

# SplitStream: High-bandwidth content distribution in a cooperative environment

Miguel Castro<sup>1</sup>

Peter Druschel<sup>2</sup>

Anne-Marie Kermarrec<sup>1</sup>

Animesh Nandi<sup>2</sup>

Antony Rowstron<sup>1</sup>

Atul Singh<sup>2</sup>

<sup>1</sup>Microsoft Research, 7 J J Thomson Avenue, Cambridge, CB3 0FB, UK.

<sup>2</sup>Rice University, 6100 Main Street, MS-132, Houston, TX 77005, USA.

## Abstract

In tree-based multicast systems, a relatively small number of interior nodes carry the load of forwarding multicast messages. This works well when the interior nodes are dedicated infrastructure routers. But it poses a problem in cooperative end-system multicast, where participants expect to contribute resources proportional to the benefit they derive from using the system. Moreover, many participants may not have the network capacity and availability required of an interior node in high-bandwidth multicast applications. SplitStream is a high-bandwidth content distribution system based on end-system multicast. It distributes the forwarding load among all the participants, and is able to accommodate participating nodes with different bandwidth capacities. We sketch the design of SplitStream and present some preliminary performance results.

## 1 Introduction

End-system or application-level multicast [2, 11, 21, 7, 18, 14, 1] has become an attractive alternative to IP multicast. Instead of relying on a multicast infrastructure in the network (which is not widely available), the participating hosts pool their resources to route and distribute multicast messages using only unicast network services. In this paper, we are particularly concerned with application-level multicast in *cooperative* environments, where participants contribute resources in exchange for using the service. In such environments, participants expect that the forwarding load be shared among all members.

Unfortunately, conventional tree-based multicast is inherently not well matched to a cooperative environment. The reason is that in any efficient (i.e. low-depth) multicast tree, a small number of interior nodes carry the burden of splitting and forwarding multicast traffic, while a large number of leaf nodes contribute no resources. This conflicts with the expectation that all members should share the forwarding load. The problem is further aggravated in high-bandwidth applications like video or bulk file distribution, where many nodes may not even have the capacity and availability required of an interior node in a conventional multicast tree. SplitStream is designed to address

these problems.

SplitStream enables efficient cooperative distribution of high-bandwidth content, whilst distributing the forwarding load among the participating nodes. SplitStream can also accommodate nodes with different network capacities and with asymmetric bandwidth on the inbound and outbound network paths. Subject to these constraints, it balances the forwarding load across all the nodes.

The key idea is to *split* the multicast content into  $k$  stripes, and multicast each stripe in a separate multicast tree. Participants join as many trees as there are stripes they wish to receive. The aim is to construct this *forest* of multicast trees such that an interior node in one tree is a leaf node in all the remaining trees. In this way, the forwarding load can be spread across all participating nodes. We show that it is possible, for instance, to efficiently construct a forest in which the inbound and outbound bandwidth requirements of each node are the same, while maintaining low delay and link stress across the system.

SplitStream also offers improved robustness to node failure and sudden node departures. Since ideally, any given node is an interior node in only one tree, its failure can cause the temporary loss of at most one of the stripes. With appropriate data encodings such as erasure coding [3] of bulk data or multiple description coding (MDC) [13, 15] of streaming media, applications can thus mask or mitigate the effects of node failures, even while the affected tree is being repaired.

The key challenge in the design of SplitStream is to efficiently construct a forest of multicast trees that distributes the forwarding load, subject to the bandwidth constraints of the participating nodes, in a decentralized, scalable, and self-organizing manner. SplitStream relies on a structured peer-to-peer overlay network called Pastry [19], and on Scribe [7], an application-level multicast system built upon this overlay to construct and maintain these trees.

The rest of this paper is organized as follows. Section 2 outlines the SplitStream approach in more detail. A brief description of Pastry and Scribe is given in Section 3. We sketch the design of SplitStream in Section 4. Section 5 describes related work and Section 6 concludes.

## 2 The SplitStream approach

In this section, we give a more detailed overview of SplitStream’s approach to cooperative, high-bandwidth content distribution.

**Tree-based multicast** In all multicast systems based on a single tree, participating nodes are either interior nodes or leaf nodes. The interior nodes carry all the burden of forwarding multicast messages. In a  $k$ -level balanced tree with arity  $f$ , the number of interior nodes is  $\frac{f^{k+1}-1}{f-1}$  and the number of leaf nodes is  $f^k$ . Thus, the fraction of leaf nodes increases with  $f$ . For example, more than half of the nodes are leaves in a binary tree, and over 90% of nodes are leaves in a tree with arity 16. In the latter case, the forwarding load is carried by less than 10% of the nodes; whilst all nodes have equal inbound bandwidth, the internal nodes have an outbound bandwidth requirement of 16 times the inbound bandwidth. Even in a binary tree, which would be impractically deep in most circumstances, the outbound bandwidth required by the interior nodes is twice that of their inbound bandwidth.

**SplitStream** SplitStream is designed to overcome the inherently unbalanced forwarding load in conventional tree-based multicast systems. SplitStream strives to distribute the forwarding load over all participating nodes, and respects different capacity limits of individual participating nodes. SplitStream achieves this by splitting the multicast content into multiple stripes, and using separate multicast trees to distribute each stripe.

Figure 1 illustrates how SplitStream balances the forwarding load among the participating nodes. In this simple example, the original content is split into two stripes and multicast in separate trees. For simplicity, let us assume that the original content has a bandwidth requirement of  $B$ , and that each stripe has half the bandwidth requirement of the original content. Each node other than the source subscribes to both stripes, inducing an inbound bandwidth requirement of  $B$ . As shown in Figure 1 each node is an internal node in only one tree and forwards the stripe to two children, yielding an outbound bandwidth requirement of no more than  $B$ .

In general, the content is split into  $k$  stripes. Participating nodes may subscribe to a subset of the stripes, thus controlling their inbound bandwidth requirement in increments of  $B/k$ . Similarly, participating nodes may control their outbound bandwidth requirement in increments of  $B/k$  by limiting the number of children they adopt. Thus, SplitStream can accommodate nodes with different bandwidths, and nodes with unequal inbound and outbound network capacities. SplitStream is able to satisfy all participating nodes as long as the total number of stripes to which all nodes wish to subscribe does not exceed the total number of children that all nodes are willing to adopt.

Applications may choose any content encoding that produces stripes with even bandwidth requirements. In practice, applications may choose to encode content in a man-

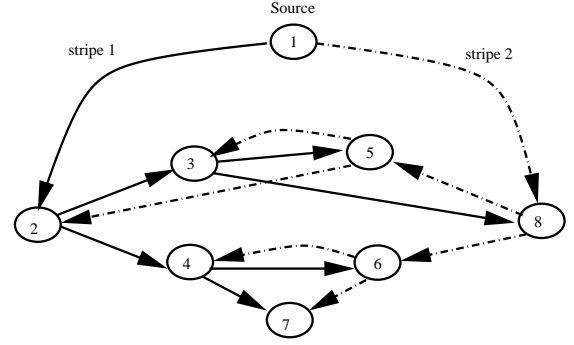


Figure 1: A simple example illustrating the basic approach of SplitStream. Original content is split into two stripes. An independent multicast tree is constructed for each stripe such that a node is an interior node in one multicast tree and a leaf in the other.

ner that requires bandwidth greater than  $B/k$  per stripe, in return for the ability to reconstitute the content from less than  $k$  stripes. For example, a media stream could be encoded using MDC so that the video can be reconstituted from any subset of the  $k$  stripes, with video quality proportional to the number of stripes received. Such an encoding also allows low-bandwidth clients to receive the video at lower quality.

As another example, erasure codes could be used to stripe file data, thus allowing the files to be reconstituted from any sufficiently large subset of the stripe blocks. For instance, a participant may subscribe to all stripes, but reconstitute the file as soon as a sufficient number of stripe blocks has arrived, discarding later arriving stripe blocks.

Using redundant content encodings also improves the resilience of the system to node failures or sudden departures of participants. Since a node failure affects at most one of the stripe trees, erasure codes can mask such failures. In the case of MDC encoded video, a node failure may at worst cause an intermittent drop in video quality while the affected tree is being repaired.

## 3 Background: Pastry and Scribe

In this section, we briefly sketch Pastry, a scalable, self-organizing, structured p2p overlay network, and Scribe, a scalable application-level multicast system based on Pastry. Both systems are key building blocks in the design of SplitStream.

**Pastry** In Pastry, nodes and objects are assigned random identifiers (called *nodeIds* and *keys*, respectively) from a large, sparse id space. Keys and nodeIds are 128 bits in length and can be thought of as a sequence of digits in base  $2^b$  ( $b$  is a configuration parameter with a typical value of 3 or 4). Given a message and a key, Pastry routes the message to the node with the nodeId that is numerically closest to the key, which is called the key’s *root*.

In order to route messages, each node maintains a rout-

ing table and a leaf set. A node’s routing table has about  $\log_{2^b} N$  rows and  $2^b$  columns. The entries in row  $n$  of the routing table refer to nodes whose nodeIds share the first  $n$  digits with the local node’s nodeId; the  $(n+1)$ th nodeId digit of a node in column  $m$  of row  $n$  equals  $m$ . The column in row  $n$  corresponding to the value of the  $(n+1)$ th digits of the local node’s nodeId remains empty. Routing in Pastry requires that at each routing step, a node normally forwards the message to a node whose nodeId shares with the key a prefix that is at least one digit longer than the prefix that the key shares with the present node’s id. If no such node is known, the message is forwarded to a node whose nodeId shares a prefix with the key as long as the current node, but is numerically closer to the key than the present node’s id.

Each Pastry node maintains a set of neighboring nodes in the nodeId space (called the leaf set), both to ensure reliable message delivery, and to store replicas of objects for fault tolerance. The expected number of routing hops is less than  $\log_{2^b} N$ . The Pastry overlay construction observes proximity in the underlying Internet. Each routing table entry is chosen to refer to a node with low network delay, among all nodes with an appropriate nodeId prefix. As a result, one can show that Pastry routes have a *low delay penalty*: the average delay of Pastry messages is only approximately twice the IP delay between source and destination [5]. Similarly, one can show the *local route convergence* of Pastry routes: the routes of messages route to the same key from nearby nodes tend to converge at a nearby intermediate node. Both of these properties are important for the construction of efficient multicast trees, described below. A full description of Pastry can be found in [19].

**Scribe** Scribe is an application-level multicast system built upon Pastry. A pseudo-random Pastry key, known as the *groupId*, is chosen for each multicast group. A multicast tree associated with the group is formed by the union of the Pastry routes from each group member to the *groupId*’s root (which is also the root of the multicast tree). Messages are multicast from the root to the members using reverse path forwarding [9].

The properties of the Pastry overlay ensure that the multicast trees are efficient. The delay to forward a message from the root to each group member is low due to the low delay penalty of Pastry routes. Pastry’s local route convergence ensures that the load imposed on the physical network is small because most message replication occurs at intermediate nodes that are close in the network to the leaf nodes in the tree.

Group membership management in Scribe is decentralized and highly efficient, because it leverages the existing, proximity-aware Pastry overlay. Adding a member to a group merely involves routing towards the *groupId* until the message reaches a member of the tree, followed by adding the route traversed by the message to the group multicast tree. As a result, Scribe can efficiently support large numbers of groups, arbitrary numbers of group mem-

bers, and groups with highly dynamic membership.

The latter property, combined with an anycast [6] primitive recently added to Scribe, can be used to perform distributed resource discovery. As we will show in the next section, SplitStream uses this mechanism to discover nodes with spare forwarding capacity. A full description and evaluation of Scribe multicast can be found in [7]. Scribe anycast is described in [6].

## 4 Building SplitStream

In this section, we sketch the design of SplitStream.

**Building independent trees** SplitStream uses a separate Scribe multicast tree for each of the  $k$  stripes. SplitStream exploits the properties of Pastry routing to ensure the desired independence. Recall that Pastry normally forwards a message towards nodes whose nodeIds share progressively longer prefixes with the message’s key. Since a Scribe tree is formed by the routes from all members to the *groupId*, the nodeIds of all interior nodes share some number of digits with the tree’s *groupId*. Therefore, we can ensure that  $k$  Scribe trees have a disjoint set of interior nodes simply by choosing *groupIds* for the trees that all differ in the most significant digit.

Setting  $k = 2^b$  ensures that each participating node has an equal chance of becoming an interior node in some tree. If  $k$  is chosen such that  $k = 2^i$  and  $i \leq b$ , then it is still possible to ensure this fairness by exploiting certain properties of the Pastry routing table, but we omit the details to conserve space. Without loss of generality, we assume that  $k = 2^b$  in the rest of this paper.

**Limiting node degree** The resulting forest of Scribe trees satisfies the independence requirement and the nodes’ constraints on the inbound bandwidth, but it does not necessarily satisfy the individual nodes’ outgoing bandwidth constraints. Let us first consider the inbound bandwidth. A node’s inbound bandwidth is proportional to the number of stripes to which the node subscribes. Note that a node has to subscribe to at least one stripe, the one whose *stripeId* shares a prefix with its nodeId, because the node may have to serve as an interior node for that stripe.

The number of children that may attempt to attach to a node is bounded by its indegree in the Pastry overlay, which is influenced by the physical network topology. In general, this number may exceed the number of children a node is able to support. For a SplitStream node to limit its outbound network bandwidth, it must limit its outdegree in the SplitStream forest, i.e., the total number of children it takes on.

Scribe has a built-in mechanism to limit a node’s outdegree. When a node that has reached its maximal outdegree receives a request from a prospective child, it provides the prospective child with a list of its current children. The prospective child then seeks to be adopted by the child with lowest delay, and so on recursively. In Scribe, this procedure is guaranteed to terminate because a leaf node

is required to take on at least one child.

Unfortunately, this procedure is not guaranteed to work in SplitStream. The reason is that a leaf node in one tree may be an interior node in another stripe tree, and may have already reached its outdegree limit with respect to that stripe tree.

**Balancing trees** SplitStream uses the following algorithm to resolve the case where a node that has reached its outdegree limit receives a join request from a prospective child. First, the node adopts the prospective child regardless of the outdegree limit. Then, it evaluates its new set of children to select a child to reject. This selection is made in an attempt to maximize the efficiency of the SplitStream forest.

First, the node looks for children that are subscribed to stripes whose stripeIds do not share a prefix with the local node's nodeId. (How the node could have acquired such a child in the first place will become clear in a moment). If multiple such nodes exist, one is chosen randomly. If no such child exists then the current node is an interior node for only one stripe tree, and it selects the child whose nodeId has the shortest prefix match with that stripeId. If multiple such nodes exist, one is chosen randomly. The chosen child is then notified that it has been orphaned.

The orphaned child then seeks to locate a new parent by sending an anycast message to a special Scribe group called the *spare capacity group*. All SplitStream nodes whose number of children is below their limit join this group. The anycast message is delivered to a leaf node in the spare capacity group tree that is near the orphan in the physical network. This node checks whether it receives any of the stripes to which the orphaned child seeks to subscribe. If so, it verifies that the orphan is not an ancestor in the corresponding stripe tree, which would create a cycle. If both tests succeed for some stripe, the node takes on the orphan as a child; if as a result, the node has now reached its outdegree limit, it leaves the spare capacity group. If one of the tests fails, the node forwards the message to its parent, starting a depth-first search (DFS) of the spare capacity group tree until an appropriate member is found.

This procedure is guaranteed to locate an appropriate parent for the orphan if one exists. Moreover, the properties of Scribe trees and the DFS of the spare capacity tree ensure that the parent is near the orphan in the physical network, among all prospective parents. This provides low delay and low link stress in the physical network. However, the algorithm as described may sacrifice tree independence, because the new parent may be already an interior node in another stripe tree. Thus, its failure may cause the temporary loss of more than one stripe for some nodes.

It is possible to minimize this partial loss of independence at the expense of higher delay, link stress, and cost of the forest construction. However, complete independence is generally only feasible if there is some excess forwarding capacity, where the total outdegree of all nodes exceeds the total indegree of all nodes. One approach to preserving independence is to add a third test during the DFS in the

spare capacity group tree, which verifies that the prospective parent's nodeId shares a prefix with the stripeId to which the orphan subscribes. This ensures independence, but may require a more extensive exploration of the spare capacity group tree, may yield a parent that is more distant in the physical network, and may not always locate a parent in the absence of sufficient excess forwarding capacity. One may balance these concerns by limiting the scope of the DFS, and relax the third test if no parent was found within that scope. SplitStream allows applications to control this tradeoff between independence, delay, link stress, total required forwarding capacity and overhead of forest construction according to its needs.

**Preliminary results** We have performed a preliminary performance evaluation of SplitStream, by running 50,000 SplitStream nodes over an emulated network with 5050 core routers based on the Georgia Tech network topology generator. We constructed a SplitStream forest with 16 stripes, and assigned per-node inbound and outbound bandwidth limits that follow a distribution measured among Gnutella clients in May 2001 [20].

The results are very encouraging. During the SplitStream forest construction (50,000 nodes, 16 stripes), the mean and median number of control messages handled by each node were 33 and 57, respectively. When multicasting a message in each stripe, the medians of the relative average delay penalty (RAD) and the relative maximum delay penalty (RMD), compared to IP multicast, were 2.33 and 3.64, respectively. These values are about 1.5 and 2 times higher, respectively, than the values measured in a single Scribe tree on the same topology. This increase reflects the principal cost of balancing the forwarding load across all participants in SplitStream.

We also considered the degree of independence in the SplitStream forest. Without any of the independence-preserving techniques described above, and with a highly constrained bandwidth allocation (outbound bandwidth not to exceed inbound bandwidth at any node), we found that over 95% of the nodes had independent (i.e., node disjoint) paths to the source in 13 or more of the 16 stripes to which they subscribed. Thus, even in pessimal cases, the loss of independence is modest. A more comprehensive evaluation of SplitStream will be presented in a forthcoming full paper.

## 5 Related work

Many application-level multicast systems have been proposed recently, e.g. [8, 14, 18, 21, 7, 1]. All are based on a single multicast tree.

Several systems exist that use end-system multicast for video distribution, notably Overcast [14] and SpreadIt [10]. Both systems create a single multicast tree. Overcast relies on dedicated servers, whilst both SpreadIt and SplitStream utilise the participating clients. However, unlike SpreadIt, SplitStream distributes forwarding load over

all participants using multiple multicast trees, thereby reducing the bandwidth demands on individual peers.

Nguyen and Zakhor [16] propose streaming video from multiple sources concurrently, thereby exploiting path diversity and increasing tolerance to packet loss. They subsequently extend the work in [16] to use Forward Error Correction [3] encodings. The work assumes that the client is aware of the set of servers from which to receive the video. SplitStream constructs multiple endsystem-based multicast trees in a decentralized fashion and is therefore more scalable.

CoopNet [17], like SplitStream, utilises multiple trees and stripes video using MDC. Each stripe is delivered to the client using a different source. When a server is overloaded, clients are 'redirected' to other clients, thereby creating a distribution tree routed at the server. There are two fundamental differences between CoopNet and SplitStream: (i) CoopNet uses a centralised algorithm (running on the server) to build the trees whilst SplitStream is completely decentralised; and (ii) CoopNet does not attempt to manage the bandwidth contribution of individual nodes; however, it is possible to add this capability to CoopNet.

In [4], algorithms and content encodings are described that enable parallel downloads and increase packet loss resilience in richly connected, collaborative overlay networks by exploiting downloads from multiple peers. SplitStream provides a complete system for content distribution in collaborative overlay networks. It explicitly stripes content and creates a tree for each stripe. Also, SplitStream's main goal is to spread the forwarding load across all participants.

Fcast [12] is a reliable file transfer protocol based on IP multicast. It combines a Forward Error Correction [3] encoding and a data carousel mechanism. Instead of relying on IP multicast, Fcast could be easily built upon SplitStream, for example, to provide software updates cooperatively.

## 6 Conclusions

We have sketched the design of SplitStream, a high-bandwidth content distribution system based on endsystem multicast in cooperative environments. Preliminary performance results are very encouraging. The system is able to distribute the forwarding load among the participating nodes, subject to individual node bandwidth limits. When combined with redundant content encoding, SplitStream yields resilience to node failures and unannounced departures, even while the affected multicast tree is repaired. The overhead of the forest construction is modest and well balanced, and the resulting increase in delay penalty and link stress is modest, when compared to a conventional tree-based endsystem multicast system.

We are currently exploring various optimizations in constructing the SplitStream forest, guided by application needs. A forthcoming paper will present comprehensive

results, including results of experiments using the Planet-Lab Internet testbed.

## References

- [1] S. Banerjee, B. Bhattacharjee, and C. Kommareddy. Scalable application layer multicast. In *ACM SIGCOMM*, Aug. 2002.
- [2] K. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu, and Y. Minsky. Bimodal multicast. *ACM TOCS*, 17(2):41–88, May 1999.
- [3] R. Blahut. *Theory and Practice of Error Control Codes*. Addison Wesley, MA, 1994.
- [4] J. Byers, J. Considine, M. Mitzenmacher, and S. Rost. Informed content delivery across adaptive overlay networks. In *SIGCOMM'2002*, Pittsburgh, PA, USA, Aug. 2002.
- [5] M. Castro, P. Druschel, Y. C. Hu, and A. Rowstron. Exploiting network proximity in peer-to-peer overlay networks, 2002. Technical report MSR-TR-2002-82.
- [6] M. Castro, P. Druschel, A.-M. Kermarrec, and A. Rowstron. Scalable peer-to-peer anycast for distributed resource management, 2002. Submitted to IPTPS'03.
- [7] M. Castro, P. Druschel, A.-M. Kermarrec, and A. Rowstron. SCRIBE: A large-scale and decentralized application-level multicast infrastructure. *IEEE JSAC*, 20(8), Oct. 2002.
- [8] Y. Chu, S. Rao, and H. Zhang. A case for end system multicast. In *ACM Sigmetrics*, pages 1–12, June 2000.
- [9] Y. K. Dalal and R. Metcalfe. Reverse path forwarding of broadcast packets. *CACM*, 21(12):1040–1048, 1978.
- [10] H. Deshpande, M. Bawa, and H. Garcia-Molina. Streaming live media over a peer-to-peer network, Apr. 2001. Stanford University, CA, USA.
- [11] P. Eugster, S. Handurukande, R. Guerraoui, A.-M. Kermarrec, and P. Kouznetsov. Lightweight probabilistic broadcast. In *DSN*, July 2001.
- [12] J. Gemmell, E. Schooler, and J. Gray. Fcast multicast file distribution. *IEEE Network*, 14(1):58–68, Jan 2000.
- [13] V. K. Goyal. Multiple description coding: Compression meet the network. *IEEE Signal Processing Magazine*, 18(5):74–93, Sept. 2001.
- [14] J. Jannotti, D. Gifford, K. Johnson, M. Kaashoek, and J. O'Toole. Overcast: Reliable multicasting with an overlay network. In *OSDI 2000*, San Diego, CA, 2000.
- [15] A. Mohr, E. Riskin, and R. Ladner. Unequal loss protection: Graceful degradation of image quality over packet erasure channels through forward error correction. *IEEE JSAC*, 18(6):819–828, June 2000.
- [16] T. Nguyen and A. Zakhor. Distributed video streaming with forward error correction. In *Packet Video Workshop*, Pittsburgh, USA., 2002.
- [17] V. Padmanabhan, H. Wang, P. Chou, and K. Sripanidkulchai. Distributing streaming media content using cooperative networking. In *NOSSDAV*, Miami Beach, FL, USA, May 2002.
- [18] S. Ratnasamy, M. Handley, R. Karp, and S. Shenker. Application-level multicast using content-addressable networks. In *NGC*, Nov. 2001.
- [19] A. Rowstron and P. Druschel. Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In *IFIP/ACM Middleware 2001*, Heidelberg, Germany, Nov. 2001.
- [20] S. Saroiu, P. K. Gummadi, and S. D. Gribble. A measurement study of peer-to-peer file sharing systems. In *MMCN*, San Jose, CA, Jan. 2002.
- [21] S. Zhuang, B. Zhao, A. Joseph, R. Katz, and J. Kubiawicz. Bayeux: An architecture for scalable and fault-tolerant wide-area data dissemination. In *NOSSDAV*, June 2001.