

# On the Emergence of Shortest Paths by Reinforced Random Walks

Daniel R. Figueiredo<sup>1</sup> and Michele Garetto<sup>2</sup>

<sup>1</sup>*Computer Science and Syst. Eng. Dept., Federal University of Rio de Janeiro (UFRJ), Brazil*

<sup>2</sup>*Computer Science Department, University of Torino, Italy*

## Abstract

The co-evolution between network structure and functional performance is a fundamental and challenging problem whose complexity emerges from the intrinsic interdependent nature of structure and function. Within this context, we investigate the interplay between the efficiency of network navigation (i.e., path lengths) and network structure (i.e., edge weights). We propose a simple and tractable model based on iterative biased random walks where edge weights increase over time as function of the traversed path length. Under mild assumptions, we prove that biased random walks will eventually only traverse shortest paths in their journey towards the destination. We further characterize the transient regime proving that the probability to traverse non-shortest paths decays according to a power-law. We also highlight various properties in this dynamic, such as the trade-off between exploration and convergence, and preservation of initial network plasticity. We believe the proposed model and results can be of interest to various domains where biased random walks and de-centralized navigation have been applied.

## 1 Introduction

The interplay between network structure (nodes, edges, weights) and network function (high level features enabled by the network) is a fundamental and challenging problem present in a myriad of systems ranging from biology to economics and sociology. In many complex systems network structure and network function co-evolve interdependently: while network structure constraints functional performance, the drive for functional efficiency pressures the network structure to change over time. Within this tussle, *network activity* (i.e., basic background processes running on the network) plays a key role in tying function and structure: in one hand, function execution often requires network activity, while in the other hand network structure often constraints network activity.

Given the complexity of co-evolution, simple and tractable models are often used to understand and reveal interesting phenomena. In this paper, we focus on *network navigation*, proposing and analyzing a simple model that captures the interplay between function and structure. Our case-study embodies repetition, plasticity, randomization, valuation and memory which are key ingredients for evolution: repetition and memory allow for learning; plasticity and randomization for exploring new possibilities; valuation for comparing alternatives. Moreover, in our case-study co-evolution is enabled by a single and simple network activity process: *biased random walks*, where time-varying edge weights play the role of memory.

Network navigation (also known as routing) refers to the problem of finding short paths in networks and has been widely studied due to its importance in various contexts. Efficient network navigation can be achieved by running centralized or distributed algorithms. Alternatively, it can also be achieved when running simple greedy algorithms over carefully crafted network topologies. But can efficient navigation emerge without computational resources and/or specifically tailored topologies?

A key contribution of our work is to answer affirmatively the above question by means of Theorem 1, which states that under mild conditions efficient network navigation always emerges through the repetition of extremely simple network activity. More clearly, a biased random walk will eventually only take

paths of minimum length, independently of network structure and initial weight assignment. Beyond its long term behavior, we also characterize the system transient regime, revealing interesting properties such as the power-law decay of longer paths, and the (practical) preservation of initial plasticity on edges far from ones on the shortest paths. The building block for establishing the theoretical results of this paper is the theory of Pólya urns, applied here by considering a network of urns.

We believe the proposed model and its analysis could be of interest to various domains where some form of network navigation is present and where random walks are used as the underlying network activity, such as computer networking [16, 8, 14], animal movement in biology [3, 21], memory recovery in the brain [18, 1, 19]. Moreover, our results can enrich existing theories such as Ant Colony Optimization (ACO) meta-heuristic [6, 5], Reinforcement Learning (RL) theory [22], and Edge Reinforced Random Walks (ERRW) theory [4, 17] – see related work in Section 3.

## 2 Model

We consider an arbitrary (fixed) network  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of vertices and  $\mathcal{E}$  is a set of directed edges among the vertices. We associate a weight (a positive real value)  $w_{i,j}$  to every directed edge  $(i, j) \in \mathcal{E}$ . Edge weights provide a convenient and flexible abstraction for structure, specially when considering evolution. Finally, a pair of fixed nodes  $s, d \in \mathcal{V}$  are chosen to be the source and destination, respectively. But how to go from  $s$  to  $d$ ?

We adopt a very simple network activity model to carry out the function of navigation: weighted random walks (WRW). Specifically, a sequence of random walks, indexed by  $n = 1, 2, \dots$ , is executed on the network one after the other. Each WRW starts at  $s$  and steps from node to node until it hits  $d$ . At each visited node, the WRW randomly follows an outgoing edge with probability proportional to its edge weight. We assume that weights on edges remain constant during the execution of a single WRW, and that decisions taken at different nodes are independent from each other.

Once the WRW reaches the destination and stops, edge weights are updated, thus impacting the behavior of the next WRW in the sequence. In particular, edges on the path followed by the WRW are rewarded (reinforced) by increasing their weights with a positive amount which depends on the length of the path taken (expressed in number of hops). Let  $f : \mathbb{N}^+ \rightarrow \mathcal{R}^+$  be some positive function of the path length, hereinafter called the reward function.

We consider two different ways in which edges are reinforced:

- **single-reward model:** each edge belonging to the path followed by the WRW is rewarded once, according to function  $f(\cdot)$ .
- **multiple-reward model:** each edge belonging to the path followed by the WRW is rewarded according to function  $f(\cdot)$  for each time the edge was traversed.

Throughout the paper, we will interpret  $n$ , the number of random walks that have gone from  $s$  to  $d$  as a discrete time step. Thus, by co-evolution of the system we actually mean what happens to the network structure (i.e., weights) and navigation (i.e., path lengths) as  $n \rightarrow \infty$ .

Let  $w_{i,j}[n]$  be the weight on edge  $(i, j)$  at time  $n$  (right after the execution of the  $n$ -th WRW but before the  $(n + 1)$ -th WRW starts). Let  $\mathcal{P}_n$  denote the sequence of edges (i.e., the path) traversed by the  $n$ -th WRW and  $L_n = |\mathcal{P}_n|$  the path length (in number of hops)<sup>1</sup>.

After reaching the destination, the weight of any distinct edge  $(i, j)$  in  $\mathcal{P}_n$  is updated according to the rule:

$$w_{i,j}[n] = w_{i,j}[n - 1] + u_{i,j}(\mathcal{P}_n) \cdot f(L_n) \quad (1)$$

---

<sup>1</sup>In this paper we assume that the cost to traverse any edge is equal to 1, but results can be immediately generalized to the case in which a generic (positive) cost  $c_{i,j}$  is associated to each edge  $(i, j)$ .

where  $u_{i,j}(\mathcal{P}_n) = 1$  under the single-reward model, whereas  $u_{i,j}(\mathcal{P}_n)$  equals the number of times edge  $(i, j)$  appears in  $\mathcal{P}_n$ , under the multiple reward model. We also allow for the event that the WRW does not reach the destination because it ‘gets lost’ in a part of the network from which  $d$  is no longer reachable. In this case, we assume that no edge is updated by the WRW who fails to reach  $d$ .

Note that our model has the desirable ingredients for co-evolution: edge set  $\mathcal{E}$  and initial weights provide plasticity and WRW provides randomization, which allows for exploring alternative paths; edge weights provide memory and the sequence of WRW provides repetition, which enables learning; path length taken by WRW provides valuation, which allows for comparing alternative paths. Moreover, note that functional performance induces structural changes through network activity as navigation (traversed path) changes edge weights, while network structure constraints function, as edge weights influence observed path lengths. Thus, our model captures the essence of co-evolution. But will efficient navigation emerge? In particular, which paths  $\mathcal{P}_n$  are taken as  $n$  increases?

### 3 Related work

The problem of finding shortest paths in networks is, of course, a well understood problem in graph theory and computer science, for which efficient algorithms are available, both centralized (e.g., Dijkstra) and distributed (e.g., Bellman-Ford). Our approach follows in the second category (distributed), as it does not require knowledge of the topology, however we stress that our goal is not to propose yet another way to compute shortest paths in network (actually, the convergence of our process is slower than that of Bellman-Ford), but to show that shortest paths can naturally emerge from the repetition of a simple and oblivious network activity which does not require computational/memory resources on the nodes. As such, our model is more tailored to biological systems, rather than technological networks.

The celebrated work of Kleinberg [13] was probably the first to show that efficient navigation is indeed feasible through a simple greedy strategy based solely on local information, but under the stringent assumption that the network exhibits a very particular structure. Greedy algorithms can also lead to efficient network navigation under distributed hash tables (DHTs), but again this requires the network to exhibit a very particular topology [10].

The idea of reinforcing edges along paths followed by random walks is surely reminiscent of Ant Colony Optimization (ACO), a biologically-inspired meta-heuristic for exploring the solution space of complex optimization problems which can be reduced to finding good paths through graphs [6, 5]. Although some versions of ACO can be proved to converge to the global optimum, their analysis turns out to be complicated and mathematically non-rigorous, especially due to *pheromone evaporation* (i.e., weights on edges decrease in the absence of reinforcement). Moreover, like most meta-heuristics, it is very difficult to estimate the theoretical speed of convergence. In contrast to ACO, our model is simpler and has the modest goal of revealing shortest paths in a given network, instead of exploring a solution space. Moreover, we do not introduce any evaporation, and we exploit totally different techniques (the theory of Pólya urns) to establish our results, including the transient behavior (convergence) of the system.

In Reinforcement Learning (RL), the problem of finding an optimal policy through a random environment has also been tackled using Monte Carlo methods that reinforce actions based on earned rewards, such as the  $\epsilon$ -soft policy algorithm [22]. Under a problem formulation with no terminal states and expected discounted rewards, it can be rigorously shown that an iterative algorithm converges to the optimal policy [23]. However, in general and more applicable scenarios, the problem of convergence to optimal policies is still an open question, with most algorithms settling for an approximate solution. Although lacking the notion of action set, our model is related to RL in the sense that it aims at finding paths accumulating the minimum cost, through an unknown environment, using a Monte-Carlo method. Our convergence results (convergence in probability) and the techniques used in the analysis (Pólya urns) are fundamentally different from what is commonly found in RL theory, and could thus be useful to tackle problems in this area.

Edge Reinforcement Random Walks (ERRW) is a mathematical modeling framework consisting of a weighted graph where weights evolve over time according to steps taken by a random walker [4, 17]. In ERRW, a single random walk moves around without having any destination and without being restarted.

Moreover, an edge weight is updated immediately after traversal of the edge, according to functions based on local information. Mathematicians have studied theoretical aspects of ERRW such as the convergence of the network structure (relative weights) and the recurrence behavior of the walker (whether it will continue to visit every node in the long run, or get trapped in one part). Similarly to our model, a key ingredient in the analysis of ERRW is the Pólya urn model, specially on directed networks. In contrast to us, ERRW model was not designed to perform any particular function and thus does not have an objective. Our model is substantially different from traditional ERRW, and we believe it could suggest a concrete application as well as new directions to theoreticians working on ERRW.

Animal movement is a widely studied topic in biology to which probabilistic models have been applied, including random walk based models [3, 21]. In particular, in the context of food foraging, variations of ERRW models have been used to capture how animals search and traverse paths to food sources. A key difference in such variations is a *direction vector*, an information external to the network (but available on all nodes) that provides hints to the random walk. Such models have been used to show the emergence of relatively short paths to food sources, as empirically observed with real (monitored) animals. In contrast, we show that shortest paths (and not just short) can emerge even when external information is not available.

Understanding how neurons in the brain connect and fire to yield higher level functions like memory and speech is a fundamental problem that has recently received much attention and funding [20, 9]. Within this context, random walk based models have been proposed and applied [19, 1] along with models where repeated network activity modifies the network structure [18]. In particular, the latter work considers a time varying weighted network model under a more complex rule (than random walks) for firing neurons to show that the network structure can arrange itself to perform better function. We believe our work can provide building blocks in this direction since our simple model for a time varying (weighted) network also self-organizes to find optimal paths.

Biased random walks have also been applied to a variety of computer networking architectures [16, 8, 14], with the goal of designing self-organizing systems to locate, store, replicate and manage data in time-varying scenarios. We believe our model and findings could be of interest in this area as well.

## 4 Main finding

Let  $L_{\min}$  be the length of the shortest path in graph  $\mathcal{G}$  connecting source node  $s$  to destination node  $d$ . Denote by  $\mathcal{P}_n$  the path taken by the  $n$ -th WRW, and by  $\mathcal{P}$  an arbitrary path from  $s$  to  $d$ , of length  $L_{\mathcal{P}}$ .

**Theorem 1.** *Given a weighted directed graph  $\mathcal{G}$ , a fixed source-destination pair  $s$ - $d$  (such that  $d$  is reachable from  $s$ ), an initial weight assignment (such that all initial weights are positive), consider an arbitrary path  $\mathcal{P}$  from  $s$  to  $d$ . Under both the single-reward model and the multiple-reward model, provided that the reward function  $f(\cdot)$  is a strictly decreasing function of the path length, as the number  $n$  of random walks performed on the graph tends to infinity, we have:*

$$\lim_{n \rightarrow \infty} \mathbb{P}\{\mathcal{P}_n = \mathcal{P}\} = \begin{cases} c(\mathcal{P}), & \text{if } L_{\mathcal{P}} = L_{\min} \\ 0, & \text{if } L_{\mathcal{P}} > L_{\min} \end{cases}$$

where  $c(\mathcal{P})$  is a random variable taking values in  $(0, 1]$ , that depends on the specific shortest path  $\mathcal{P}$ .

The above theorem essentially says that *all* shortest paths are taken with non-vanishing probability, while *all* non-shortest paths are taken with vanishing probability, as  $n \rightarrow \infty$ . Note however that the probability that a specific shortest path is taken is a random variable, in the sense that it depends on the ‘system run’ (system sample path).

**Remark 1.** *The asymptotic property stated in Theorem 1 is very robust, as it holds for any directed graph, any strictly decreasing function  $f(\cdot)$ , and any (valid) initial weights on the edges. Note instead that the (asymptotic) distribution of  $c(\mathcal{P})$ , for a given shortest path  $\mathcal{P}$ , as well as the convergence rate to it, depends strongly on the update function  $f(\cdot)$ , on the graph structure, and on the initial conditions on the edges.*

**Remark 2.** We will see in the proof of Theorem 1 that the assumption of having a strictly decreasing function  $f(\cdot)$  can be partially relaxed, allowing the reward function to be non-increasing for  $L > L_{\min}$ .

## 5 Preliminaries

### 5.1 Definitions

The following definitions for nodes and edges play a central role in our analysis.

**Definition 1 (decision point).** A decision point is a node  $i \in \mathcal{V}$ , reachable by  $s$ , that has more than one outgoing edge that can reach  $d$ .

**Remark 3.** Clearly, we can restrict our attention to nodes that are decision points, since all other nodes are either never reached by random walks originating in  $s$ , have zero or one outgoing edge (having no influence on the random walk behavior), or their outgoing edges are never reinforced since the destination cannot be reached from them.

**Definition 2 ( $\alpha$ -edge and  $\beta$ -edge).** An outgoing edge of decision point  $i$  is called an  $\alpha$ -edge if it belongs to some shortest path from  $i$  to  $d$ , whereas it is called a  $\beta$ -edge if it does not belong to any shortest path from  $i$  to  $d$ .

Note that every outgoing edge of a decision point is either an  $\alpha$ -edge or  $\beta$ -edge. Let  $q_\alpha(i, j, n)$  denote the probability that the random walk, at time  $n$ , takes a shortest path from  $i$  to  $d$  after traversing the  $\alpha$ -edge  $(i, j)$ . Let  $q_\beta(i, j, n)$  denote the probability that the random walk, at time  $n$ , will *not* return back to node  $i$  after traversing the  $\beta$ -edge  $(i, j)$ . Note that the above probabilities depend, in general, on the considered edge, on the network structure and on the set of weights at time  $n$ .

**Definition 3 ( $\alpha^*$ -edge and  $\beta^*$ -edge).** An  $\alpha^*$ -edge is an  $\alpha$ -edge such that, after traversing it, the random walk takes a shortest path to  $d$  with probability 1, and thus  $q_\alpha(i, j, n) = 1$ . A  $\beta^*$ -edge is a  $\beta$ -edge such that, after traversing it, the random walk does not return to node  $i$  with probability 1, and thus  $q_\beta(i, j, n) = 1$ .

Note that  $\alpha^*$ -edge and  $\beta^*$ -edge can occur due solely to topological constraints. In particular, we have an  $\alpha^*$ -edge whenever the random walk, after traversing the edge, can reach  $d$  only through paths of minimum length. In a cycle-free network, all  $\beta$ -edges are necessarily  $\beta^*$ -edges.

### 5.2 The single decision point

As a necessary first step, we will consider the simple case in which there is a single decision point in the network. The thorough analysis of this scenario provides a basic building block towards the analysis of the general case.

We start considering the simplest case in which there are two outgoing edges (edge 1 and edge 2) from the decision point, whose initial weights are denoted by  $w_1[0]$  and  $w_2[0]$ , respectively. Let  $L_1$  and  $L_2$  denote the (deterministic) length of the path experienced by random walks when traversing edge 1 and edge 2, respectively. Correspondingly, let  $\Delta_1 = f(L_1)$  and  $\Delta_2 = f(L_2)$  denote the rewards given the edge 1 and edge 2, respectively.

The mathematical tool used here to analyze this system, especially its asymptotic properties, are Pólya urns [15]. The theory of Pólya urns is concerned with the evolution of the number of balls of different colors (let  $K$  be the number of colors) contained in an urn from which we repeatedly draw one ball uniformly at random. If the color of the ball withdrawn is  $i$ ,  $i = 1, \dots, K$ , then  $A_{i,j}$  balls of color  $j$  are added to the urn,  $j = 1, \dots, K$ , in addition to the ball withdrawn, which is returned to the urn. In general,  $A_{i,j}$  can be deterministic or random, positive or negative. Let  $\mathbf{A}$  be the matrix with entries  $A_{i,j}$ , usually referred to as the *schema* of the Pólya urn.

We observe that a decision point can be described by a Pólya urn, where the outgoing edges represent colors, the edge weight is the number of balls in the urn<sup>2</sup>, and entries  $A_{i,j}$  correspond to edge reinforcements according to taken path lengths (through function  $f(\cdot)$ ). In the simple case with only two edges, we obtain the following schema:

$$\mathbf{A} = \begin{pmatrix} \Delta_1 & 0 \\ 0 & \Delta_2 \end{pmatrix} \quad (2)$$

We first consider the situation in which  $\Delta_1 = \Delta_2$ , which occurs when both edges are part of a shortest path, and thus, both edges are  $\alpha^*$ -edges. A classical result in Pólya urns states that the normalized weight of edge 1 (similarly for edge 2), i.e., the weight on edge 1 divided by the sum of the weights, tends in distribution to a beta distribution:

$$\frac{w_1[n]}{w_1[n] + w_2[n]} \xrightarrow{\mathcal{D}} \beta \left( \frac{w_1[0]}{\Delta_1}, \frac{w_2[0]}{\Delta_2} \right) \quad (3)$$

Note that in this simple case the above beta distribution completely characterizes the asymptotic probability of traversing the shortest path comprising edge 1 (or edge 2). Hence, we obtain a special case of the general result stated in 1, where the random variable  $c(\mathcal{P})$  is a beta distribution which depends both on the update function and the initial weights. Informally, we say that both shortest paths will always ‘survive’, as they will be asymptotically used with a (random) non-zero probability, independent of the sample path taken by the system.

The above result can be directly generalized to the case of  $K$  outgoing edges, all belonging to shortest paths. Indeed, let’s denote the asymptotic normalized weight of edge  $i$  by  $r_i$ :

$$r_i = \lim_{n \rightarrow \infty} \frac{w_i[n]}{\sum_{j=1}^K w_j[n]}$$

Moreover, let  $\alpha_i = \frac{w_i[0]}{\Delta_i}$ . Then it is known that the joint probability density function of the  $r_i$ ’s tends to a Dirichlet distribution with parameters  $\{\alpha_i\}$ .

A useful property of the Dirichlet distribution is aggregation: if we replace any two edges with initial weights  $w_i, w_j$  by a single edge with initial weight  $w_1 + w_2$ , we obtain another Dirichlet distribution where the ‘combined’ edge is associated to parameter  $\alpha_i + \alpha_j$ , i.e., if  $\mathbf{r} = (r_1, \dots, r_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$  then  $\mathbf{r}' = (r_1, \dots, r_i + r_j, \dots, r_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_i + \alpha_j, \dots, \alpha_K)$ . Note that the marginal distribution with respect to any of the edges is, as expected, a beta distribution, i.e.,

$$r_i \sim \beta \left( \alpha_i, \sum_{j=1, j \neq i}^K \alpha_j \right)$$

Let’s now consider outgoing edges that lead to paths of different lengths, starting from the simple situation in which we have just two edges. Without lack of generality, let’s assume that  $\Delta_1 > \Delta_2$  in which case edge 1 is an  $\alpha^*$ -edge and edge 2 is a  $\beta^*$ -edge. The analysis of the corresponding Pólya urn model uses a technique known as *Poissonization* [15]. The basic idea is to embed the discrete-time evolution of the urn in continuous time, associating to each ball in the urn an independent exponential timer with parameter 1. When a timer ‘fires’, the associated ball is drawn, and we immediately perform the corresponding ball additions (starting a new timer for each added ball). The memoryless property of the exponential distribution guarantees that the time at which a ball is drawn is a renewal instant for the system. Moreover, competition among the timers running in parallel exactly produces the desired probability to extract a ball of a given color at the next renewal instant. This means that, if  $t_n$  is the (continuous) time at which the  $n$ -th timer fires, at time  $t_n$  the number of balls in the continuous-time system has *exactly* the same distribution as the number of balls in the original discrete-time system after  $n$  draws. It follows that the asymptotic behavior (as  $t \rightarrow \infty$ ) of the continuous-time system coincides with the asymptotic behavior of the discrete-time system (as  $n \rightarrow \infty$ ), but the continuous-time system is more amenable to analysis, thanks to the independence of all Poisson processes in the urn.

The Poissonization technique leads to the following fundamental result: Let  $\mathbf{w}(t)$  be the (column) vector of edge weights at time  $t$  in the continuous-time system. We have (Theorem 4.1 in [15]):  $\mathbf{E}[\mathbf{w}(t)] =$

<sup>2</sup>Although Pólya urn models have been traditionally developed considering an integer number of balls for each color, analogous results hold in the case of real numbers, when all  $A_{i,j}$  are positive (as in our case).

$e^{\mathbf{A}^T t} \mathbf{w}(0)$ . The above result can be extended to the case in which the entries of schema  $\mathbf{A}$  are independent random variables (independent among them and from one draw to another) by simply substituting  $\mathbf{A}^T$  with  $\mathbf{E}[\mathbf{A}^T]$ :

$$\mathbf{E}[\mathbf{w}(t)] = e^{\mathbf{E}[\mathbf{A}^T]t} \mathbf{w}(0) \quad (4)$$

i.e., by considering a schema in which random entries are replaced by their expectations. This extension will be particularly useful in our context.

## 6 Asymptotic analysis

In this section we prove Theorem 1 first constrained to directed acyclic graphs (DAG), then relaxing to general topologies under the multiple-reward model and finally to the single-reward model.

### 6.1 The DAG case

Let  $\mathcal{G}$  be a directed acyclic graph (DAG) and note that in this case edges are either  $\alpha$ -edges or  $\beta^*$ -edges. Moreover, the absence of cycles forbids traversing an edge more than once, so the single-reward model coincides with the multiple-reward model.

We first introduce the following key lemma.

**Lemma 1.** *Consider a decision point having one or more  $\alpha^*$ -edges and one or more  $\beta^*$ -edges. The normalized weight of any  $\beta^*$ -edge vanishes to zero as  $n \rightarrow \infty$ .*

*Proof.* Let  $\hat{L}$  denote the length of the shortest path from the decision point to  $d$ . Note that this path length is realized by the random walk after following an  $\alpha^*$ -edge. Observe that  $\alpha^*$ -edges can be merged together into a single virtual  $\alpha^*$ -edge whose weight, denoted by  $\hat{w}$ , is defined as the sum of the weights of the merged  $\alpha^*$ -edges. Similarly, we will merge all  $\beta^*$ -edges into a single virtual  $\beta^*$ -edge of weight  $\hat{w}$ , defined as the sum of the weights of the merged  $\beta^*$ -edges.

Let  $\{Z_n, n \geq 1\}$  be the stochastic process corresponding to  $Z_n = \frac{\hat{w}[n]}{\hat{w}[n] + \hat{w}[n]}$ , i.e.,  $Z_n$  is the normalized weight of the virtual merged  $\beta^*$ -edge after the  $n$ -th random walk. We are going to show that  $\lim_{n \rightarrow \infty} Z_n = 0$  which implies that the asymptotic probability to follow any  $\beta^*$ -edge goes to zero as well. The proof is divided into two parts. First, we show that  $\lim_{n \rightarrow \infty} Z_n$  exists almost surely, namely,  $Z_n$  converges to a given constant  $z \in [0, 1]$ . Second, we will show that  $z$  can only be equal to 0. For the first part, we will use Doob's Martingale Convergence Theorem [7], after proving that  $Z_n$  is a super-martingale. Since  $\{Z_n\}$  is discrete time, and  $0 \leq Z_n \leq 1$ , it suffices to prove that  $\mathbb{E}[Z_{n+1} | \mathcal{F}_n] \leq Z_n$ , where the filtration  $\mathcal{F}_n$  corresponds to all available information after the  $n$ -th walk. Now, the normalized weight, at time  $n + 1$ , of any  $\beta^*$ -edge is stochastically dominated by the normalized weight, at time  $n + 1$ , of the same  $\beta^*$ -edge assuming that it belongs to a path of length  $\hat{L} + 1$ . This is essentially the reason why we can merge all  $\beta^*$ -edges into a single virtual  $\beta^*$ -edge belonging to a path of length  $\hat{L} + 1$ . Hence,  $\mathbb{E}[Z_{n+1} | \mathcal{F}_n] \leq \mathbb{E}[Z'_{n+1} | \mathcal{F}_n]$ , where  $Z'_{n+1}$  is the aggregate normalized weight of the virtual  $\beta^*$ -edge. We proceed by considering what can happen when running the  $(n + 1)$ -th walk. Two cases are possible: i) either the random walk does not reach the decision point, in which case  $Z'_{n+1} = Z_n$  since edge weights are not updated, or ii) it reaches the decision point having accumulated a (random) hop count  $\ell_{n+1}$ . In the second case, we can further condition on the value taken by  $\ell_{n+1}$  and prove that  $\mathbb{E}[Z'_{n+1} | \mathcal{F}_n, \ell_{n+1}] \leq Z_n, \forall \ell_{n+1}$ :

$$\begin{aligned} \mathbb{E}[Z'_{n+1} | \mathcal{F}_n, \ell_{n+1}] &= Z_n \frac{\hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1)}{\hat{w}(n) + \hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1)} + (1 - Z_n) \frac{\hat{w}(n)}{\hat{w}(n) + \hat{w}(n) + f(\ell_{n+1} + \hat{L})} = \\ &= Z_n \left[ \frac{\hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1)}{\hat{w}(n) + \hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1)} + \frac{\hat{w}(n)}{\hat{w}(n)} \frac{\hat{w}(n)}{\hat{w}(n) + \hat{w}(n) + f(\ell_{n+1} + \hat{L})} \right] \leq \\ &= Z_n \left[ \frac{\hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1) + \hat{w}(n)}{\hat{w}(n) + f(\ell_{n+1} + \hat{L} + 1) + \hat{w}(n)} \right] = Z_n \quad (5) \end{aligned}$$

where the inequality holds because  $f(\cdot)$  is assumed to be non-increasing.

At last, unconditioning with respect to  $\ell_{n+1}$ , whose distribution descends from  $\mathcal{F}_n$ , and considering also the case in which the random walk does not reach the decision point, we obtain  $\mathbb{E}[Z'_{n+1}|\mathcal{F}_n] \leq Z_n$  and thus  $\mathbb{E}[Z_{n+1}|\mathcal{F}_n] \leq Z_n$ . So far we have proven that  $Z_n$  converges to a constant  $z \in [0, 1]$ . To show that necessarily  $z = 0$ , we employ the Poissonization technique recalled in Section 5.2, noticing again that  $Z_n$  is stochastically dominated by  $Z'_n$ . For the process  $Z'_n$ , we have:

$$\mathbf{A}^T = \begin{pmatrix} f(\ell + \hat{L}) & 0 \\ 0 & f(\ell + \hat{L} + 1) \end{pmatrix}$$

where  $\ell$  is the (random) hop count accumulated at the decision point. We will show later that the normalized weight of any edge in the network converges asymptotically almost surely. Hence,  $\ell$  has a limit distribution, that we can use to compute expected values of the entries in the above matrix:

$$\mathbb{E}[\mathbf{A}^T] = \begin{pmatrix} \mathbb{E}_\ell[f(\ell + \hat{L})] & 0 \\ 0 & \mathbb{E}_\ell[f(\ell + \hat{L} + 1)] \end{pmatrix} = \begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix}$$

obtaining that  $a > d$  when  $f(\cdot)$  is strictly decreasing. At this point, we can just apply known results of Pólya urns' asymptotic behavior (see Theorem 3.21 in [11]), and conclude that the normalized weight of the  $\beta^*$ -edge must converge to zero. Alternatively, we can apply (4) and observe that in this simple case

$$\begin{pmatrix} \mathbf{E}[\hat{w}(t)] \\ \mathbf{E}[\dot{w}(t)] \end{pmatrix} = \begin{pmatrix} e^{at} & 0 \\ 0 & e^{dt} \end{pmatrix} \begin{pmatrix} \hat{w}(0) \\ \dot{w}(0) \end{pmatrix} \quad (6)$$

Therefore the (average) weight of the  $\alpha^*$ -edge increases exponentially faster than the (average) weight of the  $\beta^*$ -edge.  $\square$

Lemma 1 provides the basic building block to prove Theorem 1.

*Proof of Theorem 1 (DAG case).* We sequentially consider the decision points of the network according to the partial topological ordering given by the hop-count distance from the destination. Simply put, we start considering decision points at distance 1 from the destination, then those at distance 2, and so on, until we hit the source node  $s$ . We observe that Lemma 1 can be immediately applied to decision points at distance 1 from the destination. Indeed, these decision points have one (or more)  $\alpha^*$ -edge, with  $\hat{L} = 1$ , connecting them directly to  $d$ , and zero or more  $\beta^*$ -edges connecting them to nodes different from  $d$ . Then, Lemma 1 allows us to conclude that, asymptotically, the normalized weight of the virtual  $\alpha^*$ -edge will converge to 1, whereas the normalized weight of all  $\beta^*$ -edges will converge to zero. This fact essentially allows us to *prune* the  $\beta^*$ -edges of decision nodes at distance 1, and re-apply Lemma 1 to decision points at distance 2 (and so on). Note that after the pruning, an  $\alpha$ -edge of a decision point at distance 2 necessarily becomes an  $\alpha^*$ -edge. As a consequence of the progressive pruning of  $\beta$ -edges, we remove from the graph all edges which do not belong to shortest paths from a given node  $i$  to  $d$  (when we prune a  $\beta$ -edge, we contextually remove also edges that can only be traversed by following the pruned edge, and notice that by so doing we can also remove some  $\alpha$ -edge).

When the above iterative procedure hits the source node  $s$ , we are guaranteed that only shortest paths from  $s$  to  $d$  remain in the residual graph (and all of them). As a consequence, over the residual graph, a random walk starting from  $s$  can only reach  $d$  through a shortest path. Note that the normalized weight of any edge  $(i, j)$  belonging to a shortest path will converge to a random variable  $z_{i,j}$  bounded away from zero. Hence the asymptotic probability to follow any given shortest path  $\mathcal{P}$ , given by the product of normalized weights of its edges, will converge as well to a random variable  $c(\mathcal{P})$  bounded away from zero. Conversely, any path which is not a shortest path cannot 'survive'. Indeed, any such path must traverse at least one decision point and take at least one  $\beta^*$ -edge. However, the above iterative procedure will eventually prune all  $\beta^*$ -edges belonging to the considered non-shortest path, which therefore cannot survive.  $\square$

## 6.2 The multiple reward model in general network

We now consider the case of an arbitrary directed graph possibly with nodes exhibiting (even multiple) self-loops. Moreover, we first focus on the multiple-reward model which is more challenging to analyze, and discuss the single-reward model in Section 6.3.

Essentially, we follow the same reasoning as in the DAG case, by first proving a generalized version of Lemma 1.

**Lemma 2.** *Consider a decision point having one or more  $\alpha^*$ -edges and one or more  $\beta$ -edges. The normalized weight of any  $\beta$ -edge vanishes to zero as  $n \rightarrow \infty$ .*

*Proof.* Similarly to the proof of Lemma 1, we merge all  $\alpha^*$ -edges into a single virtual  $\alpha^*$ -edge with total weight  $\hat{w}$ . Moreover, we merge all  $\beta$ -edges into a single virtual  $\beta$ -edge with weight  $\hat{w}$ , defined as the sum of the weights of the merged  $\beta$ -edges. Such virtual  $\beta$ -edge can be interpreted as the best *adversary* against the virtual  $\alpha^*$ -edge. Clearly, the best  $\beta$ -edge is an outgoing edge that (possibly) brings the random walk back to the decision point over the shortest possible cycle, i.e., a self-loop. It is instead difficult, *a priori*, to establish which is the best possible value of its parameter  $q_\beta(n)$ , i.e., the probability (in general dependent on  $n$ ) the makes the virtual  $\beta$ -edge the best competitor of the virtual  $\alpha^*$ -edge. Therefore, we consider arbitrary values of  $q_\beta(n) \in [0, 1]$  (technically, if  $q_\beta(n) > 0$  then the  $\beta$ -edge cannot be a self-loop, but we optimistically assume that loops have length 1 even in this case). In the following, to ease the notation, let  $q = q_\beta(n)$ . Similarly to the DAG case, we optimistically assume that if the random walk reaches the destination without passing through the  $\alpha^*$ -edge, the overall hop count will be  $\ell_{n+1} + i + \hat{L} + 1$ , where  $\ell_{n+1}$  is the hop count accumulated when first entering the decision point, while  $i \geq 0$  denotes the number of (self) loops. Instead, if the random walk reaches the destination by eventually following the  $\alpha^*$ -edge, the overall hop count will be  $\ell_{n+1} + i + \hat{L}$ . In any real situation, the normalized cumulative weight  $Z_n$  of  $\beta$ -edges is stochastically dominated by the weight of the virtual best adversary, having normalized weight  $Z'_n$ . We have:

$$\mathbb{E}[Z'_{n+1} | \mathcal{F}_n, \ell_{n+1}] = Z_n \left[ \sum_{i=0}^{\infty} [(1-q)Z_n]^i \left( \frac{\hat{w}(n)}{\hat{w}(n)} \frac{\hat{w}(n) + i\Delta(i)}{\hat{w}(n) + i\Delta(i) + \hat{w}(n) + \Delta(i)} + q \frac{\hat{w}(n) + i\Delta'(i)}{\hat{w}(n) + i\Delta'(i) + \hat{w}(n)} \right) \right] \quad (7)$$

where  $\Delta(i) = f(\ell_{n+1} + i + \hat{L})$  and  $\Delta'(i) = f(\ell_{n+1} + i + \hat{L} + 1)$ .

Now, it turns out that the term in square brackets of the latter expression is smaller than or equal to one for any value of  $\hat{w}(n)$ ,  $\hat{w}(n)$ ,  $\hat{L}$ ,  $\ell_{n+1}$ ,  $q$  and non-increasing function  $f(\cdot)$ . This property can be easily checked numerically, but a formal proof requires some effort (see App. A). As a consequence,  $\mathbb{E}[Z'_{n+1} | \mathcal{F}_n, \ell_{n+1}] \leq Z_n$ . At last, unconditioning with respect to  $\ell_{n+1}$ , whose distribution descends from  $\mathcal{F}_n$ , and considering also the case in which the  $(n+1)$ -th random walk does not reach the decision point, we obtain  $\mathbb{E}[Z'_{n+1} | \mathcal{F}_n] \leq Z_n$  and thus  $\mathbb{E}[Z_{n+1} | \mathcal{F}_n] \leq Z_n$ . Hence, we have that  $Z_n$  converges to a constant  $z \in [0, 1]$ .

To show that necessarily  $z = 0$ , we employ the Poissonization technique as in Section 5.2, noticing again that  $Z_n$  is stochastically dominated by  $Z'_n$ . For the process  $Z'_n$ , we have:

$$\mathbf{E}[\mathbf{A}^T] = \begin{pmatrix} a & b \\ 0 & d \end{pmatrix} \quad (8)$$

The entries in the above matrix have the following meaning:

- $a$  is the average reward given to the  $\alpha^*$ -edge if we select the  $\alpha^*$ -edge;
- $b$  is the average reward given to the  $\alpha^*$ -edge if we select the  $\beta$ -edge;
- $d$  is the average reward given to the  $\beta$ -edge if we select the  $\beta$ -edge;

Note that the average reward given to the  $\beta$ -edge if we select the  $\alpha^*$ -edge is zero.

Luckily, the exponential of a 2x2 matrix in triangular form is well known [2] (see also [12] for limit theorems of triangular Pólya urn schemes). In particular, when  $a \neq d$  we obtain:

$$\begin{pmatrix} \mathbf{E}[\hat{w}(t)] \\ \mathbf{E}[\dot{w}(t)] \end{pmatrix} = \begin{pmatrix} e^{at} & \frac{b}{d-a}(e^{dt} - e^{at}) \\ 0 & e^{dt} \end{pmatrix} \begin{pmatrix} \hat{w}(0) \\ \dot{w}(0) \end{pmatrix} \quad (9)$$

The special case in which  $a = d$  will be considered later (see Section 6.3).

To show that necessarily  $z = 0$ , we reason by contradiction, assuming that  $Z'(n)$  converges to  $z > 0$ . This implies that

$$\dot{w}(t) = \frac{z}{1-z} \hat{w}(t) + o(\hat{w}(t)) \quad (10)$$

Moreover, we will assume that a large enough number of walks has already been performed such that, for all successive walks, the probability to follow the  $\beta$ -edge is essentially equal to  $z$ . Specifically, let  $n^*$  be a large enough time step such that the normalized weight of the  $\beta$ -edge is  $z - \epsilon < Z'(n) < z + \epsilon$  for all  $n > n^*$ . We can then ‘restart’ the system from time  $n^*$ , considering as initial weights  $\dot{w}(n^*)$  and  $\hat{w}(n^*)$  (the specific values are not important).

Taking expectation of (10) and plugging in the expressions of the average weights in (9), we have that the following asymptotic<sup>3</sup> relation must hold:

$$e^{dt} \dot{w}(n^*) \sim_e \frac{z}{1-z} \left( e^{at} \hat{w}(n^*) + \frac{b}{d-a} (e^{dt} - e^{at}) \dot{w}(n^*) \right)$$

Clearly, the above relation does not hold if  $d < a$ . If  $d > a$ , the relation is satisfied when

$$\frac{b}{d-a} = \frac{1-z}{z} \Leftrightarrow \frac{d-a}{d-a+b} = z \quad (11)$$

Interestingly, we will see that (11) is verified when the reward function is constant, suggesting that in this case the  $\beta$ -edge can indeed ‘survive’ the competition with the  $\alpha^*$ -edge. Instead, we will show that  $\frac{d-a}{d-a+b} < z - \epsilon$ , for any strictly decreasing function  $f(\cdot)$ , proving that the normalized weight of the  $\beta$ -edge cannot converge to any  $z > 0$ .

For simplicity, we will consider first the the case in which  $\ell_{n+1}$ , the hop count accumulated by the random walk while first entering the decision point, is not random but deterministic and equal to  $\ell$ . Under the above simplification, we have:

$$a = f(\ell + \hat{L}) \quad (12)$$

$$b = (1-q)(1-z) \sum_{i=0}^{\infty} [z(1-q)]^i f(\ell + i + 1 + \hat{L}) \quad (13)$$

$$d = \sum_{i=0}^{\infty} [z(1-q)]^i \left[ q(i+1)f(\ell + i + 1 + \hat{L}) + (1-q)(1-z)(i+1)f(\ell + i + 1 + \hat{L}) \right] \quad (14)$$

In the special case in which the reward function is constant (let this constant be  $C$ ), we obtain:

$$a = C \quad (15)$$

$$b = C \frac{(1-q)(1-z)}{1-z+qz} \quad (16)$$

$$d = C \frac{1}{1-z+qz} \quad (17)$$

It is of immediate verification that (15),(16),(17) satisfy (11) for any  $q \in [0, 1)$  (the case  $q = 1$  corresponds to having  $a = d$ , which is considered separately in Section 6.3).

<sup>3</sup>Given two functions  $f(n)$  and  $g(n)$ , we write  $f(n) \sim_e g(n)$  if  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$ .

To analyze what happens when  $f(\cdot)$  is a decreasing function, we adopt an iterative approach. We consider a sequence of reward functions  $\{f_k(\cdot)\}_k$ , indexed by  $k = 0, 1, 2, \dots$ , defined as follows. Let  $L = \ell + \hat{L}$  be the minimum path length experience by random walks traversing the decision point. We define:

$$f_k(L+i) = \begin{cases} f(L+i) & \text{if } 0 \leq i \leq k \\ f(L+k) & \text{if } i > k \end{cases} \quad (18)$$

In words, function  $f_k(\cdot)$  matches the actual reward function  $f(\cdot)$  up to hop count  $L+k$ , while it takes a constant value (equal to  $f(L+k)$ ) for larger hop count. See Figure 1.

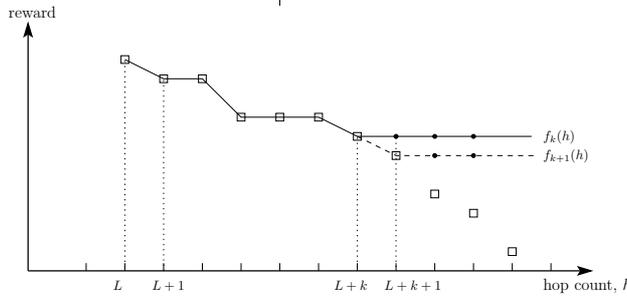


Figure 1: Example of reward functions  $f_k(\cdot)$  and  $f_{k+1}(\cdot)$ . Values taken by the actual reward function  $f(k)$  are denoted by squares. Values taken by function  $f_k(\cdot)$  (function  $f_{k+1}(\cdot)$ ) are connected by solid (dashed) line.

In our proof, we will actually generalize the result in Theorem 1, allowing the reward function to be non-increasing for values larger than  $L$ . To simplify the notation, let  $f(L) = C$ . For  $i = 1, 2, \dots$ , let  $f(L+i) = C - \delta_i$ , with  $\delta_i > 0$ , and  $\delta_i \geq \delta_{i-1}$ .

Let  $a_k, b_k, d_k$  ( $a_{k+1}, b_{k+1}, d_{k+1}$ ) be the entries of matrix (8) when we assume that rewards are given to edges according to function  $f_k(\cdot)$  (function  $f_{k+1}(\cdot)$ ), with  $k \geq 0$ . As a first step, we can show that (11) does not hold already for  $k = 0$ , i.e., for a reward function which is equal to  $C$  for hop count  $h = L$ , and equal to  $C - \delta_1$  for any  $h > L$ . Indeed, in this case we have:

$$\begin{aligned} a_0 &= C \\ b_0 &= (1-q)(1-z) \sum_{i=0}^{\infty} [z(1-q)]^i (C - \delta_1) \\ &= \frac{(1-q)(1-z)}{1-z+zq} (C - \delta_1) \\ d_0 &= (1-z+zq) \sum_{i=0}^{\infty} [z(1-q)]^i (i+1)(C - \delta_1) \\ &= \frac{1}{1-z+zq} (C - \delta_1) \end{aligned}$$

It can be easily check that  $\frac{d_0 - a_0}{d_0 - a_0 + b_0} < z - \epsilon$  for any  $0 < \epsilon < \frac{\delta_1(1-z)(1-z+zq)}{(1-q)C - \delta_1}$ .

To show that (11) cannot hold for the actual reward function  $f(\cdot)$ , it is then sufficient to prove the inductive step

$$\frac{b_k}{d_k - a_k} \leq \frac{b_{k+1}}{d_{k+1} - a_{k+1}}$$

Note, indeed, that the sequence of functions  $\{f_k(\cdot)\}_k$  tends point-wise to  $f(\cdot)$ . Now, for any  $k$  for which  $\delta_{k+1} = \delta_k$  there is nothing to prove, since in this case  $\frac{b_k}{d_k - a_k} = \frac{b_{k+1}}{d_{k+1} - a_{k+1}}$ . So let's suppose that  $\delta_{k+1} > \delta_k$ .

We have  $a_k = a_{k+1} = C$ . We can write  $b_k$  as:

$$\begin{aligned} b_k &= \hat{b} + (1-q)(1-z) \sum_{i=k+1}^{\infty} [z(1-q)]^i (C - \delta_k) \\ &= \hat{b} + (1-q)(1-z)(C - \delta_k) \frac{[z(1-q)]^{k+1}}{1-z+zq} \end{aligned}$$

where

$$\hat{b} = (1-q)(1-z) \sum_{i=0}^k [z(1-q)]^i (C - \delta_{i+1})$$

We can write  $b_{k+1}$  as:

$$\begin{aligned} b_{k+1} &= \hat{b} + (1-q)(1-z) \sum_{i=k+1}^{\infty} [z(1-q)]^i (C - \delta_{k+1}) \\ &= \hat{b} + (1-q)(1-z)(C - \delta_{k+1}) \frac{[z(1-q)]^{k+1}}{1-z+zzq} \end{aligned}$$

Similarly, we have:

$$\begin{aligned} d_k &= \hat{d} + (C - \delta_k)[z(1-q)]^{k+1} \left( k + 1 + \frac{1}{1-z+zzq} \right) \\ d_{k+1} &= \hat{d} + (C - \delta_{k+1})[z(1-q)]^{k+1} \left( k + 1 + \frac{1}{1-z+zzq} \right) \end{aligned}$$

where

$$\hat{d} = \sum_{i=0}^k [z(1-q)]^i (i+1)(1-z+zzq)(C - \delta_{i+1})$$

We will assume that both  $d_k > a_k$  and  $d_{k+1} > a_{k+1}$ , otherwise the result is trivial (if  $d_k < a_k$ , then also  $d_{k+1} < a_{k+1}$ , since  $a_{k+1} = a_k$ ,  $d_{k+1} < d_k$ . If  $d_{k+1} < a_{k+1}$ , the normalized ratio of the  $\beta$ -edge can only tend to zero). Under this assumption, we can show that

$$\frac{b_k}{d_k - a_k} < \frac{b_{k+1}}{d_{k+1} - a_{k+1}} \quad (19)$$

Indeed, plugging in the expressions of  $a_k, b_k, d_k, a_{k+1}, b_{k+1}, d_{k+1}$ , after some algebra we reduce inequality (19) to:

$$\hat{b}[(k+1)(1-z+zzq) + 1] + C(1-q)(1-z) > \hat{d}(1-q)(1-z)$$

At last, recalling the definitions of  $\hat{b}$  and  $\hat{d}$ , we obtain that the above inequality is satisfied if

$$\begin{aligned} &\sum_{i=0}^k [z(1-q)]^i [(k+1)(1-z+zzq) + 1](C - \delta_{i+1}) > \\ &\sum_{i=0}^k [z(1-q)]^i (i+1)(1-z+zzq)(C - \delta_{i+1}) \end{aligned}$$

which is clearly true, since  $k \geq i$  when  $i$  varies from 0 to  $k$ .

We now provide a sketch of the proof for the case in which the random walk arrives at the decision point having accumulated a random hop count  $\ell_{n+1}$ . After long enough time, we can assume that the probability distribution of  $\ell_{n+1}$  has converged to a random but fixed distribution that no longer depends on  $n$ . Indeed, such distribution depends only on normalized edge weights, which in the long run converge to constant values. Let  $p_m = \mathbb{P}\{\ell_{n+1} = \ell_{\min} + m\}$ ,  $m \geq 0$ , where  $\ell_{\min}$  is the minimum hop count that can be accumulated at the decision point. We can use  $\{p_m\}_m$  to compute expected values of  $a, b, d$  as defined in (12), (13), (14), and apply again the Poissonization technique to compute asymptotic values of edge weights.

Specifically, letting  $C = f(\ell_{\min} + \hat{L})$ , we obtain:

$$\begin{aligned} \mathbb{E}[a] &= \sum_{m=0}^{\infty} p_m (C - \delta_m) \\ \mathbb{E}[b] &= (1-q)(1-z) \sum_{m=0}^{\infty} \sum_{i=0}^{\infty} [z(1-q)]^i (C - \delta_{m+i+1}) \\ \mathbb{E}[d] &= (1-z+zzq) \sum_{m=0}^{\infty} \sum_{i=0}^{\infty} [z(1-q)]^i (i+1)(C - \delta_{m+i+1}) \end{aligned}$$

Similarly to before, we prove by contradiction that (11) cannot hold, through an iterative approach based on the sequence of reward functions  $\{f_k(\cdot)\}_k$ . As basic step of the induction, we take the reward function  $f_0(\cdot)$  equal to  $C - \delta_1$  for any hop count larger than  $\ell_{\min} + \hat{L}$ . Hence, we have  $\delta_{m+i+1} = \delta_1$ ,  $\forall m, i \geq 0$ . It follows that  $E[b]$  is exactly the same as in (13), and  $E[d]$  is exactly the same as in (14). The only quantity that is different is  $\mathbb{E}[a] = p_0 C + (1 - p_0)(C - \delta_1) = C - \delta_0$ , where  $\delta_0 < \delta_1$  as long as  $p_0 > 0$ . Therefore, whenever there is a non-null probability  $p_0$  to reach the decision point with minimum hop count, the basic induction step proven before still holds here, by redefining  $C$  and  $\delta_1$  as  $C - \delta_0$  and  $\delta_1 - \delta_0$ , respectively. One can also prove that the generic iterative step still holds, by following the same lines as in the basic case. Indeed, one can verify that

$$\frac{\mathbb{E}[b_k]}{\mathbb{E}[d_k] - \mathbb{E}[a_k]} \leq \frac{\mathbb{E}[b_{k+1}]}{\mathbb{E}[d_{k+1}] - \mathbb{E}[a_{k+1}]}$$

when  $\mathbb{E}[d_{k+1}] > \mathbb{E}[a_{k+1}]$ . This concludes the proof of Lemma 2.  $\square$

*Proof of Theorem 1 (general case).* The proof is exactly the same as in the DAG case, with the difference that we employ Lemma 2 instead of Lemma 1 to iteratively prune  $\beta$ -edges from the decision points, leaving only paths from  $s$  to  $d$  of minimum length.  $\square$

### 6.3 The single reward model in general network

We conclude the asymptotic analysis considering the single-reward model in a general directed network. Given the analysis for the multiple-reward model, the single-reward model is almost immediate. Indeed, the expressions for  $a$  and  $b$  (respectively in (12) and (13)) are left unmodified, as well as their averages  $\mathbb{E}[a]$  and  $\mathbb{E}[b]$  with respect to hop count accumulated at the decision point. Instead, we have

$$\mathbb{E}[d] = (1 - z + zq) \sum_{m=0}^{\infty} \sum_{i=0}^{\infty} [z(1 - q)]^i (C - \delta_{m+i+1}) \quad (20)$$

which is clearly smaller than the  $\mathbb{E}[d]$  obtained under the multiple-reward model. Hence, the basic step of the induction used to prove Lemma 2 follows immediately from the consideration that  $\frac{\mathbb{E}[d] - \mathbb{E}[a]}{\mathbb{E}[d] - \mathbb{E}[a] + \mathbb{E}[b]}$  is an increasing function of  $\mathbb{E}[d]$ . Moreover, simple algebra shows that the iterative step holds also in the case of single-reward, allowing us to extend the validity of Lemma 2, and thus Theorem 1.

Last, it is interesting to consider the case of single reward model and constant reward function,  $f(\cdot) = C$ . We have in this case:

$$a = C \quad (21)$$

$$b = C \frac{(1 - q)(1 - z)}{1 - z + qz} \quad (22)$$

$$d = C \quad (23)$$

Since  $a = d$ , the matrix exponential takes a different form with respect to 9, that now reads:

$$\begin{pmatrix} \mathbf{E}[\hat{w}(t)] \\ \mathbf{E}[\dot{w}(t)] \end{pmatrix} = e^{at} \begin{pmatrix} 1 & b \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{w}(0) \\ \dot{w}(0) \end{pmatrix} \quad (24)$$

We can show by contradiction that the normalized weight of the  $\beta$ -edge cannot tend to any  $z > 0$ . Indeed, assuming to restart the system after a long enough number of walks  $n^*$  such that  $\hat{w}(n^*) \approx \frac{1-z}{z} \dot{w}(n^*)$ , we should have:

$$e^{at} \dot{w}(n^*) \sim_e \frac{z}{1-z} (e^{at} \hat{w}(n^*) + e^{at} b \dot{w}(n^*))$$

which can only be satisfied if  $b = 0$ . Interestingly,  $b$  equals 0 when  $q = 1$ , i.e., when the  $\beta$ -edge becomes a  $\beta^*$ -edge. This means that, asymptotically, the probability that the random walk makes any loop must vanish to zero. We conclude that, in the case of a constant single reward model, many paths can survive (including non-shortest paths), but not those containing loops. In other words, surviving edges must belong to a DAG. Simulation results, omitted here due to lack of space, confirm this prediction.

## 7 Transient analysis

Beyond the asymptotic behavior, it is interesting to consider the evolution of edge weights over time. In particular, since all non-shortest paths are taken with vanishing probability, what law governs the decay rate of such probabilities? How does the decay rate depend on system parameters, such as network topology and reward function? Such questions are directly routed to analogous questions regarding how normalized edge weights evolve over time, as the probability of taking a given path is simply the product of the probabilities of taking its edges. Thus, we investigate the transient behavior of normalized edge weights.

### 7.1 Single decision point

We again start by considering the case of a single decision point with two outgoing edges (edge 1 and edge 2), whose initial weights are denoted by  $w_1[0]$  and  $w_2[0]$ , respectively. Let  $\Delta_1 = f(L_1)$  and  $\Delta_2 = f(L_2)$  be the rewards associated to edge 1 and edge 2, and  $L_{1,2}$  the corresponding path lengths.

As discussed in Section 5.2, the dynamics of this discrete time system can be usefully embedded into continuous time using the Poissonization technique, which immediately provides the transient behavior of the system in the simple form (4). To complete the analysis, the solution in continuous time  $t$  should be transformed back into discrete time  $n$ . Unfortunately, this operation can be done exactly only in the trivial case of just one edge. With two (or more) edges, we can resort to an approximate (yet quite accurate) heuristic called depoissonization, which can be applied to all Pólya urn models governed by invertible ball addition matrices [15]. In this simple topology, assuming  $\Delta_1 > \Delta_2$ , the approximation consists in assuming that all extractions that have occurred by time  $t$  are associated to the winning edge only (this becomes more and more true with the passing of time), which permits deriving the following approximate relation between  $n$  and  $\bar{t}_n$ , where  $\bar{t}_n$  is the average time at which the  $n$ -th ball is drawn:

$$n \approx \frac{w_1[0]}{\Delta_1} e^{\Delta_1 \bar{t}_n} \quad (25)$$

from which one obtains  $\bar{t}_n \approx \left(\log \frac{n\Delta_1}{w_1[0]}\right) / \Delta_1$ . Using this approximate value of  $\bar{t}_n$  into (4), we can approximate the expected values of edge weights after  $n$  walks as:

$$\begin{pmatrix} \mathbf{E}[w_1[n]] \\ \mathbf{E}[w_2[n]] \end{pmatrix} \approx \begin{pmatrix} e^{\Delta_1 \bar{t}_n} & 0 \\ 0 & e^{\Delta_2 \bar{t}_n} \end{pmatrix} \begin{pmatrix} w_1[0] \\ w_2[0] \end{pmatrix} = \begin{pmatrix} \Delta_1 n \\ w_2[0] \left(\frac{n\Delta_1}{w_1[0]}\right)^{\frac{\Delta_2}{\Delta_1}} \end{pmatrix} \quad (26)$$

The above approximation is not quite accurate for small values of  $n$ . In particular, the normalized weight of edge 2, according to (26), can be even larger than the initial value  $\frac{w_2[0]}{w_1[0]+w_2[0]}$ . For this reason, for small values of  $n$ , we improve the approximation by assuming that the (average) normalized weight of edge 2 cannot exceed its initial value at time 0. Indeed, we can easily find analytically the maximum value of  $n$ , denoted by  $n^*$ , for which we bound the normalized weight of edge 2 to the value  $\frac{w_2[0]}{w_1[0]+w_2[0]}$ . It turns out that  $n^* = \frac{w_1[0]}{\Delta_1}$ . Note that  $n^*$  depends solely on parameters of the first edge.

Our final approximation for the (average) normalized weight of edge 2 is then:

$$\mathbf{E} \left[ \frac{w_2[n]}{w_1[n] + w_2[n]} \right] \approx \begin{cases} \frac{w_2[0]}{w_1[0] + w_2[0]} & \text{if } n \leq n^* \\ \frac{1}{1 + \frac{w_1[0]}{w_2[0]} \left(\frac{n\Delta_1}{w_1[0]}\right)^{1 - \frac{\Delta_2}{\Delta_1}}} & \text{if } n > n^* \end{cases} \quad (27)$$

The expression for the (average) normalized weight of edge 1 is then easily derived as the complement of the above.

The value of  $n^*$  can be used to separate the transient regime into two parts: we call the first one, for  $n \leq n^*$ , the *exploration* phase, because during this initial interval there is still no clear winner between the competing edges, and random walks explore all possibilities with lots of variability in the selected edges. Instead, we call the second one, for  $n > n^*$ , the *convergence* phase, where the winning edge starts

to emerge and dominate the competition, whereas the loosing edge inexorably decays. The behavior of this phase is much more deterministic than the initial one, especially because at this point edges have accumulated quite a lot of weight, which individual random walks cannot significantly modify from one walk to another. These two phases and decays are illustrated numerically in Section 8 (Figure 5).

Interestingly, from (27) we see that the probability to select edge 2 decays asymptotically to zero (as  $n \rightarrow \infty$ ) according to the power law  $n^{\frac{\Delta_2}{\Delta_1}-1}$ . In particular, the larger the ratio between  $\Delta_1$  and  $\Delta_2$ , the faster the decay, which cannot however be faster than  $n^{-1}$ .

## 7.2 General network: recursive method

We propose two different approaches to extend the transient analysis to a general network. Our goal is to approximate the evolution of the average weight  $\mathbb{E}[w_{i,j}[n]]$  of individual edges over time  $n$  (where the average is with respect to all sample paths of the system).

The first approach is computationally more expensive but conceptually simple and surprisingly accurate in all scenarios that we have tested (see Section 8). It is based on the simple idea of making a step-by-step, recursive approximation of  $\mathbb{E}[w_{i,j}[n]]$  by just taking the average of (1):

$$\mathbb{E}[w_{i,j}[n]] = \mathbb{E}[w_{i,j}[n-1]] + \mathbb{E}[\Delta_{i,j}[n]] \quad (28)$$

where the approximation lies in the computation of  $\mathbb{E}[\Delta_{i,j}[n]]$ , which is the expected reward given to edge  $(i, j)$  after executing the  $n$ -th walk. This quantity can be (approximately) evaluated using just the set of values  $\{\mathbb{E}[w_{i,j}[n-1]]\}_{i,j}$  obtained at step  $(n-1)$ .

Indeed, note that  $\mathbb{E}[\Delta_{i,j}[n]]$  requires to compute the distribution of the lengths of paths from  $s$  to  $d$  containing edge  $(i, j)$ . Note that we do not need the complete enumeration of these paths, but just the distribution of their length. For this, standard techniques of Markov Chain analysis can dramatically reduce the computational burden, as we will see. The fundamental approximation that we make while evaluating this distribution is the following. First observe that the (averaged) probability to follow a given path at time  $n$  is exactly the product of (averaged) independent probabilities to select individual edges. Unfortunately, the probability to select any given edge corresponds to its (averaged) *normalized* weight at time  $n-1$ :

$$\mathbb{E}[r_{i,j}[n]] = \mathbb{E} \left[ \frac{w_{i,j}[n-1]}{\sum_k w_{i,k}[n-1]} \right]$$

which cannot be evaluated exactly, since we do not know the (joint) probability density function of weights. So we approximate  $\mathbb{E}[r_{i,j}[n]]$  by the ratio of averages:

$$\mathbb{E}[r_{i,j}[n]] \approx \frac{\mathbb{E}[w_{i,j}[n-1]]}{\mathbb{E}[\sum_k w_{i,k}[n-1]]}$$

which is instead completely known if we have values  $\{\mathbb{E}[w_{i,j}[n-1]]\}_{i,j}$ . In essence, this approximation consists in using the ratio of expectations as the expectation of a ratio.

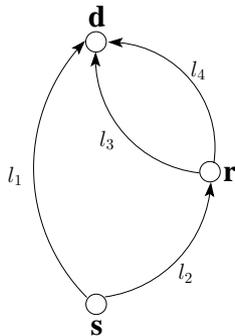


Figure 2: Example of topology comprising two decision points

In order to illustrate this recursive approach, consider the topology in Fig. 2, comprising two decision points: the source node  $s$  and the relay node  $r$ . The arcs shown in Fig. 2 do not represent individual edges but paths (i.e., sequence of nodes) with lengths denoted by  $l_i$ ,  $i = 1, \dots, 4$ , in hops. We also denote by  $w_i$ , with some abuse of notation, the weight associated to the first edge of the corresponding path  $l_i$ .

The approximate transient analysis of this graph is obtained by the following set of recursive equations:

$$\begin{cases} \mathbb{E}[w_1[n]] &= \mathbb{E}[w_1[n-1]] + \mathbb{E}[r_1[n]]\Delta_1 \\ \mathbb{E}[w_2[n]] &= \mathbb{E}[w_2[n-1]] + \mathbb{E}[r_2[n]]\Delta_2[n] \\ \mathbb{E}[w_3[n]] &= \mathbb{E}[w_3[n-1]] + \mathbb{E}[r_2[n]]\mathbb{E}[r_3[n]]\Delta_3 \\ \mathbb{E}[w_4[n]] &= \mathbb{E}[w_4[n-1]] + \mathbb{E}[r_2[n]]\mathbb{E}[r_4[n]]\Delta_4 \end{cases} \quad (29)$$

where  $\Delta_1 = f(l_1)$ ,  $\Delta_2[n] = \mathbb{E}[r_3[n]]f(l_2 + l_3) + \mathbb{E}[r_4[n]]f(l_2 + l_4)$ ,  $\Delta_3 = f(l_2 + l_3)$ ,  $\Delta_4 = f(l_2 + l_4)$ . In the above equations we have denoted the (approximated) normalized weights as  $\mathbb{E}[r_i[n]]$ . For example,  $\mathbb{E}[r_1[n]] \approx \frac{\mathbb{E}[w_1[n-1]]}{\mathbb{E}[w_1[n-1]] + \mathbb{E}[w_2[n-1]]}$ , and similarly for the other values  $\mathbb{E}[r_i[n]]$

The recursive approach can be applied to an arbitrary graph, but in general it requires to recompute, at each time  $n$  (in the case of the single reward model): i) the distribution of path lengths from  $s$  to  $j$  passing through edge  $(i, j)$ , which is an outgoing edge of decision point  $i$ ; ii) the distribution of path lengths from  $j$  to the destination  $d$ . The above distributions can be computed numerically by solving the transient of discrete-time Markov chains with proper absorbing states, but we do not provide the details here. Since the overall procedure can be computationally quite expensive in large graphs, we present in the next section a different, much simpler approach which captures the asymptotic law by which average edge weights decay.

### 7.3 General network: asymptotic power-law decay

The asymptotic analysis in Section 6 shows that the normalized weight of all  $\beta$ -edges (or  $\beta^*$ -edges) vanishes to zero as  $n \rightarrow \infty$ . Can we analytically predict the asymptotic law for such decay? The answer is affirmative, and the results offer fundamental insights into how the network structure evolves over time.

We start defining a key concept associated to decision points.

**Definition 4 (clock of a decision point).** *The clock  $c_i[n]$  of a decision point  $i$  is the expected number of random walks that reach  $i$  by time  $n$ :*

$$c_i[n] := \sum_{j=1}^n \mathbb{P}\{\text{random walk } j \text{ hits } i\}$$

The clock of a decision point dictates how fast the dynamics of its outgoing edges evolve with respect to the reference time  $n$ . As a corollary of Theorem 1, the clock of all decision points traversed by at least one shortest path is  $\Theta(n)$ , since any shortest path is asymptotically used with non-zero probability. However, decision points not traversed by shortest paths have clock  $o(n)$ , and if we put them in sequence we get decision points with increasingly slower clocks. Nevertheless, we can show that the clock of any decision point is  $\omega(1)$ .

Consider, for example, the simple topology in Fig. 2, and suppose that  $l_1$  is the only shortest path and that  $l_3 < l_4$ . Since  $l_4$  will be asymptotically used a vanishing fraction of times as compared to  $l_3$  (restricting our attention to the set of random walks passing through  $r$ , i.e., the clock of  $r$ ), we can asymptotically consider decision point  $s$  as the sole decision point of the network with two outgoing paths of lengths  $l_1$  and  $l_2 + l_3$ . Hence, we can just apply (27) to compute the power law decay of the  $\beta$ -edge leading to the path with length  $l_2 + l_3$ :

$$\mathbf{E} \left[ \frac{w_2[n]}{w_1[n] + w_2[n]} \right] = \Theta(n^{\frac{\Delta_2}{\Delta_1} - 1}) \quad (30)$$

where  $\Delta_1 = f(l_1)$  and  $\Delta_2 = f(l_2 + l_3)$ .

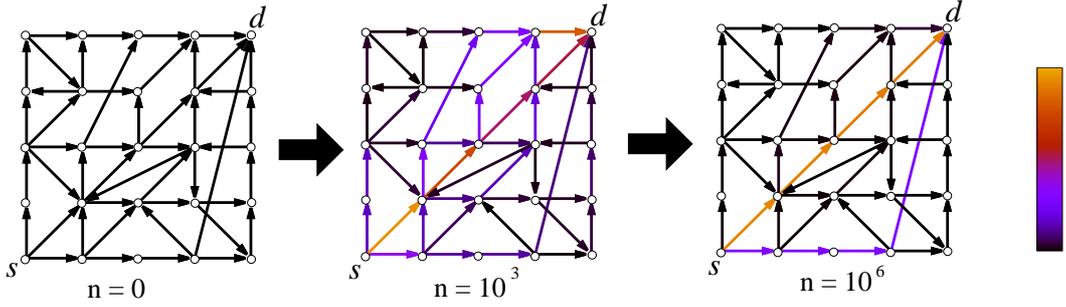


Figure 3: From random walks to short walks: (a) At  $n = 0$ , all weights are identical; (b) at  $n = 10^3$  an edge weight structure starts to emerge along shorter paths; (c) at  $n = 10^6$  edge weights along the (two) shortest paths are dominant.

Moreover, we can again use (27) to compute the scaling order of the (average) clock of decision point  $r$ :

$$\begin{aligned} \mathbb{E}[c_r[n]] &= \mathbb{E} \left[ \sum_{j=1}^n \frac{w_2[j]}{w_1[j] + w_2[j]} \right] = \sum_{j=1}^n \mathbb{E} \left[ \frac{w_2[j]}{w_1[j] + w_2[j]} \right] \\ &\approx \sum_{j=1}^n \frac{1}{1 + \frac{w_1[0]}{w_2[0]} \left( \frac{j\Delta_1}{w_1[0]} \right)^{1 - \frac{\Delta_2}{\Delta_1}}} = \Theta \left( \int_0^n x^{\frac{\Delta_2}{\Delta_1} - 1} dx \right) = \Theta \left( n^{\frac{\Delta_2}{\Delta_1}} \right) \quad (31) \end{aligned}$$

Note that the clock of  $r$  is both  $o(n)$  and  $\omega(1)$ . At last, we can compute the power law decay of the first  $\beta$ -edge of path  $l_4$ , by applying again (27) to decision point  $r$ , with the caveat of plugging in the clock of  $r$  in place of  $n$ :

$$\mathbf{E} \left[ \frac{w_4[n]}{w_3[n] + w_4[n]} \right] = \Theta \left( \left( n^{\frac{\Delta_2}{\Delta_1}} \right)^{\frac{\Delta_4}{\Delta_3} - 1} \right) = \Theta \left( n^{\frac{\Delta_2}{\Delta_1} \left( \frac{\Delta_4}{\Delta_3} - 1 \right)} \right)$$

where  $\Delta_3 = f(l_2 + l_3)$  and  $\Delta_4 = f(l_2 + l_4)$ .

A simple algorithm, that we omit here, can recursively compute the clock of all decision points (in scaling order), and the power law decay exponent of all decaying edges, starting from the source and moving towards the destination. We will discuss the implications of our results in a significant example presented later in Section 8.4.

## 8 Validation and insights

We present a selection of interesting scenarios explored numerically through simulations to confirm our approximate transient analysis and offer insights into the system behavior.

### 8.1 Emergence of shortest paths

We start by considering a 25-nodes network containing a few loops, evolving under the multiple reward model. Nodes are arranged in a 5x5 grid, with the source located at the bottom left corner and the destination at the top right corner. The initial weight on any edge is 1, and the reward function is  $f(L) = 1/L$ . Figure 3 (in color, better seen on screen) shows three snapshots of one system run, at times  $n = 0$ ,  $n = 10^3$ ,  $n = 10^6$ , where magnitude of edge weights is converted into a color code according to a heat-like palette. We observe that, by time  $n = 10^6$ , edge weights along the two shortest paths are dominant. Note that one shortest path (along the diagonal) appears to be stronger than the other (i.e., more likely to be used) but this changes from one run to another, since the asymptotic probability to use a specific shortest path is a random variable (recall Theorem 1).

## 8.2 Non-monotonous behavior of random walks

Interestingly, although edge weights increase monotonically, normalized edge weights (which are the quantities actually steering the random walk through the network) can exhibit non-monotonous behavior, even when we consider their expected values (across system runs). We illustrate this on the simple topology of Fig. 2, using segment lengths  $l_1 = 7$ ,  $l_2 = l_3 = 3$ ,  $l_4 = 18$ . Fig. 4 shows the transient of the normalized weights of the four outgoing edges, comparing simulation results (obtained averaging 1,000 runs) and the analytical approximation based on the recursive approach (29). Here initial weights are equal to 1,  $f(L) = L^{-2}$ .

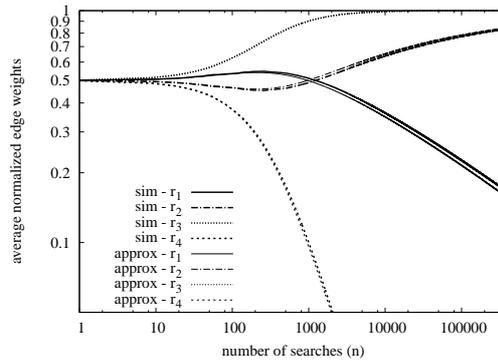


Figure 4: Transient behavior of normalized weights in a simple topology with two decision points. Comparison between simulation and analytical approximation based on the recursive approach.

Besides showing the surprising accuracy of the recursive approximation, the plot in Fig. 4 confirms that normalized edge weights can be non-monotonous (see  $r_1$  or  $r_2$ ). Note that this might appear to contradict a fundamental result that we have obtained while proving Lemma 1 (or Lemma 2), namely, the fact that the (average) normalized weight of a  $\beta$ -edge competing against an  $\alpha^*$ -edge is non-increasing. However, this result cannot be applied to the first decision point (the source  $s$ ) since the assumptions of Lemma 1 do not hold here (there are no  $\alpha^*$ -edges going out of  $s$ ).

## 8.3 Trade-off between exploration and convergence

What happens when we change the reward function  $f(L)$ ? How is the transient of a network affected by taking a reward function that decreases faster or slower with the hop count? We investigate this issue in the case of a single decision point, considering the family of reward functions  $f(L) = L^{-\phi}$  where we vary the exponent  $\phi > 0$ .

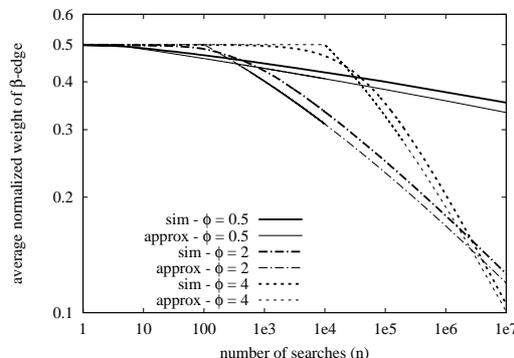


Figure 5: Transient behavior of a single decision between two paths of length 10 and 11, update function  $L^{-\phi}$ , initial weights 1.

We consider a simple topology in which the source is connected to the destination by two edge-independent paths of length 10 and 11. The first edges of these paths have weights  $w_1[n]$  and  $w_2[n]$ ,

respectively, with initial value 1. Figure 5 shows the transient behavior of the average normalized weight  $\mathbb{E}[w_2[n]/(w_1[n] + w_2[n])]$ , for three different values of  $\phi = 0.5, 2, 4$ , comparing simulation results (averaging 1,000 runs) with our analytical approximation (27).

We observe an interesting trade-off between exploration and convergence. Note the role of the threshold  $n^*$  introduced in Section 7.1, here equal to  $n^* = \frac{w_1[0]}{\Delta_1} = 10^\phi$ . Thus, the duration of the exploration phase grows exponentially with  $\phi$ . Moreover, the exponent of the asymptotic power law decay (30) is equal to  $\Delta_2/\Delta_1 - 1 = (10/11)^\phi - 1$ . Thus, larger  $\phi$  leads to larger exponent and thus faster convergence. Therefore, as  $\phi$  increases (corresponding to a reward function that decreases much more rapidly with the hop count), convergence is asymptotically faster but exploration requires much more time. Intuitively, reward functions that decay too fast with path lengths require many repetitions of the WRW before some initial structure emerges. However, when a structure does emerge, they will quickly drive the WRW deeper and deeper into it.

## 8.4 Slowing-down clocks

Consider the network illustrated in Figure 6, comprising of an long sequence of decision points indexed by  $1, 2, \dots$ . Each decision point is directly connected to the destination  $d$  and to the next decision point in the sequence. The source coincides with node 1.

This scenario provides interesting insights into the impact of slowing-down clocks on the evolution of the system structure, and will also illustrate the calculation of the asymptotic power-law decay of  $\beta$ -edges.

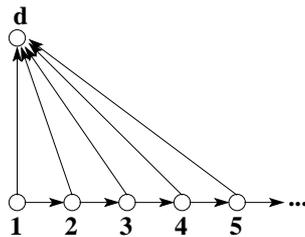


Figure 6: Network with multiple decision points in sequence.

Note that, asymptotically, random walks reaching decision point  $i$  will end up going directly to  $d$  instead of the next decision point. This means that, asymptotically, each decision point can be studied in isolation, considering two outgoing edges: an  $\alpha^*$ -edge belonging to a path of length  $i$ , and a  $\beta^*$ -edge belonging to a path of length  $i + 1$ .

Following this reasoning, we can iteratively compute the power law decay of all  $\beta^*$ -edges, along with the scaling law for the clocks of the respective decision points, using the formulas introduced in Section 7.3. Let  $e_\beta^i$  be the scaling exponent of the outgoing  $\beta^*$ -edge of node  $i$ , and  $e_c^i$  be the scaling exponent of node  $i$ 's clock. Using (30) we have  $e_\beta^1 = \frac{f(2)}{f(1)} - 1$ , and from (31) we get  $e_c^2 = \frac{f(2)}{f(1)}$ . Subsequently, we can derive  $e_\beta^2 = \frac{f(2)}{f(1)} \left( \frac{f(3)}{f(2)} - 1 \right) = \frac{f(3) - f(2)}{f(1)}$  and from this obtain  $e_c^3 = 1 + \frac{f(3) - f(2)}{f(1)}$ . We then have  $e_\beta^3 = \frac{(f(1) + f(3) - f(2))(f(4) - f(3))}{f(1)f(3)}$ , and so on<sup>4</sup>.

In Figures 7 and 8 we compare simulation and analytical results for the first four decision points shown in Fig. 6, considering initial weights equal to 1,  $f(L) = L^{-1}$ . Analytical predictions for the power-law exponents are represented by segments placed above the corresponding simulation curve (note the log-log scale). Besides showing the accuracy of the analytical prediction, results in Fig. 7 and 8 illustrates an important fact: the network structure (i.e., weights on edges) is left essentially unmodified as we move away from the shortest path. Note that this is not quite evident from the math, which predicts that the clock of any decision point in the sequence diverges. In practice, clocks of decision points sufficiently far from the shortest paths evolves so slowly that we can essentially ignore the perturbations caused by

<sup>4</sup>We lack a general closed-form expression for  $e_\beta^i$  or  $e_c^i$ .

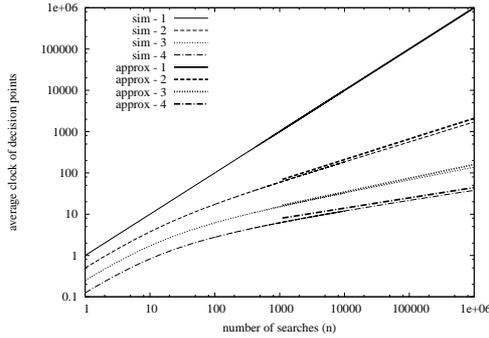


Figure 7: Clocks associated to the first four decision points of the network in Fig. 6.

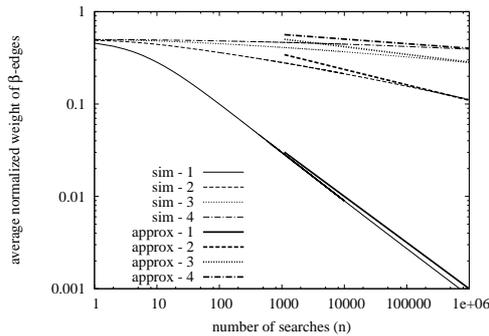


Figure 8: Normalized weight of  $\beta$ -edges going out of the first four decision points of the network in Fig. 6.

the random walks. Hence, sufficiently far regions of a large network preserve their initial ‘plasticity’, allowing them to be used for other purposes.

## 8.5 Transient analysis of the complete graph

As a final interesting case, we consider the complete graph with  $m$  nodes, where the shortest path has length 1 and every node has cycles of all lengths. Without lack of generality, let the source and destination correspond to nodes 1 and  $m$ , respectively.

The asymptotic decay exponent of  $\beta$ -edges can be easily computed following the approach in Section 7.3. Due to symmetry, there are essentially two types of decision point to analyze: the source node (node 1), and any other node different from the source and the destination, for example, node 2.

Node 1 will have, asymptotically, one surviving edge traversed by the unique shortest path of length 1, and  $m - 2$  decaying  $\beta$ -edges traversed by paths whose average length will converge to 2. Hence, the decay exponent of any  $\beta$ -edge going out of  $s$ , such as edge  $(1, 2)$ , is  $\frac{f(2)}{f(1)} - 1$ . The clock of decision point 2 will then run with scaling exponent  $\frac{f(2)}{f(1)}$ .

Node 2 will have, asymptotically, one surviving edge traversed by paths of length 2, and  $(m - 2)$  decaying edges traversed by paths of average length tending to 3. Each of them will decay with power law exponent  $\frac{f(3) - f(2)}{f(1)}$ .

For the complete graph, we have also run the recursive method introduced in Section 7.2, which required us to numerically solve, at each time step, the transient behavior of different discrete-time Markov chains with structure similar to that of the complete graph, modified by the introduction of proper absorbing states to obtain the path length distributions needed by the recursive formulas.

The recursive approach provides a more detailed prediction of the system behavior (at the cost of

higher computational complexity). In particular, it allows to distinguish, among the  $(m - 2)$   $\beta$ -edges going out of decision point 2, the special case of edge  $(2, 1)$ . Intuitively, such edge will accumulate more weight than the  $\beta$ -edge connecting node 2 to, say, node 3, because node 1 is a different decision point with respect to all other decision points (in particular, it has the smallest average residual path length to reach the destination). Considering the completely symmetric structure of the rest of the graph, it turns out that there are, essentially, 5 types of edges having different transient behavior<sup>5</sup>: 1) the  $\alpha$ -edge  $(1, m)$ ; 2) the  $\beta$ -edge  $(1, 2)$ ; 3) the  $\alpha$ -edge  $(2, m)$ ; 4) the  $\beta$ -edge  $(2, 1)$ ; 5) the  $\beta$ -edge  $(2, 3)$ ;

The results obtained by the recursive method are compared against simulations in Figure 9, for  $m = 50$  nodes, initial weights equal to 1,  $f(L) = L^{-1}$ , and single reward model. Besides confirming the surprising accuracy of the recursive approximation, results in Figure 9 suggest that, except for the outgoing edges of the source node, all other edges are marginally affected by the reinforcement process. Indeed, there are so many edges in this network (i.e., available structure) that the ‘perturbation’ necessary to discover and consolidate the shortest path between two particular nodes practically does not significantly affect any edge which is not directly connected to the source.

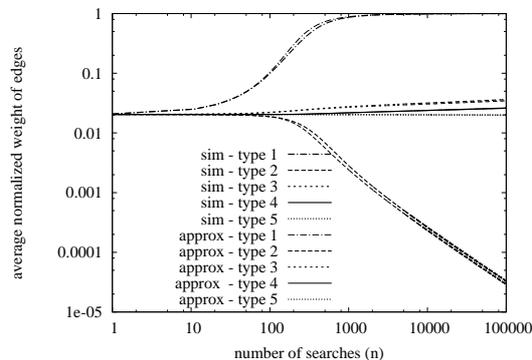


Figure 9: Transient behavior of the complete graph with 50 nodes, comparing simulation and recursive approximation, for 5 different edge types.

## 9 Conclusion and future work

Focusing on the important problem of network navigation, we have introduced and analyzed a novel, simple model capturing the co-evolution of network structure and function performance. We have shown how the repetition of a simple network activity process (WRW with edge reinforcement) is able to build over time a network structure that always leads itself to navigate through shortest paths, in a surprisingly robust manner. Many variations and extensions of the proposed model are possible, which could shed light on how information is efficiently found and/or stored in biological systems lacking the computational and storage resources required to run sophisticated routing algorithms.

## Appendix A Complement to the proof of Lemma 2

The algebraic property related to (7) that we need to prove holds even if the reward function is constant. Therefore, we will prove it under the assumption that  $\Delta(i) = C, \forall i$ . The proof can be extended to a general non-increasing reward function using the sequence of functions  $\{f_k(\cdot)\}_k$  and the iterative approach introduced in Section 6.2, but we omit this extension here. Moreover, we will first consider the simpler case in which  $g = 0$ . Introducing the following normalized variables:  $\alpha = \frac{\hat{w}}{w}$  and  $\sigma = \frac{C}{w}$ , we essentially need to show that

$$\frac{1}{\alpha} \sum_{i=0}^{\infty} \left( \frac{\alpha}{\alpha + 1} \right)^i \left[ \frac{\alpha + i\sigma}{\alpha + 1 + (i + 1)\sigma} \right] \leq 1$$

<sup>5</sup>to avoid complex notation, we take decision point 2 as representative of any decision point different from 1, and decision point 3 as representative of any decision point other than 1 and 2.

for any  $\alpha > 0$  and  $\sigma > 0$ . Observe that the above expression is exactly equal to 1 for  $\sigma = 0$ . We make another change of variable, introducing  $x = \frac{\alpha+1}{\sigma}$ . After some simple algebra, our target reduces to show that:

$$g(\alpha, x) = \frac{x + \alpha + 1}{\alpha} \sum_{k=x+1}^{\infty} \left( \frac{\alpha}{\alpha + 1} \right)^{k-x} \frac{1}{k} \geq 1$$

for any  $\alpha > 0$  and  $x > 0$ . We know that,  $\forall \alpha > 0$ ,  $\lim_{x \rightarrow \infty} g(\alpha, x) = 1$ . Therefore, for arbitrarily small  $\epsilon > 0$ , there exists an  $x_\epsilon$  such that, for  $x > x_\epsilon$ ,  $g(\alpha, x) \geq 1 - \epsilon$ .

We can show that, if  $g(\alpha, x) \geq 1 - \epsilon$ , then  $g(\alpha, x - 1) > 1 - \epsilon$ , for any  $\alpha > 0$  and  $\epsilon \geq 0$ . Indeed, we have:

$$\begin{aligned} g(\alpha, x - 1) &= \frac{x + \alpha}{\alpha} \sum_{k=x}^{\infty} \left( \frac{\alpha}{\alpha + 1} \right)^{k-x+1} \frac{1}{k} \\ &= \frac{x + \alpha}{\alpha} \frac{\alpha}{\alpha + 1} \sum_{k=x}^{\infty} \left( \frac{\alpha}{\alpha + 1} \right)^{k-x} \frac{1}{k} \\ &= \frac{x + \alpha}{\alpha + 1} \left[ g(\alpha, x) \frac{\alpha}{x + \alpha + 1} + \frac{1}{x} \right] \\ &\geq (1 - \epsilon) \frac{x + \alpha}{x + \alpha + 1} \frac{\alpha}{\alpha + 1} + \frac{1}{x} \frac{x + \alpha}{\alpha + 1} > 1 - \epsilon \end{aligned}$$

as can be easily checked. Note that, by recursion, if  $g(\alpha, x) \geq 1 - \epsilon$ , then  $g(\alpha, x - m) > 1 - \epsilon$  for any  $m$  such that  $x - m > 0$ . Armed with this result, we can easily prove that  $g(\alpha, x) > 1$  for any  $x > 0$ . Indeed, suppose, by contradiction, that  $g(\alpha, x_\sigma) < 1$  at a given point  $x_\sigma$ . Then we can write  $g(\alpha, x_\sigma) = 1 - 2\epsilon$ . Now, we build a sequence of values  $\{x_\sigma + 1, x_\sigma + 2, \dots, x_\sigma + m, \dots\}$  which eventually enters the stripe  $[1 - \epsilon, 1 + \epsilon]$ , since  $\lim_{x \rightarrow \infty} g(\alpha, x) = 1$ . Therefore there exists a sufficiently large  $m \geq 1$  such that  $g(\alpha, x_\sigma + m) > 1 - \epsilon$ . But then it must be  $g(\alpha, x_\sigma) > 1 - \epsilon$ , which contradicts the hypothesis that  $g(\alpha, x_\sigma) = 1 - 2\epsilon$ . It remains to prove that  $g(\alpha, x)$  cannot be identically equal to 1 at all points  $x > 0$ . Again, this can be proven by contradiction: suppose that  $g(\alpha, x) = 1$  at any  $x$ . Considering a generic point  $x_\sigma + 1$ ,  $g(\alpha, x_\sigma + 1) = 1$  implies that  $g(\alpha, x_\sigma) > 1$ , which contradicts the hypothesis. The more general case in which  $q > 0$  can be treated essentially in the same way, but requires more tedious algebra. In this case, we need to show that:

$$g(\alpha, x, q) = \sum_{k=x+1}^{\infty} \left[ \frac{\alpha(1-q)}{\alpha+1} \right]^{k-x} \frac{1}{k} \left[ \frac{\alpha+1+x}{\alpha(1-q)} + xq\alpha \right] \geq 1 - \alpha q$$

Evaluating the expression of  $g(\alpha, x - 1, q)$ , one can again show that, if  $g(\alpha, x, q) \geq 1 - \epsilon$ , then  $g(\alpha, x - 1, q) > 1 - \epsilon$ , for any  $\alpha > 0$ ,  $q \geq 0$ ,  $\epsilon \geq 0$ , and repeat the arguments adopted in the case  $q = 0$ .

## References

- [1] Joshua T Abbott, Joseph L Austerweil, and Thomas L Griffiths. Human memory search as a random walk in a semantic network. In *NIPS*, pages 3050–3058, 2012.
- [2] D. S. Bernstein and W. So. Some explicit formulas for the matrix exponential. *IEEE Transactions on Automatic Control*, 38(8):1228–1232, Aug 1993.
- [3] Edward A Codling, Michael J Plank, and Simon Benhamou. Random walk models in biology. *Journal of the Royal Society Interface*, 5(25):813–834, 2008.
- [4] Burgess Davis. Reinforced random walk. *Probability Theory and Related Fields*, 84(2):203–229, 1990.
- [5] Marco Dorigo and Christian Blum. Ant colony optimization theory: A survey. *Theoretical computer science*, 344(2):243–278, 2005.
- [6] Marco Dorigo and Thomas Stützle. *Ant Colony Optimization*. Bradford Company, 2004.

- [7] Rick Durrett. *Probability: theory and examples*. Cambridge University Press, 2010.
- [8] Huilong Huang, John H Hartman, and Terril N Hurst. Data-centric routing in sensor networks using biased walk. In *IEEE Conf. on Sensor and Ad Hoc Communications and Networks (SECON)*, pages 1–9, 2006.
- [9] Human brain project, 2013.
- [10] I. Stoica et al. Chord: a scalable peer-to-peer lookup protocol for internet applications. *IEEE/ACM Transactions on Networking*, 11(1):17–32, 2003.
- [11] Svante Janson. Functional limit theorems for multitype branching processes and generalized pólya urns. *Stochastic Processes and their Applications*, 110(2):177 – 245, 2004.
- [12] Svante Janson. Limit theorems for triangular urn schemes. *Probability Theory and Related Fields*, 134(3):417–452, 2005.
- [13] Jon Kleinberg. The small-world phenomenon: An algorithmic perspective. In *ACM Symposium on Theory of Computing, STOC’00*, pages 163–170, 2000.
- [14] Kenji Leibnitz, Naoki Wakamiya, and Masayuki Murata. Biologically inspired self-adaptive multi-path routing in overlay networks. *Communications of the ACM*, 49(3):62–67, 2006.
- [15] Hosam M. Mahmoud. *Pólya Urn models*. Chapman & Hall/CRC, 2008.
- [16] Andrea Passarella. A survey on content-centric technologies for the current internet: CDN and P2P solutions. *Computer Communications*, 35(1):1–32, 2012.
- [17] Robin Pemantle et al. A survey of random processes with reinforcement. *Probab. Surv*, 4(0):1–79, 2007.
- [18] Sadra Sadeh, Claudia Clopath, and Stefan Rotter. Emergence of functional specificity in balanced networks with synaptic plasticity. *PLoS Comput Biol*, 11(6):1–27, 06 2015.
- [19] Marco Saerens, Youssef Achbany, François Fouss, and Luh Yen. Randomized shortest-path problems: Two related models. *Neural Computation*, 21(8):2363–2404, 2009.
- [20] Sebastian Seung. *Connectome: How the brain’s wiring makes us who we are*. Houghton Mifflin Harcourt, 2012.
- [21] Peter E Smouse, Stefano Focardi, Paul R Moorcroft, John G Kie, James D Forester, and Juan M Morales. Stochastic modelling of animal movement. *Phil. Trans. Royal Soc. of London B: Biological Sciences*, 365(1550):2201–2211, 2010.
- [22] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- [23] John N Tsitsiklis. On the convergence of optimistic policy iteration. *The Journal of Machine Learning Research*, 3:59–72, 2003.