



# A Simple Fluid Model for the Analysis of SQUIRREL

Florence Clévenot, Philippe Nain

## ► To cite this version:

Florence Clévenot, Philippe Nain. A Simple Fluid Model for the Analysis of SQUIRREL. RR-4911, INRIA. 2003. inria-00071669

**HAL Id: inria-00071669**

**<https://inria.hal.science/inria-00071669>**

Submitted on 23 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# *A Simple Fluid Model for the Analysis of SQUIRREL*

Florence Clévenot — Philippe Nain

**N° 4911**

Août 2003

\_\_\_\_\_ THÈME 1 \_\_\_\_\_



*rapport  
de recherche*



## A Simple Fluid Model for the Analysis of SQUIRREL

Florence Clévenot\*, Philippe Nain†

Thème 1 — Réseaux et systèmes  
Projet Mistral

Rapport de recherche n° 4911 — Août 2003 — 24 pages

**Abstract:** Peer-to-peer (P2P) systems are complex to analyze, due to their large number of users who connect intermittently, and to the frequency of requests for files or Web objects. In this paper we propose a mathematical model where request streams are represented as fluid flows and apply it to analyze Squirrel, a recent P2P cooperative Web cache. Our fluid model provides a low-complexity means to estimate the performance of Squirrel (hit probability, latency) and exhibits some key qualitative properties of this system. A comparison with discrete-event simulation validates the accuracy of our model.

**Key-words:** Peer-to-peer systems, content distribution networks, performance evaluation, stochastic processes, fluid models, Palm calculus

\* INRIA Sophia Antipolis. E-mail: Florence.Clevenot@sophia.inria.fr

† INRIA Sophia Antipolis. E-mail: Philippe.Nain@sophia.inria.fr

## Un Modèle Fluide Simple pour l'Analyse de SQUIRREL

**Résumé :** Les systèmes peer-to-peer (P2P) mettent en œuvre deux types d'événements: les connexions/déconnexions des utilisateurs et les requêtes pour des fichiers ou documents multimédia. Dans cet article nous proposons un modèle mathématique simple assimilant ces requêtes à des flux continus, modulés par les arrivées et départs aléatoires des usagers, et nous l'appliquons à Squirrel, un système de cache Web coopératif. Notre modèle permet de calculer les performances (probabilité de hit) de ce système avec une faible complexité, et met en évidence les propriétés qualitatives du système. Nous validons enfin la pertinence du modèle par comparaison avec une simulation à événements discrets.

**Mots-clés :** Systèmes peer-to-peer, évaluation de performances, modèles fluides stochastiques, calcul de Palm

## 1 Introduction

Beginning with the seminal work of Anick, Mitra and Sondhi in 1982 [1], stochastic fluid models have been successfully applied to a variety of packet-switching systems over the past 20 years (e.g., see [9, 10, 17, 18, 3, 19]). In these papers, detailed models of system behavior, which involve random arrivals of packets of discrete size to network nodes, are replaced by macroscopic models that substitute fluid flows for packet streams. The rates of the fluid flows are typically modulated by a stochastic process (such as a Markov process), thereby resulting in a “stochastic fluid model”. Although the resulting stochastic fluid models ignore the detailed, packet-level interactions, often they are mathematically tractable and accurate, providing significant insight into the qualitative properties of the original system.

We believe that stochastic fluid models are promising tools for modeling content distribution systems. Instead of replacing packet streams with fluid (as is done in the packet network modeling literature), we propose to replace content – such as images, MP3s, text files, video clips – with fluid. We have used this approach in [6] to study the performance (hit rate) of a cluster of Web caches. In this paper we use a stochastic fluid model to investigate the performance of Squirrel [14], a novel peer-to-peer (P2P) cooperative Web cache. The principle of Squirrel is to replace a corporate dedicated Web cache, by making client desktop machines cooperate in a peer-to-peer fashion to act globally as an efficient distributed Web cache.

Alternative approaches include Markovian analysis and event-driven (or trace-driven) simulations. These approaches have their own merit and we do not want to systematically oppose fluid models to more traditional approaches. Our take-home message is that simple (macroscopic) fluid models, whenever they apply, may give fairly accurate qualitative and quantitative results, and this at a low numerical complexity. Content distribution networks appear to be good candidates to illustrate our approach since they typically involve a large number of users and many parameters. These characteristics in turn imply large state spaces and a high numerical complexity when one uses detailed (microscopic) models such as Markovian models and simulations.

In Section 2 we provide an overview of Squirrel. Our fluid model is introduced in Section 3 and we use it in Section 4 to compute the main performance of Squirrel (hit probability, latency, etc.). In particular, we provide a simple expression for the hit probability, whose complexity is linear in the number of nodes in the Squirrel network. We show in Section 5 that our model provides substantial insight into performance issues of P2P cooperative Web caches such as Squirrel. Our analysis shows that two key parameters largely determine the performance of the system. In Section 6 we compare results obtained with the fluid model to results obtained with a discrete-event simulation of Squirrel. We find that the fluid model is both qualitatively and quantitatively accurate. We conclude in Section 7 with possible extensions of our fluid model.

## 2 Overview of Squirrel

Squirrel [14] is a decentralized, peer-to-peer Web cache that uses Pastry [20] as a location and routing protocol. When a client requests an object it first sends a request to the Squirrel proxy running on the client's machine. If the object is uncacheable then the proxy forwards the request directly to the origin Web server. Otherwise it checks the local cache, like every Web browser would do, in order to exploit locality and reuse. If a fresh copy of the object is not found in this cache, then Squirrel tries to locate one on some other node. To do so, it uses the distributed hash-table and the routing functionalities provided by Pastry. First, the URL of the object is hashed to give a 128-bit object identity (a number called *object-Id*) from a circular list; then the routing procedure of Pastry forwards the request to the node with the identity (called *node-Id*; this number is assigned randomly by Pastry to a participating node) the closest to *object-Id*. This node then becomes the *home node* for this object. Squirrel then proposes two schemes from this point on: *home-store* and *directory* schemes.

In the home-store scheme, objects are stored both at client caches and at its home node. The client cache may either have no copy of the requested object or a stale copy. In the former case the client issues a GET request to its home-node, and it issues a *conditional* GET (cGET) request in the latter case. If the home-node has a fresh copy of an object then it forwards it to the client or it sends the client a not-modified message depending on which action is appropriate. If the home-node has no copy of the object or has a stale copy in its cache, then it issues a GET or a cGET request, respectively, to the origin server. The origin server then either forwards a cacheable copy of the object or sends a not-modified message to the home-node. Then, the home-node takes the appropriate action with respect to the client (i.e. send a not-modified message or a copy of the object).

In the directory scheme the home-node for an object maintains a small directory of pointers to nodes that have recently accessed the object. A request for this object is sent randomly to one of these nodes. We will not go deeper into the description of this scheme since from now on we will only focus on the home-store scheme. We do so mainly because the latter scheme has been shown to be overall less attractive than the directory scheme [14]. In addition, the home-store scheme is more amenable to a fluid analysis than the directory scheme.

In a Squirrel network (a corporate network, a university network, etc.), like in any peer-to-peer system, clients arrive and depart the system at random times. There are two kinds of failures (or departures): abrupt and announced failures. Each failure has a different impact on the performance of Squirrel. An abrupt failure will result in a loss of objects. To see this, assume that node  $i$  is the home-node for object  $O$ . If node  $i$  fails, then a new home-node for object  $O$  has to be found by Pastry, as explained above, the next time object  $O$  is requested. Assume that the copy of object  $O$  was fresh when node  $i$  failed and consider the first GET request issued for  $O$  after the failure of node  $i$ . The GET request is therefore forwarded to the new home-node for object  $O$  (say node  $j$ ); this request will result in a miss if  $j$  has no copy of  $O$  or if its copy is stale. In this case, the failure of node  $i$  will yield a degradation in the performance since node  $j$  will have to contact the origin server to get a new copy

of object 0 or a not-modified message, as appropriate. If a node is able to announce its departure and to transfer its content to its immediate neighbors in the node-Id space before leaving Squirrel (announced failure), then no content is lost when the node leaves.

When a node joins Squirrel then it automatically becomes the home node for some objects but does not store those objects yet (see details in [14]). In case a request for one of those objects is issued, then its two neighbors in the node-Id space transfer a copy of the object, if any. Therefore, we can consider that there is no performance degradation in Squirrel due to a node arrival, since the transfer time between two nodes is supposed to be at least one order of magnitude smaller than the transfer time between any given node and the origin server.

From now on, the terms “node” and “client” will be used interchangeably.

### 3 A Model for Squirrel

#### 3.1 Modeling the client dynamics

To capture the dynamic behavior of the Squirrel nodes we use the following model: we assume that there are  $N < \infty$  clients who join and leave Squirrel independently of each other. The time until a given client joins (resp. leaves) is exponentially distributed with rate  $\lambda > 0$  (resp.  $\mu > 0$ ). If we denote by  $N(t) \in \{0, 1, 2, \dots, N\}$  the number of participating (i.e. connected) clients at time  $t \geq 0$ , then  $\{N(t)\}_t$  is a birth and death process, known in the literature as the *Ehrenfest* (or *Engset*) model [15].

Let  $N^\infty$  denote the stationary number of participating clients. Setting  $\rho = \lambda/\mu$ , we have [15, p. 17]

$$\mathbb{P}[N^\infty = i] = \binom{N}{i} \frac{\rho^i}{(1 + \rho)^N}, \quad 1 \leq i \leq N. \quad (1)$$

From now on we will consider that the process  $\{N(t)\}_t$  is in steady-state at time  $t = 0$ . Let  $0 \leq T_1 < T_2 < \dots$  denote the successive *jump times* of the (stationary) process  $\{N(t)\}_t$ . Introduce  $N_n \stackrel{\text{def}}{=} N(T_n+)$ , the stationary number of participating clients just *after* the occurrence of the  $n$ -th event/jump (i.e. join or leave of a client). It is shown in Appendix C that  $\pi_i \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} P(N_n = i)$ , the steady-state probability that there are  $i$  participating clients just after a jump, is given by

$$\begin{aligned} \pi_i &= \binom{N-1}{i-1} \frac{i + \rho(N-i)}{2i(1+\rho)^{N-1}} \rho^{i-1}, \quad 1 \leq i \leq N \\ \pi_0 &= \frac{1}{2(1+\rho)^{N-1}}. \end{aligned} \quad (2)$$

#### 3.2 Modeling the content dynamics

We assume that each participating client issues requests for objects at a constant rate  $\sigma > 0$ , and therefore model the request process of each client by a fluid flow. The motivation for this



fluid approximation is that requests occur at a much faster time scale (typically hundreds per second) than the nodes join/leave events (once a day or even less). Therefore, we expect the long-run average performance of the fluid model to be similar to that of the real, discrete-time system, where requests occur with any distribution with mean inter-request time  $1/\sigma$ .

In particular, our model represents the instantaneous set of cached objects in the whole Squirrel network by a global amount of fluid called  $X(t)$  at time  $t$ . This quantity of fluid will increase when objects are downloaded in the network from the origin server and added to their home node, i.e. whenever there is a cache miss. On the other hand, the amount of fluid will decrease as cached objects become stale. We assume that cached objects have the same constant time-to-live in cache, given by  $1/\theta$ ; this assumption is made both for the sake of simplicity and also because most caches use default time-to-live for objects without any specified expiration date (about 70% of requested objects [7]). The usual default value is 24 hours (see [7] for more details).

We assume that each node can store an unlimited number of objects. Though individual nodes would probably not dedicate too much memory to the collaborative cache, even reasonable cache sizes are sufficient to avoid losses due to a full cache; one reason for this is that cached objects become stale fast enough to avoid continuous increase of the content. For centralized caches, the largest size needed to avoid most capacity misses is dictated by the clients request rates [8] and is fairly small. We expect this property to be applicable to the Squirrel decentralized cache system.

As a result, the variation of the fluid is proportional to the miss rate and to the expiration rate  $\theta$  when the nodes population is constant<sup>1</sup>. If we call  $\mathbb{P}[\text{hit}|i, x]$  the hit probability when there are  $i$  connected nodes containing the fluid  $x$ , then the variation rate of the amount of fluid is

$$\frac{dx}{dt} = i\sigma (1 - \mathbb{P}[\text{hit}|i, x]) - \theta x$$

where  $i\sigma$  is the overall request rate in the Squirrel network when there are  $i$  connected nodes.

We now define an appropriate model for the hit probability function  $\mathbb{P}[\text{hit}|i, x]$ . Let us first call  $c$  the total number of objects that can be requested (in our model, the total amount of existing fluid in the universe). Since  $x$  is the quantity of cached fluid, a very simple model for the hit probability is  $\mathbb{P}[\text{hit}|i, x] = \frac{x}{c}$ . However, this linear function does not take into account the fact that some objects may be requested more often than others and thus are more likely to be present in the network. Since the popularity of Web objects follows a Zipf-like distribution [4], we can also model  $\mathbb{P}[\text{hit}|i, x]$  as a concave function of the type  $\mathbb{P}[\text{hit}|i, x] = \left(\frac{x}{c}\right)^\beta$ , which reflects the fact that once the most popular objects are present in the system, the fluid increase becomes slower.

It remains to specify how the node join and leave events impact the performance of our fluid model. We have seen in Section 2 that join events will probably not affect the performance of the system. On the other hand, we consider all failures (leaves) to be abrupt failures; this assumption is discussed in Section 4.4. Therefore, when a node leaves its share

<sup>1</sup> Assuming, of course, that at least one node is present.

of objects is lost to the system. If we assume that the requests are well balanced across all nodes of the network (property of the Pastry hashing technique), then a fraction  $1/i$  of the total amount of fluid is lost when a leave occurs if there were  $i$  nodes connected just before this leave event. This value has been confirmed empirically in [14].

For the sake of generality we introduce two mappings,  $\Delta_d(i)$  and  $\Delta_u(i)$ , that give the fluid *reduction* generated by a leave and a join, respectively, given that  $i$  nodes were connected before this leave/join event. For Squirrel,  $\Delta_d(i) = \frac{i-1}{i}$  and  $\Delta_u(i) = 1$  as discussed above. In other words, if the amount of fluid is  $x$  and that  $i$  nodes are connected before a leave (resp. join) then the amount of fluid just after this event will be  $x\Delta_d(i)$  (resp.  $x\Delta_u(i)$ ).

A glossary of the model parameters is provided in Table 1.

Table 1: System Parameters

$N$	Maximum number of nodes
$\lambda$	Birth rate of each Squirrel node
$\mu$	Death rate of each Squirrel node
$\rho$	$\lambda/\mu$
$\pi$	Stationary distribution of $\{N_n\}_n$
$\sigma$	Request rate per client
$\theta$	Expiration rate of cached objects
$c$	Total number of objects in the universe (i.e. total amount of fluid)
$\Delta_d(i)$	Fluid reduction after a node failure when there were $i \geq 1$ connected nodes. Default value: $(i-1)/i$ (cf. [14])
$\Delta_u(i)$	Fluid reduction after a node join when there were $i \geq 0$ connected nodes. Default value: 1 (cf. [14])

## 4 Performance analysis of the Squirrel p2p cache system

In this section we provide a simple closed-form expression for the hit probability of the Squirrel system, under the main assumption that all objects are all equally popular. Although somewhat unrealistic, this assumption leads to a clearer analysis and highlights the parameters interactions of the system. For practical numerical results, we show how this assumption can be relaxed in Section 4.3. The end-to-end latency reduction offered by the Squirrel system, which might be a more meaningful metric than the hit probability, can easily be derived from the following results as shown in Section 4.2. Finally, we discuss the possible sources of inaccuracy of this model in Section 4.4 and try to identify remedies whenever possible.

#### 4.1 Hit probability analysis

Under the equal popularity assumption, the hit probability is a linear function of the amount of cached fluid, as shown in Section 3.2. Our first task will be to characterize the fluid process  $\{X(t)\}_t$ .

Recall the definition of  $N(t)$  and  $N_n$ , the number of connected nodes at time  $t$  and at time  $T_n+$  (i.e. just after the  $n$ -th jump time), respectively (see Section 3.1). We assume that the sample-paths of  $\{N(t)\}_t$  and  $\{X(t)\}_t$  are left continuous. Hence,  $X_n \stackrel{\text{def}}{=} X(T_n+)$  is the amount of cached fluid just *after* the  $n$ -th jump in the process  $\{N(t)\}_t$ . For easy reference, the main definitions and notation have been collected in Table 2.

Table 2: Variables

$X(t)$	Total amount of fluid in the system
$N(t)$	Number of connected nodes at time $t$
$\{T_n\}_n$	Jump times of the process $\{N_t\}_t$
$N_n = N(T_n^+)$	Number of connected nodes just after the $n$ -th jump.
$X_n = X(T_n^+)$	Total amount of fluid just after the $n$ -th jump.
$Y_n = X(T_{n+1}^-)$	Total amount of fluid just <i>before</i> the $n+1$ -th jump. (at the end of the $n$ -th period)
$v_i$	$\lim_{n \rightarrow \infty} \mathbb{E}[Y_n   N_n = i] / c$
$\eta_i$	$c / (1 + \alpha/i)$
$\gamma$	$\sigma / (\mu c)$
$\alpha$	$(\theta c) / \sigma$

The fluid process is defined as described in Section 3.2: between two consecutive jumps  $(T_n, T_{n+1})$  of  $\{N(t)\}_t$  the fluid increases at rate

$$\frac{d}{dt}X(t) = \sigma N_n \left(1 - \frac{X(t)}{c}\right) - \theta X(t) \quad (3)$$

provided that  $N_n > 0$ . Integrating (3) gives

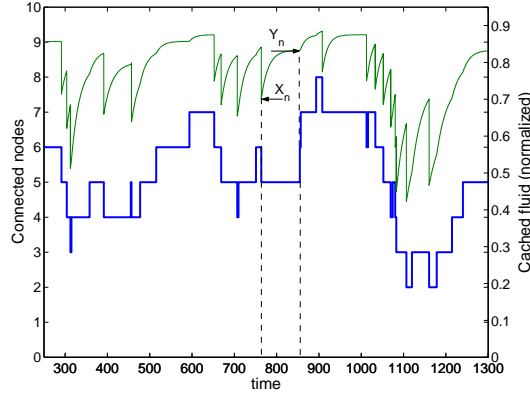
$$X(t) = \frac{\sigma N_n}{\frac{\sigma N_n}{c} + \theta} + \left(X_n - \frac{\sigma N_n}{\frac{\sigma N_n}{c} + \theta}\right) e^{-(t-T_n)(\theta + \frac{\sigma N_n}{c})}$$

for  $T_n < t < T_{n+1}$  provided that  $N_n > 0$ . If  $N_n = 0$  then  $X(t) = 0$  for  $T_n < t < T_{n+1}$ .

If  $T_n$  corresponds to a node leave (resp. join) then the amount of cached fluid is reduced as follows

$$X_n = \Delta_d(N_n)X(T_n-) \text{ (resp. } X_n = \Delta_u(N_n)X(T_n-))$$

Therefore,  $\{X(t)\}_t$  is a piecewise (exponential) process, with randomness at jump times  $\{T_n\}_n$ . A sample path of the process  $\{(N(t), X(t))\}_t$  is represented on Fig. 1.

Figure 1: Sample path of  $\{(N(t), X(t))\}_t$ .

For the sake of convenience we introduce

$$\eta_i \stackrel{\text{def}}{=} \frac{c}{1 + \frac{\theta c}{i\sigma}}. \quad (4)$$

We can now re-write the solution of (3) as

$$X(t) = \eta_{N_n} + (X_n - \eta_{N_n}) e^{-(t - T_n) \frac{\sigma N_n}{\eta_{N_n}}}. \quad (5)$$

The process  $\{(N(t), X(t))\}_t$  is an irreducible Markov process on the set  $\{0, 0\} \cup \{1, 2, \dots, N\} \times [0, c]$ . Denote by  $X$  the stationary regime of  $\{X(t)\}_t$ .

We define the steady-state hit probability  $p_H$  as

$$p_H = \frac{\mathbb{E}[X]}{c} \quad (6)$$

We give a simple formula for  $p_H$  in Proposition 4.1. This formula is expressed in terms of the following new parameters that will play a key role in the understanding of the system behavior<sup>2</sup>:

$$\alpha \stackrel{\text{def}}{=} \frac{\theta c}{\sigma} \quad \text{and} \quad \gamma \stackrel{\text{def}}{=} \frac{\sigma}{\mu c}. \quad (7)$$

**Proposition 4.1** *Assume that for  $i=0, \dots, N-1$ ,*

$$0 \leq \Delta_u(i) \Delta_d(i+1) \leq 1. \quad (8)$$

<sup>2</sup>The system is defined in terms of 6 parameters:  $N, c, \rho, \mu, \theta, \sigma$ ; definitions in (7) will allow us to express the hit probability only in terms of 4 parameters, namely,  $N, \rho, \alpha$  and  $\gamma$ , as shown in Proposition 4.1.

The hit probability  $p_H$  is given by

$$p_H = \frac{1}{(1+\rho)^N} \sum_{i=1}^N \binom{N}{i} \rho^i v_i \quad (9)$$

where the vector  $\mathbf{v} = (v_1, \dots, v_N)^T$  is the unique solution of the linear equation

$$A \mathbf{v} = \mathbf{b} \quad (10)$$

with  $\mathbf{b} = (b_1, \dots, b_N)^T$  a vector whose components are given by  $b_i = \gamma i$  for  $1 \leq i \leq N$ , and  $A = [a_{i,j}]_{1 \leq i,j \leq N}$  a  $N \times N$  tridiagonal matrix whose non-zero elements are

$$\begin{aligned} a_{i,i} &= \alpha\gamma + (\gamma + 1)i + \rho(N - i), & 1 \leq i \leq N \\ a_{i,i-1} &= -i\Delta_u(i-1), & 2 \leq i \leq N \\ a_{i,i+1} &= -\rho(N-i)\Delta_d(i+1), & 1 \leq i \leq N-1. \end{aligned}$$

**Proof.** The idea of the proof is to first compute the expected amount of cached fluid just before a jump in the process  $\{N(t)\}_t$  conditioned on the value of  $N(t)$  just before this jump, and then invoke Palm calculus to deduce the expected amount of cached fluid at *any time*.

Let  $Y_n$  be the amount of cached fluid just before the  $(n+1)$ -st jump in the process  $\{N(t)\}_t$  (i.e.  $Y_n = X(T_{n+1}-)$ ). We first compute  $v_i \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} (1/c) \mathbb{E}[Y_n | N_n = i]$  for  $1 \leq i \leq N$ . We show in Appendix A that  $v_i$  satisfies the following recursive equation:

$$\begin{aligned} v_i (\rho(N-i) + \alpha\gamma + (\gamma + 1)i) &= i\Delta_u(i-1)v_{i-1} \\ &+ \rho(N-i)\Delta_d(i+1)v_{i+1} + i\gamma \end{aligned} \quad (11)$$

for  $i = 1, 2, \dots, N$ , or equivalently (10) in matrix form with  $\mathbf{v} = (v_1, \dots, v_N)$ .

The uniqueness of the solution of (10) is shown in Appendix D. The vector  $\mathbf{v}$  in (10) gives the conditional stationary expected amount of cached fluid just before jump epochs (up to a multiplicative constant).

However, the hit probability  $p_H$  in (6) is defined in terms of the stationary expected amount of cached fluid at *arbitrary* epochs. The latter metric can be deduced from the former by using Palm calculus, through the identity (see e.g. [2, Formula (4.3.2)])

$$\mathbb{E}[X] = \Lambda \mathbb{E}^0 \left[ \int_0^{T_1} X(t) dt \right] \quad (12)$$

where  $\mathbb{E}^0$  denotes the expectation with respect to the Palm distribution<sup>3</sup>,  $T_1$  denotes the time of the first jump after 0, and where  $\Lambda$  denotes the global rate of the Engset model, i.e.

$$\Lambda = \frac{1}{\mathbb{E}^0[T_1]}. \quad (13)$$

---

<sup>3</sup>The Palm distribution is the distribution of the process  $\{X_t\}_t$  assuming that a jump occurs at time 0 and that the system is in steady-state at time 0.

From now on we assume that the system is in steady-state at time 0. Under the Palm distribution we denote by  $N_{-1}$  and  $Y_{-1}$  the number of connected nodes and the amount of cached fluid respectively, just before time 0 (i.e. just before the jump to occur at time 0).

We first compute  $1/\Lambda$ . Using (2) we find

$$\begin{aligned} \frac{1}{\Lambda} &= \sum_{i=0}^N \pi_i \mathbb{E}^0[T_1 | N_0 = i] = \frac{1}{\mu} \sum_{i=0}^N \frac{\pi_i}{\rho(N-i) + i} \\ &= \frac{1 + \rho}{2N\rho\mu}. \end{aligned} \quad (14)$$

We show in Appendix B that

$$\mathbb{E}[X] = \frac{c}{(1 + \rho)^N} \sum_{i=1}^N \binom{N}{i} \rho^i v_i. \quad (15)$$

Dividing both sides of (15) by  $c$ , we get (9), which concludes the proof.  $\diamond$

Conditions (8) in Proposition 4.1 ensure that the system (10) has a unique solution (see Appendix C). They are satisfied for the home store scheme (since  $\Delta_u(i)\Delta_d(i+1) = i/(i+1)$ ).

**Remark 4.1** *Since  $A$  is a tridiagonal matrix, (10) can be solved in only  $\mathcal{O}(N)$  operations, once the mappings  $\Delta_u$  and  $\Delta_d$  are specified.*

## 4.2 Latency reduction

The expected delay to fetch a document can easily be derived from the hit probability as follows: let  $T_e$  and  $T_i$  be the external and internal latency, respectively. The internal latency is the average delay induced by the local network, and thus is experienced by clients even in case of home node hits. The external latency is caused by network bottlenecks and Web server delays outside the organization, and is added to the internal latency in case of a miss when an object has to be retrieved from the origin Web server by its home node and is then sent to the client through the local network. Typically,  $T_e$  accounts for most of the total latency in the absence of caching (e.g. 88% for geographically located network [16]). The total expected delay with Squirrel is

$$\mathbb{E}[T] = T_i p_H + (T_i + T_e)(1 - p_H).$$

The Squirrel cache system reduces the average delay by saving the external latency whenever there is a hit. The relative latency reduction observed with Squirrel is thus

$$\frac{T_i + T_e - \mathbb{E}[T]}{T_i + T_e} = p_H \frac{T_e}{T_i + T_e}.$$

### 4.3 Zipf popularity

A concave model for the hit probability such as  $\mathbb{P}[\text{hit}|i, x] = (x/c)^\beta$  makes the differential equation (3) nonlinear and without a known closed-form solution. An alternative approach is to consider  $K$  classes of decreasing popularity: the first class contains the  $c_1$  most popular objects, and the  $K$ -th class the  $c_K$  least popular ones. With  $c$  being the total number of existing objects, we have  $\sum_{k=1}^K c_k = c$ . Let us call  $o_k$  the index of the least popular object of class  $k$  in the ordered set of objects, that is,

$$o_k = \sum_{l=1}^k c_l. \quad (16)$$

Using the Zipf-like popularity distribution [4], the probability of each class of fluid can easily be obtained as

$$p_k = \mathbb{P}[\text{request for class } k] = \frac{o_k^{1-\beta} - o_{k-1}^{1-\beta}}{c^{1-\beta}} \quad (17)$$

for class  $k$ , where  $\beta$  is the skew factor of the Zipf distribution.

The request rate for each class is given by  $\sigma_k = \sigma p_k$  for  $k = 1, 2, \dots, K$ , where  $\sigma$  is the average request rate per node (regardless of object popularity).

The hit probability  $\hat{p}_H$  is now defined as the weighted sum of the conditional hit probabilities within each class, namely,

$$\hat{p}_H = \sum_{k=0}^K \frac{\mathbb{E}[X_k]}{c_k} p_k \quad (18)$$

where  $X_k$  is the (stationary) amount of cached fluid of class  $k$ . Similar to the derivation of the hit probability  $p_H$  in Proposition 4.1, we find that

$$\hat{p}_H = \sum_{k=0}^K p_k \frac{1}{c_k(1+\rho)^N} \sum_{i=1}^N \binom{N}{i} \rho^i v_i^{(k)} \quad (19)$$

where the vector  $\mathbf{v}^{(k)} = (v_1^{(k)}, \dots, v_N^{(k)})^T$  is the unique solution of the linear equation

$$A^{(k)} \mathbf{v}^{(k)} = \mathbf{b}^{(k)} \quad (20)$$

with  $\mathbf{b}^{(k)} = (b_1^{(k)}, \dots, b_N^{(k)})^T$  a vector whose components are given by  $b_i^{(k)} = \frac{i\sigma_k}{\mu}$  for  $i = 1, 2, \dots, N$ , and  $A^{(k)} = [a_{i,j}^{(k)}]_{1 \leq i,j \leq N}$  a  $N \times N$  tridiagonal matrix whose non-zero elements are

$$\begin{aligned} a_{i,i}^{(k)} &= \frac{\theta}{\mu} + \frac{\sigma_k i}{\mu c_k} + i + \rho(N-i), & 1 \leq i \leq N \\ a_{i,i-1}^{(k)} &= -i\Delta_u(i-1), & 2 \leq i \leq N \\ a_{i,i+1}^{(k)} &= -\rho(N-i)\Delta_d(i+1), & 1 \leq i \leq N-1. \end{aligned}$$

#### 4.4 Discussion and extensions

We now discuss some specific features that were not explicitly taken into account in the analysis of Section 4.1, apart from the popularity of documents.

The first remark is that the model assumes that every requested object is saved in the cooperative cache when downloaded a first time from the origin server. However, a non-negligible fraction (around 28%, cf. [7]) of the requested objects is in practice non-cacheable (explicit non-cacheable, expiration date before current date, dynamically generated documents such as cgi-forms, etc.). We can take into account the uncacheability in our model as follows: let  $u$  be the fraction of objects that are uncacheable. So far, we have considered that the fluid increases after each miss, thereby implicitly assuming that all objects are cacheable. The uncacheability can be incorporated in our model by considering that only a fraction  $1 - u$  of misses will yield a fluid increase. This gives rise to the following equation

$$\begin{aligned} \frac{d}{dt}X(t) &= (1 - u)\sigma N_n \left(1 - \frac{X(t)}{c}\right) - \theta X(t) \\ &= (1 - u)\sigma N_n - \left(\frac{(1 - u)\sigma N_n}{c} + \theta\right) X(t) \end{aligned} \quad (21)$$

for  $T_n < t < T_{n+1}$  and  $N_n \in \{1, 2, \dots, N\}$ , since only requests for cacheable objects will lead to a fluid increase. Therefore, uncacheable objects can be added to the model simply by modifying the request rate accordingly.

Secondly, the impact of node join and leave events, modeled through the mappings  $\Delta_u$  and  $\Delta_d$ , may be slightly different from the values described in Section 3. Indeed, two factors need to be taken into account. On the one hand, some nodes may announce their intention to disconnect, thereby avoiding performance degradation (see Section 2). This would change the value of  $\Delta_d(i)$ . Though Proposition 4.1 provides an expression for general values of  $\Delta_d(i)$ , this would require a re-estimation of  $\Delta_d(i)$ . Since it would not exceed one, condition (8) would still be satisfied. On the other hand, the individual Squirrel caches may be stored either on disk or in memory. In the first case, a node  $i$  may join with a previously stored set of documents that have not been removed from its disk cache when the node went down. This would possibly add fluid into the network, if the corresponding objects have not been retrieved by the system while  $i$  was down and if  $i$  has not announced its last departure. The problem would not only be to re-estimate  $\Delta_u(i)$ , but also that  $\Delta_u(i)$  might be greater than one, making condition (8) more difficult to verify. However, we expect that node  $i$  will stay down for a minimum time that will be orders of magnitude greater than requests inter-arrival times (the reboot time is typically a few minutes). Meanwhile, most of the objects stored in node  $i$  will be requested again and added to their new home nodes. As a result, when  $i$  will join again the Squirrel network with its own full cache, it will probably not add any fluid in the system, thus guaranteeing that  $\Delta_u(i) \leq 1$ , and thereby the validity of assumption (8).

Finally, formula (9) involves binomial coefficients  $\binom{N}{i}$  and an exponential in  $N$ . Therefore, computing  $p_H$  accurately for very large values of  $N$  might reveal difficult. Nonetheless,



we would like first to mention that though we have occasionally encountered such problems, Proposition 4.1 is clearly tractable for the order of magnitude of several thousands of nodes, where a simulation would be untractable for high-confidence results. In addition, for much larger values of  $N$ , we believe that the node dynamics can be approximated by an  $M/M/\infty$  model instead of an Engset model. This extension is a work in progress.

## 5 Qualitative insight in the Squirrel system

Proposition 4.1 shows that the performance of the Squirrel system exhibits only four degrees of freedom:  $N, \rho, \gamma, \alpha$  while our model introduced six parameters  $(N, \lambda, \mu, \sigma, \theta, c)$ . We now examine the relative importance of these new parameters and how they characterize the Squirrel system behavior.

We first rapidly examine the influence of  $\rho$ . Fig. 2 shows that while there is a sharp drop of the hit probability for very small values of  $\rho$  (smaller than one), the performance is almost constant when  $\rho$  increases. Therefore, except when it is close to zero,  $\rho$  has very

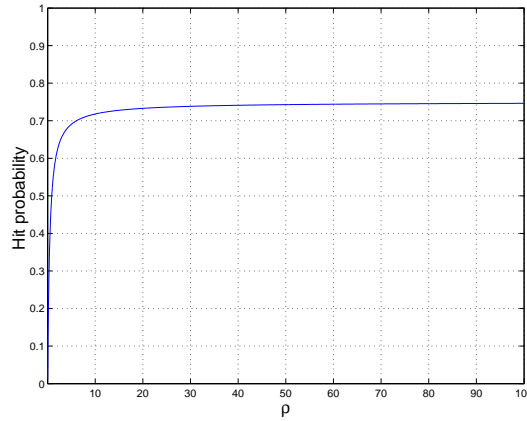


Figure 2: Impact of  $\rho$  (with  $N = 3$ ,  $\alpha = 1$  and  $\gamma = 2$ ).

little influence on the performance of the Squirrel system. Also, it is very unlikely that  $\rho$  will be really small, since it would mean a non-negligible probability of reaching a state when all nodes are down (a very unrealistic situation). Under these circumstances, the limiting factors for the hit probability will be parameters  $\gamma$  and  $\alpha$ .

In Fig. 3 we examine the comparative influence of  $\gamma$  and  $\alpha$  on the hit probability. We find that for fixed  $\alpha$ , the hit probability as a function of  $\gamma$  follows a concave shape, and can reach almost one when  $\alpha = 0$ . (This is consistent with our observation that  $\rho$  does not limit the hit probability when greater or equal to one.) Recall that  $\gamma = \sigma/(\mu c)$  where  $\sigma$  is

the individual request rate of the nodes. This concave shape in  $\gamma$  reminds us the log-like<sup>4</sup> performance of a centralized Web cache (or Web cache cluster) as described in [21, 11, 8].

However, we observe that the hit probability is high on a very narrow domain ( $\alpha \leq 1$ ,  $\gamma \geq 10$ ). Indeed, there is a strong attenuation in  $\alpha$ , so  $\gamma$  only has a real impact on the performance when  $\alpha$  is small.

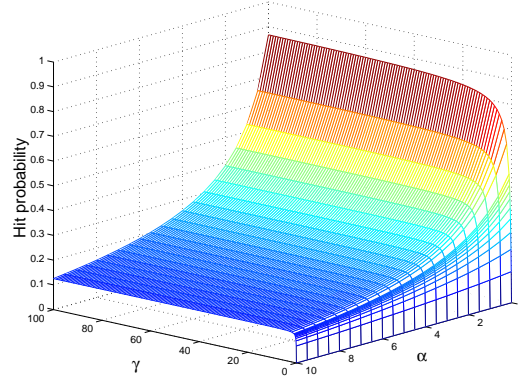


Figure 3: Impact of  $\gamma$  and  $\alpha$  on the hit probability (with  $N = 3$  and  $\rho = 1$ ). (Note that  $\alpha$  is decreasing.)

These observations suggest possible methods to improve the performance of the Squirrel system. The best possible improvement will be to reduce parameter  $\alpha = \theta c / \sigma$ . Since the total number of existing objects,  $c$ , cannot be modified, there are two options:

- Reduce  $\theta$  as much as possible: change the default value of the HTTP parameter CONF\_MAX<sup>5</sup> (usually 24 hours in HTTP/1.1) in the freshness calculation heuristic for example, especially since most cGET requests (e.g. 90%) are responded with Not-Modified message [7].
- Increase  $\sigma$  (which determines the rate at which objects are retrieved to the Squirrel network), for instance by using prefetching techniques. We believe that prefetching can be incorporated to the fluid model, which will allow us to quantify the gain of using it.

<sup>4</sup>The hit rate is either a logarithm or a small power of the global request rate.

<sup>5</sup>In the freshness heuristic, the lifetime is  $\min(\text{CONF\_MAX}, \text{CONF\_PERCENT} \times (\text{DATE\_LASTMODIFIED}))$ .

## 6 Experimental results

In this section we compare quantitatively our macroscopic fluid model with a discrete-event driven simulation of the Squirrel home-store system. Request arrivals are Poisson and object time-to-live are taken to be all constant and all identical. We also assume that nodes follow the same time-evolution as in the fluid model, i.e. an Engset model. The external latency is taken into account whereas the internal latency is considered to be zero (corresponding to instantaneous internal transfers). Simulation results are given with 99% confidence intervals.

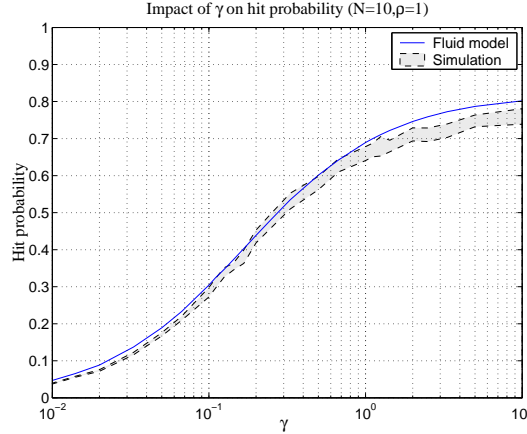


Figure 4: Fluid model vs discrete-event simulation. ( $N = 10$ ,  $\rho = 1$  and  $\alpha = \gamma$ ).

Fig. 4 displays the hit probability as a function of  $\gamma$  with  $\rho = 1$  and  $\alpha = 1$ . We observe that the fluid model curves closely follow the same shapes as the discrete-event simulations and therefore mimics the simulated system behavior very accurately. We conclude that the model is robust to assumptions such as the request rate distribution (which we assumed constant in Section 4.1), and though microscopic features such as objects replication and local hits (requests not forwarded to home node) are being ignored, the fluid model provides an accurate approximation for the actual performance of the Squirrel system.

Furthermore, we would like to emphasize that each simulation, even for very small networks (10 nodes), ran for several hours (typically 20 hours or more) on a Pentium 4 running at 2.66GHz. Even if our code may not be fully optimized, it is clear that in comparison with the instantaneous results provided by the fluid model, simulation is very slow and limited to very small network sizes.

We show in Fig. 5 how the hit probability would look like for large networks, since simulation of such systems would be either too slow or statistically irrelevant. Since Fig. 4 validated the accuracy of our model for small network sizes, we expect the results for large networks to be as relevant – though we do not have simulation results to demonstrate it. We observe the same shape as in Fig. 4, though on a larger range (thanks to the low complexity

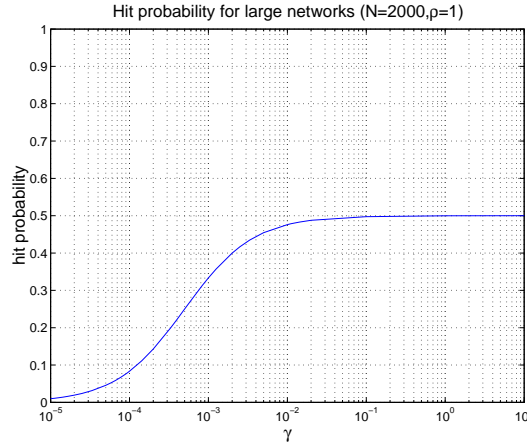


Figure 5: Hit probability for large networks ( $N = 2000$  and  $\alpha = 1000.0$ ).

of the model results), which suggests that Squirrel scales with the same type of behavior, and that the characteristics observed in Section 5 should be valid for large networks.

## 7 Conclusion and perspectives

In this paper we have proposed and studied a fluid model for the performance analysis of the Squirrel cooperative cache system. To cope with the large number of users that join and leave the cache system randomly, we have approximated the request streams of the individual nodes by a fluid flow. Our resulting stochastic fluid model turns out to be mathematically tractable, and has allowed us to provide a simple and very low-complexity procedure for computing the hit probability. Moreover, the analysis has emphasized the key characteristics of the Squirrel system and allows a better understanding of its performance. Comparison with simulation results has shown that the hit probability provided by the solution to the model is an accurate approximation of the actual hit rate and has validated the qualitative conclusions driven by the model results.

Future work will focus on extending the model to handle prefetching techniques and larger populations of peers. Also, the accuracy and tractability of this fluid model suggest that it could be used to analyze other content distribution systems, such as P2P file sharing applications and CDNs.

## Acknowledgment

The authors thank K. W. Ross for useful discussions during the course of this work.

## A Proof of equation (11)

Recall that  $v_i = \lim_{n \rightarrow \infty} \frac{\mathbb{E}[Y_n | N_n = i]}{c}$  for  $1 \leq i \leq N$ . With (5) we have

$$\begin{aligned}
 v_i &= \frac{1}{c} \lim_{n \rightarrow \infty} \mathbb{E}[Y_n | N_n = i] \\
 &= \frac{1}{c} \lim_{n \rightarrow \infty} \mathbb{E} \left[ \eta_i + (X_n - \eta_i) e^{-(T_{n+1} - T_n) \frac{\alpha i}{\eta_i}} | N_n = i \right] \\
 &= \frac{1}{c} \lim_{n \rightarrow \infty} \left( \eta_i \times \frac{\alpha \gamma + \gamma i}{\alpha \gamma + \gamma i + \rho(N - i) + i} \right. \\
 &\quad \left. + \frac{(\rho(N - i) + i) \mathbb{E}[X_n | N_n = i]}{\rho(N - i) + i + \alpha \gamma + \gamma i} \right). \tag{22}
 \end{aligned}$$

To derive (22) we have used the fact that, given  $N_n = i$ , the random variables  $X_n$  and  $T_{n+1} - T_n$  are independent, and  $T_{n+1} - T_n$  is exponentially distributed with parameter  $(N - i)\lambda + \mu i$ .

Let us now evaluate  $\lim_{n \rightarrow \infty} \mathbb{E}[X_n | N_n = i]$ . Conditioning on  $N_{n-1}$  we have

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \mathbb{E}[X_n | N_n = i] &= \\
 &\quad \lim_{n \rightarrow \infty} \mathbb{E}[X_n | N_n = i, N_{n-1} = i - 1] \\
 &\quad \times \mathbb{P}[N_{n-1} = i - 1 | N_n = i] \\
 &\quad + \lim_{n \rightarrow \infty} \mathbb{E}[X_n | N_n = i, N_{n-1} = i + 1] \\
 &\quad \times \mathbb{P}[N_{n-1} = i + 1 | N_n = i] \mathbb{1}_{[i < N]} \\
 &= \Delta_u(i - 1) \lim_{n \rightarrow \infty} \mathbb{E}[Y_{n-1} | N_{n-1} = i - 1] \\
 &\quad \times \frac{\pi_{i-1}}{\pi_i} \frac{\rho(N - i + 1)}{\rho(N - i + 1) + i - 1} \\
 &\quad + \Delta_d(i + 1) \lim_{n \rightarrow \infty} \mathbb{E}[Y_{n-1}, | N_{n-1} = i + 1] \\
 &\quad \times \frac{\pi_{i+1}}{\pi_i} \frac{i + 1}{\rho(N - i - 1) + i + 1} \mathbb{1}_{[i < N]} \\
 &= c \frac{i \Delta_u(i - 1) v_{i-1} + \Delta_d(i + 1) v_{i+1} \rho(N - i)}{\rho(N - i) + i} \tag{23}
 \end{aligned}$$

by using (2) and the definition of  $v_i$ . Finally, dividing both sides by  $c$  and introducing (23) into (22) yields (11).  $\diamond$

## B Proof of equation (15)

Let us determine  $\mathbb{E}[X]$  from the  $v_i$ s. We use the Palm formula and condition on the value of  $N_0$ . From (12), (5), (14) we find

$$\begin{aligned}
 \mathbb{E}[X] &= \Lambda \sum_{i=1}^N \pi_i \mathbb{E}^0 \left[ \int_0^{T_i} \left( \eta_i + (X_0 - \eta_i) e^{-t \frac{\sigma_i}{\eta_i}} \right) dt \mid N_0 = i \right] \\
 &= \Lambda \left[ \sum_{i=1}^N \pi_i \eta_i \mathbb{E}^0[T_1 \mid N_0 = i] + \sum_{i=1}^N \frac{\pi_i \eta_i}{\sigma_i} \right. \\
 &\quad \times \mathbb{E}^0 \left[ (X_0 - \eta_i) \left( 1 - e^{-T_1 \frac{\sigma_i}{\eta_i}} \right) \mid N_0 = i \right] \Big] \\
 &= \Lambda \left[ \sum_{i=1}^N \pi_i \eta_i \frac{1}{\lambda(N-i) + \mu i} + \sum_{i=1}^N \frac{\pi_i \eta_i}{\sigma_i} \right. \\
 &\quad \times (\mathbb{E}^0[X_0 \mid N_0 = i] - \eta_i) \\
 &\quad \times \left( 1 - \mathbb{E}^0 \left[ e^{-T_1 \frac{\sigma_i}{\eta_i}} \mid N_0 = i \right] \right) \Big] \tag{24}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{\Lambda}{\mu} \left[ \sum_{i=1}^N \pi_i \eta_i \frac{1}{\rho(N-i) + i} \right. \\
 &\quad \left. + \sum_{i=1}^N \pi_i \frac{\mathbb{E}^0[X_0 \mid N_0 = i] - \eta_i}{\rho(N-i) + i + \alpha\gamma + \gamma i} \right] \\
 &= \frac{2N\rho}{1+\rho} \sum_{i=1}^N \pi_i \left[ \eta_i \frac{1}{\rho(N-i) + i} + \right. \\
 &\quad \left. \frac{\mathbb{E}^0[X_0 \mid N_0 = i] - \eta_i}{\rho(N-i) + \alpha\gamma + (\gamma+1)i} \right]. \tag{25}
 \end{aligned}$$

By definition,  $\mathbb{E}^0[X_0 \mid N_0 = i] = \lim_{n \rightarrow \infty} \mathbb{E}[X_n \mid N_n = i]$ , which has been computed in (23). By combining (23) and (11) we obtain

$$\mathbb{E}^0[X_0 \mid N_0 = i] = \frac{(\rho(N-i) + i + \alpha\gamma + i\gamma)v_i - i \frac{\sigma}{\mu}}{\rho(N-i) + i}.$$

Plugging this value of  $\mathbb{E}^0[X_0 \mid N_0 = i]$  into the r.h.s. of (25), and using (2), yields after some straightforward algebra:

$$\mathbb{E}[X] = \frac{c}{(1+\rho)^N} \sum_{i=1}^N \binom{N}{i} \rho^i v_i.$$

which is nothing but (15).  $\diamond$

## C Stationary distribution of the Engset model at jump times

In this section we compute the limiting distribution of the Markov chain  $\{N_n, n \geq 1\}$ . Let  $P = [p_{i,j}]_{0 \leq i,j \leq N}$  be its transition probability matrix. We have  $p_{i,i+1} = \rho(N-i)/(\rho(N-i)+i)$  for  $0 \leq i \leq N-1$ ,  $p_{i,i-1} = i/(\rho(N-i)+i)$  for  $1 \leq i \leq N$  and  $p_{i,j} = 0$  when  $|i-j| \neq 1$ .

Since this Markov chain<sup>6</sup> has a finite-state space and is irreducible, it is positive recurrent [5, Cor. 5.3.19, 5.3.22]. Therefore, it possesses a unique stationary distribution  $\pi = (\pi_0, \dots, \pi_N)$  given by the (unique) solution of the equation  $\pi = \pi P$  such that  $\sum_{i=0}^N \pi_i = 1$  [12, page 208].

We proceed by induction to compute  $\pi$ . From the equation  $\pi = \pi P$  we find that  $\pi_1 = (\rho(N-1)+1)\pi_0$  and  $\pi_2 = \frac{\rho(N-2)+2}{2}\rho(N-1)\pi_0$ . This suggests that

$$\pi_j = \frac{\rho(N-j)+j}{j} \frac{\rho^{j-1}}{(j-1)!} \frac{(N-1)!}{(N-j)!} \pi_0 \quad (26)$$

for  $j = 1, 2, \dots, N$ . Assume that (26) holds for  $j = 1, 2, \dots, i < N-1$ . Let us show that it still holds for  $j = i+1$ . We have

$$\begin{aligned} \pi_{i+1} &= \frac{\rho(N-(i+1))+i+1}{i+1} \\ &\quad \times \left( \pi_i - \frac{\rho(N-(i-1))}{\rho(N-(i-1))+i-1} \pi_{i-1} \right) \\ &= \frac{\rho(N-i-1)+i+1}{i+1} \left( \frac{\rho(N-i)+i}{i!} \frac{\rho^{i-1}}{(N-i)!} \right. \\ &\quad \left. - \frac{\rho(N-i+1)}{\rho(N-i+1)+i-1} \times \frac{i-1+\rho(N-i+1)}{i-1} \right. \\ &\quad \left. \times \frac{\rho^{i-2}}{(i-2)!(N-i+1)!} \right) (N-1)! \pi_0 \\ &= \frac{\rho(N-i-1)+i+1}{i+1} \frac{\rho^i (N-1)!}{i! (N-i-1)!} \pi_0 \end{aligned} \quad (27)$$

where (27) follows from the induction hypothesis. The constant  $\pi_0$  is computed by using the normalizing condition  $\sum_{i=0}^N \pi_i = 1$ ; we find  $\pi_0 = 1/(2(1+\rho)^{N-1})$  as announced in (2). Plugging this value of  $\pi_0$  into (26) gives the general expression of (2).  $\diamond$

## D Uniqueness of the solution of (10)

The linear system (10) defined in Proposition 4.1 admits a unique solution if and only if  $\det(A) \neq 0$ . Since  $A$  is a tridiagonal matrix we can use the LU decomposition [13, Sec. 3.5]

<sup>6</sup>Note that this chain is periodic with period 2.

$A = LU$  with

$$L = \begin{pmatrix} l_1 & 0 & \cdots & 0 \\ \beta_2 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \beta_n & l_n \end{pmatrix}$$

and

$$U = \begin{pmatrix} 1 & u_1 & \cdots & 0 \\ 0 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & u_{n-1} \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

where  $l_i$ 's and  $u_i$ 's are defined as follows:

$$\begin{aligned} a_{1,1} &= l_1 \\ a_{i,i} &= l_i + a_{i,i-1}u_{i-1}, \quad i = 2, \dots, N \\ l_i u_i &= a_{i,i+1}, \quad i = 1, \dots, N-1. \end{aligned}$$

Both matrices  $L$  and  $U$  being bidiagonal matrices it follows that  $\det(A) \neq 0$  if and only if  $l_i \neq 0$  for  $i = 1, 2, \dots, N$ .

We use an induction argument to show that  $l_i \neq 0$  for  $i = 1, 2, \dots, N$ . We have  $l_1 = \gamma + \rho(N-1) + 1$ . Assume that  $l_i > \gamma + \rho(N-i)$  for  $i = 1, 2, \dots, n < N-1$  and let us show that  $l_{n+1} > \gamma + \rho(N-n-1)$ . We have

$$\begin{aligned} l_{n+1} &= a_{n+1,n+1} - \frac{a_{n+1,n}a_{n,n+1}}{l_n} \\ &= \gamma + \rho(N-n-1) \\ &\quad + (n+1) \frac{l_n - \rho(N-n)\Delta_u(n)\Delta_d(n+1)}{l_n} \\ &> \gamma + \rho(N-n-1) \end{aligned}$$

by using the induction hypothesis together with the fact that  $0 \leq \Delta_u(n)\Delta_d(n+1) \leq 1$ .  $\diamond$

## References

- [1] ANICK, D., MITRA, D., AND SONDDHI, M. M. Stochastic theory of data-handling systems with multiple sources. *Bell Systems Technical Journal* 61 (1982), 1871–1894.
- [2] BACCELLI, F., AND BRÉMAUD, P. *Elements of Queueing Theory: Palm-Martingale Calculus and Stochastic Recurrences*. Springer Verlag, 1994.



- [3] BOORSTYN, R., BURCHARD, A., LIEBEHERR, J., AND OOTTAMAKORN, C. Statistical service assurances for traffic scheduling algorithms. *IEEE Journal on Selected Areas in Communications, Special Issue on Internet QoS 18* (December 2000), 2651–2664.
- [4] BRESLAU, L., CAO, P., FAN, L., PHILLIPS, G., AND SHENKER, S. Web caching and Zipf-like distributions: Evidence and implications. In *Proc. IEEE INFOCOM '99* (New York, 1999), pp. 126–134.
- [5] ÇINLAR, E. *Introduction to Stochastic Processes*. Prentice Hall, 1975.
- [6] CLÉVENOT, F., NAIN, P., AND ROSS, K. W. Stochastic fluid models for cache clusters. Tech. rep., INRIA, Sophia Antipolis, 2003.
- [7] COHEN, E., AND KAPLAN, H. The age penalty and its effect on cache performance. In *Proc. 3rd USENIX Symposium on Internet Technologies and Systems (USITS)* (San Francisco, California, 2001), pp. 73–84.
- [8] DUSKA, B., MARWOOD, D., AND FEELEY, M. The measured access characteristics of World Wide Web client proxy caches. In *Proc. USENIX Symp. on Internet Technologies and Systems* (Monterey, California, 1997), pp. 23–35.
- [9] ELWALID, A., AND MITRA, D. Fluid models for the analysis and design of statistical multiplexing with loss priorities on multiple classes of bursty traffic. In *Proc. IEEE INFOCOM '92* (Florence, Italy, 1992), pp. 415–425.
- [10] ELWALID, A., AND MITRA, D. Effective bandwidth of general markovian traffic sources and admission control of high speed networks. *IEEE/ACM Transactions on Networking* 1 (June 1993), 329–343.
- [11] GRIBBLE, S., AND BREWER, E. System design issues for internet middleware services: Deductions from a large client trace. In *Proc. USENIX Symp. on Internet Technologies and Systems* (Monterey, California, 1997), pp. 207–218.
- [12] GRIMMETT, G., AND STIRZAKER, D. *Probability and Random Processes*. Oxford University Press, 1992.
- [13] HORN, R. A., AND JOHNSON, C. R. *Matrix Analysis*. Cambridge University Press, 1985.
- [14] IYER, S., ROWSTRON, A., AND DRUSCHEL, P. Squirrel: A decentralized, peer-to-peer Web cache. In *Proc. of ACM Symposium on Principles of Distributed Computing (PODC 2002)* (Monterey, California, 2002), pp. 213–222.
- [15] KELLY, F. P. *Reversibility and Stochastic Networks*. Wiley, Chichester, 1979.
- [16] KROEGER, T., LONG, D. D. E., AND MOGUL, J. C. Exploring the bounds of web latency reduction from caching and prefetching. In *Proc. USENIX Symp. on Internet Technologies and Systems* (Monterey, California, 1997), pp. 13–22.

- [17] KUMARAN, K., AND MANDJES, M. Multiplexing regulated traffic streams: design and performance. In *Proc. IEEE INFOCOM '01* (Anchorage, 2001), pp. 527–536.
- [18] LOPRESTI, F., ZHANG, Z., TOWSLEY, D., AND KUROSE, J. Source time scale and optimal buffer/bandwidth trade-off for regulated traffic in an ATM node. In *Proc. IEEE INFOCOM '97* (Kobe, Japan, 1997), pp. 676–683.
- [19] REISSLEIN, M., ROSS, K. W., AND RAJAGOPAL, S. A framework for guaranteeing statistical QoS. *IEEE/ACM Transactions on Networking* 10, 1 (February 2002), 27–42.
- [20] ROWSTRON, A., AND DRUSCHEL, P. Pastry: Scalable distributed object location and routing for large-scale peer-to-peer systems. In *Proc. Int. Conf. on Distributed Systems Platforms (Middleware)* (Heideberger, Germany, November 2001).
- [21] WOLMAN, A., VOELKER, G., SHARMA, N., CARDWELL, N., KARLIN, A., AND LEVY, H. On the scale and performance of cooperative Web proxy caching. In *17th ACM Symp. on Operating Systems Principles (SOSP '99)* (Kiawah Island, South Carolina, 1999), pp. 16–31.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Overview of Squirrel</b>	<b>4</b>
<b>3</b>	<b>A Model for Squirrel</b>	<b>5</b>
3.1	Modeling the client dynamics . . . . .	5
3.2	Modeling the content dynamics . . . . .	5
<b>4</b>	<b>Performance analysis of the Squirrel p2p cache system</b>	<b>7</b>
4.1	Hit probability analysis . . . . .	8
4.2	Latency reduction . . . . .	11
4.3	Zipf popularity . . . . .	12
4.4	Discussion and extensions . . . . .	13
<b>5</b>	<b>Qualitative insight in the Squirrel system</b>	<b>14</b>
<b>6</b>	<b>Experimental results</b>	<b>16</b>
<b>7</b>	<b>Conclusion and perspectives</b>	<b>17</b>
<b>A</b>	<b>Proof of equation (11)</b>	<b>18</b>
<b>B</b>	<b>Proof of equation (15)</b>	<b>19</b>

<b>C Stationary distribution of the Engset model at jump times</b>	<b>20</b>
<b>D Uniqueness of the solution of (10)</b>	<b>20</b>



---

Unité de recherche INRIA Sophia Antipolis  
2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes  
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399