



**HAL**  
open science

# AS-cast: Lock Down the Traffic of Decentralized Content Indexing at the Edge

Adrien Lebre, Brice Nédelec, Alexandre van Kempen

► **To cite this version:**

Adrien Lebre, Brice Nédelec, Alexandre van Kempen. AS-cast: Lock Down the Traffic of Decentralized Content Indexing at the Edge. [Research Report] RR-9418, Inria Rennes. 2021, pp.1-21. hal-03333669v2

**HAL Id: hal-03333669**

**<https://inria.hal.science/hal-03333669v2>**

Submitted on 13 Sep 2021 (v2), last revised 12 Jan 2023 (v3)

**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.



# AS-cast : Lock Down the Traffic of Decentralized Content Indexing at the Edge

Adrien Lèbre, Brice Nédelec, Alexandre Van Kempen

**RESEARCH  
REPORT**

**N° 9418**

September 2021

Project-Team STACK





## AS-cast : Lock Down the Traffic of Decentralized Content Indexing at the Edge

Adrien Lèbre, Brice Nédelec, Alexandre Van Kempen

Project-Team STACK

Research Report n° 9418 — September 2021 — 21 pages

**Abstract:** Although the holy grail to store and manipulate data in Edge infrastructures is yet to be found, state-of-the-art approaches demonstrated the relevance of replication strategies that bring content closer to consumers: The latter enjoy better response time while the volume of data passing through the network decreases overall. Unfortunately, locating the closest replica of a specific content requires indexing *every live* replica along with its *location*. Relying on remote services enters in contradiction with the properties of Edge infrastructures as locating replicas may effectively take more time than actually downloading content. At the opposite, maintaining such an index at every node would prove overly costly in terms of memory and traffic, especially since nodes can create and destroy replicas at any time.

In this paper, we abstract content indexing as distributed partitioning: every node only indexes its closest replica, and connected nodes with a similar index compose a partition. Our decentralized implementation AS-cast is (i) efficient, for it uses partitions to lock down the traffic generated by its operations to relevant nodes, yet it (ii) guarantees that every node eventually acknowledges its partition despite concurrent operations. Our complexity analysis supported by simulations shows that AS-cast scales well in terms of generated traffic and termination time. As such, AS-cast can constitute a new building block for geo-distributed services.

**Key-words:** Content indexing, logical partitioning, decentralized algorithm, Edge

RESEARCH CENTRE  
RENNES – BRETAGNE ATLANTIQUE

Campus universitaire de Beaulieu  
35042 Rennes Cedex

## AS-cast : Confiner le trafic de l'indexation décentralisée de contenu au niveau du Edge

**Résumé :** Bien qu'une solution parfaite pour le stockage et la manipulation de données au niveau des infrastructures Edge reste encore à déterminer, les approches de l'état de l'art ont d'ors et déjà démontré la pertinence des stratégies de réplication qui rapprochent le contenu des consommateurs : ces derniers bénéficient de meilleurs temps de réponse et le volume de données transitant par le réseau diminue globalement. Malheureusement, localiser la réplique la plus proche d'un contenu spécifique nécessite d'indexer *chacune* des répliques *existantes* avec leurs informations de *localisation*. Se fier à des services distants entre en contradiction avec les propriétés des infrastructures Edge de telle sorte que localiser des répliques prend davantage de temps que de télécharger le contenu. À l'opposé, chaque nœud ne peut maintenir l'index des toutes les répliques existantes car cela s'avérerait extrêmement coûteux en termes de mémoire et de trafic, en particulier lorsque les nœuds sont libres de créer ou détruire des répliques à n'importe quel moment. Dans cet article, nous abstrayons l'indexation de contenu par un problème de partitionnement réparti : chacun des nœuds maintient l'index de sa réplique la plus proche, et les nœuds connectés possédant un index similaire font partie d'une même partition. Notre implémentation décentralisée nommée AS-cast est (i) efficace, car elle utilise les partitions afin de confiner le trafic généré par son fonctionnement aux seuls nœuds concernés, tout en (ii) garantissant le fait que chaque nœud finisse par identifier la partition à laquelle il appartient, et ce malgré les opérations concurrentes. Notre analyse en complexité appuyée par des simulations montre que AS-cast passe à l'échelle en termes de trafic généré et de temps de terminaison. De ce fait, AS-cast peut constituer les nouvelles fondations pour construire des services géo-distribués.

**Mots-clés :** Indexation de contenu, partitionnement logique, algorithme décentralisé, Edge

## 1 Introduction

In recent years, the data storage paradigm shifted from centralized in the cloud to distributed at the edges of the network. This change aims at keeping the data close to (i) its producers since data may be too expensive or sensitive to be transmitted through the network; and (ii) its consumers so data may quickly and efficiently reach them [12, 15, 35]. To favour this transition, new designs for data management across Edge infrastructures have been investigated [8, 9, 16, 18]. They enable strategies to confine traffic by writing data locally and replicating content according to effective needs. However, locating content remains challenging. Retrieving a content location may actually take more time than retrieving the content itself. Indeed, these systems, when not using a centralized index hosted in a Cloud, rely on distributed hash tables (DHT) [29] spread across the different nodes composing the infrastructure. When a client wants to access specific content, it requests a remote node to provide at least one node identity to retrieve this content from. After it retrieved the content, the client can create another replica to improve the performance of future accesses, but it must contact the content indexing service again to notify of such a change.

These approaches contradict with the objectives of Edge infrastructures that aim at reducing the impact of latency as well as the volume of data passing through the network. First, accessing a remote node to request content location(s) raises hot spots and availability issues. But most importantly, it results in additional delays [3, 11] that occur even before the actual download started. Second, the client gets a list of content locations at the discretion of content indexing services. Without information about these locations, it often ends up downloading from multiple replica hosts, yet only keeping the fastest answer. In turns, clients either waste network resources, or face lower response time.

To tackle aforementioned limitations, every node that might request or replicate content must also host its own content indexing service in a fully decentralized fashion [25]. At any time, it can immediately locate the closest replica of specific content. A naive approach would be that every node indexes and ranks *every live* replica along with its *location* information. When creating or destroying a replica, a node would notify all other nodes by efficiently broadcasting its operation [6, 17, 31]. Unfortunately, this also contradicts with Edge infrastructure objectives, for such protocol does not confine the traffic generated to maintain its indexes. A node may acknowledge the existence of replicas at the other side of the network while there already exists a replica next to it.

Instead, a node creating a replica should notify all and only nodes that have no closer replica in the system. This would create interconnected sets of nodes, or *partitions*, gathered around a *source* being their respective replica. A node deleting its replica should notify all members of its partition so they can identify and rally their new closest partition. A periodic advertisement protocol already provides both operations for routing purposes [19]. However, its functioning requires (i) to generate traffic even when the system is quiescent, and (ii) to finely tune physical-time-based intervals that depend on network topology parameters such as network diameter.

The contribution of this paper is threefold:

- We highlight the properties that guarantee decentralized consistent partitioning in dynamic infrastructures. We demonstrate that concurrent creation and removal of partitions may impair the propagation of control information crucial for consistent partitioning. Yet, nodes are able to purge stale information forever using only neighbor-to-neighbor communication, hence leaving room for up-to-date information propagation, and eventually consistent partitioning.
- We provide an implementation entitled AS-cast that uses aforementioned principles to adapt its partitioning to creations and deletions of partitions even in dynamic systems where nodes join, leave, and crash at any time. AS-cast's efficiency relies on a communication primitive

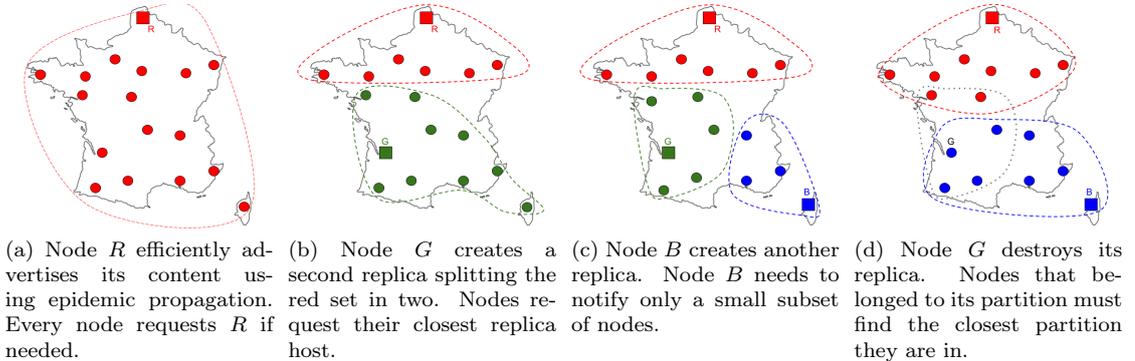


Figure 1: The french RENATER topology. Partitions grow and shrink depending on creations and removals of replicas.

called `scoped-broadcast` that enables epidemic dissemination of messages as long as receiving nodes verify an application-dependant predicate.

- We evaluate AS-cast with simulations. Our experiments empirically support our complexity analysis and the scalability of AS-cast. Using scoped broadcast, AS-cast quickly disseminate messages to a subset of relevant nodes. The higher the number of partitions, the smaller AS-cast’s overhead. Most importantly, AS-cast does not generate traffic in quiescent systems. Our experiments also highlight the relevance of AS-cast in autonomous systems. AS-cast would allow to lock down the traffic generated by content indexing, even in contexts where *physical* partitioning may occur. As such, AS-cast can constitute a new building block for geo-distributed services.

The rest of this paper is organized as follows. Section 2 illustrates the motivation and problem behind our proposal. Section 3 describes dynamic consistent partitioning and its implementation along with its complexity. Section 4 presents our evaluations. Section 5 reviews related work. Finally, Section 6 concludes and discusses future work.

## 2 Motivation and problem

Numerous studies addressed content indexing in geo-distributed infrastructures ranging from centralized services to fully decentralized approaches [23]. We advocate for the latter where nodes host the service themselves so they do not need to request remote – possibly far away third-party – entities to retrieve their content location. In this paper, we tackle the content indexing issue as a dynamic logical partitioning problem. This section motivates our positioning and explains the shortcomings of existing implementations.

**Dissemination:** Figure 1 depicts an infrastructure comprising 17 interconnected nodes spread across France. In Figure 1a, a single Node  $R$  hosts the content, so every other node downloads from this node when it needs it. In that regard, Node  $R$  only needs to disseminate a message notifying all other nodes of its new content. Node  $R$  can use uniform reliable broadcast [17] and epidemic propagation [13] to guarantee that (i) every node eventually knows the content location (ii) by efficiently using neighbor-to-neighbor communication.

**Location:** Then, in Figure 1b, another Node  $G$  creates a replica of this content. Similarly to Node  $R$ , it notifies other nodes of this new replica. However, to avoid that north nodes request its replica, and south nodes request the northern one, nodes hosting a replica must add *location information* along with their notifications. As consequence, every node eventually knows every

replica location and can download from its closest host. Red nodes would request Node  $R$  while green nodes would request Node  $G$ .

**Scoped broadcast:** Then, in Figure 1c, another Node  $B$  creates a replica of this content. Similarly to Node  $R$  and Node  $G$ , Node  $B$  can notify other nodes of this new replica. However, the set of nodes that could actually use this new replica is actually much smaller than the network size. Uniform reliable broadcast is not designed for such a context and would generate a lot of unnecessary traffic. Instead, nodes need a communication primitive that propagates notifications within a scope by evaluating a predicate, starting at its broadcaster (the *source*). In other terms, nodes propagate notifications as long as they consider them useful based on location information they carry. We call such primitive *scoped broadcast* (see Section 3.1), for messages transitively reach a subset of interconnected nodes (the *scope*). Using this primitive, nodes can lock down the traffic of content indexing to relevant nodes.

**Logical partitioning:** Every node ends up with at least one known replica that is its closest. The set of interconnected nodes with the same closest replica is a *partition*. Every node belongs to one, and only one, partition (see Section 3.2).

**Dynamic partitioning and adaptive scoped broadcast:** Finally, in Figure 1d, Node  $G$  destroys its replica. Every nodes that belonged to its green partition must choose another partition to be in. While it makes sense for Node  $G$  to scoped broadcast its removal, Node  $B$  and Node  $R$  cannot afford to continuously advertise their replica to fill the gap left open by Node  $G$ . Instead, we want to trigger once again scoped broadcast at bordering nodes of red and blue partitions. The scope of scoped broadcast changes over receipts by nodes. Dynamic partitioning raises additional challenges related to concurrent operations where removed partitions would block the propagation of other partitions (see Section 3.3).

Next Section details the properties of scoped broadcast and dynamic partitioning in dynamic networks, and provides an implementation called AS-cast along with its complexity.

### 3 Adaptive scoped broadcast

To assign and maintain nodes to their best partition according to replica creations and removals, as well as dynamic infrastructure changes, we designed and implemented AS-cast. AS-cast stands for A d a p t i v e s c o p e d s c o p e d b r a d c a s t. It relies on a primitive that allows a node to broadcast a message within a limited scope. We first use this primitive to guarantee consistent partitioning when a node can only create a new replica within the system. We highlight the issue when a process can also destroy a replica, and provide a second algorithm that handles replica removals as well as dynamic changes of the infrastructure. This section describes the communication primitive called scoped broadcast, then discusses the properties that guarantee consistent partitioning in our context, and finally details our implementation AS-cast and its complexity.

#### 3.1 Scoped broadcast

In this paper, we consider Edge infrastructures with heterogeneous nodes interconnected by communication links. Nodes involved in the management of content may crash but are not byzantine. Finally, nodes can reliably communicate through asynchronous message passing to other known nodes called neighbors.

**Definition 1** (Edge infrastructure). *An Edge infrastructure is a connected graph  $G(V, E)$  of vertices  $V$  and bidirectional edges  $E \subseteq V \times V$ . A path  $\pi_{ij}$  from Node  $i$  to Node  $j$  is a sequence of vertices  $[i, k_1, k_2, \dots, k_n, j]$  with  $\forall m : 0 \leq m \leq n, \langle \pi_{ij}[m], \pi_{ij}[m+1] \rangle \in E$ .*

We define scoped broadcast as a communication primitive that propagates a message around its broadcaster within an application-dependant scope:

**Definition 2** (Scoped broadcast (S-cast)). *When Node  $x$  scoped broadcasts  $b_x(m)$  a message  $m$ , every correct node  $y$  within a scope receives  $r_y(m)$  and delivers it  $d_y(m)$ . The scope depends on the state  $\sigma$  of each node, the metadata  $\mu$  piggybacked by each message, and a predicate  $\phi$  verified from node to node:  $b_x(m) \wedge r_y(m) \implies \exists \pi_{xy} : \forall z \in \pi_{xy}, \phi(\mu_z, \sigma_z)$ .*

This definition encompasses more specific definitions of related work [22, 28, 36]. It also highlights that epidemic propagation and scoped broadcast have well-aligned preoccupations. More precisely, it underlines the transitive relevance of messages, thus nodes can stop forwarding messages as soon as the corresponding predicate becomes unverified. We use S-cast to dynamically and efficiently modify the state of each node depending on all partitions that exist in the system.

### 3.2 Consistent partitioning

At any time, a node can decide to become a *source*, hence creating a new partition in the system by executing an **Add** operation. This partition includes at least its source plus neighboring nodes that estimate they are closer to this source than any other one. Such distance (or *weight*) is application-dependant: in the context of maintaining distributed indexes, this would be about link latency.

**Definition 3** (Partitioning). *Let  $S \subseteq V$  be the set of sources, and  $P_s$  be the partition including at least Node  $s$ , each node belongs to at most one partition  $\forall p, q \in V, \forall s_1, s_2 \in S : p \in P_{s_1} \wedge q \in P_{s_2} \implies p \neq q \vee s_1 = s_2$ , and there exists at least one path  $\pi_{ps}$  of nodes that belong to this partition  $\forall q \in \pi_{ps} : q \in P_s$ .*

Definition 2 and Definition 3 share the transitive relevance of node states. However, we further constrain the partitioning in order to guarantee the existence of exactly one consistent partitioning that nodes eventually converge to.

**Definition 4** (Consistent partitioning (CP)). *Let  $W_{xy} = W_{yx}$  be the positive symmetric weight between  $x$  and  $y$ ,  $\Pi_{xz}$  be the shortest path from  $x$  to  $z$  the weight of which  $|\Pi_{xz}|$  is lower than any other path weight, the only consistent partitioning  $\mathcal{P}$  is a set of partitions  $P_{s \in S}$  such that each node belongs to a logical partition comprising its closest source:  $\forall p \in P_{s_1} : \nexists P_{s_2}$  such that  $|\Pi_{s_1 p}| > |\Pi_{s_2 p}|$ .*

Unfortunately, nodes do not share a common global knowledge of the network state. For nodes to reach consistent partitioning together, a Node  $s$  adding a partition must send a notification  $\alpha_s$  to every node that is closer to it than any other source. Since epidemic dissemination and scoped broadcast have well-aligned preoccupations, we assume implementations relying on message forwarding from neighbor-to-neighbor.

**Theorem 1** (Best eventual forwarding (BEF)  $\implies$  CP). *Assuming reliable communication links where a correct node  $q$  eventually receives the message  $m$  sent to it by a node  $p$  ( $s_{pq}(m) \implies r_q(m)$ ), nodes eventually reach consistent partitioning if each node eventually forwards its best known partition.*

*Proof.* When a node  $s_1$  becomes a source, it belongs to its own partition, for there exists no better partition than its own:  $\forall p \in V : |\Pi_{s_1 s_1}| < |\Pi_{s_1 p}|$ . It delivers, hence forwards such an  $\alpha$  message to its neighbors. Since communication links are reliable, neighboring nodes eventually receive such a notification  $\forall \langle s_1, q \rangle \in E, s_1 \in S \iff \square r_q(\alpha_{s_1}^{w_{s_1 q}})$ . Most importantly, whatever the order of

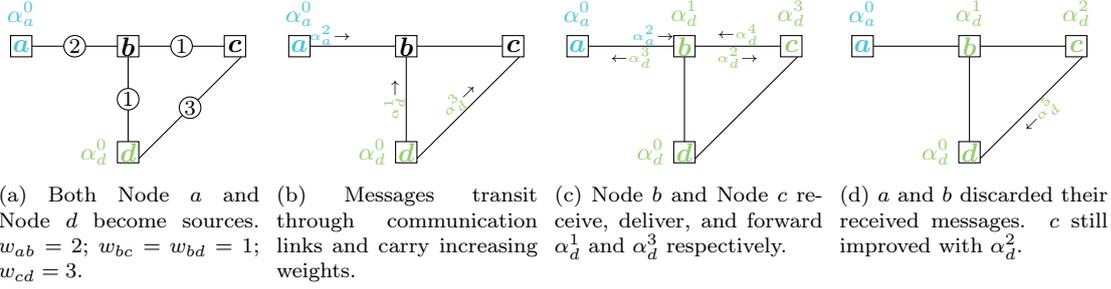


Figure 2: Efficient consistent partitioning using S-cast. Partition  $P_a$  includes  $a$  while Partition  $P_d$  includes  $b, c,$  and  $d$ . Node  $c$  and Node  $d$  never acknowledge the existence of Source  $a$ , for Node  $b$  stops the propagation of the latter's notifications.

received messages, every node  $q'$  in this neighborhood – such that there exists no better partition  $s_2$  than the received one – delivers and forwards it:  $\forall s_2 \in S : |\Pi_{q's_1}| < |\Pi_{q's_2}| \implies \square d_{q'}(\alpha_{s_1}^{w_{s_1 q'}})$ . By transitivity, the message originating from  $s_1$  reaches all such nodes through their shortest paths:  $\forall q'' \in V, s_1, s_2 \in S : |\Pi_{q''s_1}| < |\Pi_{q''s_2}| \implies \square d_{q''}(\alpha_{s_1}^{|\Pi_{s_1 q''}|})$ . Since there exists only one best sum of weights per node that can never be retracted, nodes eventually reach consistent partitioning.  $\square$

---

**Algorithm 1:** Adding a partition by Node  $p$ .

---

```

1  $O_p, W_p$  // set of neighbors and weights
2  $s \leftarrow \emptyset$  // best source of partition ( $\sigma$ )
3  $d \leftarrow \infty$  // smallest distance to  $s$  ( $\sigma$  and  $\mu$ )
4 func Add ( ) //  $\alpha_p^0$ 
5 | receiveAdd( $\emptyset, p, 0$ ) //  $b_p(\alpha_p^0)$ 
6 func receiveAdd( $q, s', d'$ ) //  $r_p(\alpha_{s'}^{d'})$  from  $q$ 
7 | if  $d' < d$  then // ( $\phi$ )
8 | |  $s \leftarrow s'$  //  $d_p(\alpha_{s'}^{d'})$ 
9 | |  $d \leftarrow d'$ 
10 | | foreach  $n \in O_p \setminus q$  do //  $f_{pn}(\alpha_{s'}^{d'})$ 
11 | | | send $_n(s', d' + W_{pn})$  //  $s_{pn}(\alpha_{s'}^{d' + w_{pn}})$ 

```

---

Algorithm 1 shows the instructions that implement consistent partitioning when (i) weights are scalar values, (ii) nodes only add new partitions to the system, (iii) and nodes never crash nor leave the system. Figure 2 illustrates its behavior on a system comprising 4 nodes  $a, b, c,$  and  $d$ . Node  $a$  and Node  $d$  become the sources of their partition. They S-cast a notification add message:  $\alpha_a^0$  and  $\alpha_d^0$ . They initialize their own state with the lowest possible bound 0 (see Line 5), and send a message to each of their neighbors by accumulating the corresponding edge weight (see Line 11). In Figure 2c, Node  $b$  receives  $\alpha_a^1$ . Since it improves its own partition distance, it keeps it and forwards it to its neighbors. In Figure 2d, Node  $b$  discards  $\alpha_a^2$ , for it does not improve its partition distance. Nodes  $c$  and  $d$  will never acknowledge that Source  $a$  exists. Ultimately, nodes discard last transiting messages. Despite the obvious lack of traffic optimization, the system reaches consistent partitioning.

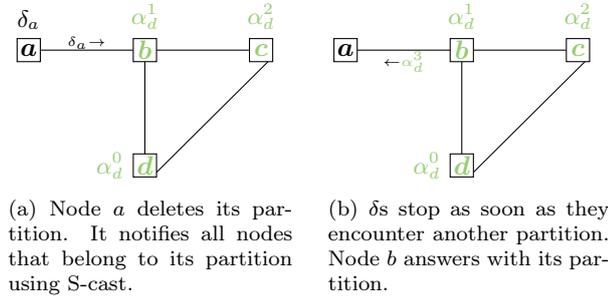


Figure 3: Efficient removal of a partition using S-cast. Node  $a$  eventually knows that it belongs to Partition  $P_d$ .

While only adding logical partitions to the distributed system is straightforward, removing them introduces additional complexity.

### 3.3 Dynamic consistent partitioning

At any time, a source can revoke its self-appointed status of source by executing a `Del` operation, hence deleting its partition from the system. All nodes that belong to this partition must eventually choose another partition to belong to. In Figure 3, two partitions exist initially:  $P_a$  and  $P_d$  that respectively include  $\{a\}$ , and  $\{b, c, d\}$ . In Figure 3a, Node  $a$  deletes its partition. It notifies all neighboring nodes – here only Node  $b$  – that may belong to its partition using S-cast. Upon receipt, Node  $b$  discards the delete notification  $\delta_a$  since  $\delta_a$  does not target the former's partition  $P_d$ . Node  $b$  sends its own best partition  $\alpha_d^1$  that may be the best for Node  $a$ . Eventually, every node belongs to Partition  $P_d$ . In this scenario, they reach consistent partitioning.

Delete operations add a new notion of order between events, and most importantly between message deliveries. Delete operations implicitly state that all preceding events become obsolete, and that all messages originating from these preceding events convey stale control information.

**Definition 5** (Happens-before  $\rightarrow$  [27]). *The transitive, irreflexive, and antisymmetric happens-before relationship defines a strict partial order between events. Two messages are concurrent if none happens before the other.*

**Definition 6** (Last win order and staleness). *When a node  $p$  broadcasts two messages  $b_p(m) \rightarrow b_p(m')$ , no node  $q$  can deliver  $m$  if it has delivered  $m'$ :  $\nexists q \in V$  with  $d_q(m') \rightarrow d_q(m)$ , for  $m$  convey stale control information.*

Unfortunately, stale control information as stated in Definition 6 may impair the propagation of both (i) notifications about actual sources, and (ii) notifications about deleted partitions. Figure 4 highlights such consistency issues with dynamic partitions, even in contexts where nodes have FIFO links, i.e., where nodes receive the messages in the order of their sending. Both Node  $a$  and Node  $d$  add then delete their partition concurrently. Node  $c$  receives, delivers, and forwards  $\alpha_d^3$  followed by  $\delta_a$ . In the meantime, Node  $b$  receives, delivers, and forwards  $\alpha_d^2$ . In Figure 4c, both Node  $c$  and Node  $d$  deliver the message about Partition  $P_a$ , for they did not have any known partition at receipt time. On the contrary, Node  $b$  delivers  $\alpha_d^1$ , for it improves its distance to a known source. Node  $b$  then blocks Node  $a$ 's removal notification  $\delta_a$ . It never reaches Node  $c$  nor Node  $d$ . Also, Node  $c$  does not deliver  $\alpha_d$  and  $\delta_d$  since it already delivered

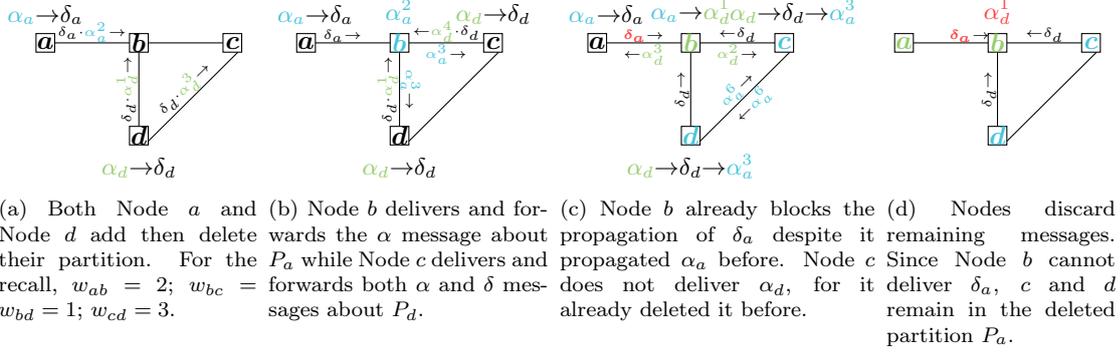


Figure 4: Naive propagation of  $\alpha$  and  $\delta$  messages is insufficient to guarantee consistent partitioning.

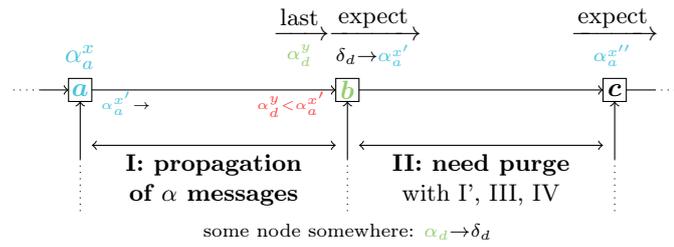
it. The system converges to an inconsistent state where some nodes assume they belong to a partition while it does not exist anymore. Naive propagation of  $\alpha$  and  $\delta$  messages is insufficient to guarantee consistent partitioning when any node can add or delete its partition at any time. Nodes need to exploit their local knowledge gathered over message receipts.

One could guarantee consistent partitioning by always propagating the  $\delta$  messages corresponding to the  $\alpha$  messages it propagated before. In Figure 4, it means that as soon as Node  $b$  forwards  $\alpha_a^2$ , it assumes that its neighbors Node  $c$  and Node  $d$  may need the notification of removal  $\delta_a$  if it exists. However, such a solution also implies that nodes generate traffic not only related to their current partition, but also related to partitions they belong to in the past. This would prove overly expensive in dynamic systems where nodes join and leave the system, create and delete partitions, at any time. Instead, we propose to use the delivery order at each node to detect possible inconsistencies and solve them. Together, nodes eventually remove all stale control information (transiting messages and local states) of the system leaving room for propagation of messages about up-to-date partitions.

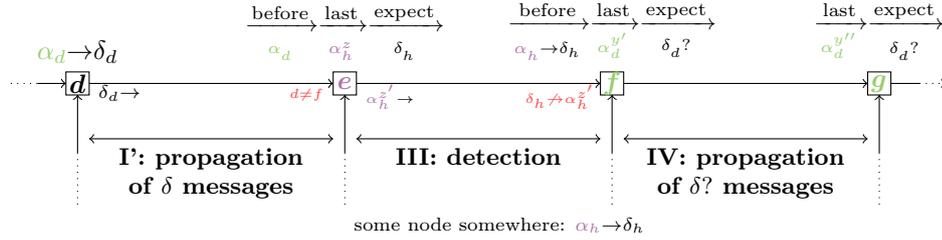
**Corollary 1** (purge + BEF  $\implies$  Dynamic consistent partitioning (DCP)). *When each node can **Add** and **Delete** its partition at any time, consistent partitioning requires (i) eventual purging of stale messages ( $b_p(m) \rightarrow b_p(m') \implies \square(d_q(m) \implies d_q(m) \rightarrow d_q(m''))$ ); (ii) and best eventual forwarding.*

*Proof.* Nodes deliver and forward the best sum of weights they received. Figure 5a depicts the only case where a node would block the epidemic propagation: the node received and delivered a message from a partition that has since then been deleted but with a better sum of weights than newly received ones (in the example  $\alpha_d^y < \alpha_a^{x'}$ ). Removing *forever* all such stale messages about deleted partitions would allow nodes to propagate their best up-to-date partition again, eventually reaching consistent partitioning together, as stated by Theorem 1.  $\square$

**Theorem 2** (three-phase purge). *Purging stale control information from its source to every node that may have delivered it last requires three phases: (i) propagation of delete notifications, (ii) detection of possible blocking conditions, and (iii) propagation of possibly deleted but blocked notifications.*



(a) Stale  $\alpha$  messages may stop other  $\alpha$  messages from reaching all nodes ( $b$  and  $c$ ) that require it along the shortest path. Consistent partitioning requires eventual purging of stale messages (see Figure 5b).



(b) Stale  $\alpha$  messages may stop  $\delta$  messages from reaching processes with the corresponding partition. Consistent partitioning requires that all nodes propagate  $\delta$  messages as often as possible (I); detect possible inconsistencies when parents'  $\alpha$  messages contradict history or state (III); notify processes of possible inconsistencies (IV).

Figure 5: Concurrent operations, delivery orders, and stale control information may impair consistent partitioning. Dynamic consistent partitioning requires eventual purging of stale messages, and best eventual forwarding. *before*, *last*, and *expect* respectively refer to the delivery history, the last delivered  $\alpha$  message, and the message expected to reach consistent partitioning.

*Proof.* Figure 5b depicts all possible cases that need appropriate purging. A chain of nodes  $[d, e, f, g]$  delivered and forwarded a message  $\alpha_d$ . Some nodes may have delivered messages about other partitions before or after  $\alpha_d$ .

**nodes with last  $\alpha_d$ :** Nodes such as Node  $d$  where the last and best partition is still  $P_d$ . By propagating from node to node the corresponding  $\delta_d$  message, these nodes purge  $\alpha_d$ .

**nodes with  $\alpha_d^y \rightarrow \alpha_h^z$  with last  $\alpha_h^z$ , for  $\alpha_h^z < \alpha_d^y$ :** Nodes such as Node  $e$  that do not directly suffer from the non-delivery of  $\delta_d$ . They need to purge their own  $P_h$  if need be. Nonetheless, they deliver and forward  $\alpha_h$  that at least one following node will receive:

**nodes with  $\delta_h \rightarrow \alpha_d$ :** Nodes such as Node  $f$  that do not belong to the preceding category, for they already delivered  $\delta_h$ . Nevertheless, best eventual forwarding guarantees that they receive  $\alpha_h$ , which contradicts their history or state. In such a case, the detecting nodes purge their own  $\alpha_d$ . The detecting nodes must also notify subsequent nodes that their current partition *may be* deleted. Indeed their parent may have blocked the corresponding  $\delta$  messages and subsequent nodes must purge their possibly stale state, for the sake of consistency. We note such notification  $\delta_d?$  to emphasize the uncertainty of a  $\delta$  message.

**nodes with last  $\alpha_d$  receiving  $\delta_d?$  from their parent:** Nodes such as Node  $g$  that have  $\alpha_d$  but preceding delivered messages are not important. Such nodes must trust detecting nodes by delivering and forwarding  $\delta_d?$ . It is worth noting that these messages are subject to all aforementioned blocking conditions, but are also solved by aforementioned mechanisms.

Propagation I' terminates: a node does not deliver a  $\delta$  message if it has already delivered it. Propagation IV terminates: a node does not deliver a looping message. These three phases propagation-detection-propagation ensure that nodes eventually purge stale  $\alpha$  messages from the system, leaving room for Propagation I.  $\square$

### 3.4 Implementation and complexity

Algorithm 2 provides the few instructions of AS-cast that implement three-phase purge and best eventual forwarding to enable dynamic consistent partitioning.

To implement last win order as stated in Definition 6, each node maintains a vector of versions that associates the respective known local counter of each known source, or has-been source. Each message  $m$  carries its source  $s$  along with its version  $c$  that we subscript  $m_{s,c}$ . Each node delivers and forwards (i) only newest messages (see Line 10) that (ii) improve their best known partition (see Line 17). This vector enables both Propagation I and Propagation I'. Its size increases linearly with the number of sources that the node ever known, which is the number of nodes in the system  $\mathcal{O}(V)$  in the worst case. Nevertheless, following the principles of scoped broadcast, we expect that every node only acknowledges a small subset of sources.

This vector of versions constitutes a summary of deliveries. AS-cast uses it to detect possible inconsistent partitioning as depicted in Figure 5b. A detecting node then propagates a  $\delta$  message conveying the unique identifier  $s$  and version  $c$  of the partition along with a path  $\pi$  comprising the identity of nodes that forwarded the corresponding  $\alpha$  message. Such  $\delta_{s,c}^\pi$  message implements the uncertain  $\delta?$  message of Figure 5b. It fulfils the double role of removing stale control information without modifying the local counter of the source (since the partition may still exist), and removing loops to ensure termination (since nodes do not modify their local state to reflect such a delivery, messages piggyback monotonically such increasing control information).

$\delta_{s,c}$  messages used in Propagation I' are more efficient than  $\delta_{s,c}^\pi$  messages used in Propagation IV, both in propagation time and traffic. First, the former propagate using every communication link available in the partition, while the latter propagate following the dissemination tree of the corresponding  $\alpha$  to purge. Second, the former only piggyback source and counter while the latter piggyback the path taken by the corresponding  $\alpha$  to purge. In the worst case,

**Algorithm 2:** AS-cast at Node  $p$  in static networks.

---

```

1  $O_p, W_p$  // set of neighbors and weights
2  $A_{s,c}^{d,\pi} \leftarrow \alpha_{\infty,0}^{\infty,\emptyset}$  // best  $\alpha$  so far
3  $V \leftarrow \emptyset; V[p] \leftarrow 0$  // vector of versions

4 func Add ( )
5 | receiveAdd( $p, \alpha_{p,V[p]+1}^{0,\emptyset}$ )

6 func Del ( )
7 | receiveDel( $p, \delta_{p,V[p]+1}$ )

8 func receiveAdd( $q, \alpha_{s',c'}^{d',\pi'}$ )
9 | if  $q = \pi[|\pi| - 1]$  and //  $\alpha$  and  $A$  come from  $q$ ?
10 |   ( $c' < V[s']$  or // history?
11 |   ( $\alpha_{s',c'}^{d',\pi'} > A_{s,c}^{d,\pi}$  and  $p \notin \pi'$ )) then // state?
12 |   | receiveDel( $q, \delta_{s,c}^{\pi}$ ) // III detect
13 | else if  $A_{s,c}^{d,\pi} < \alpha_{s',c'}^{d',\pi'}$  and  $p \notin \pi'$  then
14 |   |  $V[s'] \leftarrow c'$ 
15 |   |  $A_{s,c}^{d,\pi} \leftarrow \alpha_{s',c'}^{d',\pi'}$ 
16 |   | foreach  $n \in O_p$  do
17 |   |   | send $_n(\alpha_{s',c'}^{d'+W_{pq,\pi' \cup p}})$  // I propagate

18 func receiveDel( $q, \delta_{s',c'}^{\pi'}$ )
19 | if (( $\delta_{s',c'}^{\pi'}$  and  $\alpha_{s',c'-1} = A_{s,c}$ ) or //  $\delta$  received
20 |   ( $\delta_{s',c'}^{\pi'}$  and  $\alpha_{s',c'}^{\pi'} = A_{s,c}^{\pi}$ )) and //  $\delta?$  received
21 |    $p \notin \pi'$  then
22 |   |  $V[s'] \leftarrow \max(V[s'], c')$ 
23 |   |  $A_{s,c}^{d,\pi} \leftarrow \alpha_{\infty,0}^{\infty,\emptyset}$ 
24 |   | foreach  $n \in O_p$  do
25 |   |   | if  $\delta_{s',c'}^{\pi'}$  then send $_n(\delta_{s',c'}^{\pi'})$  // I' propagate
26 |   |   | else send $_n(\delta_{s',c'}^{\pi' \cup p})$  // IV propagate
27 |   | else if  $A_{s,c}^{d,\pi} \neq \alpha_{\infty,0}^{\infty,\emptyset}$  then
28 |   |   | send $_q(\alpha_{s,c}^{d+W_{pq,\pi \cup p}})$  // compete using  $A$ 

```

---

a path includes every node in the system but converges to the average diameter of partitions. Nevertheless, message paths and AS-cast synergize with each other: Paths tend to be smaller as the number of sources in the system increases. Bloom filters could introduce another trade-off by further decreasing the size of path data structures at the cost of premature stopped propagation that could impair consistent partitioning [37].

As stated in Theorem 1, dynamic consistent partitioning not only requires eventual purging of stale control information, but also the retrieval of the best up-to-date partitions. In that regard, both kinds of  $\delta$  messages have dual use: since they already reach the borders of their partition when they remove stale control information, they also trigger a competition when reaching such bordering nodes (see Line 28). This simply consists in sending the current partition through the communication link from which the node received the  $\delta$ . Upon receipt of this answer, nodes act normally by propagating their changes when they improve.

AS-cast guarantees dynamic consistent partitioning by making extensive use of node-to-node communications. It requires (i) purging stale control information, hence propagating  $\delta$  messages; and (ii) retrieving the closest source, hence propagating  $\alpha$  messages. In terms of number of messages, in the average-case, a node  $i$  chosen uniformly at random among all nodes creates a logical partition. Its messages  $\alpha_i$  spread through the network until reaching nodes that belong to another partition closer to them. This splits partitions in half in average. Overall, the  $a^{\text{th}}$  new partition comprises  $\mathcal{O}(\frac{|V|}{2^{\lfloor \log_2 a \rfloor}})$  nodes. This decreases every new partition until reaching one node per new partition: its source. The average number of messages per node is  $\mathcal{O}(\frac{|O|}{2^{\lfloor \log_2 a \rfloor}})$ , where  $|O|$  is the average number of neighbors per node. Deleting the  $a^{\text{th}}$  partition generates twice as many messages as creating the  $a^{\text{th}}$  partition:  $\delta$  messages travel through the partition, then  $\alpha$  messages compete to fill the gap. Overall, the communication complexity shows that AS-cast scales well to the number of partitions.

---

**Algorithm 3:** AS-cast at Node  $p$  in dynamic networks.

---

```

1 func edgeUp( $q$ )                                     // new link to  $q$ 
2   | if  $A_{s,c}^{d,\pi} \neq \alpha_{\infty,0}^{\infty,\emptyset}$  then
3   |   |  $\text{send}_q(\alpha_{s,c}^{d+W_{pq},\pi \cup p})$            // compete with  $q$ 
4 func edgeDown( $q$ )                                    // link to  $q$  removed
5   | if  $q = \pi[\lceil \pi \rceil - 1]$  then
6   |   |  $\text{receiveDel}(q, \delta_{s,c}^\pi)$                  // IV:  $\delta?$  propagation

```

---

Algorithm 3 extends AS-cast to enable dynamic consistent partitioning in dynamic networks where nodes can join or leave the system without notification. We consider that removing a node is equivalent to removing all its edges, and adding a node is equivalent to add as many edges as necessary.

Adding an edge between two nodes may only improve the shortest path of one of these nodes. Therefore, it triggers a competition between the two nodes only. If one or the other node improves, the propagation of  $\alpha$  messages operates normally. Removing an edge between two nodes may invalidate the shortest path of one of involved nodes plus downstream nodes. As a side effect, removing an edge may also impair the propagation of a partition delete. To implement this behavior, AS-cast uses its implementation of  $\delta?$  messages. This prove costly but enables AS-cast even in dynamic networks subject to physical partitioning.

Next Section provides the details of our simulations that assess the proposed implementation.

## 4 Experimentation

This section describes the results of our experimentation evaluating the properties of AS-cast that enables adaptive scoped broadcast in dynamic networks. It answers the following questions: What is the overhead of reaching consistent partitioning using AS-cast in large scale networks? Does AS-cast is relevant in Edge infrastructures where nodes are geo-distributed into clusters? We performed the experimentation using PeerSim, a simulator written in Java that allows researchers to evaluate and assess distributed algorithms in large scale networks [30]. The code is available on the Github platform at <https://anonymous.4open.science/r/peersim-partition-5592>.

### 4.1 Confirm the complexity analysis stating that communication and time overheads depend on the current partitioning of the system.

**Description:** We build a network comprising 10k nodes. First, we chain nodes together, then we add another communication link per node to another random node. Since links are bidirectional, each node has 4 communication links on average. We set the latency of links between 20 and 40 ms at random following a uniform distribution. We set the weight of links between 5 and 15 at random following a uniform distribution. Weights and latency are different, hence the first  $\alpha$  received by a node may not be that of its shortest path. It implies more concurrent operations, hence more traffic to reach consistent partitioning.

We evaluate dynamic consistent partitioning. First, we create 100 partitions one at a time: nodes reach consistent partitioning before we create each new partition. Second, we remove every partition one at a time, in the order of their creation, hence starting from the first and oldest partition that had been added.

We focus on the complexity of AS-cast. We measure the average number of messages generated per node per second during the experiment; and the time before reaching consistent partitioning after adding or removing a partition.

**Results:** Figure 6 shows the results of this experiment. The top part shows the average traffic per node per second divided between  $\alpha$  and  $\delta$  messages. The bottom part shows the time before reaching consistent partitioning.

Figure 6 confirms that AS-cast’s overhead depends on the size of partitions. The first partition is the most expensive, for  $\alpha$ ’s must reach every node which takes time and generate more traffic. AS-cast quickly converges towards consistent partitioning in only 350 milliseconds. The last and 100<sup>th</sup> partition added around 21 seconds is the least expensive. By using scoped broadcast, control information only reaches a small subset of the whole network.

Figure 6 also confirms that AS-cast’s delete operations are more expensive than add operations. Indeed, the top part of the figure shows that after 21 seconds, when the experiment involves removals, traffic includes both  $\alpha$ ’s and  $\delta$ ’s. The latter aims at removing stale information and triggering competition while the former aims at updating shortest paths. As corollary, the convergence time increases, for  $\delta$  messages must reach the partition borders before sound competitors propagate their partition. It is worth noting that this delete operation involves concurrency: removals continue to propagate while the competition is already partially triggered and propagating.

Figure 6 shows that the overhead of adding the 1<sup>st</sup> partition does not correspond to the overhead of deleting this 1<sup>st</sup> partition. When adding it, messages must reach all nodes while when removing it, messages must reach a small subset of this membership. AS-cast’s overhead actually depends on current partitioning as opposed to past partitionings.

Finally, Figure 6 highlights that after 49 seconds, i.e., after the 100 adds and the 100 deletes, nodes do not generate traffic anymore. Being reactive, AS-cast has no overhead when there is

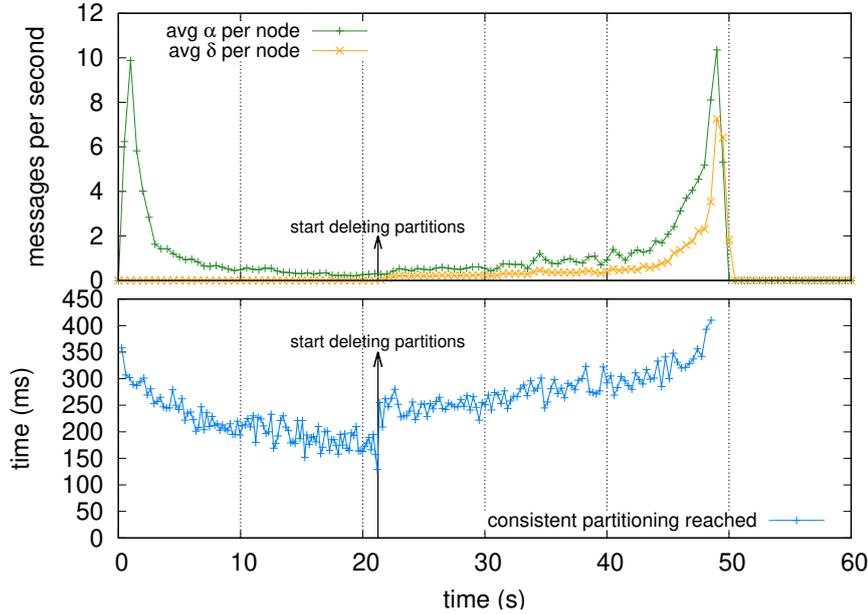


Figure 6: Overhead of consistent partitioning using AS-cast.

no operation in the system. In other words, AS-cast’s overhead actually depends on its usage.

#### 4.2 Confirm that AS-cast locks down the traffic of decentralized content indexing in the context of dynamic inter-autonomous system communications.

**Description:** We build a network by duplicating the GÉANT topology [26] – a real infrastructure spanning across Europe – and by connecting these two clusters with a high latency link: 200 ms simulating cross-continental communications such as Europe-America. The experiments comprise  $2 \times 271 = 542$  nodes and we derive intra-cluster latency from their nodes’ geographic location.

We evaluate the traffic of AS-cast by measuring the average number of messages per node over the experiments. In the first experiment, at 50 ms, only one node becomes source, hence there is only one partition for the whole distributed system. In the second experiment, at 50 ms, two nodes become sources, one per cluster. Afterwards both scenarios are identical. At 850 ms, we remove the link between the two clusters. At 1.7 s, we insert back this link.

**Results:** Figure 7 shows the results of this experimentation. The top part displays the one source experiment while the bottom part displays the one source per cluster experiment.

Figure 7 confirms that concurrent Adds may reach consistent partitioning faster. In particular, the top part of Figure 7 depicts a slow down in traffic around 500 ms due to the high latency inter-continental link. The first plateau corresponds to the source’s autonomous system acknowledging this source, while the second plateau corresponds to the other autonomous system catching up. The inter-continental link acts as a bottleneck that does not exist when each cluster has its own source.

Figure 7 highlights that removing the inter-continental link generates additional traffic. The cluster cut off from the source must purge all control information about this source. Most

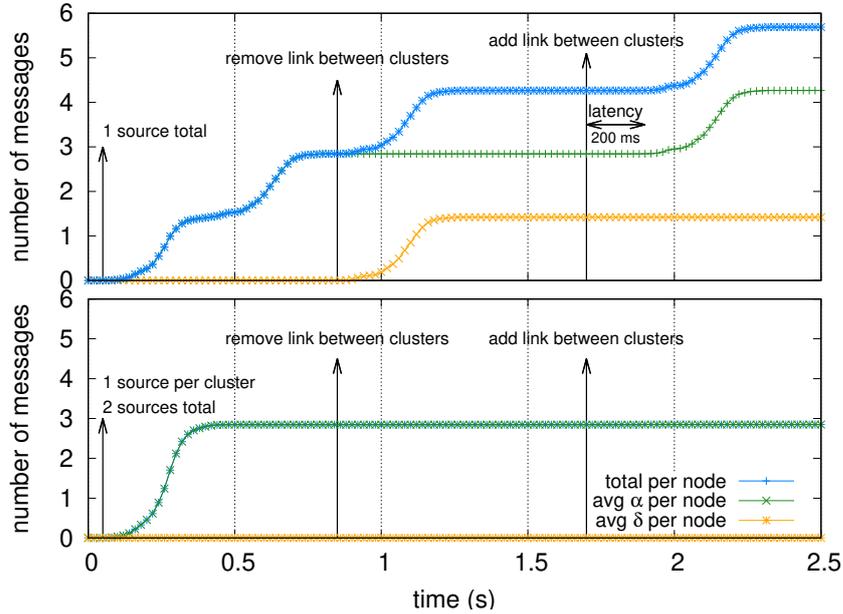


Figure 7: Overhead of AS-cast in the partitioning of 2 clusters.

importantly, Figure 7 shows that AS-cast still operates well when *physical* partitions appear. Nodes from the cluster without source do not belong to any partition while nodes from the other cluster belong to the same partition.

Figure 7 shows that, when adding back the inter-continental link, the two clusters can communicate again. In the experiment involving one source for two clusters, it generates traffic. After a 200 milliseconds delay corresponding to the link latency, the cut off cluster starts to communicate  $\alpha$  messages again. Eventually, all nodes belong to the same partition. However, in the experiment involving one source *per* cluster, the new link does not trigger noticeable traffic, for nodes already know their closest source located in their respective cluster.

Overall, this experimentation empirically demonstrates the benefits of AS-cast in the context of geo-distributed infrastructures, e.g., inter-connected autonomous systems where nodes and communications links may be added and removed at any time. AS-cast coupled with replication strategies designed for Edge infrastructures would allow locking down the traffic of replica indexing, while providing every node with quick access to replicated content.

Next Section reviews state-of-the-art approaches for content indexing in geo-distributed infrastructures and explains their shortcomings.

## 5 Related Work

Content indexing services in geo-distributed infrastructures allow nodes to retrieve specific content while leveraging the advantages of replication. These systems mostly rely on dedicated location services host by possibly remote third-party nodes; but cutting out the middleman requires that each node maintains its own index in a decentralized fashion.

**Third-party:** Dedicated services possibly hosted by remote third-party nodes are popular, for they propose the simplest mean to deploy such service. They must maintain (i) the list of

current replicas along with their respective location; and (ii) the network topology to determine the closest replica for every requesting node.

A central server that registers both these information facilitates the computation, for it gathers all knowledge in one place [1, 2, 16, 33]. However, it comes with well-known issues such as load balancing, or single point of failure.

Distributing this index among nodes circumvents these issues [4, 5, 10, 23, 38], but still raises locality issues where nodes (i) request to a possibly far away location service the content location, and then (ii) request the actual content from far away replicas. For instance, using distributed hash tables (DHT) [5, 10, 23], each node stores a part of the index defined by an interval between hash values. The hash values are the keys to retrieve the associated content location. Before downloading any content, a node asks for its location using its key. After a series of round-trips between DHT servers possibly distant from each others, the node gets a list of available replicas. Contrarily to AS-cast, such services do not ensure to include the closest replica.

In addition, they often do not include distance information along with replica location. Determining where resides the closest replica for every requesting node necessarily involves some knowledge about the *current* topology. Maintaining a consistent view of an ever changing topology across a network is inherently complicated [7, 32]. As consequence, the requesting node ends up downloading from multiple replica hosts, yet only keeping the fastest answer. Nodes either waste network resources, or face lower response time.

**Fully decentralized:** Cutting out the middleman by hosting a location service on each and every node improves response time of content retrieval. However, it still requires every node to index every live replica as well as their respective distance in order to rank them. Named Data Networking (NDN) [21] broadcasts information about cache updates to all nodes in the system. Having the entire knowledge of all the replica locations along with distance information carried into messages, each node easily determine where the closest copy of specific content resides, without contacting any remote service. Conflict-free replicated datatype (CRDT) [34] for set data structures could also implement such a location service. Nodes would eventually have a set of all replicas, assuming eventual delivery of messages. Such solutions imply that every node receives every message, which contradicts with the principles of Edge infrastructures that aim at reducing the volume of data passing through the network. On the opposite, each node running AS-cast only registers its closest replica. This allows AS-cast to use scoped broadcast as a communication primitive to lock down the traffic generated by content indexing based on distances.

**Scoped flooding:** Some state-of-the-art approaches confine the dissemination of messages to interested nodes too.

Distance-based membership protocols such as T-Man [24] make nodes converge to a specific network topology depending on targeted properties. They periodically communicate with their neighbors to compare their state. Following this exchange, they add and remove communication links to other nodes. Over exchanges, they converge towards a topology ensuring the targeted properties. While at first glance membership protocols and AS-cast could share common pre-occupations, AS-cast does not aim at building any topology and never modifies neighbors of nodes.

The most closely related approaches to AS-cast come from information-centric networking (ICN) [14, 19]. They periodically advertise their cache content, and advertisements stop as soon as they reach uninterested nodes. However their operation requires that (i) they constantly generate traffic even in quiescent systems where nodes do not add nor destroy replicas; and (ii) they define physical-clock-based timeouts the value of which depends on network topology properties such as network diameter. Unfitting timeouts lead to premature removals of live replicas; or slow convergence where nodes wrongly believe that their closest replica is live while it was destroyed.

On the opposite, AS-cast quickly informs each node of its closest replica even in large networks; and when the system becomes quiescent, AS-cast has no overhead.

## 6 Conclusion

With the advent of the Edge and IoT era, major distributed services proposed by our community should be revised in order to mitigate as much as possible traffic between nodes. In this paper, we tackled the problem of content indexing, a key service for storage systems. We proposed to address this problem as a dynamic logical partitioning where partitions grow and shrink to reflect content and replicas locations, as well as infrastructure changes. Using an efficient communication primitive called scoped broadcast, each node composing the infrastructure eventually knows the node hosting the best replica of specific content. The challenge resides in handling concurrent operations that may impair the propagation of messages, and in turns, lead to inconsistent partitioning. We highlighted the properties that solve this problem and proposed an implementation called AS-cast. Our simulations confirmed that nodes quickly reach consistent partitioning together, and that the generated traffic remains locked down to partitions.

As future work, we plan on extending the logical partitioning protocol to enable lazy retrieval of partitions. Indeed, using AS-cast, every node partakes in the propagation of indexes of content and replicas. We would like to leverage the hierarchical properties of interconnected autonomous systems to limit the propagation of such indexes only to interested systems.

We also would like to evaluate AS-cast within a concrete storage system such as InterPlanetary File System (IPFS) [20]. This would allow us to assess the relevance of AS-cast in a real system subject to high dynamicity, and compare it against its current DHT-based indexing system that does not include distance in its operation.

## References

- [1] Alexander Afanasyev, Cheng Yi, Lan Wang, Beichuan Zhang, and Lixia Zhang. SNAMP: Secure namespace mapping to scale NDN forwarding. In *2015 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, pages 281–286, April 2015.
- [2] Vinay Aggarwal, Anja Feldmann, and Christian Scheideler. Can ISPS and P2P users cooperate for improved performance? *SIGCOMM Comput. Commun. Rev.*, 37(3):29–40, July 2007.
- [3] Alemnew Shefaraw Asrese, Steffie Jacob Eravuchira, Vaibhav Bajpai, Pasi Sarolahti, and Jörg Ott. Measuring web latency and rendering performance: Method, tools & longitudinal dataset. *IEEE Transactions on Network and Service Management*, 16(2):535–549, June 2019.
- [4] Olivier Beaumont, Anne-Marie Kermarrec, Loris Marchal, and Etienne Rivière. VoroNet: A scalable object network based on Voronoi tessellations. In *2007 IEEE International Parallel and Distributed Processing Symposium*, pages 1–10, 2007.
- [5] Juan Benet. IPFS - content addressed, versioned, P2P file system. *CoRR*, abs/1407.3561, 2014.
- [6] Kenneth P. Birman, Mark Hayden, Ozgur Ozkasap, Zhen Xiao, Mihai Budiu, and Yaron Minsky. Bimodal multicast. *ACM Transactions on Computer Systems*, 17(2):41–88, May 1999.

- 
- [7] Yuri Breitbart, Minos Garofalakis, Cliff Martin, Rajeev Rastogi, Sangeetha Seshadri, and Avi Silberschatz. Topology discovery in heterogeneous IP networks. In *Proceedings IEEE INFOCOM 2000. Conference on Computer Communications. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies (Cat. No.00CH37064)*, volume 1, pages 265–274 vol.1, 2000.
- [8] Bastien Confais, Adrien Lebre, and Benoît Parrein. An object store service for a fog/edge computing infrastructure based on IPFS and a scale-out NAS. In *2017 IEEE 1st International Conference on Fog and Edge Computing (ICFEC)*, pages 41–50. IEEE, 2017.
- [9] Bastien Confais, Adrien Lebre, and Benoît Parrein. Performance analysis of object store systems in a fog and edge computing infrastructure. In *Transactions on Large-Scale Data- and Knowledge-Centered Systems XXXIII*, pages 40–79. Springer, 2017.
- [10] Matteo D’Ambrosio, Christian Dannewitz, Holger Karl, and Vinicio Vercellone. MDHT: A hierarchical name resolution service for information-centric networks. In *Proceedings of the ACM SIGCOMM Workshop on Information-centric Networking, ICN ’11*, pages 7–12, New York, NY, USA, 2011. ACM.
- [11] Trinh Viet Doan, Ljubica Pajevic, Vaibhav Bajpai, and Jörg Ott. Tracing the path to YouTube: A quantification of path lengths and latencies toward content caches. *IEEE Communications Magazine*, 57(1):80–86, 2019.
- [12] Ustav Drolia, Katherine Guo, Jiaqi Tan, Rajeev Gandhi, and Priya Narasimhan. Cachier: Edge-caching for recognition applications. In *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*, pages 276–286, June 2017.
- [13] Patrick T. Eugster, Rachid Guerraoui, Anne-Marie Kermarrec, and Laurent Massoulié. Epidemic information dissemination in distributed systems. *Computer*, 37(5):60–67, 2004.
- [14] Jose Joaquin Garcia-Luna-Aceves, Jesus Elohim Martinez-Castillo, and Rolando Menchaca-Mendez. Routing to multi-instantiated destinations: Principles, practice, and applications. *IEEE Transactions on Mobile Computing*, 17(7):1696–1709, 2018.
- [15] Peizhen Guo, Bo Hu, Rui Li, and Wenjun Hu. FoggyCache: Cross-device approximate computation reuse. In *Proceedings of the 24th Annual International Conference on Mobile Computing and Networking, MobiCom ’18*, pages 19–34, New York, NY, USA, 2018. ACM.
- [16] Harshit Gupta and Umakishore Ramachandran. Fogstore: A geo-distributed key-value store guaranteeing low latency for strongly consistent access. In *Proceedings of the 12th ACM International Conference on Distributed and Event-based Systems, DEBS ’18*, pages 148–159, New York, NY, USA, 2018. ACM.
- [17] Vassos Hadzilacos and Sam Toueg. A modular approach to fault-tolerant broadcasts and related problems. Technical report, USA, 1994.
- [18] Jonathan Hasenburger, Martin Grambow, and David Bermbach. Towards a replication service for data-intensive fog applications. In *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, pages 267–270, 2020.
- [19] Ehsan Hemmati and Jose Joaquin Garcia-Luna-Aceves. A new approach to name-based link-state routing for information-centric networks. In *Proceedings of the 2nd ACM Conference on Information-Centric Networking, ACM-ICN ’15*, page 29–38, New York, NY, USA, 2015. Association for Computing Machinery.

- 
- [20] Sebastian Henningsen, Martin Florian, Sebastian Rust, and Björn Scheuermann. Mapping the interplanetary filesystem, 2020.
- [21] A K M Mahmudul Hoque, Syed Obaid Amin, Adam Alyyan, Beichuan Zhang, Lixia Zhang, and Lan Wang. Nlsr: Named-data link state routing protocol. In *Proceedings of the 3rd ACM SIGCOMM Workshop on Information-centric Networking, ICN '13*, pages 15–20, New York, NY, USA, 2013. ACM.
- [22] Hung-Chang Hsiao and Chung-Ta King. Scoped broadcast in dynamic peer-to-peer networks. In *Proceedings of the 29th Annual International Computer Software and Applications Conference (COMPSAC 2005)*, volume 1, pages 533–538, Los Alamitos, CA, USA, jul 2005. IEEE Computer Society.
- [23] Sitaram Iyer, Antony Rowstron, and Peter Druschel. Squirrel: A decentralized peer-to-peer web cache. In *Proceedings of the Twenty-first Annual Symposium on Principles of Distributed Computing, PODC '02*, pages 213–222, New York, NY, USA, 2002. ACM.
- [24] Márk Jelasity and Ozalp Babaoglu. T-man: Gossip-based overlay topology management. volume 3910, pages 1–15, 05 2006.
- [25] Anne-Marie Kermarrec and François Taïani. Want to scale in centralized systems? think P2P. *Journal of Internet Services and Applications*, 6(1):1–12, 2015.
- [26] Simon Knight, Hung X. Nguyen, Nickolas Falkner, Rhys Bowden, and Matthew Roughan. The internet topology zoo. *IEEE Journal on Selected Areas in Communications*, 29(9):1765–1775, 2011.
- [27] Leslie Lamport. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM*, 21(7):558–565, July 1978.
- [28] Min-Yen Lue, Chung-Ta King, and Ho Fang. Scoped broadcast in structured P2P networks. In *Proceedings of the 1st International Conference on Scalable Information Systems, InfoScale '06*, page 51–es, New York, NY, USA, 2006. Association for Computing Machinery.
- [29] Petar Maymounkov and David Mazières. Kademlia: A peer-to-peer information system based on the xor metric. In Peter Druschel, Frans Kaashoek, and Antony Rowstron, editors, *Peer-to-Peer Systems*, pages 53–65, Berlin, Heidelberg, 2002. Springer Berlin Heidelberg.
- [30] Alberto Montresor and Márk Jelasity. PeerSim: A scalable P2P simulator. In *Proceedings of the 9th International Conference on Peer-to-Peer (P2P'09)*, pages 99–100, Seattle, WA, September 2009.
- [31] Michel Raynal. *Distributed algorithms for message-passing systems*, volume 500. Springer, 2013.
- [32] RFC2328. OSPF Version 2. <https://tools.ietf.org/html/rfc2328>, 1998.
- [33] Jan Seedorf, Sebastian Kiesel, and Martin Stiernerling. Traffic localization for P2P-applications: The ALTO approach. In *2009 IEEE Ninth International Conference on Peer-to-Peer Computing*, pages 171–177, Sept 2009.
- [34] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. Conflict-free Replicated Data Types. Research Report RR-7687, July 2011.

- 
- [35] Weisong Shi, Jie Cao, Quan Zhang, Youhuizi Li, and Lanyu Xu. Edge computing: Vision and challenges. *IEEE Internet of Things Journal*, 3(5):637–646, 2016.
  - [36] Liang Wang, Suzan Bayhan, Jörg Ott, Jussi Kangasharju, Arjuna Sathaseelan, and Jon Crowcroft. Pro-diluvian: Understanding scoped-flooding for content discovery in information-centric networking. In *Proceedings of the 2nd ACM Conference on Information-Centric Networking*, ACM-ICN '15, page 9–18, New York, NY, USA, 2015. Association for Computing Machinery.
  - [37] Andrew Whitaker and David Wetherall. Forwarding without loops in icarus. In *2002 IEEE Open Architectures and Network Programming Proceedings.*, pages 63–75, 2002.
  - [38] Junjie Xie, Chen Qian, Deke Guo, Minmei Wang, Shouqian Shi, and Honghui Chen. Efficient indexing mechanism for unstructured data sharing systems in edge computing. In *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*, pages 820–828, 2019.



**RESEARCH CENTRE  
RENNES – BRETAGNE ATLANTIQUE**

Campus universitaire de Beaulieu  
35042 Rennes Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399