

On-demand caches and content-oriented networks

Stochastic content-service systems

Nicaise E. CHOUNGMO FOFACK

joint work with Philippe Nain, Sara Alouf, Don Towsley, Giovanni Neglia

Orange Labs Networks
Issy-les-Moulineaux, France
nicaise.choungmofofack@orange.com

Orange Labs



Outline

- 1 Introduction
- 2 Single cache analysis
- 3 Cache network analysis
- 4 Conclusion

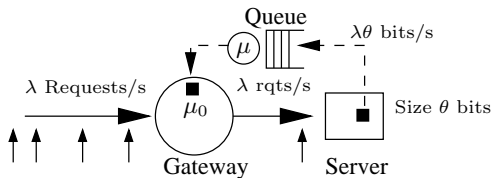
Outline

- 1 Introduction
- 2 Single cache analysis
- 3 Cache network analysis
- 4 Conclusion

Context: Data communication networks

Traditional Client—Server architectures

- 1 Host-to-Host communication model
- 2 Data flows
 - Kelly or Jackson's Queueing models

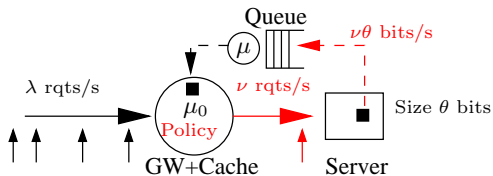


(a) Client-Server model

Context: Content-oriented networks

Client—Cache network—Server architectures

- 1 Host-to-Content communication model
- 2 Data acceleration, Server load reduction, Popular contents, Self-adaptation
 - Caching (Storage & Policy) + Queueing models
 - “On-demand” policies: Least Recently Used (LRU), FIFO, Random, Timeout

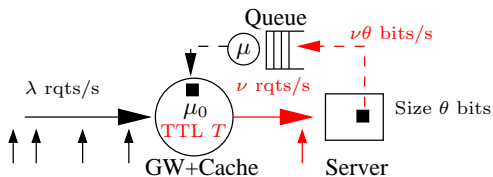


(b) Client-(Cache)-Server model, $\nu \leq \lambda$

Context: Model of Content-oriented networks

Client—Time-To-Live (TTL)-based Cache network—Server models

- (TTL)-based Caching + Queueing models
 - Simple, Tractable, and Extensible



(c) Client-(TTL Cache)-Server model, $\nu \leq \lambda$

Link with existing cache models

“Space” representations of LRU, FIFO, Random policies

- Describe the position of an item in the cache
- Exact analysis via Markov chains (King 1972, Gelenbe 1973)
- Approximations (Dan and Towsley 1990, Flajolet et al. 1992, Jelenkovic 1999)

“Time” representations of cache replacement policies

- Describe the remaining time of an item in the cache
- Approximate results of LRU, FIFO, Random policies via Time-To-Live (TTL)-based models (Che et al. 2002, Fricker et al. 2012, Martina et al. 2013)
- Exact results on TTL-based caches (Hou et al. 2004, Fofack et al. 2012, 2014, Fofack and Alouf 2013, Berger et al. 2014)

Link with existing cache models

“Space” representations of LRU, FIFO, Random policies

- Describe the position of an item in the cache
- Exact analysis via Markov chains (King 1972, Gelenbe 1973)
- Approximations (Dan and Towsley 1990, Flajolet et al. 1992, Jelenkovic 1999)

“Time” representations of cache replacement policies

- Describe the remaining time of an item in the cache
- Approximate results of LRU, FIFO, Random policies via Time-To-Live (TTL)-based models (Che et al. 2002, Fricker et al. 2012, Martina et al. 2013)
- Exact results on TTL-based caches (Hou et al. 2004, Fofack et al. 2012, 2014, Fofack and Alouf 2013, Berger et al. 2014)

Link with existing cache models

“Space” representations of LRU, FIFO, Random policies

- Describe the position of an item in the cache
- Exact analysis via Markov chains (King 1972, Gelenbe 1973)
- Approximations (Dan and Towsley 1990, Flajolet et al. 1992, Jelenkovic 1999)

“Time” representations of cache replacement policies

- Describe the remaining time of an item in the cache
- Approximate results of LRU, FIFO, Random policies via Time-To-Live (TTL)-based models (Che et al. 2002, Fricker et al. 2012, Martina et al. 2013)
- Exact results on TTL-based caches (Hou et al. 2004, Fofack et al. 2012, 2014, Fofack and Alouf 2013, Berger et al. 2014)

TTL approach: Towards a unifying cache model?

Simple and accurate models for existing cache networks

- “*Che approximation*” of LRU caches \rightarrow *Deterministic (Det.)* \mathcal{R} -TTL model
- FIFO and Random caches \rightarrow *Det.* and *Exponential* Σ -TTL models
- Amazon ElastiCache, Squid web caches \rightarrow *Det.* $\min(\Sigma, \mathcal{R})$ -TTL models

TTL-based caches within real systems

- (Modern) Domain Name System caches run *Det.* Σ -TTL policy.
- (OpenFlow) Software-defined switches run *Det.* $\min(\Sigma, \mathcal{R})$ -TTL policy.

Why TTL models are further interesting?

New caching replacement policies: *TTL as control parameter* for

- User QoE metrics: Delay, Downloading time
- Server load, Network QoS, Power save, etc.
- Content-service differentiation (Premium offer, Real-Time Apps, etc.)
- Elastic storage provisioning and management (Cloud, Data center)

Why TTL models are further interesting?

New caching replacement policies: *TTL as control parameter* for

- User QoE metrics: Delay, Downloading time
- Server load, Network QoS, Power save, etc.
- Content-service differentiation (Premium offer, Real-Time Apps, etc.)
- Elastic storage provisioning and management (Cloud, Data center)

Not covered here: Questions are welcome!

- Optimization and control
- Behaviour of TTL-based cache and $G/M/1$ queue in tandem

Challenges of the performance evaluation

In this talk (PhD thesis)

- 1 Workload model
 - Independent Reference Model (IRM) or Poisson assumption
 - (Markov-) Renewal, (non-) Stationary request models
- 2 Cache model and performance
 - per-content/demand metrics of interest (hit & occupancy probabilities)
 - global cache performance (hit ratio)
 - Filtering (cache misses)
- 3 Extension to cache networks
 - **Filtering (cache misses)**
 - Splitting and Aggregating (cache routing)

Assumptions (very easy to relax!)

- Downloading delays \ll Request time scales (or Infinite bandwidth)
- Infinite cache capacity (only TTL causes a file eviction)

Outline

- 1 Introduction
- 2 Single cache analysis
 - TTL-based concept
 - Description of basic TTL-based policies
 - Single TTL-based cache and single file
- 3 Cache network analysis
- 4 Conclusion

A nice property of TTL-based models

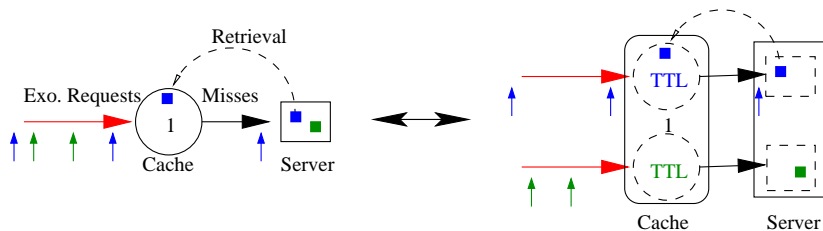


Figure : From capacity-driven to TTL-based caches

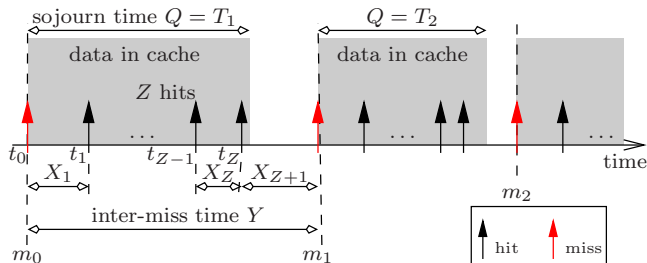
Infinite capacity \implies Decoupling \implies Focus on SINGLE content item¹

¹Choungmo Fofack (Nicaise). – On models for performance analysis of a core cache network and power save of a wireless access network. –Ph.D. thesis, Univ. of Nice Sophia Antipolis, <http://tel.archives-ouvertes.fr/tel-00968894>, Feb. 2014.

Σ -TTL policy

Algorithm

- On cache miss: add content, assign timer T to it
- On cache hit: use remaining value of T
- Only when** timer T expires: remove content

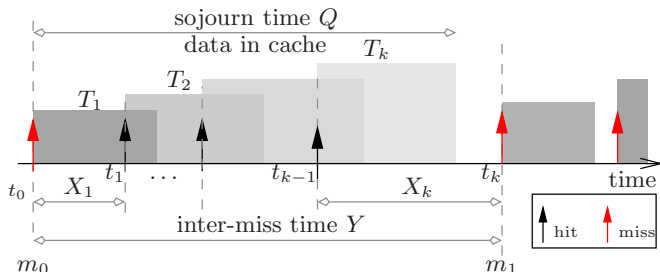


(a) Σ -TTL policy, $r = 0$

\mathcal{R} -TTL policy

Algorithm

- On cache miss: add content, assign timer T to it
- On cache hit: re-initialize timer T
- Only when** timer T expires: remove content



(b) \mathcal{R} -TTL policy, $r = 1$

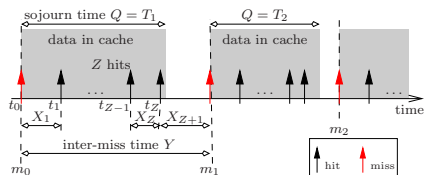
min(Σ, \mathcal{R})-TTL policy

Algorithm: a straightforward generalization

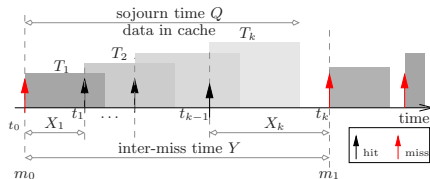
- On cache miss: apply \mathcal{R} -TTL policy with probability r
- Otherwise: apply Σ -TTL policy with probability $1 - r$

where $r = \mathbb{P}^0(Q^{(1)} < Q^{(0)})$ and

$Q^{(1)}$ and $Q^{(0)}$ are sojourn times in \mathcal{R} -TTL and Σ -TTL models respectively.



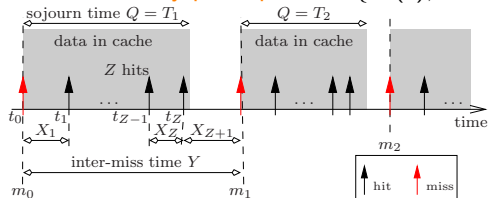
(c) Σ -TTL with prob. $1 - r$



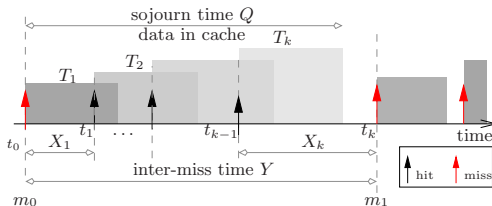
(d) \mathcal{R} -TTL with prob r

Analysis of $G/G/\Sigma$ and $G/G/\mathcal{R}$ -TTL cache models

- File requests form a **stationary point process**, $\{\mathcal{N}(t), t \geq 0\}$



(e) Case $r = 0$, $G/G/\Sigma$ -TTL



(f) Case $r = 1$, $G/G/\mathcal{R}$ -TTL

TTL-based caches under Poisson requests

- $F(t) = \mathbb{P}(X < t) = 1 - e^{-\lambda t}$, $F^*(s) = \mathbb{E}[e^{-sX}] = \frac{\lambda}{\lambda + s}$
- $M(t) = \mathbb{E}[\mathcal{N}(t)] = \lambda t$, $T(t) = \mathbb{P}(T < t)$, $T^*(s) = \mathbb{E}[e^{-sT}]$
- e.g.: Exp. $T^*(s) = \frac{\mu}{\mu + s}$, $\mu = 1/\mathbb{E}[T]$; Det. $T^*(s) = e^{-sD}$, $T = D$;

TTL-based caches under Poisson requests

- $F(t) = \mathbb{P}(X < t) = 1 - e^{-\lambda t}$, $F^*(s) = \mathbb{E}[e^{-sX}] = \frac{\lambda}{\lambda + s}$
- $M(t) = \mathbb{E}[\mathcal{N}(t)] = \lambda t$, $T(t) = \mathbb{P}(T < t)$, $T^*(s) = \mathbb{E}[e^{-sT}]$
- e.g.: Exp. $T^*(s) = \frac{\mu}{\mu + s}$, $\mu = 1/\mathbb{E}[T]$; Det. $T^*(s) = e^{-sD}$, $T = D$;

Proposition (Exact Performance and Miss stream, Fofack et al. 2012, 2013)

Hit and Occupancy probability (P.A.S.T.A): $H_P = O_P$

$$H_P^{(0)} = 1 - \frac{1}{1 + \lambda \mathbb{E}[T]}; \quad H_P^{(1)} = 1 - T^*(\lambda)$$

Miss stream is a **Renewal process**, and $G^{(r)*}(s) = \mathbb{E}[e^{-sY}]$ the LST of the CDF $G^{(r)}(t) = \mathbb{P}(Y < t)$ is given by

$$G^{(0)*}(s) = 1 - (1 - F^*(s)) \times \frac{\lambda}{s}(1 - T^*(s)), \quad \text{if } r = 0 \quad (1)$$

$$G^{(1)*}(s) = F^*(s)(1 - T^*(s)), \quad \text{if } r = 1 \quad (2)$$

$$G^{(r)*}(s) = (1 - r)G^{(0)*}(s) + rG^{(1)*}(s), \quad \text{if } r = \mathbb{P}^0(Q^{(1)} < Q^{(0)}) \quad (3)$$

General performance metrics of basic TTL-based caches

Theorem (Fofack et al. 2014, *PhD thesis*)

Under general stationary assumption, \mathbb{P}^0 the Palm probability, \mathbb{E}^0 the expectation w.r.t. \mathbb{P}^0 ,

$$H_P = r \mathbb{E}^0 [F(T_1)] + (1 - r) \left(1 - (1 + \mathbb{E}^0 [\mathcal{N}(T_1)])^{-1} \right) \quad (4)$$

$$O_P = (\mathbb{E}^0 [X_1]^{-1} \times (1 - H_P)) \times \mathbb{E}^0 [Q], \quad \text{"Little's Law for caches"} \quad (5)$$

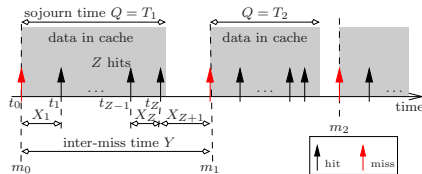
where $\mathbb{E}^0 [F(T_1)] = \mathbb{P}^0 (X_1 < T_1)$, and Q is the sojourn time:

$$Q = \begin{cases} Q^{(0)} = T_1, & \text{if } r = 0 \\ Q^{(1)} = T_1 \times \mathbf{1}\{X_1 > T_1\} + (X_1 + \tilde{Q}) \times \mathbf{1}\{X_1 < T_1\}, & \text{if } r = 1 \\ Q^{(r)} = \min(Q^{(0)}, Q^{(1)}), & \text{for } \min(\Sigma, \mathcal{R})\text{-TTL model if } r = \mathbb{P}^0(Q^{(1)} < Q^{(0)}) \end{cases} \quad (6)$$

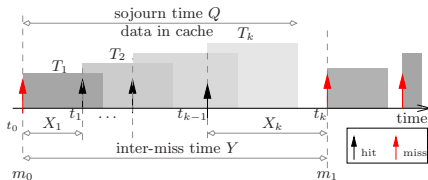
Miss process under general stationary request correlations

- Only an approximation of the CDF of the first inter-miss time.

$$Y_1 = Q + \tilde{X} \approx \mathbb{E}^0[Q] + \tilde{X}$$



(g) Σ -TTL, \tilde{X} residual request time



(h) \mathcal{R} -TTL, $\tilde{X} = X - T | X > T$

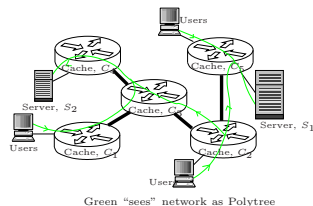
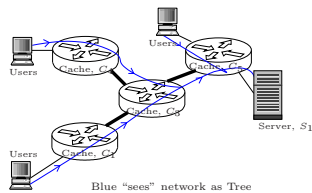
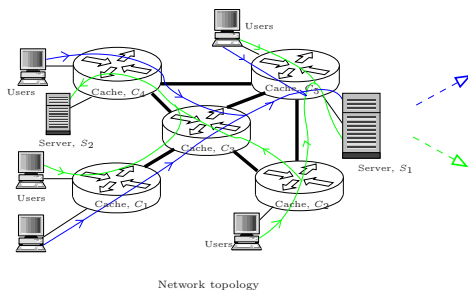
Outline

- 1 Introduction
- 2 Single cache analysis
- 3 Cache network analysis**
- 4 Conclusion

Exact analysis of “Pure” TTL-based cache networks

Negligible content/chunk size w.r.t cache size (DNS cache hierarchy)

- Assumption: **Markov-Renewal traffic model**
- Filtering, Aggregating & Splitting produce **Markov Renewal Processes**



Per-content TTL-based cache network analysis

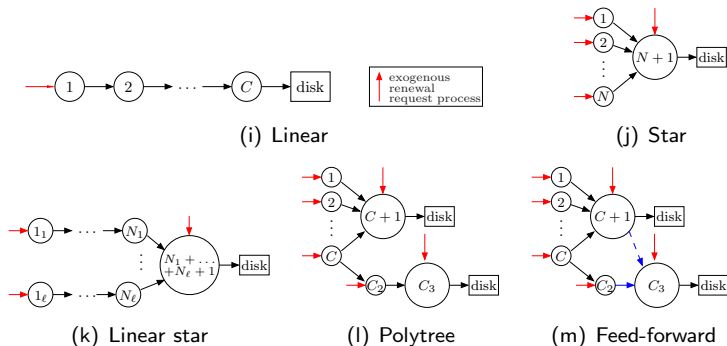
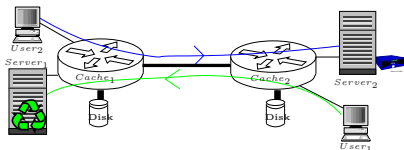


Figure : Per-content directed acyclic graph routing topologies

Sequential methodology

- Apply sequentially the single TTL-based cache analysis

“Constrained” TTL-based cache networks



Constraints on cache capacities (Video-on demand systems)

- Issue: **Dependency among cache states and TTLs**
- Why?: Saturation of capacity constraints $\sum_{i=1}^{K_n} O_{P,i,n} = C_n, \forall n$
- Example: Heterogeneous networks of LRU, FIFO, RND caches

Recursive methodology, Fofack et al. 2014 (Valuetools)

- Accurate and Polynomial time algorithm (with quadratic speed of convergence)

(Possible) Applications of TTL-based models

On real caching systems

1 Recently...

- Dynamic Page Caching mechanism (Akamai)
- WebRTC protocol (Google)

2 Past few years...

- New concepts: Information-Centric Networking (CCN, NDN architectures)
- (Mature) technos: push-based (CDNs, WWW) and pull-based (DNS, P2P)

In other systems

- Green networks and Smart grids (Idle mode, Power saving protocols)
- Economics (Product warranty), Physics (Geiger-Müller counters)

Outline

- 1 Introduction
- 2 Single cache analysis
- 3 Cache network analysis
- 4 Conclusion**

Take home message

1 Cache networks

- Think “TTL model” for cache analysis
- Proved results on TTL caches are valid for most caches!
e.g. If $F(t)$ is concave then Deterministic TTL performs the best.
- Deploy your “TTL model” using LRU, RND, FIFO or Pra-TTL

2 New research opportunities

- Cache network optimization and control (Current work)
- Spatio-temporal diversity in content access on mobile networks?

Take home message

1 Cache networks

- Think “TTL model” for cache analysis
- Proved results on TTL caches are valid for most caches!
e.g. If $F(t)$ is concave then Deterministic TTL performs the best.
- Deploy your “TTL model” using LRU, RND, FIFO or Pra-TTL

2 New research opportunities

- Cache network optimization and control (Current work)
- Spatio-temporal diversity in content access on mobile networks?

Thank you!

Miss process under general independent requests

- Arrivals are i.i.d \implies **Renewal process**
- $F(t) = \mathbb{P}(X < t)$, $M(t) = \mathbb{E}[\mathcal{N}(t)]$, $T(t) = \mathbb{P}(T < t)$

Proposition (Exact CDF of the inter-miss times, Fofack et al. 2012, 2013)

Hit and Occupancy probability: $H_P \neq O_P = \lambda(1 - H_P)\mathbb{E}[Q]$

$$H_P^{(0)} = 1 - \frac{1}{1 + \mathbb{E}[M(T)]}; \quad H_P^{(1)} = \mathbb{E}[F(T)]$$

Miss stream is a **Renewal process**, and the CDF $G^{(r)}(t) = \mathbb{P}(Y < t)$ is given by

$$G^{(0)}(t) = F(t) - \int_0^t (1 - T(x)) dM(x) (1 - F(t - x)) \quad (7)$$

$$G^{(1)}(t) = F(t) - \int_0^t (1 - T(x)) dF(x) + \int_0^t (1 - T(x)) dF(x) G^{(1)}(t - x) \quad (8)$$

$$G^{(r)}(t) = (1 - r)G^{(0)}(t) + rG^{(1)}(t), \quad \text{for } \min(\Sigma, \mathcal{R})\text{-TTL model} \quad (9)$$

Miss process under Markov-correlated requests

- Arrivals occur at jumps of a DTMC $\{(t_k, \xi_k)_{k \geq 0}\}$ on $\mathcal{S} = \{1, 2, \dots, J\}$,
- Requests form a **stationary Markov Renewal Process (MRP)**,
 $[\mathbf{F}(t)]_{i,j} := \mathbb{P}(X_k < t, \xi_{k+1} = j | \xi_k = i)$, $T_i(t) = \mathbb{P}(T < t | \xi_k = i)$

Proposition (Exact Kernel of the miss process, *PhD thesis*)

Miss stream is a *Markov Renewal Process*, with the kernel $\mathbf{G}^{(r)}(t)$ given by

$$\mathbf{G}^{(0)}(t) = \mathbf{F}(t) - \int_0^t d\mathbf{R}(x)(\mathbf{I} - \mathbf{F}(t-x)) \quad (10)$$

$$\mathbf{G}^{(1)}(t) = \mathbf{F}(t) - \mathbf{L}(t) + \int_0^t d\mathbf{L}(x)\mathbf{G}^{(1)}(t-x) \quad (11)$$

$$\mathbf{G}^{(r)}(t) = (1-r)\mathbf{G}^{(0)}(t) + r\mathbf{G}^{(1)}(t), \quad \text{for } \min(\Sigma, \mathcal{R})\text{-TTL model} \quad (12)$$

$$[\mathbf{L}(t)]_{i,j} := \int_0^t (1 - T_i(x)) dF_{i,j}(x), \quad [\mathbf{R}(t)]_{i,j} := \int_0^t (1 - T_i(x)) dM_{i,j}(x)$$

Proofs

Hint: under (Markov) Renewal Assumption.

$$Y = \begin{cases} X_1 + \dots + X_{N(T_1)+1} & \text{if } r = 0 \\ X_1 \mathbf{1}(X_1 > T_1) + (X_1 + \tilde{Y}) \mathbf{1}(X_1 < T_1) & \text{if } r = 1. \end{cases}$$