

Stackelberg Routing on Horizontal Queueing Networks

Walid Krichene, Jack Reilly, Saurabh Amin, Alexandre Bayen

Abstract—In order to address inefficiencies of Nash equilibria for horizontal queueing congestion networks, we study the Stackelberg routing game on parallel networks: assuming a coordinator has control over a fraction of the flow, what is a good way of routing that compliant flow so that the induced Nash equilibrium is closer to the social optimum than the initial Nash equilibrium?

We study Stackelberg routing for a new class of latency functions, that models congestion on horizontal queueing networks. We show that in this setting, optimal Stackelberg strategies can be computed in polynomial time in the size of the network. We introduce a candidate Stackelberg strategy, the *non-compliant first* strategy, and prove it to be optimal. We apply these results by modeling a transportation network in which a coordinator can choose the routes of fraction of compliant agents, while the rest of the agents choose their routes selfishly.

Some proofs are omitted for brevity, and are available for the reviewer’s convenience in [6].

I. INTRODUCTION

A. Congestion Games and Stackelberg Routing

Nash equilibria of congestion games (or user optimal assignments) have been extensively studied [8], [9], [11] and are known to be inefficient compared to the system optimal assignments, where a coordinator assigns flow as to minimize a system-wide cost function [2], [13].

In order to address this inefficiency, Stackelberg routing games have been proposed, in which a fraction of the agents is centrally controlled, while the rest is free to act selfishly [2], [10]. The objective of the coordinator is to assign the compliant flow in a manner that minimizes a system-wide cost function in anticipation of the rest of the flow’s selfish response.

Congestion games and Stackelberg routing on parallel networks have been studied extensively for the class of non-decreasing latency functions, and it is known that computing the optimal Stackelberg strategy is NP-hard in the number of links [10]. This led to considering polynomial time approximate strategies such as Largest Latency First and Scale [10], and several bounds have been shown on the efficiency of these strategies. While this class of latency functions provides a good model of congestion for a considerable range of networks, such as communication networks, it does not accurately model horizontal queueing congestion, such as congestion on transportation networks [4], [7]. A new class of latency functions is introduced in [6] to model congestion on horizontal queueing networks. We study Stackelberg routing for this new class of latency, which leads to novel theoretic and algorithmic results.

B. Motivating Application: Compliant and Non-compliant Drivers on Highway Networks

Advances in technology have made it possible to interact with individual drivers on a traffic network and exchange information through GPS-enabled smartphone applications or before and after-market vehicular navigation systems. This offers an opportunity to not only provide the driver with relevant traffic information and collect anonymized data that help improve traffic estimation, but also to provide routing advice that can improve the overall efficiency of the network by relieving congestion. However, when providing routing advice, one needs to take into account the possible impact of rerouting drivers on the traffic conditions of the network, and the response of other drivers to this change in traffic conditions. This fits into the framework of Stackelberg routing, in which a fraction of the population of drivers is assumed to be *compliant* to routing suggestions, and the rest of the drivers are considered to be *non-compliant*. We call compliant a driver who is connected to a central coordinator, through a smartphone application or a navigation system, and is willing (or has an external incentive) to follow routing suggestions provided by the coordinator. Other drivers (drivers who are not connected or who are simply not willing to follow alternative routes) are described as non-compliant. These two populations form respectively the leader and the followers in the Stackelberg routing game.

In numerous transportation networks, two highly populated areas can be connected by disjoint highways (we consider one such example in the numerical results section). Therefore, Stackelberg routing strategies on simple parallel networks are of practical importance to traffic planners [3]. However, due to the limitations of congestion models discussed above, the existing literature on Stackelberg routing is not readily applicable to practical traffic networks. The present work addresses these limitations by considering a new class of latency functions that better models congestion for horizontal queueing network, in particular traffic networks.

C. Contributions

We study the Stackelberg routing game on parallel networks for the class of *horizontal queueing congestion* latencies, and show that optimal Stackelberg strategies can be computed in $O(N^2)$ time, where N is the number of links in the network. This result contrasts with the class of non-decreasing latency functions, for which computing the optimal Stackelberg strategy is NP-hard [10]. We define in particular the *non-compliant first* (NCF) strategy and prove it to be optimal.

We then apply these results to model a real transportation network, and identify ranges of the flow demand and compliance rates where optimal Stackelberg routing are most efficient, and quantify the decrease in inefficiency achieved by the NCF strategy.

These results are an encouraging and necessary step towards a scalable, accurate model for optimal route assignment on horizontal queueing networks with partial compliance.

D. Organization of the Article

In Section II, we start by defining the congestion game and the class of horizontal queueing congestion latencies, then review some properties of Nash equilibria for horizontal queueing networks. The main results on optimal Stackelberg routing are presented in Sections III and IV, where a polynomial time algorithm is presented for computing a provably optimal Stackelberg strategy on N link parallel networks. This is followed by numerical results that illustrate the effects of optimal Stackelberg routing in Section V. We conclude with a summary of our results and directions for future work in Section VI.

II. PRELIMINARIES

A. The Model: Congestion games and the class of Horizontal Queueing Latencies

We consider a non-atomic [12] congestion game on a parallel network with a single source, a single sink (or destination) and N parallel edges (or links) indexed by $n \in \{1, \dots, N\}$. The network is subject to a constant positive flow demand r at the source. We will denote by (N, r) an instance of a network with N parallel links subject to demand r .

A flow assignment for the instance (N, r) is a vector $x \in \mathcal{R}_+^N$ such that $\sum_{n=1}^N x_n = r$ where x_n is the flow on link n . We will denote by $\text{Supp}(x)$ the support of x , i.e. the set of links that hold strictly positive flow $\{n \in \{1, \dots, N\} | x_n > 0\}$.

Every (non-atomic) agent chooses a route to go from the source to the destination, and in this simple setting, every agent simply chooses a link. All agents on link n experience the same latency, and we assume that the total flow x_n affects the latency on link n . As detailed in [6], in the case of horizontal queueing networks, the latency also depends on whether the link is in *free-flow* (the density is below a critical density) or *congested* (the density is above the critical density). Intuitively, a given flow x_n corresponds to two different configurations:

- either few agents moving fast (the density is low and the link is in free-flow), in which case the latency is low,
- or many agents moving slowly (the density is high and the link is congested), in which case the latency is high.

Let $m_n \in \{0, 1\}$ be the congestion state of link n :

$$m_n = \begin{cases} 0 & \text{if } n \text{ is in free-flow} \\ 1 & \text{if } n \text{ is congested} \end{cases}$$

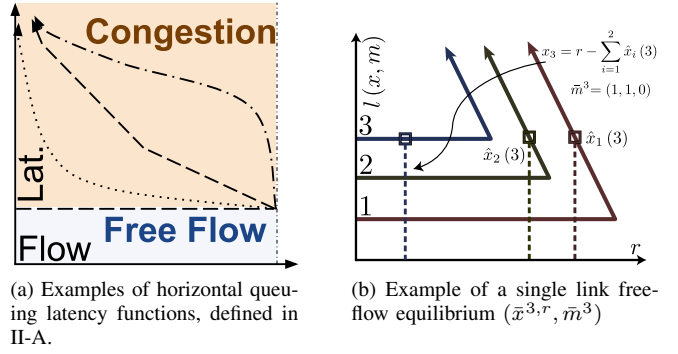


Fig. 1: Horizontal queueing congestion latencies and single link free-flow equilibria.

Then the latency on the link n is given by

$$l_n : [0, x_n^{\max}] \times \{0\} \cup (0, x_n^{\max}] \times \{1\} \rightarrow \mathcal{R}_+ \\ (x_n, m_n) \mapsto l_n(x_n, m_n)$$

We further require the latency function to satisfy the following properties:

- The latency is constant in free-flow: $\forall x_n \in [0, x_n^{\max}]$, $l_n(x_n, 0) = a_n > 0$.
- $\lim_{x_n \rightarrow x_n^{\max}} l_n(x_n, 1) = l_n(x_n^{\max}, 0) = a_n$
- $x_n \mapsto l_n(x_n, 1)$ is decreasing on $(0, x_n^{\max})$.

This defines the class of horizontal queueing latencies, introduced in [6]. The last property can be justified intuitively by the following argument: as the density on the link increases, the agents slow down (to avoid collision), therefore the flow decreases and the latency increases.

Some examples of latency functions in this class are illustrated in Figure 1a.

We further assume, to simplify our discussion, that the free-flow latencies are different, and that links are ordered by increasing free-flow latencies:

$$a_1 < a_2 < \dots < a_N$$

Total System Cost: Every non-atomic agent on link n experiences latency $l_n(x_n, m_n)$, therefore the total latency experienced by the flow x_n is $C_n(x_n, m_n) = x_n l_n(x_n, m_n)$ and the total system cost (latency experienced by the total flow) is simply

$$C(x, m) = \sum_{n=1}^N C_n(x_n, m_n) \quad (1)$$

B. Nash Equilibria

In this section, we briefly review some properties of Nash equilibria for congestion games on parallel networks for the class of horizontal queueing latencies. These properties will be useful in studying the Stackelberg routing game. For a more detailed discussion and proofs, we refer the reader to [6].

Consider a network instance (N, r) . A pure Nash equilibrium, which we simply refer to as Nash equilibrium, is an

assignment in which every non-atomic agent cannot improve her/his latency by switching to another link:

Definition 1: Nash Equilibrium

An assignment $(x, m) \in \mathcal{R}_+^N \times \{0, 1\}^N$ feasible for the network instance (N, r) is a Nash equilibrium, if $\forall n$

$$x_n > 0 \Rightarrow \forall k \leq N, l_n(x_n, m_n) \leq l_k(x_k, m_k)$$

In particular, all links in the support of x have the same latency, and if l_0 denote the common latency on the support, then the total system cost incurred by the network is simply $C(x, m) = r l_0$.

Single link free-flow equilibria: Let $\text{NE}(N, r)$ denote the set of Nash Equilibria for network instance (N, r) . For the class of horizontal queuing latency functions, there may exist multiple Nash equilibria with different costs. For an example of a network instance where this is the case, see [6]. The set of Nash equilibria can then be partitioned into single link free-flow equilibria (equilibria such that the last link in the support is in free-flow) and congested equilibria (such that all links in the support are congested). One can show that these are indeed the only possible equilibria.

We now focus our attention on single link free-flow equilibria. The following lemma characterizes the congestion state vectors for these equilibria:

Lemma 1: Congestion states under single link free-flow equilibrium

Let $(x, m) \in \text{NE}(N, r)$. Assume that $\exists j \in \text{Supp}(x)$ such that $m_j = 0$. Then $m = (1, \dots, \overset{j-1}{1}, \overset{j}{0}, \dots, 0)$ and $\text{Supp}(x) = \{1, \dots, j\}$.

The lemma states that if some link k in the support of a Nash equilibrium is in free-flow, this completely determines the congestion state vector of the equilibrium: links $\{1, \dots, k-1\}$ are in the support and are congested, and links $\{k+1, \dots, N\}$ are not in the support. Note that this also determines the flow vector: since link k is in free flow and in the support, its latency is $l_k(x_k, 0) = a_k$. Therefore every link in the support has the same latency, in particular $\forall n \in \{1, \dots, k-1\}$, $l_n(x_n, 1) = a_k$. The unique flow that satisfies this equality is referred to as *congestion flow*. More precisely,

Definition 2: Congestion flow

For $1 \leq n < k \leq N$, the congestion flow $\hat{x}_n(k)$ is defined as the unique flow in $(0, x_n^{\max})$ that satisfies

$$l_n(\hat{x}_n(k), 1) = a_k \quad (2)$$

Note that the congestion flow $\hat{x}_n(k) = l_n(\cdot, 1)^{-1}(a_k)$ is a decreasing function of k since a_k is increasing in k and $l_n(\cdot, 1)^{-1}$ is decreasing.

One can then show that all single link free-flow equilibria are of the form $(\bar{x}^{k,r}, \bar{m}^k)$ where

$$\bar{m}^k := (1, \dots, \overset{k-1}{1}, \overset{k}{0}, \dots, 0) \quad (3)$$

$$\bar{x}^{k,r} := (\hat{x}_1^1(k), \dots, \hat{x}_{k-1}^{k-1}(k), r - \sum_{n=1}^{k-1} \hat{x}_n(k), 0, \dots, 0) \quad (4)$$

and $(\bar{x}^{k,r}, \bar{m}^k)$ is an equilibrium if and only if it is a feasible assignment, i.e. $r - \sum_{n=1}^{k-1} \hat{x}_n(k) \in [0, x_k^{\max}]$.

An example of single link free-flow equilibrium is shown in Figure 1b.

Lemma 2: Existence of a single link free-flow Nash equilibrium

$\forall r \in [0, \max_{1 \leq k \leq N} \{x_k^{\max} + \sum_{n=1}^{k-1} \hat{x}_n(k)\}]$, there exists a single link free-flow Nash equilibrium for the instance (N, r) .

This lemma in fact shows that if the set of Nash equilibria is non empty, then it contains a single link free-flow equilibrium.

C. Best Nash Equilibrium

In order to study the inefficiency of Nash equilibria, and the improvement of performance that we can achieve using a Stackelberg game (in which a fraction of the total flow is controlled by a central authority), we focus our attention on *best Nash equilibria* and *price of stability* as a measure of their inefficiency (see for example [1]). A *best Nash equilibrium* (BNE) is defined to be a Nash equilibrium of least total latency $\text{BNE}(N, r) = \arg \min_{(x, m) \in \text{NE}(N, r)} C(x, m)$.

The following theorem characterizes best Nash equilibria.

Theorem 1: Characterization of Best Nash Equilibria [6]

For a parallel network instance (N, r) , the unique best Nash equilibrium is the single-link free-flow equilibrium that has smallest support

$$\text{BNE}(N, r) = \arg \min_{(x, m) \in \text{NE}(N, r)} \{\max[\text{Supp}(x)]\}$$

Theorem 1 provides a simple characterization of the best Nash equilibrium for any network instance (N, r) , and shows in particular that the best Nash equilibrium can be computed in $O(N^2)$ where N is the size of the network: the best Nash equilibrium can be computed by simply enumerating all candidate single link free-flow equilibria $(\bar{x}^{k,r}, \bar{m}^k)$, starting from the smallest support ($k = 0$). There are N such candidate equilibria, corresponding to the congestion states $(0, \dots, 0)$ up to $(1, \dots, 1, 0)$, and each candidate equilibrium is a vector in \mathcal{R}^N that can be computed in $O(N)$, which corresponds to a worst-case time complexity of $O(N^2)$.

III. STACKELBERG ROUTING

In order to address the inefficiency of Nash equilibria due to selfish routing and lack of coordination, we assume that a fraction α of the flow is centrally controlled, and we investigate possible strategies for improving the equilibria of the network. Leader-follower routing games have been considered in the transportation literature [2], [5]. However, latency functions considered in the previous literature do not model decrease in flow on a link as a result of density buildup, while the class of latency functions introduced in [6] better models horizontal queuing. In this section, we setup the problem and introduce useful definitions. In the next section, we show that optimal Stackelberg strategies for the class of horizontal queuing latency can be computed in polynomial time, and give a constructive algorithm for computing such an optimal strategy.

A. Stackelberg routing game

We consider the following problem: given a network instance (N, r) under constant flow demand r , assume a coordinator (a central authority) has control over a fraction α of the flow: the corresponding agents are *compliant* and willing to let the coordinator choose their routes. The coordinator (who plays the role of the leader in the Stackelberg game) assigns the compliant flow αr according to a Stackelberg strategy s that is a feasible flow assignment for the instance $(N, \alpha r)$, i.e. s satisfies: $s_n \leq x_n^{\max} \forall n \leq N$ and $\sum_n s_n = \alpha r$.

We assume that the remaining flow $(1 - \alpha)r$ represents selfish agents (who play the role of followers in the Stackelberg game), who will choose their routes after the Stackelberg strategy s is revealed. This induces an assignment $(t(s), m(s))$ of the selfish flow at Nash equilibrium, and we assume that the assignment s of compliant agents is *not affected* after introducing the non-compliant flow on the network.

Since s may induce multiple Nash equilibria, we define the assignment $(t(s), m(s))$ to be the best such equilibrium (as defined in Section II-C). To characterize this Nash equilibrium, which we refer to as the *induced equilibrium* by strategy s , we note that the flow on link n is simply $s_n + t_n(s)$, and we have for all $n \in \text{Supp}(t_n(s))$ and for all $k \in \{1, \dots, N\}$:

$$l_n(s_n + t_n(s), m_n(s)) \leq l_k(s_k + t_k(s), m_k(s))$$

Equivalently, all links that are in the support of the selfish flow assignment $t(s)$ have a common latency l_0 in the induced equilibrium, and links that are not in the support have latency greater than or equal to l_0 .

This can be summarized in the following definition of a Stackelberg strategy. Let (N, r, α) denote an instance of the Stackelberg game played on the network instance (N, r) with compliance rate α .

Definition 3: Stackelberg Strategy

Consider an instance (N, r, α) of the Stackelberg routing game. A Stackelberg strategy for this instance is an assignment s of the compliant flow αr that is feasible for the instance $(N, \alpha r)$, and which induces *best Nash equilibrium* $(t(s), m(s))$ of the non-compliant flow, as defined in Section II-C, such that $s + t(s)$ is feasible for the instance (N, r) and $\exists l_0 > 0$ such that

$$\begin{aligned} \forall n \in \text{Supp}(t(s)), l_n(s_n + t_n(s), m_n(s)) &= l_0 \\ \forall n \notin \text{Supp}(t(s)), l_n(s_n, m_n(s)) &\geq l_0 \end{aligned}$$

This extends the definition usually used in the congestion games literature, see for example [10].

We will denote by $S(N, r, \alpha) \subset \mathcal{R}^N$ the set of Stackelberg strategies for network instance (N, r, α) .¹

We next show that the induced Nash equilibrium has one link in free-flow:

¹Note that a feasible flow assignment s of compliant flow may fail to induce a Nash equilibrium (t, m) and therefore is not considered to be a Stackelberg strategy.

Lemma 3: Characterization of the induced Nash Equilibrium

Let $s \in S(N, r, \alpha)$ be a Stackelberg strategy for the Stackelberg instance (N, r, α) , and $(t(s), m(s))$ its induced best Nash equilibrium. Then the last link in the support of $t(s)$ is in free-flow, i.e. $m_{\max \text{Supp}(t(s))} = 0$.

Proof: Note that $(t(s), m(s))$ is the best Nash equilibrium for the instance $(N, \alpha r)$ and latencies

$$\begin{aligned} \tilde{l}_n : [0, x_n^{\max} - s_n] \times \{0, 1\} &\longrightarrow \mathcal{R}_+ \\ (x_n, m_n) &\longmapsto l_n(s_n + x_n, m_n) \end{aligned}$$

Latencies \tilde{l}_n satisfy the assumptions of the horizontal queuing latencies class. Therefore, by Theorem 1, we immediately have the result. \blacksquare

B. Optimal Stackelberg strategies

In this section we solve for optimal Stackelberg strategies, i.e. Stackelberg strategies that induce Nash equilibria of minimal cost.

Definition 4: Optimal Stackelberg strategy

An optimal Stackelberg strategy s^* is a solution to the optimization problem

$$s^* = \arg \min_{s \in S(N, r, \alpha)} C(s + t(s), m(s))$$

Here $(t(s), m(s))$ is the equilibrium induced by s .

We also introduce a definition that will be useful in proving the main result.

Definition 5: At least i -congested link

Consider a network under feasible flow assignment (x, m) . Link n is said to be at least i -congested ($i \geq n + 1$) under assignment (x, m) if it is congested ($m_n = 1$) and its latency is at least a_i

$$l_n(x_n, m_n) \geq a_i$$

This is equivalent to $m_n = 1$ and $x_n \leq \hat{x}_n(i)$. If the above holds with equality, we say that the link is exactly i -congested.

Note that if $j \geq i \geq n + 1$, then if link n is at least j -congested under (x, m) , then it is also at least i -congested under (x, m) .

And if (t, m) is a single link free-flow equilibrium, and $i = \max \text{Supp}(t)$, then all links $n \in \{1, \dots, i - 1\}$ are exactly i -congested.

IV. COMPUTING THE OPTIMAL STACKELBERG STRATEGY

In this section, we show the following result: the optimal Stackelberg strategy can be computed in polynomial time for parallel networks with N links for the class of horizontal queuing congestion functions defined in II-A. This result

To see this, consider the following 2-link network such that link 1 is faster $a_1 < a_2$ and has larger capacity $x_1^{\max} > x_2^{\max}$. Now assume that the network is subject to flow demand $r = x_1^{\max} + \epsilon$ and most of the flow is compliant $\alpha r = x_1^{\max}$. Consider the following assignment: $s = (x_1^{\max}, 0)$.

Assuming that the assignment of compliant agents is not affected by introducing the non-compliant flow, we have for any assignment t of non-compliant flow, $t_1 = 0$ and $t_2 > 0$. Therefore t is not at Nash equilibrium since $\text{Supp}(t) = \{2\}$ and $l_2(s_2 + t_2, m_2) > l_1(s_1, 0)$ (non compliant agents are forced to use less efficient link 2).

contrasts with the class of non-decreasing latency functions where the optimal Stackelberg strategy is shown to be NP-hard to compute, see [10].

The optimal Stackelberg strategy in our case corresponds to:

- First computing the best Nash equilibrium of non-compliant agents alone, $(\bar{t}, \bar{m}) = \text{BNE}(N, (1-\alpha)r)$
- Then assigning the compliant flow by filling the remaining links (i.e. those that are not congested under (\bar{t}, \bar{m})), up to maximum capacity, starting with the faster links.

Intuitively, the best induced Nash equilibrium $(t(s), m(s))$ of any Stackelberg strategy s will be more congested than the best Nash equilibrium (\bar{t}, \bar{m}) of instance $(N, (1-\alpha)r)$. So if we can find a strategy \bar{s} that induces equilibrium (\bar{t}, \bar{m}) and that has minimal cost, then one expects this strategy to be optimal. Next, we detail this idea by defining a candidate Stackelberg strategy \bar{s} that will later be shown to be optimal.

A. A candidate Stackelberg strategy: Non-Compliant First

Let (\bar{t}, \bar{m}) denote the *best Nash equilibrium* for the instance $(N, (1-\alpha)r)$. Let $k = \max \text{Supp}(\bar{t})$ be the last link in the support of \bar{t} . Then we have from Equations (4) and (3), $\bar{m} = (1, \dots, 1, 0, \dots, 0)$ and

$$\bar{t} = \left(\hat{x}_1(k), \dots, \hat{x}_{k-1}(k), (1-\alpha)r - \sum_{n=1}^{k-1} \hat{x}_n(k), 0, \dots, 0 \right)$$

i.e. links $\{1, \dots, k-1\}$ are k -congested, and link k is in free-flow. Figure 2a shows best Nash equilibrium (\bar{t}, \bar{m}) on a sample network, where the latency in congestion $l_n(\cdot, 1)$ is taken to be affine for simplicity.

We now define Stackelberg strategy \bar{s} as the optimal assignment (i.e. of least cost) of compliant flow αr that induces equilibrium (\bar{t}, \bar{m}) . It is easy to see that \bar{s} is simply given by assigning the compliant flow to remaining links $\{k, k+1, \dots, N\}$ successively, each up to maximum capacity. The strategy \bar{s} will assign $x_k^{\max} - \bar{t}_k$ on link k , then x_{k+1}^{\max} on link $k+1$, x_{k+2}^{\max} on link $k+2$ and so on. Let $l = \min\{n | \alpha r - (\sum_{n=k}^{l-1} x_n^{\max} - \bar{t}_k) \geq 0\}$ be the least efficient link used by the Stackelberg assignment. Then \bar{s} is given by

$$\bar{s} = \left(0, \dots, 0, x_k^{\max} - \bar{t}_k, x_{k+1}^{\max}, \dots, x_{l-1}^{\max}, \alpha r - \left(\sum_{n=k}^{l-1} x_n^{\max} - \bar{t}_k \right), 0, \dots, 0 \right) \quad (5)$$

Equivalently, the total assignment $\bar{x} = \bar{s} + \bar{t}$ is given by

$$\bar{x} = \left(\hat{x}_1(k), \dots, \hat{x}_{k-1}(k), x_k^{\max}, x_{k+1}^{\max}, \dots, x_{l-1}^{\max}, r - \sum_{n=1}^{k-1} \hat{x}_n(k) - \sum_{n=k}^{l-1} x_n^{\max}, 0, \dots, 0 \right) \quad (6)$$

and the corresponding latencies are

$$\bar{l} = (a_k, \dots, a_k, a_{k+1}, \dots, a_l, 0, \dots, 0) \quad (7)$$

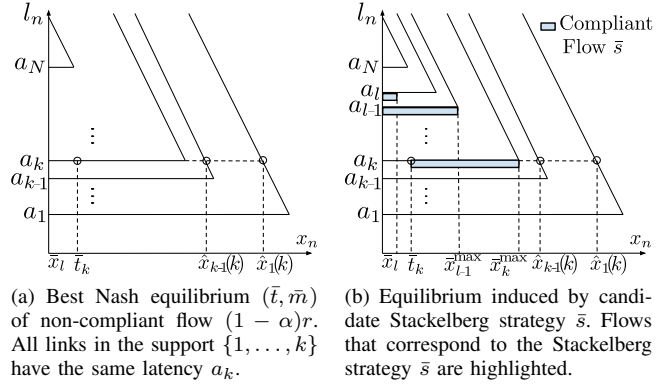


Fig. 2: Non-compliant first (NCF) strategy. In this example, latencies are taken to be affine in congestion for simplicity, the results hold for the general class of horizontal queuing latencies.

We will denote by $\text{NCF}(N, r, \alpha) = \bar{s}$ the *non-compliant first strategy* for the Stackelberg instance (N, r, α) .

Figure 2b shows the total flow $\bar{x}_n = \bar{s}_n + \bar{t}_n$ on each link. Links $\{1, \dots, k-1\}$ are exactly k -congested, links $\{k, \dots, l-1\}$ are in free-flow and at maximum capacity, and the remaining flow goes on link l .

In the next section we show that strategy \bar{s} is indeed an optimal Stackelberg strategy.

B. The Non-Compliant First strategy is optimal

Theorem 2: Optimal Stackelberg Strategy

$\bar{s} = \text{NCF}(N, r, \alpha)$ is an optimal Stackelberg strategy for the Stackelberg instance (N, r, α) .

Proof: Let $s \in S(N, r, \alpha)$ be a Stackelberg strategy for the Stackelberg instance (N, r, α) and (t, m) be the best induced Nash equilibrium for the non-compliant flow. We will show that $C(x, m) \geq C(\bar{x}, \bar{m})$, where $x = s + t$ and $\bar{x} = \bar{s} + \bar{t}$.

The proof proceeds as follows: we first show that links $\{1, \dots, l-1\}$ are more congested under assignment (x, m) than under (\bar{x}, \bar{m}) , in the following sense: these links have worse latency $l_n(x_n, m_n) \geq l_n(\bar{x}_n, \bar{m}_n)$, and hold less flow $x_n \leq \bar{x}_n$. Then we conclude by lower bounding the cost $C(x, m)$.

Let $k' = \max \text{Supp}(t)$ be the link with largest free-flow latency, in the support of the non-compliant flow. By Lemma 3, we have $m_{k'} = 0$, i.e. link k' is in free-flow under assignment $(x, m) = (s + t, m)$. We start by showing that $k' \geq k$ where $k = \max \text{Supp}(\bar{t})$.

Lemma 4: Comparing the supports of induced equilibria

The last link in the support of $t(s)$ has higher free-flow latency than the last link in the support of \bar{t} : $k' \geq k$.

Intuitively, since (\bar{t}, \bar{m}) is the *best Nash equilibrium* of the non-compliant agents when they are not sharing the network with any other flow, the cost of this assignment (\bar{t}, \bar{m}) is less than the cost of any equilibrium after introducing additional flow s .

Proof: First note that $(s + t(s), m)$ restricted to $\text{Supp}(t(s))$

is at Nash equilibrium. Then since link k' is in free-flow we have $l_{k'}(s_{k'} + t_{k'}(s), m_{k'}) = a_{k'}$, and since $k' \in \text{Supp}(t(s))$, we have by Definition 3 that any other link has higher latency. In particular, $\forall i \in \{1, \dots, k' - 1\}$, $l_i(s_i + t_i(s), m_i) \geq a_{k'}$, thus $s_i + t_i(s) \leq \hat{x}_i(k')$. Therefore we have $\sum_{n=1}^{k'} s_n + t_n(s) \leq \sum_{n=1}^{k'-1} \hat{x}_n(k') + x_{k'}^{\max}$. But $\sum_{n=1}^{k'} (s_n + t_n(s)) \geq \sum_{n \in \text{Supp}(t)} t_n(s) = (1 - \alpha)r$ since $\text{Supp}(t) \subset \{1, \dots, k'\}$. Therefore

$$(1 - \alpha)r \leq \sum_{n=1}^{k'-1} \hat{x}_n(k') + x_{k'}^{\max}$$

and by Lemma 2 applied to a network with k' links, this guarantees the existence of a single-link free-flow Nash Equilibrium for the network instance $(k', (1 - \alpha)r)$. Let $(\tilde{t}, \tilde{m}) \in \mathcal{R}^{k'} \times \{0, 1\}^{k'}$ be such an equilibrium. The cost of (\tilde{t}, \tilde{m}) is $(1 - \alpha)r l_0$ where $l_0 \leq a_{k'}$ is the free-flow latency of the last link in the support of \tilde{t} . Thus $C(\tilde{t}, \tilde{m}) \leq (1 - \alpha)r a_{k'}$.

Then $((\tilde{t}_1, \dots, \tilde{t}_{k'}, 0, \dots, 0), (\tilde{m}_1, \dots, \tilde{m}_{k'}, 0, \dots, 0)) \in \mathcal{R}^N \times \{0, 1\}^N$ is clearly a Nash equilibrium for the instance $((1 - \alpha)r, N)$, and has the same cost $C(\tilde{t}, \tilde{m}) \leq (1 - \alpha)r a_{k'}$. Since by definition (\tilde{t}, \tilde{m}) is the *best Nash equilibrium* for the instance $((1 - \alpha)r, N)$ and has cost $(1 - \alpha)r a_k$, we must have $(1 - \alpha)r a_k \leq (1 - \alpha)r a_{k'}$, i.e. $a_k \leq a_{k'}$. This completes the proof of the Lemma. \blacksquare

Using the lemma, we can now show that links $\{1, \dots, l - 1\}$ are more congested under assignment (x, m) than candidate assignment (\bar{x}, \bar{m}) .

Since $k' \in \text{Supp}(t)$, we have from Definition 3 of a Stackelberg strategy and its induced equilibrium, that the latency on k' is smaller than the latency on any other link. Thus $\forall n \in \{1, \dots, k' - 1\}$, $l_n(x_n, m_n) \geq l_{k'}(x_{k'}, m_{k'}) \geq a_{k'}$, i.e. $\forall n \in \{1, \dots, k' - 1\}$, n is at least k' -congested under assignment (x, m) . We also have by definition of the candidate assignment (\bar{x}, \bar{m}) and the resulting latencies given by Equation (7), $\forall n \in \{1, \dots, k - 1\}$, n is exactly k -congested under assignment (\bar{x}, \bar{m}) . Thus using the fact that $k' \geq k$, we have $\forall n \in \{1, \dots, k - 1\}$, $l_n(x_n, m_n) \geq a_{k'} \geq a_k = l_n(\bar{x}_n, \bar{m}_n)$, and $x_n \leq \hat{x}_n(k') \leq \hat{x}_n(k) = \bar{x}_n$, obtained by inverting the latency function $l_n(\cdot, m_n)$.

We have from Equation (6) that $\forall n \in \{k, \dots, l - 1\}$, n is in free-flow and at maximum capacity under assignment (\bar{x}, \bar{m}) (i.e. $\bar{x}_n = x_n^{\max}$ and $l_n(\bar{x}_n) = a_n$). Thus $\forall n \in \{k, \dots, l - 1\}$, $l_n(x_n, m_n) \geq a_n = l_n(\bar{x}_n, \bar{m}_n)$ and $x_n \leq x_n^{\max} = \bar{x}_n$. Therefore we have

$$l_n(x_n, m_n) \geq l_n(\bar{x}_n, \bar{m}_n) \quad \forall n \in \{1, \dots, l - 1\} \quad (8)$$

$$x_n \leq \bar{x}_n \quad \forall n \in \{1, \dots, l - 1\} \quad (9)$$

Note that $\forall n \in \{1, \dots, k\}$, $l_n(\bar{x}_n, \bar{m}_n) = a_k \leq a_l$, and $\forall n \in \{k, \dots, l - 1\}$, $l_n(\bar{x}_n, \bar{m}_n) = a_n \leq a_l$, thus we have

$$l_n(\bar{x}_n, \bar{m}_n) \leq a_l \quad \forall n \in \{1, \dots, l - 1\} \quad (10)$$

Also note that each link $n \in \{l, \dots, N\}$ has latency at least a_n (the latency on a link is always greater than the free-flow latency) and $a_n \geq a_l$, thus

$$l_n(x_n, m_n) \geq a_l \quad \forall n \in \{l, \dots, N\} \quad (11)$$

We can now lower-bound the cost of the assignment (x, m) where $x = s + t$ and (t, m) is the *best Nash equilibrium* induced by s . We have

$$\begin{aligned} C(x, m) &= \sum_{n=1}^N x_n l_n(x_n, m_n) \\ &= \sum_{n=1}^{l-1} x_n l_n(x_n, m_n) + \sum_{n=l}^N x_n l_n(x_n, m_n) \\ &\geq \sum_{n=1}^{l-1} x_n l_n(\bar{x}_n, \bar{m}_n) + \sum_{n=l}^N x_n a_l \end{aligned}$$

using (8) and (11). Then rearranging the terms we have

$$C(x, m) \geq \sum_{n=1}^{l-1} (x_n - \bar{x}_n) l_n(\bar{x}_n, \bar{m}_n) + \sum_{n=1}^{l-1} \bar{x}_n l_n(\bar{x}_n, \bar{m}_n) + \sum_{n=l}^N x_n a_l$$

Then by (9) and (10) we have $\forall n \in \{1, \dots, l - 1\}$, $x_n - \bar{x}_n \leq 0$ and $l_n(\bar{x}_n, \bar{m}_n) \leq a_l$, thus

$$\begin{aligned} C(x, m) &\geq \sum_{n=1}^{l-1} (x_n - \bar{x}_n) a_l + \sum_{n=1}^{l-1} \bar{x}_n l_n(\bar{x}_n, \bar{m}_n) + \sum_{n=l}^N x_n a_l \\ &= a_l \left(\sum_{n=1}^N x_n - \sum_{n=1}^{l-1} \bar{x}_n \right) + \sum_{n=1}^{l-1} \bar{x}_n l_n(\bar{x}_n, \bar{m}_n) \\ &= a_l \left(r - \sum_{n=1}^{l-1} \bar{x}_n \right) + \sum_{n=1}^{l-1} \bar{x}_n l_n(\bar{x}_n, \bar{m}_n) \end{aligned}$$

But $a_l \left(r - \sum_{n=1}^{l-1} \bar{x}_n \right) = \bar{x}_l l_l(\bar{x}_l, \bar{m}_l)$ since $\text{Supp}(\bar{x}) = \{1, \dots, l\}$ and $l_l(\bar{x}_l, \bar{m}_l) = a_l$. Therefore

$$\begin{aligned} C(x, m) &\geq \bar{x}_l l_l(\bar{x}_l, \bar{m}_l) + \sum_{n=1}^{l-1} \bar{x}_n l_n(\bar{x}_n, \bar{m}_n) \\ &= C(\bar{x}, \bar{m}) \end{aligned}$$

Therefore the NCF strategy is an optimal Stackelberg strategy, and it can be computed in polynomial time since it is generated in linear time after computing the best Nash equilibrium BNE $(N, (1 - \alpha)r)$, which was shown to be quadratic in N .

Finally, we note that the NCF strategy is, in general, not the unique optimal Stackelberg strategy, but the set of optimal Stackelberg strategies can be easily generated from the NCF strategy. We do not describe this procedure for brevity.

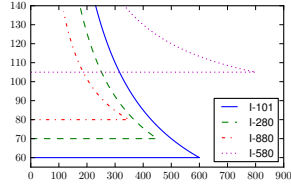
V. NUMERICAL RESULTS

A. Optimal Stackelberg routing on an example network

In this section, we apply the previous results to a scenario of freeway traffic from the San Francisco Bay Area. Four parallel highways are chosen starting in San Francisco and ending in San Jose: I-101, I-280, I-880 and I-580 (shown in Figure 3a). We analyze the inefficiency of Nash equilibria



(a) Map of parallel highway network showing four parallel highways connecting San Francisco to San Jose.



(b) Latency (minutes) vs. demand (cars/minute) for parallel highway routes.

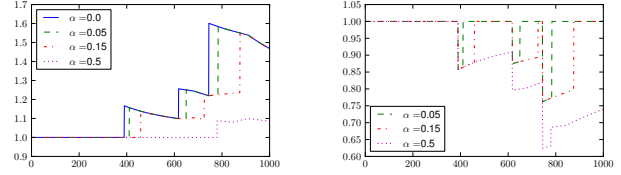
Fig. 3: Example highway network connecting San Francisco to San Jose.

due to selfish routing and lack of coordination, and show how optimal Stackelberg routing strategies (non-compliant first strategy) can improve these conditions. We first use price of stability [1] to measure the improvement in performance achieved by optimal Stackelberg routing. Price of Stability is defined as the ratio between the cost of the best Nash equilibrium, and the cost of the social optimum: $POS(N, r) = \frac{C(\text{BNE}(N, r))}{C^*}$ where $C^* = \min_{(x, m)} C(x, m)$.

Figure 3b shows the latency functions for the highway network, which we derived assuming a triangular fundamental diagram for each highway stretch (the resulting latency functions are decreasing harmonic in congestion. See [6] for a detailed derivation of latency functions from a triangular fundamental diagram). Under free-flow conditions, I-101 is the fastest route available between San Francisco and San Jose. But when I-101 becomes congested, other routes represent viable alternatives. To analyze how congestion increases with demand, we compute total network latency profiles (Equation (1)) for optimal Stackelberg strategies as a function of demand. In addition, to show how congestion improves as a function of fraction of compliance α , we compute network latency profiles over a range of compliance rates.

The numerical results are summarized in Figure 4. The price of stability plot in Figure 4a shows that even with a small compliance rate, Stackelberg routing can decongest a given link n for a fixed flow demand, when the Nash equilibrium is slightly above maximum capacity on link n . This shows the significant benefits of Stackelberg routing, especially around the critical regions of flow demand where the support of the best Nash equilibrium changes.

We also note that for a fixed compliance rate α , Stackelberg routing can delay the congestion of a particular link n , i.e. increase the critical flow demand $r^{(n)}$ above which link n becomes congested (formally, $r^{(n)} = \inf_r \{r | n \text{ is congested under BNE}(N, r)\}$). Let $r^{(n, \alpha)}$ denote the critical flow for the Stackelberg instance (N, r, α) , i.e. $r^{(n, \alpha)} = \inf_r \{r | n \text{ is congested under } (t(\bar{s}), m(\bar{s})), \bar{s} = \text{NCF}(N, r, \alpha)\}$. Then the delay range (the increase in critical flow demand $r^{(n, \alpha)} - r^{(n)}$) is proportional to the compliance rate. Observing the behavior around the 600-800 cars/minute, it takes about 30 extra cars/minute to congest the second



(a) Price of stability vs. demand (cars/minute) for different compliance rates.

(b) Value of altruism vs. demand (cars/minute) for different compliance rates.

Fig. 4: Results for network efficiency on parallel highway link example.

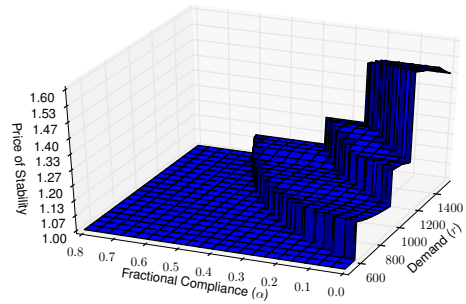


Fig. 5: Price of stability profile

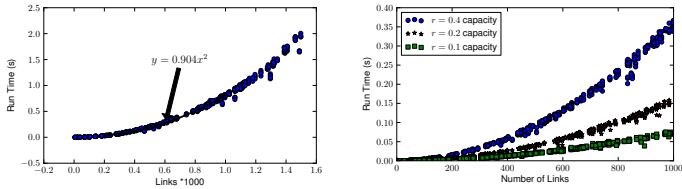
link with 5% compliance, about 100 cars/minute for 15% compliance, while for 50% compliance, the second link appears to only congest near network capacity.

While price of stability describes how inefficient the best Nash equilibrium is with respect to the social optimal strategy, another metric, *value of altruism* [2], was introduced to show how performant a particular Stackelberg strategy is with respect to the best Nash equilibrium (0% compliance). Value of altruism is defined by:

$$VOA(N, r, \alpha) = \frac{\min_{s \in S(N, r, \alpha)} C(s + t(s), m(s))}{C(\text{BNE}(N, r))}$$

Value of altruism for the example network is shown in Figure (4b). The figure illustrates the “delaying” effect of optimal Stackelberg routing. A small compliance fraction (5-15% on this network) can decongest a given link n when the demand is around the critical demand for that link ($r^{(n)} \leq r < r^{(n, \alpha)}$), but the link will be congested again for higher values of demand ($r \geq r^{(n, \alpha)}$) resulting in a value of altruism of 1.

The importance of efficient computation of optimal Stackelberg strategies can be seen by the complete demand-compliance fraction profile generated in Figure 5. If the amount of compliance is viewed as a cost to some central controller, then the tiers visible in Figure 5 (take $r = 1000$, $.3 \leq \alpha \leq .5$) can be seen as regions of potential waste. If



(a) Quadratic fit to computation time-network size relationship (b) Increase of computational time with demand.

Fig. 6: Computation time of optimal Stackelberg strategy relative to network size.

a central controller can predict a demand in a range around 1000 cars/minute and a maximum compliance fraction of .5, then the controller can reference the two-dimensional profile and reduce the compliance fraction to a region around .4, since any compliance above 30% and less than 50% does not improve the performance of the network.

B. Scaling of optimal Stackelberg strategy algorithm on size of network

To illustrate the performance of the algorithm as the size of the network scales up, the computation time of the optimal Stackelberg strategy was measured for 500 randomly generated networks. The number of links in a network was chosen between 3 and 1500 and the latency functions of each link correspond to randomly generated triangular fundamental diagrams. The compliance rate was arbitrarily chosen to be 40% and the demand was chosen to be 70% that of the maximum capacity of the network at best Nash equilibrium. The results are shown in Figure 6a.

As shown in Section III-B, the worst-case complexity of computing optimal Stackelberg assignments is quadratic in the size of the network, which is verified experimentally as shown in Figure 6a.

Figure 6b shows that the computation time of the optimal Stackelberg strategy increases as the demand increases. This is due to the fact the best Nash equilibrium is computed using sequential search: the algorithm tests if a Nash equilibrium exists for a particular support, and if it fails to find such an equilibrium, increases the size of the support. As the demand increases, the algorithm will have to check for larger supports, which explains the increase in computation time.

VI. DISCUSSION AND OPEN PROBLEMS

In order to address the inefficiency of Nash equilibria on horizontal queuing networks, we considered the Stackelberg routing game where a central coordinator has control over a fraction α of compliant agents. We proved that for the class of horizontal queuing congestion latencies introduced in [6], the *non-compliant first* (NCF) strategy is optimal, and that it can be computed in quadratic time in the size of the network. We illustrated these results using a benchmark network for which we computed the decrease in inefficiency that can be achieved using optimal Stackelberg routing. This example showed that when the demand is near critical flows $r^{(n)}$,

optimal Stackelberg routing can achieve a significant increase in efficiency even for small values of compliance rate α .

These results show that careful routing of a small compliant population can significantly improve the efficiency of the network. It is also worth noting that these results indicate that for specific demand and compliance ranges, Stackelberg routing can be completely ineffective. Therefore identifying the ranges for which optimal Stackelberg routing does improve the efficiency of the network is crucial for effective planning and control.

This work offers several directions of future research: the work presented here only considers parallel networks under static conditions (constant flow demand r , and static equilibria): one question is how one may dynamically steer the system from one equilibrium to a better one. For example, consider the case in which the agents are stuck in a congested equilibrium, and assume a coordinator has control over a fraction of the flow. Can the coordinator steer the system to a single link free-flow equilibrium? And what is the minimal compliance rate needed to achieve this?

Another question is how robust are the NCF strategy results? Do they hold for general network topologies? The extension of our results to general network topologies is still an open problem.

REFERENCES

- [1] E. Anshelevich, A. Dasgupta, J. Kleinberg, E. Tardos, T. Wexler, and T. Roughgarden. The Price of Stability for Network Design with Fair Cost Allocation. *45th Annual IEEE Symposium on Foundations of Computer Science*, pages 295–304, 2004.
- [2] A. Aswani and C. Tomlin. Game-theoretic routing of GPS-assisted vehicles for energy efficiency. In *American Control Conference (ACC), 2011*, pages 3375–3380. IEEE, 2011.
- [3] Caltrans. US 101 South, corridor system management plan, 2010.
- [4] C. F. Daganzo. The cell transmission model: A dynamic representation of highway traffic consistent with the hydrodynamic theory. *Transportation Research Part B: Methodological*, 28(4):269–287, 1994.
- [5] D. Joksimovic, M. C. J. Bliemer, and P. H. L. Bovy. Different policy objectives of the road-pricing problem: a game-theoretic approach. *Pricing in road transport: a multi-disciplinary perspective*, pages 151–169.
- [6] W. Krichene, J. Reilly, S. Amin, and Bayen A. M. Nash equilibria and stackelberg routing on horizontal queueing networks, <https://sites.google.com/site/wkrichene/research/stackelberg>.
- [7] M. J. Lighthill and G. B. Whitham. On kinematic waves. II. A theory of traffic flow on long crowded roads. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 229(1178):317, 1955.
- [8] H. K. Lo and W. Y. Szeto. A cell-based variational inequality formulation of the dynamic user optimal assignment problem. *Transportation Research Part B: Methodological*, 36(5):421–443, 2002.
- [9] C. Papadimitriou and G. Valiant. A new look at selfish routing. *Innovations in Computer Science (ICS)*, 2010.
- [10] T. Roughgarden. Stackelberg scheduling strategies. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pages 104–113. ACM, 2001.
- [11] T. Roughgarden and E. Tardos. How bad is selfish routing? *Journal of the ACM (JACM)*, 49(2):236–259, 2002.
- [12] D. Schmeidler. Equilibrium points of nonatomic games. *Journal of Statistical Physics*, 7(4):295–300, 1973.
- [13] A. K. Ziliaskopoulos. A linear programming model for the single destination system optimum dynamic traffic assignment problem. *Transportation science*, 34(1):37, 2000.