# PA-Cache: Evolving Learning-Based Popularity-Aware Content Caching in Edge Networks

Qilin Fan, *Member, IEEE,* Xiuhua Li, *Member, IEEE,* Jian Li, *Member, IEEE,*
Qiang He, *Senior Member, IEEE,* Kai Wang, *Member, IEEE,* and Junhao Wen

*Abstract*—As ubiquitous and personalized services are growing boomingly, an increasingly large amount of traffic is generated over the network by massive mobile devices. As a result, content caching is gradually extending to network edges to provide low-latency services, improve quality of service, and reduce redundant data traffic. Compared to the conventional content delivery networks, caches in edge networks with smaller sizes usually have to accommodate more bursty requests. In this paper, we propose an evolving learning-based content caching policy, named PA-Cache in edge networks. It adaptively learns time-varying content popularity and determines which contents should be replaced when the cache is full. Unlike conventional deep neural networks (DNNs), which learn a fine-tuned but possibly outdated or biased prediction model using the entire training dataset with high computational complexity, PA-Cache weighs a large set of content features and trains the multi-layer recurrent neural network from shallow to deeper when more requests arrive over time. We extensively evaluate the performance of our proposed PA-Cache on real-world traces from a large online video-on-demand service provider. The results show that PA-Cache outperforms existing popular caching algorithms and approximates the optimal algorithm with only a 3.8% performance gap when the cache percentage is 1.0%. PA-Cache also significantly reduces the computational cost compared to conventional DNN-based approaches.

*Index Terms*—Edge Caching, Popularity Prediction, Deep Learning, Quality of Service.

Q. Fan, X. Li and J. Wen are with Key Laboratory of Dependable Service Computing in Cyber Physical Society of Ministry of Education, Chongqing University, Chongqing 401331, China, with State Key Laboratory of Power Transmission Equipment and System Security and New Technology, Chongqing University, Chongqing 401331, China, and with the School of Big Data & Software Engineering, Chongqing University, Chongqing 401331, China (e-mail: fanqilin@cqu.edu.cn; lixiuhua1988@gmail.com; jhwen@cqu.edu.cn).

J. Li is with the Department of Electrical and Computer Engineering, Binghamton University, State University of New York, Binghamton, NY 13902, USA (e-mail: lij@binghamton.edu).

Q. He is with School of Software and Electrical Engineering, Swinburne University of Technology, Melbourne, 3122, Australia (e-mail: qhe@swin.edu.au).

K. Wang is with School of Computer Science and Technology, Harbin Institute of Technology, Weihai 264209, China, and with Research Institute of Cyberspace Security, Harbin Institute of Technology, Weihai 264209, China (e-mail: dr.wangkai@hit.edu.cn).

## I. INTRODUCTION

Today's Internet has seen an explosion of mobile data traffic in data-consuming application services such as mobile video services requested from a wide variety of mobile devices [1]. According to the Cisco VNI report [2], the overall average mobile traffic will increase at a compound annual growth rate of 46% between 2017 and 2022, and it is estimated that video traffic will account for about 80% of the global mobile data traffic. To enable ubiquitous mobile video services, cloud computing has been widely acknowledged and deployed as it can provide adequate computing and storage resources. However, with the sky-rocketing network traffic load and more stringent requirements of users, cloud-based mobile video services are facing new challenges such as large transmission latency and limited bandwidth resources.

To tackle the above challenges, edge computing has emerged as a new and evolving paradigm to accommodate future mobile video services, which provides computing and caching services close to end-users at the network edges [3], [4]. In particular, the deployment of edge caching can minimize the redundant data traffic, thereby leading to a significant reduction in service latency and elimination in bandwidth wastage.

Edge caching strategy attempts to learn the pattern of content requests in some fashions [5], ensuring the availability of contents as high as possible in the cache nodes (e.g., base stations (BSs) and edge routers). Generally, the requested content is searched for in a cache node. If it is unavailable, a miss occurs, and the requested content is fetched from an upstream server (typically with higher latency and more expensive transmission cost). The content is then stored in the cache node and finally transmitted to the user. Besides, compared to the total size of contents, the capacity of cache nodes is usually limited and much smaller. If the new content is required to be cached when the cache is full, several cached contents may have to be evicted. Therefore, caching algorithms can also be described by the employed eviction strategy. When caching contents at network edges, user requests for the same content will be locally served. This can effectively reduce the redundant data traffic and greatly improve the quality of service [6].

Compared to the conventional content delivery networks (CDNs), edge caching has its unique characteristics [7], [8]: (i) *Limited resources.* The cloud usually has a large number of diverse resources. However, the edge cache with limited

computing and storage resources enables only a small fraction of contents to be cached and low-complexity tasks to be executed; and (ii) *Bursty requests*. The requests from edge networks usually vary a lot over time.

Nowadays, most caching system still utilize recency-based [9], frequency-based [10], size-based [11], or combinations of them. The limitation is that they might work well for some access patterns but poorly for others. Recently, learning-based caching algorithms have been proposed either to determine which contents should be evicted when the cache is full or decide whether or not to admit a content upon a request by learning content popularity. However, the content requests of edge networks are time-varying and bursty. On the one hand, it is difficult for shallow machine learning models to capture complex patterns. On the other hand, by using the entire training dataset, conventional deep neural networks (DNNs) would learn a fine-tuned but possibly outdated or biased prediction model with high computation complexity, making is difficult to support the application at edge caches with limited computing capability.

In this paper, we propose a novel popularity-aware content caching policy, namely PA-Cache, in edge networks. Different from previous approaches, PA-Cache weighs a large set of content features to learn the time-varying content popularity adaptively. Furthermore, to overcome the high computational cost of conventional DNN, PA-Cache takes advantage of a shallow network with fast convergence at the beginning, and then the powerful representation of DNN when more requests arrive over time. In this way, PA-Cache achieves high scalability and high computation efficiency. The contributions of this paper can be summarized as follows:

- We investigate the issue of popularity-aware content caching and design modules and operations of popularity-aware cache nodes in edge networks. We apply a learning-based approach to tackle this problem.
- We amend the multi-layer recurrent neural network (RNN) architecture by attaching every hidden layer representation to an output regression to predict the temporal content popularity. We utilize a hedge strategy, which enables adaptive training of DNN in an evolving setting. The popularity-aware cache node makes appropriate cache replacement decisions to maximize the long-term cache hit rate based on the estimated popularity.
- We conduct extensive experiments on a real-world dataset derived from iQiYi, which is the largest online video-on-demand (VoD) service provider in China. Trace-driven evaluation results demonstrate the effectiveness and superiority of PA-Cache over several candidate algorithms.

The rest of this paper is organized as follows. The related work is briefly introduced in Section II. Section III provides a system overview. Section IV gives a formal description of the cache replacement problem. Section V presents the design of PA-Cache algorithm in detail. In Section VI, the performance of our proposed algorithm by trace-driven experiments is evaluated. Section VII concludes this paper.

## II. RELATED WORK

Content caching algorithms have been studied for many decades. Existing work can be divided into the following two main branches.

### A. Rule-Based Algorithms

The first branch is named rule-based algorithms. These cache eviction algorithms rely on one of the most widely used features (i.e., recency, frequency and size) or their combinations. LRU-K [12] is a combination of least recently user (LRU) [9] and least frequently used (LFU) [10], in which the cache remembers the time of last K occurrences instead of the last occurrence for each content. ARC [13] adaptively divided the cache space into two segments and maintains contents that have been referenced exactly once and at least twice, respectively. S4LRU [14] partitioned the cache into four lists, and each list was an LRU queue. Vietri *et al.* [15] proposed a LeCaR algorithm which managed two histories (i.e., recency and frequency) of metadata for each cache entry, and their weights were adaptively updated by regret minimization. Bahn *et al.* [16] evaluated the content based on its past accesses to estimate the likelihood of re-access, defined a least-unified value metric and normalizes it by the cost per unit size. Berger *et al.* [17] formulated the caching as a min-cost flow problem when considering variable content sizes.

To evaluate the performance of classic cache eviction algorithms, Martina *et al.* [18] proposed a unified and flexible approach by extending and generalizing a fairly decoupling principle (the so-called Che's approximation [19]). These algorithms follow heuristic rules and are easy to implement. However, most of the analysis was derived under an independent reference model, assuming that content popularity follows a fixed Zipf law [20], i.e., $p_i \propto i^{-\alpha}, \alpha > 0$, where $p_i$ refers to the request probability of $i$-th most popular content. The performance of a cache eviction algorithm under synthetic data traces is found to be quite different from that under real data traces [21]. Therefore, these algorithms might hardly accommodate the dynamic content access pattern in edge networks.

### B. Machine Learning-Based Algorithms

The second branch relies on machine learning algorithms to optimize the content caching strategy. This branch is further subdivided into two categories.

The first category investigates the use of "model-free" reinforcement learning (RL) in which the system starts without any prior assumption about the traffic pattern to converge towards the optimal caching strategy [22]–[28]. Such system learns to make decisions from experience through interactions with the environment, in which good actions are enhanced via a reward function. Kirilin *et al.* [23] trained a feedforward neural network (FNN) using a Monte Carlo method, which computed the admission probability to maximize the cache hit rate. Sadeghi *et al.* [24] developed a two-timescale deep Q-network approach in a hierarchical cache network to learn an online caching policy. Wang *et al.* [25] formulated the

content caching procedure as a Markov decision process and utilized double deep Q-network which stabilized the learning in discrete huge spaces for training the caching process. Fan *et al.* [26] leveraged the benefits of the RNN as well as the Deep Q-Network to maximize the cache efficiency by jointly learning request features, caching space dynamics, and making decisions. Guan *et al.* [27] designed a content-aware cache admission strategy. It consists of a tree-structured learning model for mining the critical features and a UCB-tree algorithm for making cache admission decisions dynamically. Yan *et al.* [28] leveraged RL and Belady to learn content popularity for both content admission and content eviction. However, RL-based algorithms require a tremendous number of learning samples and suffer from large-delayed rewards (cache hits). This can result in slow reaction times in highly dynamic environments [29]. Furthermore, as RL-based algorithms are quite sensitive to hyperparameters and random seeds, it is difficult to configure and maintain these algorithms [30].

The second category utilizes supervised learning, which learns key attribute features in requests and predicts the popularity or arrival times of contents to facilitate efficient caching [31]–[34]. Fedchenko *et al.* [31] employed FNN to predict content popularity and made caching decisions accordingly. Pang *et al.* [32] utilized deep long short-term memory (LSTM) to calculate the probability for each content to arrive. Song *et al.* [33] explored gradient boosting machine model to predict log (time-to-next-request). Chen *et al.* [34] designed a DNN-based popularity evolving model to estimate the near future popularity of IoT data. However, on the one hand, shallow machine learning models would be difficult to learn complex patterns. On the other hand, by using the entire training dataset, conventional DNNs would learn a fine-tuned but possibly outdated or biased prediction model with high computation complexity. Only a few papers [35], [36] proposed to use online learning-based algorithms as an adaptation method to predict content popularity by clustering the contents according to their similarities of access patterns. However, the selection of hand-crafted features has a significant impact on its performance gain. Moreover, the prediction model, which is essentially based on the sample average approximation method, is still shallow.

## III. SYSTEM OVERVIEW

In this section, we firstly introduce the cache node structure in the considered edge network. Then, we describe the main operations of the edge node.

### A. Cache Node Structure

The edge caching service is implemented and provided collectively by multiple cache nodes. Fig. 1 illustrates the modules of a considered popularity-aware cache node in an edge network. It provides both storage and computation capacity for caching service. As high-performance CPUs and GPUs are enabled (e.g., NVIDIA Jetson TX21[1]), deep learning tasks can be processed in the cache node. Typically, a

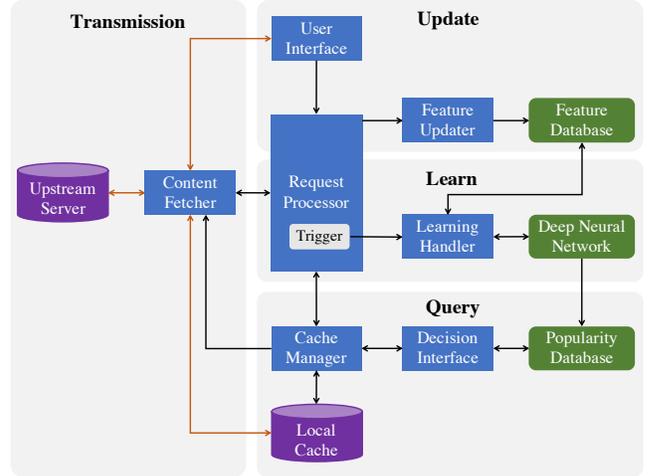[1]https://developer.nvidia.com/embedded-computing



Fig. 1. Modules and operations of a single popularity-aware cache node in an edge network.

cache node contains basic models: *Request Processor*, *Cache Manager*, *User Interface*, *Content Fetcher*, *Local Cache*. A popularity-aware cache node also integrates with *Feature Updater*, *Learning Handler*, *Decision Interface*, two databases (i.e., *Feature Database*, *Popularity Database*), as well as a *Deep Neural Network* to enable popularity awareness.

1) The *Feature Updater* module updates features from raw data of content requests, including contextual features (e.g., the number of video requests) and semantic features (e.g., video type) of content, and storing them in the *Feature Database* with a unified format.

2) The *Learning Handler* module trains a *Deep Neural Network* to predict content popularity in the next time step from the data in *Feature Database*.

3) The *Decision Interface* module determines which cached content should be evicted when a cache miss is declared.

### B. Operations

The main operations of a popularity-aware cache node consist of three procedures:

1) **Update.** The *Request Processor* receives a request and initiates an update operation based on the metadata of the received request. Then, the *Feature Updater* keeps track of the raw data, extracts the latest features, and writes them to the *Feature Database*.

2) **Learn.** The *Request Processor* periodically triggers a *Learning Handler* to learn a prediction model every $\phi$ hours. Parameter $\phi$ requires a proper selection for balancing the high accuracy against the low computational cost for the prediction model. The *Learning Handler* extracts the raw features and ground-truth popularity from the *Feature Database*. Then, it takes the normalized features together with the ground-truth popularity as the input to train a *Deep Neural Network* in an evolving manner. The predicted content popularity calculated by *Deep Neural Network* is recorded in the *Popularity Database*.

3) **Query.** The *Cache Manager* examines if the requested content is in the *Local Cache*. When it is locally available, the *Content Fetcher* fetches the content from the *Local Cache* and

serves the user. Otherwise, the *Cache Manager* sends a query to the *Decision Interface* that determines which cached content should be evicted according to *Popularity Database* and get its response. The *Cache Manager* evicts the least popular content and notifies the *Content Fetcher* to fetch the requested content from the *Upstream Server*, store it in the *Local Cache* and serve the user.

## IV. PROBLEM STATEMENT

In this paper, let $\mathcal{C} = \{1, 2, \ldots, C\}$ denote a given set of contents that are distributed to multiple end-users associated with the cache node. Since many replication strategies in the VoD system fragment the contents into equally sized chunks [37], we assume that all contents are unit-sized. It reduces the complexities and inefficiencies of continually allocating and de-allocating storage space to contents with non-uniform sizes. However, our work could also be extended to the arbitrary sizes considering popularity per unit size at the eviction phase like previous works [16], [38]. Let the cache capacity be $s$, that is, the cache node can accommodate up to $s$ contents. The sequence of requests for content is denoted as $\mathcal{K} = \{1, 2, \ldots, K\}$. Each request $k \in \mathcal{K}$ in this sequence is composed of three elements: the requested content $c_k$, the timestamp $t_k$ and the $d$-dimensional feature vector of the requested content $\boldsymbol{f}_k$.

For each request $k$, we have to examine whether it can be served by the cache. For this purpose, let $\boldsymbol{Z}_k = [Z_k^1, Z_k^2, \ldots, Z_k^C]$ be the indication vector of cache status at time $t_k$, where $Z_k^c \in \{0, 1\}$ is an indicator of whether content $c$ is in the cache or not (i.e., $Z_k^c = 1$ indicates that $c$ is available in the local cache and current request can be served, and 0 otherwise.)

When content $c_k$ is hit in the cache, the cache status vector stays the same: $\boldsymbol{Z}_{k+1} = \boldsymbol{Z}_k$. Otherwise, when content $c_k$ suffers from a miss, the cache node will retrieve it from a specific upstream server according to the traffic assignment criteria [39]. In this situation, the cache node will remove an old content $c^{evict}$ to make room for the new content $c_k$. Formally, the cache state transition can be modelled as follows:

$$Z_{k+1}^c = \begin{cases} 0, & c = c^{evict}, \\ 1, & c = c_k, \\ Z_k^c, & \text{otherwise.} \end{cases} \quad (1)$$

Whenever a request $k$ arrives, a caching policy $\pi$ maps the current cache status vector $\boldsymbol{Z}_k$, the requested content $c_k$ and the feature vector of the requested content $\boldsymbol{f}_k$ to the new cache status vector $\boldsymbol{Z}_{k+1}$. The cache status vector can be updated based on $\pi$ as follows:

$$\boldsymbol{Z}_{k+1} = \pi(\boldsymbol{Z}_k | c_k, \boldsymbol{f}_k). \quad (2)$$

We use the term $H_\pi(K)$ as a cache hit rate metric to evaluate the efficiency of the caching system. It is defined as the ratio of requests that are served from the local cache to $K$ requests, which is given by:

$$H_\pi(K) = \frac{1}{K} \sum_k Z_k^{c_k}. \quad (3)$$

TABLE I
DETAILED FEATURES

| Category | Feature |
|---|---|
| contextual feature | access times in previous time step |
| | age = current time - publish time |
| semantic feature | type |
| | length |
| | area |
| | language |
| | score |
| | number of comments |
| | director |
| | performer |

Furthermore, $H_\pi$ is introduced to represent the long-term average cache hit rate when the number of requests goes to infinity over time by adopting caching policy $\pi$, which is written as:

$$H_\pi = \lim_{K \to \infty} H_\pi(K). \quad (4)$$

In the considered edge network, our objective is to maximize $H_\pi$ based on the constraint of cache capacity, and then find a policy $\pi$ for generating a series of popularity-aware replacement actions. Consequently, the corresponding problem can be formulated as:

$$\max_\pi \quad H_\pi \quad (5)$$

$$s.t. \quad \sum_c Z_k^c \le s, \quad \forall k, \quad (6)$$

$$Z_k^c \in \{0, 1\}, \quad \forall k, \forall c. \quad (7)$$

## V. PA-CACHE ALGORITHM DESIGN

In this section, we present the PA-Cache, a popularity-aware content caching approach that makes appropriate cache replacement decisions to maximize the long-term cache hit rate based on the estimated popularity. We begin by introducing the basic idea of the optimal replacement policy. We then present a deep learning-based approach for predicting the content popularity. Finally, we describe how the PA-Cache makes online replacement decisions.

### A. Basic Idea

Van Roy *et al.* [40] presented a proof showing that the replacement policy, named MIN [41] proposed by Belady is an optimal policy of the above problem. The MIN policy replaces the content in the cache, which has the longest time to be visited next time. As the future information is required in advance, it is an idealistic algorithm which is unimplementable in a real system. However, Belady's MIN algorithm gives a performance upper bound for content caching algorithms.

It is a challenging task to imitate the MIN algorithm using learning-based approaches directly. To pick out the content whose next request time is farthest in the future, learning-based approaches are required to predict the next request time of all contents in the cache precisely. Besides, running a predictor

for all contents in the cache upon each request will consume significant computing resources and time.

Therefore, in this paper, we aim to approximate the MIN algorithm by predicting content popularity within a specific time interval/step to address the aforementioned issues. We propose a popularity-aware content caching algorithm, named PA-Cache, which consists of two phases: *offline content popularity prediction* (in Section V-B) and *online replacement decision* (in Section V-C).

### B. Offline Content Popularity Prediction

As motivated earlier, we consider an evolving prediction task. The time period is partitioned into consecutive time steps indexed by $t = 1, \dots, T$. Our goal of evolving deep learning is to learn a function $F : \mathbb{R}^{m \times d} \to \mathbb{R}^m$ based on a sequence of training samples $\mathcal{D} = \{(\boldsymbol{x}_1, \boldsymbol{y}_1), \dots, (\boldsymbol{x}_t, \boldsymbol{y}_t), \dots, (\boldsymbol{x}_T, \boldsymbol{y}_T)\}$, that arrive sequentially, where $\boldsymbol{x}_t \in \mathbb{R}^{m \times d}$ represents the input features at time step $t$, $m$ is the number of instances, and $d$ is the feature dimension. $(\boldsymbol{x}_t^i, y_t^i)$ is the $i^{th}$ sample at time step $t$. The corresponding ground-truth content popularity is denoted as $\boldsymbol{y}_t \in \mathbb{R}^m$ while the predicted content popularity is denoted as $\hat{\boldsymbol{y}}_t \in \mathbb{R}^m$. The performance of learnt model is evaluated in terms of the cumulative prediction error of $m$ mini-batch instances.

*1) Feature Selection:* We consider two main types of features: the contextual features that would be time-varying; and the semantic features that are invariant with time. The detailed features are given in Table I. As the categorical features such as type and language are non-numeric features, which would not have a natural rank-ordering, we transform them into one-hot coding vectors, which is widely utilized in word embedding area [42]. Furthermore, as the numeric features (e.g., access times in previous time step, length) might have a wide range of values, their values are normalized into the range of $[0, 1]$.

*2) Evolving Deep Learning Model:* Given an input $\boldsymbol{x}_t$, the content popularity prediction task can be conducted by a conventional multi-layer recurrent neural network (RNN) with $L$ hidden layers. However, using such a model for evolving learning faces several challenges [43]: (i) *Model selection.* The network depth must be fixed in advance and cannot be changed. However, selecting an appropriate depth is a daunting task, especially in an evolving environment. For a small number of instances, a shallow network would be preferred as it converges quickly, while for a large number of instances, a deep network would be better at generalization and accurate prediction; (ii) *Convergence.* Critical issues of deep architecture consist of exploding or vanishing gradients, saddle point and diminishing feature reuse. These issues will result in slow or unstable convergence and be further exaggerated in the evolving setting.

To address these issues, we amend the multi-layer RNN architecture by attaching every hidden layer representation to an output regression for evolving learning through hedge backpropagation (HBP) [43]. HBP automatically determines how and when to accommodate network depth in an evolving manner. An evolving deep learning framework using HBP is depicted in Fig. 2.
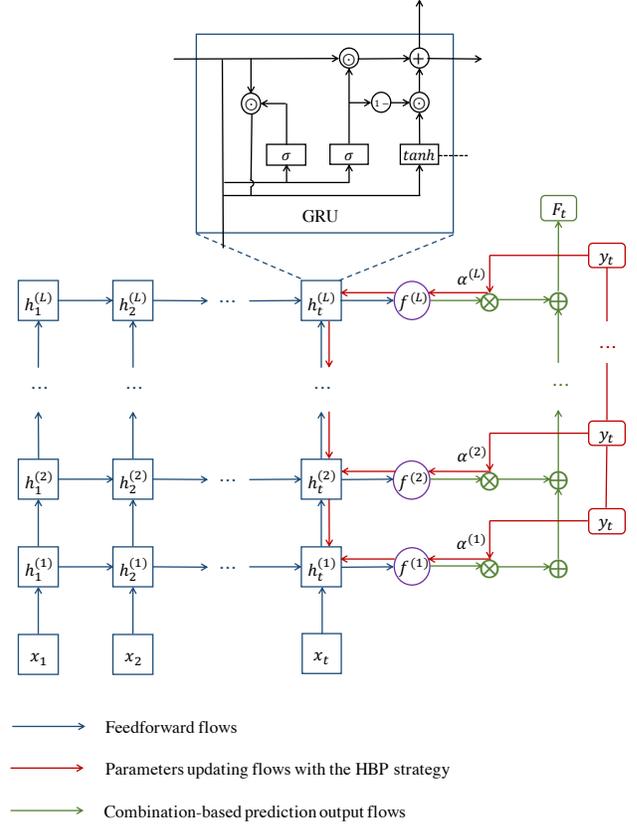


Fig. 2. Evolving deep learning framework using HBP.

In this paper, we consider a multi-layer RNN with $L$ hidden layers, where $L$ denotes the maximum depth of the network. Feature vectors are fed into a gated recurrent unit (GRU) layer, which is capable of capturing the long-term dependencies in sequential data and learning the time-varying patterns of user requests. The standard GRU architecture is based on two multiplicative gate units. The prediction function $F(\boldsymbol{x}_t)$ for the proposed evolving DNN is given by:

$$F(\boldsymbol{x}_t) = \sum_{l=1}^{L} \alpha^{(l)} f^{(l)}(\boldsymbol{x}_t), \quad \text{where} \tag{8}$$

$$f^{(l)}(\boldsymbol{x}_t) = \boldsymbol{\Theta}^{(l)} \boldsymbol{h}_t^{(l)},$$
$$\boldsymbol{r}_t^{(l)} = sigm(\boldsymbol{W}_r^{(l)} \boldsymbol{h}_t^{(l-1)} + \boldsymbol{U}_r^{(l)} \boldsymbol{h}_{t-1}^{(l)} + \boldsymbol{b}_r^{(l)}),$$
$$\boldsymbol{z}_t^{(l)} = sigm(\boldsymbol{W}_z^{(l)} \boldsymbol{h}_t^{(l-1)} + \boldsymbol{U}_z^{(l)} \boldsymbol{h}_{t-1}^{(l)} + \boldsymbol{b}_z^{(l)}),$$
$$\boldsymbol{h}_t^{(l)} = \boldsymbol{z}_t^{(l)} \odot \boldsymbol{h}_{t-1}^{(l)} + (\boldsymbol{1} - \boldsymbol{z}_t^{(l)}) \odot \tilde{\boldsymbol{h}}_t^{(l)},$$
$$\tilde{\boldsymbol{h}}_t^{(l)} = tanh(\boldsymbol{W}_h^{(l)} \boldsymbol{h}_t^{(l-1)} + \boldsymbol{U}_h^{(l)} (\boldsymbol{r}_t^{(l)} \odot \boldsymbol{h}_{t-1}^{(l)}) + \boldsymbol{b}_h^{(l)}),$$
$$\boldsymbol{h}_t^{(0)} = \boldsymbol{x}_t,$$

where $\boldsymbol{r}_t, \boldsymbol{z}_t, \boldsymbol{h}_t, \tilde{\boldsymbol{h}}_t$ are referred to as reset gate, update gate, hidden state vector and candidate hidden state vector respectively. The reset gate is used to specify how much previous knowledge should be ignored. The update gate helps the model to control the information that flows into memory. The candidate hidden state, also named intermediate memory unit, combines the knowledge from hidden states of the

previous time step and previous layer. It can facilitate the computation of subsequent hidden state. The parameters for training are the matrices $\{\boldsymbol{W}_r, \boldsymbol{W}_z, \boldsymbol{W}_h\}$ (the feedforward connection weights), $\{\boldsymbol{U}_r, \boldsymbol{U}_z, \boldsymbol{U}_h\}$ (the recurrent weights) and the vectors $\{\boldsymbol{b}_r, \boldsymbol{b}_z, \boldsymbol{b}_h\}$ (the bias). $sigm(\cdot)$ is the sigmoid activation function, which limits $\boldsymbol{r}_t$ and $\boldsymbol{z}_t$ to take values ranging from 0 and 1. The element-wise tensor multiplication is denoted by $\odot$.

Different from the conventional DNN, in which only the hidden state vector of last layer $\boldsymbol{h}^{(L)}$ is used by the regression to calculate the final predicted value, we utilize a weighted combination of regressions learned based on the multiple hidden state vectors from $\boldsymbol{h}^{(1)}, \ldots, \boldsymbol{h}^{(L)}$ in this paper. Two sets of new parameters to be learnt are introduced, i.e., $\boldsymbol{\Theta}^{(l)}$ and $\alpha^{(l)}$. Each regression in intermediate layer $f^{(l)}(\boldsymbol{x}_t)$ is parameterized by $\boldsymbol{\Theta}^{(l)}$. The final prediction of this model is a linear weighted sum of regressions $f^{(1)}(\boldsymbol{x}_t), \ldots, f^{(L)}(\boldsymbol{x}_t)$ , where the weight of each regression is denoted by $\alpha^{(l)} > 0$. The loss function of this model is defined as: $\mathcal{L}(F(\boldsymbol{x}_t), \boldsymbol{y}_t) = \sum_{l=1}^{L} \alpha^{(l)} \mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t)$. During the evolving learning procedure, parameters $\alpha^{(l)}$, $\boldsymbol{\Theta}^{(l)}$, $\boldsymbol{W}^{(l)}$, $\boldsymbol{U}^{(l)}$, $\boldsymbol{b}^{(l)}$ are required to be learnt.

We employ the HBP strategy [44] to learn $\alpha^{(l)}$. In the beginning, all weights $\alpha^{(l)}$ are uniformly split, i.e., $\alpha^{(l)} = \frac{1}{L}, l = 1, \ldots, L$. At each iteration, the regression of layer $l$ makes a prediction $f^{(l)}(\boldsymbol{x}_t)$. The weight of the regression then is updated as follows:

$$\alpha_{t+1}^{(l)} \leftarrow \alpha_t^{(l)} \beta^{\min(\mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \kappa)}, \tag{9}$$

where $\beta \in (0, 1)$ is the discount factor, $\kappa$ is the threshold parameter for smoothing "noisy" data. Thus, a regression's weight decays of a factor of $\beta^{\min(\mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \kappa)}$ in each iteration. At the end of each iteration, weights $\alpha$ are normalized to the interval $[0, 1]$, i.e., $\sum_{l=1}^{L} \alpha^{(l)} = 1$.

We adopt gradient descent methods to learn the parameters $\boldsymbol{\Theta}^{(l)}$ for all regressions, where the input to the $l^{th}$ regression is $\boldsymbol{h}^{(l)}$. It is identical with the update of the weights of the output layer in the conventional feedforward neural networks. This update is given by:

$$\begin{aligned} \boldsymbol{\Theta}_{t+1}^{(l)} &\leftarrow \boldsymbol{\Theta}_t^{(l)} - \eta \nabla_{\boldsymbol{\Theta}_t^{(l)}} \mathcal{L}(F(\boldsymbol{x}_t), \boldsymbol{y}_t) \\ &\leftarrow \boldsymbol{\Theta}_t^{(l)} - \eta \alpha^{(l)} \nabla_{\boldsymbol{\Theta}_t^{(l)}} \mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t). \end{aligned} \tag{10}$$

Updating the feature representation parameters $\boldsymbol{W}^{(l)}, \boldsymbol{U}^{(l)}$, $\boldsymbol{b}^{(l)}$ is not trivial. Unlike the original backpropagation scheme, in which the error derivatives are backpropagated from the output layer to each hidden layer, the error derivatives in this paper are backpropagated from each regression $f^{(l)}(\boldsymbol{x}_t)$. Thus, by combining the gradient descent methods and the dynamic objective function $\mathcal{L}(F(\boldsymbol{x}_t), \boldsymbol{y}_t) = \sum_{l=1}^{L} \alpha^{(l)} \mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t)$, the update rule for $\boldsymbol{W}^{(l)}, \boldsymbol{U}^{(l)}, \boldsymbol{b}^{(l)}$ is given as follows:

$$\boldsymbol{W}_{t+1}^{(l)} \leftarrow \boldsymbol{W}_t^{(l)} - \eta \sum_{j=l}^{L} \alpha^{(j)} \nabla_{\boldsymbol{W}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \tag{11}$$

$$\boldsymbol{U}_{t+1}^{(l)} \leftarrow \boldsymbol{U}_t^{(l)} - \eta \sum_{j=l}^{L} \alpha^{(j)} \nabla_{\boldsymbol{U}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \tag{12}$$

---

**Algorithm 1** Popularity prediction using evolving deep learning.

**Input:** Input features: $\boldsymbol{x}_t$; Revealed popularity: $\boldsymbol{y}_t$; Parameters for evolving DNN: $\{\beta, \eta, \zeta, \kappa, \alpha_t, \boldsymbol{\Theta}_t, \boldsymbol{W}_t, \boldsymbol{U}_t, \boldsymbol{b}_t\}$
**Output:** Predicted popularity: $\hat{\boldsymbol{y}}_t$
1: $\hat{\boldsymbol{y}}_t = F(\boldsymbol{x}_t) = \sum_{l=1}^{L} \alpha_t^{(l)} f^{(l)}(\boldsymbol{x}_t)$;
2: Calculate $\mathcal{L}_t^{(l)} = \mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \forall l = 1, \ldots, L$ by Eq. (14);
3: Update $\boldsymbol{\Theta}_{t+1}^{(l)}, \boldsymbol{W}_{t+1}^{(l)}, \boldsymbol{U}_{t+1}^{(l)}, \boldsymbol{b}_{t+1}^{(l)}, \forall l = 1, \ldots, L$ by Eq. (10) - Eq. (13);
4: Update $\alpha_{t+1}^{(l)} = \alpha_t^{(l)} \beta^{\min(\mathcal{L}_t^{(l)}, \kappa)}, \forall l = 1, \ldots, L$;
5: Smoothing $\alpha_{t+1}^{(l)} = \max(\alpha_{t+1}^{(l)}, \frac{\zeta}{L}), \forall l = 1, \ldots, L$;
6: Normalize $\alpha_{t+1}^{(l)} = \frac{\alpha_{t+1}^{(l)}}{Z_{t+1}}$ where $Z_{t+1} = \sum_{l=1}^{L} \alpha_{t+1}^{(l)}$;

---

$$\boldsymbol{b}_{t+1}^{(l)} \leftarrow \boldsymbol{b}_t^{(l)} - \eta \sum_{j=l}^{L} \alpha^{(j)} \nabla_{\boldsymbol{b}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t), \tag{13}$$

where $\nabla_{\boldsymbol{W}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t))$, $\nabla_{\boldsymbol{U}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t))$, $\nabla_{\boldsymbol{b}_t^{(l)}} \mathcal{L}(f^{(j)}(\boldsymbol{x}_t), \boldsymbol{y}_t))$ are computed via backpropagation from error derivatives of $f^{(j)}(\boldsymbol{x}_t)$. The summation part in Eq. (11)-(13) starts at $j = l$ as the shallower regression does not rely on the parameters $\boldsymbol{W}, \boldsymbol{U}, \boldsymbol{b}$ of deeper layers to make predictions.

Since shallower models are usually inclined to converge faster than deeper models [45], the weights of deeper regression might be diminished to a very small value by the HBP strategy. This will lead to a slow convergence in deeper regressions. Therefore, a smoothing parameter $\zeta \in (0, 1)$ is introduced to set a minimum weight for each regression. After the weight of the regressions are updated in each iteration according to Eq.(9), we further set the weights as $\alpha^{(l)} \leftarrow \max(\alpha^{(l)}, \frac{\zeta}{L})$, where $\zeta$ guarantees that each regression will be selected with at least probability $\frac{\zeta}{L}$. This balances the trade-off between exploration (all regressions at every depth will affect the backpropagation update) and exploitation.

*3) Loss Function:* Typically, DNN tracks the mean square error (MSE) as a metric when fitting the model. However, this metric is unable to make an accurate prediction when applied to the data that follows the heavy-tailed distribution, which is common in VoD systems [46]. It is because popular contents with high access times only account for a small fraction, but the arithmetic mean will be biased towards them. Hence, the overall performance of this metric would not be so satisfactory for the majority of contents. To eliminate this drawback, a mean relative squared error (MRSE) [47] is employed here, which is written as:

$$\mathcal{L}(f^{(l)}(\boldsymbol{x}_t), \boldsymbol{y}_t) = \frac{1}{m} \sum_{i=1}^{m} (\frac{f^{(l)}(\boldsymbol{x}_t^i)}{y_t^i} - 1)^2. \tag{14}$$

Algorithm 1 outlines the popularity prediction using evolving deep learning.

### C. Online Replacement Decision

When a request $k$ of content $c_k$ arrives, we first extract its latest features of the requested content and write them into

**Algorithm 2** PA-Cache algorithm.

**Input:** Request $k$

1: Update features vector for $c_k$ in *Feature Database*;
2: **if** $Z_{k+1}^{c_k} == 0$ **then**
3:      Remove the head element $c^{evict}$ from $Q$;
4:      Fetch $c_k$ from the upstream server;
5:      Insert $c_k$ and its estimated popularity into $Q$;
6: **end if**
7: Serve the client with $c(k)$;
8: **if** $t_k \bmod \phi == 0$ **then**
9:      Re-predict the popularity for all contents by Algorithm 1;
10:      Rebuild the priority queue $Q$;
11: **end if**

the *Feature Database* module. Then, the requested content is searched for in the local cache. If it is available in the cache, the user is directly served by the content replica from the cache node. Otherwise, PA-Cache removes the least popular content $c^{evict}$ in the cache to leave space for the new content $c_k$. *Popularity Database* module provides the estimated popularity which is updated by *Deep Neural Network* periodically described in Section V-B. Afterward, $c_k$ is fetched from the upstream server, stored in the cache, and transmitted to the user. PA-Cache manages a priority queue $Q$, which stores the cached contents along with their predicted popularities. The head element of $Q$ is regarded as the least popular content which can be quickly retrieved under this data structure. Each eviction operation will update $Q$ accordingly. To keep the prediction of content popularity up-to-date, PA-Cache triggers a *Learning Handler* periodically after every $\phi$ hours to update the prediction of content popularity. The PA-Cache algorithm is presented in Algorithm 2.

## VI. TRACE-DRIVEN EVALUATION RESULTS

In this section, we conduct extensive evaluations with real-world traces to evaluate the performance of PA-Cache.

### A. Dataset

The experiments are simulated on a real-world dataset derived from iQiYi[2], which is the largest online VoD service provider in China. This dataset contains 0.3 million unique videos watched by 2 million users over 2 weeks and has been widely used in previous works [32], [48]. The relevant information is recorded for each trace item as follows: (i) The device identifier (anonymized), which is unique for each device and will be used to label different users. (ii) Request time, which records the timestamp when the user request arrives; (iii) Video content, which involves the video name and some basic information, e.g., score, number of comments. In addition, we implement a crawler to collect more data to complement the features of videos (e.g., the area, type, language, length, publish date, director, performer).

[2]http://www.iqiyi.com

TABLE II
PARAMETER VALUES.

| Parameter | Value | Description |
|---|---|---|
| $C$ | 10,000 | Number of contents |
| $K$ | 446629 | Number of requests |
| $p$ | 0.1%-5.0% | Cache percentage |
| $L$ | 10 | Maximum capacity of DNN |
| $m$ | 128 | Mini-batch size |
| $\beta$ | 0.99 | Discount factor |
| $\kappa$ | 100 | Smoothing parameter of "noisy" data |
| $\zeta$ | 0.1 | Smoothing parameter of minimum weight |
| $\eta$ | 10 | Learning rate |
| $\phi$ | 1h | Time window |

### B. Algorithm Implementation

We implement a discrete event simulator based on the framework depicted in Fig. 1. To calculate and compare it with existing algorithms under restricted computational resources, we randomly sample $C = 10,000$ videos from the dataset 9 times. The cache percentage, which is the ratio between the cache size and the total number of unique contents, is considered to range from 0.1% to 5.0%. By default, it is set $p = 1.0\%$. The *Learning Handler* module trains a DNN in Fig. 2 after every 1 hour. The traces are divided into two periods. The first one is named *warm-up* period. It spans the first seven days of request traces, representing the input of evolving DNN. The second one is *test* period, which begins after the *warm-up* period. Unless explicitly clarified, the experimental results presented are all obtained under the above settings. PA-Cache is run on a PC with an Intel(R) Core(TM) i5-7360U CPU @ 2.30GHz, NVIDIA GTX 1080 GPU, and 8 GB RAM using MXNet framework. In this paper, we train a 10-layer evolving DNN with 512 units in the first layer and 16 units in the last hidden layer. The number of units from the $2^{nd}$ layer to the $9^{th}$ layer decreases gradually. Some key parameters, along with their descriptions and values, are listed in Table II.

### C. Benchmarks

We compare our algorithm PA-Cache against the following baselines:

- **Optimal** [41]. The cache runs an optimal and offline strategy for replacing the content, which has the longest time to be visited next time.
- **FNN-Caching** [31]. The cache employs FNN to predict content popularity and accordingly makes caching decisions.
- **NA-Caching** [26]. The cache uses deep reinforcement learning from its own experiences to make caching decisions based on the features of dynamic requests and caching space.
- **PopCaching** [35]. The cache clusters different contents into hypercubes based on the similarity between their access patterns and predicts their popularity when making replacement decisions.
- **LeCaR** [15]. The cache maintains two history entries based on recency policy and frequency policy. The weight of each entry is updated by regret minimization. It is used
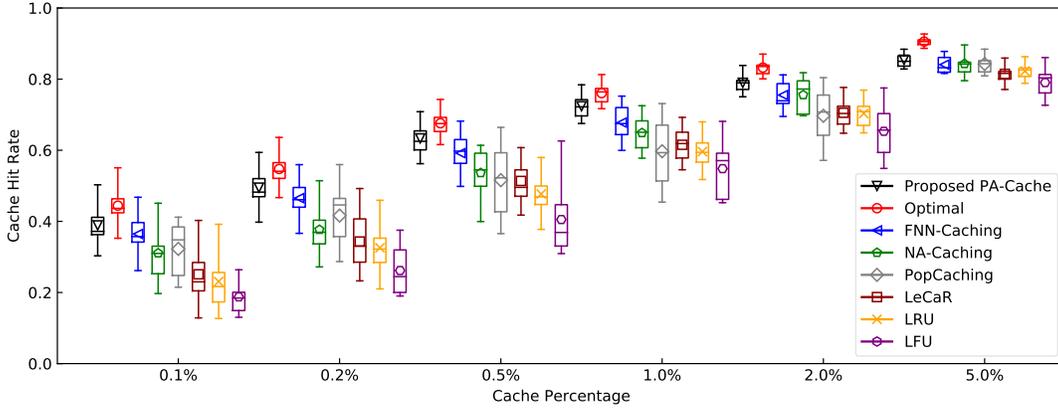
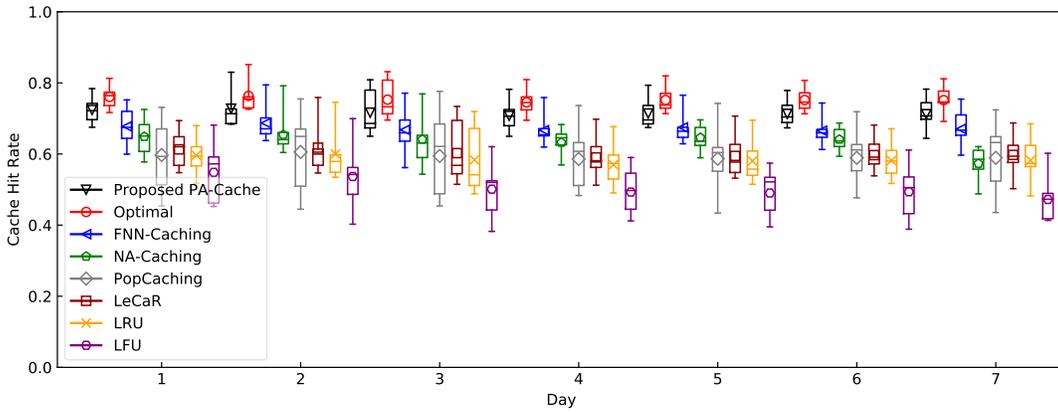Fig. 3. Cache hit rate under different cache percentages.



Fig. 4. Cache hit rate over time.

to determine which policy to be applied for the cache eviction.

- **LRU** [9]. The cache manages an ordered queue which records the recent accesses of all the cached contents. When the cache is full, the least recently accessed content will be replaced by the newly requested content.

- **LFU** [10]. The cache is implemented as a priority queue: the one that has been requested least frequently is replaced by the new content when the cache is full.

### D. Performance Comparison

The cache hit rate is depicted by the boxplot in Fig. 3, where the x-axis represents the cache percentage. The boxplot shows the means through the hollow markers, the medians through the bars, and the maximum and minimum values through the whiskers. It illustrates that PA-Cache significantly outperforms LRU, LFU, and their combination LeCaR in all cases. In particular, when the cache percentage is 0.1%, PA-Cache's average performance advantage over LFU, LRU, and LeCaR exceed 107.1%, 68.4%, and 54.8%, respectively. This is because these rule-based algorithms make eviction decisions based on simple metrics such as access time or frequency without considering the context information of contents. Moreover, when the content is evicted from the cache (e.g., LRU, LFU)

or history (e.g., LeCaR), all information about it is missing and cannot be utilized for future decision making.

Because the performance of the PopCaching algorithm is highly dependent on the hand-crafted features, we generate the content context vector several times and select the parameter setting, which achieves a $75^{th}$ percentile. It is observed that PA-Cache achieves a higher hit rate than FNN-Caching, NA-Caching, and PopCaching, especially when the cache percentage is small ($\leq 1.0\%$). For example, PA-Cache outperforms PopCaching, NA-Caching, and FNN-Caching by 22.8%, 18.1%, and 7.0% in the case where cache percentage is 0.5%. These results demonstrate that previous caching algorithms perform well under a certain scale of cache percentages, whereas PA-Cache can adapt to any cache percentages. For instance, when the cache percentage ranges from 0.2% to 1.0%, PA-Cache increases the hit rate by 7.9%-12.7% compared to PopCaching. Our interpretation is that PopCaching relies on the average sampling method for prediction, which is remarkable over hot contents but does not get an edge on lukewarm or cold contents.

Fig. 4 depicts how the cache hit rate of the algorithms varies over time. On average, we can see that PA-Cache outperforms FNN-Caching, NA-Caching, PopCaching, LeCaR, LRU, LFU by 6.4%, 13.2%, 20.9%, 19.0%, 22.5%, 42.1%, respectively. PA-Cache approximates the Optimal with only a
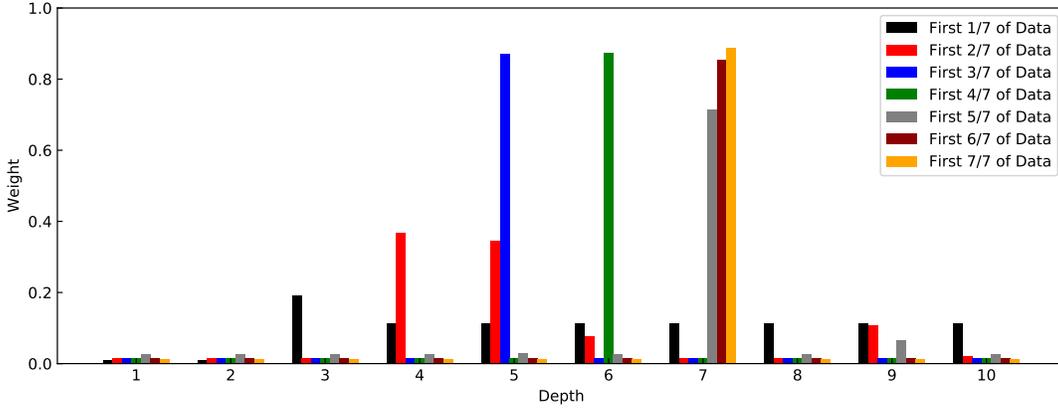
Fig. 5. Evolution of weight distribution of the DNN over time on the training set.

3.8% performance gap. This can be attributed to our model's prediction error. Besides, For the average hit rate, PA-Cache maintains a more stable cache hit rate with a smaller variance of $4.2 \times 10^{-5}$ compared to that of other methods (e.g., $5.0 \times 10^{-5}$ for FNN-Caching, $6.6 \times 10^{-4}$ for NA-Caching, $4.3 \times 10^{-5}$ for PopCaching, $8.7 \times 10^{-5}$ for LeCaR, $8.8 \times 10^{-5}$ for LRU, $6.4 \times 10^{-4}$ for LFU). In particular, NA-Caching and LFU exhibit the highest variances. This is because LFU predicts content popularity based on historical frequency without decay, which would be outdated over time. NA-Caching suffers from large-delayed rewards for some contents, which might lead to convergence to poor action choices and slow policy learning. On the contrary, PA-Cache could quickly adapt to changes in the workload over time.

Fig. 5 illustrates the evolution of the weight distribution of the DNN over time on the training set for the sampled dataset, which achieves the median performance. Initially (first 1/7 data), the maximum weight has gone to the shallower regression in the $3^{rd}$ layer. Then (first 2/7 data), the maximum weights have moved to the regressions in the $4^{th}$ and $5^{th}$ layers. As more data arrives sequentially, deeper regressions will pick up higher weights. For instance, in the segment with the first 6/7 data, the $7^{th}$ layer has obtained the highest weight. It shows that our evolving DNN is capable of selecting an appropriate model automatically.

We compare the convergence behavior of our evolving DNN with a conventional 10-layer DNN in Fig. 6. It shows the variation of loss as the batch number increases. We can see that the loss of the 10-layer conventional DNN converges to a local optimum after about the $120^{th}$ batches, while our evolving DNN converges much more quickly. It means that our evolving DNN could benefit from the fast convergence of its shallow networks at the beginning. Moreover, Fig. 7 depicts the boxplots for 800 loss values obtained during the training after 200 batches by the conventional 10-layer DNN and our evolving DNN. In this plot, the boxes relate to the interquartile range; small circles denote the outliers; the upper and lower whiskers represent loss values outside the middle 50%. We can observe that both the 25th and 75th percentile of evolving DNN are lower than those of conventional DNN.
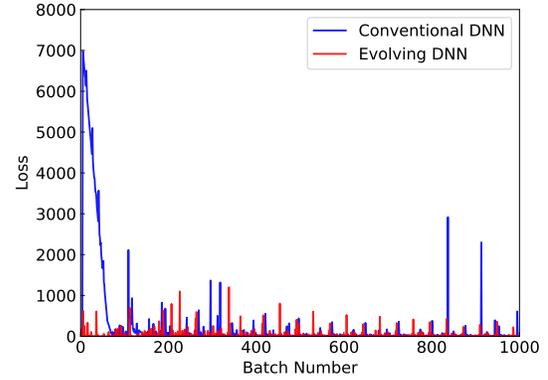


Fig. 6. Convergence behavior of the conventional DNN and our evolving DNN over batch number.
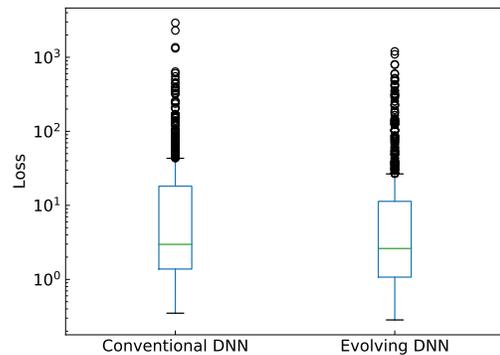


Fig. 7. Boxplots for 800 loss values during training after 200 batches by the conventional DNN and our evolving DNN.

The median value of evolving DNN and conventional DNN is 2.97 and 2.61, respectively, represented by the bar in the box. The boxplots indicate that the evolving DNN obtains both smoother variance and lower median loss. It illustrates that our evolving DNN keeps the merits of a powerful representation of the deeper network.
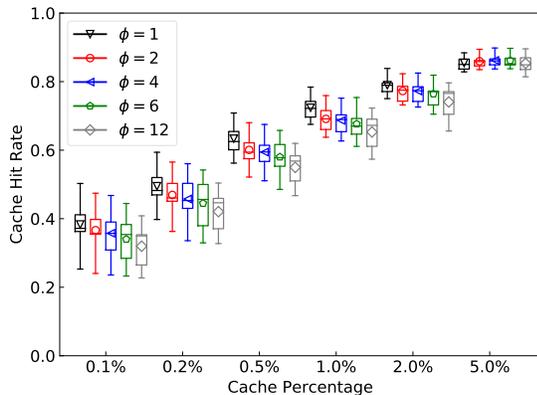
Fig. 8. Cache hit rate under different cache percentages for different $\phi$.

### E. Practical Consideration

The first consideration is the frequency of re-estimation for content popularity. Although the system is capable of correcting deviation between the predicted popularity and the real one and adapt to time-varying demands more quickly if updating the evolving DNN more frequently, each update incurs additional overhead. Fig. 8 illustrates how different values of parameter $\phi$ affect the caching performance under different cache percentages. It is observed that the smaller the $\phi$ is, the higher the cache hit rate achieves In almost all cases, other than the case in which the cache percentage is 5.0%. When the cache percentage ranges from 0.1% to 2.0%, the cache hit rate with $\phi = 1$ is 2.2%-5.4% higher than that with $\phi = 2$, but the performances under $\phi = 2$ and $\phi = 4$ are very similar. After the cache percentage increases to $\phi = 5.0\%$, the cache hit rate is almost constant even if a different $\phi$ is selected. This observation indicates that the edge cache could select an appropriate parameter $\phi$ based on its storage capacity. That is, when the storage capacity is small, the cache node could optimize the parameter by seeking a trade-off between prediction accuracy and computational cost, whereas, for the large capacity, our algorithm scales to time windows with different sizes for updating.

The historical information is available for existing contents, but new contents generated from content providers may be added to the system continually. Although the historical information of $\phi$ hours in our PA-Cache, much smaller than PopCaching or FNN-Caching, is enough for prediction. The cold start problem of new contents is still crucial for the performance improvement of the caching system. We can employ a small LRU-cache space [49] to handle the traffic demand of new contents for which we do not estimate.

### VII. Conclusion

In this paper, we have presented the design, implementation, and evaluation of PA-Cache, a novel popularity-aware edge content caching algorithm in edge networks. It makes adaptive caching decisions with the aim to maximize the cache hit rate in the long term. PA-Cache can effectively tackle the daunting task of content caching upon time-varying traffic demands.

It has combined the strength of multi-layer RNN in learning the comprehensive representations of requested contents and the hedge backpropagation strategy that automatically determines how and when to accommodate network's depth in an evolving manner. The experiments on real-world traces have been conducted, and the performance of PA-Cache has been evaluated. Trace-driven evaluation results have demonstrated the effectiveness and superiority of our proposed PA-Cache in terms of the hit rate and computational cost compared to existing popular caching algorithms.

### References

[1] M. I. A. Zahed, I. Ahmad, D. Habibi, Q. V. Phung, and L. Zhang, "A cooperative green content caching technique for next generation communication networks," *IEEE Transactions on Network and Service Management*, vol. 17, no. 1, pp. 375–388, Mar. 2020.

[2] M. Carrie and R. David, "Worldwide global datasphere iot device and data forecast, 2019-2023," *Market Forcast*, 2019.

[3] S. Shen, Y. Han, X. Wang, and Y. Wang, "Computation offloading with multiple agents in edge-computing–supported iot," *ACM Transactions on Sensor Networks*, vol. 16, no. 1, pp. 1–27, Dec. 2019.

[4] K. Wang, H. Yin, W. Quan, and G. Min, "Enabling collaborative edge computing for software defined vehicular networks," *IEEE Network*, vol. 32, no. 5, pp. 112–117, Sep./Oct. 2018.

[5] X. Xia, F. Chen, Q. He, J. Grundy, M. Abdelrazek, and H. Jin, "Online collaborative data caching in edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 2, pp. 281–294, Aug. 2020.

[6] F. Zafari, J. Li, K. K. Leung, D. Towsley, and A. Swami, "Optimal energy consumption for communication, computation, caching, and quality guarantee," *IEEE Transactions on Control of Network Systems*, vol. 7, no. 1, pp. 151–162, Mar. 2019.

[7] F. Wang, F. Wang, J. Liu, R. Shea, and L. Sun, "Intelligent video caching at network edge: A multi-agent deep reinforcement learning approach," in *Proc. of IEEE INFOCOM*, Toronto, Canada, Jul. 2020, pp. 2499–2508.

[8] X. Xia, F. Chen, Q. He, J. Grundy, M. Abdelrazek, and H. Jin, "Cost-effective app data distribution in edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 2, pp. 31–44, Jul. 2020.

[9] M. Ahmed, S. Traverso, P. Giaccone, E. Leonardi, and S. Niccolini, "Analyzing the performance of LRU caches under non-stationary traffic patterns," *arXiv preprint arXiv:1301.4909*, 2013.

[10] A. Jaleel, K. B. Theobald, S. C. Steely Jr, and J. Emer, "High performance cache replacement using re-reference interval prediction (RRIP)," *ACM SIGARCH Computer Architecture News*, vol. 38, no. 3, pp. 60–71, Jun. 2010.

[11] M. Abrams, C. R. Standridge, G. Abdulla, S. Williams, and E. A. Fox, "Caching proxies: limitations and potentials," Department of Computer Science, Virginia Polytechnic Institute & State University, Tech. Rep., 1995.

[12] E. J. O'neil, P. E. O'neil, and G. Weikum, "The LRU-K page replacement algorithm for database disk buffering," *ACM SIGMOD Record*, vol. 22, no. 2, pp. 297–306, Jun. 1993.

[13] N. Megiddo and D. S. Modha, "ARC: A self-tuning, low overhead replacement cache," in *Proc. of USENIX FAST*, San Francisco, USA, Mar. 2003, pp. 115–130.

[14] Q. Huang, K. Birman, R. van Renesse, W. Lloyd, S. Kumar, and H. C. Li, "An analysis of facebook photo caching," in *Proc. of ACM SOSP*, Pennsylvania, USA, Nov. 2013, p. 167–181.

[15] G. Vietri, L. V. Rodriguez, W. A. Martinez, S. Lyons, J. Liu, R. Rangaswami, M. Zhao, and G. Narasimhan, "Driving cache replacement with ML-based LeCaR," in *Proc. of USENIX HotStorage*, Boston, USA, Jul. 2018, pp. 1–6.

[16] H. Bahn, K. Koh, S. H. Noh, and S. Lyul, "Efficient replacement of nonuniform objects in web caches," *Computer*, vol. 35, no. 6, pp. 65–73, Jun. 2002.

[17] D. S. Berger, N. Beckmann, and M. Harchol-Balter, "Practical bounds on optimal caching with variable object sizes," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 2, no. 2, pp. 1–38, Jun. 2018.

[18] V. Martina, M. Garetto, and E. Leonardi, "A unified approach to the performance analysis of caching systems," in *Proc. of IEEE INFOCOM*, Toronto, Canada, May. 2014, pp. 2040–2048.

[19] H. Che, Y. Tung, and Z. Wang, "Hierarchical web caching systems: Modeling, design and experimental results," *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 7, pp. 1305–1314, Sep. 2002.

[20] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and Zipf-like distributions: Evidence and implications," in *Proc. of IEEE INFOCOM*, New York, USA, Mar. 1999, pp. 126–134.

[21] J. Li, S. Shakkottai, J. C. Lui, and V. Subramanian, "Accurate learning or fast mixing? dynamic adaptability of caching algorithms," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 6, pp. 1314–1330, Jun. 2018.

[22] C. Zhong, M. C. Gursoy, and S. Velipasalar, "A deep reinforcement learning-based framework for content caching," in *Proc. of CISS*, Princeton, USA, Mar. 2018, pp. 1–6.

[23] V. Kirilin, A. Sundarrajan, S. Gorinsky, and R. K. Sitaraman, "RL-Cache: Learning-based cache admission for content delivery," *IEEE Journal on Selected Areas in Communications*, pp. 2372–2385, Oct. 2020.

[24] A. Sadeghi, G. Wang, and G. B. Giannakis, "Deep reinforcement learning for adaptive caching in hierarchical content delivery networks," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 4, pp. 1024–1033, Dec. 2019.

[25] X. Wang, C. Wang, X. Li, V. C. Leung, and T. Taleb, "Federated deep reinforcement learning for internet of things with decentralized cooperative edge caching," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9441–9455, Oct. 2020.

[26] Q. Fan, X. Li, S. Wang, S. Fu, X. Zhang, and Y. Wang, "NA-Caching: An adaptive content management approach based on deep reinforcement learning," *IEEE Access*, vol. 7, pp. 152 014–152 022, Oct. 2019.

[27] Y. Guan, X. Zhang, and Z. Guo, "CACA: Learning-based content-aware cache admission for video content in edge caching," in *Proc. of ACM Multimedia*, Nice, France, Oct. 2019, pp. 456–464.

[28] G. Yan and J. Li, "RL-Bélády: A unified learning framework for content caching," in *Proc. of ACM Multimedia*, Seattle, USA, Oct. 2020, pp. 1009–1017.

[29] D. S. Berger, "Towards lightweight and robust machine learning for CDN caching," in *Proc. of ACM HotNets*, Washington, USA, Nov. 2018, pp. 134–140.

[30] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *Proc. of AAAI*, New Orleans, USA, Feb. 2018, pp. 3207–3214.

[31] V. Fedchenko, G. Neglia, and B. Ribeiro, "Feedforward neural networks for caching: enough or too much?" *ACM SIGMETRICS Performance Evaluation Review*, vol. 46, no. 3, pp. 139–142, Jan. 2019.

[32] H. Pang, J. Liu, X. Fan, and L. Sun, "Toward smart and cooperative edge caching for 5g networks: A deep learning based approach," in *Proc. of IEEE/ACM IWQoS*, Banff, Canada, Jun. 2018, pp. 1–6.

[33] Z. Song, D. S. Berger, K. Li, and W. Lloyd, "Learning relaxed belady for content distribution network caching," in *Proc. of USENIX NSDI*, Santa Clara, USA, Feb. 2020, pp. 529–544.

[34] B. Chen, L. Liu, M. Sun, and H. Ma, "Iotcache: Toward data-driven network caching for internet of things," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 10 064–10 076, Dec. 2019.

[35] S. Li, J. Xu, M. Van Der Schaar, and W. Li, "Popularity-driven content caching," in *Proc. of IEEE INFOCOM*, San Francisco, USA, Apr. 2016, pp. 1–9.

[36] S. Müller, O. Atan, M. van der Schaar, and A. Klein, "Context-aware proactive content caching with service differentiation in wireless networks," *IEEE Transactions on Wireless Communications*, vol. 16, no. 2, pp. 1024–1036, Feb. 2016.

[37] L. Gkatzikis, V. Sourlas, C. Fischione, and I. Koutsopoulos, "Low complexity content replication through clustering in content-delivery networks," *Computer Networks*, vol. 121, pp. 137–151, Jul. 2017.

[38] Q. Fan, H. Yin, G. Min, S. Wang, Y. Lyu, and X. Zhang, "Content and network aware replication and scheduling mechanism for user generated content videos," in *Proc. of IEEE UIC*, Leicester, UK, Aug. 2019, pp. 739–746.

[39] Q. Fan, H. Yin, L. Jiao, Y. Lv, H. Huang, and X. Zhang, "Towards optimal request mapping and response routing for content delivery networks," *IEEE Transactions on Services Computing*, pp. 1–1, Jan. 2018.

[40] B. Van Roy, "A short proof of optimality for the min cache replacement algorithm," *Information Processing Letters*, vol. 102, no. 2-3, pp. 72–73, Apr. 2007.

[41] L. A. Belady, "A study of replacement algorithms for a virtual-storage computer," *IBM Systems Journal*, vol. 5, no. 2, pp. 78–101, 1966.

[42] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. of NIPS*, Nevada, USA, Dec. 2013, pp. 3111–3119.

[43] D. Sahoo, Q. Pham, J. Lu, and S. C. Hoi, "Online deep learning: learning deep neural networks on the fly," in *Proc. of IJCAI*, Stockholm, Sweden, Jul. 2018, pp. 2660–2666.

[44] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, Aug. 1997.

[45] G. Larsson, M. Maire, and G. Shakhnarovich, "Fractalnet: Ultra-deep neural networks without residuals," *arXiv preprint arXiv:1605.07648*, 2016.

[46] C. Zhang, Y. Xu, Y. Zhou, and X. Fu, "On the "familiar stranger" phenomenon in a large-scale vod system," in *Proc. of IEEE ICC Workshops*, Atlanta, USA, May. 2017, pp. 928–933.

[47] H. Pinto, J. M. Almeida, and M. A. Gonçalves, "Using early view patterns to predict the popularity of youtube videos," in *Proc. of ACM WSDM*, Rome, Italy, Feb. 2013, pp. 365–374.

[48] G. Ma, Z. Wang, M. Zhang, J. Ye, M. Chen, and W. Zhu, "Understanding performance of edge content caching for mobile video streaming," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 5, pp. 1076–1089, Mar. 2017.

[49] D. Applegate, A. Archer, V. Gopalakrishnan, S. Lee, and K. Ramakrishnan, "Optimal content placement for a large-scale vod system," *IEEE/ACM Transactions on Networking*, vol. 24, no. 4, pp. 2114–2127, Aug. 2016.