# Marginal Cost Pricing with a Fixed Error Factor in Traffic Networks

Guni Sharon
Texas A&M University
College Station, Texas, USA
guni@tamu.edu

Stephen D. Boyles
The University of Texas at Austin
Austin, Texas, USA
sboyles@mail.utexas.edu

Shani Alkoby
The University of Texas at Austin
Austin, Texas, USA
shani.alkoby@gmail.com

Peter Stone
The University of Texas at Austin
Austin, Texas, USA
pstone@cs.utexas.edu

## ABSTRACT

It is well known that charging *marginal cost tolls* (MCT) from self interested agents participating in a congestion game leads to optimal system performance, i.e., minimal total latency. However, it is not generally possible to calculate the correct marginal costs tolls precisely, and it is not known what the impact is of charging incorrect tolls. This uncertainty could lead to reluctance to adopt such schemes in practice. This paper studies the impact of charging MCT with some fixed factor error on the system's performance. We prove that under-estimating MCT results in a system performance that is at least as good as that obtained by not applying tolls at all. This result might encourage adoption of MCT schemes with conservative MCT estimations. Furthermore, we prove that no local extrema can exist in the function mapping the error value, $r$, to the system's performance, $T(r)$. This result implies that accurately calibrating MCT for a given network can be done by identifying an extremum in $T(r)$ which, consequently, must be the global optimum. Experimental results from simulating several large-scale, real-life traffic networks are presented and provide further support for our theoretical findings.

## KEYWORDS

Routing games; Congestion games; Marginal-cost pricing; Traffic flow; Flow optimization

## INTRODUCTION

Self interested agents that are routed in a congestible network, such as vehicles in a road network or packets in a data network, impose a *user equilibrium* (UE) that is often far worse than the *system optimum* (SO) flow [25]. Charging *marginal cost tolls* (MCT), in which each agent is charged a toll equivalent to the damage it inflicts on all other agents, results in a UE that achieves SO performance [1, 2, 21].

Calculating the MCT for a given agent, on a given path, i.e., the damage that the agent in question inflicts on other agents by traversing the path in question, is very challenging without making several restrictive assumptions (e.g., well-defined and known latency

functions) that do not hold in most traffic models and certainly not in real-life traffic. Recent work [27, 28] suggested a model free technique, denoted Δ-tolling, for approximating MCT.

Since Δ-tolling, or any tolling scheme that approximates MCT for that matter, is not guaranteed to result in the exact MCT, no optimality guarantees can be given regarding the system's performance. In fact, applying tolls different from MCT might result in a system performance that is worse than not applying tolls at all. This fact might deter public officials from implementing any tolling scheme that is not guaranteed to impose the exact MCT.

This paper examines the impact of imposing inaccurate MCT on the system's performance. Specifically, we provide conditions under which the system's performance will not be worse than applying no tolls, i.e., the system will not be worse off by imposing the tolling scheme. This paper establishes that charging a toll that is off by a factor, $r$, from the true MCT will not hurt the system's performance if $0 \leq r \leq 1$ (i.e., if MCT is underestimated by a constant factor). Moreover, this paper proves that the function mapping $r$ to the system's performance (total travel time) has a single (global) minimum and no local extrema. This fact implies that calibrating schemes for evaluating MCT e.g., Δ-tolling, can be carried out by identifying a minimum, which is guaranteed to be the global optimum.

Finally, experimental results from a traffic simulator are presented for different traffic scenarios. The experimental results match our theoretical claims by showing that, across various traffic scenarios, a global optimal flow is achieved for $r = 1$ and no extrema exist elsewhere.

## PRELIMINARIES

This paper assumes a standard flow model that is common in the routing and congestion games literature [21, 25, 31]. The terminology for this model follows Sharon et al. ([2018]) and is given next.

### The flow model

The flow model in this work is composed of a directed graph $G(V, E)$, where each link $e \in E$ is affiliated with a latency function. Additionally, the flow model requires a demand function $R(s, t) \to \mathbb{R}^+$ mapping a pair of vertices, $s, t \in V^2$, to a non-negative real number

representing the required amount of flow between source, $s$, and target, $t$.[1] A traffic flow scenario is a $\{G, R\}$ pair.

The variable $f_p$ represents the flow volume assigned to a path, $p$. Similarly, $f_e$ is the flow volume assigned to link $e$. Note that, a flow assignment to all paths implies a unique assignment to all links. By contrast, a flow assignment to all links does not necessarily imply a unique assignment to all paths. As an example, assume that link $e1$ is assigned a flow of $f_{e1}$. Further assume that $e1$ is part of two paths, $p1$ and $p2$, the flow assignment requires that $f_{e1} = f_{p1} + f_{p2}$ which might produce a range of possible flow assignments for $p1$ and $p2$. Hereafter we use the term *flow* or $f$ to represent a unique links flow assignment (which might be non-unique with respect to paths flow assignment).

A flow is defined as *valid* if:

- $f_p \geq 0$ for all paths $p$, that is, no path is assigned negative flow.
- the flow on each link ($f_e$) equals the summation of flows on all paths of which $e$ is a part. That is, $f_e = \sum_{p \in \mathcal{P}_e} f_p$ where $\mathcal{P}_e$ is the set of acyclic paths that include link $e$.

DEFINITION 1 (FEASIBLE FLOW). *A flow is defined as feasible if it is valid and the traffic demand is satisfied, that is, $\sum_{p \in \mathcal{P}_{st}} f_p = R(s, t)$ for all demand pairs $(s, t)$, where $\mathcal{P}_{st}$ is the set of acyclic paths leading from $s$ to $t$.*

Each link $e \in E$ has a latency function $l_e(f_e)$ which, given a flow volume ($f_e$), returns the latency (travel time) on $e$. The following regularity conditions on the latency function are a standard assumption in the transportation literature [20]

ASSUMPTION 1. *The latency function $l_e(f_e)$ is non-negative, convex, and its derivative, with regards to $f_e$ is positive for each link $e \in E$.*

The above assumption implies that, travel-time cannot be negative, more vehicles results in larger travel time, and that the $i^{th}$ vehicle causes a larger increase in travel time compared to the $j^{th}$ iff $i > j$.

The latency of a path, $p$, for a given flow, $f$, is defined as $l_p(f) = \sum_{e \in p} l_e(f_e)$. A feasible flow $f$ is defined as a *user equilibrium* (*UE*) if for every $s, t \in V^2$ and $p_a, p_b \in \mathcal{P}_{st}$ with $f_{p_a} > 0$ it holds that $l_{p_a}(f) \leq l_{p_b}(f)$ (see Lemma 2.2 in [25]). In other words, at *UE*, no amount of flow can be rerouted to a path with lower latency when the rest of the flow is fixed.

Define the total travel time associated with a link $e$ as $T_e(f_e) = l_e(f_e)f_e$. The total system travel time, for a given flow $f$, is $T(f) = \sum_{e \in E} T_e(f_e)$.

A feasible flow $f$ is defined as a *system optimum* (*SO*) if $T(f)$ is minimal over the set of feasible flows. We use $T(UE)$ to denote the total travel time at the UE solution. Similarly, $T(SO)$ denotes the total travel time at the SO solution.

Following the fact that, under Assumption 1, $T_e(f_e)$ is convex for any link $e$, it is easy to show that $T(f)$ is strictly convex in $f$. As a result, unique *UE* and *SO* flows exist [1, 4].

## Applying tolls

A recent body of work [3, 10, 28, 30, 33] assumed that each link in the network ($e \in E$) is assigned a toll value, $\tau_e$. The goal of such

tolls is to affect the route choice of self interested agents. Such work assumes that drivers are willing to sustain time delays in return for monetary gain (or avoiding monetary loss). This line of work requires translating time delays into monetary value using the agents' *value of time* (VOT). VOT represents the agents' monetary evaluation of a single unit of time.

Following previous work dealing with non-atomic flow [1, 2, 21, 24] we make the following assumptions and definition.

ASSUMPTION 2. *The agents' are homogeneous with regards to their time evaluation (VOT).*

DEFINITION 2 (GENERALIZED-COST *UE* (GUE)). *Let $\tau_p$ be the toll associated with path $p$ (the sum of the tolls on its constituting links i.e., $\sum_{e \in p} \tau_e$). A feasible flow $f$ is a GUE if for every $s, t \in V^2$ and $p_a, p_b \in \mathcal{P}_{st}$ with $f_{p_a} > 0$ it holds that $l_{p_a}(f) \times VOT + \tau_{p_a} \leq l_{p_b}(f) \times VOT + \tau_{p_b}$. In other words, at GUE, no amount of flow can be rerouted to a path with lower generalized cost (latency multiplied by VOT plus toll) when the rest of the flow is fixed.*

ASSUMPTION 3. *A solution for a traffic scenario follows the generalized-cost UE principle.*

Note that the above definition of *GUE* requires homogeneous VOT (Assumption 2). Nonetheless, *GUE* for heterogeneous VOT can be formulated as a dynamic user equilibrium (DUE) [17]. Though we expect that the main contributions of this paper extend naturally to that case, for clarity of presentation, we leave consideration of such models for future work.

A traffic scenario is said to be *toll-optimized* if the set of tolls ($\tau$) causes the *SO* and *GUE* solutions to align. Specifically, a sufficient (yet not necessary) condition for an optimized system is that $\tau$ equals the set of marginal cost tolls, $\tau^{MCT}$ [1, 4].

DEFINITION 3 (MARGINAL COST TOLL). *In marginal cost tolling (MCT) each agent (infinitesimally portion of the flow) is charged a toll equivalent to the damage it inflicts on the system. When the latency functions are differentiable, the MCT for link $e$ equals $f_e \frac{\partial l_e}{\partial f_e}$ That is, the increase in travel time caused by adding one more unit of flow to link $e$ (i.e., $\frac{\partial l_e}{\partial f_e}$) multiplied by all the flow that suffers from this increase (i.e., $f_e$). We use $\tau_e^{MCT}$ to denote the marginal cost toll for link $e$.*

Assuming that the latency functions are known and differentiable is not practical in many traffic models e.g., the cell transmission model [5, 6] or microsimulation models [8, 11, 32]. Such an assumption is certainly not practical for real-life traffic networks. Consequently, Sharon et al. [27, 28] introduced Δ-tolling, a model-free method for approximating MCT when the latency function is unknown. Despite showing reductions in total system travel time across markedly different traffic models, Δ-tolling, or any mechanism that approximates MCT for that matter, is not guaranteed to be toll-optimized. This fact poses a major problem since applying tolls that are different than MCT might result in arbitrarily worse total system travel time compared to that at the *UE*.

This paper makes a first attempt to examine the impact of applying inaccurate MCT. Specifically, it provides conditions under which the system performance (total system travel time) will be no worse than that at the *UE* solution. Providing such conditions is not trivial since

---

[1]The demand between any source and target, $R(s, t)$, can be viewed as an infinitely divisible set of agents (also known as a non-atomic flow [23]).
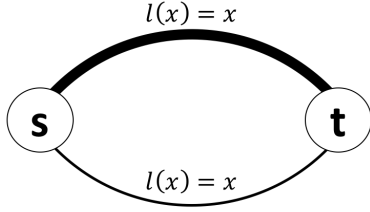
**Figure 1: A symmetrical network where *SO* and *UE* naturally align. Applying inaccurate MCT on one link will result in a *GUE* that is worse compared to the *SO* and *UE* solutions.**

slight errors regarding the marginal cost toll on specific links can add up and dramatically affect the price affiliated with many paths in a given road network.

The network presented in Figure 1 illustrates a possible effect of inaccurate MCT. In this symmetrical network the flow, $R(s, t)$, would split evenly between the top and bottom links in both the *SO* and *UE* solutions. In such a solution the MCT equals $0.5R(s, t)$ on both links. Increasing/decreasing this toll value on one link while keeping it constant for the other would throw the system out of balance and result in a new *GUE* that is worse than both the *SO* and *UE*.

Providing bounds for arbitrary errors in the value of MCT across a network is challenging, as illustrated in the above example. As a result, this work focuses on scenarios where the MCT error is of constant factor across the network.

## INACCURATE MARGINAL COST TOLLS

We consider a scenario where the tolls assigned to all links in a network are off by some factor from the MCT. Such a scenario might represent a systemic error in evaluating the $\beta$ parameter in $\Delta$-tolling (see [28] for exact details). Another relevant scenario is one in which MCT can accurately be computed in units of time delays. In such cases, a systemic error in the evaluation of the agents' VOT would result in a constant factor MCT error.

DEFINITION 4 (MCT-ERRORED SCENARIO). *A scenario is said to be MCT-errored if the toll affiliated with every link, $e \in E$, equals $r \cdot \tau_e^{MCT}$ for some error factor $r \geq 0$.*

Define the *GUE* flow for an MCT-errored scenario with error $r$ as $f^r$. As a result, $T(f^r)$ denotes the total system travel time for the *GUE* flow. Since $f^r$ is a function of $r$, we use $T(r)$ instead of $T(f^r)$ for brevity.

## BOUNDING THE SYSTEM'S PERFORMANCE

The following section presents the main contribution of this work i.e., provable bounds on the system's performance (total system travel time) as a function of the error factor $r$. We begin with several supporting lemmas.

LEMMA 1. *A GUE flow, $f$, for an MCT-errored system minimizes*

$$r \sum_{e \in E} \left[ f_e l_e(f_e) \right] + (1 - r) \sum_{e \in E} \left[ \int_0^{f_e} l_e(z) dz \right] \qquad (1)$$

*subject to $f$ being feasible (see Definition 1).*

PROOF. Combining this objective function with the feasibility constraints results in the following convex program (convexity is proven in Theorem 1):

$$Min. \quad r \sum_{p \in \mathcal{P}} \left[ f_p l_p(f_p) \right] + (1 - r) \sum_{p \in \mathcal{P}} \left[ \int_0^{f_p} l_p(z) dz \right]$$

Subject to:

$$\sum_{p \in \mathcal{P}_{st}} f_p = R(s, t) \qquad \forall s, t \qquad (2)$$

$$f_p \geq 0 \qquad \forall p \qquad (3)$$

Notice that the objective function in the above convex program includes a summation over paths. This is in contrast to Equation 1 which includes a summation over links. This discrepancy is made possible by the flow constraint which is defined by $f_e = \sum_{p \in \mathcal{P}_e} f_p$.

The appropriate Lagrange function for this convex program (ignoring the non-negativity constraint) is:

$$\mathcal{L}(f, \lambda) = r \sum_{p \in \mathcal{P}} \left[ f_p l_p(f) \right] + (1 - r) \sum_{p \in \mathcal{P}} \left[ \int_0^{f_p} l_p(z) dz \right] + \sum_{s, t \in V^2} \left[ \lambda_{st} \left( R(s, t) - \sum_{p \in \mathcal{P}_{st}} f_p \right) \right]$$

Incorporating the non-negativity constraint (given in Equasion 3) results in the following KKT optimality conditions:

$$f_p \geq 0 \qquad \forall p$$
$$(4)$$

$$\frac{\partial \mathcal{L}}{\partial f_p} \geq 0 \equiv l_p(f) + \mathbf{r} f_p l'_p(f) \geq \lambda_{st} \qquad \forall s, t \in V^2, \ p \in \mathcal{P}_{st}$$
$$(5)$$

$$f_p \frac{\partial \mathcal{L}}{\partial f_p} = f_p \left( l_p(f) + \mathbf{r} f_p l'_p(f) - \lambda_{st} \right) = 0 \quad \forall s, t \in V^2, \ p \in \mathcal{P}_{st}$$
$$(6)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{st}} = 0 \qquad \forall s, t \in V^2$$
$$(7)$$

Notice that Conditions 4 - 7 imply GUE for an MCT-errored scenario. The condition given in Equation 4 enforces non-negative path flows. The condition given in Equation 5 enforces that $\lambda_{st}$ is the minimal generalized cost (latency, $l_p(f)$, plus errored marginal-cost toll, $r f_p l'_p(f)$) over all paths leading from $s$ to $t$. The condition given in Equation 6 enforces that if a path is used ($f_p > 0$) its generalized cost must be equal to $\lambda_{st}$.

□

Next, we turn to prove that any solution that satisfies the GUE criterion results in the same system travel time. Specifically, we show that all flow assignments satisfying the above optimality conditions must be the same solution.

THEOREM 1. *A GUE flow for an MCT-errored scenario exists and is unique.*

PROOF. In order to prove this lemma it is sufficient to show that the objective function given in Lemma 1 (Equation 1) is strictly convex in the flow assignment ($f$). The Hessian matrix for Equation

1 ($H \in \mathbb{R}^{|E| \times |E|}$) is diagonal, where each entry on the diagonal (representing one link, $e \in E$) equals:

$$(r+1)\frac{\partial l_e}{\partial f_e} + r f_e \frac{\partial l_e^2}{\partial f_e^2} \tag{8}$$

For any link, $e$, the value of equation 8 is strictly positive since:

- $r \geq 0$, see Definition 4.
- $\frac{\partial l_e}{\partial f_e} > 0$, see Assumption 1.
- $f_e \geq 0$, see Definition 1.
- $\frac{\partial l_e^2}{\partial f_e^2} \geq 0$, see Assumption 1.

A diagonal matrix with strictly positive entries along its diagonal is positive definite. As a result, Equation 1 is strictly convex. □

Given that a unique *GUE* flow that minimizes equation 1 exists, we now turn to evaluate its impact on total system travel time for three key $r$ values: 0, 1, and $\infty$.

LEMMA 2. $T(0) = T(UE)$

PROOF. Setting $r = 0$ in Equation 1 results in the minimization of

$$\sum_{e \in E} \int_0^{f_e} l_e(z)dz$$

subject to the feasibility constraint. This minimization problem results in the UE flow [1]. □

LEMMA 3. $T(1) = T(SO)$

PROOF. Setting $r = 1$ in Equation 1 results in the minimization of

$$\sum_{e \in E} f_e l_e(f_e)$$

subject to the feasibility constraint. This minimization problem translates to minimizing total system travel time i.e., an SO flow [1]. □

LEMMA 4. $T(\infty) = T(f^\infty)$ *where $f^\infty$ is a UE solution for a scenario in which the latency affiliated with every path, p, equals* $f_p \frac{\partial l_p}{\partial f_p}$.

PROOF. Dividing Equation 1 by a positive constant (specifically $r$) preserves the minimizing assignment and yields

$$\sum_{e \in E} [f_e l_e(f_e)] + \frac{1-r}{r} \sum_{e \in E} \left[ \int_0^{f_e} l_e(z)dz \right] \tag{9}$$

Since $\lim_{r \to \infty}(1-r)/(r) = -1$, Equation 9 converges to

$$\sum_{e \in E} [f_e l_e(f_e)] - \sum_{e \in E} \left[ \int_0^{f_e} l_e(z)dz \right] \tag{10}$$

The KKT optimality conditions for minimizing Equation 10 under the feasibility constraints include:

$$f_p \geq 0 \qquad\qquad \forall p \tag{11}$$
$$f_p l_p'(f_p) \geq \lambda_{st} \qquad\qquad \forall st, \ p \in \mathcal{P}_{st} \tag{12}$$
$$f_p(f_p l_p'(f_p) - \lambda_{st}) = 0 \qquad\qquad \forall st, \ p \in \mathcal{P}_{st} \tag{13}$$

from which the *UE* definition (see Section ) holds if the latency function for any path $p$ is replaced by $f_p \frac{\partial l_p}{\partial f_p}$. □
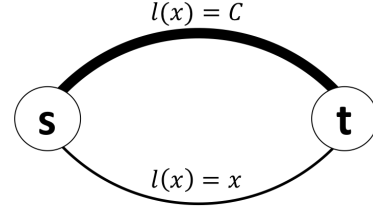


**Figure 2: A network where setting $r = \infty$ results in an arbitrary worse system performance compared to both the UE and SO solutions.**

Lemma 4 implies that at $r = \infty$ the system performance (total system travel time) can be arbitrarily worse than $T(SO)$ or $T(UE)$. As an example, consider the network depicted in Figure 2. The latency on the bottom link equals the fraction of flow that is assigned to it. If, for instance, 25% of the flow is assigned to the bottom link then the travel time on that link equals 0.25. The latency on the top link equals a constant, $C$, regardless of the amount of flow that is assigned to it. For $C \geq 2$ the SO and UE align and $T(SO) = T(UE) = 1 \cdot R(s,t)$. Since the latency on the top link is not a function of the flow, $MCT = x\frac{\partial l_p}{\partial x} = 0$ for the top link while $MCT = x\frac{\partial l_p}{\partial x} \geq 0$ for the bottom link. As a result, at $r = \infty$, 100% of the flow from $s$ to $t$ would travel the top link while 0% would travel the bottom link. Such a flow would result in total system travel time $= C \cdot R(s,t)$. It is easy to see that as $C$ increases so does the difference between $T(\infty)$ and $T(SO)$ or $T(UE)$, potentially to infinity.

Given that no bound on the system's performance can be given for $r = \infty$ we turn to examine bounds on other values of $r$. We start by examining values of $r$ that fall between zero and one.

LEMMA 5. *Any two error values $0 \leq r_1 < r_2 < 1$ satisfy $T(r_1) \geq T(r_2)$.*

PROOF. For simplicity of presentation we use $U(r)$ to denote

$$\sum_{e \in E} \left[ \int_0^{f_e^r} l_e(z)dz \right]$$

Any *GUE* flow $f^r$ must minimize Equation 1 (Lemma 1). That is, subject to being feasible, $f^r$ minimizes the expression $rT(r) + (1-r)U(r)$. Minimizing Equation 1 under $r_1$ requires that

$$r_1 T(r_2) + (1-r_1)U(r_2) \geq r_1 T(r_1) + (1-r_1)U(r_1)$$

and as a result

$$r_1(T(r_2) - T(r_1)) \geq (1-r_1)(U(r_1) - U(r_2)) \tag{14}$$

Similarly, minimizing Equation 1 under $r_2$ requires that

$$r_2(T(r_2) - T(r_1)) \leq (1-r_2)(U(r_1) - U(r_2)) \tag{15}$$

Assume, in contradiction to the lemma, that $T(r_2) - T(r_1) > 0$. Since $1 - r_2 > 0$ and $r_2 > 0$, Equation 15 would require $U(r_1) - U(r_2) > 0$. Since all the components of Equations 14 and 15 are strictly positive,

we can rewrite them as:

$$\frac{r_1}{1-r_1} \geq \frac{U(r_1) - U(r_2)}{T(r_2) - T(r_1)} \tag{16}$$

$$\frac{r_2}{1-r_2} \leq \frac{U(r_1) - U(r_2)}{T(r_2) - T(r_1)} \tag{17}$$

From Equations 16 and 17 we obtain

$$\frac{r_1}{1-r_1} \geq \frac{r_2}{1-r_2} \tag{18}$$

Since the function $f(r) = r/(1-r)$ is continuous and strictly increasing for $r < 1$ then Equation 18 must satisfy $r_1 \geq r_2$ in contradiction to the lemma's premise. □

Next we turn to examine the behavior of error values that are greater than one.

LEMMA 6. *Any two error values* $1 < r_1 < r_2$ *satisfy* $T(r_1) \leq T(r_2)$.

PROOF. Assume, in contradiction to the lemma, that $T(r_2) - T(r_1) < 0$. Since $1 - r_1 < 0$ and $r_1 > 1 > 0$, Equation 14 requires $U(r_1) - U(r_2) > 0$. Even though the signs of $(T(r_2) - T(r_1))$ and $(1 - r_1)$ and $(1 - r_2)$ are in contrast to the case presented in Lemma 5, rearranging Equations 14 and 15 still result in Equations 16 and 17 which leads to the inequality in Equation 18. Since the function $f(r) = r/(1-r)$ is continuous and strictly increasing for $r > 1$ then Equation 18 must satisfy $r_1 \geq r_2$ in contradiction to the lemma's premise. □

Following Lemma 5 and 6 we can now provide bounds for an MCT-errored system.

THEOREM 2. *If* $0 \leq r \leq 1$ *then* $T(r) \leq T(UE)$.

PROOF. $T(0) = T(UE)$ (Lemma 2) and $T(r)$ is non increasing in the interval $[0, 1)$ (Lemma 5). Also $T(1) = T(SO) \leq T(UE)$ (Lemma 3). □

THEOREM 3. *If* $r \geq 1$ *then* $T(r) \leq T(f^\infty)$ *when* $f^\infty$ *is a UE solution for a scenario where the latency on every path, p, equals* $f_p l'_p(f_p)$.

PROOF. $T(\infty) = T(f^\infty)$ when $f^\infty$ is a UE solution for a scenario where the latency for every path, $p$, equals $f_p l'_p(f_p)$ (Lemma 4). $T(r)$ is non decreasing for $r > 1$ (Lemma 6). Also $T(1) = T(SO) \leq T(\infty)$ (Lemma 3). □

Theorem 2 implies that when underestimating MCT by a constant factor, $0 \geq r < 1$, the systems performance cannot be worse that the one obtain by the UE solution, $T(UE)$.

Theorem 3 implies that when overestimating MCT by a constant factor, $r > 1$, the systems performance cannot be worse then $T(\infty)$. However since $T(\infty)$ can be arbitrary worse than $T(UE)$ and $T(SO)$, this bound is not as useful as the one provided for the previous case, $0 < r < 1$.

## EMPIRICAL STUDY

In order to validate our theoretical findings, we simulated different traffic scenarios while varying the MCT error factor ($r$). The total system performance (total system travel time) was measured for each setting and the trends were compared to the above theoretical claims.

### Traffic scenario

Each simulated traffic scenario is defined by two attributes:

(1) The road network, $G(V, E)$, specifying the set of vertices and links where each link is affiliated with a length, capacity, and speed limit, these link attributes are used to set the link's latency function. Following standard practice, networks are partitioned into traffic analysis zones (TAZs) and each zone contains a vertex belonging to $V$ called the centroid. All traffic originating and terminating within the zone is assumed to enter and leave the network at the centroid.
(2) A trip table specifying the traffic demand between pairs of centroids. The demand, $R(s, t)$, between vertices, $s, t \in V^2$, other than centroids, is set to zero.

Following Sharon et al. ([2018]) the following benchmark scenarios were chosen: Sioux Falls, Eastern Massachusetts, Anaheim, Chicago Sketch, Philadelphia, and Chicago Regional. All traffic scenarios are available at: https://github.com/bstabler/TransportationNetworks. Figure 3 depicts three representative network topologies (Sioux Falls, Eastern Massachusetts, Anaheim).

Table 1 presents the scenario specifications i.e., number of vertices, links, and zones for the traffic network that is affiliated with each scenario. Total demand, summed over all $\{s, t\}$ pairs, as specified by the affiliated trip table are also provided (as "Total Demand"). The same table also presents total system travel times for different error values, these results are discussed later.

### The Traffic Model

The *GUE* solutions for the above scenarios were computed using Algorithm B [7]. For all scenarios, the model assumes that travel times follow the *Bureau of Public Roads* (BPR) function [19] with the commonly used parameters $\alpha = 0.15$ and $\beta = 4$. Since computing the *GUE* solution requires solving a convex program (see 1), we only solve it to a limited precision.

To measure convergence, given an assignment of agents to paths, we define the average excess cost (AEC) as the average difference between the travel times on paths taken by the agents and their shortest alternative path. The algorithm is terminated when the AEC is less than $1 \times 10^{-6}$ minutes.

### Results

In addition to the scenario specifications, Table 1 also presents the system's performance (total system travel time) for five different error values ($r = \{0, 0.5, 1, 2, \infty\}$). The *SO* solution ($r = 1$) provides the best performance (minimal total system travel time), as expected. The performance for $r = \infty$ is slightly better than that at the *UE* solution ($r = 0$) in some cases, e.g., Sioux Falls and Philadelphia, but might be significantly worse in others, e.g., Eastern Massachusetts where $T(\infty)$ was outperformed by $T(UE)$ by 15%.
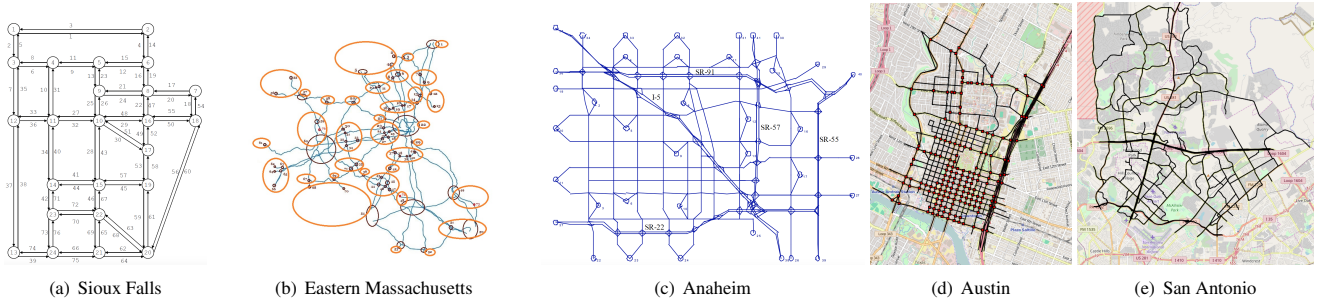
(a) Sioux Falls


(b) Eastern Massachusetts


(c) Anaheim


(d) Austin


(e) San Antonio

**Figure 3: Traffic networks used in the experiments.**

| Scenario | Vertices | Links | Zones | Total Demand | $T(UE)$ | $T(0.5)$ | $T(SO)$ | $T(2)$ | $T(\infty)$ |
|---|---|---|---|---|---|---|---|---|---|
| Sioux Falls | 24 | 76 | 24 | 360,600 | 7,480,223 | 7,205,048 | 7,194,256 | 7,198,091 | 7,222,857 |
| Eastern MA | 74 | 258 | 74 | 65,576 | 28,181 | 27,411 | 27,324 | 27,392 | 32,460 |
| Anaheim | 416 | 914 | 38 | 104,694 | 1,419,913 | 1,397,216 | 1,395,015 | 1,398,631 | 1,549,075 |
| Chicago S | 933 | 2,950 | 387 | 1,260,907 | 18,377,331 | 17,991,235 | 17,953,268 | 17,994,192 | 19,630,440 |
| Chicago R | 12,982 | 39,018 | 1790 | 1,360,427 | 33,656,969 | 32,078,668 | 31,942,957 | 32,096,038 | 38,190,675 |
| Philadelphia | 13,389 | 40,003 | 1525 | 18,503,872 | 335,647,096 | 325,211,099 | 324,268,465 | 325,176,216 | 335,296,306 |

**Table 1: The system performance (total system travel time) given as "$T(x)$" for different scenarios along with network specifications, for each scenario: number of vertices, links, zones, and total demand ($\sum_{st} R(s, t)$).**
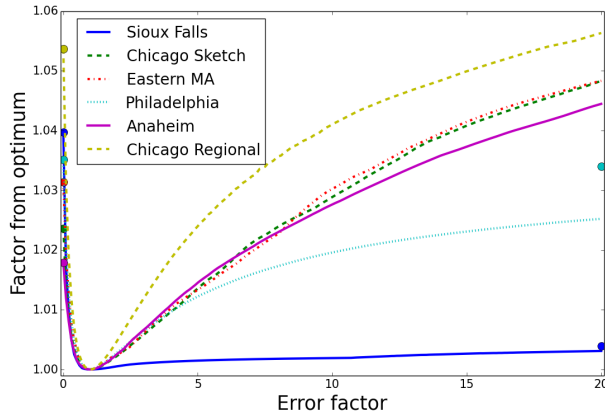


**Figure 4: Normalized total system travel time as a function of the error factor (r) for six benchmark traffic scenarios.**



**Figure 5: Normalized total system travel time (with 95% confidence intervals) as a function of $\beta$ in $\Delta$-*tolling* (representing the error factor) for three benchmark traffic scenarios.**

Results for applying half and double the true MCT are also provided ($T(0.5)$ and $T(2)$ respectively). Results for these values are mixed where in some cases $T(0.5)$ performs slightly better than $T(2)$ and vice versa in others. Nonetheless, $r = 0.5$ has a clear advantage over $r = 2$ since, unlike $T(2)$, the value of $T(0.5)$ is bounded by $T(UE)$ for any scenario (Theorem 2).

Figure 4 presents normalized values for total system travel time as a function of the error factor $r$. The total system travel time values (y-axis) for each curve are normalized according to $T(SO)$ e.g., a total system travel time value of 2 correlates to double $T(SO)$ for the relevant curve (scenario). Consequently, $T(1) = T(SO) = 1$
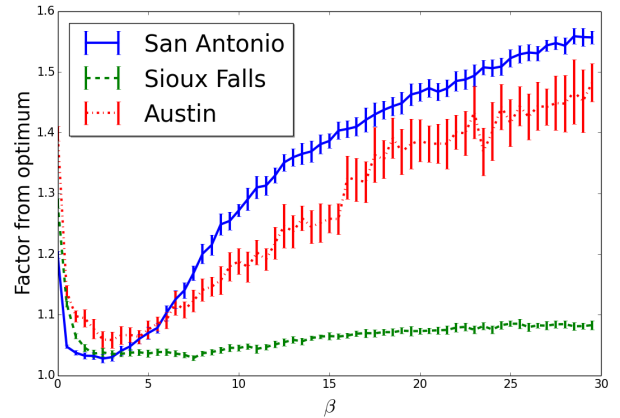
in all the curves. The data points were computed for the range $r = [0, 20]$ with a step size of 0.1. Each of the curves starts with a dot representing $T(UE)$. Additionally, dots on the right border of the plot represent $T(\infty)$. Such dots are presented only for the Sioux Falls and Philadelphia scenarios as $T(\infty)$ is out of the presented total system travel time range for the rest (exact values are available in Table 1). As predicted by Lemmas 5 and 6 the curves are non-increasing in the range $[0, 1]$ and non-decreasing in the range $[1, \infty]$.

## Dynamic traffic assignment

The traffic model that is assumed in this paper, though common in the traffic literature, does not apply to many real-life traffic scenarios. In order to broaden the impact of this research, we turn to investigate the performance of an MCT-errored scenario in a more realistic traffic flow model. Specifically, we test our findings in a **dynamic traffic assignment** setting. A dynamic traffic assignment model combines a traffic model with time-varying network states with a route choice principle (drivers choose routes to minimize some combination of their travel time and toll cost).

Dynamic traffic assignment iterates between finding shortest paths, assigning vehicles, and evaluating travel times through simulation, to find a route assignment near dynamic user equilibrium [14]. DTA models can be used to perform many simulations of city network traffic in a reasonable time. DTA models commonly use the kinematic wave theory of traffic flow, which models traffic as a compressible fluid [15, 22]. The kinematic wave theory models several important aspects of traffic behavior including the formation and dissipation of congestion waves over time due to bottlenecks. The kinematic wave model involves a system of partial differential equations which are solved numerically given initial and boundary conditions. One common solution method is the cell transmission model (CTM) [5, 6], which is a Godunov scheme [9] for the kinematic wave theory. The CTM can be used with a variety of intersection models [29], including traffic signals and autonomous reservation schemes [13]. Using such intersection models, CTM, unlike the static model defined by Assumptions 1, 2, and 3, takes into account inter-link effects, making CTM more realistic on the one hand but intractable for large networks.

In order to further mimic a realistic setting, drivers were assigned heterogeneous evaluation of time. The time evaluation per driver was randomly drawn from a Dagum distribution with parameters $\hat{a} = 22020.6$, $\hat{b} = 2.7926$, and $\hat{c} = 0.2977$, reflecting the distribution of personal income in the United States [16]. These settings were chosen to be identical to those presented in previous work [27, 28].

Three traffic scenarios (depicted in Figure 3) were evaluated using the CTM framework.

- **Sioux Falls** - [12] — this scenario is widely used in the transportation research literature [13], and consists of 76 directed links, 24 nodes (intersections) and 28,835 trips spanning 3 hours.
- **Downtown Austin** - [14] — this network consists of 1,247 directed links, 546 nodes and 62,836 trips spanning 2 hours during the morning peak.
- **Uptown San Antonio** - this network consists of 1,259 directed links, 742 nodes and 223,479 trips spanning 3 hours during the morning peak.

Since there is no closed form equation for computing MCT in DTA for the general case, the $\Delta$-*tolling* mechanism was used to approximate MCT. The $\beta$ parameter in $\Delta$-*tolling* acts as a proportional parameter for $\Delta$-*tolling* (for more details see [28]) and, thus, was used to represent different error values ($r$).[2] Figure 5 is similar in structure to Figure 4, representing total system travel time as a function of the MCT error (represented by different $\beta$ values) but

[2]In the reported experiments the $R$ parameter was set to $10^{-4}$ for $\Delta$-*tolling* following the best performing value reported by Sharon et al. [2017b].

for DTA scenarios. Since DTA is not deterministic with regard to the VOT assigned to each driver, an average of 20 runs is presented per data-point with 95% confident intervals.

DTA does not follow the assumptions made in the above theoretical analysis (Assumptions 1 and 3). As a result, Lemmas 5 and 6 and Theorems 2 and 3 do not hold. Nonetheless, the general trend where the system performance improves until some optimal point and then deteriorates can still be observed suggesting that the general conclusions drawn in this work are relevant to real-world traffic.

## DISCUSSION

Lemmas 5 and 6 and Theorems 2 and 3 lead to the following theoretical conclusions:

- Underestimating MCT by a constant factor across a traffic network would result in a system performance that is not worse than the no-toll user equilibrium.
- When calibrating a parameter that is a multiplier of the true MCT, a value that is locally optimal is guaranteed to be globally optimal.

The presented empirical results suggest that these conclusions extend to realistic traffic models. The implications of these conclusions might be substantial when installing a new tolling scheme with a tunable parameter, $\theta$ where the value of $\theta$ correlates to a fixed error in MCT. As stated in Section , this can occur when calibrating the expected drivers' value of time or the $\beta$ parameter in $\Delta$-*tolling* [28] as done in *Enhanced $\Delta$-tolling* [18]. The calibration process in such cases amounts to detecting a local minimum (which is guaranteed to be the global minimum).

## SUMMARY AND FUTURE WORK

This paper considers a traffic scenario in which marginal-cost tolls (MCT) with a fixed factor error is imposed on all drivers. The system performance is analyzed with regards to the error rate and performance bounds are provided as a function of the error value. Three main claims are proven:

(1) If the error factor is lower than 1 (MCT is underestimated) the system will not perform worse than if no tolls were applied.
(2) As the error factor increases from 0 to 1 the system's performance will not deteriorate.
(3) As the error factor increases from 1 to infinity the system's performance will not improve.

These claims can allow the tuning of MCT-based tolling schemes while insuring quality of service along the tuning process. There are many other conceivable errors besides a multiplicative, system-wide factor on the true MCT. Consequently, future work ought to examine scenarios with other assumptions on the toll error, such as when the assessed toll is within some bounded interval around the MCT.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Martin Beckmann, Charles B. McGuire, and Christopher B. Winsten. 1956. *Studies in the Economics of Transportation*. Yale University Press.

[2] Dietrich Braess. 1968. Über ein Paradoxon aus der Verkehrsplanung. *Unternehmensforschung* 12, 1 (1968), 258–268.

[3] Haipeng Chen, Bo An, Guni Sharon, Josiah P. Hanna, Peter Stone, Chunyan Miao, and Yeng Chai Soh. 2018. DyETC: Dynamic Electronic Toll Collection for Traffic Congestion Alleviation. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI-18)*.

[4] Stella C. Dafermos and Frederick T. Sparrow. 1969. The traffic assignment problem for a general network. *Journal of Research of the National Bureau of Standards B* 73, 2 (1969), 91–118.

[5] Carlos F. Daganzo. 1994. The cell transmission model: A dynamic representation of highway traffic consistent with the hydrodynamic theory. *Transportation Research Part B: Methodological* 28, 4 (1994), 269–287.

[6] Carlos F. Daganzo. 1995. The cell transmission model, part II: network traffic. *Transportation Research Part B: Methodological* 29, 2 (1995), 79–93.

[7] Robert B. Dial. 2006. A path-based user-equilibrium traffic assignment algorithm that obviates path storage and enumeration. *Transportation Research Part B: Methodological* 40, 10 (2006), 917–936.

[8] Kurt Dresner and Peter Stone. 2008. A multiagent approach to autonomous intersection management. *Journal of artificial intelligence research* 31 (2008), 591–656.

[9] Sergei Konstantinovich Godunov. 1959. A difference method for numerical calculation of discontinuous solutions of the equations of hydrodynamics. *Matematicheskii Sbornik* 89, 3 (1959), 271–306.

[10] Stephen D. Boyles Peter Stone Josiah P. Hanna, Guni Sharon. 2019. Selecting Compliant Agents for Opt-in Micro-Tolling. In *Proceedings of the 33nd AAAI Conference on Artificial Intelligence (AAAI-19)*.

[11] Daniel Krajzewicz, Georg Hertkorn, Christian Rössel, and Peter Wagner. 2002. SUMO (Simulation of Urban MObility)-an open-source traffic simulation. In *Proceedings of the 4th middle East Symposium on Simulation and Modelling (MESM20002)*. 183–187.

[12] Larry J LeBlanc, Edward K Morlok, and William P Pierskalla. 1975. An efficient approach to solving the road network equilibrium traffic assignment problem. *Transportation research* 9, 5 (1975), 309–318.

[13] Michael W Levin and Stephen D Boyles. 2015. Intersection auctions and reservation-based control in dynamic traffic assignment. *Transportation Research Record: Journal of the Transportation Research Board* 2497 (2015), 35–44.

[14] Michael W Levin, Matt Pool, Travis Owens, Natalia Ruiz Juri, and S Travis Waller. 2015. Improving the convergence of simulation-based dynamic traffic assignment methodologies. *Networks and Spatial Economics* 15, 3 (2015), 655–676.

[15] Michael James Lighthill and Gerald Beresford Whitham. 1955. On kinematic waves II. A theory of traffic flow on long crowded roads. *Proc. R. Soc. Lond. A* 229, 1178 (1955), 317–345.

[16] P Łukasiewicz, K Karpio, and A Orłowski. 2012. The models of personal incomes in USA. *Acta Physica Polonica A* 121, 2B (2012).

[17] Hani Mahmassani and Robert Herman. 1984. Dynamic user equilibrium departure time and route choice on idealized traffic arterials. *Transportation Science* 18, 4 (1984), 362–384.

[18] Hamid Mirzaei, Guni Sharon, Stephen Boyles, Tony Givargis, and Peter Stone. 2018. Enhanced Delta-tolling: Traffic Optimization via Policy Gradient Reinforcement Learning. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 47–52.

[19] Ren Moses and Enock T. Mtoi. 2017. Calibration and Evaluation of Link Congestion Functions. *Journal of Transportation Technologies* 4, 2 (2017).

[20] Michael Patriksson. 2015. *The traffic assignment problem: models and methods*. Courier Dover Publications.

[21] Arthur Pigou. 1920. *The Economics of Welfare*. Palgrave Macmillan.

[22] Paul I. Richards. 1956. Shock waves on the highway. *Operations research* 4, 1 (1956), 42–51.

[23] Robert W. Rosenthal. 1973. A class of games possessing pure-strategy Nash equilibria. *International Journal of Game Theory* 2, 1 (1973), 65–67.

[24] Tim Roughgarden. 2004. Stackelberg scheduling strategies. *SIAM J. Comput.* 33, 2 (2004), 332–350.

[25] Tim Roughgarden and Éva Tardos. 2002. How bad is selfish routing? *Journal of the ACM (JACM)* 49, 2 (2002), 236–259.

[26] Guni Sharon, Michael Albert, Tarun Rambha, Stephen Boyles, and Peter Stone. 2018. Traffic Optimization For a Mixture of Self-interested and Compliant Agents. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI-18)*.

[27] Guni Sharon, Josiah P. Hanna, Tarun Rambha, Michael W. Levin, Michael Albert, Stephen D. Boyles, and Peter Stone. 2017. Real-time Adaptive Tolling Scheme for Optimized Social Welfare in Traffic Networks. In *Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS-2017)*.

[28] Guni Sharon, Michael W. Levin, Josiah P. Hanna, Tarun Rambha, Stephen D. Boyles, and Peter Stone. 2017. Network-wide adaptive tolling for connected and automated vehicles. *Transportation Research Part C* 84 (September 2017), 142–157. https://doi.org/10.1016/j.trc.2017.08.019

[29] Chris MJ Tampère, Ruben Corthout, Dirk Cattrysse, and Lambertus H Immers. 2011. A generic class of first order node models for dynamic macroscopic simulation of traffic flows. *Transportation Research Part B: Methodological* 45, 1 (2011), 289–309.

[30] Hai Yang, Qiang Meng, and Der-Horng Lee. 2004. Trial-and-error implementation of marginal-cost pricing on networks in the absence of demand functions. *Transportation Research Part B: Methodological* 38, 6 (2004), 477–493.

[31] Hai Yang, Xiaoning Zhang, and Qiang Meng. 2007. Stackelberg games and multiple equilibrium behaviors on networks. *Transportation Research Part B: Methodological* 41, 8 (2007), 841–861.

[32] Qi Yang and Haris N. Koutsopoulos. 1996. A microscopic traffic simulator for evaluation of dynamic traffic management systems. *Transportation Research Part C: Emerging Technologies* 4, 3 (1996), 113–129.

[33] Bojian Zhou, Michiel Bliemer, Hai Yang, and Jie He. 2015. A trial-and-error congestion pricing scheme for networks with elastic demand and link capacity constraints. *Transportation Research part B: Methodological* 72 (2015), 77–92.