Deep Reinforcement Learning for Adaptive Caching in Hierarchical Content Delivery Networks

Alireza Sadeghi[®], *Student Member, IEEE*, Gang Wang[®], *Member, IEEE*, and Georgios B. Giannakis[®], *Fellow, IEEE*

Abstract—Caching is envisioned to play a critical role 2 in next-generation content delivery infrastructure, cellular 3 networks, and Internet architectures. By smartly storing the most 4 popular contents at the storage-enabled network entities during 5 off-peak demand instances, caching can benefit both network 6 infrastructure as well as end users, during on-peak periods. 7 In this context, distributing the limited storage capacity across 8 network entities calls for decentralized caching schemes. Many 9 practical caching systems involve a parent caching node con-10 nected to multiple leaf nodes to serve user file requests. To model 11 the two-way interactive influence between caching decisions at the 12 parent and leaf nodes, a reinforcement learning (RL) framework 13 is put forth. To handle the large continuous state space, a scal-14 able deep RL approach is pursued. The novel approach relies 15 on a hyper-deep Q-network to learn the Q-function, and thus 16 the optimal caching policy, in an online fashion. Reinforcing the 17 parent node with ability to learn-and-adapt to unknown poli-18 cies of leaf nodes as well as spatio-temporal dynamic evolution 19 of file requests, results in remarkable caching performance, as 20 corroborated through numerical tests.

21 *Index Terms*—Caching, deep RL, deep Q-network, next-22 generation networks, function approximation.

23

I. INTRODUCTION

²⁴ **I** N LIGHT of the tremendous growth of data traffic ²⁵ **I** over both wireline and wireless communications, next-²⁶ generation networks including future Internet architectures, ²⁷ content delivery infrastructure, and cellular networks stand ²⁸ in need of emerging technologies to meet the ever-increasing ²⁹ data demand. Recognized as an appealing solution is *caching*, ³⁰ which amounts to storing reusable contents in geographi-³¹ cally distributed storage-enabled network entities so that future ³² requests for those contents can be served faster. The ratio-³³ nale is that unfavorable shocks of peak traffic periods can be ³⁴ smoothed by proactively storing 'anticipated' highly popular ³⁵ contents at those storage devices and during off-peak peri-³⁶ ods [1], [2]. Caching popular content is envisioned to achieve

Manuscript received February 26, 2019; revised June 30, 2019; accepted August 10, 2019. This work was supported in part by NSF grants 1711471, 1514056, and 1901134. The associate editor coordinating the review of this article and approving it for publication was H. T. Dinh. (*Corresponding author: Gang Wang.*)

The authors are with the Digital Technology Center, University of Minnesota, Minneapolis, MN 55455 USA, and also with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: sadeghi@umn.edu; gangwang@umn.edu; georgios@umn.edu).

Digital Object Identifier 10.1109/TCCN.2019.2936193

major savings in terms of energy, bandwidth, and cost, in 37 addition to user satisfaction [1]. 38

To fully unleash its potential, a content-agnostic caching 39 entity has to rely on available observations to learn what 40 and when to cache. Toward this goal, contemporary machine 41 learning and artificial intelligence tools hold the promise to 42 empower next-generation networks with 'smart' caching con-43 trol units, that can learn, track, and adapt to unknown dynamic 44 environments, including space-time evolution of content popu-45 larities and network topology, as well as entity-specific caching 46 policies. 47

Deep neural networks (DNNs) have lately boosted 48 the notion of "learning from data" with field-changing 49 performance improvements reported in diverse artificial intelli-50 gence tasks [3]. DNNs can cope with the 'curse of dimension-51 ality' by providing compact low-dimensional representations 52 of high-dimensional data [4]. Combining deep learning with 53 **RL**, deep (D) RL has created the first artificial agents to 54 achieve human-level performance across many challenging 55 domains [5], [6]. As another example, a DNN system was 56 built to operate Google's data centers, and shown able to con-57 sistently achieve a 40% reduction in energy consumption for 58 cooling [7]. This system provides a general-purpose frame-59 work to understand complex dynamics, which has also been 60 applied to address other challenges including, e.g., dynamic 61 spectrum access [8], multiple access and handover con-62 trol [9], [10], as well as resource allocation in fog-radio access 63 networks [11], [12] or software-defined networks [13], [14]. 64

A. Prior Art on Caching

Early approaches to caching include the least recently used 66 (LRU), least frequently used (LFU), first in first out (FIFO), 67 random replacement (RR) policies, and their variants. Albeit 68 simple, these schemes cannot deal with the dynamics of con-69 tent popularities and network topologies. Recent efforts have 70 gradually shifted toward developing learning and optimization 71 based approaches that can 'intelligently' manage the cache 72 resources. For unknown but time-invariant content populari-73 ties, multi-armed bandit online learning was pursued in [15]. 74 Yet, these methods are generally not amenable to online 75 implementation. To serve delay-sensitive requests, a learning 76 approach was developed in [16] using a pre-trained DNN to 77 handle a non-convex problem reformulation. 78

In realistic networks however, popularities exhibit dynamics, which motivate well the so-termed *dynamic* caching. A

2332-7731 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

65

81 Poisson shot noise model was adopted to approximate the 82 evolution of popularities in [17], for which an age-based ⁸³ caching solution was developed in [18]. RL based methods ⁸⁴ have been pursued in [19], [20], [21], [22]. Specifically, a 85 Q-learning based caching scheme was developed in [19] to 86 model global and local content popularities as Markovian 87 processes. Considering Poisson shot noise popularity dynam-88 ics, a policy gradient RL based caching scheme was devised ⁸⁹ in [20]. Assuming stationary file popularities and service costs, ⁹⁰ a dual-decomposition based Q-learning approach was pursued ⁹¹ in [21]. Albeit reasonable for discrete states, these approaches 92 cannot deal with large continuous state-action spaces. To cope ⁹³ with such spaces, DRL approaches have been considered ⁹⁴ for content caching in, e.g., [6], [22], [23], [24], [25], [26]. 95 Encompassing finite-state time-varying Markov channels, a 96 deep Q-network approach was devised in [22]. An actor-critic 97 method with deep deterministic policy gradient updates was 98 used in [23]. Boosted network performance using DRL was ⁹⁹ documented in several other applications, such as connected vehicular networks [24], and smart cities [25]. 100

The aforementioned works focus on devising caching poli-101 102 cies for a single caching entity. A more common setting ¹⁰³ in next-generation networks however, involves a network of interconnected caching nodes. It has been shown that con-104 105 sidering a network of connected caches jointly can further ¹⁰⁶ improve performance [27], [28]. For instance, leveraging 107 network topology and the broadcast nature of links, the coded 108 caching strategy in [27] further reduces data traffic over a 109 network. This idea has been extended in [29] to an online 110 setting, where popularities are modeled Markov processes. Collaborative and distributed online learning approaches have 111 ¹¹² been pursued [28], [30], [31]. Indeed, today's content delivery 113 networks such as Akamai [32], have tree network structures. 114 Accounting for the hierarchy of caches has become a common ¹¹⁵ practice in recent works; see also [33], [34], [35]. Joint rout-¹¹⁶ ing and in-network content caching in a hierarchical cache 117 network was formulated in [33], for which greedy schemes with provable performance guarantees can be found in [34]. 118

¹¹⁹ We identify the following challenges that need to be ¹²⁰ addressed when designing practical caching methods for next-¹²¹ generation networks.

122 c1) Networked Caching: Caching decisions of a node, in
 a network of caches, influences decisions of all other
 nodes. Thus, a desired caching policy must adapt to the
 network topology and policies of neighboring nodes.

c2) Complex Dynamics: Content popularities are random
 and exhibit unknown space-time, heterogeneous, and
 often non-stationary dynamics over the entire network.

c3) Large Continuous State Space: Due to the shear size of
 available content, caching nodes, and possible realiza tions of content requests, the decision space is huge.

132 B. This Work

Prompted by the recent interest in hierarchical caching, this paper focuses on a two-level network caching, where a partise ent node is connected to multiple leaf nodes to serve end-user file requests. Such a two-level network constitutes the buildtise popular tree hierarchical cache networks in,



Fig. 2. A hierarchical tree network cache system.

e.g., [32]. To model the interaction between caching decisions ¹³⁸ of parent and leaf nodes along with the space-time evolution ¹³⁹ of file requests, a scalable DRL approach based on hyper deep ¹⁴⁰ Q-networks (DQNs) is developed. As corroborated by extensive numerical tests, the novel caching policy for the parent ¹⁴² node can adapt itself to local policies of leaf nodes and spacetime evolution of file requests. Moreover, our approach is ¹⁴⁴ simple-to-implement, and performs close to the optimal policy. ¹⁴⁵

II. MODELING AND PROBLEM STATEMENT

146

Consider a two-level network of interconnected caching ¹⁴⁷ nodes, where a parent node is connected to N leaf nodes, ¹⁴⁸ indexed by $n \in \mathcal{N} := \{1, \ldots, N\}$. The parent node is con- ¹⁴⁹ nected to the cloud through a (typically congested) back-haul ¹⁵⁰ link; see Fig. 1. One could consider this network as a part of ¹⁵¹ a large hierarchical caching system, where the parent node is ¹⁵² connected to a higher level caching node instead of the cloud; ¹⁵³ see Fig. 2. In a content delivery network for instance, edge ¹⁵⁴ servers (a.k.a. points of presence or PoPs) are the leaf nodes, ¹⁵⁵ and a fog server acts as the parent node. Likewise, (small) base ¹⁵⁶ stations in a 5G cellular network are the leaf nodes, while a ¹⁵⁷ serving gate way (S-GW) may be considered as the parent ¹⁵⁸ node; see also [36, p. 110].

All nodes in this network store files to serve file requests. ¹⁶⁰ Every leaf node serves its locally connected end users, by ¹⁶¹ providing their requested files. If a requested content is locally ¹⁶² available at a leaf node, the content will be served immediately ¹⁶³ at no cost. If it is not locally available due to limited caching ¹⁶⁴ capacity, the content will be fetched from its parent node, at ¹⁶⁵ a certain cost. Similarly, if the file is available at the parent ¹⁶⁶ node, it will be served to the leaf at no cost; otherwise, the ¹⁶⁷ file must be fetched from the cloud at a higher cost. ¹⁶⁸

To mitigate the burden with local requests on the network, 169 each leaf node stores 'anticipated' locally popular files. In 170 ¹⁷¹ addition, this paper considers that each parent node stores ¹⁷² files to serve requests that are *not* locally served by leaf nodes. ¹⁷³ Since leaf nodes are closer to end users, they frequently receive ¹⁷⁴ file requests that exhibit rapid temporal evolution at a *fast* ¹⁷⁵ timescale. The parent node on the other hand, observes aggre-¹⁷⁶ gate requests over a large number of users served by the *N* ¹⁷⁷ leaf nodes, which naturally exhibit smaller fluctuations and ¹⁷⁸ thus evolve at a *slow* timescale.

This motivated us to pursue a two-timescale approach to managing such a network of caching nodes. To that end, let $\tau = 1, 2, ...$ denote the slow time intervals, each of which is further divided into *T* fast time slots indexed by t = 1, ..., T; see Fig. 3 for an illustration. Each fast time slot may be, use cach slow time interval is a period of say 4-5 minutes. We assume that the network state remains unchanged during the each fast time slot *t*, but can change from *t* to t + 1.

Consider a total of F files in the cloud, which are col-188 189 lected in the set $\mathcal{F} = \{1, \ldots, F\}$. At the beginning of each 190 slot t, every leaf node n selects a subset of files in \mathcal{F} to 191 prefetch and store for possible use in this slot. To deter-192 mine which files to store, every leaf node relies on a local ¹⁹³ caching policy function denoted by π_n , to take (cache or no-¹⁹⁴ cache) action $\boldsymbol{a}_n(t+1,\tau) = \pi_n(\boldsymbol{s}_n(t,\tau))$ at the beginning 195 of slot t + 1, based on its state vector \boldsymbol{s}_n at the end of slot t. We assume this action takes a negligible amount of 196 197 time relative to the slot duration; and define the state vector 198 $\boldsymbol{s}_n(t,\tau) \coloneqq \boldsymbol{r}_n(t,\tau) \coloneqq [r_n^1(t,\tau) \cdots r_n^F(t,\tau)]^\top$ to collect the ¹⁹⁹ number of requests received at leaf node n for individual files 200 over the duration of slot t on interval τ . Likewise, to serve file 201 requests that have not been served by leaf nodes, the parent ²⁰² node takes action $a_0(\tau)$ to store files at the beginning of every ²⁰³ interval τ , according to a certain policy π_0 . Again, as aggrega-204 tion smooths out request fluctuations, the parent node observes 205 slowly varying file requests, and can thus make caching deci-²⁰⁶ sions at a relatively slow timescale. In the next section, we 207 present a two-timescale approach to managing such a network 208 of caching nodes.

209 III. TWO-TIMESCALE PROBLEM FORMULATION

File transmission over any network link consumes resources, including, e.g., energy, time, and bandwidth. Hence, serving requested file that is not locally stored at a node, incurs a cost. Among possible choices, the present paper considers the following cost for node $n \in \mathcal{N}$, at slot t + 1 of interval τ

where $\mathbf{c}_n(\cdot) \coloneqq [c_n^1(\cdot) \cdots c_n^F(\cdot)]^\top$ concatenates the cost for serving individual files per node n; symbol \odot denotes entrywise vector multiplication; entries of \mathbf{a}_0 and \mathbf{a}_n are either 1 (cache, hence no need to fetch), or, 0 (no-cache, hence fetch); and 1 stands for the all-one vector. Specifically, the second summand in (1) captures the cost of the leaf node fetching files for end users, while the first summand corresponds to that of the parent fetching files from the cloud.





We model user file requests as Markov processes with ²²⁶ unknown transition probabilities [19]. Per interval τ , a ²²⁷ reasonable caching scheme for leaf node $n \in \mathcal{N}$ ²²⁸ could entail minimizing the expected cumulative cost; ²²⁹ that is, ²³⁰

$$\pi_{n,\tau}^* \coloneqq \underset{\pi_n \in \Pi_n}{\operatorname{arg\,min}} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}^\top \boldsymbol{c}_n(\pi_n(\boldsymbol{s}_n(t,\tau)), \boldsymbol{r}_n(t+1,\tau), \boldsymbol{a}_0(\tau)) \right]$$
(2) 23(2)

where Π_n represents the set of all feasible policies for ²³³ node *n*. Although solving (2) is in general challenging, efficient near-optimal solutions have been introduced in several ²³⁵ recent contributions; see, e.g., [15], [19], [20], and references therein. In particular, a RL based approach using tabular ²³⁷ *Q*-learning was pursued in our precursor [19], which can ²³⁸ be employed here to tackle this fast timescale optimization. ²³⁹ The remainder of this paper will thus be on designing the ²⁴⁰ caching policy π_0 for the parent node, that can learn, track, ²⁴¹ and adapt to the leaf node policies as well as user file ²⁴² requests. ²⁴³

At the beginning of every interval $\tau + 1$, the parent node 244 takes action to refresh its cache; see Fig. 4. To design a 245 caching policy that can account for both leaf node poli- 246 cies and file popularity dynamics, the parent collects local 247 information from all leaf nodes. At the end of interval τ , each 248 leaf node $n \in \mathcal{N}$ first obtains the time-average state vector 249 $\bar{s}_n(\tau) := (1/T) \sum_{t=1}^T s_n(t, \tau)$, and subsequently forms and 250 forwards the per-node vector $\bar{s}_n(\tau) \odot (1 - \pi_n(\bar{s}_n(\tau)))$ to its 251 parent node. This vector has nonzero entries the average number of file requests received during interval τ , and zero entries 253 if π_n stores the corresponding files. Using the latter, the parent node forms its 'weighted' state vector as 255

$$\boldsymbol{s}_0(\tau) \coloneqq \sum_{n=1}^N w_n \bar{\boldsymbol{s}}_n(\tau) \odot (\boldsymbol{1} - \pi_n(\bar{\boldsymbol{s}}_n(\tau))) \tag{3} 256$$

where the weights $\{w_n \geq 0\}$ control the influence of leaf ²⁵⁷ nodes $n \in \mathcal{N}$ on the parent node's policy. Similarly, having ²⁵⁸ received at the end of interval $\tau + 1$ the slot-averaged costs ²⁵⁹

$$\bar{\boldsymbol{c}}_n\Big(\pi_0(\boldsymbol{s}_0(\tau)); \{\boldsymbol{a}_n(t,\tau+1), \boldsymbol{r}_n(t,\tau+1)\}_{t=1}^T\Big)$$
 260

$$= \frac{1}{T} \sum_{t=1}^{T} \boldsymbol{c}_n(\boldsymbol{a}_n(t,\tau+1), \boldsymbol{r}_n(t,\tau+1), \boldsymbol{a}_0(\tau+1)) (4)$$

²⁶² from all leaf nodes, the parent node determines its cost

263
$$\boldsymbol{c}_{0}(\boldsymbol{s}_{0}(\tau), \pi_{0}(\boldsymbol{s}_{0}(\tau)))$$

264 $= \sum_{n=1}^{N} w_{n} \bar{\boldsymbol{c}}_{n} \left(\boldsymbol{\pi}_{0}(\boldsymbol{s}_{0}(\tau)); \{ \boldsymbol{a}_{n}(t, \tau+1), \boldsymbol{r}_{n}(t, \tau+1) \}_{t=1}^{T} \right) \right).$
265 (5)

²⁶⁶ Having observed $\{s_0(\tau'), a_0(\tau'), c_0(s_0(\tau'-1), a_0(\tau')\}_{\tau'=1}^{\tau},$ ²⁶⁷ the objective is to take a well-thought-out action $a_0(\tau+1)$ for ²⁶⁸ next interval. This is tantamount to finding an optimal policy ²⁶⁹ function π_0^* to take a caching action $a_0(\tau+1) = \pi_0^*(s_0(\tau))$. ²⁷⁰ Since the requests $\{r_n(t,\tau)\}_{n,t}$ are Markovian and present ²⁷¹ actions $\{a_n(t,\tau)\}_{n,t}$ also affect future costs, the optimal ²⁷² RL policy for the parent node will minimize the expected ²⁷³ cumulative cost over all leaf nodes in the long term, namely

$$\pi_{0}^{*:=} \underset{\pi_{0} \in \Pi_{0}}{\operatorname{arg\,min}} \mathbb{E}\left[\sum_{\tau=1}^{\infty} \gamma^{\tau-1} \mathbf{1}^{\top} \boldsymbol{c}_{0}(\boldsymbol{s}_{0}(\tau), \pi_{0}(\boldsymbol{s}_{0}(\tau)))\right]$$
(6)

where Π_0 represents the set of all feasible policies; expectation is over $\{\boldsymbol{r}_n(t,\tau)\}_{n,t}$, as well as (possibly random) parent $\boldsymbol{a}_0(\tau)$ and leaf actions $\{\boldsymbol{a}_n(t,\tau)\}_{n,t}$; and the discount factor $\gamma \in [0,1)$ trades off emphasis of current versus future costs. It is evident from (6) that the decision taken at a given state $\boldsymbol{s}_0(\tau)$, namely $\boldsymbol{a}_0(\tau+1) = \pi(\boldsymbol{s}_0(\tau))$, influences the next state $\boldsymbol{s}_0(\tau+1)$ through $\pi_{n,\tau}(\cdot)$ in (3), as well as the cost $\boldsymbol{c}_0(\cdot)$ in (5). Therefore, problem (6) is a discounted infinite time horizon Markov decision process (MDP). Finding the optimal policy of an MDP is NP-hard [37]. To cope with this complexity of solving (6), an adaptive RL approach is pursued next.

286 IV. ADAPTIVE RL-BASED CACHING

RL deals with action-taking policy function estimation in an environment with dynamically evolving states, so as to minimize a long-term cumulative cost. By interacting with the environment (through successive actions and observed states and costs), RL seeks a policy function (of states) to draw actions from, in order to minimize the average cumulative cost as in (6) [38]. To proceed, we start by giving some basics on Markov decision processes (MDPs) [38, p. 310].

295 A. MDP Formulation

MDPs provide a well-appreciated model for decision mak-296 ²⁹⁷ ing tasks. It is represented by a tuple $\langle S, A, C, P \rangle$, where S 298 denotes a set of possible states, A a set of feasible actions, C299 a set of costs, and \mathcal{P} the set of state transition probabilities. 300 In our problem of interest, at the end of each interval τ , the ³⁰¹ parent node is in state $s_0(\tau) \in S$ and takes a caching decision $a_0(\tau+1) \in \mathcal{A}$ for the next interval $\tau+1$, according to its ³⁰³ local caching policy $\boldsymbol{a}_0(\tau+1) = \pi_0(\boldsymbol{s}_0(\tau))$. Upon proceeding to interval $\tau + 1$, every leaf node *n* serves its locally connected 305 users. Let vector $\boldsymbol{r}_n(t,\tau+1)$ collect the number of received 306 requests at node n during slot t across files. We model the 307 temporal evolution of this vector with Markov dynamics as 308 $\boldsymbol{r}_n(t+1,\tau+1) = [\boldsymbol{r}_n(t,\tau+1) + \boldsymbol{\Delta}_n(t,\tau+1)]^+$, where 309 $\Delta_n(t, \tau + 1)$ is a multivariate Gaussian random noise; $|\cdot|$ ³¹⁰ and $(\cdot)^+$ denote the entry-wise floor and $\max\{\cdot, 0\}$ operators. In this context, the incurred cost $c_0(s_0(\tau), \pi_0(s_0(\tau))) \in C$ (see (5)), is a random vector with an unknown distribution. ³¹² As requests $\{r_n\}$ are random, caching decisions $\{a_n\}$ are ³¹³ also random. In addition, decision a_0 of the parent node influences those of leaf nodes via (2), as well as the next state of ³¹⁵ the parent node in (3). That is, the current decision of parent ³¹⁶ node probabilistically influences its next state. To formalize ³¹⁷ this, we use $P^a_{ss'} \in \mathcal{P}$ to denote the *unknown* state transition ³¹⁸ probability from state *s* to *s'* upon taking action *a* ³¹⁹

$$P_{ss'}^{a} = \Pr\{s(\tau+1) = s' \mid s(\tau) = s, \ a = \pi(s)\}$$
(7) 320

where the subscript 0 referring to the parent node is dropped $_{321}$ for brevity. As a_0 probabilistically influences the next state as $_{322}$ well as the incurred cost, our problem is indeed an MDP. The $_{323}$ rest of this paper targets finding an optimal policy solving (6). $_{324}$

B. RL Based Caching

starting from initial state $s_0(0)$ as

Towards developing an RL solver for (6), define the so- $_{326}$ termed value function to indicate the quality of policy π_0 , $_{327}$

325

328

341

$$V_{\pi_0}(\boldsymbol{s}_0(0)) \coloneqq \mathbb{E}\left[\sum_{\tau=1}^{\infty} \gamma^{\tau-1} \mathbf{1}^{\top} \boldsymbol{c}_0(\boldsymbol{s}_0(\tau), \pi_0(\boldsymbol{s}_0(\tau)))\right]$$
(8) 326

which represents the average cost of following policy π_0 ³³⁰ to make caching decisions starting from $s_0(0)$. Upon finding $V_{\pi_0}(\cdot)$ for all $\pi_0 \in \Pi_0$, one can readily obtain π_0^* that ³³² minimizes $V_{\pi_0}(s_0(0))$ over all possible initial states $s_0(0)$. ³³³

For brevity, the time index and the subscript 0 referring to the parent node will be dropped whenever it is clear from the context. To find $V_{\pi}(\cdot)$, one can rely on the Bellman equation, which basically relates the value of a policy at one state to values of the remaining states [38, p. 46]. Leveraging (8) and (7), the Bellman equation for value function $V_{\pi}(\cdot)$ is given by 339

$$V_{\pi}(\boldsymbol{s}) = \mathbb{E}\left[\mathbf{1}^{\top}\boldsymbol{c}(\boldsymbol{s},\pi(\boldsymbol{s})) + \gamma \sum_{\boldsymbol{s}'} P_{\boldsymbol{s}\boldsymbol{s}'}^{\pi(\boldsymbol{s})} V_{\pi}(\boldsymbol{s}')\right] \quad (9) \text{ }_{340}$$

where the average immediate cost can be found as

$$\mathbb{E}\Big[\mathbf{1}^{\top}\boldsymbol{c}(\boldsymbol{s},\pi(\boldsymbol{s}))\Big] = \sum_{\boldsymbol{s}'} P_{\boldsymbol{s}\boldsymbol{s}'}^{\pi(\boldsymbol{s})} \mathbf{1}^{\top}\boldsymbol{c}\big(\boldsymbol{s},\pi(\boldsymbol{s})|\boldsymbol{s}'\big).$$

If $P_{ss'}^{a}$ were known $\forall a \in \mathcal{A}, \forall s, s' \in \mathcal{S}$, finding $V_{\pi}(\cdot)$ ³⁴³ would be equivalent to solving the system of linear equations (9). Indeed, if one could afford the complexity of ³⁴⁵ evaluating $V_{\pi}(\cdot)$ for all $\pi \in \Pi_0$, the optimal policy π^* is the ³⁴⁶ one that minimizes the value function for all states. However, ³⁴⁷ the large (possibly infinite) number of policies in practice ³⁴⁸ discourages such an approach. An alternative is to employ ³⁴⁹ the so-termed policy iteration algorithm [38, p. 64] outlined ³⁵⁰ next. Define first the state-action value function, also called ³⁵¹ Q-function for policy π ³⁵²

$$Q_{\pi}(\boldsymbol{s}, \boldsymbol{a}) := \mathbb{E}\Big[\mathbf{1}^{\top} \boldsymbol{c}(\boldsymbol{s}, \boldsymbol{a})\Big] + \gamma \sum_{\boldsymbol{s}'} P^{\boldsymbol{a}}_{\boldsymbol{s}\boldsymbol{s}'} V_{\pi}(\boldsymbol{s}').$$
 (10) 353

This function captures the expected immediate cost of starting ³⁵⁴ from state *s*, taking the first action to be *a*, and subse- ³⁵⁵ quently following policy π to take future actions onwards. ³⁵⁶ The only difference between the value function in (8) and that ³⁵⁷ ³⁵⁸ of *Q*-function in (10) is that the former takes the first action ³⁵⁹ according to policy π , while the latter starts with **a** as the first ³⁶⁰ action, which may not necessarily be taken when adhering ³⁶¹ to π . Having defined the *Q*-function, we are ready to present ³⁶² the policy iteration algorithm, in which every iteration *i* entails ³⁶³ the following two steps:

Policy Evaluation: Find $V_{\pi^i}(\cdot)$ for the current policy π^i by solving the system of linear equations in (9).

³⁶⁶ *Policy Improvement:* Update the current policy greedily as

367
$$\pi^{i+1}(\boldsymbol{s}) = \arg\min_{\boldsymbol{\alpha}\in\mathcal{A}} Q_{\pi^i}(\boldsymbol{s},\boldsymbol{\alpha}).$$

To perform policy evaluation, we rely on knowledge of $P_{ss'}^{\pi^i(s)}$. However, this is impractical in our setup that involves dynamic evolution of file requests, and unknown caching policies for the leaf nodes. This calls for approaches that target directly the optimal policy π^* , without knowing $P_{ss'}^{a}$, $\forall a, s, s'$. One such approach is *Q*-learning [38, p. 107]. In the ensuing section, we first introduce a *Q*-learning based adaptive caching scheme, which is subsequently generalized in Section V by invoking a DQN.

377 C. Q-Learning Based Adaptive Caching

The *Q*-learning algorithm finds the optimal policy π^* by 379 estimating $Q_{\pi^*}(\cdot, \cdot)$ 'on-the-fly.' It turns out that π^* is the 380 greedy policy over the optimal *Q*-function [38, p. 64], that is

$$\pi^*(\boldsymbol{s}) = \arg\min_{\boldsymbol{\alpha}\in\mathcal{A}} Q_{\pi^*}(\boldsymbol{s},\boldsymbol{\alpha}), \quad \forall \boldsymbol{s}\in\mathcal{S}$$
(11)

where Q_{π^*} is estimated using Bellman's equation for the Qfunction. This is possible because the V-function is linked with the Q-function under π^* through (see [38, p. 51] for details)

$$V_{\pi^*}(\boldsymbol{s}) = \min_{\boldsymbol{\alpha} \in \mathcal{A}} Q_{\pi^*}(\boldsymbol{s}, \boldsymbol{\alpha}), \quad \forall \boldsymbol{s}.$$
(12)

³⁸⁶ Substituting (12) into (10), Bellman's equation for the Q-³⁸⁷ function under π^* is expressed as

388
$$Q_{\pi^*}(\boldsymbol{s}, \boldsymbol{a}) = \mathbb{E}\Big[\mathbf{1}^\top \boldsymbol{c}(\boldsymbol{s}, \boldsymbol{a})\Big] + \gamma \sum_{\boldsymbol{s}'} P^{\boldsymbol{a}}_{\boldsymbol{s}\boldsymbol{s}'} \min_{\boldsymbol{\alpha}} Q_{\pi^*}(\boldsymbol{s}', \boldsymbol{\alpha})$$
 (13)

³⁸⁹ which plays a key role in many RL algorithms. Examples ³⁹⁰ include *Q*-learning [39], and SARSA [38], where one relies ³⁹¹ on (13) to update estimates of the *Q*-function in a stochas-³⁹² tic manner. In particular, the *Q*-learning algorithm follows an ³⁹³ exploration-exploitation procedure to take some action \boldsymbol{a} in ³⁹⁴ a given state \boldsymbol{s} . Specifically, it chooses the action minimiz-³⁹⁵ ing the current estimate of $Q_{\pi^*}(\cdot, \cdot)$ denoted by $\hat{Q}_{\tau}(\cdot, \cdot)$, with ³⁹⁶ probability (w.p.) $1 - \epsilon_{\tau}$, or, it takes a random action $\boldsymbol{a} \in \mathcal{A}$ ³⁹⁷ otherwise; that is,

398
$$\boldsymbol{a} = \begin{cases} \arg\min_{\boldsymbol{\alpha} \in \mathcal{A}} \hat{Q}_{\tau}(\boldsymbol{s}, \boldsymbol{\alpha}), \text{ w.p. } 1 - \epsilon_{\tau} \\ \boldsymbol{\alpha} \in \mathcal{A} \\ \operatorname{random} \boldsymbol{a} \in \mathcal{A}, \quad \operatorname{w.p.} \epsilon_{\tau}. \end{cases}$$

³⁹⁹ After taking action \boldsymbol{a} , moving to some new state \boldsymbol{s}' , and incur-⁴⁰⁰ ring cost \boldsymbol{c} , the *Q*-learning algorithm adopts the following loss ⁴⁰¹ function for the state-action pair $(\boldsymbol{s}, \boldsymbol{a})$

$${}_{402} \quad \mathcal{L}(\boldsymbol{s}, \boldsymbol{a}) = \frac{1}{2} \left(\mathbf{1}^{\top} \boldsymbol{c}(\boldsymbol{s}, \boldsymbol{a}) + \gamma \min_{\boldsymbol{\alpha} \in \mathcal{A}} \hat{Q}_{\tau}(\boldsymbol{s}', \boldsymbol{\alpha}) - \hat{Q}_{\tau}(\boldsymbol{s}, \boldsymbol{a}) \right)^{2}.$$

$${}_{403} \qquad (14)$$



Fig. 5. Deep Q-network.

The estimated Q-function for a *single* state-action pair is subsequently updated, by following a gradient descent step to 405 minimize the loss in (14), which yields the update 406

$$\hat{Q}_{\tau+1}(\boldsymbol{s}, \boldsymbol{a}) = \hat{Q}_{\tau}(\boldsymbol{s}, \boldsymbol{a}) - \beta \frac{\partial \mathcal{L}(\boldsymbol{s}, \boldsymbol{a})}{\partial \hat{Q}_{\tau}(\boldsymbol{s}, \boldsymbol{a})}$$
(15) 407

where $\beta > 0$ is some step size. Upon evaluating the gradient 408 and merging terms, the update in (15) boils down to 409

$$\hat{Q}_{ au+1}(oldsymbol{s},oldsymbol{a}) = (1-eta)\hat{Q}_{ au}(oldsymbol{s},oldsymbol{a})$$
 410

+
$$\beta \left[\mathbf{1}^{\top} \boldsymbol{c}(\boldsymbol{s}, \boldsymbol{a}) + \gamma \min_{\boldsymbol{\alpha} \in \mathcal{A}} \hat{Q}_{\tau}(\boldsymbol{s}', \boldsymbol{\alpha}) \right].$$
 411

 Three remarks are worth making at this point.
 412

 Remark 1: As far as the fast-timescale caching strategy of
 413

 leaf nodes is concerned, multiple choices are possible, includ 414

 ing, e.g., LRU, LFU, FIFO, [40], the multi-armed bandit
 415

 scheme [15], and even RL ones [19], [20].
 416

Remark 2: The exploration-exploitation step for taking ⁴¹⁷ actions guarantees continuously visiting state-action pairs, and ⁴¹⁸ ensures convergence to the optimal *Q*-function [37]. Instead ⁴¹⁹ of the ϵ -greedy exploration-exploitation step, one can employ ⁴²⁰ the upper confidence bound scheme [41]. Technically, any ⁴²¹ exploration-exploitation scheme should be greedy in the limit ⁴²² of infinite exploration (GLIE) [42, p. 840]. An example obey-⁴²³ ing GLIE is the ϵ -greedy algorithm [42, p. 840] with $\epsilon_{\tau} = 1/\tau$. ⁴²⁴ It converges to an optimal policy, albeit at a very slow rate. On ⁴²⁵ the other hand, using a constant $\epsilon_{\tau} = \epsilon$ approaches the optimal ⁴²⁶ $Q^*(\cdot, \cdot)$ faster, but its exact convergence is not guaranteed as ⁴²⁷ it is not GLIE.

Clearly, finding $Q_{\pi^*}(\boldsymbol{s}, \boldsymbol{a})$ entails estimating a function 429 defined over state and action spaces. In several applications 430 however, at least one of the two vector variables is either con- 431 tinuous or takes values from an alphabet of high cardinality. 432 Revisiting every state-action pair in such settings is impossi- 433 ble, due to the so-called curse of dimensionality – a typical 434 case in practice. To deal with it, function approximation tech- 435 niques offer as a promising solution [38]. These aim at finding 436 the original Q-function over all feasible state-action pairs, by 437 judicious generalization from a few observed pairs. Early func- 438 tion approximators for RL design good hand-crafted features 439 that can properly approximate $V_{\pi^*}(\cdot)$, $Q_{\pi^*}(\cdot, \cdot)$, or, $\pi^*(\cdot)$ [5]. 440 Their applicability has only been limited to domains, where 441 such features can be discovered, or, to state spaces that are 442 low dimensional [38]. 443

Deep learning approaches on the other hand, have recently 444 445 demonstrated remarkable potential in applications such as 446 object detection, speech recognition, and language transla-447 tion, to name a few. This is because DNNs are capable extracting compact low-dimensional features from high-448 Of 449 dimensional data. Wedding DNNs with RL results in 'deep RL' that can effectively deal with the curse of dimensional-450 451 ity, by eliminating the need of hand-crafting good features. These considerations have inspired the use of DNNs to esti-452 ⁴⁵³ mate either the *Q*-function, value function, or, the policy. The ⁴⁵⁴ most remarkable success in playing AI games adopted DNNs estimate the *Q*-function, what is termed DQN in [5]. 455 to

Prompted by this success, the next leverages the power of
 function approximation through DNNs to develop an adaptive
 aching scheme for the parent node.

V. ADAPTIVE DQN-BASED CACHING

To find the optimal caching policy π^* for the parent node, the success of function approximators for RL (e.g., [5]), motivated us to pursue a parametric approach to estimating $Q(\boldsymbol{s}, \boldsymbol{a})$ with a DNN (see (11)). This DNN has as input pairs of vectors $(\boldsymbol{s}, \boldsymbol{a})$, and as scalar output the corresponding estimated $Q(\boldsymbol{s}, \boldsymbol{a})$ values. Clearly, the joint state-action space of the sought Q-function has cardinality $|\mathcal{A} \times \mathcal{S}| = |\mathcal{A}| |\mathcal{S}|$. To reduce the search space, we shall take a parametric DQN approach [5] that we adapt to our setup, as elaborated next.

Consider a deep feedforward NN with *L* fully connected lay-470 ers, with input the $F \times 1$ state vector $\mathbf{s}(\tau)$ as in (3), and $F \times 1$ 471 output cost vector $\mathbf{o}(\tau+1)$ [5]. Note that input does not include 472 the action vector \mathbf{a} . Each hidden layer $l \in \{2, \ldots, L-1\}$ 473 comprises n_l neurons with rectified linear unit (ReLU) acti-474 vation functions $h(z) := \max(0, z)$ for $z \in \mathbb{R}$; see, e.g., [43]. 475 Neurons of the *L*-th output layer use a softmax nonlinearity¹ 476 to yield for the *f*-th entry of \mathbf{o} , an estimated long term cost 477 o_f that would incur, if file *f* is not stored at the cache.

If this cache memory can store up to M (< F) files at the parent node, the M largest entries of the DQN output $o(\tau+1)$ are chosen by the decision module in Fig. 5 to obtain the action vector $a(\tau + 1)$. We can think of our DQN as a 'soft policy' function estimator, and $o(\tau+1)$ as a 'predicted cost' or a 'soft action' vector for interval $\tau + 1$, whose 'hard thresholding' wields $a(\tau + 1)$. It will turn out that excluding a from the DQN input and picking it up at the output lowers the search space from $|\mathcal{A}||S|$ to |S|.

To train our reduced-complexity DQN amounts to finding a 488 set of weights collected in the vector $\boldsymbol{\theta}$ that parameterizes the 489 input-output relationship $\boldsymbol{o}(\tau + 1) = Q(\boldsymbol{s}(\tau); \boldsymbol{\theta}_{\tau})$. To recur-490 sively update $\boldsymbol{\theta}_{\tau}$ to $\boldsymbol{\theta}_{\tau+1}$, consider two successive intervals 491 along with corresponding states $\boldsymbol{s}(\tau)$ and $\boldsymbol{s}(\tau+1)$; the action 492 $\boldsymbol{a}(\tau+1)$ taken at the beginning of interval $\tau+1$; and, the cost 493 $\boldsymbol{c}(\tau+1)$ revealed at the end of interval $\tau+1$. The instantaneous 494 approximation of the optimal cost-to-go from interval $\tau+1$ is 495 given by $\boldsymbol{c}(\tau+1) + \gamma Q(\boldsymbol{s}(\tau+1); \boldsymbol{\theta}_{\tau})$, where $\boldsymbol{c}(\tau+1)$ is the 496 immediate cost, and $Q(\boldsymbol{s}(\tau+1); \boldsymbol{\theta}_{\tau})$ represents the predicted 497 cost-to-go from interval $\tau+2$ that is provided by our DQN with θ_{τ} , and discounted by γ . Since our DQN offers $Q(s(\tau); \theta_{\tau})$ 498 as the predicted cost for interval $\tau + 1$, the prediction error of 499 this cost as a function of θ_{τ} is given by 500

$$\boldsymbol{\delta}(\boldsymbol{\theta}_{\tau}) \coloneqq \begin{bmatrix} \text{target cost-to-go from interval } \tau+1 \\ \boldsymbol{c}(\tau+1) + \gamma Q(\boldsymbol{s}(\tau+1); \boldsymbol{\theta}_{\tau}) & -Q(\boldsymbol{s}(\tau); \boldsymbol{\theta}_{\tau}) \end{bmatrix} \quad \text{50}$$

$$\tilde{\odot} (\mathbf{1} - \mathbf{a}(\tau + 1))$$
(16) 502

and has non-zero entries for files not stored at interval $\tau + 1$. 503 Using the so-termed experience $\boldsymbol{E}_{\tau+1} \coloneqq [\boldsymbol{s}(\tau), \boldsymbol{a}(\tau + 504 1), \boldsymbol{c}(\tau+1), \boldsymbol{s}(\tau+1)]$, and the ℓ_2 -norm of $\boldsymbol{\delta}(\boldsymbol{\theta}_{\tau})$ as criterion 505

$$\mathcal{L}(oldsymbol{ heta}_ au) = \|oldsymbol{\delta}(oldsymbol{ heta}_ au)\|_2^2$$
 (17) 506

510

520

the sought parameter update minimizing (17) is given by the 507 stochastic gradient descent (SGD) iteration as 508

$$\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_{\tau} - \beta_{\tau} \nabla \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{\tau}}$$
 (18) 505

where $\beta_{\tau} > 0$ denotes the learning rate.

Since the dimensionality of θ can be much smaller than ⁵¹¹ |S||A|, the DQN is efficiently trained with few experi-⁵¹² ences, and generalizes to unseen state vectors. Unfortunately, ⁵¹³ DQN model inaccuracy can propagate in the cost prediction ⁵¹⁴ error in (16) that can cause instability in (18), which can ⁵¹⁵ lead to performance degradation, and even divergence [44]. ⁵¹⁶ Moreover, (18) leverages solely a single most recent experience $E_{\tau+1}$. These limitations will be mitigated as elaborated ⁵¹⁸ next. ⁵¹⁹

A. Target Network and Experience Replay

NN function approximation, along with the loss (17) and ${}_{521}$ the update (18), often result in unstable RL algorithms [5]. ${}_{522}$ This is due to: i) correlated experiences used to update the ${}_{523}$ DQN parameters θ ; and, ii) the influence of any change in ${}_{524}$ policy on subsequent experiences and vice versa.

Possible remedies include the so-called *experience replay* ⁵²⁶ and *target network* to update the DQN weights. In experience replay, the parent node stores all past experiences E_{τ} ⁵²⁷ in $\mathcal{E} := \{E_1, \ldots, E_{\tau}\}$, and utilizes a batch of *B* uniformly ⁵²⁹ sampled experiences from this data set, namely $\{E_{i\tau}\}_{i=1}^B \sim$ ⁵³⁰ $U(\mathcal{E})$. By sampling and replaying previously observed experiences, experience replay can overcome the two challenges. On the other hand, to obtain decorrelated target values in (16), a ⁵³³ second NN (called target network) with structure identical to ⁵³⁴ the DQN is invoked with parameter vector θ^{Tar} . Interestingly, ⁵⁵⁶ θ^{Tar} can be periodically replaced with θ_{τ} every *C* training ⁵³⁶ iterations of the DQN, which enables the target network to ⁵³⁷ smooth out fluctuations in updating the DQN [5].

With a randomly sampled experience $E_{i_{\tau}} \in \mathcal{E}$, the 539 prediction error with the target cost-to-go estimated using the 540 target network (instead of the DQN) is 541

$$\boldsymbol{\delta}^{\mathrm{Tar}}(\boldsymbol{\theta}; \boldsymbol{E}_{i_{\tau}}) := \begin{bmatrix} \boldsymbol{c}(i_{\tau}+1) + \gamma Q \big(\boldsymbol{s}(i_{\tau}+1); \boldsymbol{\theta}^{\mathrm{Tar}} \big) - Q(\boldsymbol{s}(i_{\tau}); \boldsymbol{\theta}) \end{bmatrix} \qquad 542$$

$$(1 - a(i_{\tau} + 1)).$$
 (19) 543

Different from (16), the target values here are found through 544 the target network with weights θ^{Tar} . In addition, the error 545 in (16) is found by using the most recent experience, while 546

459

¹Softmax is a function that takes as input a vector $\boldsymbol{z} \in \mathbb{R}^{F}$, and normalizes it into a probability distribution via $\sigma(\boldsymbol{z})_{f} = \mathrm{e}^{z_{f}} / (\sum_{f=1}^{F} \mathrm{e}^{z_{f}}), \forall f$.

Algorithm 1: Deep RL for Adaptive Caching **Initialize:** s(0), $s_n(t,\tau)$, $\forall n$, θ_{τ} , and θ^{Tar} 1 for $\tau = 1, 2, ...$ do Take action $\boldsymbol{a}(\tau)$ via exploration-exploitation 2 $\int \text{Best files via } Q(\boldsymbol{s}(\tau-1);\boldsymbol{\theta}_{\tau}) \text{ w.p. } 1 - \epsilon_{\tau}$ $\boldsymbol{a}(\tau) =$ 3 random $\boldsymbol{a} \in \mathcal{A}$ w.p. ϵ_{τ} for $t = 1, \ldots, T$ do 4 for n = 1, ..., N do 5 Take action \boldsymbol{a}_n using local policy 6 $\boldsymbol{a}_n(t,\tau) = \begin{cases} \pi_n(\boldsymbol{s}(t-1,\tau)) \text{ if } t \neq 1 \\ \pi_n(\boldsymbol{s}(T,\tau-1)) \text{ if } t = 1 \end{cases}$ 7 Requests $\boldsymbol{r}_n(t,\tau)$ are revealed 8 Set $\boldsymbol{s}_n(t,\tau) = \boldsymbol{r}_n(t,\tau)$ 9 Incur $\boldsymbol{c}_n(\cdot)$, see (1) 10 11 end end 12 Leaf nodes 13 Set $\bar{\boldsymbol{s}}_n(\tau) \coloneqq (1/T) \sum_{t=1}^T \boldsymbol{s}_n(t,\tau)$ 14 Send $\bar{\boldsymbol{s}}_n(\tau) \odot (\boldsymbol{1} - \pi_n(\bar{\boldsymbol{s}}_n))$ to parent node 15 Send $\bar{c}_n(\cdot)$ see (4), to parent node 16 Parent node 17 Set $\boldsymbol{s}(\tau) \coloneqq \sum_{n=1}^{N} w_n \bar{\boldsymbol{s}}_n(\tau) \odot (\boldsymbol{1} - \pi_n(\bar{\boldsymbol{s}}_n))$ Find $\boldsymbol{c}(\boldsymbol{s}(\tau-1), \boldsymbol{a}(\tau))$ 18 19 Save $(\boldsymbol{s}(\tau-1), \boldsymbol{a}(\tau), \boldsymbol{c}(s(\tau-1), \boldsymbol{a}(\tau)), \boldsymbol{s}(\tau))$ in \mathcal{E} 20 Uniformly sample B experiences from \mathcal{E} 21 Find $\nabla \mathcal{L}^{\text{Tar}}(\boldsymbol{\theta})$ for these samples, using (20) 22 Update $\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_{\tau} - \beta_{\tau} \nabla \mathcal{L}^{\mathrm{Tar}}(\boldsymbol{\theta})$ 23 If $mod(\tau, C) = 0$, then update $\theta^{Tar} = \theta_{\tau}$ 24 25 end



Fig. 6. Convergence of DQN to the target network.



 $_{547}$ the experience here is randomly drawn from past experiences $_{548}$ in \mathcal{E} . As a result, the loss function becomes

$$\mathcal{L}^{\mathrm{Tar}}(\boldsymbol{\theta}) = \mathbb{E} \left\| \boldsymbol{\delta}^{\mathrm{Tar}}(\boldsymbol{\theta}; \boldsymbol{E}) \right\|_{2}^{2}$$
(20)

⁵⁵⁰ where the expectation is taken with respect to the uniformly ⁵⁵¹ sampled experience E. In practice however, only a batch of B⁵⁵² samples is available and used to update θ_{τ} , so the expectation ⁵⁵³ will be replaced with the sample mean. Finally, following a ⁵⁵⁴ gradient descent step over the sampled experiences, we have

55
$$\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_{\tau} - \beta_{\tau} \nabla \mathcal{L}^{\mathrm{Tar}}(\boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{\tau}}.$$

5

560

⁵⁵⁶ Both the experience replay and the target network help sta-⁵⁵⁷ bilize the DQN updates. Incorporating these remedies, Alg. 1 ⁵⁵⁸ tabulates our deep RL based adaptive caching scheme for the ⁵⁵⁹ parent node.

VI. NUMERICAL TESTS

In this section, we present several numerical tests to assess he performance of our deep RL based caching schemes which is represented in Alg. 1.

The first experiment considers a simple setup, consisting of a parent node that directly serves user file requests. A total of F = 50 files were assumed in the cloud, each having the same size, while $M_0 = 5$ files can be stored in the cache, amounting to 10% of the total. The file popularities were drawn randomly

Fig. 7. Impact of C on DQN convergence in a static setting.

from [0, 1], invariant across all simulated time instants, but 569 unknown to the caching entity. A fully connected feed-forward 570 NN of 3 layers was implemented for DQN with each layer 571 comprising 50 hidden neurons. For simplicity, we assume that 572 the dataset \mathcal{E} can store a number \mathcal{R} of most recent experiences, 573 which is also known as the replay memory. It was set to $\mathcal{R} = 574$ 10 in our test, along with mini-batch size B = 1 to update the 575 target network every C = 10 iterations. In addition, hyperparameter values $\gamma = 0.8$, learning rate $\beta_{\tau} = 0.01$, as well as 577 $\epsilon_{\tau} = 0.4$ were used throughout our numerical tests. 578

Figure 6 shows the convergence of the DQN parameter θ_{τ} 579 to that of the target network θ^{Tar} . Since θ^{Tar} is updated with 580 θ_{τ} per *C* iteration, the error $\|\theta_{\tau} - \theta^{\text{Tar}}\|_2$ vanishes period-581 ically. While the error changes just in a few iterations after 582 the θ^{Tar} update, it is constant in several iterations before the 583 next θ^{Tar} update, suggesting that θ_{τ} is fixed across those iterations. This motivates investigating the impact of *C* on θ_{τ} 's 585 convergence. Figure 7 shows the convergence of θ_{τ} to θ^{Tar} 586 for C = 20, 5, 3, 2. Starting with C = 20, vector θ_{τ} is not 587 updated in most iterations between two consecutive updates of 588 θ^{Tar} . Clearly, the smaller *C* is, the faster the convergence for 589 θ_{τ} is achieved. Indeed, this is a direct result of having static 590



Fig. 8. Hyper deep-Q network for scalable caching.

⁵⁹¹ file popularity. Based on our numerical tests, having small *C* ⁵⁹² values is preferable in dynamic settings too.

For the second experiment, a similar setup having solely the parent node, plus $M_0 = 5$ and F = 50, was considered. Dynamic file popularities were generated with time evolutions modeled by Markov process described earlier in the paper, and again they are unknown to the caching node.

We first consider that the parent node is connected to 598 $_{599} N = 5$ leaf nodes, and every leaf node implements a local 600 caching policy π_n with capacity of storing $M_n = 5$ files. 601 Every leaf node receives user file requests within each fast $_{602}$ time slot, and each slow-timescale interval consists of T = 2slots. File popularity exhibits different Markovian dynam-603 604 ics locally at leaf nodes. Having no access to local policies $\{\pi_n\}$, the parent node not only should learn file populari-605 606 ties along with their temporal evolutions, but also learn the 607 caching policies of leaf nodes. To endow our approach with scalability to handle $F \gg 1$, we advocate the following 609 hyper Q-network implementation. Files are first split into Ksmaller groups of sizes F_1, \ldots, F_K with $F_k \ll F$. This ⁶¹¹ yields the representation $\mathbf{s}^{\top}(\tau) := [\mathbf{s}^{1\top}(\tau), \dots, \mathbf{s}^{K^{\top}}(\tau)],$ ₆₁₂ where $s^k \in \mathbb{R}^{F_k}$. By running K DQNs in parallel, every 613 DQN-k now outputs the associated predicted costs of input files through $\boldsymbol{o}^k(\tau) \in \mathbb{R}^{F^k}$. Concatenating all these outputs, 615 one obtains the predicted output cost vector of all files as 616 $\boldsymbol{o}^{\top}(\tau+1) \coloneqq [\boldsymbol{o}^{1\top}(\tau+1), \dots, \boldsymbol{o}^{K\top}(\tau+1)];$ see Section V ₆₁₇ for discussion on finding action a from vector o, and also 618 Fig. 8 for an illustration of our hyper Q-network.

The ensuing numerical tests consider F = 1,000 files with K = 25 and $\{F_k = 20\}$, where only $M_0 = 50$ can be stored. To further assess the performance, we adopt a *non-case causal* optimal policy as a benchmark, which unrealistically assumes knowledge of future requests and stores the most frequently requested files. In fact, this is the best policy that best policy that is schemes including, e.g., the LRU, LFU, and FIFO [45] are also simulated. A difference between LRU, LFU, FIFO, and whenever



Fig. 9. Instantaneous reduced cost for different policies.



Fig. 10. Instantaneous reduced cost for different policies.

a request is received, while our scheme refreshes the cache 629 only at the end of every time slot. 630

By drawing 100 samples² randomly from all 2,000 time ⁶³¹ intervals, the instantaneous reduced cost and the empirical ⁶³² cumulative distribution function (CDF) obtained over these ⁶³³ 100 random samples for different policies are plotted in Fig. 9 ⁶³⁴ and Fig. 10, respectively. These plots further verify how ⁶³⁵ the DRL policy performs relative to the alternatives, and in ⁶³⁶ particular very close to the optimal policy. ⁶³⁷

LRU, LFU, and FIFO make caching decisions based on 638 instantaneous observations, and can refresh the cache many 639 times within each slot. Yet, our proposed policy as well as 640 the optimal one here learns from all historical observations 641 to cache, and refreshes the cache only once per slot. Because 642 of this difference, the former policies outperform the latter at 643 the very beginning of Fig. 9, but they do not adapt to the 644 underlying popularity evolution and are outperformed by our 645 learning-based approach after a number of slots. The merit of 646 our approach is further illustrated by the CDF of the reduced 647 cost depicted in Fig. 10. 648

²During these samples, DRL policy is enforced to exploit its learned caching policy.



Fig. 11. Performance of all policies for N = 10 leaf nodes.



Fig. 12. Performance of all policies for N = 1,000 leaf nodes.

In the last test, we increased the number of leaf nodes from 10 in the previous experiment to N = 1,000. Figures 11 and 12 showcase that the DRL performance approaches that of the optimal policy as the number of nodes increases. This is likely because the more leaf nodes, the smoother the popularity fluctuations, and therefore the easier it becomes for DRL to learn the optimal policy.

VII. CONCLUSION

656

Caching highly popular contents across distributed caching 657 658 entities can substantially reduce the heavy burden of content 659 delivery in modern data networks. This paper considered a 660 network caching setup, comprising a parent node connected 661 to several leaf nodes to serve end user file requests. Since the 662 leaf nodes are closer to end users, they observe file requests in a fast timescale, and could thus serve requests, either through 664 their locally stored files or by fetching files from the parent 665 node. Observing aggregate requests from all leaf nodes, the 666 parent node makes caching decisions in a slower-time scale. 667 Such a two-timescale caching task was tackled in this paper 668 using a RL approach. An efficient caching policy leveraging 669 deep RL was introduced, and shown capable of learning-and-670 adapting to dynamic evolutions of file requests, and caching policies of leaf nodes. Simulated tests corroborated its impressive performance. This work also opens up several directions for future research, including multi-armed bandit online learning [46], and distributed deep RL using recurrent NNs [47] for future spatio-temporal file request prediction.

REFERENCES

- G. S. Paschos, G. Iosifidis, M. Tao, D. Towsley, and G. Caire, "The 677 role of caching in future communication systems and networks," *IEEE 678 J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1111–1125, Jun. 2018. 679
- [2] E. Bastug, M. Bennis, and M. Debbah, "Living on the edge: The role 680 of proactive caching in 5G wireless networks," *IEEE Commun. Mag.*, 681 vol. 52, no. 8, pp. 82–89, Aug. 2014.
- [3] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep Learning*.
 Cambridge, MA, USA: MIT Press, 2016.
- [4] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A 685 review and new perspectives," *IEEE Trans. Pattern Anal. Mach. Intell.*, 686 vol. 35, no. 8, pp. 1798–1828, Aug. 2013. 687
- [5] V. Mnih *et al.*, "Human-level control through deep reinforcement 688 learning," *Nature*, vol. 518, no. 7540, p. 529, Feb. 2015. 689
- [6] N. C. Luong *et al.*, "Applications of deep reinforcement learning in 690 communications and networking: A survey," *IEEE Commun. Surveys* 691 *Tuts.*, to be published.
- [7] J. Gao and R. Jamidar, "Machine learning applications for data center 693 optimization," Menlo Park, CA, USA, Google, White Paper, Oct. 2014. 694
- [8] O. Naparstek and K. Cohen, "Deep multi-user reinforcement learning for dynamic spectrum access in multichannel wireless networks," in *Proc.* 696 *Glob. Commun. Conf.*, Singapore, Dec. 2017, pp. 1–7. 697
- Y. Yu, T. Wang, and S. C. Liew, "Deep-reinforcement learning multiple access for heterogeneous wireless networks," in *Proc. Int. Conf. Commun.*, Kansas City, MO, USA, May 2018, pp. 1–7.
- Z. Wang, L. Li, Y. Xu, H. Tian, and S. Cui, "Handover control in wireless 701 systems via asynchronous multiuser deep reinforcement learning," *IEEE 702 Internet Things J.*, vol. 5, no. 6, pp. 4296–4307, Dec. 2018.
- Y. Sun, M. Peng, and S. Mao, "Deep reinforcement learning-based mode selection and resource management for green fog radio access networks," *TeEE Internet Things J.*, vol. 6, no. 2, pp. 1960–1971, Apr. 2019.
- [12] Y. Dong, M. Z. Hassan, J. Cheng, M. J. Hossain, and V. C. M. Leung, 707 "An edge computing empowered radio access network with UAVmounted FSO fronthaul and backhaul: Key challenges and approaches," 709 *IEEE Wireless Commun.*, vol. 25, no. 3, pp. 154–160, Jul. 2018. 710
- [13] C. Yu, J. Lan, Z. Guo, and Y. Hu, "DROM: Optimizing the routing 711 in software-defined networks with deep reinforcement learning," *IEEE 712 Access*, vol. 6, pp. 64533–64539, 2018.
- [14] Z. Guo, W. Chen, Y.-F. Liu, Y. Xu, and Z.-L. Zhang, "Joint 714 switch upgrade and controller deployment in hybrid software-defined 715 networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 5, pp. 1012–1028, 716 Mar. 2019. 717
- P. Blasco and D. Gündüz, "Learning-based optimization of cache content 718 in a small cell base station," in *Proc. IEEE Int. Conf. Commun.*, Sydney, 719 NSW, Australia, Jun. 2014, pp. 1897–1903.
- [16] L. Lei, L. You, G. Dai, T. X. Vu, D. Yuan, and S. Chatzinotas, "A deep 721 learning approach for optimizing content delivering in cache-enabled 722 HetNet," in *Proc. Int. Symp. Wireless Commun. Syst.*, Bologna, Italy, 723 Aug. 2017, pp. 449–453. 724
- [17] S. Traverso, M. Ahmed, M. Garetto, P. Giaccone, E. Leonardi, and 725 S. Niccolini, "Temporal locality in today's content caching: Why it matters and how to model it," ACM SIGCOMM Comput. Commun. Rev., 727 vol. 43, no. 5, pp. 5–12, Nov. 2013. 728
- M. Leconte, G. S. Paschos, L. Gkatzikis, M. Draief, S. Vassilaras, and 729
 S. Chouvardas, "Placing dynamic content in caches with small popula-730 tion," in *Proc. Int. Conf. Comput. Commun.*, San Francisco, CA, USA, 731
 Apr. 2016, pp. 1–9. 732
- [19] A. Sadeghi, F. Sheikholeslami, and G. B. Giannakis, "Optimal and 733 scalable caching for 5G using reinforcement learning of space-time 734 popularities," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, 735 pp. 180–190, Feb. 2018.
- [20] S. O. Somuyiwa, A. György, and D. Gündüz, "A reinforcement-learning 737 approach to proactive caching in wireless networks," *IEEE J. Sel. Areas 738 Commun.*, vol. 36, no. 6, pp. 1331–1344, Jun. 2018. 739
- [21] A. Sadeghi, F. Sheikholeslami, A. G. Marques, and G. B. Giannakis, 740 "Reinforcement learning for adaptive caching with dynamic storage 741 pricing," *IEEE J. Sel. Areas Commun.*, to be published. 742

676

743 [22] Y. He et al., "Deep-reinforcement-learning-based optimization for cache-

enabled opportunistic interference alignment wireless networks," IEEE 744 745

Trans. Veh. Technol., vol. 66, no. 11, pp. 10433-10445, Nov. 2017. C. Zhong, M. C. Gursoy, and S. Velipasalar, "A deep reinforcement 746 [23]

- learning-based framework for content caching," in Proc. Conf. Inf. Sci. 747 748 Syst., Princeton, NJ, USA, Mar. 2018, pp. 1-6.
- 749 [24] Y. He, N. Zhao, and H. Yin, "Integrated networking, caching, and com-750 puting for connected vehicles: A deep reinforcement learning approach," IEEE Trans. Veh. Technol., vol. 67, no. 1, pp. 44-55, Jan. 2018. 751
- 752 [25] Y. He, F. R. Yu, N. Zhao, V. C. M. Leung, and H. Yin, "Software-defined 753 networks with mobile edge computing and caching for smart cities: A big data deep reinforcement learning approach," IEEE Commun. Mag., 754
- vol. 55, no. 12, pp. 31-37, Dec. 2017. 755
- 756 [26] H. Zhu, Y. Cao, W. Wang, T. Jiang, and S. Jin, "Deep reinforcement learning for mobile edge caching: Review, new features, and open 757 758 issues," IEEE Netw., vol. 32, no. 6, pp. 50-57, Nov./Dec. 2018.
- M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," 759 [27] IEEE Trans. Inf. Theory, vol. 60, no. 5, pp. 2856-2867, May 2014. 760
- 761 [28] S. C. Borst, V. Gupta, and A. Walid, "Distributed caching algorithms for content distribution networks," in Proc. Int. Conf. Comput. Commun., 762 763 San Diego, CA, USA, Mar. 2010, pp. 1478-1486.
- 764 [29] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," IEEE/ACM Trans. Netw., vol. 24, no. 2, pp. 836-845, 765 Apr. 2016. 766
- 767 [30] J. Dai, Z. Hu, B. Li, J. Liu, and B. Li, "Collaborative hierarchical 768 caching with dynamic request routing for massive content distribution," in Proc. Int. Conf. Comput. Commun., Orlando, FL, USA, Mar. 2012, 769 pp. 2444-2452. 770
- 771 [31] W. Wang, D. Niyato, P. Wang and A. Leshem, "Decentralized caching for content delivery based on blockchain: A game theoretic perspective,' 772 773 in Proc. IEEE Intl. Conf. Commun., Kansas City, MO, USA, May 2018, pp. 1-6. 774
- 775 [32] E. Nygren, R. K. Sitaraman, and J. Sun, "The Akamai network: A platform for high-performance Internet applications," ACM SIGOPS Oper. 776 Syst. Rev., vol. 44, no. 3, pp. 2-19, 2010. 777
- 778 [33] M. Dehghan et al., "On the complexity of optimal request routing and content caching in heterogeneous cache networks," IEEE/ACM Trans. 779 Netw., vol. 25, no. 3, pp. 1635-1648, Jun. 2017. 780
- 781 [34] S. Shukla, O. Bhardwaj, A. A. Abouzeid, T. Salonidis, and T. He, "Proactive retention-aware caching with multi-path routing for wire-782 less edge networks," IEEE J. Sel. Areas Commun., vol. 36, no. 6, 783 pp. 1286-1299, Jun. 2018. 784
- L. Tong, Y. Li, and W. Gao, "A hierarchical edge cloud architecture for 785 [35] 786 mobile computing," in Proc. IEEE Int. Conf. Comput. Commun., 2016, pp. 1–9. 787
- 788 [36] E. Dahlman, S. Parkvall, and J. Skold, 4G: LTE/LTE-Advanced for Mobile Broadband. Amsterdam, The Netherlands: Academic, 2013. 789
- C. H. Papadimitriou and J. N. Tsitsiklis, "The complexity of Markov 790 [37] decision processes," Math. Oper. Res., vol. 12, no. 3, pp. 441-450, 1987. 791
- R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction. 792 [38] Cambridge, MA, USA: MIT Press, 2018. 793
- 794 [39] C. J. C. H. Watkins and P. Dayan, "Q-learning," Mach. Learn., vol. 8, nos. 3-4, pp. 279-292, May 1992. 795
- R. Fagin, "Asymptotic miss ratios over independent references," J. 796 [40] Comput. Syst. Sci., vol. 14, no. 2, pp. 222-250, 1977. 797
- 798 [41] P. Auer, "Using confidence bounds for exploitation-exploration tradeoffs," J. Mach. Learn. Res., vol. 3, pp. 397-422, Nov. 2002. 799
- 800 [42] S. J. Russell and P. Norvig, Artificial Intelligence: A Modern Approach. Upper Saddle River, NJ, USA: Prentice-Hall, 2016. 801
- 802 [43] G. Wang, G. B. Giannakis, and J. Chen, "Learning ReLU networks on 803 linearly separable data: Algorithm, optimality, and generalization," IEEE Trans. Signal Process., vol. 67, no. 9, pp. 2357-2370, May 2019. 804
- D. Silver *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, p. 484, Jan. 2016. 805 [44] 806
- 807 [45] A. Dan and D. F. Towsley, "An approximate analysis of the LRU and
- FIFO buffer replacement schemes," ACM SIGMETRICS Perform. Eval. 808 Rev., vol. 18, no. 1, pp. 143-152, 1990. 809
- B. Li, T. Chen, and G. B. Giannakis, "Bandit online learning with 810 [46] unknown delays," in Proc. Int. Conf. Artif. Intell. Stat., Naha, Japan, 811 Apr. 2019, pp. 1-10. 812
- 813 [47] A. Narayanan, S. Verma, E. Ramadan, P. Babaie, and Z.-L. Zhang, "DeepCache: A deep learning based framework for content caching," 814
- in Proc. ACM Workshop Netw. Meets AI & ML, Budapest, Hungary, 815 Aug. 2018, pp. 48-53. 816



Alireza Sadeghi (S'16) received the B.Sc. 817 degree (Hons.) in electrical engineering from the 818 Iran University of Science and Technology, Tehran, 819 Iran, in 2012, and the M.Sc. degree in electrical 820 engineering from the University of Tehran in 821 2015. He is currently pursuing the Ph.D. degree 822 with the Department of Electrical and Computer 823 Engineering, University of Minnesota, Minneapolis, 824 MN, USA. In 2015, he was a Visiting Scholar with 825 the University of Padua, Padua, Italy. His research 826 interests include machine learning, optimization, 827

and signal processing with applications to networking. He was a recipient 828 of the ADC Fellowship awarded by the Digital Technology Center of the 829 University of Minnesota, and the Student Travel Awards from the IEEE 830 Communications Society and the National Science Foundation. 831



Gang Wang (M'18) received the B.Eng. degree 832 in electrical engineering and automation from the 833 Beijing Institute of Technology, Beijing, China, in 834 2011, and the Ph.D. degree in electrical and com- 835 puter engineering from the University of Minnesota, 836 Minneapolis, MN, USA, in 2018. 837

He is currently a Post-Doctoral Associate with the 838 Department of Electrical and Computer Engineering, 839 University of Minnesota. His research interests 840 focus on the areas of statistical signal processing, 841 optimization, and deep learning with applications to 842

data science and smart grids. He was a recipient of the National Scholarship 843 in 2013, the Guo Rui Scholarship in 2015, and the Innovation Scholarship 844 (First Place in 2017), all from China, as well as the Best Conference Papers 845 at the 2017 European Signal Processing Conference and 2019 IEEE Power 846 and Energy Society General Meeting. 847



Georgios B. Giannakis (F'97) received the Diploma 848 degree in electrical engineering from the National 849 Technical University of Athens, Greece, in 1981, 850 and the first M.Sc. degree in electrical engineer- 851 ing, the second M.Sc. degree in mathematics, and 852 the Ph.D. degree in electrical engineering from the 853 University of Southern California in 1983, 1986, and 854 1986, respectively. 855

From 1982 to 1986, he was with the University 856 of Southern California. He was a Faculty Member 857 with the University of Virginia from 1987 to 1998. 858

Since 1999, he has been a Professor with the University of Minnesota, where he holds an ADC Endowed Chair, a University of Minnesota McKnight 860 Presidential Chair in ECE, and serves as the Director of the Digital Technology 861 Center. His general interests span the areas of statistical learning, commu-862 nications, and networking-subjects on which he has published over 450 863 journal papers, 750 conference papers, 25 book chapters, 2 edited books, and 864 2 research monographs (H-index 142). He has co-invented 32 patents. Current 865 research focuses on data science and network science with applications to the 866 Internet of Things, social, brain, and power networks with renewables. He 867 was a (co-)recipient of nine best journal paper awards from the IEEE Signal 868 Processing (SP) and Communications Societies, including the G. Marconi 869 Prize Paper Award in Wireless Communications. He also received Technical 870 Achievement Awards from the SP Society in 2000, the EURASIP in 2005, the 871 Young Faculty Teaching Award, the G. W. Taylor Award for Distinguished 872 Research from the University of Minnesota, and the IEEE Fourier Technical 873 Field Award (inaugural recipient in 2015). He has served the IEEE in a num- 874 ber of posts, including that of a Distinguished Lecturer for the IEEE-SPS. He 875 is a fellow of EURASIP. 876