

# Adaptive TTL-Based Caching for Content Delivery

Soumya Basu<sup>1</sup>, Aditya Sundarrajan<sup>2</sup>, Javad Ghaderi<sup>3</sup>, Sanjay Shakkottai<sup>1</sup>, and Ramesh Sitaraman<sup>2,4</sup>

UT Austin<sup>1</sup>, UMass Amherst<sup>2</sup>, Columbia University<sup>3</sup>, and Akamai Technologies<sup>4</sup>

## Content Delivery Networks (CDN)

- CDNs deliver **millions of requests** from content provider to users.
- CDNs bring **content closer to end users** by caching content locally.
- CDNs handle **heterogeneous and ephemeral** contents, e.g. webpage, video.

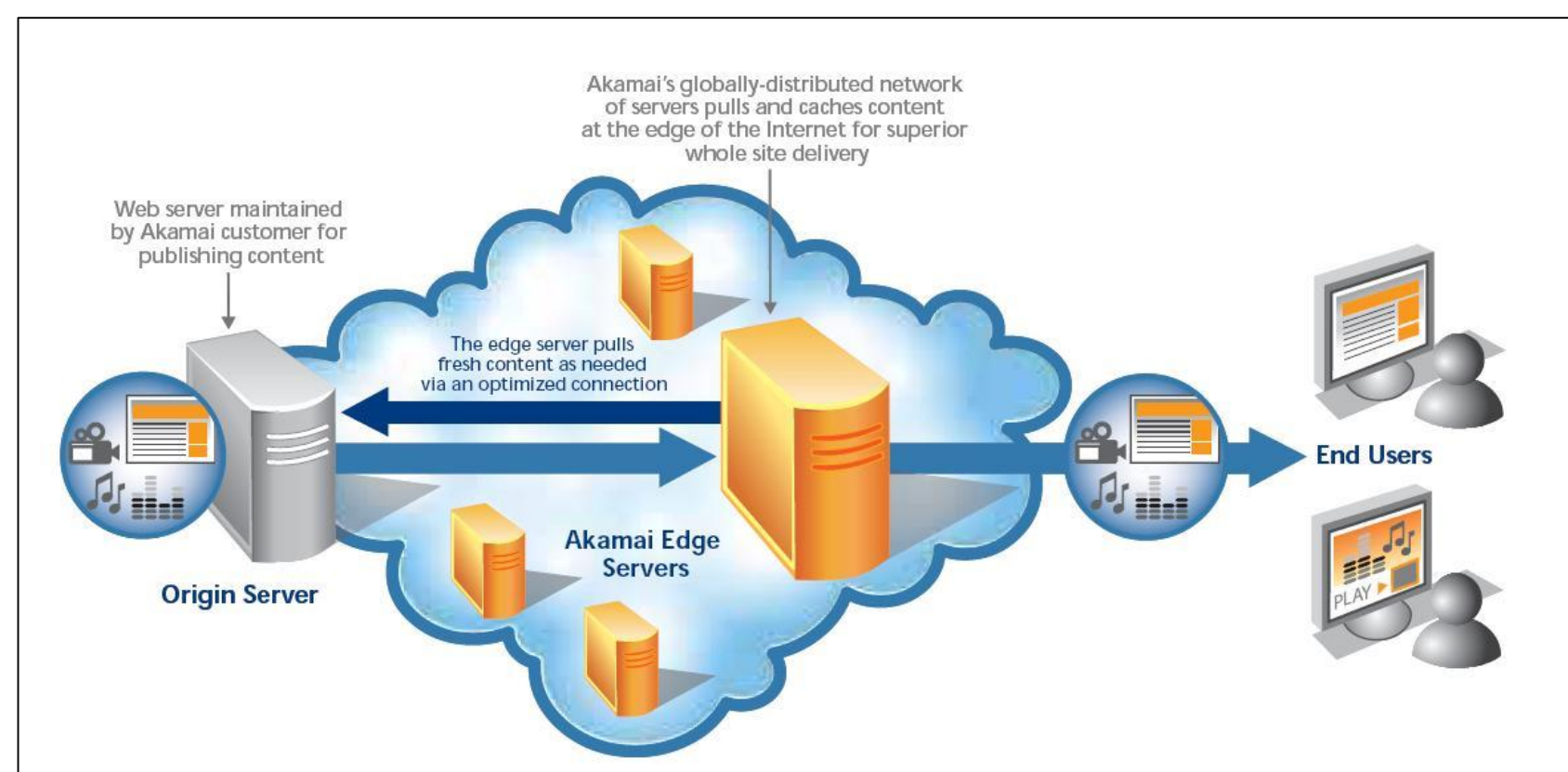


Fig 1: Content Delivery Networks

## CDN Objectives and Cache Design

- **Primary Objective:** Guarantying Quality of Service (QoS) to the end users.



- **Ambitious Objective:** Enabling Pricing Schemes for content providers.



\* Explained in Caching Terminology (to the right)

## Challenges in Cache Design

### Content Request Process

1. Millions of Objects spanning thousands of types
2. Correlated Arrivals with complex Inter-arrival distribution
3. Complex and non-independent content popularity
4. Non-stationary arrivals – (i) One-hot wonders, (ii) Flash Crowds

- **Cache Design:** Trace-based methods and Theoretical methods

### Theoretical methods

$$\text{Design Parameter} = f\left(\begin{matrix} \text{Hit rate, size rate,} \\ \text{Request Process param.} \end{matrix}\right)$$

### Trace-based Model

Exhaustive search for parameters  
Evaluation on traces via simulation

- |                                                    |                                                    |
|----------------------------------------------------|----------------------------------------------------|
| 1) Very few results with <b>non-stationarity</b>   | 1) <b>Expensive</b> in large scale, such as CDN    |
| 2) <b>Knowledge of Model</b> for closed form       | 2) Erroneous for <b>non-stationary</b> traces      |
| 3) Request parameters <b>high-dimensional</b>      | 3) Inaccurate for traces with <b>diverse types</b> |
| 4) Learning is <b>complex</b> and <b>expensive</b> | 4) <b>Long traces</b> necessary for accuracy       |

## Overcoming Challenges: Adaptive Caching

- **Model-oblivious** and **Target-driven, Online Adaptation** of the parameters
- **Time-to-live (TTL)** caches for adaptation with hit-rate guarantees
- **Circumventing non-convexity:** Achieving size-rate, not optimizing
- Higher degrees of freedom for size-rate control: **Two level TTL caches.**
- **Adaptive Filtering** of Non-stationary content to limit wastage of size

## Dynamic TTL Cache (d-TTL)

### Single Level TTL Cache

- One 'TTL',  $\theta_t$  for each type  $t$
- On Miss, cache with TTL,  $\theta_t$
- On Hit, reset the TTL to  $\theta_t$
- On timer **Expiry**, evict object

### Single Level TTL Adaptation

- Hit-rate target for each type  $t$ ,  $h_t < 1$
- Adaptive TTL,  $\theta_t(l)$  on  $l^{\text{th}}$  request
- On  $l^{\text{th}}$  request: Type  $t$  object,  $\alpha \in (0.5, 1)$ 
  - If Miss,  $\theta_t(l) = \theta_t(l-1) + \frac{1}{l^\alpha} h_t$
  - If Hit,  $\theta_t(l) = \theta_t(l-1) + \frac{1}{l^\alpha} (h_t - 1)$

## Caching Terminology

**Meta-data:** Id of an Object.

**Type:** E.g. data, audio, video

**Hit(i):** Object in cache (i),  $i=1,2$

**Miss:** Object not in (any) cache.

**Virtual Hit:** Object not in both

cache, but object-id in cache (2).

$$\text{Hit rate} := \frac{\# \text{ Cache hit}}{\# \text{ Requests}}$$

$$\text{Size rate} := \frac{\text{Avg. cache size (Gb)}}{\text{Arrival rate (Gbps)}}$$

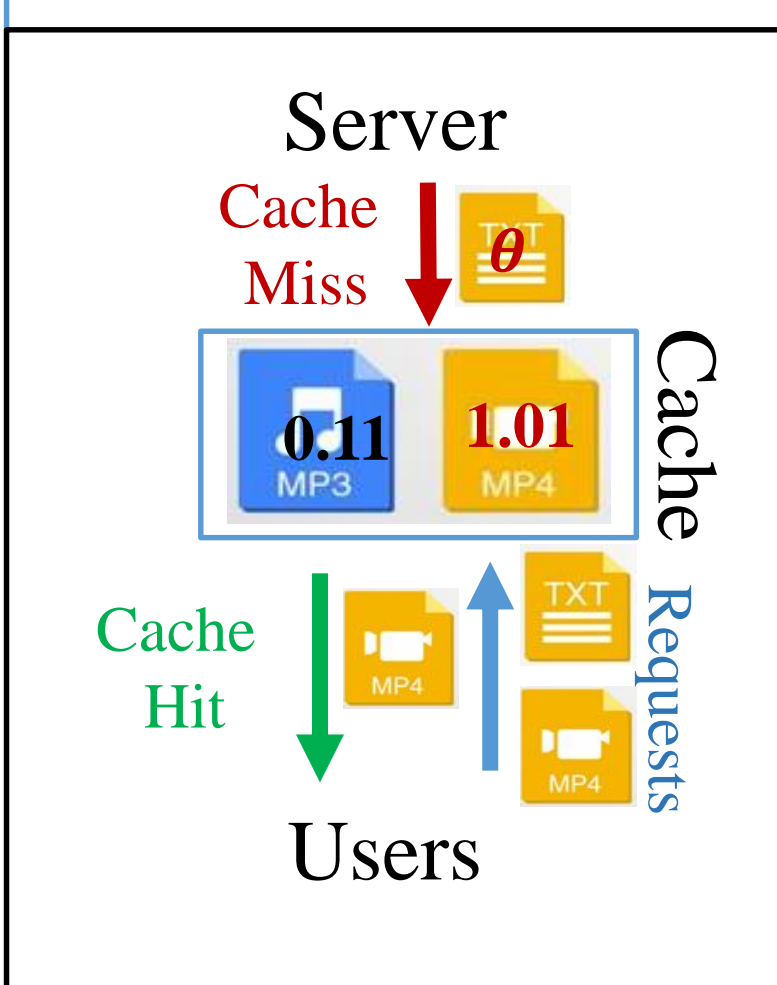


Fig 2: d-TTL Cache

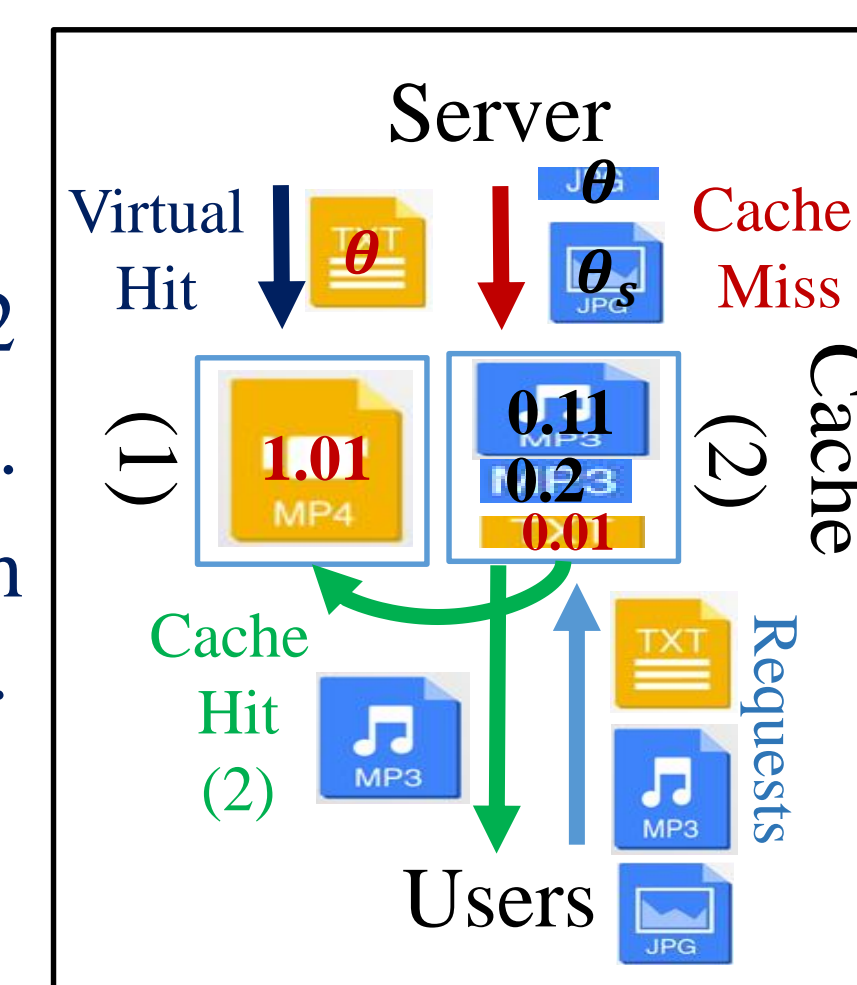


Fig 3: f-TTL Cache

## Filtering TTL Cache (f-TTL)

### Two Level TTL Cache

- Two level of caches: Cache (1) and (2)
- 'TTL' pair,  $(\theta_t, \theta_t^s)$  for each type  $t$
- 'TTL' pair satisfies  $\theta_t \geq \theta_t^s$  for all  $t$
- On Miss,
  - i) Cache object in (2) with timer  $\theta_t^s$
  - ii) Cache meta-data in (2) with timer  $\theta_t$
- On Hit in cache (1), reset the timer to  $\theta_t$
- On Hit in cache (2),
  - i) Cache object in (1) with timer  $\theta_t$
  - ii) Evict object and meta-data from (2)
- On Virtual hit,
  - i) Cache object in (1) with timer  $\theta_t$
  - ii) Evict meta-data from (2)
- On timer **Expiry**, evict object/meta-data

### Two Level TTL Adaptation

- Hit-rate target for each type  $t$ ,  $h_t < 1$
- Size-rate target for each type  $t$ ,  $s_t$
- Adaptive TTL,  $(\theta_t(l), \theta_t^s(l))$  on  $l^{\text{th}}$  request.
- On  $l^{\text{th}}$  request: Type  $t$  object,  $\alpha \in (0.5, 1)$ 
  - If Miss or Virtual hit,
 
$$\theta_t(l) = \theta_t(l-1) + \frac{1}{l^\alpha} h_t$$
  - If Hit, denote the timer value as  $\psi > 0$ 

$$\theta_t^s(l) = \theta_t^s(l-1) + \frac{1}{l} (s_t - \theta_t^s(l-1))$$

$$\theta_t(l) = \theta_t(l-1) + \frac{1}{l^\alpha} (h_t - 1)$$

$$\theta_t^s(l) = \theta_t^s(l-1) + \frac{1}{l} (s_t + \psi - \theta_t^s(l-1))$$
  - $\theta_t^s(l) = \min\{\theta_t(l), \theta_t^s(l)\}$

## Performance on Akamai Traces

- Akamai trace: Duration, **9 days**, #Requests, **504m**, #Objects, **25m**.
- **Targeted hit rate** for d-TTL and f-TTL: 40%, 50%, 60%, 70% and 80%.
- **Targeted size rate** for f-TTL: 50% of size rate of d-TTL.
- **ACCURACY:** Error in Achieved hit rate: **1.3%**, Achieved size rate: **2%**.

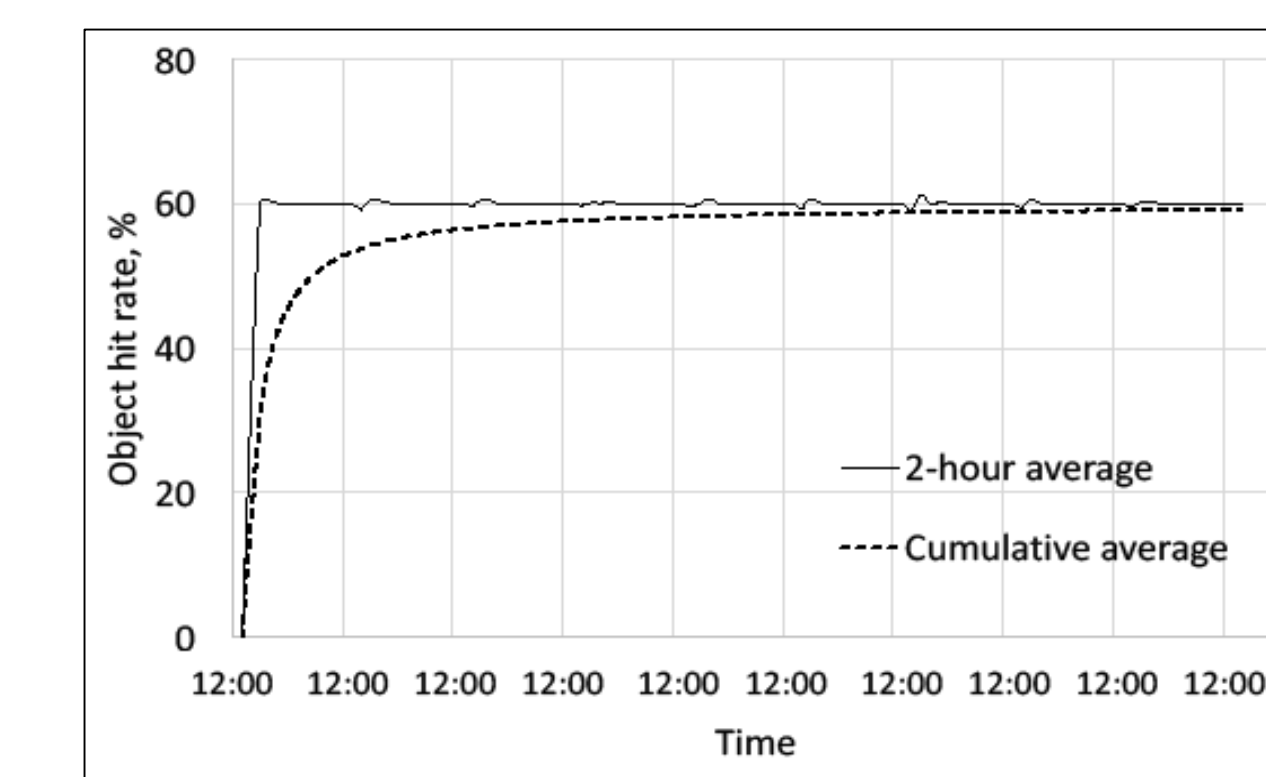


Fig 4: d-TTL Convergence Plot

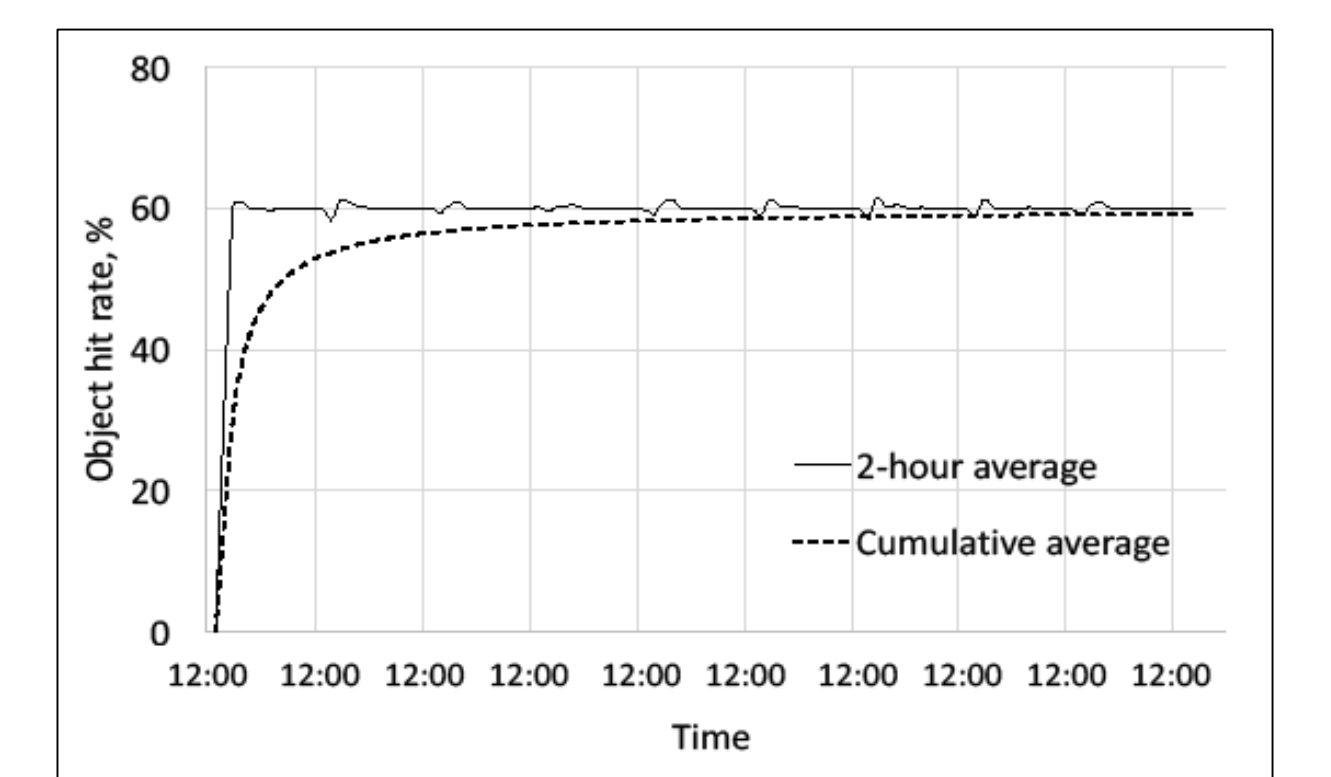


Fig 5: f-TTL Convergence Plot

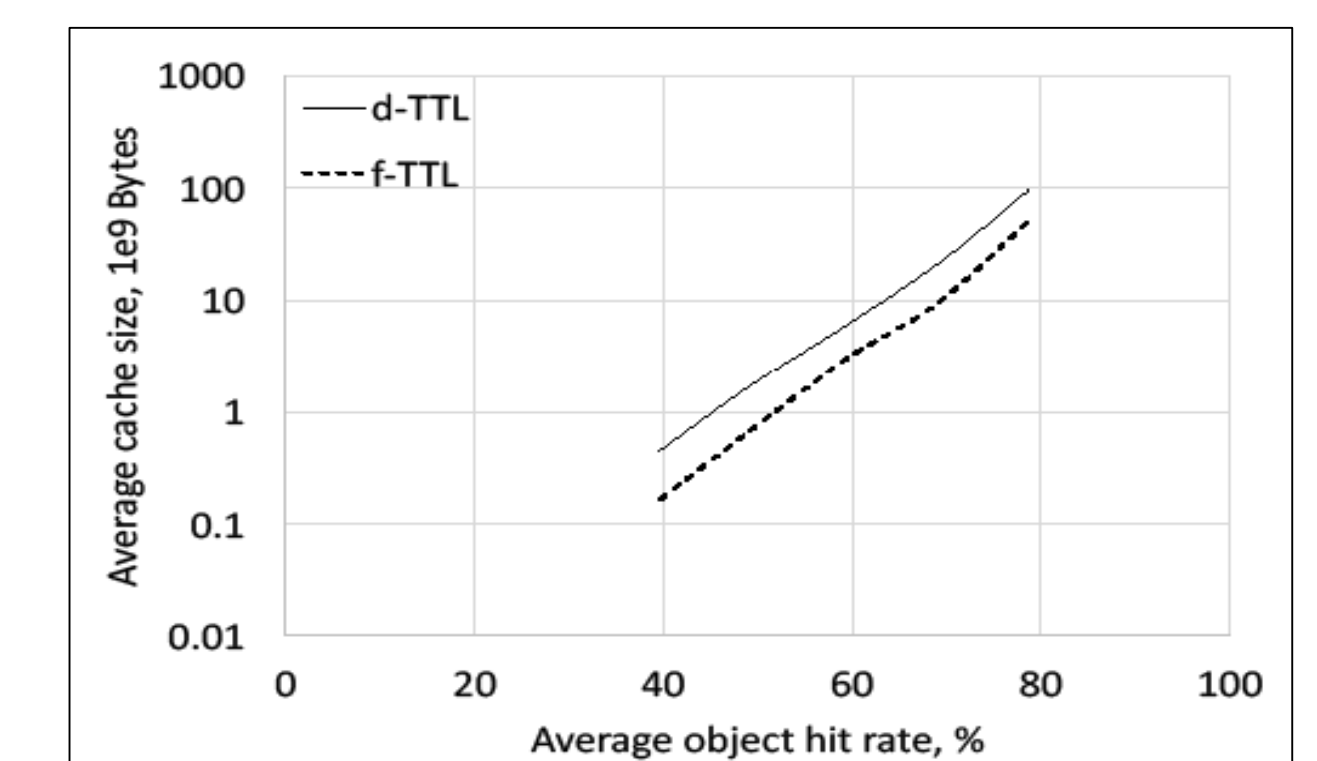


Fig 6: Hit rate vs Average Cache Size curve

## Adaptation Insights

- The **larger** the value of  $\theta_t$  the **higher** the hit-rate.
- On **cache hit decrease**  $\theta_t$  value and on **cache miss increase**.
- The second 'TTL',  $\theta_t^s$ , enables **filtering** of **rare** objects.
- Lower  $\theta_t^s$  + hit-rate target  $\Rightarrow$  **Smaller** cache (2) but **larger** cache (1).
- Thumb Rule: **Filtering reduces total size** under high non-stationarity.

## Performance Guarantee

### System Model

- Finite no. of types and finite no. of 'recurrent' objects of each type
- Arrival of 'rare' objects at a **non zero rate**, e.g. flash-crowd, one-hot objects
- Content request modeled as **Markovian Arrival Process**
- Inter arrival time with any **absolutely contd. pdf** with **connected support**

### Performance

- d-TTL\* and f-TTL\* **attains the target hit-rate** asymptotically almost surely
- f-TTL\* **attains size-rate  $\leq$  the target** or **collapses** to  $\theta_t^s = 0$  a.a.s.

### Proof Techniques

- **Stochastic approximation** technique used for TTL adaptation
- **Projected ODE** based proof technique to show convergence of d-TTL
- **Two timescale Actor-Critic** framework for convergence of f-TTL

\* A modified version of the current d-TTL and f-TTL algorithm