

# Universal Bufferless Routing

Costas Busch<sup>1</sup>, Malik Magdon-Ismail<sup>1</sup>, and Marios Mavronicolas<sup>2</sup>

<sup>1</sup> Department of Computer Science  
Rensselaer Polytechnic Institute  
Troy, NY 12180, USA

<sup>2</sup> Department of Computer Science  
University of Cyprus  
Nicosia CY-1678, Cyprus

**Abstract.** Given an arbitrary network, and a routing problem with congestion  $C$  and dilation  $D$ , a long standing open problem is to show the existence of *bufferless* routing algorithms with optimal performance guarantees (routing time close to the lower bound  $\Omega(C+D)$ ). Our main result is a new *deterministic* technique that constructs a universal bufferless algorithm by emulating a universal buffered algorithm. The heart of the emulation is to replace packet buffering with packet circulation on regions of the network. The cost of the emulation on the routing time is proportional to the square of the node buffer size used by the buffered algorithm. We apply this emulation to a simple randomized buffered algorithm to obtain a *distributed*, universal bufferless algorithm with routing time  $O((C+D) \cdot \log^3(n+N))$ , which is within poly-logarithmic factors from the optimal, where  $n$  is the size of the network and  $N$  is the number of packets. The *bufferless competitive ratio* is the ratio of the best achievable bufferless routing time, to the best achievable buffered routing time. We give the first non-trivial bound of  $O(\log^3(n+N))$  for the bufferless competitive ratio for arbitrary routing problems.

## 1 Introduction

*Packet routing* has received a large amount of attention over the past decade on account of its importance to applications ranging from parallel and distributed algorithms to communication networks. The task is to deliver packets from their sources to their destinations along specified paths in a given network. A packet routing algorithm is *universal* if it can be applied to any routing problem on any network topology. For a given set of paths, the *routing time* (denoted  $rt$ ) is the time at which the last packet reaches its destination. Universal algorithms with optimal or near-optimal routing time are known if packets may be buffered along their paths, [20, 23, 24]

A long standing and important open problem is to give universal *bufferless* routing algorithms with near optimal performance guarantees. In this paper, we will present a distributed bufferless routing algorithm that is optimal up to poly-logarithmic factors. We introduce a new technique for developing bufferless algorithms based upon emulating buffered algorithms. Applying this technique to a simple randomized buffered protocol gives the advertised result.

*Preliminaries.* A routing problem  $Q = (G, \Pi, P)$  on the graph  $G$  with  $n$  nodes consists of a set of  $N$  packets  $\Pi = \{\pi_1, \pi_2, \dots, \pi_N\}$  that are to be routed on their respective paths  $P = \{p_1, p_2, \dots, p_N\}$ , where  $p_i$  is a path in  $G$ . We will represent paths either as a sequence of edges, or as a sequence of nodes, and the length of a path  $|p|$  is the number of edges in the path. The *edge-congestion*  $C$  is the maximum number of packets that use an edge in  $G$ , the *node-congestion*  $\overline{C}$  is the maximum number of packets that use a node in  $G$ , and the *dilation*  $D$  is the maximum path length in  $P$ .

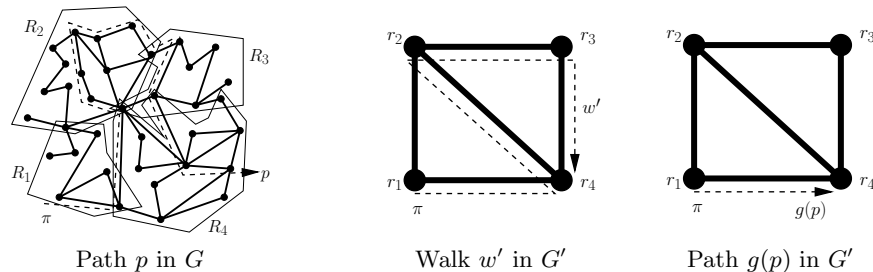
We assume a synchronous routing model, in which time is divided into a sequence of discrete time steps. An edge may be traversed by at most one packet in either direction during a time step. A well known lower bound on the routing time in this model is given by  $\Omega(C + D)$ , and so the optimal routing time  $rt^* = \Omega(C + D)$ . In a buffered algorithm, packets may either traverse edges or be buffered at a node. In a bufferless algorithm, a packet must traverse an available edge at every time step.

*An Impossibility Result.* If all packets must follow the paths specified in  $P$ , without collisions or buffering, then the only degree of freedom for a bufferless routing algorithm is the injection times of the packets. Such a routing paradigm is known as *direct routing*, [3, 13]. In this case, it is shown in [13] that there exist routing problems for which bufferless routing times better than a  $\sqrt{N}$  factor from optimal are not possible. Thus, if the paths remain unchanged, then near-optimal universal bufferless algorithms do not exist (where near optimal means within poly-logarithmic factors from the lower bound  $C + D$ ). Thus, to obtain near-optimal bufferless schedules, we must allow packets to deviate from their paths. However, we still measure performance with respect to  $C$  and  $D$  of the original paths. The justification of this is that if the paths  $P$  themselves are optimal, i.e., they minimize  $C + D$ , then we obtain bufferless routing times that are near-optimal for the given sources and destinations. We do not discuss how to obtain the optimal paths, but rather how to send the packets to their destinations given the paths.

*Contributions.* Our main result is a deterministic technique for bufferless emulation of buffered algorithms. Given a near-optimal universal buffered algorithm that routes problems with simple paths, and uses buffers of size  $\gamma$ , we give a universal bufferless algorithm, which emulates the buffered algorithm. The cost of the emulation on the routing time is  $O(\gamma^2 \cdot \log n)$ .

We apply this emulation result to a simple randomized buffered algorithm that uses  $O(\log(n + N))$  buffers to obtain a bufferless routing algorithm with routing time  $O((C + D) \cdot \log^3(n + N))$  with high probability, which approximates within poly-logarithmic factors the optimal routing time for the given paths. If all the nodes know the network topology, and the values of  $C$  and  $N$ , then the bufferless algorithm is distributed, i.e., routing decisions are made locally at each node.

*Overview of the Approach.* The main idea behind the bufferless emulation of a buffered algorithm is to use regions in the network in order to emulate buffer



**Fig. 1.** An example of a region graph.

space. We decompose the graph into connected regions each containing approximately  $\gamma$  edges. The regions form a region graph, on which the nodes are regions. Now, a buffered algorithm executes as if the regions were the nodes. In the buffered algorithm, a packet is either buffered in a node (region) or “hops” from node to node. The path of each packet in the original graph is translated to a path on the region graph. The buffer needed at each node is at most  $\gamma$ . Figure 1 illustrates the general idea of decomposing the graph into regions and then mapping a packet’s path to the graph in which every node corresponds to a region.

The buffered algorithm on the region graph is emulated by a bufferless algorithm on the original graph. If in the buffered algorithm a packet needs to be buffered in a node (region), then, in the emulation the packet “circulates” in the respective region by moving from one edge of the region to the next. A packet circulates until the buffered algorithm prescribes that the packet makes its next hop, in which case the packet moves to the respective adjacent region. Since the buffered algorithm requires  $\gamma$  buffer space per node (region), there is enough room to circulate all the packets in the  $\gamma$  edges of the region in a bufferless fashion.

*Related Work.* There are no previously known results for universal bufferless routing with near-optimal routing time guarantees. However, near-optimal bufferless routing has been obtained for specific bufferless routing models and architectures, which we summarize. In *hot-potato routing*, packets are deflected along available links in a collision [5]. Our model of bufferless routing is essentially the hot-potato routing model, with packets being deflected along particular available edges specified by the emulation (i.e., not on an arbitrary available edge as is typically done in hot-potato algorithms). Hot-potato routing algorithms have been extensively studied for a variety of architectures such as the mesh and torus [4, 6, 10, 12, 16, 19], hypercubes [9, 16, 18], trees [14, 25], vertex-symmetric networks [21], and leveled networks [8, 11]. Typically, by allowing packets to deviate from their paths slightly, one obtains routing times that are within poly-logarithmic factors of optimal. In *direct routing*, packets follow their

paths without buffering and without any collisions, [3, 13]. Wormhole routing is similar to direct routing, but here, packets occupy more than one edge [15, 17]. A dual to direct routing is *time constrained routing*, where the task is to schedule as many packets as possible within a given time frame [1]. In *matching routing*, packets are swapped at adjacent nodes, and permutation problems on trees have been studied in [2, 26].

There are two variants of buffered algorithms. Those that use buffers on every edge (*edge-buffers*) and those that use buffers in every node (*node-buffers*). For non-bounded degree networks, these variants are distinct. The existence of optimal, universal *buffered* routing algorithms using constant size edge-buffers was first established by Leighton, Maggs and Rao [20]. Thereafter, the main focus has been on constructive algorithms with optimal,  $O(C + D)$ , routing time, [7, 22–24]. These algorithms use large (proportional to the congestion  $C$ ) buffers. Leighton *et al.* [20] improve this result, requiring only edge-buffers of size  $O(\log ND)$  to obtain routing time  $O(C + D \log ND)$ . Cypher *et al.* [15] give an algorithm with edge-buffers of size  $O(\log CD)$  and slightly better routing time. Our bufferless algorithm is based on emulating a universal buffered algorithm. However, the existing results, though powerful, do not suit our purpose because we need algorithms where the node-buffers are small (logarithmic), and so we offer a simple randomized algorithm that satisfies the conditions for bufferless emulation.

*Paper Outline.* We first discuss how to decompose a graph into connected regions of approximately a given size (Section 2). We then show how these regions are used for bufferless emulation of a buffered algorithm (Section 3). Finally we apply the emulation to a randomized buffered algorithm (Section 4) to obtain near-optimal universal bufferless routing (Section 5). We conclude with a discussion (Section 6). Due to lack of space, several proofs have moved to the full version of the paper.

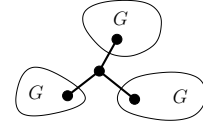
## 2 Regions

We first discuss how to decompose a connected graph  $G$  into connected components of approximately a specified size. Such a decomposition will be required by the bufferless emulation algorithm. Specifically, let  $G = (V, E)$  be an undirected connected graph. Let  $F$  be a subset of the edges in  $E$ . The subgraph induced by  $F$  is the graph  $H = (U, F)$ , where  $U$  is the union of all vertices in  $V$  that are incident with edges in  $F$ . We say that the edge set  $F$  is *connected* if the induced subgraph  $H$  is connected. A *connected decomposition* of  $G$  is a partition of the edges in  $E$  into disjoint sets  $E_1, E_2, \dots, E_k$  such that  $\cup_{i=1}^k E_i = E$  and every  $E_i$  is connected. We refer to the  $E_i$ 's as the *connected edge sets* or *regions* in the decomposition, and denote the number of edges in  $E_i$  as the size of  $E_i$ ,  $|E_i|$ . Notice that the subgraphs,  $H_1 = (V_1, E_1), \dots, H_k = (V_k, E_k)$  induced by the edge sets may have overlapping vertex sets. We say that  $E_i$  is *connected* to  $E_j$  if and only if  $V_i \cap V_j \neq \emptyset$ . Notice that if  $E_i$  is connected to  $E_j$ , then  $E_i \cup E_j$  is a connected edge set.

An  $[\alpha, \beta]$ -partition of  $G$  (if it exists) is a connected decomposition of  $G$ ,  $\{E_1, \dots, E_k\}$ , such that  $\alpha \leq |E_i| \leq \beta$  for  $i = 1, \dots, k$ . Notice that if  $\alpha \approx \beta$ , then an  $[\alpha, \beta]$ -partition decomposes  $G$  into connected edge sets of size approximately equal to  $\alpha$ . We now show that such approximate decompositions are possible for any connected graph.

**Theorem 1 (Existence of a  $[k, 3k - 3]$ -partition).** *Let  $G = (V, E)$  be a connected graph. For any  $k$ , where  $1 < k \leq |E|$ , there exists a  $[k, 3k - 3]$ -partition of  $G$ .*

The following example proves that the result of Theorem 1 is tight. For a given  $k$ , let  $G$  be any connected graph with  $k - 2$  edges, and connect 3 such graphs in a wheel configuration as shown on the right. It is easy to see that the only decomposition in which every edge set has  $\geq k$  edges is the entire graph itself, which has  $3k - 3$  edges.



The proof in Theorem 1 is constructive, hence it can be directly converted to an algorithm. One can show that the time complexity of the algorithm to compute a decomposition is in  $O(|E|^2)$ .

*Region Graph.* Consider a connected graph  $G = (V, E)$ , with  $n$  nodes. Take an  $[\alpha, \beta]$ -partition of  $G$ , which gives regions (connected edge sets)  $R_1, R_2, \dots, R_k$ . Let the subgraphs induced by these regions have vertex sets  $U_1, U_2, \dots, U_k$ . The *region graph*  $G' = (V', E')$ , has a vertex set  $V' = \{r_1, r_2, \dots, r_k\}$  where each vertex  $r_i$  corresponds to the region  $R_i$  of  $G$ . Two vertices  $r_i, r_j$  are adjacent in  $G$ , i.e.,  $(r_i, r_j) \in E'$  if and only if  $U_i \cap U_j \neq \emptyset$ , i.e., the corresponding regions have intersecting vertex sets. An example of a region graph is given in Figure 1.

*Routing Problems on Region Graph.* Let  $Q = (G, \Pi, P)$  denote a routing problem with edge-congestion  $C$ , node-congestion  $\overline{C}$  and dilation  $D$ . Let  $\{R_1, \dots, R_k\}$  be an  $[\alpha, \beta]$ -partition of  $G$ . Every edge in  $G$  belongs to exactly one region. Let  $G' = (V', E')$  be the corresponding region graph. The mapping  $f : E \rightarrow V'$  is defined for every  $e \in E$  by  $f(e) = r_i$  if and only if  $e \in R_i$ . Consider a path  $p \in P$ , with  $p = (e_1, e_2, \dots, e_l)$ . We define a function  $g$  which maps a path in  $G$  to a path in  $G'$  as follows. For any path  $p = (e_1, e_2, \dots, e_l)$  in  $G$ , consider the walk in  $G'$  given by  $w' = (f(e_1), f(e_2), \dots, f(e_l))$ .  $g(p)$  is the path obtained after removing all the cycles in  $w'$ ,  $g(p) = (f(e_{i_1}), f(e_{i_2}), \dots, f(e_{i_l}))$ .

We now transform the routing problem  $Q$  on the original graph into a routing problem  $Q' = (G', \Pi, P')$  on the region graph, in which the paths in  $G'$  are given by the transformed paths,  $P' = \{p'_1, p'_2, \dots, p'_N\}$  where  $p'_i = g(p_i)$ ,  $\forall p_i \in P$ . Let  $C', \overline{C}'$  and  $D'$  denote the edge-congestion, the node-congestion and the dilation of the paths in  $P'$ . For any routing problem, the edge-congestion is bounded by the node-congestion. A path uses node  $r_i$  only if it contains edges in  $R_i$ . By construction,  $|R_i| \leq \beta$ , so the number of edges in  $P$  that use  $R_i$  is at most  $\beta C$ , thus  $\overline{C}' \leq \beta C$ . Since  $|g(p)| \leq |p|$  for any path  $p$  in  $G$ , we have the following lemma.

**Lemma 1.**  $C' \leq \overline{C}' \leq \beta C; D' \leq D$ .

*Euler Tours on Regions.* We define Euler tours with respect to the directed representation  $G^D = (V, E^D)$  of the undirected graph  $G$ : each (undirected) edge  $(u, v) \in E$  is replaced by two directed edges  $(u, v), (v, u) \in E^D$ . Let  $R_i^D$  denote the region of  $G^D$  that corresponds to the region  $R_i$  in  $G$ . Since the in-degree equals the out-degree of every node in  $R_i^D$ ,  $R_i^D$  has an Euler tour. Let  $\psi_i = (v_1, v_2, \dots, v_1)$  denote an Euler tour in  $R_i^D$ . Note that  $\psi_i$  is walk in  $R_i$ . We will refer to  $\psi_i$  as the “Euler tour” of  $R_i$  (an abuse of notation, since  $\psi_i$  is not an Euler tour of  $R_i$ ). Note that for an  $[\alpha, \beta]$ -partition of  $G$ , every Euler tour  $\psi_i$  satisfies  $2\alpha \leq |\psi_i| \leq 2\beta$ .

### 3 Emulation

Let  $G = (V, E)$  be a connected graph with  $n$  nodes and let  $\{R_1, \dots, R_k\}$  be an  $[\alpha, \beta]$ -partition of  $G$  with corresponding region graph  $G' = (V', E')$ . For routing problem  $Q = (G, \Pi, P)$  in  $G$ , we obtain the corresponding routing problem  $Q' = (G', \Pi, P')$  in  $G'$ . Let  $(s_i, d_i)$  denote the source and destination of each packet  $\pi_i \in \Pi$ , and let  $S = \{(s_1, d_1), (s_2, d_2), \dots, (s_N, d_N)\}$ . Let  $Q_s = (G, \Pi, S)$  denote the routing problem in  $G$  in which the packets need to be delivered from their sources to their destination, without necessarily following the paths in  $P$ .

The general idea behind our approach is to design a bufferless routing Algorithm  $B$  to solve the routing problem  $Q_s$ . The bufferless algorithm will depend on a buffered Algorithm  $A$  to solve the routing problem  $Q'$  in  $G'$ . The bufferless algorithm will then *emulate* the running of Algorithm  $A$  in  $G'$  to solve  $Q_s$  in  $G$ .

#### 3.1 Buffered Routing in $G'$ – Algorithm $A$

Our bufferless algorithm in  $G$  will emulate a buffered algorithm  $A$  in  $G'$ . Algorithm  $A$  solves routing problem  $Q'$  in  $G'$  and uses node-buffers of size at most  $\gamma$  to do so. We require algorithm  $A$  to receive at most  $\gamma$  packets at every time step. It is then possible to divide the execution of Algorithm  $A$ , into a sequence of *phases*, in which each phase has the following two properties:

- (i) Each phase is a fixed time period consisting of at least one time step;
- (ii) During each phase, each packet traverses at most one edge in  $G'$ , and each node receives at most  $\gamma$  packets from adjacent nodes or through injection.

A trivial division of the execution of Algorithm  $A$  into phases that satisfies these two properties is to take each phase to be a single time step. In Section 4, we give a specific buffered Algorithm  $A_1$  in which each phase contains  $O(\log(n + N))$  time steps. During a single phase of Algorithm  $A$ , a packet  $\pi$  may perform one of four actions (in  $G'$ ):

- (i) Remain in the buffer of its current node. **[Buffering]**
- (ii) Move from its current node to a neighboring node. **[Packet Transfer]**
- (iii) Be injected into the network at its source node. **[Injection]**
- (iv) Move to and be absorbed in its destination node. **[Absorbtion]**

### 3.2 Bufferless Routing in $G$ – Algorithm $B$

Algorithm  $B$  emulates the phases of Algorithm  $A$  (which is faster than emulating the individual time steps of Algorithm  $A$ ). Algorithm  $B$  emulates the buffering of packets and their transfer from node to node using an  $[\alpha, \beta]$ -partition of  $G$ , where  $\alpha = 2\gamma$ . (We assume that  $2\gamma \leq |E|$  and by Theorem 1, we can set  $\beta = 6\gamma - 3$ .) In Algorithm  $A$ , when a packet is buffered in a node  $r_i$  of  $G'$ , then Algorithm  $B$  emulates this by letting the packet circulate in the edges of region  $R_i$  in  $G$ . When in Algorithm  $A$  a packet is transferred from node  $r_i$  to node  $r_j$  of  $G'$ , in Algorithm  $B$  the packet is transferred from region  $R_i$  to region  $R_j$  in  $G$ . Similarly, algorithm  $B$  handles the packet injection and absorption. Next we describe the emulation in more detail.

*Phases and Rounds.* Let  $\Phi$  denote the number of phases in Algorithm  $A$ . In Algorithm  $B$ , time is divided into  $\Phi$  phases. Each phase of  $B$  emulates a phase of  $A$ . In order to perform the emulation of a phase, Algorithm  $B$  further divides each phase into  $\Sigma$  rounds, where  $\Sigma$  is defined below. The duration of each round is  $T_r = 4\beta^2 + 4\beta$  time steps. Thus the bufferless algorithm runs for  $\Phi \cdot \Sigma \cdot T_r$  time steps in total.

For the duration of a round, a region is either in the *sending* or the *receiving* state – we say that the region is sending, or receiving. In the emulation, when a packet has to be transferred from one region to the next, the first region should be sending while the other receiving. We guarantee that for any pair of adjacent nodes there is a round in each phase in which one region is sending and the other is receiving (and vice-versa), as follows.

In order to determine if a region is sending or receiving, we first obtain a vertex coloring of  $G'$ . Let  $\delta_i$  denote the color (non-negative integer in binary representation) assigned to node  $r_i$  in  $G'$  (which will also be the color of region  $R_i$ ), and let  $\delta$  denote the maximum color we obtain from the vertex coloring. Note that  $\delta \leq n'$ , where  $n' = |V'| \leq |E|/\alpha$ . Let  $\sigma$  denote the number of bits in  $\delta$ ,  $\sigma = \lceil \log \delta \rceil \leq \lceil \log n' \rceil$ . By pre-padding with zeros, we assume that every  $\delta_i$  has  $\sigma$  bits. We define the *state parameter*  $\mathbf{x}_i$  for region  $R_i$  to be the  $2\sigma$ -bit integer  $\bar{\delta}_i \delta_i$ , where  $\bar{\delta}_i$  is the binary complement of  $\delta_i$ . We use the notation  $\mathbf{x}_i(k)$  to denote the  $k$ -th bit of  $\mathbf{x}_i$ . We set  $\Sigma = 2\sigma \leq 2\lceil \log n' \rceil$ , i.e., each phase in Algorithm  $B$ , consists of  $2\sigma$  rounds,  $\omega_1, \omega_2, \dots, \omega_{2\sigma}$ . During round  $\omega_k$ , if  $\mathbf{x}_i(k) = 0$  then region  $R_i$  is sending, otherwise, if  $\mathbf{x}_i(k) = 1$ , then region  $R_i$  is receiving. Our assignment of colors ensures that during every phase, a region can send or receive from each of its neighbors.

**Lemma 2.** *If  $R_i$  and  $R_j$  are adjacent, then during every phase  $\phi$ , there is at least one round  $\omega_s$  ( $\omega_r$ ) in which  $R_i$  is sending (receiving) and  $R_j$  is receiving (sending).*

*Proof.* Since  $R_i$  and  $R_j$  are adjacent,  $\delta_i$  and  $\delta_j$  must differ at some bit  $s$ ,  $0 \leq s \leq \sigma - 1$ . Thus, rounds  $s$  and  $s + \sigma$  satisfy the requirements, since  $\overline{\mathbf{x}_i(s + \sigma)} = \mathbf{x}_i(s) = \mathbf{x}_j(s) = \mathbf{x}_j(s + \sigma)$ .

*Packet Circulation.* Packet circulation is a basic function for the emulation. During packet circulation, a packet  $\pi$  repeatedly follows the Euler tour of the region  $R_i$  that it is in: at each time step, packet  $\pi$  follows the next edge in the Euler tour; when  $\pi$  reaches the end of the Euler tour it continues from the beginning of the tour, and so on. At the time step in which packet  $\pi$  traverses an edge  $e \in \psi_i$ , we say that  $e$  is the *current* edge of  $\pi$ .

At each round of a phase, a region is either sending or receiving. The speed at which a packet circulates in its region depends on whether the region is sending or receiving. If the region is receiving, then the packet follows the Euler tour in the normal fashion.

If the region is sending, then the packet moves at an effectively slower speed as follows. At time step 0 (the beginning of the round), suppose that  $\pi$  is at node  $u$  with current edge  $e = (u, v) \in \psi_i$ . At time step 0, packet  $\pi$  follows its current edge  $(u, v)$  and at time step 1,  $\pi$  appears in node  $v$ . At time step 1, suppose that its new current edge in  $\psi_i$  is  $(v, w)$ ; the packet *does not* follow its new current edge in  $\psi_i$ , but instead it follows edge  $(v, u)$  from  $v$  back to  $u$ , and thus at time step 2, it appears back in node  $u$ . Thus after two time steps, the packet has effectively not moved. We call such an operation an *oscillation*, and we say that packet  $\pi$  oscillates on its current edge in the Euler path. The time period of the oscillation is 2 time steps. The packet continues in this fashion for subsequent time steps, so at even time steps  $t = 2i$ , it appears in node  $u$ , and at odd time steps  $t = 2i + 1$  it appears in node  $v$ , for  $i \geq 0$ . The packet performs  $\beta$  such oscillations on its current edge  $e$ , and so after  $2\beta$  time steps, the packet appears at  $u$  and follows edge  $e$  for the last time. At time step  $T_s = 2\beta + 1$ , the packet is now at  $v$  and at this point it stops oscillating on edge  $e$  and begins oscillating on its new current edge  $(v, w) \in \psi_i$ . Thus, after  $T_s$  time steps, the packet advances by one edge in the Euler path of  $\psi_i$ . Consequently, since  $|\psi_i| \leq 2\beta$ , after  $2\beta T_s = 4\beta^2 + 2\beta$  time steps, a packet circulating in region  $R_i$  has oscillated at least once on every edge of  $\psi_i$ .

**Lemma 3.** *After  $4\beta^2 + 2\beta < T_r$  time steps, a packet circulating in a sending region  $R_i$  has oscillated at least once on every edge in  $\psi_i$ .*

Suppose that the directed edge  $e = (u, v) \in \psi_i$ , is an edge in the Euler path of a receiving region  $R_i$ . If at time step  $t$ , no packet has edge  $e$  as its current edge, then we say that  $e$  is *empty*. At each time step, we say that an empty edge is associated with an *empty slot*. Empty slots are similar to packets in that they too circulate – as the packets in a receiving region circulate (forward) in  $\psi_i$ , the empty slots circulate (backward) in  $\psi_i$  at the same rate. They continue to circulate until some packet occupies the empty edge.

*Emulation of Buffering.* Suppose that packet  $\pi$  is buffered at node  $r_i$  of  $G'$  during the execution of phase  $\phi$  of Algorithm A. Assume that in Algorithm B, packet  $\pi$  is in region  $R_i$  of  $G$ . Packet  $\pi$  will circulate in  $R_i$  through the entire phase  $\phi$ .

**Lemma 4.** *If packet  $\pi$  is in  $R_i$  at the end of phase  $\phi - 1$  of bufferless Algorithm B, and in phase  $\phi$  of buffered Algorithm A it is buffered in node  $r_i$ , then in phase  $\phi$  of bufferless Algorithm B, it can be buffered in region  $R_i$  using circulation.*

*Emulation of Packet Transfer.* Suppose that in phase  $\phi$  of Algorithm A, packet  $\pi$  moves from node  $r_i$  to node  $r_j$ . Assume that at the beginning of phase  $\phi$  in Algorithm B, packet  $\pi$  is in region  $R_i$ . During phase  $\phi$  in Algorithm B,  $\pi$  will move from  $R_i$  to  $R_j$  as follows. Packet  $\pi$  will circulate in  $R_i$  until a round  $\omega$  of  $\phi$  in which  $R_i$  is sending and  $R_j$  is receiving (the existence of such a round is guaranteed by Lemma 2).

Since  $r_i$  and  $r_j$  are adjacent in  $G'$ , there exists a node  $u$  which is common to  $R_i$  and  $R_j$ . Since node  $u$  is in  $R_i$ , there exists an edge  $e_i = (u_i, u) \in \psi_i$  on the Euler path of  $R_i$ . Similarly, there exists an edge  $e_j = (u, u_j) \in \psi_j$  on the Euler tour of  $R_j$ . During round  $\omega$ , packet  $\pi$  circulates (in slow mode) in region  $R_i$  along the Euler tour  $\psi_i$ . At some particular slow time step  $\tau$  of the round, the current edge of  $\pi$  will be  $e_i$ . During the course of its  $T_s > \beta$  oscillations on edge  $e_i$ , the packet will appear at the common node  $u$  at the  $\beta + 1$  times  $\tau + 1, \tau + 3, \dots, \tau + 2\beta + 1$ . If at any of these times, the edge  $e_j \in \psi_j$  is an empty slot, i.e., not the current edge of any packet circulating (in normal mode) in  $R_j$ , then  $\pi$  switches from oscillation on edge  $e_i$ , making  $e_j$  its new current edge.  $\pi$  now continues to circulate in  $R_j$  at normal speed. Note that  $\pi$  will have completed its circulation on edge  $e_i$  in at most  $4\beta^2 + 2\beta$  time steps, thus  $\pi$  will enter  $R_j$  within the first  $4\beta^2 + 2\beta$  time steps of round  $\omega$ .

We now show that during round  $\omega$ , for at least one of the time steps  $\tau + 1, \tau + 3, \dots, \tau + 2\beta + 1$ , the edge  $e_j \in \psi_j$  will be an empty slot. Remember that empty slots circulate in  $R_j$  at the rate of one edge per time-step. Thus, if an empty slot is not occupied by any packet during its circulation, then every edge in  $\psi_j$  will become an empty slot at least once during a consecutive  $2\beta$  time steps. In particular, edge  $e_j$  will become an empty slot at least once in the time steps  $\tau + 1, \tau + 2, \tau + 3, \dots, \tau + 2\beta + 1$ . A problem arises if  $e_j$  becomes empty at time  $\tau + k$  where  $k$  is even, because then packet  $\pi$  will not be at node  $u$ , able to utilize this edge. This problem is solved if there is a second *consecutive* empty slot in  $R_j$  that will also not be occupied by any other packet during its circulation. This second empty slot must also appear at least once in the time steps  $\tau + 1, \tau + 2, \tau + 3, \dots, \tau + 2\beta + 1$ , and since both these empty slots cannot appear at  $\tau + k$  for  $k$  even, we are assured that  $\pi$  will be able to transfer into  $R_j$ .

From the previous phase, suppose that there are at most  $\gamma$  packets circulating in  $R_j$ . During the current phase, at most  $\gamma$  more packets will enter  $R_j$ , by definition of the buffered Algorithm A. In the worst case, all the  $\gamma - 1$  packets other than  $\pi$  that will enter have already entered, and none of the packets that are to leave this region in this phase have left yet. In this case there are at most  $2\gamma - 1$  packets that could be circulating in  $R_j$  during round  $\omega$ . Since  $\alpha = 2\gamma$  and there are at least in  $2\alpha = 4\gamma$  edges  $\psi_j$ , we conclude that there are at least  $2\gamma + 1$  empty slots during round  $\omega$ . By the pigeonhole principle, at least two of these empty slots must be consecutive, and we have the following lemma.

**Lemma 5.** *Suppose that in phase  $\phi - 1$  of bufferless Algorithm B, at most  $\gamma$  packets are circulating in region  $R_j$ , and that packet  $\pi$  is circulating in the adjacent region  $R_i$ . Suppose that in buffered Algorithm A, packet  $\pi$  moves from  $r_i$  to  $r_j$  in phase  $\phi$ . Then during phase  $\phi$  of bufferless Algorithm B, packet  $\pi$  can be transferred (using circulation) from region  $R_i$  to  $R_j$ .*

*Emulation of Injection.* Suppose that  $\pi$  is a packet that is to be injected into the network in Algorithm A. Let  $p$  be the path of  $\pi$  in  $G$ , and let  $e$  be the first edge in this path, and  $u$  the injection node. Suppose that  $e \in R_i$  – note also that  $u \in R_i$ . In this case,  $\pi$  is injected into node  $r_i$  in  $G'$ . Suppose that  $\pi$  is injected into  $r_i$  during phase  $\phi$  of buffered Algorithm A. Then  $\pi$  will be injected into  $R_i$  in phase  $\phi$  of bufferless Algorithm B during the last round in which  $R_i$  is receiving. After injection, it will circulate in  $R_i$  until the end of phase  $\phi$ . Let  $e = (u, v)$  be an edge on the Euler path  $\psi_i$  of  $R_i$ . We know that from the previous analysis of packet transfer that if  $R_i$  had at most  $\gamma$  packets circulating in phase  $\phi - 1$ , then  $e$  will be an empty slot at least  $2\gamma + 1$  times during every receiving round. At the time that  $e$  becomes empty,  $\pi$  is injected into the network and  $e$  becomes its current edge.  $\pi$  then continues to circulate in  $R_i$ . Note that at least  $\gamma$  packets could be injected into  $R_i$  from the *same* injection node during a single receiving round.

**Lemma 6.** *Suppose that in phase  $\phi - 1$  of bufferless Algorithm B, at most  $\gamma$  packets are circulating in region  $R_i$ . Suppose that packet  $\pi$  has first edge  $e \in R_i$  and that during phase  $\phi$  of buffered Algorithm A, packet  $\pi$  is injected into node  $r_i$ . Then during phase  $\phi$  of bufferless Algorithm B, packet  $\pi$  can be injected into  $R_i$ . Further, at least  $\gamma$  packets can be injected into the same node during a single receiving round.*

*Emulation of Absorption.* Suppose that packet  $\pi$  moves from node  $r_i$  to its destination node  $r_j$  in phase  $\phi$  in buffered Algorithm A. We use the packet transfer emulation to first move the packet from region  $R_i$  to  $R_j$  in phase  $\phi$ . This takes at most  $4\beta^2 + 2\beta$  time steps. Then the packet circulates in the receiving region at normal speed until it reaches its destination node, at which point it is absorbed. Since the packet completes the Euler tour for  $R_j$  in at most  $2\beta$  time steps, the number of time steps to move and be absorbed is  $4\beta^2 + 4\beta \leq T_r$ , giving the following lemma.

**Lemma 7.** *Suppose that in phase  $\phi - 1$  of bufferless Algorithm B, at most  $\gamma$  packets are circulating in region  $R_j$ , and that packet  $\pi$  is circulating in the adjacent region  $R_i$ . Suppose that in phase  $\phi$  of buffered Algorithm A, packet  $\pi$  is absorbed in  $r_j$ . Then, during phase  $\phi$  of bufferless Algorithm B, packet  $\pi$  can be absorbed at its destination node in region  $R_j$ .*

### 3.3 Analysis of Emulation by Bufferless Algorithm B

First, we prove that Algorithm B correctly emulates Algorithm A. We then analyse the routing time of Algorithm B in  $G$  in terms of the routing time of Algorithm A in  $G'$ .

*Correctness.* Assume that  $\alpha = 2\gamma \leq |E|$  in order to guarantee the existence of the  $[\alpha, \beta]$ -partition. Algorithm  $B$  correctly emulates algorithm  $A$  if at the end of every phase  $\phi$ :

- i. In Algorithm  $A$ , packet  $\pi$  is in node  $r_i$  iff in Algorithm  $B$  it is circulating in region  $R_i$
- ii. In algorithm  $A$  packet  $\pi$  is injected (absorbed) at node  $r_i$ , if and only if in Algorithm  $B$  packet  $\pi$  is injected (absorbed) into region  $R_i$ .

We show by induction on  $\phi$  that Algorithm  $B$  correctly emulates Algorithm  $A$ . Observe that when  $\phi = 1$ , Algorithm  $A$  can only inject packets into nodes. The conditions of Lemma 6 are satisfied, and since at most  $\gamma$  packets are injected into a node in  $G'$ , Algorithm  $B$  can successfully inject these packets into the corresponding regions. Suppose that Algorithm  $B$  correctly emulates Algorithm  $A$  up to phase  $\phi_0 \geq 1$ . At the end of phase  $\phi_0$ , there are at most  $\gamma$  packets circulating in any region  $R_i$  since every packet  $\pi$  in node  $r_i$  in the execution of Algorithm  $A$  is in region  $R_i$  in the execution of Algorithm  $B$ . Thus, the conditions of Lemmas 4, 5, 6, and 7 are satisfied for every packet  $\pi$ . Every action that  $\pi$  could make in phase  $\phi_0 + 1$  of Algorithm  $A$  can now be emulated in phase  $\phi_0 + 1$  of Algorithm  $B$ . By induction, we have the following theorem.

**Theorem 2 (Correctness of Emulation).** *Algorithm  $B$  correctly emulates in  $G$  every phase in the execution of Algorithm  $A$  in  $G'$ . Each packet in Algorithm  $B$  follows a path from its source to destination, hence Algorithm  $B$  solves routing problem  $Q_s$  without buffers.*

*Routing Time.* Let  $rt_B(Q_s)$  be the routing time for Algorithm  $B$  to solve routing problem  $Q_s$ . Let  $\Phi_A(Q')$  be the number of phases used by Algorithm  $A$  to solve routing problem  $Q'$ . Since Algorithm  $B$  emulates Algorithm  $A$  phase for phase, the number of phases of algorithm  $B$  is also  $\Phi_A(Q')$ . The routing time is therefore given by  $\Phi_A \cdot \Sigma \cdot T_r$ . Since  $T_r = 4\beta^2 + 4\beta$ ,  $\beta = 6\gamma - 3$  and  $\Sigma = 2\lceil \log \delta \rceil$ , we obtain:

**Theorem 3 (Routing Time of Emulation).**  $rt_B(Q_s) = \Theta(\Phi_A(Q') \cdot \gamma^2 \cdot \log \delta)$ .

Since  $\delta \leq |E|/\alpha = O(n^2)$ , from Theorem 3, we have that  $rt_B(Q_s) = O(\Phi_A(Q') \cdot \gamma^2 \cdot \log n)$

## 4 A Randomized Buffered Algorithm

We give a buffered algorithm that can be used to obtaining bufferless routing on arbitrary networks. Since the per-node buffer size enters into the routing time of the bufferless emulation, it is necessary to have buffered algorithms that limit the amount of per-node buffering. We refer to this algorithm as Algorithm  $A_1$ .

Algorithm  $A_1$  is a randomized routing algorithm for routing problems with simple paths, in arbitrary networks. Let  $Q' = (G', \Pi, P')$  be a routing problem

---

**Algorithm 1** Buffered Algorithm  $A_1$ 

---

- 1: Divide time into phases of length  $\gamma$  time steps.
  - 2: **for** Each packet  $\pi$  **do**
  - 3:    $\pi$  selects uniformly at random an injection phase  $\phi_\pi$  between phases 1 and  $12\overline{C'}/\gamma$ ;
  - 4:   Packet  $\pi$  is injected at the first time step of phase  $\phi_\pi$ ;
  - 5:   Packet  $\pi$  follows its path at the speed of one edge per phase;
- 

with acyclic paths  $P'$  on an arbitrary graph  $G' = (V', E')$ . Let  $\overline{C'}$  be the node-congestion and  $D'$  the dilation. Let  $N$  be the number of packets and  $n'$  the size of  $V'$ . Algorithm  $A_1$  uses buffers of size  $\gamma = 6 \log(n' + 2N)$ .

We show that with high probability, Algorithm  $A_1$  successfully routes the packets, and at the same time satisfies the requirements in Section 3.1.

**Theorem 4 (Routing Time of Algorithm  $A_1$ ).** *With probability at least  $1 - O(1/(n' + 2N))$ , Algorithm  $A_1$  solves routing problem  $Q'$  in at most  $12\overline{C'}/\gamma + D'$  phases. The node-buffer size required is  $\gamma = 6 \log(n' + 2N)$ .*

## 5 A Universal Bufferless Algorithm

We use buffered Algorithm  $A_1$  to construct bufferless Algorithm  $B_1$  for arbitrary networks. Algorithm  $B_1$  emulates Algorithm  $A_1$ . The buffer size used by algorithm  $A_1$  is  $\gamma = 6 \log(n' + 2N)$ . Since  $n' \leq |E|/\alpha$ , in order to guarantee the existence of an  $[\alpha, \beta]$ -partition, we assume that  $\alpha \leq |E|$ . Since  $\alpha = 2\gamma$ , we assume that  $12 \log(|E|/\alpha + 2N) \leq |E|$ . It is sufficient that  $2N \leq 2^{|E|/12} - |E|$ .

Suppose  $2N \leq 2^{|E|/12} - |E|$ . Since  $n' \leq n^2/2$ ,  $\gamma \leq 6 \log(n^2/2 + 2N)$ , independent of  $G'$ . Combining Theorems 3 and 4, and the fact that in the emulation,  $\Phi_{B_1}(Q_s) = \Phi_{A_1}(Q')$ , we obtain that  $rt_{B_1}(Q_s) = O((12\overline{C'}/\gamma + D') \cdot \log \delta \cdot \log^2(n + 2N))$ . Using Lemma 1 and the facts that  $\beta \leq 6\gamma$  and  $\delta \leq n' = O(n^2)$ , we obtain that  $rt_{B_1}(Q_s) = O((C + D) \cdot \log n \cdot \log^2(n + N))$ , with probability at least  $1 - O(1/(n' + N))$ .

Consider now the case when  $2N \geq 2^{|E|/12} - |E|$ . We can send the  $N$  packets of routing problem  $Q_s$  on  $G$  to their destinations one after the other. Each packet takes time  $O(D)$  to be delivered to its destination, and thus the total routing time to send all the packets is  $O(DN)$ . Clearly,  $C \geq N/|E|$ , and thus  $C \geq (2^{|E|/12} - |E|)/2|E|$ . Since  $|E| = O(\log(N))$  and  $D \leq |E|$ , the routing time is  $ND \leq CD|E| = O(C \log^2(N))$ . This simple algorithm can easily be converted to a distributed algorithm with the same routing time.

Combining the above results for both cases of the number of packets, we obtain the main result of this paper:

**Theorem 5 (Routing Time of Bufferless Algorithm).**  $rt_{B_1}(Q_s) = O((C + D) \cdot \log n \cdot \log^2(n + N))$ , with probability at least  $1 - O(1/(n' + N))$ .

## 6 Discussion

We have presented a distributed algorithm for routing packets in bufferless networks. Our algorithm is based on the emulation of algorithms with buffers. We partition the original graph into regions, and construct a respective region graph. Each region serves the purpose of a buffer. We then consider an algorithm with buffers on the region graph, and emulate this algorithm by circulating the packets in the regions, and thus avoiding the need of buffers. With this technique, the resulting routing time of our algorithm is  $O((C + D) \cdot \log^3(n + N))$ , which is poly-logarithmic factors away from the optimal for the given paths.

For a particular (source-destination) routing problem, we can define the *bufferless competitive ratio* as the ratio between the best possible routing time of a bufferless algorithm and the best possible routing time of a buffered algorithm. Our result shows that the bufferless competitive ratio is at most  $\log^3(n + N)$  for *any* routing problem. A interesting open problem is to improve this bound.

## References

1. Micah Adler, Sanjeev Khanna, Rajmohan Rajaraman, and Adi Rosen. Time-constrained scheduling of weighted packets on trees and meshes. In *Proceedings of 11th ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, pages 1–12, 1999.
2. N. Alon, F.R.K. Chung, and R.L.Graham. Routing permutations on graphs via matching. *SIAM Journal on Discrete Mathematics*, 7(3):513–530, 1994.
3. Stephen Alstrup, Jacob Holm, Kristian de Lichtenberg, and Mikkel Thorup. Direct routing on trees. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 98)*, pages 342–349, 1998.
4. A. Bar-Noy, P. Raghavan, B. Schieber, and H. Tamaki. Fast deflection routing for packets and worms. In *Proceedings of the Twelfth Annual ACM Symposium on Principles of Distributed Computing*, pages 75–86, Ithaca, New York, USA, August 1993.
5. P. Baran. On distributed communications networks. *IEEE Transactions on Communications*, pages 1–9, 1964.
6. A. Ben-Dor, S. Halevi, and A. Schuster. Potential function analysis of greedy hot-potato routing. *Theory of Computing Systems*, 31(1):41–61, January/February 1998.
7. Petra Berenbrink and Christian Scheideler. Locally efficient on-line strategies for routing packets along fixed paths. In *Proceedings of the Tenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 112–121, N.Y., January 17–19 1999. ACM-SIAM.
8. Sandeep N. Bhatt, Gianfranco Bilardi, Geppino Pucci, Abhiram G. Ranade, Arnold L. Rosenberg, and Eric J. Schwabe. On bufferless routing of variable-length message in leveled networks. *IEEE Trans. Comput.*, 45:714–729, 1996.
9. J. T. Brassil and R. L. Cruz. Bounds on maximum delay in networks with deflection routing. *IEEE Transactions on Parallel and Distributed Systems*, 6(7):724–732, July 1995.
10. A. Broder and E. Upfal. Dynamic deflection routing on arrays. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing*, pages 348–358, May 1996.

11. C. Busch.  $\tilde{O}$ (Congestion + Dilation) hot-potato routing on leveled networks. In *Proceedings of the Fourteenth ACM Symposium on Parallel Algorithms and Architectures*, pages 20–29, August 2002.
12. C. Busch, M. Herlihy, and R. Wattenhofer. Hard-potato routing. In *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, pages 278–285, May 2000.
13. Costas Busch, Malik Magdon-Ismael, Marios Mavronicolas, and Paul Spirakis. Direct routing. In *Proceedings of the 12th Annual European Symposium on Algorithms ESA 2004*, Bergen, Norway, September 2004.
14. Costas Busch, Malik Magdon-Ismael, Marios Mavronicolas, and Roger Wattenhofer. Near-optimal hot-potato routing on trees. In *Proceedings of the 10th International Conference on Parallel and Distributed Computing (Euro-par)*, September 2004.
15. Robert Cypher, Friedhelm Meyer auf der Heide, Christian Scheideler, and Berthold Vöcking. Universal algorithms for store-and-forward and wormhole routing. In *Proceedings of the 28th ACM Symp. on Theory of Computing (STOC)*, pages 356–365, 1996.
16. U. Feige and P. Raghavan. Exact analysis of hot-potato routing. In IEEE, editor, *Proceedings of the 33rd Annual Symposium on Foundations of Computer Science*, pages 553–562, Pittsburgh, PN, October 1992.
17. Ronald I. Greenberg and Hyeong-Cheol Oh. Universal wormhole routing. *IEEE Transactions on Parallel and Distributed Systems*, 8(3):254–262, 1997.
18. B. Hajek. Bounds on evacuation time for deflection routing. *Distributed Computing*, 1:1–6, 1991.
19. Ch. Kaklamanis, D. Krizanc, and S. Rao. Hot-potato routing on processor arrays. In *Proceedings of the 5th Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 273–282, Velen, Germany, June 30–July 2, 1993.
20. F. T. Leighton, B. M. Maggs, and S. B. Rao. Packet routing and job-scheduling in  $O(\text{congestion} + \text{dilation})$  steps. *Combinatorica*, 14:167–186, 1994. (preliminary version appears in FOCS 1988).
21. Friedhelm Meyer auf der Heide and Christian Scheideler. Routing with bounded buffers and hot-potato routing in vertex-symmetric networks. In *Proceedings of the Third Annual European Symposium on Algorithms*, volume 979 of LNCS, pages 341–354, Corfu, Greece, 25–27 September 1995.
22. Friedhelm Meyer auf der Heide and Berthold Vöcking. Shortest-path routing in arbitrary networks. *Journal of Algorithms*, 31(1):105–131, April 1999.
23. Rafail Ostrovsky and Yuval Rabani. Universal  $O(\text{congestion} + \text{dilation} + \log^{1+\epsilon} N)$  local control packet switching algorithms. In *Proceedings of the 29th Annual ACM Symposium on the Theory of Computing*, pages 644–653, New York, May 1997.
24. Yuval Rabani and Éva Tardos. Distributed packet switching in arbitrary networks. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing*, pages 366–375, Philadelphia, Pennsylvania, 22–24 May 1996.
25. Alan Roberts, Antonios Symvonis, and David R. Wood. Lower bounds for hot-potato permutation routing on trees. In *Proceedings of the 7th Int. Coll. Structural Information and Communication Complexity, SIROCCO*, pages 281–295, 20–22 June 2000.
26. L. Zhang. Optimal bounds for matching routing on trees. In *Proceedings of the 8th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 445–453, 1997.