

Optimal dynamic proactive caching via reinforcement learning

Alireza Sadeghi, Fatemeh Sheikholeslami, and Georgios Giannakis

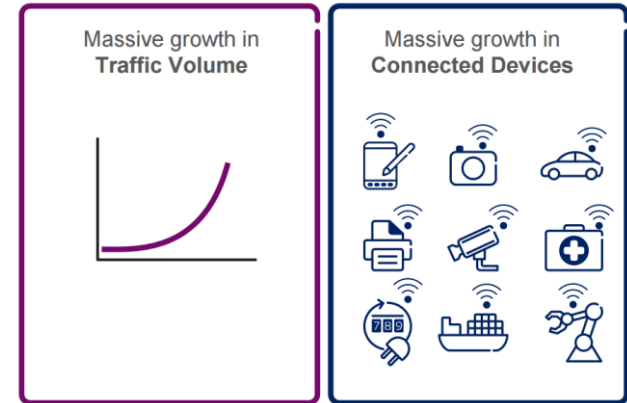
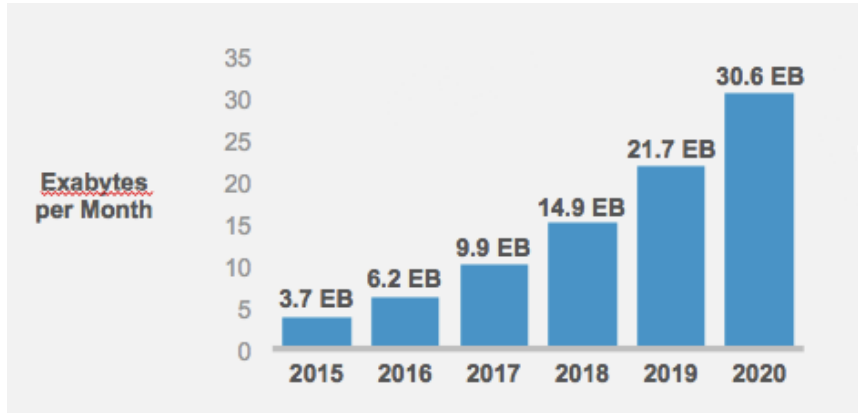
April 14, 2017

Acknowledgment: ARO W911NF-15-1-0492 and NSF-EARS 1343248.

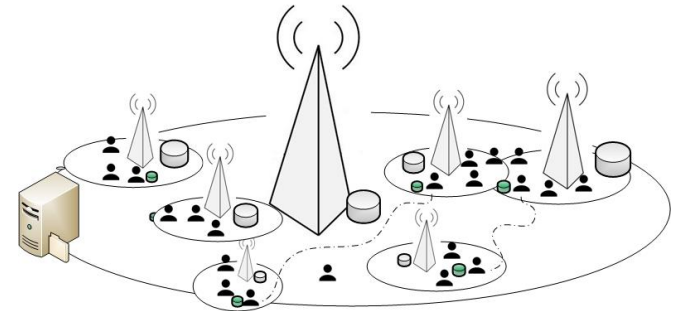


Evolution of wireless networks

- 8-fold growth of global mobile data traffic between 2015 and 2020



- Heterogeneous network architecture (HetNet)
 - Access technologies in heterogeneous subnets
- 60% of data is reusable, a.k.a. **contents**
 - Utilization of storage units at small base stations (SBs)
 - Challenge: *what* and *when* to cache?



Caching in wireless networks

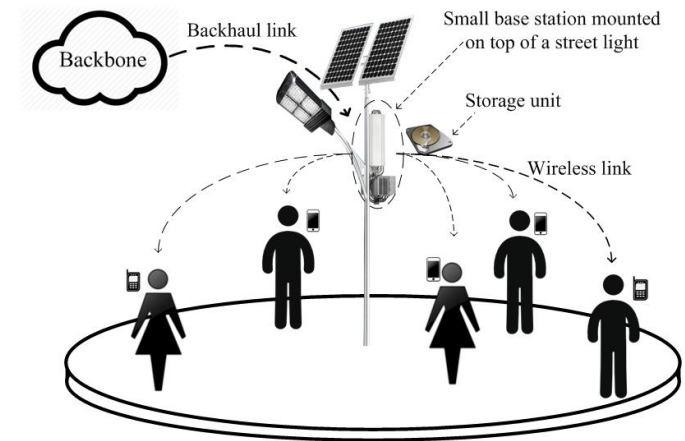
❑ Memory-enabled SBs

- Cache during off-peak hours
- Reduce load on backhauls during peak traffic periods
- Reduce cost for providing service with high QoS

❑ Generally unknown content popularity profiles

❑ Prior art

- Multi-armed bandit (MAB) formulation [D. Belasco et al'14]
- Distributed and convexified MAB [A. Sengupta et al'14] \Rightarrow Unknown static popularity profile
- Dynamic user demand [Kim et al'17] \Rightarrow Unknown dynamic popularity profile

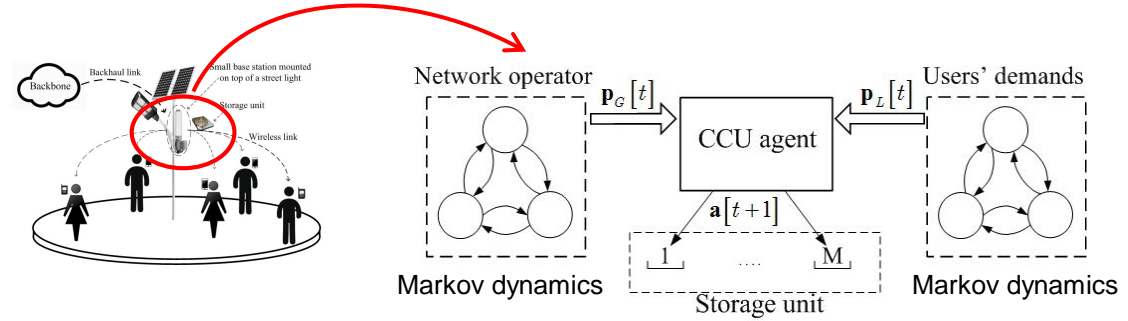


Proposed approach

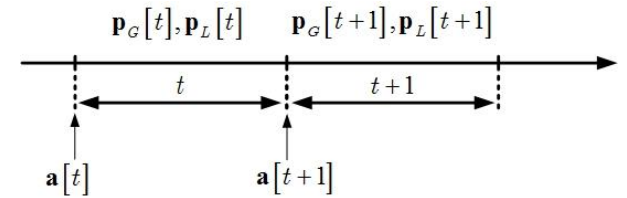
- ✓ Caching via reinforcement learning over files with spatio-temporally dynamic popularities

System model

Discrete-time network



- SB with caching control unit (CCU)
- Total number of F contents w/ unit size in backbone
- Storage capacity of M files in SBs



- Local popularity profile
$$\left[\mathbf{p}_L[t] \right]_f = \frac{\# \text{ of local requests for file } f \text{ at time interval } t}{\text{Total } \# \text{ of local requests at time interval } t};$$
- Global popularity profiles $\mathbf{p}_G[t]$
- Action vector $\mathbf{a}[t] \in \{0, 1\}^F$: SB caches f^{th} content if $\mathbf{a}_f[t] = 1$
- State vector $\mathbf{s}[t] := [\mathbf{p}_G^T[t], \mathbf{p}_L^T[t], \mathbf{a}^T[t]]^\top$
- Policy $\pi(\cdot)$ is a mapping from state space to action space $\Rightarrow \mathbf{a}[t+1] = \pi(\mathbf{s}[t])$

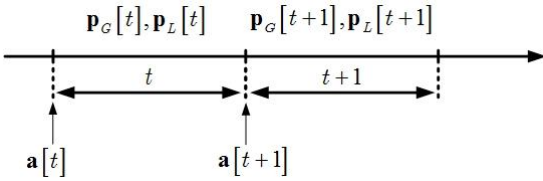
Goal: Given $\{\mathbf{s}[\tau]\}_{\tau=0}^{\tau=t}$ and observed costs, optimize policy $\pi(\cdot)$

Problem formulation

□ Recall

□ Costs

- Refreshing the cached contents
- Fetching requested non-cached files
- Tracking global popularities


$$\begin{aligned}h(\mathbf{a}[t], \mathbf{a}[t-1]) &:= \lambda_1 \mathbf{a}^\top[t] (\mathbf{1} - \mathbf{a}[t-1]) \\g(\mathbf{s}[t]) &:= \lambda_2 (\mathbf{1} - \mathbf{a}[t])^\top \mathbf{p}_L[t] \\f(\mathbf{s}[t]) &:= \lambda_3 (\mathbf{1} - \mathbf{a}[t])^\top \mathbf{p}_G[t]\end{aligned}$$

$$\Rightarrow C(\mathbf{s}[t-1], \mathbf{a}[t] | \mathbf{p}_G[t], \mathbf{p}_L[t]) := h(\mathbf{a}[t], \mathbf{a}[t-1]) + g(\mathbf{s}[t]) + f(\mathbf{s}[t])$$

□ Expected discounted cost

$$V_\pi(\mathbf{s}[\tau]) := \lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=\tau}^T \gamma^{t-\tau} C(\mathbf{s}[t], \pi(\mathbf{s}[t])) \right]$$

□ Goal: Find the optimal policy

$$\pi^* = \arg \min_{\pi \in \Pi} V_\pi(\mathbf{s}_0)$$

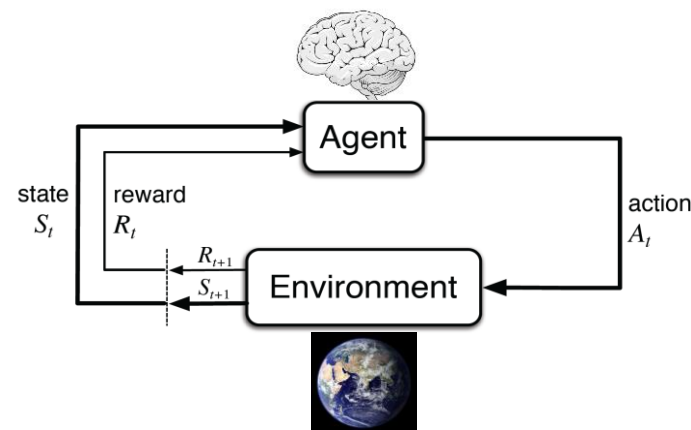
□ Viable approaches

- ✓ Adaptive dynamic programming
- ✓ Q-learning
- ✓ SARSA

Reinforcement learning (RL)

Agent-environment interactions

- **State:** mathematical representation of environment
- **Action:** decision made by the agent
- **Reward:** scalar feedback, how well agent is doing



State value function (under policy π)

- Immediate reward + discounted future rewards

Objective

- ✓ Find optimal policy $\pi(\cdot)$ such that $V_\pi(s)$ is maximized for all states

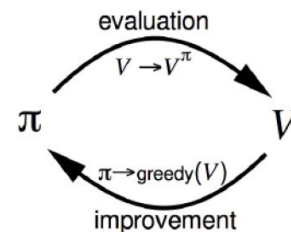
□ Bellman equation
$$V_\pi(s) = \mathbb{E}[C(s, \pi(s))] + \gamma \sum_{s' \in \mathcal{S}} \mathbf{T}(s'; s, \pi(s)) V_\pi(s'), \forall s, s'.$$

Q-learning

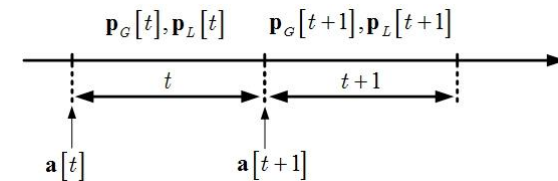
- State-action value function
$$Q^\pi(s, a) := \mathbb{E}[C(s, a)] + \gamma \sum_{s' \in \mathcal{S}} \mathbf{T}(s'; s, a) V^\pi(s')$$

Optimality

$$\pi^*(s) = \arg \min_a Q^*(s, a), \quad \forall s \in \mathcal{S},$$



Q-learning for proactive caching



Algorithm 1: Proactive-caching via Q-learning at CCU

1 **Initialize** $\mathbf{s}[0]$ randomly and $\hat{Q}(\mathbf{s}, \mathbf{a}) = 0 \forall \mathbf{s}, \mathbf{a}$

2 **for** $t=1, 2, \dots$ **do**

3 Take action $\mathbf{a}[t]$ chosen probabilistically by

$$\mathbf{a}[t] = \begin{cases} \arg \min_{\mathbf{a}} \hat{Q}(\mathbf{s}[t-1], \mathbf{a}) & \text{w.p. } 1 - \epsilon_t \\ \text{random } \mathbf{a} \in \mathcal{A} & \text{w.p. } \epsilon_t \end{cases}$$

} Exploration-exploitation tradeoff

4 $\mathbf{p}_L[t]$ and $\mathbf{p}_G[t]$ are revealed

5 Set $\mathbf{s}[t] = [\mathbf{p}_L[t], \mathbf{p}_G[t], \mathbf{a}[t]]$

6 Incur cost $C(\mathbf{s}[t-1], \mathbf{a}[t] | \mathbf{p}_G[t], \mathbf{p}_L[t])$

7 Update

} Observations are revealed

$$\begin{aligned} \hat{Q}_t(\mathbf{s}[t-1], \mathbf{a}[t]) &\leftarrow (1 - \beta_t) \hat{Q}_{t-1}(\mathbf{s}[t-1], \mathbf{a}[t]) + \beta_t \\ &\times [C(\mathbf{s}[t-1], \mathbf{a}[t] | \mathbf{p}_G[t], \mathbf{p}_L[t]) + \gamma \min_{\mathbf{a}} \hat{Q}_{t-1}(\mathbf{s}[t], \mathbf{a})] \end{aligned}$$

} Stochastic update of Q-table

Convergence

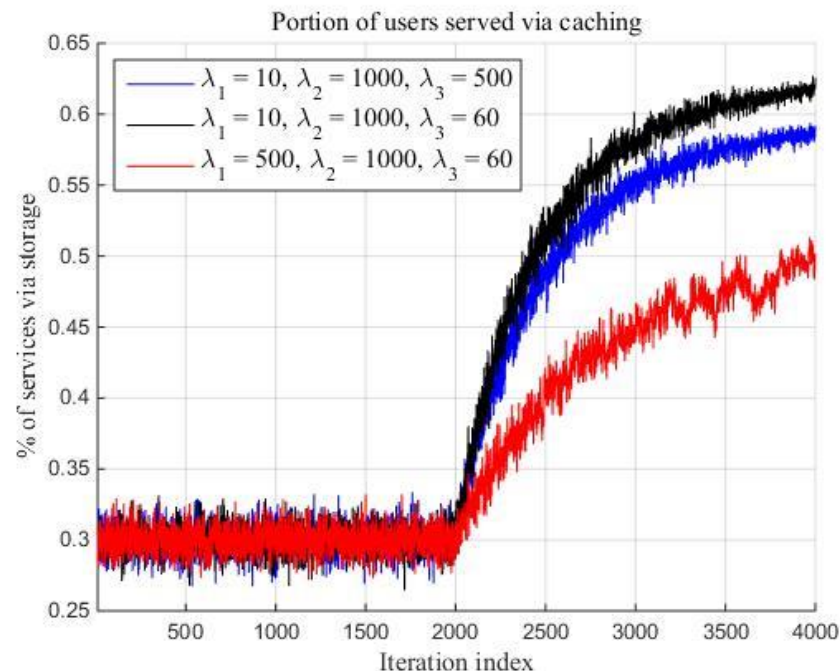
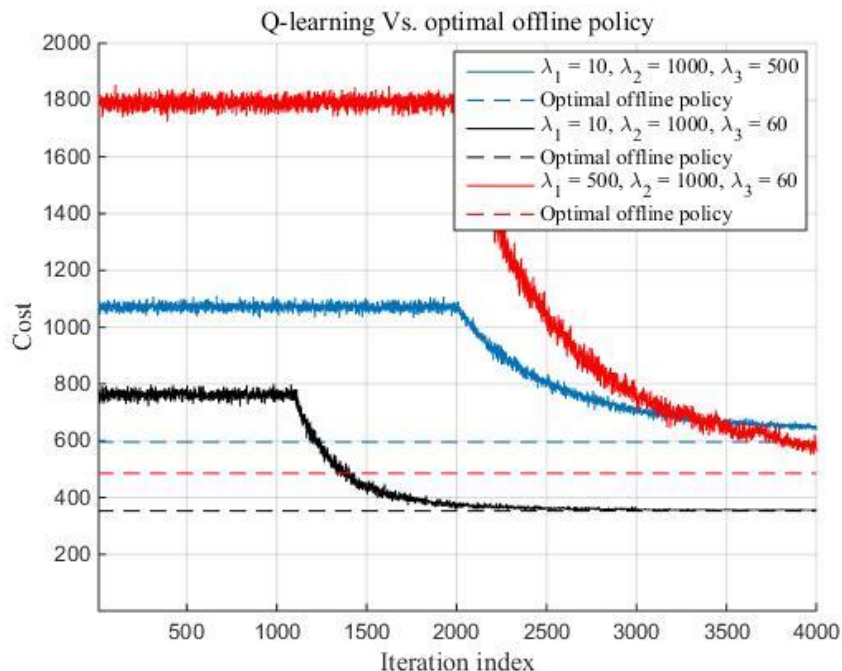
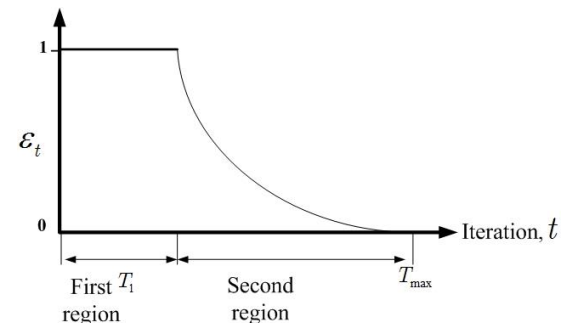
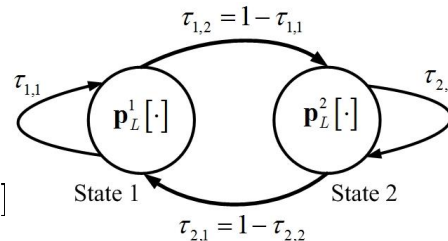
- ✓ If
1. All state-action pairs are continuously visited,
 2. Step-size β_t satisfies $\sum_{t=1}^{\infty} \beta_t = \infty$ and $\sum_{t=1}^{\infty} \beta_t^2 < \infty$

then policy will converge to the optimal policy, i.e., $\pi \rightarrow \pi^*$ w.p. 1.

Simulation tests

□ Consider a total of $F=10$, and $M=3$

- Two-state Markov chain for modeling $\mathbf{p}_L[t]$
- State transition probabilities $\tau := \begin{bmatrix} 0.35 & 0.65 \\ 0.75 & 0.25 \end{bmatrix}$
- Similarly for $\mathbf{p}_G[t]$ with $\tau' := \begin{bmatrix} 0.6 & 0.4 \\ 0.45 & 0.55 \end{bmatrix}$
- $\beta_t = 0.3$ and $\gamma = 0.6$



Thank you!