



UNIVERSITÀ DEGLI STUDI DI TRENTO

DEPARTMENT OF INFORMATION AND COMMUNICATION TECHNOLOGY

38050 Povo — Trento (Italy), Via Sommarive 14
<http://dit.unitn.it/>

STOCHASTIC GRAPH PROCESSES FOR PER-
FORMANCE EVALUATION OF CONTENT DE-
LIVERY APPLICATIONS IN OVERLAY NET-
WORKS

D. Carra, R. Lo Cigno and E.W. Biersack

February 2006 – Ver. 1.0

Technical Report # DIT-06-013

Stochastic Graph Processes for Performance Evaluation of Content Delivery Applications in Overlay Networks

Damiano Carra

Dip. di Informatica e Telecomunicazioni
Università di Trento, Trento, Italy
carra@dit.unitn.it

Renato Lo Cigno

Dip. di Informatica e Telecomunicazioni
Università di Trento, Trento, Italy
locigno@dit.unitn.it

Ernst W. Biersack

Institut EURECOM
Sophia Antipolis, France
erbi@eurecom.fr

Abstract— We consider the problem of distributing a content of finite size to a group of users connected through an overlay network that is built by a peer-to-peer application. The goal is the fastest possible diffusion of the content until it reaches all the peers. Applications like Bit-Torrent or SplitStream are examples where the problem we study is of great interest.

In order to represent the content diffusion process, we model the system as a stochastic graph process and define the constraints the graph evolution is subject to. The evolution of the graph is a semi-Markov process where the sojourn times are the rewards of interest for the computation of the time needed to complete the file distribution. We discuss the general properties of the constrained stochastic graphs and we show preliminary results obtained with an ad-hoc Monte-Carlo technique.

I. INTRODUCTION

The peer-to-peer (P2P) networking paradigm received a lot of attention in recent years. P2P systems construct an overlay at the application layer and do not require any modification to the existing Internet. Such an ease of deployment, which is in contrast to other technologies such as IPv6 or IP-level multicast that do require modifications inside the network, makes P2P systems very attractive for supporting new communications paradigms such as ‘application level multicast’ or ‘distributed publish/subscribe’.

One of the most popular P2P application is file sharing, a variation of which can also be used for file distribution. P2P for file distribution has the appealing feature of self-scaling: the more are the users downloading the content the larger the overall amount of resources (memory, bandwidth, CPU power, ...) available for the entire system.

In this paper we consider the specific problem of how to distribute a file to a community of users organized as an overlay of peers: assuming that the content is time-critical, which is the most efficient architecture and protocol that can be used to distribute this content to the users? Applications include, for instance, distribution of virus footprints or software updates, but also instant messaging to large communities. Other ‘applications’ may include the spreading of malware (viruses, worms, bugs, ...), but also streaming and conferencing systems; however in this work we restrict the analysis

to files distributed to cooperative users willing to receive and forward it to others.

We give a formal definition of the process underlying the construction of a ‘distribution graph’ as a semi-Markov process, describing how different choices impact the structure of the stochastic process itself (and obviously the constructed graph) as well as the rewards used to derive the performances. Then, we analyze the properties of the semi-Markov processes. Finally, we discuss some results. Representation of content delivery overlay networks through stochastic graph processes allows us to give a high level description of different kind of protocols and architectures and to compare them, without focusing on the implementation details.

A. Related Work

Performance analysis in terms of the minimum time required to distribute a file using a P2P system has only in recent years received some attention. Most of the analytical work [1][2][3] focuses on a specific system and not on a generic distribution architecture, and they assume strict hypotheses, like the uniformity of access bandwidths. Only [4] tackles the problem of bandwidths heterogeneity, which is also treated in this paper.

To the best of our knowledge, very few models have been proposed that allow comparative studies of different distribution architectures. In [5], inspired by SplitStream [13] (an overlay streaming protocol) the authors have defined and analyzed linear chain and tree-based architectures assuming ideal conditions, hence in a completely deterministic situation. The work in [6] defines a model for chain-based and tree-based architectures and analyzes the system using max-plus algebra considering an infinite number of packets, calculating the long term average throughput; our analysis instead considers a finite file size and calculates the distribution of the download time of all the peers.

Stochastic graph processes, the analytical tool we use in this paper to model overlay content delivery networks, were defined in [7] with the same notation we use here, although they were known since the ‘50s. A sub case, the *random graphs* [8], were studied in detail obtaining their general properties.

The focus of the analysis in [8] is the topological properties of random graphs, whereas our aim is to take into account not only connections among nodes, but also their weights given by the bandwidths of the involved nodes (this concept is clarified in Sect. II), which give rise to the state reward structure that allows the computation of completion times. Moreover, content distribution is done building a distribution graph on top of the overlay graph (see Fig. 1) and these two levels have different properties.

II. PROBLEM FORMULATION

The aim of the service is the delivery of a given finite size content \mathcal{F} to a set of users. The only requirement of the service is content integrity and the main performance metrics are the *download time* T of the content, either for a given user i (T_i), or for the whole community (T_t), or the mean \bar{T} of all the individual download times T_i .

We assume that each node knows a subset of the whole community, i.e., a node has a finite number of neighbors.

The content \mathcal{F} is divided in C pieces called *chunks*. A chunk represents a basic unit of transmission that can be distributed independently.

Each node is characterized by its bandwidth, that can be either symmetric or asymmetric. When two nodes start exchanging a chunk of \mathcal{F} , the *effective* bandwidth used to transfer the chunk depends on the bandwidth of both nodes, or on the percentage of the bandwidth that each node wants to use for the communication (each node can simultaneously be involved in more than one exchange). Besides the bandwidth used to transfer a single chunk, we have to consider the *rate* at which chunks arrive at a node, i.e., the inter-arrival time between chunks. The rate is defined as the amount of data contained in each chunk divided by the inter-arrival time. The distribution throughput is obviously related to the rate and not to the bandwidth, and the transfer rate between two nodes may be influenced by the whole transfer path from the source to the nodes considered.

A. Formalization and general definitions

The distribution of a content within a community of users can be formalized as the propagation across a graph of nodes and edges with some (stochastically defined) characteristics. Nodes are the users and edges summarize all the characteristics of the communication network between the users.

Let \mathcal{N} be the set of nodes, i.e., the vertices of the graph, and \mathcal{A} the set of all the arcs that connect pairs of nodes, $\mathcal{A} \subseteq \mathcal{N} \times \mathcal{N}$. We only consider fully reachable networks with bidirectional connections, so that \mathcal{A} can be represented by an irreducible, symmetric adjacency matrix. \mathcal{B}_i is the set of neighbors of user i , i.e., all those nodes in \mathcal{N} that are known and directly reachable from node i , with $\bigcup_{i \in \mathcal{N}} \mathcal{B}_i = \mathcal{N}$, since the network is fully reachable. \mathcal{B}_i is represented also by row i of the adjacency matrix \mathcal{A} .

The graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ represents the overlay network created, for instance, by a P2P network (see Fig. 1). In general,

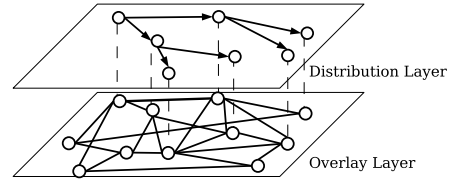


Fig. 1. Overlay and distribution graphs

$\mathcal{G}(\mathcal{N}, \mathcal{A})$ is time varying, i.e., nodes and edges can change in time, even appear or disappear. The overlay layer is the basis used to form the *distribution graph*. We define the distribution graph $\mathcal{G}^*(\mathcal{N}, \mathcal{E})$ as a directed subgraph of $\mathcal{G}(\mathcal{N}, \mathcal{A})$, with $\mathcal{E} \subseteq \mathcal{A}$. \mathcal{G}^* is a directed graph, since, from the content distribution point of view, the content propagates from the source to the destinations. How to obtain \mathcal{G}^* from \mathcal{G} is given by the rules implemented in the specific content distribution protocol. In general, we can assume that the distribution graph \mathcal{G}^* is built step by step. The building process is a stochastic process and it can be modeled as a chain. Let \mathcal{N}_n^* be the set of nodes that belong to the distribution graph at step n , and $\overline{\mathcal{N}_n^*}$ its complement with respect to \mathcal{N} . The distribution graph \mathcal{G}_{n+1}^* at step $n+1$ is obtained from \mathcal{G}_n^* by adding a new edge $\in \mathcal{A}$ from one node in \mathcal{N}_n^* to one in $\overline{\mathcal{N}_n^*}$. The complete distribution graph $\mathcal{G}^*(\mathcal{N}, \mathcal{E})$ is obtained when $\mathcal{N}_n^* = \mathcal{N}$, and $\overline{\mathcal{N}_n^*}$ is the empty set. The dynamic behavior of the distribution graph can be modeled as a *stochastic graph process*. We recall here the general definition of stochastic graph processes [7], while in Sect. III, we specialize them for the analysis of content distribution.

Definition 2.1: A stochastic graph process (SGP) on a node set \mathcal{N} is a discrete time Markov chain (DTMC) whose states are graphs on \mathcal{N} .

Even if not stated in the definition, nodes can be connected only through edges that belong to \mathcal{A} (in the next definitions, we state explicitly this dependence). Adhering to the definition given in [7], the focus is the building process and the SGP evolution implies that the graph is built step by step by adding one node and one edge at each step.

Notice that in content distribution, the distribution network (or graph) is naturally built step by step, so using the graph $\mathcal{G}^*(\mathcal{N}, \mathcal{E})$ is very appropriate. The time between two steps depends on the sojourn time of the state. If sojourn times are exponentially distributed, then we obtain a continuous time Markov chain (CTMC), but, in general, this assumption is not true and in continuous time we have a semi-Markov chain.

Definition 2.2: A constrained stochastic graph process (CSGP) on a graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ is a semi-Markov chain whose states are subgraphs on \mathcal{G} .

From the semi-Markov chain we can derive an embedded DTMC by sampling the process exactly at transition instants.

A CSGP is a fairly precise representation of the process of distributing a content. In the distribution of \mathcal{F} , a node that has started to download \mathcal{F} can in turn start to upload it after entirely receiving the first chunk. We define the *eligibility time* t_i^{el} of node i as the time at which node i can start to upload to

other nodes. t_i^{el} are random variables, since they depend on the node bandwidths. Given \mathcal{G}_n^* the next state \mathcal{G}_{n+1}^* depends only on the eligibility times of nodes in \mathcal{N}_n^* , so that the transition probabilities can be easily defined step-by-step.

If node j is son of node i in \mathcal{G}_n^* , then $t_j^{\text{el}} = t_i^{\text{el}} + t_{ij}$ where t_{ij} depends only on the bandwidth of nodes i and j . We assume that the bandwidth of a node is a random variable with a known probability density function (pdf) and node i chooses its sons j uniformly at random among all its neighbors (the neighbor set \mathcal{B}_i , is defined by the connectivity of \mathcal{G}).

Eligibility times t_j^{el} influence the semi-Markov process in two different ways. In the general case of randomly varying t_{ij} , they define both the transition probabilities between states and the states dwelling times. In the particular case of deterministic t_{ij} (e.g., when the bandwidth is only determined by access links), dwelling times are deterministic, and t_j^{el} define only the state transition probabilities. Notice, however, that the file distribution is entirely described by the embedded DTMC, so that only transition probabilities are important.

The DTMC that describes a CSGP is a transient chain with adsorbing states $\in \mathcal{G}^*(\mathcal{N}, \mathcal{E})$ that are reached when $\mathcal{N}_n^* = \mathcal{N}$.

The way we defined a CSGP implies that nodes are stable and collaborative, and that the networking infrastructure is reliable enough to allow edge stability. Clearly there is the possibility of extending the analysis to cases where nodes (or edges) can disappear during the distribution process, so that \mathcal{G}_n^* is derived from \mathcal{G}_{n-1}^* not only by adding an edge and a node, but also by removing one node and all the edges relative to it (that may include several generations of sons).

A well known class of CSGP are the *random graph process* studied back in the '50s by Erdős and Renyi in [8]. In random graph processes edges are added choosing uniformly at random, a property that does not capture the 'propagating' nature of content distribution. In the following, we specialize CSGP by adding constraints that make them suitable for modeling our problem.

III. CONTENT-DELIVERY CSGP

Different distribution architectures can be defined as special cases of stochastic graph processes with additional constraints. Before introducing the 'Content-Delivery constrained stochastic graph processes,' or CD-CSGP, we give some additional definitions that will simplify the characterization of each CD-CSGP.

A. Content-Delivery Related Definitions

For each node i , we define as b_i the node bandwidth, with a known pdf; we refer to b_i^u and b_i^d as the upload and download bandwidth of node i respectively. b_i^u and b_i^d can be correlated, e.g., $b_i^u + b_i^d = k$ constant, as in a LAN based access. When downloading the content, each node receives an effective rate r_i that depends on multiple factors. Moreover, each node has a constraint on maximum and minimum number of active uploads (the outdegree of the node): k_i^{max} and k_i^{min} .

Definition 3.1 (saturated node): A node $i \in \mathcal{N}_n^*$ is called *saturated* if

- it has k_i^{max} outgoing edges that belong to \mathcal{G}_n^* or
- fully uses b_i^u to transmit chunks to neighboring nodes that belong to \mathcal{G}_n^* .

Definition 3.2 (interior subset): The subset of nodes $\in \mathcal{N}_n^*$ that at step n are saturated is called \mathcal{I}_n , the interior node subset at step n .

Definition 3.3 (leaf subset): The set of nodes $\in \mathcal{N}_n^*$ that are not interior nodes is called the leaf node subset \mathcal{L}_n at step n , with $\mathcal{L}_n = \mathcal{N}_n^* - \mathcal{I}_n$.

We consider a single node as a root of the stochastic graph. We define a distance measure based on number of hops from the root to any node i .

Definition 3.4 (step distance): The number of hops from the root to a node i following the shortest path is called *step distance* or *step depth*, $d^{(i)}$.

In a tree, $\max_i(d^{(i)})$ is the tree depth.

B. Unbalanced and Uneven Trees

The general process that leads to a tree-based distribution structures must abide to the following rules.

CD-CSGP 1: A constrained stochastic graph process on graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ is called tree-based if

- 1) \mathcal{G}_0^* is a node, called root, randomly chosen in \mathcal{N} .
- 2) \mathcal{G}_n^* is obtained from \mathcal{G}_{n-1}^* by
 - a) choosing the node i from \mathcal{N}_{n-1}^* with the smallest eligibility time: $t_i^{\text{el}} = \min_j(t_j^{\text{el}})$; if more nodes have the same eligibility time, the choice among these nodes is made randomly;
 - b) adding edges from node i to nodes randomly chosen from $\mathcal{B}_i \cap \overline{\mathcal{N}_{n-1}^*}$, until node i becomes saturated.

Figure 2 shows an example with few states of the DTMC generated by a CD-CSGP 1 process. In this case we have only two possible bandwidths (slow nodes with black circles, fast nodes with white ones, with slow bandwidth less than half of the fast bandwidth) and $k_i^{\text{max}} = 2$ and $k_i^{\text{min}} = 1$. Starting, for instance, from a state where the server is uploading to a slow and to a fast node, the fast node has the smallest eligibility time and there are only three next possible states: (i) the fast node selects a fast nodes among its neighbors, becoming saturated; alternatively, the fast node chooses a slow node so it has to select another node: (ii) the selected node is fast and we have bandwidth saturation, or (iii) the chosen node is slow and we have saturation because k_i^{max} is reached. Note that in case (i) the node becomes saturated since the rate of the content it is receiving is high. If, for instance, the rate were slow (consider the fast node under the slow node in the shadowed state), the number of children would be in any case 2, since the bandwidth used to transmit to each child is at most equal to the rate it is receiving.

The resulting tree is, in general, a structure where the nodes in the leaf set \mathcal{L}_n do not all have the same step distance from the root. The speed of growth of the different branches is not the same. And the deeper branches are those that contain faster nodes, i.e., nodes with smaller eligibility times t^{el} . We call such trees "uneven."

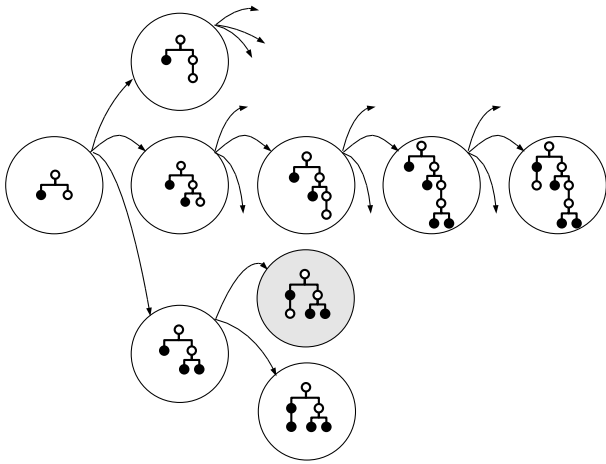


Fig. 2. Sample of the embedded DTMC for a CD-CSGP 1 process; states are graphs built on \mathcal{G} , black and white circles represent slow and fast nodes respectively, $k_i^{\max} = 2$ and $k_i^{\min} = 1$.

There is a vast literature on tree-based distribution architectures and their performance. Nevertheless, the majority of the proposals consider trees where leaves have the same distance from the root. We call such trees “*unbalanced*”¹. The difference between unbalanced and uneven trees is enormous: in an unbalanced tree, a slow node will influence the reception of all nodes in its subtree, in an uneven tree, a slow node may not even have the possibility to have children. Since we are interested in the download time, it is worth to look at a weighted graph where the weight associated to a directed edge is given by the difference between the download times of the nodes connected by the edge. Considering unbalanced trees, this representation shows the disparity in terms of download time among leaf nodes that are at the same step distance. In Fig. 3 the weight is represented as a difference in edge length. Conversely, in uneven trees, leaf nodes are at different step distances and the weighted graph gives a pictorial illustration why the tree grows in this way: a new edge is added only after a node becomes eligible and this forces a uniform growth of the *weighted* graph.

Since most of the proposed protocols use a tree architecture with an unbalanced tree, we define a special stochastic graph process for this type of tree. To do so, we consider a subset of the leaf set \mathcal{L}_n .

Definition 3.5: Let $d_n^{\max} = \max_j(d_n^{(j)})$ be the maximum step

distance of the nodes $j \in \mathcal{L}_n$. The subset $\widetilde{\mathcal{L}}_n \subseteq \mathcal{L}_n$ is defined as follows: $\widetilde{\mathcal{L}}_n = \{i \in \mathcal{L}_n \mid d_n^{(i)} < d_n^{\max}\}$.

Now we can define the process that leads to unbalanced trees.

CD-CSGP 2: A constrained stochastic graph process on graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ is called tree-based and unbalanced if

- 1) \mathcal{G}_0^* is a node, called root, randomly chosen in \mathcal{N} .
- 2) \mathcal{G}_n^* is obtained from \mathcal{G}_{n-1}^*

¹A balanced tree is a full tree having nodes all with the same outdegree (a full tree is a tree with leaf nodes at the same step distance), thus a variable outdegree implies an unbalanced structure.

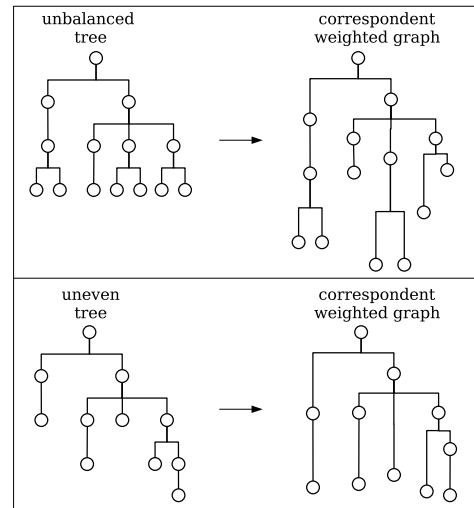


Fig. 3. Difference between unbalanced and uneven trees, considering the corresponding weighted graphs where edge length represents the download time

- a) choosing a node i from $\widetilde{\mathcal{L}}_{n-1}$ if not empty, otherwise from \mathcal{L}_{n-1} , with the smallest eligibility time; if more nodes have the same eligibility time, the choice among these nodes is made randomly;
- b) adding edges from node i to nodes randomly chosen from $\mathcal{B}_i \cap \overline{\mathcal{N}_{n-1}^*}$, until the node becomes saturated.

We will show in Sect. V the impact of unbalanced and uneven tree on the performance. Finally, we note that CD-CSGP 2 is a CD-CSGP 1 with additional constraints.

C. General Mesh Architecture

Tree based architectures allow the content to rapidly diffuse to nodes, but have known shortcomings. Each node has only one ancestor and in case of a nodes failure, the entire subtree will stop receiving data. Each node must divide the upload bandwidth among its children, so children use only a fraction of their download bandwidth for receiving chunks; if we consider the case of asymmetric capacities, where the upload bandwidth is smaller than the download bandwidth (as in the case of ADSL), the percentage of unused download bandwidth increases even further. For a bounded number of nodes, there are nodes that have received the entire file without uploading a single chunk, resulting in unfairness and poor performance.

Mesh based architectures are meant to overcome these problems. In addition, if the direction of the content diffusion within the structure is random (the diffusion has not a particular pattern) the number of leaf nodes can be made minimal (can be reduced to a single node with a perfect knowledge of the network [5]).²

Nevertheless, the mesh architecture introduces a new problem: the chunk selection strategy. A node can help another

²All this reasoning is valid if \mathcal{N} is bounded; if \mathcal{N} is unbounded, the probability to find a node already involved in the distribution process is negligible and only the tree-based architectures makes sense.

node if it has parts of \mathcal{F} not yet received by the helped node, i.e., if it has ‘fresh’ information. We are not concerned here on how freshness is checked and/or imposed (see for instance [9], [10] for works on the topic). We assume an ideal situation where if a node has received only part of \mathcal{F} and it is contacted by another node that is not already its ancestor, then all the information it can provide is either completely fresh or completely stale. Any impairment can be easily taken into account with a probabilistic approach.

Allowing the generation of mesh topologies means that a node i already included in \mathcal{G}_n^* may be contacted other nodes j , also in \mathcal{G}_n^* to receive parts of \mathcal{F} it does not yield already. If j only has information that is not fresh for i , then the ‘delivery connection’ is not established.

Let’s analyze the delivery process starting from the root, or content generator. The root can upload \mathcal{F} to its k children in exactly the same order³. If this is the case, \mathcal{F} is distributed by any successive node maintaining exactly the same order, then whenever a node j in \mathcal{G}_n^* contacts a node i also in \mathcal{G}_n^* , then the information it offers is necessarily stale, since contacts happen at the eligibility time t_j^d , that is when j has received only the first chunk.

The ‘delivery connection’ between j and i is not established and the distribution topology evolves like a single *diffusion tree* as depicted in Fig. 4a.

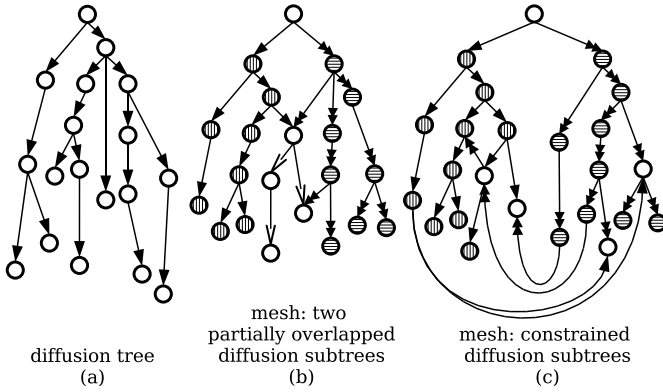


Fig. 4. A diffusion tree (a); mesh topologies obtained overlapping diffusion subtrees without constraints (b) and with constraints (c)

If, on the other hand, the root node uploads chunks to its k children in different orders, whenever the subtrees generated by the children merge, there are high chances that the two ancestors of the node where the merging happens have (at least initially) disjoint contents. This situation is represented in Fig. 4b, where the different order of the content is represented by different arrows on the edges. Each of the children generates its own *diffusion subtree* \mathcal{FG} . Subtrees can overlap since during the content diffusion a node can add an already selected node provided that they do not belong to the same subtree \mathcal{FG} .

³We use the term ‘order’ here in a very broad sense, including network coding techniques as well as simple chunk shuffling, but these details do not impact the performance of distribution architectures.

Under these conditions, we can define the process leading to a mesh-based distribution architecture.

CD-CSGP 3: A constrained stochastic graph process on graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ is called mesh-based if

- 1) \mathcal{G}_0^* is a node, called root, randomly chosen in \mathcal{N} .
- 2) \mathcal{G}_1^* consists of the root node and the k nodes randomly chosen in $\mathcal{N} - \{\text{root}\}$ called the first generation nodes, where and each of them will generate a subtree \mathcal{FG} .
- 3) \mathcal{G}_n^* is obtained from \mathcal{G}_{n-1}^*
 - a) choosing a node i from \mathcal{N}_{n-1}^* with the smallest eligibility time; if more nodes have the same eligibility time, the choice among these nodes is made randomly; \mathcal{FG}_i denotes the subtree to which i belongs;
 - b) adding edges from node i to nodes randomly chosen from $\mathcal{B}_i \cap \mathcal{N}_{n-1}^\dagger$ until the node becomes saturated, where $\mathcal{N}_{n-1}^\dagger$ is the set of nodes in the subtree \mathcal{FG}_i at step $n - 1$.

Meshes generated by CD-CSGP 3 can cause inefficiencies of the distribution process. The intuitive explanation is that in overlapping diffusion trees the ‘leaf capacity problem’ is not removed: interior nodes upload k times the amount they receive (k is the average outdegree), whereas leaf nodes do not upload at all. If we allow diffusion subtrees to overlap without constraints, the whole content will diffuse creating interior and leaf nodes, something that we would like to avoid.

For this purpose, we build diffusion trees allowing only the addition of nodes not yet reached by any chunk (untouched nodes). Only when no more untouched nodes are available, the process can start adding nodes belonging to different diffusion subtrees. With this procedure, leaf nodes of a diffusion subtree start uploading to other nodes and we obtain better fairness and improved performance. Moreover, the diffusion direction is ‘reverted’ at leaves, ensuring a better spreading of the content. This constrained mesh diffusion process is depicted in Fig. 4c, and we can formally define the CSGP that leads to these architectures. As noted above, these processes can be defined only for a bounded set of nodes.

CD-CSGP 4: A constrained stochastic graph process on graph $\mathcal{G}(\mathcal{N}, \mathcal{A})$ is called constrained mesh-based if

- 1) \mathcal{G}_0^* is a node, called root, randomly chosen in \mathcal{N} .
- 2) \mathcal{G}_1^* consists of the root node and the k nodes randomly chosen in $\mathcal{N} - \{\text{root}\}$ called the first generation nodes, where and each of them will generate a subtree \mathcal{FG} .
- 3) \mathcal{G}_n^* is obtained from \mathcal{G}_{n-1}^*
 - a) choosing a node i from \mathcal{N}_{n-1}^* with the smallest eligibility time; if more nodes have the same eligibility time, the choice among these nodes is made randomly; \mathcal{FG}_i denotes the subtree to which i belongs;
 - b) adding edges from node i to nodes randomly chosen from $\mathcal{B}_i \cap \mathcal{N}_{n-1}^*$, until the node becomes saturated and \mathcal{N}_{n-1}^* is not empty;
 - c) if \mathcal{N}_{n-1}^* is empty and node i is not saturated, adding edges from node i to nodes randomly

chosen from $\mathcal{B}_i \cap \overline{\mathcal{N}_{n-1}^\dagger}$ until the node becomes saturated, where $\mathcal{N}_{n-1}^\dagger$ is the set of nodes in the subtree $\mathcal{F}\mathcal{G}_i$ at step $n - 1$.

Two sub-cases of the above defined construction process were proposed under the name *SplitStream* [13] and *PTree* [5]: both of the schemes use a fixed outdegree k and partition the file in exactly k stripes, where each stripe is distributed along one of the k diffusion trees. The process we define is more general since we do not impose any fixed outdegree and allow that nodes upload the whole file, however in different order.

IV. CD-CSGP SOLUTION

The CD-CSGP we have defined in Sect. III can describe different behaviors of content distribution protocols and algorithms. Consider for example the two well known application layer multicast protocols, ALMI [12] and SplitStream [13], which use different topologies for the distribution. The first builds a tree structure to deliver the content, the second a mesh structure. Through CD-CSGP we can compare these two approaches, abstracting from any protocol detail.

A. Properties of CSGP and rewards

Let \mathbf{S} be the state space of the DTMC embedded in the CD-CSGP we consider; $S_k \in \mathbf{S}$ are the states of the DTMC, i.e., the graphs \mathcal{G}_n^* . We start considering trees: to compute the mean download time \bar{T} of \mathcal{F} we assign to each absorbing state $S_k \in \mathbf{S}_a$ (\mathbf{S}_a is the set of absorbing states) a reward \bar{T}_k equal to the mean download time of the nodes in the state,

$$\bar{T}_k = \frac{1}{|\mathcal{N}_n^*|} \sum_{i \in \mathcal{N}} T_i ; S_k \in \mathbf{S}_a$$

The mean download time \bar{T} is the reward of a DTMC obtained by adding a deterministic transition from all the absorbing states to an initial state represented by the empty graph \emptyset .

$$\bar{T} = \frac{\sum_{k \in \mathbf{S}_a} \bar{T}_k \pi_k}{\sum_{k \in \mathbf{S}_a} \pi_k} .$$

where π_k s are the steady state probability of the support DTMC defined above.

Another performance measure easily defined as a reward is the wasted upload bandwidth \bar{w}^u (in percentage). Let $w_i^u = 100(1 - \frac{\max(r_i^u)}{b_i^u})$ be the wasted upload bandwidth of node i , where $\max(r_i^u)$ is the maximum upload rate ever reached by the node in any visited state. Considering again the modified DTMC and letting \mathbf{S}_a be the set of absorbing states in the unmodified DTMC we have

$$\bar{w}_k^u = \frac{1}{|\mathcal{N}|} \sum_{i \in \mathcal{N}} w_i^u ; S_k \in \mathbf{S}_a$$

and

$$\bar{w}^u = \frac{\sum_{k \in \mathbf{S}_a} \bar{w}_k^u \pi_k}{\sum_{k \in \mathbf{S}_a} \pi_k} .$$

We developed a tool for the numerical solution of all the CD-CSGP we defined in this work. The tool is based on

fast Monte Carlo techniques that allow the exploration of the chain with great efficiency, computing in the meanwhile the reliability of the produced results. It implements the algorithms driving the stochastic process providing realizations of the process. It has several configuration parameters and the outputs are the rewards associated with the process. For a detailed description of the tool we refer the interested reader to [15]. Here we only give some sample results showing the feasibility and power of CD-CSGP formalization for content delivery analysis.

V. NUMERICAL RESULTS

As numerical example we consider a density function for the node bandwidth taken from [16], summarized in Table I.

TABLE I
BANDWIDTH DISTRIBUTION USED IN THE EXAMPLES

Bandwidth	% nodes
56 kbit/s	13%
640 kbit/s	23%
1.2 Mbit/s	64%

When reporting results, we normalize the data such that $\frac{|\mathcal{F}|}{\min(b_i)} = 1$ ‘round’, where $|\mathcal{F}|$ is the content size in bits and $\min(b_i)$ is the minimum bandwidth of the input pdf in bits/s. We use a number of chunks C equal to 100, but a sensitivity analysis with different values of C indicates a qualitative behavior independent of C , as long as $C \gg 1$. All results have confidence level 0.95 and confidence interval $\pm 10\%$.

We focus on the comparison between unbalanced and uneven trees; a detailed analysis of the influence of the different constraints, such as k^{\max} and k^{\min} , can be found in [15].

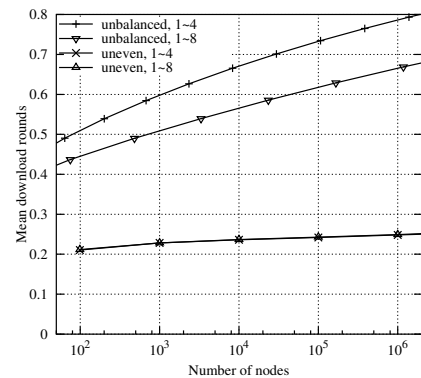


Fig. 5. \bar{T} for unbalanced and uneven trees as a function of $|\mathcal{N}|$; $k^{\min} = 1$

Figure 5 shows the results as a function of $|\mathcal{N}|$. The poorer results for unbalanced trees (CD-CSGP 2) are due to the bounds on the $d^{(i)}$. Slow nodes, especially those close to the root, impose their rate on the whole subtree, independently of the bandwidth of the nodes in the subtree. In the case of uneven trees (CD-CSGP 1), slow nodes close to the root have no time to start to upload, since the time it takes to become eligible is much more than the time it takes the fast nodes to

reach, at different levels, all the other nodes. This increase of performance for the uneven tree comes at a cost of a greater step distance.

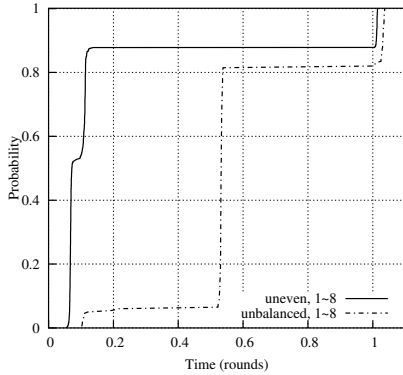


Fig. 6. CDF of T for unbalanced and uneven trees, with $|\mathcal{N}| = 10^4$.

As shown in Fig.6, our tool provides the entire CDF (or conversely the pdf) of the download time, a result that is normally not available with analytical or semi-analytical techniques. The CDF shows clearly the reason why uneven trees are more efficient and does also show that with uneven trees fast peers are far less influenced by slow ones also increasing the perceptual fairness of users.

Space forbids the presentation of a detailed set of results, that can be found in [15]. We only present a last synthetic result for mesh networks: the wasted upload bandwidth. This is a measure of how much of the node upload bandwidth is unused when the node is involved in the distribution process. Table II shows the percentage of wasted bandwidth when using a tree (uneven) and a mesh for different outdegrees and community size. The waste reduction in the mesh is due to an efficient use of upload bandwidth of the leaves.

TABLE II
COMPARISON OF TREES AND MESHES.

Outdegree	#Nodes	Levels	Upload Wasted Bandwidth	
			Uneven Tree	Mesh
1 - 8	10^5	21	46.9%	13.3%
1 - 8	10^6	24	47.5%	13.1%
2 - 8	10^5	14	66.2%	26.8%
2 - 8	10^6	18	68.9%	29.3%

VI. CONCLUSIONS

Trees have been studied intensively in the literature, however important details such as bandwidth heterogeneity and varying node outdegrees as well as different minimum and maximum outdegrees have received very little attention, probably for the difficulty in finding closed form results.

In this paper we have considered architectures for content distributions that leads to general tree topologies and mesh showing they can be described as Stochastic Graph Processes. Sample results demonstrates the power of the analytic framework we defined.

The approach based on stochastic graph processes is very attractive for the analysis of basic properties of distribution systems. To the best of our knowledge, stochastic graph processes were used only to study connectivity properties, but they were not applied in performance analysis of networks, while this work clearly demonstrated that properly adding constraints in the evolution of the stochastic graph can lead to a detailed and insightful description of the system.

ACKNOWLEDGMENT

This work has been partly supported by the European Union under the E-NEXT project FP6-506869.

REFERENCES

- [1] X. Yang and G. de Veciana, "Service Capacity of Peer-to-Peer Networks," in *Proc. IEEE INFOCOM 2004*, Hong Kong, Mar. 2004.
- [2] F. Clevenot and P. Nain, "A Simple Fluid Model for the Analysis of the Squirrel Peer-to-Peer Caching System," in *Proc. IEEE INFOCOM 2004*, Hong Kong, Mar. 2004.
- [3] D. Qiu and R. Srikant, "Modeling and Performance Analysis of BitTorrent-Like Peer-to-Peer Networks," in *Proc. ACM SIGCOMM 2004*, Portland, OR, Sept. 2004.
- [4] Z. Ge, D. R. Figueiredo, S. Jaiswal, J. Kurose, and D. Towsley, "Modeling Peer-Peer File Sharing Systems," in *Proc. IEEE INFOCOM 2003*, San Francisco, California, USA, Mar. 2003.
- [5] E. W. Biersack, P. Rodriguez, and P. Felber, "Performance Analysis of Peer-to-Peer Networks for File Distribution" in *Proc. 5th International Workshop on Quality of Future Internet Services (QofIS'04)*, Barcelona, Spain, Sept. 2004.
- [6] F. Baccelli, A. Chaintreau, Z. Liu, A. Riabov, S. Sahu "Scalability of Reliable Group Communication Using Overlays," in *Proc. IEEE INFOCOM 2004*, Hong Kong, Mar. 2004.
- [7] S. Nikolettseas, J. Reif, P. Spirakis and M. Young, "Stochastic Graphs Have Short Memory: Fully Dynamic Connectivity in Poly-Log Expected Time," in *Proc. of the 22nd ICALP*, pp. 159-170, 1995.
- [8] P. Erdős and A. Renyi, "On random graphs," *Publ. Math.* 6:290-297, 1959.
- [9] C. Gkantsidis, and P. Rodriguez, "Network Coding for Large Scale Content Distribution," in *Proc. IEEE INFOCOM 2005*, Miami, Mar. 2005.
- [10] D. Kostic, A. Rodriguez, J. Albrecht, and A. Vahdat, "Bullet: High Bandwidth Data Dissemination Using an Overlay Mesh," in *Proc. SOSP 2003*, Oct. 2003.
- [11] B. Cohen, "Incentives build robustness in BitTorrent," May 2003. Available: <http://www.bittorrent.com/documentation.html>
- [12] D. Pendarakis, S. Shi, D. Verma, and M. Waldvogel, "ALMI: An Application Level Multicast Infrastructure," in *Proc. of the 3rd Usenix Symposium on Internet Technologies & Systems (USITS)*, Mar. 2001.
- [13] M. Castro, P. Druschel, A. M. Kermarrec, A. Nandi, A. Rowstron, and A. Singh, "SplitStream: Highbandwidth Multicast in a Cooperative Environment," in *Proc. ACM Symposium on Operating Systems Principles (SOSP 03)*, The Sagamore, New York, USA, Oct. 2003.
- [14] D. Carra, R. Lo Cigno, and E. W. Biersack, "Stochastic Graph Processes for Performance Evaluation of Content Delivery Applications in Overlay Networks," Technical Report DIT-06-013, Univ. of Trento, Feb. 2006. Available: <http://www.dit.unitn.it/locigno/preprints/DIT-06-013.pdf>
- [15] D. Carra, R. Lo Cigno, and E. W. Biersack, "Fast Stochastic Exploration of P2P File Distribution Architectures," Technical Report DIT-06-014, Univ. of Trento, Feb. 2006. Available: <http://www.dit.unitn.it/locigno/preprints/DIT-06-014.pdf>
- [16] R. Gaeta, M. Griboaud, D. Manini, and M. Sereno, "Analysis of Resource Transfer in Peer-to-Peer File Sharing Applications using Fluid Models," *Performance Evaluation - Peer-to-Peer Computing Systems*, Vol. 63, Issue 3, Pages 147-264, Mar. 2006.