

On Stackelberg Routing on Parallel Networks with Horizontal Queues

Walid Krichene

Jack Reilly

Saurabh Amin

Alexandre Bayen

Abstract—In order to address inefficiencies of Nash equilibria for congestion networks with horizontal queues, we study the Stackelberg routing game on parallel networks: assuming a coordinator has control over a fraction of the flow, and that the remaining players respond selfishly, what is an optimal Stackelberg strategy of the coordinator, i.e. a strategy that minimizes the cost of the induced equilibrium?

We study Stackelberg routing for a new class of latency functions, which models congestion on horizontal queues. We introduce a candidate strategy, the *non-compliant first* strategy, and prove it to be optimal. Then we apply these results by modeling a transportation network in which a coordinator can choose the routes of a subset of the drivers, while the rest of the drivers choose their routes selfishly.

I. INTRODUCTION

A. Congestion games and Stackelberg routing

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

In order to address this inefficiency, many tools have been proposed, including congestion pricing [10], capacity allocation [6] and Stackelberg routing [12], [2], [15], [5]. In the Stackelberg routing game, a coordinator (leader) routes a fraction of the flow, then the remaining players (followers) respond selfishly [2], [12]. The objective of the coordinator is to minimize a system-wide cost function.

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 for this class of latency functions is NP-hard in the number of links. This led to considering polynomial time approximate strategies such as Largest Latency First and Scale [12], 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 congestion on networks with horizontal queues, such as transportation networks [4], [8]. A new class of latency functions,

the HQSF latency class (horizontal queues, singled-valued in free-flow) is introduced in [7] to model congestion on horizontal queues. We study Stackelberg routing for this new class of latency.

B. Motivating application: optimal routing of a subset of drivers and 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 vehicular navigation systems. These devices are used to provide the driver with relevant traffic information. They may also be used by a coordinator (a traffic control center) to provide routing advice that can improve the overall efficiency of the network. Naturally, when providing routing advice, the coordinator needs to take into account the possible response of other drivers to the resulting change in traffic conditions, hence the importance of the Stackelberg routing framework, in which a fraction of the population of drivers is assumed to be *compliant* to routing suggestions, and the rest of the drivers (*non-compliant*) are assumed to respond selfishly.

We restrict our present work to the case of parallel networks. This simple network topology is of practical importance to traffic planners [3], since numerous transportation networks can be modeled as parallel highways connecting two highly populated areas. We consider one such example in the numerical results section.

C. Contributions

We study the Stackelberg routing game on parallel networks for the HQSF latency class. We define *non-compliant first* (NCF) strategy and prove that it is an optimal Stackelberg strategy. This shows in particular that in this setting, optimal Stackelberg strategies can be computed in quadratic time in the size of the network. This result contrasts with the class of non-decreasing latency functions, for which computing the optimal Stackelberg strategy is NP-hard [12]. We then apply these results to model a real transportation network, quantify the decrease in inefficiency achieved by the NCF strategy, and identify ranges of the flow demand and compliance rates where optimal Stackelberg routing are most efficient.

D. Organization of the Article

We start by defining the Stackelberg routing game and the HQSF latency class in Section II, then review some properties of Nash equilibria for the routing game. The main results on optimal Stackelberg routing are presented in Sections III: we define the NCF strategy and prove that

Walid Krichene is with the department of Electrical Engineering and Computer Sciences, University of California at Berkeley. walid@eecs.berkeley.edu

Jack Reilly is with the department of Civil and Environmental Engineering, University of California at Berkeley. jackdreilly@berkeley.edu

Saurabh Amin is with the department of Civil and Environmental Engineering, Massachusetts Institute of Technology. amins@mit.edu

Alexandre Bayen is with the department of Electrical Engineering and Computer Sciences, and the department of Civil and Environmental Engineering, University of California at Berkeley. bayen@berkeley.edu

it is optimal. This is followed by an example network in Section IV that illustrate the effects of optimal Stackelberg routing.

II. PRELIMINARIES

A. The Model: Stackelberg routing game and the HQSF latency class

We consider 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. The state of the network is given by

- a feasible flow assignment vector $\mathbf{x} \in \mathbb{R}_+^N$ such that $\sum_{n=1}^N x_n = r$ where x_n is the flow on link n ,
- a congestion state vector $\mathbf{m} \in \{\mathbf{0}, 1\}^N$ where $m_n = \mathbf{0}$ if the link is in *free-flow* (density is below critical density) and $m_n = 1$ if the link is *congested* (density is above critical density).

In the routing game, a population of non-atomic players [14] share the network, and every non-atomic player chooses a route in order to minimize their individual latency [13]. If a player chooses link n , their latency is given by $\ell_n(x_n, m_n)$, where ℓ_n is a HQSF latency function [7] (horizontal queues, single-valued in free-flow). The latency function

$$\ell_n : D_n \rightarrow \mathbb{R}_+ \quad (1)$$

$$(x_n, m_n) \mapsto \ell_n(x_n, m_n)$$

is defined on $D_n = [\mathbf{0}, x_n^{\max}] \times \{\mathbf{0}\} \cup (\mathbf{0}, x_n^{\max}) \times \{1\}$, and satisfies the following properties:

- The latency in free-flow, $\ell_n(\cdot, \mathbf{0})$, is constant. We will denote its value by a_n and call it the *free-flow latency*.
- The latency in congestion, $x_n \mapsto \ell_n(x_n, 1)$, is decreasing from $(\mathbf{0}, x_n^{\max})$ onto $(a_n, +\infty)$.
- $\lim_{x_n \rightarrow x_n^{\max}} \ell_n(x_n, 1) = \ell_n(x_n^{\max}, \mathbf{0}) = a_n$.

As detailed in [7], in the case of horizontal queues, the latency is not uniquely determined by the flow. It depends on whether the link is in free-flow or in congestion, hence the dependency on m_n . Intuitively, a given flow x_n corresponds to two different configurations: either few drivers moving fast (the density is low and the link is in free-flow), in which case the latency is low, or many drivers moving slowly (the density is high and the link is congested), in which case the latency is high.

An example of HQSF latency function in this class is given in Fig. 1. We observe, as an immediate consequence of these assumptions, that the latency in congestion is always greater than the free-flow latency: $\forall x_n \in (\mathbf{0}, x_n^{\max}), \ell_n(x_n, 1) > a_n$.

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

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

We denote by (N, r) an instance of the routing game on a network with N links under flow demand r . Pure Nash equilibria of the game (which we will simply refer to as Nash

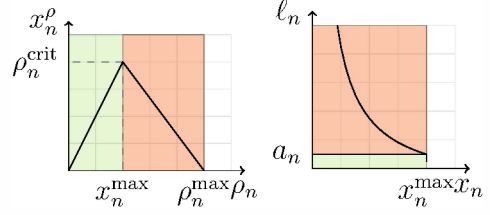


Fig. 1: Example of a triangular flux function, and the resulting latency function, as defined in section II-A.

equilibria) are assignments (\mathbf{x}, \mathbf{m}) such that every player cannot improve their latency by switching to a different link.

Definition 1: Nash equilibrium of the routing game

A feasible assignment $(\mathbf{x}, \mathbf{m}) \in \mathbb{R}_+^N \times \{\mathbf{0}, 1\}^N$ is a Nash equilibrium of the routing game instance (N, r) if $\forall k \in \text{supp}(\mathbf{x}), \forall n \in \{1, \dots, N\}, \ell_k(x_k, m_k) \leq \ell_n(x_n, m_n)$.

Here $\text{supp}(\mathbf{x})$ denotes the support of \mathbf{x} , i.e. the set of links n such that $x_n > \mathbf{0}$. As a consequence of this definition, all links in the support of \mathbf{x} have the same latency ℓ_{\bullet} , and links that are not in the support have latency greater than or equal to ℓ_{\bullet} . We will denote by $\text{NE}(N, r)$ the set of Nash equilibria of the instance (N, r) .

While a Nash equilibrium achieves minimal individual latencies, it does not minimize, in general, the *system cost* or *total cost* defined as follows:

Definition 2: Total cost

The total cost of an assignment (\mathbf{x}, \mathbf{m}) is the total latency experienced by all players

$$C(\mathbf{x}, \mathbf{m}) = \sum_{n=1}^N x_n \ell_n(x_n, m_n) \quad (2)$$

In order to study the inefficiency of Nash equilibria, and the improvement of performance that the Stackelberg routing game can achieve, we focus our attention on *best Nash equilibria* and *price of stability* as a measure of their inefficiency (see for example [1]). The *best Nash equilibrium* (BNE) is defined to be the Nash equilibrium of least total latency $\text{BNE}(N, r) = \underset{(\mathbf{x}, \mathbf{m}) \in \text{NE}(N, r)}{\text{argmin}} C(\mathbf{x}, \mathbf{m})$. It is shown in [7] that the minimizer is unique.

In the Stackelberg routing game, a coordinator (a central authority) is assumed to have control over a positive fraction α of the total flow demand r . We call α the *compliance rate*. The coordinator wants to route the *compliant flow* αr in a way that minimizes the system cost, while anticipating the response of the rest of the players, assumed to choose their routes selfishly *after* the strategy of the coordinator is revealed. We will refer to the flow of selfish players $(1 - \alpha)r$ as the *non-compliant flow*.

More precisely, the game is played as follows:

- First, the coordinator (the leader) chooses a *Stackelberg strategy*, i.e. an assignment $\mathbf{s} \in \mathbb{R}_+^N$ of the compliant flow. The assignment \mathbf{s} needs to be feasible for the instance $(N, \alpha r)$, i.e. $\sum_{n=1}^N s_n = \alpha r$.
- Then, the Stackelberg strategy \mathbf{s} of the leader is revealed, and the non-compliant players (followers)

choose their routes selfishly and form the best Nash equilibrium $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$, induced by strategy \mathbf{s}^1 . By definition, the induced equilibrium $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ satisfies

$$\begin{aligned} \forall k \in \text{supp}(\mathbf{t}(\mathbf{s})), \forall n \in \{1, \dots, N\}, \\ \ell_k(s_k + t_k(\mathbf{s}), m_k(\mathbf{s})) \leq \ell_n(s_n + t_n(\mathbf{s}), m_n(\mathbf{s})) \end{aligned} \quad (3)$$

In other words, $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ is the best Nash equilibrium for the routing game $(N, (1 - \alpha)r)$ with latency functions

$$\begin{aligned} \tilde{\ell}_n : \quad \tilde{D}_n \rightarrow \mathbb{R}_+ \\ (x_n, m_n) \mapsto \ell_n(s_n + x_n, m_n) \end{aligned} \quad (4)$$

where $\tilde{D}_n \triangleq [0, \tilde{x}_n^{\max}] \times \{0\} \cup (0, \tilde{x}_n^{\max}) \times \{1\}$ and $\tilde{x}_n^{\max} \triangleq x_n^{\max} - s_n$ (note that latencies $\tilde{\ell}_n$ also satisfy the assumptions of the HQSF latency class).

The total flow on the network is $\mathbf{s} + \mathbf{t}(\mathbf{s})$, thus the total cost is $C(\mathbf{s} + \mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$. Note that a Stackelberg strategy \mathbf{s} may induce multiple Nash equilibria in general. However, we define the assignment $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ to be the (unique) best such equilibrium.

We use the following notation:

- (N, r, α) is an instance of the Stackelberg routing game played on a parallel network with N links under flow demand r with compliance rate α . The routing game (N, r) is a special case of the Stackelberg routing game with $\alpha = 0$.
- $S(N, r, \alpha) \subset \mathbb{R}_+^N$ is the set of Stackelberg strategies for the Stackelberg instance (N, r, α) .
- $S^*(N, r, \alpha)$ is the set of optimal Stackelberg strategies defined as

$$S^*(N, r, \alpha) = \arg \min_{\mathbf{s} \in S(N, r, \alpha)} C(\mathbf{s} + \mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s})) \quad (5)$$

B. Properties of Nash Equilibria

We briefly review some properties of Nash equilibria for the routing game. For a more detailed discussion and proofs, we refer the reader to [7].

Consider a routing game instance (N, r) , and partition the set of Nash equilibria 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.

The following lemma characterizes the congestion state vectors for single-link-free-flow equilibria:

Lemma 1: Congestion states for single-link-free-flow equilibria

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

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

mines 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 $\ell_k(x_k, 0) = a_k$. Therefore every link in the support has the same latency, in particular $\forall n \in \{1, \dots, k-1\}$, $\ell_n(x_n, 1) = a_k$. The unique flow that satisfies this equality is referred to as *congestion flow*. More precisely,

Definition 1: 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

$$\ell_n(\hat{x}_n(k), 1) = a_k \quad (6)$$

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

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

$$\mathbf{m}^k \triangleq (1, \dots, 1, \overset{k}{0}, \dots, 0) \quad (7)$$

$$\mathbf{x}^{k,r} \triangleq \left(\hat{x}_1(k), \dots, \hat{x}_{k-1}(k), r - \sum_{n=1}^{k-1} \hat{x}_n(k), 0, \dots, 0 \right) \quad (8)$$

Under such an assignment $(\mathbf{x}^{k,r}, \mathbf{m}^k)$, we say that links $\{1, \dots, k-1\}$ are k -congested. An example of single-link-free-flow equilibrium is shown in Fig. 2. It is also shown in [7] that when the set of equilibria $\text{NE}(N, r)$ is nonempty, it contains at least one single-link-free-flow equilibrium.

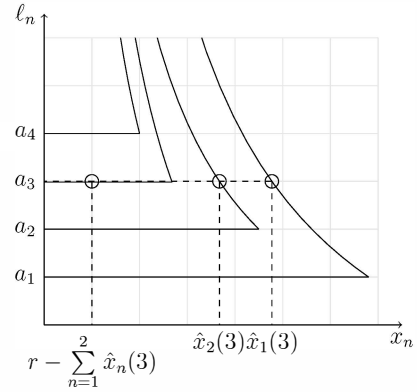


Fig. 2: Example of a single-link-free-flow equilibrium $(\mathbf{x}^{3,r}, \mathbf{m}^3)$. Links 1 and 2 are exactly 3-congested, link 3 is in free-flow, and link 4 is not in the support.

Lemma 2: Best Nash Equilibria [7]

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

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

As a consequence, the best Nash equilibrium can be computed by simply enumerating all candidate single-link-free-flow equilibria $(\mathbf{x}^{k,r}, \mathbf{m}^k)$, starting from the smallest support

($k = \mathbf{0}$). There are N such candidate equilibria, corresponding to the congestion states $(\mathbf{0}, \dots, \mathbf{0})$ up to $(1, \dots, 1, \mathbf{0})$, and each candidate equilibrium is a vector in \mathbb{R}_N that can be computed in $O(N)$, which corresponds to a worst-case time complexity of $O(N^2)$. This proves that computing the BNE is quadratic in the size of the network.

Corollary 1: Let $\mathbf{s} \in \mathbf{S}(N, r, \alpha)$ be a Stackelberg strategy for the Stackelberg instance (N, r, α) , and $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ its induced best Nash equilibrium. Then the last link in the support of $\mathbf{t}(\mathbf{s})$ is in free-flow.

Proof: This follows from Lemma 2 and the observation that $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ is the best Nash equilibrium for the instance $(N, (1-\alpha)r)$ and latencies $\tilde{\ell}_n$ given by equation (4). ■

III. OPTIMAL STACKELBERG STRATEGIES

In this section we study optimal Stackelberg strategies, i.e. Stackelberg strategies that induce Nash equilibria of minimal cost. We first define the non-compliant first (NCF) strategy, and prove that it is an optimal Stackelberg strategy, and that it can be computed in quadratic time. This result contrasts with the class of non-decreasing latency functions where the optimal Stackelberg strategy is shown to be NP-hard to compute, see [12]. The NCF strategy corresponds to:

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

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

A. A candidate Stackelberg strategy: non-compliant first

Let $(\bar{\mathbf{t}}, \bar{\mathbf{m}})$ denote the best Nash equilibrium for the routing instance $(N, (1-\alpha)r)$. Let $\bar{k} = \max \text{supp}(\bar{\mathbf{t}})$ be the last link in the support of $\bar{\mathbf{t}}$. Then we have from Equations (7) and (8),

$$\bar{\mathbf{m}} = (1, \dots, 1, \mathbf{0}, \dots, \mathbf{0}) \quad (9)$$

$$\bar{\mathbf{t}} = \left(\hat{x}_1(\bar{k}), \dots, \hat{x}_{\bar{k}-1}(\bar{k}), r - \sum_{n=1}^{\bar{k}-1} \hat{x}_n(\bar{k}), \mathbf{0}, \dots, \mathbf{0} \right) \quad (10)$$

i.e. links $\{1, \dots, \bar{k}-1\}$ are \bar{k} -congested, and link \bar{k} is in free-flow (see Fig. 3).

We now define Stackelberg strategy $\bar{\mathbf{s}}$ as the optimal assignment (i.e. of least cost) of compliant flow αr that induces equilibrium $(\bar{\mathbf{t}}, \bar{\mathbf{m}})$. It is easy to see that $\bar{\mathbf{s}}$ is simply given by assigning the compliant flow to remaining links $\{\bar{k}, \bar{k}+1, \dots, N\}$ successively, each up to maximum capacity. The strategy $\bar{\mathbf{s}}$ will assign $x_{\bar{k}}^{\max} - \bar{t}_{\bar{k}}$ on link \bar{k} , then $x_{\bar{k}+1}^{\max}$ on link $\bar{k}+1$, $x_{\bar{k}+2}^{\max}$ on link $\bar{k}+2$ and so on.

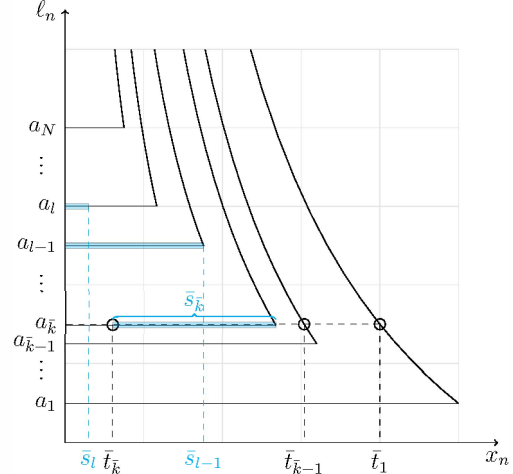


Fig. 3: Non-compliant first (NCF) strategy $\bar{\mathbf{s}}$ and its induced equilibrium. Circles show the best Nash equilibrium $(\bar{\mathbf{t}}, \bar{\mathbf{m}})$ of the non-compliant flow $(1-\alpha)r$: link \bar{k} is in free-flow, and links $\{1, \dots, \bar{k}-1\}$ are congested. The Stackelberg strategy $\bar{\mathbf{s}} = \text{NCF}(N, r, \alpha)$ is highlighted.

Let $l = \min\{i | \alpha r - (\sum_{n=\bar{k}}^{i-1} x_n^{\max} - \bar{t}_{\bar{k}}) \geq \mathbf{0}\}$ be the least efficient link used by the Stackelberg assignment. Then $\bar{\mathbf{s}}$ is given by

$$\bar{\mathbf{s}} = (0, \dots, 0, x_{\bar{k}}^{\max} - \bar{t}_{\bar{k}}, x_{\bar{k}+1}^{\max}, \dots, x_{l-1}^{\max}, \alpha r - (\sum_{n=\bar{k}}^{l-1} x_n^{\max} - \bar{t}_{\bar{k}}), 0, \dots, 0) \quad (11)$$

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

$$\bar{\mathbf{x}} = (\hat{x}_1(\bar{k}), \dots, \hat{x}_{\bar{k}-1}(\bar{k}), x_{\bar{k}}^{\max}, x_{\bar{k}+1}^{\max}, \dots, x_{l-1}^{\max}, r - \sum_{n=1}^{\bar{k}-1} \hat{x}_n(\bar{k}) - \sum_{n=\bar{k}}^{l-1} x_n^{\max}, \mathbf{0}, \dots, \mathbf{0}) \quad (12)$$

and the corresponding latencies are

$$(\mathbf{0}_{\bar{k}}, \dots, \mathbf{0}_{\bar{k}}, \mathbf{0}_{\bar{k}+1}, \dots, \mathbf{0}_l, \mathbf{0}, \dots, \mathbf{0}) \quad (13)$$

We will denote by $\text{NCF}(N, r, \alpha) = \bar{\mathbf{s}}$ the non-compliant first strategy for the Stackelberg instance (N, r, α) . Fig. 3 shows the total flow $\bar{x}_n = \bar{s}_n + \bar{t}_n$ on each link. Links $\{1, \dots, \bar{k}-1\}$ are \bar{k} -congested, links $\{\bar{k}, \dots, l-1\}$ are in free-flow and at maximum capacity, and the remaining flow goes on link l .

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

B. The Non-Compliant First strategy is optimal

Theorem 1: The NCF strategy is optimal

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

Proof: Let $\mathbf{s} \in \mathbf{S}(N, r, \alpha)$ be a Stackelberg strategy for the Stackelberg instance (N, r, α) and $(\mathbf{t}(\mathbf{s}), \mathbf{m}(\mathbf{s}))$ be its induced best Nash equilibrium for

the non-compliant flow. We will show that $C(\mathbf{x}, \mathbf{m}) \geq C(\bar{\mathbf{x}}, \bar{\mathbf{m}})$, where $\mathbf{x} = \mathbf{s} + \mathbf{t}$ and $\bar{\mathbf{x}} = \bar{\mathbf{s}} + \bar{\mathbf{t}}$.

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

Let $k = \max \text{supp}(\mathbf{t})$ be the link with largest free-flow latency, in the support of the non-compliant flow. By corollary 1, we have $m_k = 0$, i.e. link k is in free-flow under assignment $(\mathbf{x}, \mathbf{m}) = (\mathbf{s} + \mathbf{t}, \mathbf{m})$. It can be shown that $k \geq \bar{k}$ where $\bar{k} = \max \text{supp}(\bar{\mathbf{t}})$ (intuitively, since $(\bar{\mathbf{t}}, \bar{\mathbf{m}})$ is the *best Nash equilibrium* of the non-compliant players when they are not sharing the network with any other flow, the cost of this assignment $(\bar{\mathbf{t}}, \bar{\mathbf{m}})$ is less than the cost of any equilibrium after introducing additional flow \mathbf{s} . For a detailed proof, see [7]).

Since $k \in \text{supp}(\mathbf{t})$, we have by Equation (3) defining the induced equilibrium, $\forall n \in \{1, \dots, k-1\}$, $\ell_n(x_n, m_n) \geq \ell_k(x_k, m_k) \geq a_k$ (the latency on k is less than or equal to the latency on any other link). We also have by definition of the candidate assignment $(\bar{\mathbf{x}}, \bar{\mathbf{m}})$ and the resulting latencies given by Equation (13), $\forall n \in \{1, \dots, \bar{k}-1\}$, n is exactly \bar{k} -congested under assignment $(\bar{\mathbf{x}}, \bar{\mathbf{m}})$. Thus using the fact that $k \geq \bar{k}$, we have $\forall n \in \{1, \dots, \bar{k}-1\}$, $\ell_n(x_n, m_n) \geq a_k \geq a_{\bar{k}} = \ell_n(\bar{x}_n, \bar{m}_n)$, and $x_n \leq \hat{x}_n(k) \leq \hat{x}_n(\bar{k}) = \bar{x}_n$, obtained by inverting the latency function $\ell_n(\cdot, m_n)$.

We have from Equation (12) that $\forall n \in \{\bar{k}, \dots, l-1\}$, n is in free-flow and at maximum capacity under assignment $(\bar{\mathbf{x}}, \bar{\mathbf{m}})$ (i.e. $\bar{x}_n = x_n^{\max}$ and $\ell_n(\bar{x}_n) = a_n$). Thus $\forall n \in \{\bar{k}, \dots, l-1\}$, $\ell_n(x_n, m_n) \geq a_n = \ell_n(\bar{x}_n, \bar{m}_n)$ and $x_n \leq x_n^{\max} = \bar{x}_n$. Therefore we have

$$\ell_n(x_n, m_n) \geq \ell_n(\bar{x}_n, \bar{m}_n) \quad \forall n \in \{1, \dots, l-1\} \quad (14)$$

$$x_n \leq \bar{x}_n \quad \forall n \in \{1, \dots, l-1\} \quad (15)$$

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

$$\ell_n(\bar{x}_n, \bar{m}_n) \leq a_l \quad \forall n \in \{1, \dots, l-1\} \quad (16)$$

We 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

$$\ell_n(x_n, m_n) \geq a_l \quad \forall n \in \{l, \dots, N\} \quad (17)$$

We can now lower-bound the cost of the assignment (\mathbf{x}, \mathbf{m}) . We have

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

using (14) and (17). Then rearranging the terms we have

$$\begin{aligned} C(\mathbf{x}, \mathbf{m}) &\geq \sum_{n=1}^{l-1} (x_n - \bar{x}_n) \ell_n(\bar{x}_n, \bar{m}_n) + \\ &\quad \sum_{n=1}^{l-1} \bar{x}_n \ell_n(\bar{x}_n, \bar{m}_n) + \sum_{n=l}^N x_n a_l \quad (18) \end{aligned}$$

Then by (15) and (16) we have $\forall n \in \{1, \dots, l-1\}$, $x_n - \bar{x}_n \leq 0$ and $\ell_n(\bar{x}_n, \bar{m}_n) \leq a_l$, thus

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

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

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

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 $\text{BNE}(N, (1-\alpha)r)$, 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 expressed in terms of perturbations of the NCF strategy [7].

IV. 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 4a). We analyze the inefficiency of Nash equilibria 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] and value of altruism [2] to measure the improvement in performance achieved by optimal Stackelberg routing. For Stackelberg instance (N, r, α) , price of stability is defined as the ratio between the cost of the induced equilibrium, and the cost of the social optimum:

$$\text{POS}(N, r, \alpha) = \frac{C(\bar{\mathbf{s}} + \mathbf{t}(\bar{\mathbf{s}}), \mathbf{m}(\bar{\mathbf{s}}))}{C^*}$$

where $\bar{\mathbf{s}}$ is the NCF strategy $\text{NCF}(N, r, \alpha)$, and $C^* = \min_{(\mathbf{x}, \mathbf{m})} C(\mathbf{x}, \mathbf{m})$. The improvement achieved by optimal

if a compliance rate greater than 0.46 is not feasible, the controller may prefer to implement a control strategy with α close to 0.14, since further increasing the compliance rate will not improve efficiency.

B. Scaling with the size of network

To illustrate the performance of the algorithm as the size of the network scales up, we measured the computation time of the NCF strategy for 10 randomly generated networks of size $N \in [3, 1000]$, with latency functions corresponding to randomly generated triangular flux functions. We fixed the compliance rate to be $\alpha = 0.4$ and the demands to be 0.4 and 0.8 times the maximum demand $r^{\text{NE}}(N)$. The results are given in Figure 8. As shown in Section III, the worst-case complexity of computing optimal Stackelberg assignments is quadratic in the size of the network, which is verified experimentally as illustrated in Figure 8 (the dashed line is a quadratic function that fits the data).

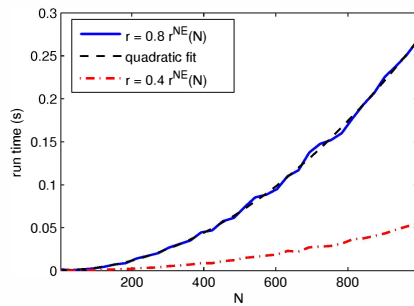


Fig. 8: Computation time of the NCF strategy for increasing network size.

Fig. 8 also 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.

V. DISCUSSION AND OPEN PROBLEMS

In order to address the inefficiency of Nash equilibria on parallel networks with horizontal queues, we considered the Stackelberg routing game where a central coordinator routes a fraction α of the total flow. We proved that for the HQSF latency class, 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 an example network for which we computed the decrease in inefficiency that can be achieved using optimal Stackelberg routing. This example showed that optimal Stackelberg routing can achieve a significant increase in efficiency even for small values of compliance rate α , especially when the demand is near critical flows $r^{(n,0)}$.

These results show that careful routing of a small compliant population can significantly improve the efficiency

of the network. The numerical results also show that for some ranges of demand and compliance ranges, Stackelberg routing can be completely ineffective (for example when the compliance rate is too low). 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 players 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] Y. A. Korilis, A. A. Lazar, and A. Orda. Achieving network optima using stackelberg routing strategies. *IEEE/ACM Transactions on Networking*, 5:161–173, 1997.
- [6] Y. A. Korilis, A. A. Lazar, and A. Orda. Capacity allocation under noncooperative routing. *IEEE Transactions on Automatic Control*, 42:309–325, 1997.
- [7] W. Krichene, J. Reilly, S. Amin, and Bayen A. M. Stackelberg routing on parallel networks with horizontal queues, <http://www.eecs.berkeley.edu/~walid/papers/stack.pdf>. 2012.
- [8] 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.
- [9] 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.
- [10] A. Ozdaglar and R. Srikant. Incentives and Pricing in Communication Networks. In *Algorithmic Game Theory*, N. Nisan, T. Roughgarden, E. Tardos, and V. Vazirani. Cambridge University Press, 2007.
- [11] C. Papadimitriou and G. Valiant. A new look at selfish routing. *Innovations in Computer Science (ICS)*, 2010.
- [12] T. Roughgarden. Stackelberg scheduling strategies. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pages 104–113. ACM, 2001.
- [13] T. Roughgarden and E. Tardos. How bad is selfish routing? *Journal of the ACM (JACM)*, 49(2):236–259, 2002.
- [14] D. Sehmeidler. Equilibrium points of nonatomic games. *Journal of Statistical Physics*, 7(4):295–300, 1973.
- [15] C. Swamy. The effectiveness of stackelberg strategies and tolls for network congestion games. In *SODA*, pages 1133–1142, 2007.