

How Bad is Selfish Routing?

Tim Roughgarden*

Éva Tardos†

Abstract

We consider the problem of routing traffic to optimize the performance of a congested network. We are given a network, a rate of traffic between each pair of nodes, and a latency function for each edge specifying the time needed to traverse the edge given its congestion; the objective is to route traffic such that the sum of all travel times—the total latency—is minimized.

In many settings, including the Internet and other large-scale communication networks, it may be expensive or impossible to regulate network traffic so as to implement an optimal assignment of routes. In the absence of regulation by some central authority, we assume that each network user routes its traffic on the minimum-latency path available to it, given the network congestion caused by the other users. In general such a “selfishly motivated” assignment of traffic to paths will not minimize the total latency; hence, this lack of regulation carries the cost of decreased network performance.

In this paper we quantify the degradation in network performance due to unregulated traffic. We prove that if the latency of each edge is a linear function of its congestion, then the total latency of the routes chosen by selfish network users is at most $\frac{4}{3}$ times the minimum possible total latency (subject to the condition that all traffic must be routed). We also consider the more general setting in which edge latency functions are assumed only to be continuous and non-decreasing in the edge congestion. Here, the total latency of the routes chosen by unregulated selfish network users may be arbitrarily larger than the minimum possible total latency; however, we prove that it is no more than the total latency incurred by optimally routing twice as much traffic.

*Department of Computer Science, Cornell University, Ithaca NY 14853. Research done while visiting the Computer Science Division, UC Berkeley, Berkeley, CA 94720. Supported by ONR grant N00014-98-1-0589. Email: timr@cs.cornell.edu.

†Department of Computer Science, Cornell University, Ithaca NY 14853. Research done in part while a visiting Miller Professor at the Computer Science Division, UC Berkeley, Berkeley, CA 94720. Partially supported by a Guggenheim Fellowship, NSF grant CCR-9700163 and ONR grant N00014-98-1-0589. Email: eva@cs.cornell.edu.

1 Introduction

A fundamental problem arising in the management of large-scale traffic and communication networks is that of routing traffic in order to optimize network performance. One problem of this type is the following: given the rate of traffic between each pair of nodes in a network, find an assignment of traffic to paths so that the sum of all travel times (the *total latency*) is minimized. A difficult aspect of this problem is that the amount of time needed to traverse a single link of a network is typically *load-dependent*, that is, link traversal time increases as the link becomes more congested.

In practice, it is often difficult or even impossible to impose optimal or near-optimal routing strategies on the traffic in a network, and thus network users are free to act according to their own interests, without regard to overall network performance. For example, existing Internet protocols place little restriction on how network traffic is routed, allowing network users to make decisions in a selfish or even malicious manner [3]. The central question of this paper is *how much does network performance suffer from this lack of regulation?*

As a first step toward formalizing this question mathematically, we assume that, in the absence of network regulation, users act in a purely selfish (but not malicious) manner. Under this assumption, we can view network users as independent agents participating in a *non-cooperative game* and expect the routes chosen by users to form a *Nash equilibrium* in the sense of classical game theory [24]. In other words, we assume that each agent uses the minimum-latency path from its source to its destination, given the link congestion caused by the rest of the network users. It is well-known that Nash equilibria do not in general optimize social welfare; perhaps the most famous example is that of “The Prisoner’s Dilemma” [8, 24]. We are then interested in comparing the total latency of a Nash equilibrium with that of the optimal assignment of traffic to paths.

This line of research was recently initiated by Koutsoupias and Papadimitriou [18], who both considered network routing as a non-cooperative game (although in a different model than ours, and only for two-node networks) and proposed the worst-case ratio of the social welfare (suitably

defined) achieved by a Nash equilibrium and by a socially optimal set of strategies as a measure of the performance degradation caused by a lack of regulation. As articulated in [18], this question studies the cost of the lack of *coordination* inherent in a non-cooperative game, as opposed to the cost of a lack of *unbounded computing power* (studied via approximation algorithms) or the cost of a lack of *complete information* (studied via on-line algorithms).

For most of the paper we assume that each agent controls a negligible fraction of the overall traffic. For example, each agent could represent a car and the network a highway system, or agents might represent individual packets in a high-bandwidth communication network; an equilibrium then represents a steady-state in the system (perhaps best achieved in a road network by daily commuters during rush hour and in a communication network by persistent or long-running applications). Under this assumption, a feasible assignment of traffic to paths in the network can be modeled as *network flow*, with the amount of flow between a pair of nodes in the network equal to the rate of traffic between the two nodes. A Nash equilibrium in the aforementioned non-cooperative game then corresponds to a flow where all flow paths between a given source and destination have the same latency (if a flow does not have this property, some agent can improve its travel time by switching from a longer flow path to a shorter one). Beckman et al. [2] showed that if the latency of each network link is a continuous nondecreasing function of the flow on the link, then a flow corresponding to a Nash equilibrium always exists and moreover all such flows have the same total latency. Thus, we can study the cost of routing selfishly via the following question: among all networks with continuous, nondecreasing link latency functions, what is the worst-case ratio between the total latency of a flow at Nash equilibrium and that of an optimal flow (i.e., a flow minimizing the total latency)?

Our Results

In networks in which the latency of each edge is a linear function of the edge congestion (a model that has been the focus of several previous papers [11, 32]), we show that a flow at Nash equilibrium has total latency at most $\frac{4}{3}$ that of the optimal flow. We give examples showing that this result is tight.

We also consider the model in which link latency functions are assumed only to be continuous and nondecreasing. We first show that the ratio between the total latency of a flow at Nash equilibrium and that of an optimal flow may be unbounded in this model. We then work toward *bicriteria* results; in particular, we compare the total latency of a Nash equilibrium flow with that of an optimal flow that routes *additional traffic* between each pair of nodes. Our main result in this setting is that for any network with continuous non-

decreasing latency functions, a flow at Nash equilibrium has total latency no more than that of an optimal flow forced to route twice as much traffic. We again give an example showing that our analysis is tight.

Finally, we examine two unrealistic assumptions made in the basic model: first, the assumption that agents can evaluate the latency of a path with arbitrary precision, and second, that there are an infinite number of agents each controlling a negligible fraction of the overall traffic. We define extensions to the basic model and use them to analyze the sensitivity of our results to these assumptions.

Related Work

Unregulated traffic has been modeled as network flow with all flow paths between a given source-destination pair having equal latency since the 1950's [2, 33] (see also Knight [14]). Beckman et al. [2], observing that such an equilibrium flow is an optimal solution to a related convex program (see also Section 2), gave existence and uniqueness results for traffic equilibria. Dafermos and Sparrow [7] were perhaps the first authors interested in computing the equilibrium efficiently, and many subsequent papers gave increasingly efficient methods for computing equilibria (see [10] for a survey); others have extended these results to more sophisticated models (see for example [1, 6, 10, 13, 21, 22, 27, 29, 30]).

In the past several decades, much of the work on this model has been inspired by a “paradox” first discovered by Braess [4] and later reported by Murchland [20]. The essence of Braess’s Paradox is captured by the example shown in Figure 1, where the edges are labeled with their latency functions (each a function of the link congestion x). Suppose one unit of traffic flow needs to be routed from s to t in the first network of Figure 1. In the unique flow at Nash equilibrium, which coincides with the optimal flow, half of the traffic takes the upper path and the other half travels along the lower path. As all agents are routed on a path of latency $\frac{3}{2}$, the flow has a total latency of $\frac{3}{2}$. Next suppose a fifth edge of latency 0 (independent of the congestion) is added to the network, with the result shown in Figure 1(b). The optimal flow is unaffected by this augmentation (there is no way to use the new link to decrease the total latency) while in the new (unique) flow at Nash Equilibrium, all traffic follows path $s \rightarrow v \rightarrow w \rightarrow t$. The total latency of this flow is 2, as is the latency experienced by each individual agent. Thus, the intuitively helpful (or at least innocuous) action of adding a new zero-latency link may negatively impact *all* of the agents!

After Braess’s Paradox was discovered (together with evidence of similarly counterintuitive and counterproductive traffic behavior following the construction of new roads in congested cities [15, 20]), researchers investi-

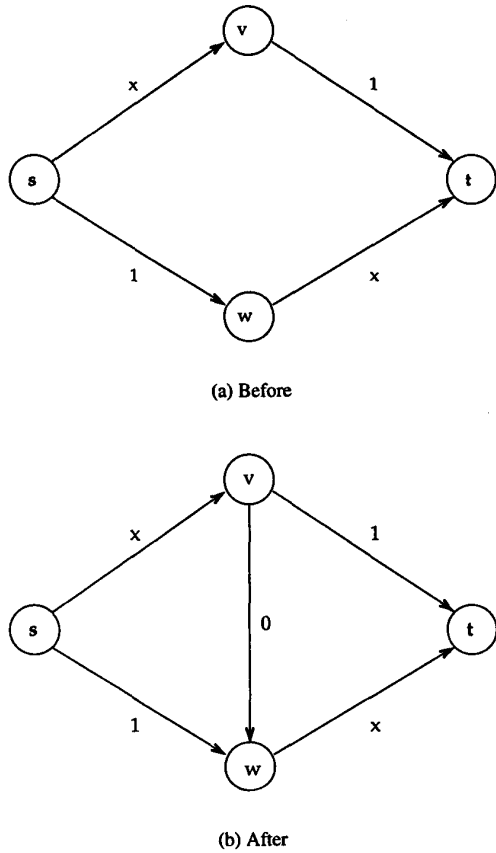


Figure 1. Braess's Paradox

gated the sensitivity of traffic equilibria to the properties of the underlying network [12], classified network topologies in which the addition of new links could degrade network performance [11, 32], discovered new types of “paradoxes” [9, 31], and showed that Braess's Paradox has an analogue in queuing theory [5]. More recently, several papers [16, 17, 23] have investigated a related model in which each agent controls a strictly positive amount of flow (and there are only finitely many agents); classes of network topologies and latency functions guaranteeing existence and uniqueness of Nash equilibria are studied by Orda et al. [23], while Korlis et al. [16, 17] study strategies for adding new edges and/or capacity to a network that guarantee an improvement in network performance. To the best of our knowledge, however, no previous work has attempted to quantify the difference in social welfare between equilibrium and optimal traffic flows.

Finally, the recent paper of Koutsoupias and Papadimitriou [18] is quite similar in spirit to ours, although their

model is fairly different. In [18], a finite number of users share a series of parallel links, and each user chooses a distribution on the set of links (specifying the probability that the agent will route all of its flow on a given link). Each agent wishes to minimize the expected congestion it will experience, while the global objective is to minimize the expected load on the most congested edge. We note that a Nash equilibrium in this model consists of a set of *mixed strategies* (i.e., agents select a distribution on paths) while in our model agents are confined to *pure strategies* (i.e., agents choose a single path); however, there is no essential distinction between pure and mixed strategies under the assumption that each agent controls a negligible amount of traffic. Different Nash equilibria may have different values in the model of [18], so the *worst-case* Nash equilibrium is compared to a globally optimal choice of distributions. Koutsoupias and Papadimitriou obtain tight results in two-node, two-link networks and partial results for two-node networks with three or more parallel links.

Organization

The paper is structured as follows. In Section 2 we give a formal definition of our network model and of flows at Nash equilibrium, and state several lemmas needed for our main results. In Section 3 we prove our main bicriteria result for networks with general edge latency functions. In Section 4 we prove a stronger and technically more involved result for networks with linear edge latency functions. Section 5 considers several extensions to the basic model.

2 Preliminaries

2.1 Model

We consider a directed network $G = (V, E)$ with vertex set V , edge set E , and k source-destination vertex pairs $\{s_1, t_1\}, \dots, \{s_k, t_k\}$. We denote the set of (simple) s_i - t_i paths by \mathcal{P}_i , and define $\mathcal{P} = \cup_i \mathcal{P}_i$. A *flow* is a function $f : \mathcal{P} \rightarrow \mathcal{R}^+$; for a fixed flow f we define $f_e = \sum_{P: e \in P} f_P$. We associate a finite and positive *rate* r_i with each pair $\{s_i, t_i\}$, the amount of flow with source s_i and destination t_i ; a flow f is said to be *feasible* if for all i , $\sum_{P \in \mathcal{P}_i} f_P = r_i$. Each edge $e \in E$ is given a load-dependent *latency* that we denote by $\ell_e(\cdot)$. The latency of a path P with respect to a flow f is then the sum of latencies of the edges on the path, denoted by $\ell_P(f) = \sum_{e \in P} \ell_e(f_e)$. For each $e \in E$, we assume that ℓ_e is non-negative, continuous, and nondecreasing. We define the *cost* $C(f)$ of a flow f in G as the total latency incurred by f , i.e.,

$$C(f) = \sum_{P \in \mathcal{P}} \ell_P(f) f_P.$$

By summing over the edges in a path P and reversing the order of summation, we may also write $C(f) = \sum_{e \in E} \ell_e(f_e) f_e$.

2.2 Flows at Nash Equilibrium

We wish to consider flows that represent an equilibrium among many non-cooperative agents—i.e., flows that behave “greedily” or “selfishly”, without regard to their overall cost. Intuitively, we expect each unit of such a flow (no matter how small) to travel along the minimum-latency path available to it, where latency is measured with respect to the rest of the flow; otherwise, this flow would reroute itself on a path with smaller latency. We formalize this idea in the next definition.

Definition 2.1 *A flow f in G is at Nash equilibrium if for all $i \in \{1, \dots, k\}$, $P_1, P_2 \in \mathcal{P}_i$, and $\delta \in [0, f_{P_1}]$, we have $\ell_{P_1}(f) \leq \ell_{P_2}(\tilde{f})$, where*

$$\tilde{f}_P = \begin{cases} f_P - \delta & \text{if } P = P_1 \\ f_P + \delta & \text{if } P = P_2 \\ f_P & \text{if } P \notin \{P_1, P_2\} \end{cases}$$

Letting δ tend to 0, continuity and monotonicity of the edge latency functions give the following useful characterization of a flow at Nash equilibrium, occasionally called a Wardrop Equilibrium [13] or Wardrop’s Principle [31, 32] in the literature, due to an influential paper of Wardrop [33].

Lemma 2.2 *A flow f is at Nash equilibrium if and only if for every $i \in \{1, \dots, k\}$ and $P_1, P_2 \in \mathcal{P}_i$ with $f_{P_1} > 0$, $\ell_{P_1}(f) \leq \ell_{P_2}(f)$.*

In particular, if f is at Nash equilibrium then all s_i - t_i flow paths (i.e., s_i - t_i paths to which f assigns a positive amount of flow) have equal latency, say $L_i(f)$. Thus, we can express the cost $C(f)$ of a flow f at Nash equilibrium in a particularly nice form.

Lemma 2.3 *If f is a feasible flow at Nash equilibrium, then*

$$C(f) = \sum_{i=1}^n L_i(f) r_i.$$

Remark. Our definition of a flow at Nash equilibrium corresponds to an equilibrium in which each agent chooses a single path of the network (a *pure strategy*), whereas in classical game theory a Nash equilibrium is typically defined via *mixed strategies* (in which an agent may choose a probability distribution over pure strategies) [24]. However, since in our model each agent carries a negligible fraction of the overall traffic, these two definitions are essentially equivalent.

2.3 A Characterization of Optimal Flows via Non-Linear Programming

We now investigate the properties of an optimal flow—i.e., a flow that minimizes total latency. Recalling that the cost of a flow f may be expressed $C(f) = \sum_{e \in E} \ell_e(f_e) f_e$, observe that the problem of finding the minimum-latency feasible flow in a network is a special case of the following non-linear program

$$\begin{aligned} & \text{Min} && \sum_{e \in E} c_e(f_e) \\ & \text{subject to:} && \\ (NLP) & && \sum_{P \in \mathcal{P}_i} f_P = r_i \quad \forall i \in \{1, \dots, k\} \\ & && f_e = \sum_{P \in \mathcal{P}: e \in P} f_P \quad \forall e \in E \\ & && f_P \geq 0 \quad \forall P \in \mathcal{P} \end{aligned}$$

where in our problem, $c_e(f_e) = \ell_e(f_e) f_e$.

For simplicity we have given a formulation with an exponential number of variables, but it is not difficult to give an equivalent compact formulation (with decision variables only on edges and explicit conservation constraints) that requires only polynomially many variables and constraints.

Next, we characterize the local optima of (NLP) . Intuitively, we expect a flow to be locally optimal if and only if moving flow from one path to another can only increase the flow’s cost. Put differently, we expect the *gradient* along any s_i - t_i flow path P (equivalently, the marginal cost of increasing flow on P or the marginal benefit of decreasing flow on P) to be at most the gradient along any other s_i - t_i path (for otherwise moving flow to a smaller-gradient path improves the objective function). Moreover, since the local and global minima of a convex function on a convex set coincide (see, e.g., [25, Thm 2.3.4]), whenever the objective function of (NLP) is convex (e.g., when each edge latency function is convex) this condition should be necessary and sufficient for a flow to be *globally* optimal.

The following lemma formalizes the preceding discussion. Letting $c'_P(f) = \sum_{e \in P} c'_e(f_e)$ (where we are assuming differentiability for simplicity only), we may apply the Karush-Kuhn-Tucker Theorem (see, e.g., [25]) to a convex program of the form (NLP) to derive the following characterization of optimal flows (see the full version [28] for details):

Lemma 2.4 *A flow f is optimal for a convex program of the form (NLP) if and only if for every $i \in \{1, \dots, k\}$ and $P_1, P_2 \in \mathcal{P}_i$ with $f_{P_1} > 0$, $c'_{P_1}(f) \leq c'_{P_2}(f)$.*

The striking similarity between the characterizations of optimal solutions to a convex program of the form (NLP)

and of flows at Nash equilibrium was noticed early on by Beckman et al. [2], and provides an interpretation of an optimal flow as a flow at Nash equilibrium *with respect to a different set of edge latency functions*. More concretely, consider a minimum-latency flow f^* for a convex program of the form (NLP). The flow f^* satisfies the conditions of Lemma 2.4, and so by Lemma 2.2 can be regarded as a flow at Nash equilibrium with respect to latency functions c' .

Now consider the special case of (NLP) where $c_e(f_e)$ has the form $\ell_e(f_e)f_e$ for each edge e . Assuming for convenience that ℓ_e is differentiable, we denote by $\ell_e^*(f_e) = (\ell_e^*(f_e)f_e)' = \ell_e(f_e) + \ell_e'(f_e)f_e$ the marginal cost of increasing flow on edge e ; any flow f^* at Nash equilibrium with respect to latency functions ℓ^* is optimal with respect to the original latency functions ℓ . Interpreting ℓ_e^* as a function with one term capturing per-unit latency and a second term accounting for the degradation in the total latency of the system, we see that the only essential difference between an optimal flow and a flow at Nash equilibrium is that the former evaluates the cost of edge use (via ℓ^*) in a way that accounts for the latency experienced by all flow using the edge, while the latter “selfishly” evaluates edge latency by the per-unit rate ℓ .

Beckman et al. [2] also exploited this similarity between the two characterizations (in particular, that a flow at Nash equilibrium can be regarded as the optimal solution of a convex program of the form (NLP)) to prove the existence and essential uniqueness of Nash equilibria.

Lemma 2.5 ([2]) *A network G with continuous, nondecreasing latency functions admits a feasible flow at Nash equilibrium. Moreover, if f, f' are flows at Nash equilibrium, then $C(f) = C(f')$.*

3 A Bicriteria Result for General Latency Functions

We have already seen (Figure 1(b)) that a flow at Nash equilibrium and a minimum-latency flow may have different costs. In the next two sections, we analyze the *ratio* of the cost of a flow at Nash equilibrium to that of the minimum-latency flow. In this section we work with general (continuous, nondecreasing) latency functions, while in Section 4 we will specialize to the case of linear latency functions.

For a network G with rate vector r and edge latency functions ℓ , admitting an optimal flow f^* and a flow at Nash equilibrium f , we denote the ratio $\frac{C(f)}{C(f^*)}$ by $\rho = \rho(G, r, \ell)$; note that ρ is well-defined by Lemma 2.5.

We begin with some simple negative results. Recall in the canonical example demonstrating Braess’s Paradox (Figure 1) a flow at Nash Equilibrium has total latency 2 while the optimal flow has total latency $\frac{3}{2}$; thus, in the above

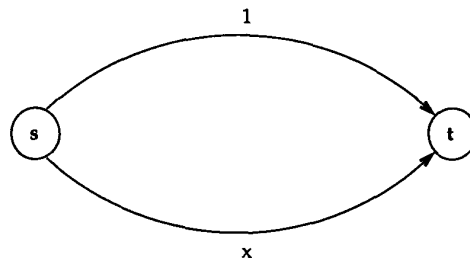


Figure 2. A Simple Bad Example

notation, $\rho = \frac{4}{3}$ in this particular instance. In fact, it is easy to construct an even simpler example (still with linear latency functions) with ratio $\frac{4}{3}$. In the network shown in Figure 2, with a single source-destination pair and rate 1, the flow at Nash equilibrium puts the entire unit of flow on the lower link (with a total latency of 1) while the minimum-latency flow spreads flow evenly across the two links, thereby incurring a cost of $\frac{3}{4}$. Thus, $\rho = \frac{4}{3}$ in this simple instance as well. (In the next section we will prove that this is the worst possible ratio for instances with linear latency functions.)

In fact, the ratio can be much worse when non-linear latency functions are allowed. For a positive integer k , we extend the example of Figure 2 by giving the lower link a latency function of x^k (other input data remains unchanged). The flow at Nash equilibrium again places the entire unit on the lower link, incurring a cost of 1, while the optimal flow assigns $(k+1)^{-1/k}$ units to the lower link and the remainder to the upper link. This solution has a total latency of $1 - k \cdot (k+1)^{-(k+1)/k}$, which tends to 0 as $k \rightarrow \infty$. Thus, assuming only continuity and monotonicity of the edge latency functions, ρ cannot be bounded above.

On the other hand, this example does not rule out interesting *bicriteria* results. Toward this end, we compare the cost of a flow at Nash equilibrium to an optimal flow feasible for *increased rates*. In the example above, an optimal flow feasible for rate $r \geq 1$ assigns the additional flow to the upper link, now incurring a cost that tends to $r - 1$ as $k \rightarrow \infty$. In particular, for any k an optimal flow feasible for twice the rate ($r = 2$) has total latency at least that of the flow at Nash equilibrium (feasible for the original rates). Our main result of this section is a proof of the surprising generalization of this result to *any* network with continuous, nondecreasing edge latencies.

Theorem 3.1 *If f is a flow at Nash equilibrium for (G, r, ℓ) and f^* is feasible for $(G, 2r, \ell)$, then $C(f) \leq C(f^*)$.*

Proof: Suppose f, f^* satisfy the hypotheses of the theorem. For $i = 1, \dots, k$, let $L_i(f)$ be the latency of an s_i - t_i flow path (of f), so that $C(f) = \sum_i L_i(f)r_i$ (see Lemma 2.3).

We seek a set of latency functions $\bar{\ell}$ that on one hand approximates the original ones (in the sense that the cost of a flow with respect to latency functions $\bar{\ell}$ is close to its original cost) and, on the other hand, allows us to easily lower bound the cost (with respect to $\bar{\ell}$) of *any* feasible flow. With this goal in mind, we define new latency functions $\bar{\ell}$ as follows:

$$\bar{\ell}_e(x) = \begin{cases} \ell_e(f_e) & \text{if } x \leq f_e \\ \ell_e(x) & \text{if } x \geq f_e \end{cases}$$

First we compare the cost of the flow f^* under the new latency functions $\bar{\ell}$ to its original cost $C(f^*)$. For any e , $\bar{\ell}_e(x) - \ell_e(x)$ is zero for $x \geq f_e$ and bounded above by $\ell_e(f_e)$ for $x < f_e$, so $x(\bar{\ell}_e(x) - \ell_e(x)) \leq \ell_e(f_e)f_e$ for all $x \geq 0$. Thus, the difference between the new cost (w.r.t. $\bar{\ell}$) and the old cost (w.r.t. ℓ) can be bounded as follows:

$$\begin{aligned} \sum_e \bar{\ell}_e(f_e^*)f_e^* - C(f^*) &= \sum_{e \in E} f_e^*(\bar{\ell}_e(f_e^*) - \ell_e(f_e^*)) \\ &\leq \sum_{e \in E} \ell_e(f_e)f_e \\ &= C(f). \end{aligned}$$

In other words, evaluating f^* with latency functions $\bar{\ell}$ (rather than ℓ) changes its cost by at most an additive $C(f)$ factor.

On the other hand, if f_0 denotes the zero flow in G , then by construction $\bar{\ell}_P(f_0) \geq L_i(f)$ for any $P \in \mathcal{P}_i$. Since $\bar{\ell}_e$ is nondecreasing for each e , it follows that $\bar{\ell}_P(f^*) \geq L_i(f)$ for any $P \in \mathcal{P}_i$. Thus, the cost of f^* with respect to $\bar{\ell}$ given by $\sum_P \bar{\ell}_P(f^*)f_P^*$ is bounded below by

$$\sum_i \sum_{P \in \mathcal{P}_i} L_i(f)f_P^* = \sum_i 2L_i(f)r_i = 2C(f).$$

Combining these two results we obtain the theorem:

$$\begin{aligned} C(f^*) &\geq \sum_P \bar{\ell}_P(f^*)f_P^* - C(f) \\ &\geq 2C(f) - C(f) = C(f). \end{aligned}$$

■

More generally, the proof of Theorem 3.1 shows that if f is at Nash equilibrium for (G, r, ℓ) and f^* is feasible for $(G, (1 + \delta)r, \ell)$, then $C(f) \leq \frac{1}{\delta}C(f^*)$.

Referring back to the bad example at the beginning of the section, we can see that Theorem 3.1 is essentially tight. More precisely, for any $\epsilon > 0$ one can take k sufficiently large to obtain an instance where the optimal flow feasible for rate $2 - \epsilon$ has cost strictly less than 1 (the cost of the flow at Nash equilibrium for the original rates) and the optimal flow feasible for rate 2 has cost at most $1 + \epsilon$.

4 A Worst-Case Ratio of $\frac{4}{3}$ for Linear Latency Functions

Perhaps the simplest model of interest is that of a network with latency functions linear in the amount of flow on an edge. In this section, we consider the case where for each edge $e \in E$, $\ell_e(f_e) = a_e f_e + b_e$ for some $a_e, b_e \geq 0$. This is the setting in which Braess's paradox was originally discovered [4, 20], and several subsequent papers focused entirely on this model [11, 32].

We have already seen (Figures 1 and 2) two examples with linear latency functions for which ρ , the ratio of the cost of a flow at Nash equilibrium and the cost of an optimal flow, is $\frac{4}{3}$. Our main result for this section (Theorem 4.5) is a matching upper bound for networks with linear latency functions. Our proof techniques build on those of the previous section, the primary extension being a more refined approach to lower bounding the cost of an optimal flow.

The results of Section 2 have particularly simple and useful forms in the special case of linear latency functions. First, the total latency of a flow f is given by $\sum_e a_e f_e^2 + b_e f_e$; since $a_e \geq 0$ for all e , $(NL P)$ is a convex (quadratic) program and thus Lemma 2.4 characterizes its optimal solutions. Also, in the notation of subsection 2.3, if $\ell_e(f_e) = a_e f_e + b_e$, then $\ell_e^*(f_e)$, the marginal cost of increasing flow on e , is simply $2a_e f_e + b_e$. For convenience, we summarize this discussion together with specialized versions of Lemmas 2.2 and 2.4 in the following lemma.

Lemma 4.1 *Suppose G is a directed network with k source-sink pairs and with edge latency functions $\ell_e = a_e f_e + b_e$ for each $e \in E$. Then,*

(a) *a flow f is at Nash equilibrium in G if and only if for each i and $P, P' \in \mathcal{P}_i$ with $f_P > 0$,*

$$\sum_{e \in P} a_e f_e + b_e \leq \sum_{e \in P'} a_e f_e + b_e$$

(b) *a flow f^* is (globally) optimal in G if and only if for each i and $P, P' \in \mathcal{P}_i$ with $f_P^* > 0$,*

$$\sum_{e \in P} 2a_e f_e^* + b_e \leq \sum_{e \in P'} 2a_e f_e^* + b_e.$$

As an aside, we note that Lemma 4.1 immediately gives the following non-trivial result regarding networks in which the latency of each edge is proportional to its congestion (i.e., $\ell_e(f_e) = a_e f_e$ for each e).

Corollary 4.2 *Let G be a network in which each edge latency function is of the form $\ell_e(f_e) = a_e f_e$. Then for any rate vector r , a flow feasible for (G, r, ℓ) is optimal if and only if it is at Nash equilibrium.*

Proof: A feasible flow for such an instance satisfies the conditions of Lemma 4.1(a) if and only if it satisfies the conditions of Lemma 4.1(b). ■

A second corollary of Lemma 4.1 will play a crucial role in our proof of the main theorem of this section.

Lemma 4.3 *Suppose (G, r, ℓ) has linear latency functions. Then if f is at Nash equilibrium for (G, r, ℓ) , the flow $f/2$ is optimal for $(G, r/2, \ell)$.*

Proof: If f satisfies the conditions of Lemma 4.1(a), then $f/2$ satisfies the conditions of Lemma 4.1(b). ■

An outline of the proof of the main theorem is as follows. It will be useful to think about creating an optimal flow for the instance (G, r, ℓ) via a two-step process: in the first step, a flow optimal for the instance $(G, r/2, \ell)$ is sent through G , and in the second step this flow is augmented to one optimal for (G, r, ℓ) (note that this augmentation may increase or decrease the amount of flow on any given arc). We will show that the first flow has cost at least $\frac{1}{4}C(f)$ and that the augmentation has cost at least $\frac{1}{2}C(f)$, where f is some flow at Nash equilibrium.

We will see in the proof of Theorem 4.5 that the first lower bound follows easily from Lemma 4.3, but the second (for the cost of the augmentation, given that the first flow has already been routed) requires more work, and in particular the following lemma. Intuitively, the lemma simply claims that the per-unit cost of increasing the amount of flow through a network is at least the gradient with respect to the current optimal flow.

Lemma 4.4 *Suppose (G, r, ℓ) is an instance with linear latency functions for which f^* is an optimal flow. Let $L_i^*(f^*)$ be such that every s_i - t_i flow path P of f^* satisfies $\ell_P^*(f^*) = L_i^*(f^*)$. Then for any $\delta > 0$, a feasible flow for the problem instance $(G, (1 + \delta)r, \ell)$ has cost at least*

$$C(f^*) + \delta \sum_{i=1}^k L_i^*(f^*) r_i.$$

Proof: $L_i^*(f^*)$, the marginal cost of increasing flow on each s_i - t_i flow path, is well-defined for each i by Lemma 4.1(b). If we knew that each L_i^* was nondecreasing in r_i , then routing α additional units of flow from s_i to t_i would cost at least $\alpha \cdot L_i^*(f^*)$ and the lemma would then follow easily by summing over s_i - t_i pairs. Although it is intuitively plausible that marginal costs are increasing in the amount of flow (it is certainly true for each edge individually), the proof requires a little work.

Formally, fix $\delta > 0$ and suppose f is feasible for $(G, (1 + \delta)r, \ell)$. In general f_e may be larger or smaller than f_e^* . For any $e \in E$, convexity of the function $\ell_e(f_e)f_e = a_e f_e^2 + b_e f_e$ implies that

$$\ell_e(f_e)f_e \geq \ell_e(f_e^*)f_e^* + (f_e - f_e^*)\ell_e^*(f_e^*).$$

In essence, this inequality states that estimating the cost of changing the flow value on edge e to f_e by $(f_e - f_e^*)\ell_e^*(f_e^*)$ (i.e., by the marginal cost of flow increase at f_e^* times the size of the perturbation) only underestimates the actual cost of an increase (when $f_e > f_e^*$) and overestimates the actual benefit of a decrease (when $f_e < f_e^*$).

Thus,

$$\begin{aligned} C(f) &= \sum_{e \in E} \ell_e(f_e)f_e \\ &\geq \sum_{e \in E} \ell_e(f_e^*)f_e^* + \sum_{e \in E} (f_e - f_e^*)\ell_e^*(f_e^*) \\ &= C(f^*) + \sum_{i=1}^k \sum_{P \in \mathcal{P}_i} \ell_P^*(f^*)(f_P - f_P^*) \\ &= C(f^*) + \sum_{i=1}^k L_i^*(f^*) \sum_{P \in \mathcal{P}_i} (f_P - f_P^*) \\ &= C(f^*) + \delta \sum_{i=1}^k L_i^*(f^*) r_i. \end{aligned}$$

■

We remark that Lemma 4.4 and its proof remain valid in much more general settings, for example when all the edge latency functions are convex.

We are now prepared to prove the main theorem.

Theorem 4.5 *If (G, r, ℓ) has linear latency functions, then $\rho(G, r, \ell) \leq \frac{4}{3}$.*

Proof: Let f be a flow in G at Nash equilibrium. Let $L_i(f)$ be the latency of an s_i - t_i flow path, so that $C(f) = \sum_i L_i(f)r_i$ (see Lemma 2.3). By Lemma 4.3, $f/2$ is an optimal solution to the instance $(G, r/2, \ell)$. Moreover, in the notation of Lemma 4.4, $\ell_e^*(f_e/2) = \ell(f_e)$ for each edge e and hence $L_i^*(f/2) = L_i(f)$ for each i (in words, marginal costs of edges and paths w.r.t. $f/2$ and latencies of edges and paths w.r.t. f coincide); this establishes the necessary connection between the cost of augmenting $f/2$ to a flow feasible for (G, r, ℓ) and the cost of a flow at Nash equilibrium, f .

Taking $\delta = 1$ in Lemma 4.4, we find that the cost of any flow f^* feasible for (G, r, ℓ) satisfies

$$\begin{aligned} C(f^*) &\geq C(f/2) + \sum_{i=1}^k L_i^*(f/2) \frac{r_i}{2} \\ &\geq C(f/2) + \frac{1}{2} \sum_{i=1}^k L_i(f) r_i \\ &= C(f/2) + \frac{1}{2} C(f). \end{aligned}$$

Finally, it's easy to lower bound the cost of $f/2$:

$$\begin{aligned} C(f/2) &= \sum_e \frac{1}{4} a_e f_e^2 + \frac{1}{2} b_e f_e \\ &\geq \frac{1}{4} \sum_e a_e f_e^2 + b_e f_e \\ &= \frac{1}{4} C(f) \end{aligned}$$

and thus $C(f^*) \geq \frac{3}{4} C(f)$. ■

We note that the analysis of this section can be extended to prove that in any network G with rate vector r where for some k , $\ell_e = a_e x^k + b_e$ for each e , $\rho(G, r, \ell) \leq (1 - k \cdot (k + 1)^{-(k+1)/k})^{-1}$. The example at the beginning of Section 3 shows that this result is tight.

5 Extensions

In this section we extend the basic model in several ways, as the model of flow considered thus far suffers from several drawbacks. First, in practice agents cannot evaluate path latency exactly, only approximately. Subsection 5.1 extends the notion of a flow at Nash equilibrium and Theorem 3.1 to this setting. Second, our basic model represents a scenario with infinitely many agents each controlling an infinitesimal amount of flow, while in practice we expect to encounter a finite number of agents, each controlling a strictly positive amount of flow. In subsection 5.2 we prove an analogue of Theorem 3.1 for the case of finitely many agents, provided each agent can route its flow fractionally over any number of paths. In subsection 5.3 we show that such an assumption is essentially necessary, in that no bicriteria result analogous to Theorem 3.1 holds when there are only finitely many agents, each of whom must route its flow on a single path. Proofs of the results in this section may be found in the full version of the paper [28].

5.1 Flows at Approximate Nash Equilibrium

It is unreasonable to expect agents to be able to evaluate the latency of different paths with arbitrary precision. We next investigate the sensitivity of our results to this assumption. We suppose that an agent can distinguish between paths that differ significantly in their latency (say by more than a $(1 + \epsilon)$ factor for some $\epsilon > 0$). Our definition of a flow at ϵ -approximate Nash equilibrium is then an obvious modification of Definition 2.1:

Definition 5.1 A flow f in G is at ϵ -approximate Nash equilibrium if for all $i \in \{1, \dots, k\}$, $P_1, P_2 \in \mathcal{P}_i$, and $\delta \in [0, f_{P_1}]$, we have $\ell_{P_1}(f) \leq (1 + \epsilon) \ell_{P_2}(\tilde{f})$, where

$$\tilde{f}_P = \begin{cases} f_P - \delta & \text{if } P = P_1 \\ f_P + \delta & \text{if } P = P_2 \\ f_P & \text{if } P \notin \{P_1, P_2\} \end{cases}$$

The analogue of Lemma 2.2 is then:

Lemma 5.2 A flow f is at ϵ -approximate Nash equilibrium if and only if for every $i \in \{1, \dots, k\}$ and $P_1, P_2 \in \mathcal{P}_i$ with $f_{P_1} > 0$, $\ell_{P_1}(f) \leq (1 + \epsilon) \ell_{P_2}(f)$.

With minor modifications to the proof of Theorems 3.1, we can then prove the following result.

Theorem 5.3 If f is at ϵ -approximate Nash equilibrium for (G, r, ℓ) and f^* is feasible for $(G, 2r, \ell)$, then $C(f) \leq (1 + \epsilon) C(f^*)$.

5.2 Finitely Many Agents: Splittable Flow

Our basic model makes the unrealistic assumption that flow is comprised of infinitely many independent agents. In this subsection we extend the basic model to the case of finitely many agents, each of whom controls a strictly positive amount of flow. In this subsection we allow an agent to split flow along any number of paths; the next subsection investigates the case where each agent must route all of its flow on a single path.

We are given a network G with continuous nondecreasing latency functions ℓ as before, and in addition k agents. We assume that agent i intends to send r_i units of flow from source s_i to destination t_i . Distinct agents may have identical source-destination pairs. We continue to denote an instance by (G, r, ℓ) , and we call the instance *finite splittable*. A flow f now consists of k functions, $f^{(i)} : \mathcal{P}_i \rightarrow \mathcal{R}^+$ for agent i . For a flow f , we denote by $C_i(f)$ the total latency experienced by agent i ; thus, $C_i(f) = \sum_{P \in \mathcal{P}_i} \ell_P(f) f_P^{(i)}$. As usual, a flow is at Nash equilibrium if no agent can decrease the latency it experiences by rerouting its flow. In this setting, a flow f is at Nash equilibrium if and only if for each i , $f^{(i)}$ minimizes $C_i(f)$ given $f^{(j)}$ for $j \neq i$. It follows from results of Rosen [26] that such a flow exists (and in fact is essentially unique) under mild convexity assumptions, for example when all latency functions are convex.

Our main result for this model is an analogue of Theorem 3.1, proved using similar ideas.

Theorem 5.4 If f is at Nash equilibrium for the finite splittable instance (G, r, ℓ) and f^* is feasible for the finite splittable instance $(G, 2r, \ell)$, then $C(f) \leq C(f^*)$.

Theorem 3.1 can thus be regarded as the limiting case of the above theorem, as the number of agents tends to infinity and the amount of flow controlled by each agent tends to 0.

5.3 Finitely Many Agents: Unsplittable Flow

In this subsection we continue our investigation of selfish routing with finitely many agents, each controlling a non-negligible amount of flow. It is easy to imagine scenarios in

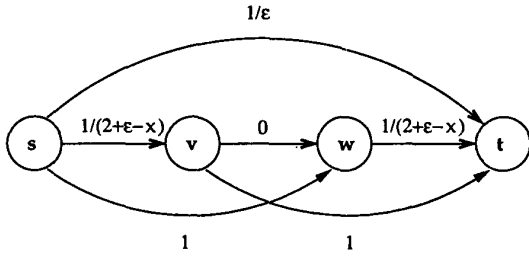


Figure 3. A Bad Example for Unsplittable Flow

which agents cannot route flow on several different paths, but instead must select a single path for routing. Our previous results have made crucial use of the “infinitely divisible” nature of flow, and we next show that this assumption is essentially necessary.

Consider an instance (G, r, ℓ) as in the previous subsection (with k agents and the i th agent controlling r_i units of flow), but with the additional constraint that each agent selects a *single* path on which to route all of its flow. We call such an instance *finite unsplittable*. Adapting the definition of the previous subsection to this new setting, a flow f (now consisting only of k paths) is at Nash equilibrium if and only if for each i , agent i routes its flow on a path minimizing $\ell_P(f)$ (with P ranging over all paths in \mathcal{P}_i), given the paths chosen by the other $k - 1$ agents.

We first consider a simple example showing that a flow at Nash equilibrium may have cost arbitrarily larger than that of an optimal flow. Consider the network given in Figure 3, and suppose there are two agents, each of whom has source s , destination t , and one unit of flow to send; $\epsilon > 0$ is arbitrary. In the optimal solution, one agent chooses path $s \rightarrow v \rightarrow t$ and the other $s \rightarrow w \rightarrow t$; the cost of this solution is less than 4 (for any $\epsilon > 0$). On the other hand, a solution with one agent choosing path $s \rightarrow v \rightarrow w \rightarrow t$ and the other routing on the $s \rightarrow t$ link is a flow at Nash equilibrium with cost greater than $\frac{1}{\epsilon}$; by choosing ϵ arbitrarily small this cost is arbitrary large, and hence arbitrarily more costly than optimal.

In light of the example at the beginning of Section 3, such a result is hardly surprising; however, this example easily extends to show that bicriteria statements analogous to the theorems of Sections 3 and 5.2 are false when we require agents to route flow unsplittably. In particular, for any positive integer q , consider the network G_q consisting of $2q + 2$ vertices arranged in a path $s, v_1, v_2, \dots, v_{2q}, t$ with edges along the path alternately having latency functions $\frac{1}{2+\epsilon-x}$ and 0, a direct $s-t$ link with constant latency function $\frac{1}{\epsilon}$, and arcs from s to v_{2i} and from v_{2i-1} to t with constant latency functions $\ell(x) = 1$ (e.g., $q = 1$ in Figure 3). As in

the previous paragraph, there is a flow at Nash equilibrium with two agents, each controlling one unit of flow, with cost greater than $\frac{1}{\epsilon}$. On the other hand, it is possible for $q + 1$ agents to each send one unit of flow through G_q at total cost at most $3q$ (the first agent uses path $s \rightarrow v_1 \rightarrow t$, the last $s \rightarrow v_{2q} \rightarrow t$, and otherwise the i th agent uses path $s \rightarrow v_{2i-2} \rightarrow v_{2i-1} \rightarrow t$). Letting ϵ tend to 0 for each fixed value of q , we see that an optimal flow can send *arbitrarily more flow at arbitrarily less cost* than a flow at Nash equilibrium.

In the above bad example, the network has latency functions with unbounded derivatives; in this situation, routing a strictly positive amount of additional flow on an edge may increase the latency of that edge by an arbitrarily large amount. This example is of particular interest as functions of the form $\ell(x) = 1/(u - x)$ have been used in several different models considered in the literature [16, 17, 19, 23] with the intention of modeling a link with capacity u . However, in networks where the largest possible change in edge latency resulting from a single agent rerouting its flow is bounded above, we can apply the results of subsection 5.1 to derive the following.

Theorem 5.5 *Suppose f is at Nash equilibrium in the finite unsplittable instance (G, ℓ, r) , and for some $\alpha \geq 1$, $\ell_e(x + r_i) \leq \alpha \cdot \ell_e(x)$ for all $i \in \{1, \dots, k\}, e \in E, x \in [0, \sum_{j \neq i} r_j]$. Then for any f^* feasible for $(G, 2r, \ell)$, $C(f) \leq \alpha \cdot C(f^*)$.*

For example, in an instance with linear latency functions (say $\ell_e(f_e) = a_e f_e + b_e$) with $b_e > 0$ for all edges e , we may apply Theorem 5.5 with $\alpha = 1 + \max_i r_i \cdot \max_e a_e/b_e$.

Acknowledgements

We would like to thank Leonard Schulman, Christos Papadimitriou, and Satish Rao for a discussion about Braess’s paradox, and for suggesting the comparison of the social value of a Nash equilibrium to the optimal social value. This discussion was the start of our research. We are also grateful to Scott Shenker for numerous insightful conversations about possible extensions of our results, and for his interest in our work.

References

- [1] H. Z. Aashtiani and T. L. Magnanti. Equilibria on a congested transportation network. *SIAM Journal on Algebraic and Discrete Methods*, 2(3):213–226, 1981.
- [2] M. Beckmann, C. B. McGuire, and C. B. Winsten. *Studies in the Economics of Transportation*. Yale University Press, 1956.

- [3] B. Braden, D. Clark, J. Crowcroft, B. Davie, S. Deering, D. Estrin, S. Floyd, V. Jacobson, G. Minshall, C. Partridge, L. Peterson, K. Ramakrishnan, S. Schenker, J. Wroclawski, and L. Zhang. Recommendations on queue management and congestion avoidance in the Internet. Network Working Group Request for Comments 2309, April 1998.
- [4] D. Braess. Über ein paradoxon der verkehrsplanung. *Unternehmensforschung*, 12:258–268, 1968.
- [5] J. E. Cohen and F. P. Kelly. A paradox of congestion in a queuing network. *Journal of Applied Probability*, 27:730–734, 1990.
- [6] S. Dafermos. Traffic equilibrium and variational inequalities. *Transportation Science*, 14(1):42–54, 1980.
- [7] S. C. Dafermos and F. T. Sparrow. The traffic assignment problem for a general network. *Journal of Research of the National Bureau of Standards, Series B*, 73B(2):91–118, 1969.
- [8] P. Dubey. Inefficiency of Nash equilibria. *Mathematics of Operations Research*, 11(1):1–8, 1986.
- [9] C. Fisk. More paradoxes in the equilibrium assignment problem. *Transportation Research*, 13B:305–309, 1979.
- [10] M. Florian. Nonlinear cost network models in transportation analysis. *Mathematical Programming Study*, 26:167–196, 1986.
- [11] M. Frank. The Braess Paradox. *Mathematical Programming*, 20:283–302, 1981.
- [12] M. A. Hall. Properties of the equilibrium state in transportation networks. *Transportation Science*, 12(3):208–216, 1978.
- [13] A. Haurie and P. Marcotte. On the relationship between Nash-Cournot and Wardrop equilibria. *Networks*, 15:295–308, 1985.
- [14] F. H. Knight. Some fallacies in the interpretation of social cost. *Quarterly Journal of Economics*, 38:582–606, 1924.
- [15] W. Knödel. *Graphentheoretische Methoden und ihre Anwendungen*. Springer-Verlag, 1969.
- [16] Y. A. Korlis, A. A. Lazar, and A. Orda. Capacity allocation under noncooperative routing. *IEEE Transactions on Automatic Control*, 42(3):309–325, 1997.
- [17] Y. A. Korlis, A. A. Lazar, and A. Orda. Avoiding the Braess paradox in noncooperative networks. *Journal of Applied Probability*, 42(3):309–325, 1999.
- [18] E. Koutsoupias and C. Papadimitriou. Worst-case equilibria. In *Proceedings of the 16th Annual Symposium on Theoretical Aspects of Computer Science*, pages 404–413, 1999.
- [19] A. A. Lazar, A. Orda, and D. E. Pendarakis. Virtual path bandwidth allocation in multiuser networks. *IEEE/ACM Transactions on Networking*, 5:861–871, 1997.
- [20] J. D. Murchland. Braess’s paradox of traffic flow. *Transportation Research*, 4:391–394, 1970.
- [21] Y. Nesterov. Stable flows in transportation networks. CORE Discussion Paper 9907, 1999.
- [22] Y. Nesterov and A. D. Palma. Stable dynamics in transportation systems. CORE Discussion Paper 00/27, 2000.
- [23] A. Orda, R. Rom, and N. Shimkin. Competitive routing in multi-user communication networks. *IEEE/ACM Transactions on Networking*, 1:510–521, 1993.
- [24] G. Owen. *Game Theory*. Academic Press, 1995. Third Edition.
- [25] A. L. Peressini, F. E. Sullivan, and J. J. Uhl. *The Mathematics of Nonlinear Programming*. Springer-Verlag, 1988.
- [26] J. B. Rosen. Existence and uniqueness of equilibrium points for concave N -person games. *Econometrica*, 33(3):520–534, 1965.
- [27] R. W. Rosenthal. A class of games possessing pure-strategy Nash equilibria. *International Journal of Game Theory*, 2:65–67, 1973.
- [28] T. Roughgarden and É. Tardos. How bad is selfish routing? <http://www.cs.cornell.edu/timr>.
- [29] Y. Sheffi. *Urban Transportation Networks: Equilibrium Analysis with Mathematical Programming Methods*. Prentice-Hall, 1985.
- [30] M. J. Smith. The existence, uniqueness and stability of traffic equilibria. *Transportation Research*, 13B:295–304, 1979.
- [31] R. Steinberg and R. E. Stone. The prevalence of paradoxes in transportation equilibrium problems. *Transportation Science*, 22(4):231–241, 1988.
- [32] R. Steinberg and W. I. Zangwill. The prevalence of Braess’ paradox. *Transportation Science*, 17(3):301–318, 1983.
- [33] J. G. Wardrop. Some theoretical aspects of road traffic research. In *Proceedings of the Institute of Civil Engineers, Pt. II*, volume 1, pages 325–378, 1952.