

Making Asynchronous Distributed Computations Robust to Noise

Keren Censor-Hillel*

Ran Gelles†

Bernhard Haeupler‡

Abstract

We consider the problem of making distributed computations robust to noise, in particular to worst-case (adversarial) corruptions of messages. We give a general distributed interactive coding scheme which simulates any asynchronous distributed protocol while tolerating an optimal corruption of a $\Theta(1/n)$ fraction of all messages while incurring a moderate blowup of $O(n \log^2 n)$ in the communication complexity.

Our result is the first *fully distributed* interactive coding scheme in which the topology of the communication network is not known in advance. Prior work required either a coordinating node to be connected to all other nodes in the network or assumed a synchronous network in which all nodes already know the complete topology of the network.

*Technion, Department of Computer Science, ckeren@cs.technion.ac.il. Supported in part by the Israel Science Foundation (grant 1696/14) and the Binational Science Foundation (grant 2015803).

†Faculty of Engineering, Bar-Ilan University, ran.gelles@biu.ac.il.

‡Carnegie Mellon University, haeupler@cs.cmu.edu. Supported in part by NSF grants CCF-1527110 and CCF-1618280.

1 Introduction

Fault tolerance is one of the central challenges in the design of distributed algorithms. Typically, computation is performed by n nodes, of which some subset may be *faulty* and not behave as expected. This includes *crash* or *Byzantine* failures. Faults can also occur as communication errors, if links suffer from, e.g., *omissions*, *alterations* or *Byzantine* errors (see, e.g., [Lyn96, AW04]).

We focus on alteration errors, in which the content of sent messages may be corrupted. Previous work in the setting of faulty channels provides fault-tolerant algorithms for several specific tasks, such as the leader election or the consensus problem (e.g., [SAA95, Sin96, GLR95, SCY98]).

In this paper, we provide a general technique that takes an asynchronous distributed protocol as an input and outputs a simulation of this protocol that is resilient to noise. Specifically, we develop several tools whose combination allows us to obtain the first *fully distributed* interactive coding scheme.

The Challenge. Once communication is unreliable, even the simplest distributed tasks, such as flooding information over the network or constructing a BFS tree, become tremendously difficult to execute correctly. For instance, the asynchronous distributed Dijkstra or Bellman-Ford algorithms [Pel00] miserably fail when messages may be corrupted. To see why, recall that in the Bellman-Ford algorithm, each node sends to all of its neighbors its distance from the root. A node then sets its neighbor that is closest to the root as its parent. However, if messages are incorrect, the distance mechanism may fail and nodes may set their parents in an arbitrary way.

Our Contribution. In any attempt to tolerate message corruptions, naturally, some bound on the noise must be given. Indeed, if a majority of the sent messages are corrupted, there is no hope to complete a computation correctly. On the other hand, when the noise falls below a certain threshold, fault tolerant computation can be obtained, for example, by employing various coding techniques.

The field of *coding for interactive communication* (see, e.g., the survey of [Gel15]) considers the case where two or more parties carry some computation by sending messages to one another over noisy channels and strives to devise *coding schemes* with good guarantees. A coding scheme is a method that is given as an input a protocol π that assumes reliable channels, and outputs a noise-resilient protocol Π that simulates π . The two main measures upon which a coding scheme is evaluated are its *noise resilience*—the fraction of noise that the resilient simulation Π can withstand—and its *overhead*—the amount of redundancy Π adds in order to tolerate faults. For networks with n nodes, it is easy to show that the maximal noise fraction that any resilient protocol can cope with is $\Theta(1/n)$ [JKL15]. Indeed, if more than $(1/n)$ -fraction of the messages are corrupted, then the noise can completely corrupt all the communication of the node that sends the least number of messages. The overhead depends on the network topology, communication model, and noise resilience, as we elaborate upon in Section 1.2.

Our main result, informally stated as follows, is a deterministic coding scheme that fortifies any asynchronous protocol designed for a noise-free setting over any network topology, such that its resilient simulation withstands the maximal $\Theta(1/n)$ -fraction of noise.

Theorem 1.1. *There exists a deterministic coding scheme that takes as an input any asynchronous distributed protocol π designed for reliable channels, and outputs an asynchronous distributed protocol Π that simulates π , is resilient to a fraction $\Theta(1/n)$ of adversarially corrupted messages, and has a multiplicative communication overhead of $O(n \log^2 n)$.*

1.1 Techniques

A Content-Oblivious BFS Construction. A key ingredient in our coding scheme is a BFS construction which is *content oblivious*. That is, in our BFS construction, the nodes send messages to each other and *ignore their content*, basing their decisions only on the order of received messages. The challenge is to be able to do this despite asynchrony and despite lack of FIFO assumptions. In a sense, our construction can be seen as a variant of the distributed Dijkstra algorithm, with the property that the nodes send “empty messages” that contain no information (alternatively, the nodes ignore the content of received messages).

Recall that the distributed Dijkstra algorithm, see, e.g., [Pel00, Chapter 5], is initiated by some node r , which governs the BFS construction layer by layer, where the construction of each layer is called a *phase*. The invariant is that after the p -th phase, the algorithm has constructed a BFS tree T_p of depth p rooted at r , where all nodes in T_p know their parent and children in T_p . The base case is $T_0 = \{r\}$, and the construction of the first layer is as follows. The node r sends an EXPLORE message to all its neighbors, who in turn set r as their parent. Each EXPLORE message is replied to with an ACK message. Once r receives ACK messages from all of its neighbors, the first phase ends and the construction of the second layer begins. Note that T_1 indeed holds r and all of its neighbors.

For the p -th phase, the root floods a message PHASE through T_{p-1} . Once a leaf in T_{p-1} receives a PHASE message, it sends EXPLORE to all of its neighbors, who in turn set their parent unless already in T_{p-1} . Each node that receives an EXPLORE replies with an ACK and an indication of its parent node, so that the exploring node learns which of its neighbors is a child and which is a sibling. Upon receiving an ACK from all of its neighbors, the node sends an ACK to its parent, which propagates it all the way to r . Once r has received ACK messages from all of its children, the phase is complete.

Our content-oblivious BFS construction imitates the above behavior while using only a “single” type of message, instead of PHASE, EXPLORE and ACK messages. Specifically, the construction begins with r sending a message (EXPLORE¹) to all of its neighbors, who in turn set r as their parent and reply with a message (ACK). When r receives a message from all of its neighbors, the first phase is complete. Then, r begins the second phase by sending another message (EXPLORE/PHASE) to all of its neighbors. This message causes a node that has already set its parent to behave like r —it sends a message to all of its neighbors (EXPLORE) *except for its parent*. After receiving a message (ACK) from all of its neighbors, it sends a message (ACK) to its parent.

One can easily verify that this approach behaves similarly to the Dijkstra algorithm described above, in the sense that every node sets its parent correctly. The only difference is when a node u sends an (EXPLORE) message to its sibling w . In the Dijkstra algorithm the sibling w replies by telling the exploring node u that they are siblings (by indicating the parent of w , which is not u). However, in our case messages contain no content and u is unable to distinguish whether w is a child or a sibling, since in both cases w should reply to the EXPLORE message in the same way.

Our insight is that serializing each phase provides a solution to the above ambiguity. That is, we let r send a message (EXPLORE/PHASE) to one child at a time, waiting to receive a message (ACK) from that child before sending a message (EXPLORE/PHASE) to the next child. This gives that if a node is expecting a message (EXPLORE) from its parent but instead it receives a message (EXPLORE) from a non-parent neighbor, then it knows that this neighbor must be a sibling. Hence, the node can mark all siblings and distinguish them from its children.

¹To ease the readability, we write in parenthesis the functionality of each sent message, but we emphasize that messages in our construction contain no content at all, and the labels of EXPLORE and ACK are given only for the analysis.

The main advantage of not basing our construction on the content of received messages is that the obtained BFS construction is *inherently tolerant against message corruptions*: the noise has no effect on the construction since the content of the communicated messages is already being ignored. Notice that in our construction, the nodes do not learn their distance from r , in contrast to what can easily be obtained in the noise-free case. However, this will suffice for our usage of the BFS tree in our coding scheme.

Interactive Coding over Sparse Subgraphs. A crucial framework we rely on in our simulation is a multiparty coding scheme for interactive communication by Hoza and Schulman [HS16], which is in turn based on ideas from [RS94]. This coding scheme allows simulating protocols over any graph $G = (V, E)$ and withstands an $\Theta(1/|E|)$ -fraction of adversarial message corruption, while incurring a *constant* communication overhead. The caveat of using this scheme for our simulation is that it applies only for *synchronous* protocols that communicate over G in a manner which we call *fully-utilized synchronous*: in each round, every node communicates one symbol over to *each* of its neighbors.

In order to obtain our coding scheme for asynchronous protocol with resilience $\Theta(1/n)$, we first convert the asynchronous input protocol π into a fully-utilized synchronous protocol defined over some subgraph $G' = (V, E')$ of G with $|E'| = \Theta(n)$. To this end, we use the BFS tree constructed by our content-oblivious method described above. Once we obtain a BFS tree \mathcal{T} , we simulate each message communicated by π via n fully-utilized synchronous rounds over the tree \mathcal{T} . During each of such n rounds, a message of π is flooded throughout \mathcal{T} until it reaches all the nodes and, in particular, its destination node. Note that in every round, all nodes send messages over all the edges of \mathcal{T} . This implies a communication overhead of $O(n^2 \log n)$: we have at most n rounds with $\Theta(n)$ messages per round. The $\log n$ term stems from adding the identity of the source node and the destination node to each flooded message.²

Using the Hoza and Schulman [HS16] coding scheme taking as an input the fully-utilized synchronous protocol defined over the topology \mathcal{T} gives a resilient simulation of π which withstands a maximal $\Theta(1/n)$ -fraction of corrupted messages. Alas, it is a synchronous simulation, while our environment is asynchronous. Hence, to complete our simulation, we need to use a *synchronizer* [Awe85].

A Root-Triggered Synchronizer. In the original error-free setting, if the input protocol to a synchronizer is guaranteed to be fully-utilized then synchronization is trivial. Each node simply attaches a round number to each of its outgoing messages and produces the outgoing messages for round $i + 1$ only after receiving messages for round i from all of its neighbors. The key difficulty is then for non-fully utilized synchronous input algorithms, in which a node cannot simply wait to receive a message for round i from all of its neighbors, as it may be the case that some of these do not exist.

In our setting, we guarantee that we produce a fully-utilized synchronous algorithm as an input to our synchronizer. However, we do not assume FIFO channels, which means that we cannot rely on the naive synchronizer, despite the promise of a fully-utilized synchronous protocol for an input. Thus, we need a different solution for synchronizing the messages, and our approach is based on having a single node responsible for triggering messages of each round only after the previous round

²Throughout this work, all logarithms are taken to base 2.

has been simulated by all nodes. To this end, our synchronizer bears similarity to the classic tree-based synchronizer of Awerbuch [Awe85], with the difference that it does not incur any message overhead because it is given a fully-utilized synchronous input.

A Spanner-Based Coding Scheme. We show that our coding technique described above can be further improved. Routing each message over a tree \mathcal{T} requires n rounds in the worst case for a message to reach its destination. A more efficient solution would be to route each message through a spanning subgraph $S = (V, E_S)$ of G in which the distance over S of every $(u, v) \in E$ is not too large. On the other hand, the Hoza-Schulman coding scheme on S has a noise resilience of $\Theta(1/|E_S|)$, and hence we require $|E_S|$ to be $O(n)$ in order to maintain a maximal resilience level of $\Theta(1/n)$. Luckily, for every G there exist sparse spanning subgraphs in which $|E_S| = O(n)$ while every two neighbors in G are at distance at most $O(\log n)$ in S ; such subgraphs are known as $O(\log n)$ -spanners [Pel00, PS89].

Flooding a message of π from u to v can be done within $O(\log n)$ rounds, in each of which $O(|E_S|) = O(n)$ messages are sent by a fully-utilized synchronous simulation of π , leading to our claimed communication overhead of $O(n \log^2 n)$. Here again, the extra $\log n$ term stems from adding identifiers to each flooded message.

However, flooding information over a spanner introduces several other difficulties. For instance, in contrast to the case of a tree, it is not guaranteed anymore that each message arrives only once to its destination—indeed, multiple paths may exist between any two nodes. Furthermore, when multiple nodes send messages, the congestion may cause super-polynomial delays if a simple flooding algorithm is used. Then, due to having multiple paths with arbitrary delays, messages may arrive to their destination out of order. Since the delay is super-polynomial in the worst case, adding a counter to each message will increase the overhead by $\omega(\log n)$ and damage the global overhead.

Instead, we provide a contention-resolution flavored technique, which consists of priority-based windows for delivering the messages. In more detail, a message flooding starts only at the beginning of an $O(\log n)$ -round window. Multiple messages that are sent during the same window may be dropped during their flooding, yet the source always learns when its message is dropped, so it can retransmit the message in the next window. A similar approach is well-known for constructing a BFS tree when no specific root is given, but our extension of this technique is more involved, since dropped messages *must be resent*.

It remains to explain how to construct the $O(\log n)$ -spanner over the noisy network to begin with. For this, we use our previously described tree-based coding scheme to simulate a distributed spanner construction, e.g., the (noiseless) construction of Derbel, Mosbah, and Zemmari [DMZ10]. While coding this part incurs a large overhead of $O(n^2 \log n)$, this overhead applies only to the part of constructing the spanner, and the global overhead of our coding scheme is still dominated by the overhead of coding the input protocol over the spanner.

1.2 Related Work

Performing computations over noisy channels is the heart of *coding for interactive communication*, initiated by Schulman [Sch92, Sch96]. A long line of work considers the 2-party case and obtains various coding schemes, as well as bounds on their capabilities in various settings and noise models [BR14, BE14, BKN14, GHS14, FGOS15, EGH16, KR13, Hae14, BGMO16, GHK⁺16]. See [Gel15] for a survey on interactive coding.

Interactive coding in the multiparty setting was first considered by Rajagopalan and Schulman [RS94] for the case of random noise, where every bit is flipped with some fixed probability. Rajagopalan and Schulman show, for any topology G , a coding scheme with an overhead of $O(\log(d + 1))$, where d is the maximal degree of G . Gelles, Moitra and Sahai [GMS14] provide an efficient extension to that scheme. Alon et al. [ABE⁺16] show a coding scheme with an overhead of $O(1)$ for d -regular graphs with degree $d = n^{\Omega(1)}$. Braverman et al. [BEGH16] demonstrate a lower bound of $\Omega(\log n)$ on the communication over a star graph. All the above works assume fully-utilized synchronous protocols, in which the protocol works in rounds and in every round all nodes communicate on all the channels connected to them. Gelles and Kalai [GK17] show that if nodes are not required to speak at every round, a lower bound of $\Omega(\log n)$ on the overhead can be proved even for coding schemes over graphs with small degree, e.g., $d = 2$.

In the case of adversarial noise, Jain, Kalai and Lewko [JKL15] show a coding scheme that is resilient to a noise fraction of $\Theta(1/n)$ and has an overhead of $O(1)$ in networks which contain a star as a subgraph. Lewko and Vitercik [LV15] improve the communication balance of that scheme. Hoza and Schulman [HS16] consider fully-utilized synchronous protocols on arbitrary graphs and show a coding with resilience $\Theta(1/|E|)$ and overhead $O(1)$. If the topology of the network G is known to all nodes, the nodes can route messages over a sparser spanning graph and decrease the number of edges used by the coding scheme. In this case, Hoza and Schulman show a coding scheme with a maximal resilience level of $\Theta(1/n)$ and an overhead of $O((|E| \log n)/n)$.

Previous work in distributed settings that allow edge failures are typically different from our setting in various aspects. Most notable are synchrony assumptions, complete communication graphs, addressing specific distributed tasks and assuming a bound on the number of links that may exhibit failures. This is in contrast to our work, which addresses an asynchronous setting with an arbitrary topology, and considers the simulation of any distributed task. In particular, all links may send corrupted messages, with the bound being the number of corruptions rather than the number of faulty links. For instance, Singh [Sin96] and Sayeed, Abu-Amara and Abu-Amara [SAA95] consider the specific task of leader election and agreement for complete networks. Gong, Lincoln, and Rushby [GLR95], Siu, Chin and Yang [SCY98] and Dasgupta [Das98] consider agreement in complete synchronous networks with both faulty nodes and faulty links.

Pelc [Pel92] shows that if the number of Byzantine-corrupted links is bounded by t , reliable communication is achievable only over graphs whose connectivity is more than $2t$. The same work also considers the case where each link is faulty with some probability. In a more recent work, Feinerman, Haeupler and Korman [FHK14] also address complete synchronous networks, and study the specific problems of broadcast and majority consensus under random errors.

Synchronizers for unreliable settings have been studied by Awerbuch et al. [APPS92], which address a dynamic setting, and by Harrington and Somani [HS94], which assume faulty nodes.

1.3 Organization

We define basic notations, our communication and noise model as well as the notion of noise resilient computations (i.e., coding schemes) in Section 2. In Section 3 we describe our content-oblivious BFS construction. A coding scheme over a spanning tree with overhead $O(n^2 \log n)$ is provided in Section 4. Finally, a coding scheme based on an underlying spanner with an improved overhead of $O(n \log^2 n)$ is provided in Section 5.

2 Preliminaries

Throughout this work we assume a network described by a graph $G = (V, E)$ with $n = |V|$ nodes and $m = |E|$ edges. Each node $u \in V$ is a party that participates in the computation and each edge $(u, v) \in E$ is a bi-directional communication channel between nodes u and v . The task of the nodes is to conduct some distributed computation given by a deterministic³ protocol π , which consists of the algorithm each node (locally) runs. In particular, the protocol dictates to each node which messages to send to which neighbor as a function of all previous communication (and possibly the node's identity, private randomness and private input, if exists). The *communication complexity* of the protocol, $CC(\pi)$, is the maximal number of bits communicated by all nodes in any instance of π . The *message complexity* of π is the maximal number of message sent by all nodes in any instance of π .

We assume that the topology of G is known only locally, namely, each node v knows only the set \mathcal{N}_v of identities of its own neighbors. However, the size of the network n is known to all nodes.

Communication Models. Our protocols are for the *Asynchronous* communication model defined below. In addition, we describe a different communication model named the *Fully-Utilized Synchronous Model*, which is common in previous interactive coding work [RS94, HS16, ABE⁺16, BEGH16]. In particular, we use coding schemes defined in the fully-utilized synchronous model (specifically, [HS16]) as primitives for encoding our asynchronous protocols (see Lemma 2.1 below).

- *Asynchronous Model.* In this setting, there are no timing assumptions. We assume each node is asleep until receiving a message. Once a message is received, the receiver wakes up, performs some local computation, transmits one or more messages to one or more adjacent nodes and goes back to sleep. Messages can be of any length. A protocol starts by waking up a single node r of its choice.
- *The Fully-Utilized Synchronous Model.* Communication in this model works in synchronous rounds, determined by a global clock. At every clock tick, every node sends one symbol (from some fixed alphabet Σ) on each and every one of the communication links connected to it. That is, at every round exactly $2m$ symbols are being communicated.

Adversary. We assume an all-powerful adversary that knows the network G , the protocol π and the private inputs of the nodes (if there are any). The adversary is able to (a) corrupt messages by changing the content of a transmitted message and (b) rush or delay the delivery of messages by an unbounded but finite amount of time. We restrict the number of messages that the adversary can corrupt, namely, we assume that the adversary can corrupt at most some fixed fraction μ of the communicated messages. We do not restrict how a message can be corrupted and, in particular, the adversary may replace a sent message M with any other message M' of any length and content. However, our coding scheme will have the invariant that each message contains a single symbol (from a given alphabet Σ), thus a message corruption will be equivalent to corrupting a single symbol. Note that the adversary is *not* allowed to inject new messages or completely delete existing messages.⁴

³While we focus here on deterministic protocols, ours result also apply to randomized Monte-Carlo protocols.

⁴This type of noise, commonly called *insertion and deletion* noise, is known to be more difficult to deal with in the interactive setting [BGMO16] and may be destructive for asynchronous protocols [FLP85].

Protocol Simulation, Resilience, and Overhead. A protocol Π is said to *simulate* π , if after the completion of Π , each node outputs the transcript it would have seen when running π assuming noiseless channels. The protocol Π is *resilient* to a μ fraction of noise, if Π succeeds in simulating π even if an all powerful adversary completely corrupts up to a fraction μ of the messages communicated by Π . The *overhead* of Π with respect to π is defined by $overhead(\Pi \mid \pi) = CC(\Pi)/CC(\pi)$.

A coding scheme $\mathcal{C} : \pi \rightarrow \Pi$ converts any input protocol π into a resilient version $\Pi = \mathcal{C}(\pi)$. The resilience of a coding scheme is the minimal resilience of any simulation generated by the coding scheme. The (asymptotic) overhead of a coding scheme considers the maximal overhead for the worst input protocol π when $CC(\pi)$ tends to infinity. Namely,

$$overhead(\mathcal{C}) = \limsup_{c \rightarrow \infty} \max_{\substack{\pi \text{ s.t.} \\ CC(\pi) = c}} overhead(\mathcal{C}(\pi) \mid \pi).$$

We are mainly interested in how the overhead scales with n and m .

A famous multiparty coding scheme in the fully-utilized synchronous model, shown by Hoza and Schulman [HS16] (based on a previous scheme [RS94]), provides a coding scheme that simulates any noiseless fully-utilized synchronous protocol π defined over some topology G with resilience $\Theta(1/m)$ and a constant overhead $O(1)$.

Lemma 2.1 ([HS16]). *In the fully-utilized synchronous model, any T -round protocol π can be simulated by a protocol $\Pi = HS(\pi)$ with round complexity $O(T)$ and communication complexity $O(CC(\pi))$ that is resilient to adversarial corruption of up to an $\Theta(1/m)$ fraction of the messages.*

3 A Distributed Content-Oblivious BFS Algorithm

In this section we show a distributed construction of a BFS tree using messages whose content can be arbitrary. We call this a *content-oblivious* construction. Our algorithm can be seen as a variant of a simple distributed layered-BFS algorithm, see, e.g., [Gal82, Pel00, Tel00].

3.1 The BFS Algorithm: Description

The BFS construction is initiated by one designated node r we call here the *root*. The construction builds the tree layer by layer. First, the root sends a message to all of its neighbors. This triggers its neighbors to set r as their parent. Each such a neighbor replies a message to r to acknowledge that it has received r 's message. Once r has received a message from all of its neighbors, it knows that the first layer is completed, and all nodes with distance 1 have set r as their parent. We call the above an EXPLORE step.

The root then begins a second EXPLORE which causes all nodes at distance 2 to set their parent and connect to the BFS tree. Specifically, the root sends a message to each of its children and waits until all children reply a message to indicate they are done. However, in contrast to previous distributed BFS algorithms, messages are sent *sequentially*—the root sends a message to its next child only after receiving the acknowledgement message from its previous child.

When a node v that has already set its parent $parent_v$ receives a message *from its parent*, it acts as a root and invokes an EXPLORE: it sends a message to all of its neighbors excluding $parent_v$ and waits until they all send a message back. Only then v sends a message to its parent to indicate its

EXPLORE process has completed. It is easy to see that when the root completes its k -th EXPLORE, all nodes within distance at most k have set their parent and connected to the BFS tree.

A special treatment is needed when a node u receives a message from a node v who is *not* the parent of u during a time at which u is not in the middle of an EXPLORE step. That is, u is not expecting any messages from its neighbors, except for its parent that may trigger it to initiate another EXPLORE step. Recalling that messages are sent to children in a sequential manner, it is easy to verify that such a message delivery may happen only when v has received a message from its own parent and is now processing its own EXPLORE. That is, such a message indicates that v is a *sibling* of u in the BFS tree (namely, v is not a parent nor a child of u in the BFS tree). Thus, upon receiving such a message, u marks v as a sibling and removes it from its list of children. To simplify the presentation, as we elaborate in Remark 1, in next exploration steps u will keep sending messages to v as if it was one of its children.

One additional property that we require from our BFS construction is that all the nodes complete the algorithm *at the same time*. As explained in the introduction, we use this construction as an initial part for our coding scheme. Furthermore, recall that in order to be noise-resilient, during the BFS construction the nodes ignore the content of the messages and their entire behavior is based on whether or not a message was received. However, once this construction is complete, the nodes send and receive messages according to the coding scheme and it is crucial that a node is able to distinguish messages that belong to the BFS construction from messages of the coding scheme.

We solve this issue by making sure that each node participates in exactly n steps of EXPLORE. Once the node has sent the n -th acknowledgement to its parent, the node knows that the next message *from the parent* belongs to the coding scheme rather than to the BFS construction.⁵ To make sure that each node participates in exactly n EXPLORE steps, regardless of its distance from r , we let every node initiate one additional EXPLORE, which we refer to as a *dummy* EXPLORE. Specifically, when a node completes its $(n - 1)$ -th EXPLORE, and *before the node sends the acknowledgement back to its parent*, it invokes another EXPLORE step. Now, just by counting the messages received from the parent, every node knows whether the BFS construction has completed or not.

The pseudocode of the BFS construction is given in Algorithm 1(a) and Algorithm 1(b).

3.2 The BFS Algorithm: Analysis

In this section we analyze Algorithm 1 and show that it satisfies the following properties.

Theorem 3.1. *For any input $G = (V, E)$ and node $r \in V$, Algorithm 1 finds a BFS tree \mathcal{T} with root r . Specifically, each node knows its parent in T and all of its adjacent edges that belong to \mathcal{T} . The algorithm communicates $O(nm)$ messages, where no payload is needed in any messages.*

Furthermore, we show that all nodes know that the BFS construction is complete, in the following sense.

Claim 3.2. *At the end of Algorithm 1 all nodes are in state DONE. Moreover, if r is in state DONE then all other nodes are in state DONE as well.*

Proof. (Theorem 3.1) Let \mathcal{T} be a graph on the nodes V defined at the end of Algorithm 1 in the following manner: If $v = \text{parent}_u$, then (u, v) is an edge in \mathcal{T} . We begin by proving that \mathcal{T} is a spanning tree. This is implied by the following claim.

⁵Note that additional messages may arrive from a sibling node for the BFS construction but still, the next message arriving from the *parent* belongs to the coding scheme rather than the BFS construction.

Algorithm 1(a) Content-oblivious BFS construction: Main Algorithm

Initialization: All nodes begin in the INIT state.

```
1: For node  $r$  designated as root:
2: Begin
3:    $parent_r \leftarrow \perp$ 
4:    $children_r \leftarrow \mathcal{N}_r$ 
5:    $count_r \leftarrow 0$ 
6:    $state_r \leftarrow \text{IDLE}$ 

7:   while  $state_r \neq \text{DONE}$  do                                 $\triangleright$  Perform  $n$  instances of EXPLORE
8:      $r$  invokes EXPLORE
9:   end while
10: End
```

Claim 3.3. *At the end of the k -th invocation of the root's EXPLORE step, all the nodes that are at distance k from r set their parent to a node with distance $k - 1$ from r and move to the state IDLE, and every node of distance larger than k from r is in state INIT.*

Proof. We prove the claim by induction on k . The base case $k = 1$ follows since in the first EXPLORE invocation all of r 's children run SETPARENT, setting r as their parent, and switch to IDLE. They send message only back to r , hence all other nodes remain in INIT.

Assume that the claim holds for the k -th invocation and consider the $(k + 1)$ -th invocation of EXPLORE by r . Messages propagating along the BFS tree cause all nodes of distance at most k to invoke EXPLORE (in some order). This triggers a message to every node of distance $k + 1$, which causes it to switch its state to IDLE and set its parent to the first node (of distance k) that sent it a message. Note that nodes of distance $k + 1$ only communicate back to their parent and do not invoke EXPLORE at this time, so nodes of distance larger than $k + 1$ remain in state INIT. At the end of the invocation each EXPLORE, the invoking node switches back to state IDLE. \square

Next, we prove that each node learns which neighbors are its children and which are not. First note that if $v = parent_u$ then u only sends v a message as a reply to a prior message received from v (i.e., an "ACK" message at the end of an EXPLORE). Therefore, whenever v receives a message from u it is in state EXPLORE, and such a message can never invoke the procedure MARKSIBLING. It follows that at the end of the algorithm $v \in children_u$.

Next, assume (u, v) is an edge in G but not in \mathcal{T} . We show that at the end of the algorithm $v \notin children_u$ and $u \notin children_v$. Let t be the first time after which both u and v have invoked SETPARENT. We claim that both u and v invoke EXPLORE after time t . This is because time t is within the execution of an EXPLORE step invoked by r and before Line 19 of that execution, and hence for every node $w \neq r$ there is a time $t_w > t$ during the execution of the loop in Lines 19–23 for r in which w invokes EXPLORE.

Finally, we note that since (u, v) is an edge in G but not in \mathcal{T} , then neither u is an ancestor of v in \mathcal{T} nor v is an ancestor of u in \mathcal{T} . This implies that when v invokes EXPLORE then u is in state IDLE, which causes it to invoke MARKSIBLING and hence $v \notin children_u$. The proof for $u \notin children_v$ is exactly the same.

Algorithm 1(b) Content-oblivious BFS construction: Message Handling Procedures

For every node u in state INIT upon receiving a message from node v

```

1: procedure SETPARENT
2:    $parent_u \leftarrow v$ 
3:    $children_u \leftarrow \mathcal{N}_u \setminus \{v\}$ 
4:    $count_u \leftarrow 0$ 
5:    $state_u \leftarrow \text{IDLE}$ 
6:   send a message to  $v$  ▷ an “ACK” message
7: end procedure

```

For every node u in state IDLE/DONE upon receiving a message from $v \neq parent_u$

```

8: procedure MARKSIBLING
9:    $children_u \leftarrow children_u \setminus \{v\}$ 
10:  send a message to  $v$  ▷ an “ACK” message
11: end procedure

```

For every node u in state IDLE upon receiving a message from $parent_u$

```

12: procedure EXPLORE
13:    $state_u \leftarrow \text{EXPLORE}$ 
14:    $count_u \leftarrow count_u + 1$ 
15:   for all  $v \in \mathcal{N}_u \setminus \{parent_u\}$  do ▷ note: for is sequential
16:     send a message to  $v$  ▷ an “Explore” message
17:     wait until a message is received from  $v$ 
18:   end for

19:   if  $count_u = n - 1$  then ▷ Extra dummy EXPLORE
20:     for all  $v \in children_u$  do
21:       send a message to  $v$ 
22:       wait until a message is received from  $v$ 
23:     end for
24:   end if

25:  send a message to  $parent_u$  ▷ an “ACK” message

26:  if  $count_u = n - 1$  then ▷ Change state to DONE if completed; otherwise, back to IDLE
27:     $state_u \leftarrow \text{DONE}$ 
28:  else
29:     $state_u \leftarrow \text{IDLE}$ 
30:  end if
31: end procedure

```

Finally let us analyze the message complexity. In Algorithm 1 each node invokes EXPLORE for n times (see also the proof of Claim 3.2 below), where during each EXPLORE it sends a message on each edge. Therefore, there are $O(n)$ messages sent on each one of the m edges, which amounts to a total message complexity of $O(nm) = O(|V| \cdot |E|)$. \square

Remark 1. *It is possible to reduce the message complexity by sending EXPLORE messages only to children_v nodes. However, this must be delayed at least one EXPLORE step, beyond the point in time where all the neighbors have completed their first EXPLORE (in order to be able to identify siblings). The new message complexity will be $O(|V|^2 + |E|)$. For simplicity, we avoid this optimization and assume EXPLORE messages are sent to all non-parent nodes all the time, incurring a message complexity of $O(|V| \cdot |E|)$.*

We now prove Claim 3.2. This property is important in particular for the next section, as it suggests that there is a point in time (known by the root), when all nodes have completed their BFS algorithm. In hindsight, this allows to distinguish messages that are part of the BFS construction, whose content is ignored, from messages of the coding scheme, whose content is meaningful and must not be ignored.

Proof. (Claim 3.2) Note that the EXPLORE procedure works in an DFS manner: a node replies an ACK to its parent only after all of its children reply an ACK to it. Similarly, the root completes an EXPLORE step after receiving an ACK from all its children, which means that they have all completed their EXPLORE steps.

Note that each node invokes exactly n EXPLORE steps due to the dummy EXPLORE step initiated in Line 19. To see this, consider the same algorithm without the extra EXPLORE in Lines 19–23 and note that nodes at distance k from the root r invoke exactly $n - k$ EXPLORE steps. Adding this extra EXPLORE step at every node makes all nodes invoke EXPLORE exactly n times. Specifically, during the n -th invocation of EXPLORE by r , every node with distance 1 from r invokes its $(n - 1)$ -th EXPLORE step, and then, *before sending an ACK to r* in Line 30, it invoke its n -th EXPLORE step. This then continues in an inductive manner all the way to the leaves.

Only once all of its children have sent an ACK and thus terminated the protocol and switched to DONE, a node replies with an ACK to its parent and changes its state to DONE. It follows that when the root receives an ACK for the n -th EXPLORE step from all of its children, all the nodes have terminated the protocol and switched state to DONE. \square

4 A Distributed Interactive Coding Scheme

In this section we show how to simulate any asynchronous protocol over a noisy network whose topology is unknown in advance. Our main theorem for this part is the following.

Theorem 4.1. *Any asynchronous protocol π over a network G can be simulated by an asynchronous protocol Π resilient to an $\Theta(1/n)$ -fraction of adversarial message corruption, and it holds that $\text{CC}(\Pi) = O(nm \log n) + \text{CC}(\pi) \cdot O(n^2 \log n)$.*

4.1 Obtaining a fully-utilized synchronous protocol from an asynchronous input protocol π

The first ingredient we need is a way to transform an asynchronous protocol into a fully-utilized synchronous protocol, in order to be able to use the Hoza-Schulman coding scheme. This trans-

formation does not need to be robust to noise, as it is not going to be executed as is, but we will rather encode the fully-utilized synchronous protocol and execute the robust version. Later, we transform it back into the asynchronous setting using a synchronizer that is robust to noise.

Recall that in a fully-utilized synchronous protocol nodes operate in rounds, where at each round every node communicates one symbol (from some fixed alphabet Σ) on each communication channel connected to it. We will assume the alphabet is large enough to convey all the information that our coding scheme needs. In particular, we assume each symbol contains $O(\log n)$ bits.

Remark 2. *In the following, we assume the network G is composed of channels with a fixed alphabet Σ of size $\text{poly}(n)$. That is, each symbol contains $O(\log n)$ bits.*

In order to avoid confusion, we will use the term “symbols” for messages sent by the coding scheme, while using “messages” to indicate the information sent by the noiseless protocol π .

The construction of our transformation into a fully-utilized synchronous protocol is given in Algorithm 2. In this construction, each node u maintains a queue of symbols that it needs to relay throughout a locally known spanning tree \mathcal{T} . The queue is initialized with the bits of any message that u needs to send according to the input protocol π , where each bit is encapsulated in a symbol that contains the bit value, the identity of the source (i.e., of u), and the identity of the destination node. Every symbol received by u is pushed into its queue, and relayed to u 's neighbors in future rounds. In particular, upon receiving the symbol $(src, dest, val)$ from a node w , the node u pushes the vector $(src, dest, val, w)$ to its queue. If u is the destination node, it does not push the symbol into its queue; instead, u collects this bit for decoding the message.

The transformation works by having each node pop a record from its queue in each round and send the obtained triplet to all of its neighbors in \mathcal{T} except for the node w from which the message was received. If the queue is empty then an empty message is sent to all neighbors in \mathcal{T} .

Note that all fragments of a message are received in order at the destination, since \mathcal{T} has no cycles. Therefore, we can assume that the protocol sends a predefined symbol that indicates the end of the message, in order to avoid an assumption of knowledge of the message length. This ensures that Line 17 is well-defined. Our transformation guarantees the following.

Lemma 4.2. *Algorithm 2 creates a fully-utilized synchronous protocol π' that simulates π , in the sense that all messages of π are sent and received. The simulation π' has a communication overhead of $O(n^2 \log n)$ with respect to π , and a message complexity of $\text{CC}(\pi) \cdot O(n^2)$.*

Proof. By construction, every node sends a symbol to all of its neighbors in each round and hence Algorithm 2 is a fully-utilized synchronous protocol. In addition, eventually every messages of π reaches its destination and hence the obtained fully-utilized synchronous protocol simulates π . For the communication overhead, note that $O(\log n)$ bits of the identities of source and destination are appended to each bit sent by π ; that is, a symbol size of $O(\log n)$ bits suffices. In addition, a delivery of a single message of π may require $O(n)$ rounds of relaying symbols sent along the tree \mathcal{T} . In each such round there are $O(n)$ symbols that are sent since the obtained protocol is a fully-utilized synchronous protocol. This implies that $O(n^2)$ symbols are communicated per each bit of π and gives a total communication overhead of $O(n^2 \log n)$.

Note that this is a worst-case analysis that assumes a single bit travels within the network at each time so that another bit is sent only after a previous bit reached its destination. If several bits are sent consecutively or if several nodes send bits simultaneously, the resulting number of messages can only decrease. \square

Algorithm 2 Simulating an asynchronous protocol π by a fully-utilized synchronous protocol π' .

Initialization: Given is a BFS tree \mathcal{T} rooted at r .

```

1: In every round, for every node  $u$ :
2: Begin
3:   for every node  $v$  do
4:     Let  $M_1 \cdots M_\ell$  be the bit representation of a message  $M$  that  $u$  has to send to  $v$  in  $\pi$ .
5:     Push  $(u, v, M_1, \perp), \dots, (u, v, M_\ell, \perp)$  into  $queue_u$ 
6:   end for
7:    $(src, dest, val, w) \leftarrow$  pop item out of  $queue_u$ 
8:   if  $(src, dest, val, w)$  is not empty then
9:     send  $(src, dest, val)$  to every  $v \in \mathcal{N}_u(\mathcal{T}) \setminus \{w\}$  and send  $\perp$  to  $w$ 
10:  else
11:    send  $\perp$  to every  $v \in \mathcal{N}_u(\mathcal{T})$ 
12:  end if
13:  For every message  $(src, dest, val)$  received from  $w$ :
14:  if  $dest \neq u$  then
15:    push  $(src, dest, val, w)$  into  $queue_u$ 
16:  else
17:    collect the bits  $val$  for decoding  $M$ 
18:  end if
19: End

```

4.2 Root-triggered synchronizers

We now describe our root-triggered synchronizer, which we use in order to execute the resilient synchronous protocol (which can be obtained by using the Hoza-Schulman coding scheme) in our asynchronous setting. We constructed a tree-based synchronizer as in Awerbuch [Awe85]. The synchronizer gets as an input a fully-utilized synchronous protocol Π' and outputs an equivalent asynchronous protocol Π that simulates Π' round by round.

We first describe our simulation of a single round of Π' over a *tree*. Our synchronizer works as follows. The root initiates the process by sending its messages, determined by Π' , to its children. This triggers its children to send their messages to their children, but not yet to their parent, and so forth, so that messages propagate all the way to the leaves. Once a leaf receives a message, it sends its message to its parent, and similarly, any node which receives a message from all of its children sends its message to its parent. This continues inductively all the way back to the root, which eventually receives messages from all of its children and complete the simulation of this round of Π' .

We build upon the above idea in order to simulate a fully-utilized synchronous algorithm Π' over an arbitrary graph S . That is, each node u has a message m_{uv} designated to each one of its neighbors $v \in \mathcal{N}_u(S)$.⁶

The pseudocode is given in Algorithm 3. We single out a node r , which we refer to as the *initiator*, which starts by sending a message to all of its neighbors in S . This triggers each neighboring node to send its messages to its neighbors, but not yet to its parent, which is now simply

⁶Later, in Section 5, we apply our root-triggered synchronizer to an input protocol on G which is fully-utilized on a spanning subgraph S of G .

the neighbor from which it receives the *first* message. This continues inductively, and only when a node receives messages from all of its neighbors it sends its message to its parent. Eventually, the initiator receives messages from all of its neighbors and completes the simulation of the round.

Algorithm 3 A root-triggered synchronizer for a fully-utilized synchronous protocol Π' over a graph S .

Initialization: All nodes begin in the INIT state.

- 1: For node r designated as initiator:
 - 2: **Begin**
 - 3: $state_r \leftarrow \text{ACTIVE}$
 - 4: $parent_r \leftarrow \perp$
 - 5: $children_r \leftarrow \mathcal{N}_r(S)$
 - 6: r sends m_{rv} to each node $v \in children_r$
 - 7: r waits to receive a message m_{vr} from every node $v \in children_r$
 - 8: $state_r \leftarrow \text{DONE}$
 - 9: **End**
 - 10: For every node u , upon receiving a message from w when in state INIT:
 - 11: **Begin**
 - 12: $state_u \leftarrow \text{ACTIVE}$
 - 13: $parent_u \leftarrow w$
 - 14: $children_u \leftarrow \mathcal{N}_r(S) \setminus \{w\}$
 - 15: u sends m_{uv} to each node $v \in children_u$
 - 16: u waits to receive a message m_{vu} from every node $v \in children_u$
 - 17: u sends m_{uw} to w
 - 18: $state_u \leftarrow \text{DONE}$
 - 19: **End**
-

We prove the following properties of Algorithm 3.

Lemma 4.3. *By the end of Algorithm 3 each node u receives the messages m_{vu} from every node $v \in \mathcal{N}_u(S)$, and all nodes are in state DONE.*

Proof. Let T denote the tree rooted at r that is induced by the edges of S that connect each node u with $parent_u$. By construction, each node $u \neq r$ sets its parent to be the first node from which it receives a messages and hence u sets exactly one node as its parent in an acyclic manner, inducing the tree T .

We prove by induction on the height of the nodes with respect to T , that each node u receives the messages m_{vu} from every node $v \in \mathcal{N}_u(S)$ and then switches its state to DONE. Note that every node sends its messages to all of its neighbors so that eventually all such messages arrive, and we only need to verify that the message from u to $parent_u$ is eventually sent.

The base case is for the leaves of T , which indeed receive messages from all of their neighbors since the only messages that get delayed are messages from nodes to their parents. Assume this holds for all nodes at height h , and consider a node u at height $h + 1$. Node u receives messages

from all of its siblings in the tree. By the induction hypothesis, every child v of u in T receives all of its messages and switches to state DONE. This implies that in between, node v sends its message m_{vu} to its parent u . When this happens for all nodes $v \in \text{children}_u$ it is the case that u receives the messages m_{vu} from every node $v \in \mathcal{N}_u(S)$ and then switches its state to DONE. \square

By having the initiator r control the simulation of each round of a simulated fully-utilized synchronous protocol Π' , we obtain synchronization for an arbitrary number of rounds.

Corollary 4.4. *Multiple consecutive invocations of Algorithm 3 simulate any input fully-utilized synchronous protocol Π' round by round, resulting in an asynchronous protocol Π that uses the same number of messages.*

4.3 The Coding Scheme

We can now complete the details of our coding scheme for asynchronous networks with unknown topology. The scheme consists of two parts. In the first part, the scheme uses the BFS construction given in Section 3 in order to obtain a BFS tree \mathcal{T} . Note that the nodes ignore the content of messages during this part, therefore an adversary that can only modify messages cannot disturb this part.

In the second part, the scheme translates π into a fully-utilized synchronous protocol π' via $O(n)$ fully-utilized synchronous rounds over \mathcal{T} . This is done using Algorithm 2. The protocol π' is still non-resilient to noise and hence is not the protocol that is executed. Instead, we add a coding layer for multiparty interactive communication, namely via the Hoza-Schulman coding scheme, whose properties are given in Lemma 2.1. This results in a fully-utilized synchronous protocol Π' that is resilient to noise, which we then execute through our root-triggered synchronizer to obtain the asynchronous resilient protocol Π .

The complete construction is given in Algorithm 4. We prove its communication overhead in the following lemma, and then we prove its correctness and resilience.

Algorithm 4 A coding scheme Π for any noiseless asynchronous input protocol π .

Initialization: All nodes begin in the INIT state.

- 1: For node r designated as initiator:
 - 2: **Begin**
 - 3: Execute Algorithm 1 with r designated as root. Let \mathcal{T} be the obtained BFS tree.
 - 4: Let π' be a fully-utilized synchronous algorithm induced by π using Algorithm 2.
 - 5: Let $\Pi' = \text{HS}(\pi')$ be the Hoza-Schulman coding scheme for π' .
 - 6: Simulate Π' using the synchronizer of Algorithm 3 over \mathcal{T} with r as the initiator.
 - 7: **End**
-

Lemma 4.5. *For any asynchronous protocol π the coding Π of Algorithm 4 has a communication complexity of*

$$\text{CC}(\Pi) = O(nm \log n) + \text{CC}(\pi) \cdot O(n^2 \log n).$$

Proof. Recall that we assume channels with a fixed alphabet size, so that each symbol contains $O(\log n)$ bits (Remark 2).

The $O(nm \log n)$ term follows from Theorem 3.1. The transformation of Algorithm 2 induces a communication overhead factor of $O(n^2 \log n)$ per bit of π , as shown in Lemma 4.2.

By Lemma 2.1 there exists a resilient fully-utilized synchronous protocol Π' that simulates π' whose message/communication complexity is linear in the message complexity of π' . Finally, Corollary 4.4 gives that the asynchronous simulation of Π' via Algorithm 2 has the same message and communication complexity as Π' .

It follows that the total overhead in communication of Algorithm 4 is $O(n^2 \log n)$, as claimed. \square

Remark 3. *Note that the BFS construction (Algorithm 1) ignores the contents of messages sent. Hence, if we relax the assumption of Remark 2, the communication complexity can be reduced by sending empty messages (without any payload) during that step. In this case the message complexity of Π remains $O(mn) + \text{CC}(\pi) \cdot O(n^2)$ yet the communication complexity effectively reduces to $\text{CC}(\Pi) = \text{CC}(\pi) \cdot O(n^2 \log n)$.*

Remark 4. *In the above, each message sent in π is split into single bits and a separate symbol is dedicated to each such bit. However, instead of communicating a single bit M_i in each symbol, nodes can aggregate blocks of $O(\log n)$ bits, so that the payload of each symbol is a single block (of π 's communication) while keeping the coding's symbol size of the magnitude $O(\log n)$.*

For some protocols, namely those which send large messages, this may result in a slight logarithmic decrease in the message complexity. This optimization, however, will not change the asymptotic overhead in the worst case, when the protocol π communicates a single bit at a time.

Lemma 4.6. *For any asynchronous protocol π the coding Π of Algorithm 4 correctly simulates π even if up to $\Theta(1/n)$ of the messages are adversarially corrupted.*

Proof. Correctness and resilience to noise are proved as follows. Theorem 3.1 proves the correctness of our content-oblivious BFS construction despite noise, since the contents of the sent messages are ignored by the nodes. We emphasize that by Corollary 3.2, all of the nodes know when to stop ignoring the content of messages for the BFS construction and start executing that synchronizer over Π' .

Lemma 4.2 proves that indeed π' is a fully-utilized synchronous transformation of π . By Lemma 2.1, we have that Π' is a fully-utilized synchronous protocol that simulates π' in a manner that is resilient to corrupting up to $\Theta(1/|\tilde{E}|)$ of the messages, where \tilde{E} is the edges over which the protocol communicates. In our case these are the edges of the BFS tree \mathcal{T} , and hence this step is resilient to an $\Theta(1/n)$ -fraction of corruptions.

Finally, Corollary 4.4 gives that Π' is executed correctly in the asynchronous setting despite noise.

We now need to sum up the maximal number of symbols that can be corrupted and the total number of communicated symbols. Recall that the noise resilience is the ratio between these two sums. Since corruption can only take place on symbols of the Hoza-Schulman coding, of which there are $\text{CC}(\pi) \cdot O(n^2)$ many, we get that the scheme is resilient to at most $\Theta(1/n) \cdot \text{CC}(\pi) \cdot O(n^2)$ corrupted symbols. The total number of symbols communicated in the scheme includes also the $O(mn)$ symbols required for constructing the BFS tree, implying that our scheme is resilient to a fraction of symbol corruption equal to

$$\frac{\Theta(1/n) \cdot \text{CC}(\pi) \cdot O(n^2)}{O(nm) + \text{CC}(\pi) \cdot O(n^2)}.$$

This is asymptotically equal to an $O(1/n)$ fraction of noise when $\text{CC}(\pi) > n$, $\text{CC}(\pi) \rightarrow \infty$. \square

Lemmas 4.5 and 4.6 directly give our main theorem for this section.

Theorem 4.1 *Any asynchronous protocol π over a network G can be simulated by an asynchronous protocol Π resilient to an $\Theta(1/n)$ -fraction of adversarial message corruption, and it holds that $\text{CC}(\Pi) = O(nm \log n) + \text{CC}(\pi) \cdot O(n^2 \log n)$.*

5 A Spanner-Based Distributed Interactive Coding Scheme

In this section we slightly improve the overhead obtained by the coding scheme of Theorem 4.1. We demonstrate a family of coding schemes with an interesting tradeoff between their overhead and resilience. The key ingredient is replacing the underlying infrastructure of the BFS tree \mathcal{T} with a sparse spanning graph S , where we can trade off the sparseness of the graph (i.e., the number of edges it contains, and as a consequence, the resilience of the obtained coding scheme) with its distance distortion (i.e., the maximal distance in S for any neighboring nodes in G , and as a consequence, the added overhead for routing messages through S in the coding scheme).

Assume u sends v a message in the input protocol π . The coding scheme of Algorithm 4 routes every such message via the BFS tree \mathcal{T} . This incurs a delay in Π' , which can be of $O(n)$ rounds: in the worst case, u and v which are neighbors in G may now be two leaves of \mathcal{T} whose distance is n . In fact, even if their distance in \mathcal{T} is smaller, the coding scheme is not aware of this fact and must propagate the message to the entire network. The only guarantee we have in this case is that the message reaches its destination after at most n rounds (of the underlying fully-utilized synchronous protocol).

In this section we suggest a way to reduce the delay factor of n by routing messages over a *spanner* rather than over the tree \mathcal{T} .

Definition 1 (*t*-Spanner). *A subgraph $S = (V, E_S)$ is a t -spanner of $G = (V, E)$ if for every $(u, v) \in E$ it holds that $\text{dist}(u, v) \leq t$ in S .*

Replacing the BFS tree \mathcal{T} with a t -spanner that has $s = |E_S|$ edges ensures that a message reaches its destination after at most t steps (instead of n). Since the noise resilience is determined by the number of edges used by the underlying fully-utilized synchronous protocol, by Lemma 2.1, we obtain a resilience of $\Theta(1/s)$. The main result of this section is the following.

Theorem 5.1. *Let π_{spanner} be an asynchronous distributed algorithm for constructing a t -spanner S with s edges in a noiseless setting. Any asynchronous protocol π over a network G with $\text{CC}(\pi) \gg \text{CC}(\pi_{\text{spanner}})$ can be simulated by a noise-resilient asynchronous protocol Π resilient to an $\Theta(1/s)$ -fraction of message corruption and it holds that $\text{CC}(\Pi) = \text{CC}(\pi) \cdot O(st \log n)$.*

Specifically, due to the existence of $O(\log n)$ -spanner with $O(n)$ edges [Awe85, PS89] (see also [Pel00, Section 16]), we can let π_{spanner} be a distributed construction of a spanner with the same parameters [DMZ10] and obtain the following corollary.

Corollary 5.2. *Let π_{spanner} be an asynchronous distributed algorithm for constructing a $\log n$ -spanner with $O(n)$ edges in a noiseless setting. Any asynchronous protocol π over a network G with $\text{CC}(\pi) \gg \text{CC}(\pi_{\text{spanner}})$ can be simulated by a noise-resilient asynchronous protocol Π resilient to an $\Theta(1/n)$ -fraction of message corruption and it holds that $\text{CC}(\Pi) = \text{CC}(\pi) \cdot O(n \log^2 n)$.*

There are several challenges that arise when we replace the BFS tree \mathcal{T} with a t -spanner S . As in the case of a tree \mathcal{T} , the nodes know neither the topology of the graph nor the shortest route between any two nodes, and the only way to propagate information to its destination is by flooding it through the network.

One difficulty stems from the fact that multiple paths may exist between any two nodes in S , while only a single path exists in \mathcal{T} . This means that we are no longer guaranteed that each message arrives to its destination only *once* and that consecutive messages arrive *in the correct order*. One possible way to overcome the above issue is to add a serial counter to each message. Unfortunately, a counter of $O(\log n)$ bits is not enough to avoid confusion. Consider for instance the case where u and v are at distance $t = O(\log n)$ of each other⁷, yet each node in between them is connected to \sqrt{n} (unique) nodes. Furthermore, assume that all nodes are currently sending messages. If we flood the information through the network similar to the case of the tree \mathcal{T} , i.e., where each node holds a queue of incoming messages which it relays one by one, then the message between u and v going through that route may be delayed by $(\sqrt{n})^t = n^{\Theta(\log n)}$. Such a large delay cannot be recorded by a counter of $O(\log n)$ bits. Therefore, if u and v are also connected via another *short* path (or alternatively, a path that is not congested), then messages from the long path may be confused with messages from the short path so that the correct order of the messages could not be retrieved. To bypass the congestion issue, our coding scheme uses a priority mechanism for contention resolution, where messages with low priority are dropped by congested nodes and are later resent by their originating node. A full description and a careful proof that no messages are dropped are given in Section 5.1.

The remaining issue is how to construct the spanner despite the noisy communication. Luckily, this can be done using Theorem 4.1, i.e., by running Algorithm 4 with an input distributed protocol $\pi_{spanner}$ for constructing spanners. For example, we can take a deterministic synchronous⁸ protocol for constructing a $O(\log n)$ -spanner of size $O(n)$, e.g., the construction of Derbel, Mosbah, and Zemmari [DMZ10] which sends messages of size $O(\log n)$ and takes $O(n)$ rounds to complete. On this input, Algorithm 4 communicates $O(mn^2)$ messages in the worst case (see also Remark 4).

Once we obtain a t -spanner, we proceed as in Section 4 by converting the input protocol π to a fully-utilized synchronous protocol over the spanner using Algorithm 2, coding the fully-utilized synchronous protocol using the Hoza-Schulman scheme, and simulating each round of the coded algorithm using the root-triggered synchronizer over the spanner via Algorithm 3.

Another issue which requires a careful attention is that when we execute two coded algorithms one after the other we must make sure that the lengths of both parts are balanced. This happens, for example, when we first run the algorithm for obtaining the spanner S and then execute the algorithm for simulating π over S . The reason for this is that otherwise the strong adversary can choose to attack the shorter algorithm using the larger budget of errors it has due to the longer algorithm. For example, assume that the coded construction of the spanner communicates x symbols while the coded simulation of π sends y symbols. If x is at most $O(y/s)$, then the adversary has enough budget to fully corrupt the first part. Requiring, say, $x = y$, makes the scheme resilient to a $\Theta(1/2s) = \Theta(1/s)$ fraction of corrupted symbols (see also footnote 10).

⁷E.g., when using an $O(\log n)$ -spanner of size $s = O(n)$ towards Corollary 5.2.

⁸In Theorem 5.1 we assume $\pi_{spanner}$ is an asynchronous protocol. Note that any synchronous protocol with time complexity τ sending messages of size σ can be thought of as an asynchronous protocol with message complexity $O(\tau m)$ and communication complexity $O(\tau m \sigma)$.

5.1 Obtaining a spanner-based fully-utilized synchronous protocol

A main difference between the coding scheme given in Section 4 and in the spanner-based coding scheme of this section is how they transform the input protocol π into a fully-utilized synchronous one. The underlying idea is similar: each node breaks the messages of π bits, encapsulates them in a symbol that contains also their source and destination, and floods them one by one through the underlying graph—a tree in the former case and a spanner in the latter.

As mentioned above, the main difficulty in the case of a spanner is the existence of several paths between every two nodes, which may cause a specific symbol to arrive multiple times at the receiver. Moreover, simply flooding the information through the spanner may cause very large delays on certain paths (super-polynomial delays⁹), due to the congestion caused by other symbols that are being relayed through the network. Hence, a receiver that gets multiple copies of symbols out of order with such long delays may not be able to reconstruct π 's message from them.

Our solution breaks the simulation into windows of $2t$ rounds each, where each node attempts to send only a single symbol at each such window. In order to avoid congestion we use a priority system: when two or more symbols are received by some node, it drops all the symbols except for the one with the highest priority—the one whose sender's identity is maximal. Then, the node relays this symbol to all of its neighbors. This procedure guarantees that at least one symbol arrives at its destination, at every $2t$ -round window. Furthermore, senders whose symbols were dropped can learn this event and resend their message during the next window. If a sender receives a symbol with a higher priority, it assumes that its own symbol was dropped during that window. We prove that this mechanism may have false negatives (i.e., a sender that resends a symbol while that symbol did arrive at the destination), but it does not have false positives (if a sender does not get an indication that it needs to resend, then its symbol is guaranteed to have arrived).

The detailed transformation into a fully-utilized synchronous protocol is given in Algorithm 5.

Towards proving the correctness of Algorithm 5, we prove some properties of the algorithm.

Lemma 5.3. *During every non-overlapping window of $2t$ rounds of π at least one message $(src, dest, val, i)$ is delivered to its destination. Furthermore, The sender of this message ends this window of $2t$ rounds with $RepeatSendMsg_u = false$.*

Proof. Since RELAY is performed according to priority, the message with the highest priority (highest src) always survives and gets relayed all the way to the destination (which takes at most t rounds). The sender of this message never receives a message with a higher priority, thus it never sets $RepeatSendMsg_u$ to *true* on Line 31. \square

Lemma 5.4. *If u sends a message and after $2t$ rounds it holds that $RepeatSendMsg_u = false$, then the message has reached its destination.*

Proof. Assume u sends a message to some node v . Additionally, assume that a node w with a higher priority ($w > u$) also attempts to send a message during the same $2t$ -round window. Without loss of generality, we assume w is the only node with priority higher than u ; also recall that nodes with lower priority have no effect on the delivery of u 's message.

For $k \leq 2t$, define B_k to be all the messages that travel during the first k rounds (of that specific $2t$ -round window) on all the edges of distance at most k from u . We consider several cases:

⁹Hence, appending $O(\log n)$ bits of a counter to each message does not suffice in order to guarantee the correct ordering of messages at the receiver's side.

Algorithm 5 Simulation of an asynchronous protocol π by a fully-utilized synchronous protocol π' using a t -spanner S .

Input: An asynchronous protocol π over G ; a t -spanner S of G (the value t is known to all nodes)

Init: for all nodes u : RepeatSendMsg $_u \leftarrow false$; NextMsg $_u$, RelayMemory $_u$, queue $_u \leftarrow \emptyset$.

```

1: In every round RND, for every node  $u$ :
2: Begin
3:   if RND  $\equiv 1 \pmod{2t}$  then
4:     Invoke SETNEXTMESSAGE
5:   end if
6:   Invoke RELAY
7: End

8: procedure SETNEXTMESSAGE (for node  $u$ )
9:   if RepeatSendMsg $_u = false$  or NextMsg $_u = \emptyset$  then
10:      $\triangleright$  Split the next message of  $\pi$  into bits, each to be delivered separately
11:     Let  $M_1 \cdots M_\ell$  be the bit representation of a message  $M$  that  $u$  has to send to  $v$  in  $\pi$ .
12:     Push  $(u, v, M_1, 1 \pmod{n}), \dots, (u, v, M_\ell, \ell \pmod{n})$  into queue $_u$   $\triangleright$  ignore if no message to send
13:     NextMsg $_u \leftarrow \text{pop}(\text{queue}_u)$ 
14:     RelayMemory $_u \leftarrow \text{NextMsg}_u$ 
15:     RepeatSendMsg $_u \leftarrow false$ 
16:   else if RepeatSendMsg $_u = true$  then  $\triangleright$  Previous bit may have not reached the destination; retry with the same NextMsg.
17:     RelayMemory $_u \leftarrow \text{NextMsg}_u$ 
18:     RepeatSendMsg $_u \leftarrow false$ 
19:   end if
20: end procedure

21: procedure RELAY (for node  $u$ )
22:    $(src', dest', val', i', w') \leftarrow \text{RelayMemory}_u$ 
23:    $\triangleright$  Relay highest priority message to all neighbors
24:   if RelayMemory $_u \neq \emptyset$  then
25:     send  $(src', dest', val', i')$  to every  $v \in \mathcal{N}_u(S) \setminus \{w'\}$  and send  $\perp$  to  $w'$ 
26:   else
27:     send  $\perp$  to every  $v \in \mathcal{N}_u(S)$ 
28:   end if
29:   RelayMemory $_u \leftarrow \emptyset$ 

30:    $\triangleright$  Receive messages from all neighbors, keep only the message with highest priority
31:   Upon receiving a message  $(src, dest, val, i)$  from  $w$ :
32:   Begin
33:     if  $src' > u$  then  $\triangleright$  Whenever a message with higher priority is received, assume own message wasn't delivered during this window of  $2t$  rounds
34:       RepeatSendMsg $_u \leftarrow true$ 
35:     end if
36:     if  $dest \neq u$  then
37:       if  $src > src'$  or RelayMemory $_u = \emptyset$  then
38:         RelayMemory $_u \leftarrow (src, dest, val, i, w)$ 
39:       end if
40:     else if  $dest = u$  then
41:       collect the bits  $val$  for decoding  $M$ 
42:       (using the index  $i$  to ignore already received bits)
43:     end if
44:   End
45: end procedure

```

1. If u 's message has the highest priority in B_{2t} then it is clear that the statement holds. This happens whenever $dist(u, w) > 4t$.
2. If $dist(u, w) \leq 2t$ then u receives w 's message during that window, and sets $RepeatSendMsg_u$ to $true$, so in this case the required conditions of the lemma do not hold to begin with.
3. The last case is when B_{2t} contains w 's message, yet $dist(u, w) > 2t$. This implies that $dist(v, w) > t$ (otherwise, via the triangle inequality, we have that $dist(u, w) \leq 2t$). This in turn implies that u 's message is the one with the highest priority in B_t , which means that u 's message is delivered to v by round t .

□

Note that the opposite direction of the statement of Lemma 5.4 does not always hold. That is, if u sends a message and after $2t$ rounds it holds that $RepeatSendMsg_u = true$, then it is possible that the message has nevertheless reached its destination.

Lemma 5.3 and Lemma 5.4 suggest that progress is made every $2t$ rounds: at least one message is being delivered and all nodes whose messages are not delivered receive an indication of this event (and retry during the next window). This leads to the correctness of the algorithm, as stated in the following lemma.

Lemma 5.5. *Algorithm 5 creates a fully-utilized synchronous protocol π' that simulates π , in the sense that all messages of π are sent and received, with a communication overhead of $O(st \log n)$, where s is the number of edges in the input t -spanner S .*

Proof. By construction, every node sends a message to all of its neighbors in S in each round and hence Algorithm 5 is a fully-utilized synchronous protocol. In addition, as implied by Lemma 5.3 and Lemma 5.4, eventually every messages of π reaches its destination and hence the obtained fully-utilized synchronous protocol simulates all messages in π . Note that a message may be resent and received multiple times at the destination. However a node u always receives from any specific node v either the next bit of the message v sends, or a re-transmission of the last bit of the same message. Therefore, adding a single parity bit of the index of the transmitted bit is enough in order to avoid confusion (In Algorithm 3 we actually add to each bit the information $(i \bmod n)$ where i is the index of the bit; yet communicating the value $(i \bmod 3)$ would have sufficed).

For the communication overhead, note that $O(\log n)$ bits of the identities of source and destination and bit index are appended to each bit sent by π , that is $\log |\Sigma| = O(\log n)$ suffices. In addition, a delivery of a single bit of π may require $O(t)$ rounds when sent over the spanner S . In each such round $O(s)$ symbols are sent in a fully-utilized synchronous protocol. This implies $O(st)$ messages are communicated per each bit of π which gives a total communication overhead of $O(st \log n)$.

Note that this is a worst-case analysis that assumes a single bit (of π) travels within the network at each time so that another bit is sent only after a previous bit reached its destination. If several bits are sent consecutively or if several nodes send bits simultaneously, the resulting number of messages can only decrease. □

Algorithm 6 Spanner-based coding scheme Π for any asynchronous noiseless protocol π .

Initialization: All nodes know (a bound on) $CC(\pi)$. $\pi_{spanner}$ is an asynchronous (noiseless) protocol for constructing a t -spanner.

- 1: For node r designated as initiator:
 - 2: **Begin**
 - 3: Execute Algorithm 4 on input $\pi_{spanner}$ with r designated as root, extending the protocol so that the total message complexity of this step is C_2 defined below. Let S be the obtained t -spanner.
 - 4: Let π' be a fully-utilized synchronous algorithm induced by π using Algorithm 5 on S .
 - 5: Let $\Pi' = HS(\pi')$ be the Hoza-Schulman coding scheme for π' . Let C_2 be the message complexity of Π' .
 - 6: Simulate Π' using the synchronizer of Algorithm 3 over S with r as the initiator.
 - 7: **End**
-

5.2 The spanner-based coding scheme

Our modified spanner-based coding scheme Π is given in Algorithm 6.

As mentioned above, when a coding scheme contains two parts (i.e., constructing the spanner and executing the coding scheme) that are being coded independently, it is necessary to make sure that the two parts are of equal length. Recall that we consider the asymptotical behavior of the coding scheme when $CC(\pi)$ tends to infinity. Hence, we assume $CC(\pi) \gg CC(\pi_{spanner})$, which implies that we need to extend the first part of our scheme where we construct the spanner graph in a resilient way. There are two possible ways to make this extension: either by artificially increasing the communication of $\pi_{spanner}$ to $\approx O(CC(\pi) \cdot st/n^2)$ bits, say, by sending zeroes after the spanner's construction has completed; or by running the Hoza-Schulman encoding for $O(C_2/n)$ rounds.

We claim that the coding scheme Π of Algorithm 6 satisfies the requirements of Theorem 5.1.

Proof. (Theorem 5.1) We first analyze the communication complexity of the scheme. We recall the assumption stated in Remark 2 that symbol of the coding scheme contains $O(\log n)$ bits; thus we can equivalently bound the message complexity. The execution of Algorithm 4 in Line 3 has a message complexity of $C_1 = O(mn) + CC(\pi_{spanner}) \cdot O(n^2)$. However, we artificially extend this step, so that it would have a message complexity of C_2 . By Lemma 5.5, the message complexity of π' given by Algorithm 5 (Line 4) is $C_2 = CC(\pi) \cdot O(st)$. Hence, the overall message complexity of Algorithm 6 is $2C_2$, and the overall communication complexity is $O(C_2 \log n) = CC(\pi) \cdot O(ts \log n)$.

Regarding the resilience of the scheme, the first part (Line 3) is resilient due to Theorem 4.1; the second part (Lines 4–6) is independently encoded via the Hoza-Schulman coding scheme, whose guarantees are given in Lemma 2.1. The first part is resilient to a $\mu_1 = \Theta(1/n)$ fraction of corrupted messages and the second part is resilient to a $\mu_2 = \Theta(1/s)$ fraction of corrupted messages. Since we balance the message complexity of the two parts so that $C_1 = C_2$ and since $s = \Omega(n)$, the new scheme is resilient to a fraction of $\mu_2/2 = \Theta(1/s)$ of corrupted messages overall.¹⁰ \square

Acknowledgement. We are grateful to Merav Parter for bringing [DMZ10] to our attention.

¹⁰By balancing the parts C_1 and C_2 in a weighted way one can obtain a slightly improved resilience of $\mu_1\mu_2/(\mu_1 + \mu_2) = \Theta(1/(n + s))$ which is, however, asymptotically equivalent to $\Theta(1/s)$.

References

- [ABE⁺16] N. Alon, M. Braverman, K. Efremenko, R. Gelles, and B. Haeupler. [Reliable communication over highly connected noisy networks](#). *Proceedings of the 2016 ACM Symposium on Principles of Distributed Computing, PODC '16*, pp. 165–173, 2016.
- [AW04] H. Attiya and J. Welch. *Distributed Computing: Fundamentals, Simulations and Advanced Topics*. John Wiley & Sons, 2004.
- [Awe85] B. Awerbuch. [Complexity of network synchronization](#). *J. ACM*, 32(4):804–823, 1985.
- [APPS92] B. Awerbuch, B. Patt-Shamir, D. Peleg, and M. E. Saks. [Adapting to asynchronous dynamic networks \(extended abstract\)](#). *STOC '92: Proceedings of the 24th Annual ACM Symposium on Theory of Computing*, pp. 557–570, 1992.
- [BKN14] Z. Brakerski, Y. T. Kalai, and M. Naor. [Fast interactive coding against adversarial noise](#). *J. ACM*, 61(6):35:1–35:30, 2014.
- [BE14] M. Braverman and K. Efremenko. [List and unique coding for interactive communication in the presence of adversarial noise](#). *Proceedings of the 55th annual IEEE Symposium on Foundations of Computer Science, FOCS '14*, pp. 236–245, 2014.
- [BEGH16] M. Braverman, K. Efremenko, R. Gelles, and B. Haeupler. [Constant-rate coding for multiparty interactive communication is impossible](#). *Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing, STOC '16*, pp. 999–1010, ACM, 2016.
- [BGMO16] M. Braverman, R. Gelles, J. Mao, and R. Ostrovsky. [Coding for interactive communication correcting insertions and deletions](#). *43rd International Colloquium on Automata, Languages, and Programming (ICALP '16), Leibniz International Proceedings in Informatics (LIPIcs)*, vol. 55, pp. 61:1–61:14, 2016.
- [BR14] M. Braverman and A. Rao. [Toward coding for maximum errors in interactive communication](#). *Information Theory, IEEE Transactions on*, 60(11):7248–7255, 2014.
- [Das98] P. Dasgupta. [Agreement under faulty interfaces](#). *Information Processing Letters*, 65(3):125 – 129, 1998.
- [DMZ10] B. Derbel, M. Mosbah, and A. Zemmari. [Sublinear fully distributed partition with applications](#). *Theory of Computing Systems*, 47(2):368–404, 2010.
- [EGH16] K. Efremenko, R. Gelles, and B. Haeupler. [Maximal noise in interactive communication over erasure channels and channels with feedback](#). *IEEE Transactions on Information Theory*, 62(8):4575–4588, 2016.
- [FHK14] O. Feinerman, B. Haeupler, and A. Korman. [Breathe before speaking: efficient information dissemination despite noisy, limited and anonymous communication](#). *ACM Symposium on Principles of Distributed Computing, PODC '14*, pp. 114–123, 2014.
- [FLP85] M. J. Fischer, N. A. Lynch, and M. Paterson. [Impossibility of distributed consensus with one faulty process](#). *J. ACM*, 32(2):374–382, 1985.

- [FGOS15] M. Franklin, R. Gelles, R. Ostrovsky, and L. J. Schulman. [Optimal coding for streaming authentication and interactive communication](#). *Information Theory, IEEE Transactions on*, 61(1):133–145, 2015.
- [Gal82] R. G. Gallager. Distributed minimum hop algorithms. Tech. Rep. LIDS-P-1175, M.I.T. Laboratory for Information and Decision Systems, 1982.
- [Gel15] R. Gelles. [Coding for interactive communication: A survey](#), 2015.
- [GHK⁺16] R. Gelles, B. Haeupler, G. Kol, N. Ron-Zewi, and A. Wigderson. [Towards optimal deterministic coding for interactive communication](#). *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1922–1936, 2016.
- [GK17] R. Gelles and Y. T. Kalai. Constant-rate interactive coding is impossible, even in constant-degree networks. *Proceedings of the 8th Conference on Innovations in Theoretical Computer Science, ITCS '17*, 2017.
- [GMS14] R. Gelles, A. Moitra, and A. Sahai. [Efficient coding for interactive communication](#). *Information Theory, IEEE Transactions on*, 60(3):1899–1913, 2014.
- [GHS14] M. Ghaffari, B. Haeupler, and M. Sudan. [Optimal error rates for interactive coding I: Adaptivity and other settings](#). *Proceedings of the 46th Annual ACM Symposium on Theory of Computing, STOC '14*, pp. 794–803, 2014.
- [GLR95] L. Gong, P. Lincoln, and J. Rushby. [Byzantine agreement with authentication: Observations and applications in tolerating hybrid and link faults](#). *Dependable Computing and Fault Tolerant Systems*, vol. 10, pp. 139–158, IEEE Computer Society, 1995.
- [Hae14] B. Haeupler. [Interactive channel capacity revisited](#). *Proceedings of the 55th annual IEEE Symposium on Foundations of Computer Science, FOCS '14*, pp. 226–235, 2014.
- [HS94] M. Harrington and A. K. Somani. [Synchronizing hypercube networks in the presence of faults](#). *IEEE Trans. Computers*, 43(10):1175–1183, 1994.
- [HS16] W. M. Hoza and L. J. Schulman. [The adversarial noise threshold for distributed protocols](#). *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 240–258, 2016.
- [JKL15] A. Jain, Y. T. Kalai, and A. Lewko. [Interactive coding for multiparty protocols](#). *Proceedings of the 6th Conference on Innovations in Theoretical Computer Science, ITCS '15*, pp. 1–10, 2015.
- [KR13] G. Kol and R. Raz. [Interactive channel capacity](#). *STOC '13: Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*, pp. 715–724, ACM, 2013.
- [LV15] A. Lewko and E. Vitercik. [Balancing communication for multi-party interactive coding](#), 2015. ArXiv preprint arXiv:1503.06381.
- [Lyn96] N. A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1996.

- [Pel92] A. Pelc. [Reliable communication in networks with byzantine link failures](#). *Networks*, 22(5):441–459, 1992.
- [Pel00] D. Peleg. *Distributed Computing: A Locality-Sensitive Approach*. Society for Industrial and Applied Mathematics, 2000.
- [PS89] D. Peleg and A. A. Schäffer. [Graph spanners](#). *Journal of Graph Theory*, 13(1):99–116, 1989.
- [RS94] S. Rajagopalan and L. Schulman. [A coding theorem for distributed computation](#). *STOC '94: Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, pp. 790–799, 1994.
- [SAA95] H. M. Sayeed, M. Abu-Amara, and H. Abu-Amara. [Optimal asynchronous agreement and leader election algorithm for complete networks with byzantine faulty links](#). *Distributed Computing*, 9(3):147–156, 1995.
- [Sch92] L. J. Schulman. [Communication on noisy channels: a coding theorem for computation](#). *Foundations of Computer Science, Annual IEEE Symposium on*, pp. 724–733, 1992.
- [Sch96] L. J. Schulman. [Coding for interactive communication](#). *IEEE Transactions on Information Theory*, 42(6):1745–1756, 1996.
- [Sin96] G. Singh. [Leader election in the presence of link failures](#). *IEEE Transactions on Parallel and Distributed Systems*, 7(3):231–236, 1996.
- [SCY98] H.-S. Siu, Y.-H. Chin, and W.-P. Yang. [Byzantine agreement in the presence of mixed faults on processors and links](#). *IEEE Transactions on Parallel and Distributed Systems*, 9(4):335–345, 1998.
- [Tel00] G. Tel. *Introduction to distributed algorithms*. Cambridge university press, 2000. Chapter 12.4. Asynchronous BFS Algorithms, pages 414–420.