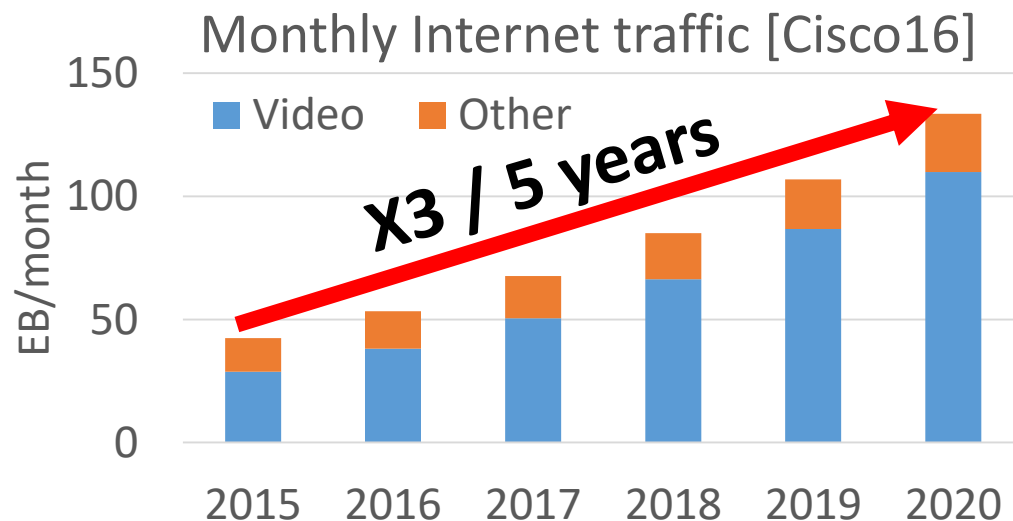  - Following rapid change in access pattern by LFU/LRU hybrid caching

- **Evaluation**
  - Traffic reduction and computational overhead compared with a sub-optimal result

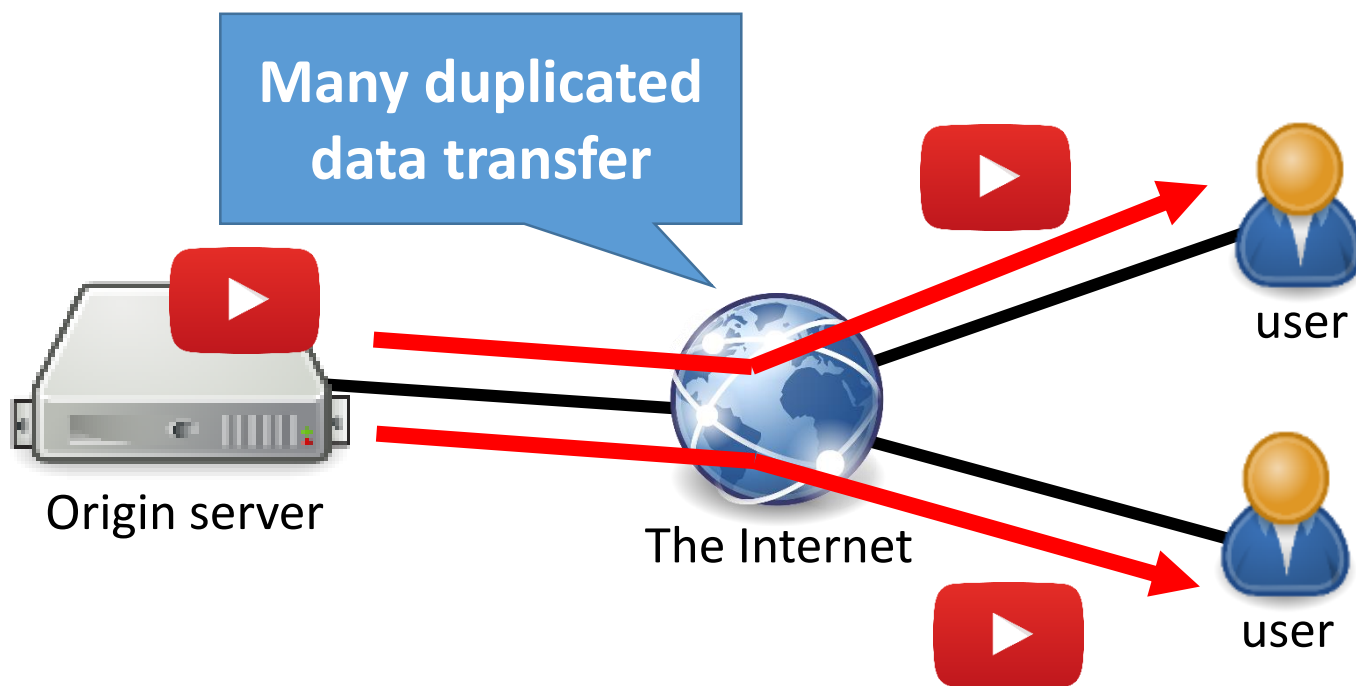- **Conclusion and Future work**

# Rapid growth of video traffic

- **<span style="color:red">Video-on-Demand</span>** (VoD) services will contribute more than **<span style="color:red">80% of internet traffic</span>** in 2020 [Cisco16]

- Such enormous traffic will cause many congested links and **degrade network performances**

- **<span style="color:red">Efficient utilization of cache servers</span>** is a key to reduce the internet traffic

Monthly Internet traffic [Cisco16]

■ Video   ■ Other

X3 / 5 years

EB/month

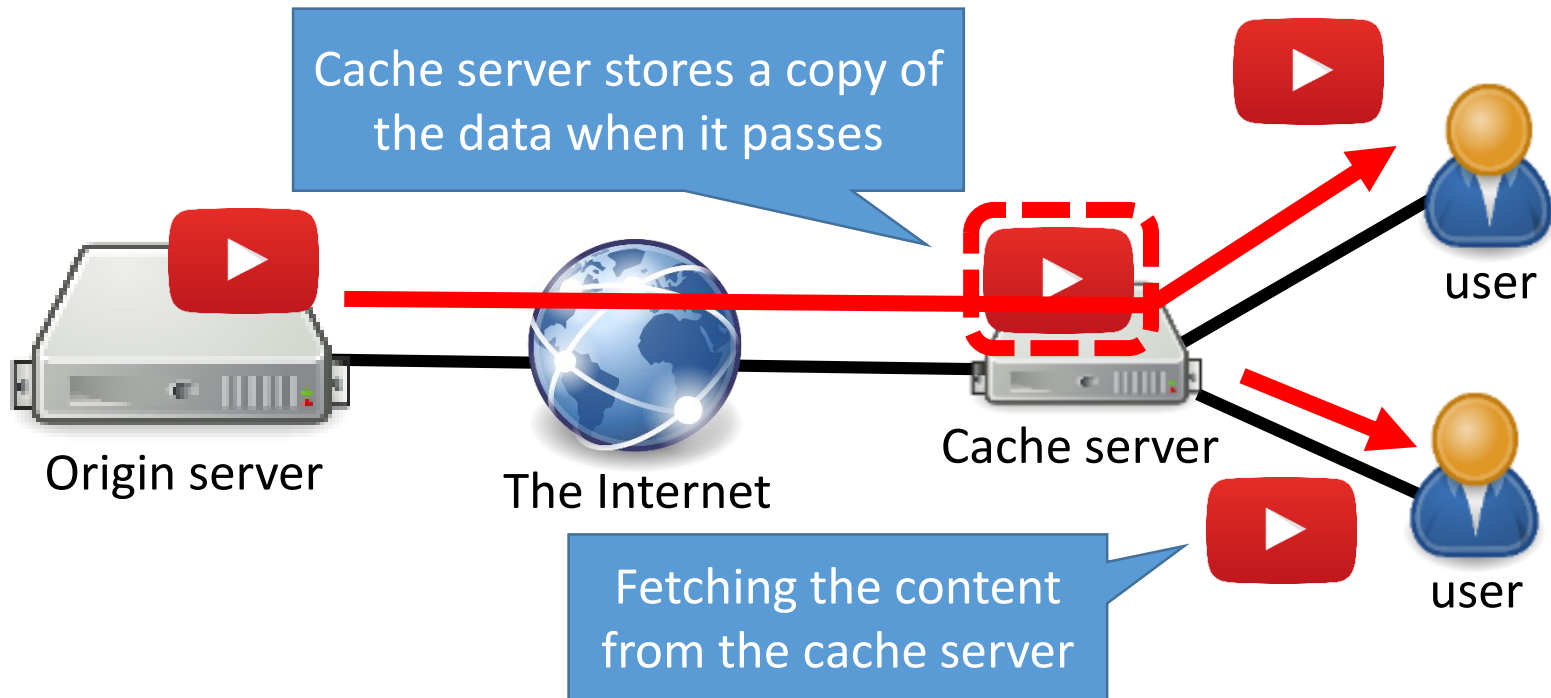150 — 100 — 50 — 0

2015  2016  2017  2018  2019  2020

# How internet traffic increases?

- The network **transfers the same data** to different users many times increasing the internet traffic
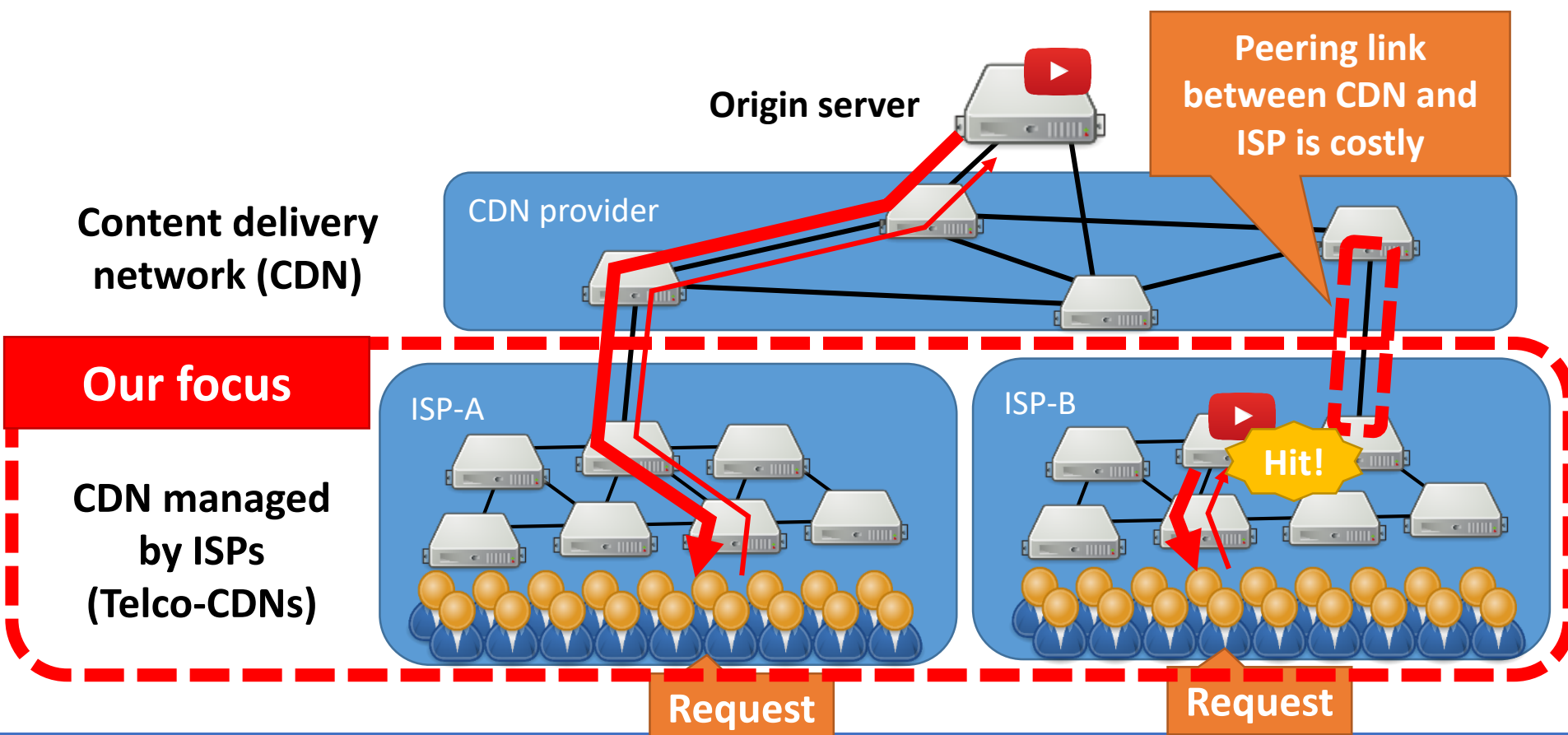
# Traffic reduction by cache servers

- The cache server **stores a copy of contents** when they pass the server

- The cache server **responds the copy to users** to reduce the traffic from the cache server to the origin
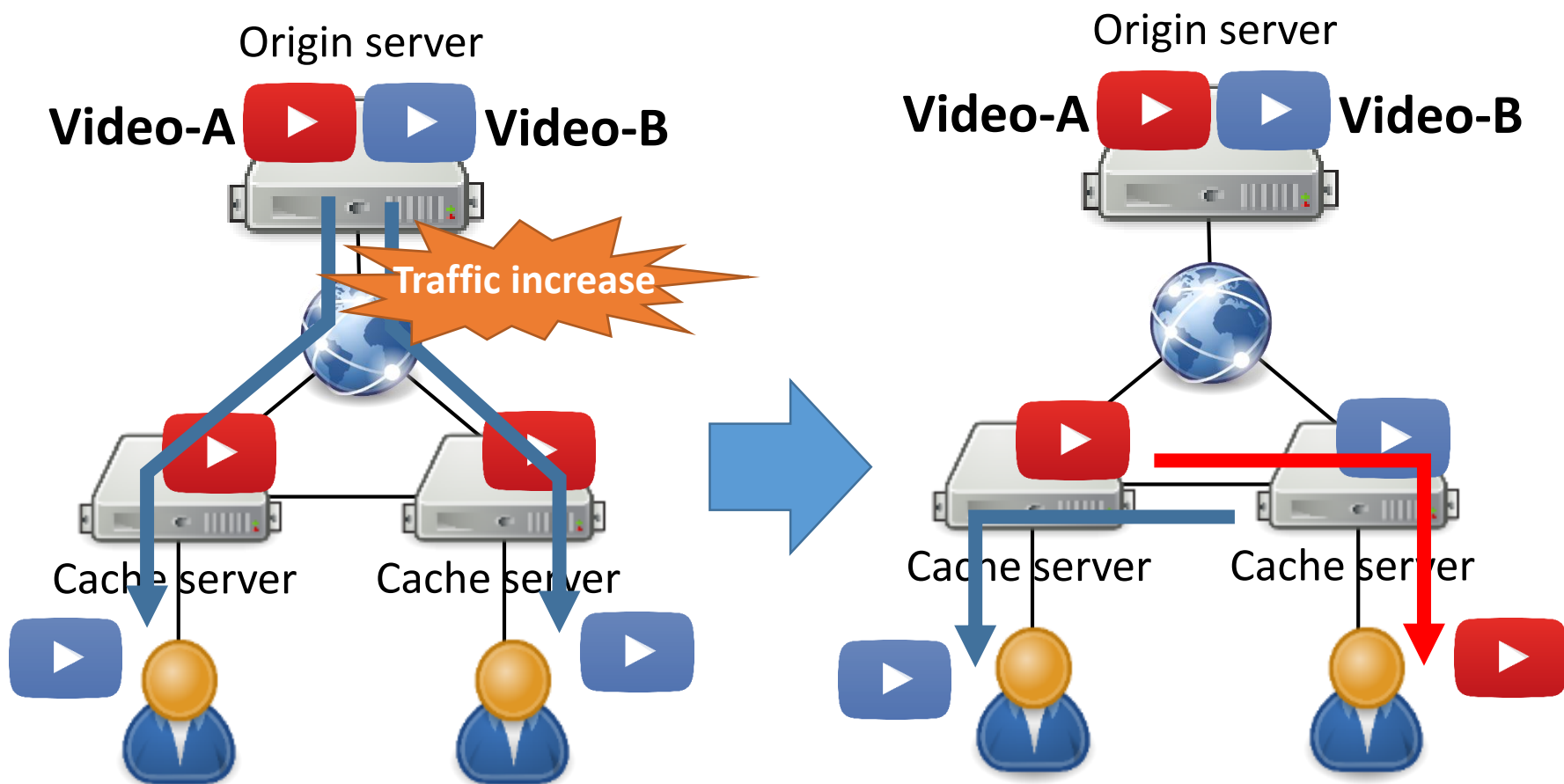
Cache server stores a copy of the data when it passes

Origin server

The Internet

Cache server

user

user

Fetching the content from the cache server

# Tiered cache networks

- It is better to **complete requests in an ISP network** to reduce traffic and communication costs
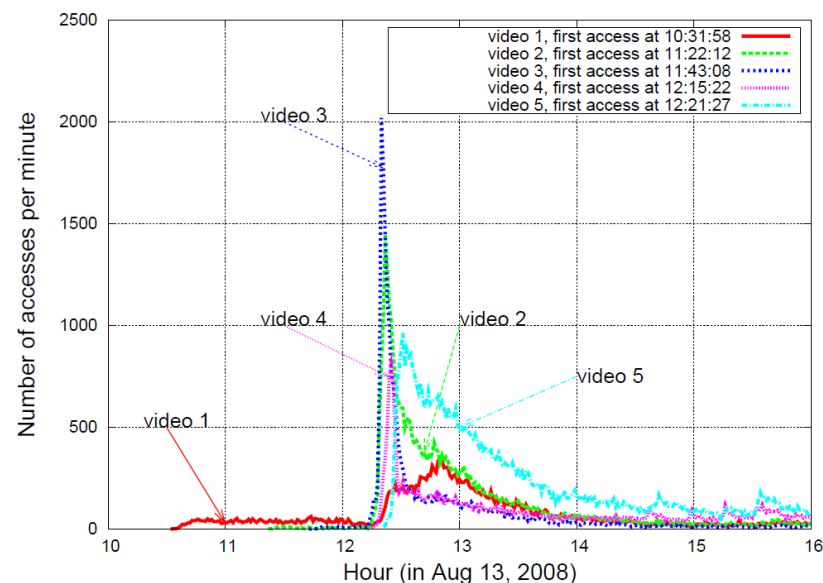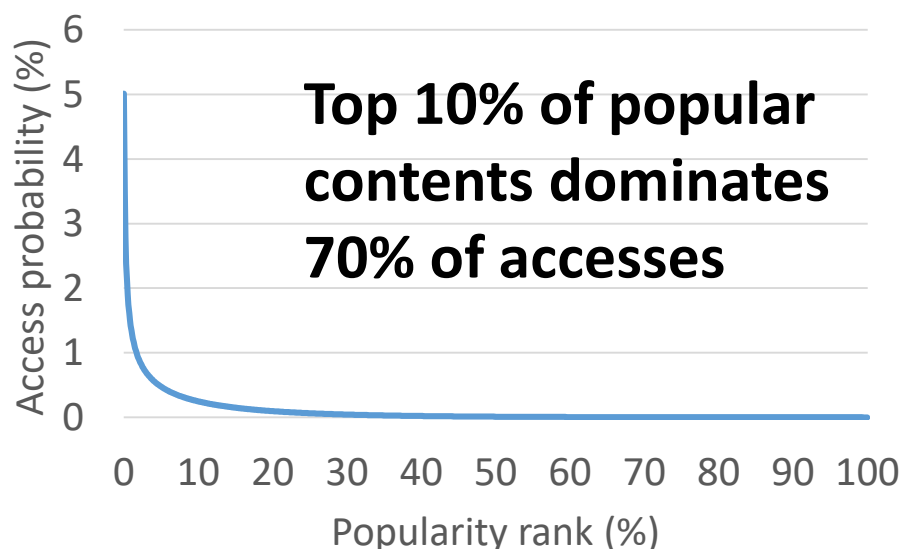
# Cooperative caching

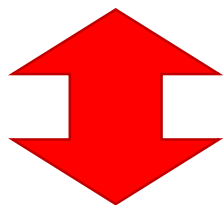- **<u>Increasing effective storage size</u>** by grouping several cache servers

# Characteristics of video accesses

- **<u>Skewed accesses</u>** [Cheng13]: Most accesses request limited popular contents

- **<u>Rapid change in contents' popularities</u>** [Yin09]: Access patterns often change widely due to news and viral communications in SNS

**Top 10% of popular contents dominates 70% of accesses**
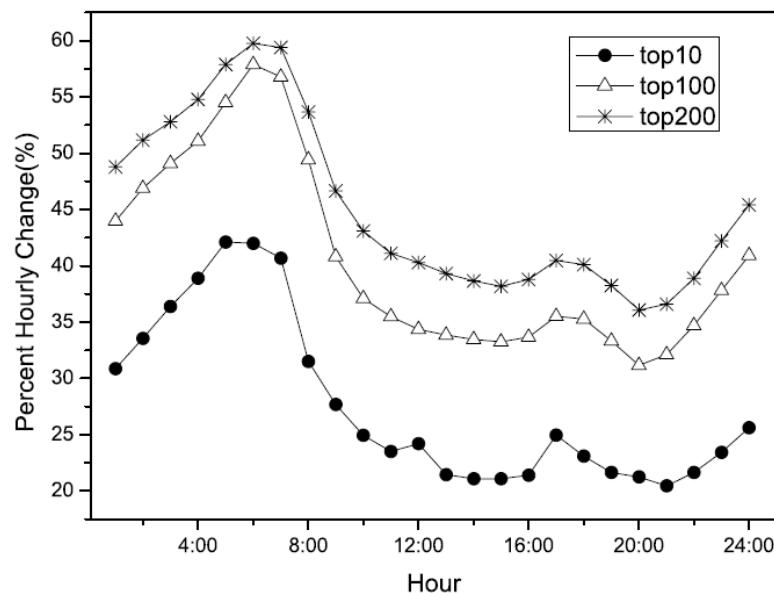
# Finding optimal cache allocation

- Calculating **sub-optimal allocations of contents** to minimize the traffic with Genetic Algorithm (GA) [Li13]
  - It takes around **10 hours' calculation**

- **Access patterns change 20-60% every hour** [Yu06]
  - Long calculation time causes mismatches in the allocation

**Hourly popularity change**

# LFU/FIFO hybrid caching [Zhou15]

- LFU/FIFO hybrid caching **improves cache hit rate** and **follow changes in access patterns**
  - LFU: Improving hit rate of each cache server
  - FIFO: Following change in access patterns

- It **does not support cooperative caching**

| LFU cache area improves hit rate by caching popular contents | FIFO cache area stores recently accessed contents |
|---|---|

LFU/FIFO Hybrid Cache

# Outline

- **Introduction**
  - Rapid growth of video traffic
  - Efficient utilization of cache servers

- **<u>Proposal</u>**
  - Adjusting cache distribution by a simple scheme
  - Following rapid change in access pattern by LFU/LRU hybrid caching

- **Evaluation**
  - Traffic reduction and computational overhead compared with a sub-optimal result

- **Conclusion and Future work**

# Efficient use of cache servers

- A key factor of an efficient cache management is a combination use of **content distribution** and **duplication of popular contents**



**Cache distribution:**
**increase effective storage size**

**Duplication of popular contents:**
**increase hit rates of servers**

# Content distribution by color tags

- **<u>Increasing cache capacity</u>** by explicitly storing contents among cache servers

- Grouping and associating cache servers and contents **with color tags with a specific color**
  - Each cache server stores contents if the color matches

**Increase cache capacity up to X4**



Cache server    Cache server    Cache server    Cache server

# Content duplication by color tags

- **Eliminating traffic** among cache servers
- **Duplicate popular contents** by **applying multiple colors** to them to increase hit rates

**Increase hit rates of cache servers**

Popular contents

Unpopular contents

Cache server  Cache server  Cache server  Cache server

# Example cache distribution

- Contents are basically distributed
- Several popular contents are duplicated

: Popular content

: Unpopular content

**Content library**

**Origin server**

# Preparing color tags

- A color tag is a set of bits, and each bit stands for a specific color

  🟥🟩🟦🟨 1111

  🟥🟩🟦⬜ 1110

  🟥🟩⬜⬜ 1100

  🟥⬜⬜⬜ 1000

- Popular contents have tags with many 1-bit to increase hit rates

| # of colors | R | G | B | Y |
|:-----------:|:-:|:-:|:-:|:-:|
| 4 | 1 | 1 | 1 | 1 |
| 3 | 1 | 1 | 1 | 0 |
| | 1 | 1 | 0 | 1 |
| | 1 | 0 | 1 | 1 |
| | 0 | 1 | 1 | 1 |
| 2 | 1 | 1 | 0 | 0 |
| | 1 | 0 | 1 | 0 |
| | 1 | 0 | 0 | 1 |
| | 0 | 1 | 1 | 0 |
| | 0 | 1 | 0 | 1 |
| | 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 | 0 |
| | 0 | 1 | 0 | 0 |
| | 0 | 0 | 1 | 0 |
| | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 |

# Coloration of cache networks

- Each cache server is preliminarily colorized with a specific color like the **four-color theorem**
    - For a case study, we colorized the network by preferring longer distances between the same colors



Origin Server

Origin Server

# Coloration of contents

- **Sorting contents** by their popularities and **set color tags in a cyclic fashion**

| # of colors | R | G | B | Y |
|---|---|---|---|---|
| 4 | 1 | 1 | 1 | 1 |
| 3 | 1 | 1 | 1 | 0 |
| | 1 | 1 | 0 | 1 |
| | 1 | 0 | 1 | 1 |
| | 0 | 1 | 1 | 1 |
| 2 | 1 | 1 | 0 | 0 |
| | 1 | 0 | 1 | 0 |
| | 1 | 0 | 0 | 1 |
| | 0 | 1 | 1 | 0 |
| | 0 | 1 | 0 | 1 |
| | 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 | 0 |
| | 0 | 1 | 0 | 0 |
| | 0 | 0 | 1 | 0 |
| | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 |

| Rank | Popularity class | Content name | Tag |
|---|---|---|---|
| 1 | High | Video01.mp4 | 1111 |
| 2 | High | Video02.mp4 | 1111 |
| ⋮ | | | |
| 11 | Mid-High | Video11.mp4 | **1110** |
| 12 | Mid-High | Video12.mp4 | 1101 |
| 13 | Mid-High | Video13.mp4 | 1011 |
| 14 | Mid-High | Video14.mp4 | 0111 |
| 15 | Mid-High | Video15.mp4 | **1110** |
| ⋮ | | | |
| 130 | Middle | Video130.mp4 | 1100 |
| 131 | Middle | Video131.mp4 | 1010 |
| 132 | Middle | Video132.mp4 | 1001 |
| 133 | Middle | Video133.mp4 | 0110 |

Set color tags from popular contents

# Following rapid access changes

- We adopt a **hybrid caching scheme** with colored LFU and no-color Modified LRU [Vleeschauwer11] areas
  - Modified LRU achieves better hit rate than LRU

- **Colored LFU area** stores contents with matching tags, while the **Modified LRU** area stores contents without matching tags

| Colored LFU cache area stores contents with matching color tags | LRU area stores any contents **regardless of tags' colors** |
|---|---|

Hybrid Cache with Red color tag

# Outline

- **Introduction**
  - Rapid growth of video traffic
  - Efficient utilization of cache servers

- **Proposal**
  - Adjusting cache distribution by a simple scheme
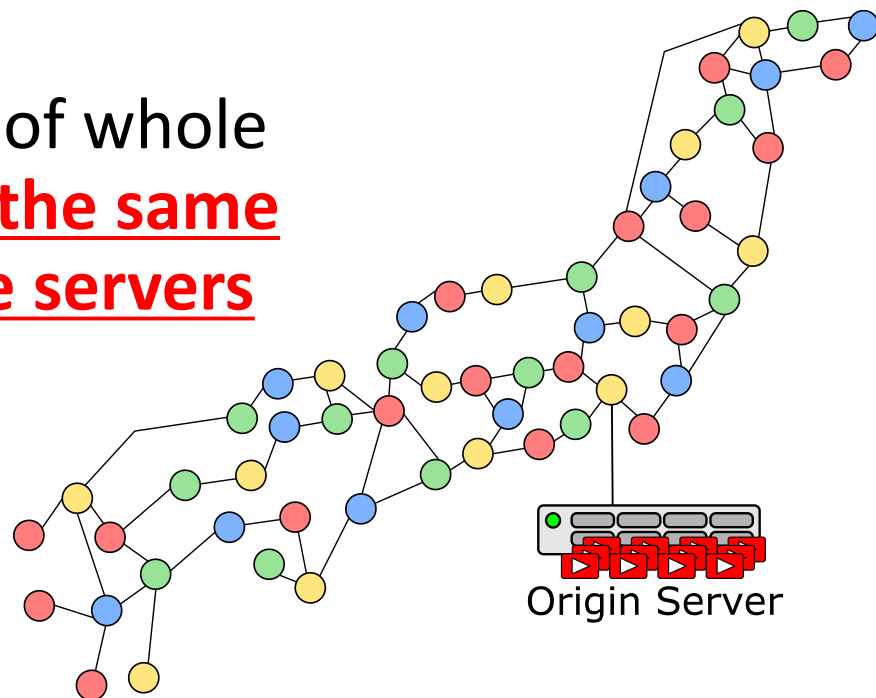  - Following rapid change in access pattern by LFU/LRU hybrid caching

- **Evaluation**
  - Traffic reduction and computational overhead compared with a sub-optimal result
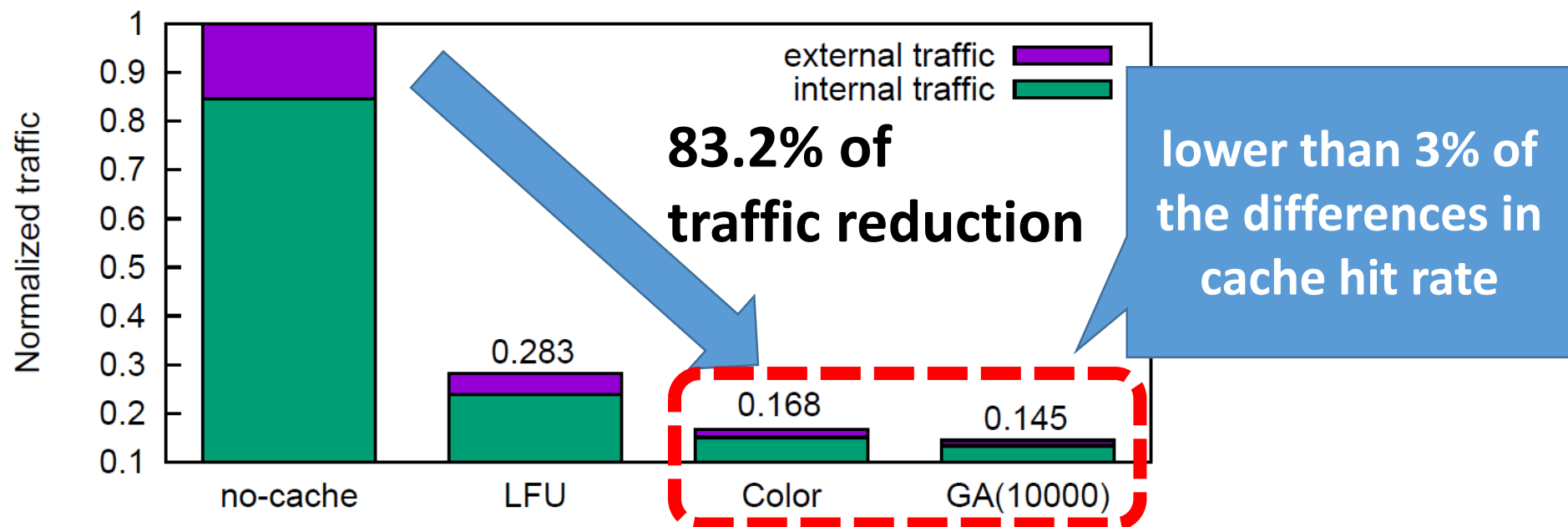
- **Conclusion and Future work**

# Evaluation environment

- **NTT-like topology** in Japan [Arteta07]

- Origin server is connected to a cache server in Tokyo

- Skewed access pattern with **YouTube video traffic** [Cheng13]

- Each cache can store 10% of whole contents, which is almost **the same capacity as Netflix's cache servers**

Origin Server

# Traffic reduction

- Proposed color-cache scheme could **achieve close to the sub-optimal result** calculated by GA

**83.2% of traffic reduction**

**lower than 3% of the differences in cache hit rate**

Normalized traffic chart:
- Legend: external traffic (purple), internal traffic (green)
- no-cache: 1
- LFU: 0.283
- Color: 0.168
- GA(10000): 0.145

# Computational overhead

- **<u>Colorization overhead is limited to a few seconds</u>** since it only have to sort and update tags in a cyclic fashion

Table: Computational time of GA

| Topology | Nodes | Generation | | | |
|----------|-------|------|------|------|-------|
| | | 1000 | 3000 | 8000 | 10000 |
| Ring | 8 | 5m33s | 16m01s | 42m05s | 52m33s |
| 2D-mesh | 25 | 34m44s | 103m11s | 274m40s | 343m17s |
| NTT | 55 | 42m08s | 127m34s | 350m38s | 440m25s |

**GA takes more than 7 hours until the conversion when using recent Core i7 CPU**

# Following dynamic accesses

- The **Colored hybrid caching** scheme could **maintain its hit rate** even when new contents are inserted

**Colored hybrid caching could limit the degradation to 2.3%**

**Colored LFU cache 90%** | **Modified LRU 10%**

Colored Hybrid Cache

Colored cache + Modified LRU
Colored cache only

Cache hit rate

**Inserted 5 popular contents (0.5% of content library) with no-color**

Request ID (x1000)

**Colored LFU cache 100%**

Colored Cache

**Single colored LFU cache drops the hit rate by 13.9%**

# Conclusion and Future work

- **<u>Proposal</u>**: A light-weight scheme to manage cache servers by focusing on **<u style="color:red">content distribution</u>** and **<u style="color:red">duplication</u>** with a simple colorization scheme

- **<u>Evaluation</u>**:
  - Colored caching scheme could **<u style="color:red">achieve close to the sub-optimal result</u>** with less than 3% of difference in hit rates
  - Computational **<u style="color:red">overhead is limited to a few seconds</u>**
  - Colored hybrid caching scheme could also **<u style="color:red">follow the rapid change in access patterns</u>** limiting the degradation to 2.3%

- **<u>Future work</u>**
  - More efficient ways to colorize cache servers and routing algorithms for further enhancing the traffic reduction