

## **Marginal congestion cost on a dynamic network with queue spillbacks**

**Kenneth A. Small, University of California at Irvine**

Dept. of Economics, University of California, Irvine, CA 92797-5100, USA

Tel: +1 949-824-5658; Fax: +1 949-824-2182

([ksmall@uci.edu](mailto:ksmall@uci.edu))

and

**Mogens Fosgerau, Technical University of Denmark & Centre for Transport Studies, Sweden**

Tel: +45 4525 6521

([mf@transport.dtu.dk](mailto:mf@transport.dtu.dk))

September 22, 2009

Keywords: congestion cost, marginal cost, dynamic congestion, queue, hypercongestion

Journal of Economic Literature code: L9, R41

### **Abstract**

We formulate an empirical model of congestion for a network where queues may form and spill back from one link to another. Its purpose is to disentangle the dynamic effect that a marginal vehicle, on a given link and at a given time, has on the distribution of travel times experienced there and on connected links. We estimate a dynamic model, based on an unusually complete and accurate dataset from Danish motorways. Each data point contains information on the vehicle flow on a link during a five-minute interval, along with the average speed experienced by those vehicles as measured by timed license-plate matches. We use the results to estimate the marginal external cost of adding a vehicle to a link's entry flow, as it is influenced by conditions on that link and on its downstream neighbor.

# Marginal congestion cost on a dynamic network with queue spillbacks

Mogens Fosgerau and Kenneth A. Small

## 1. Introduction

Congested road networks are receiving much attention as analysts and policy makers examine more sophisticated measures to manage traffic on existing facilities. These measures include ramp metering, express lanes, carpooling incentives, and pricing. The implications of such policies, especially express lanes and pricing, are mostly understood either from models of a single road link or from simulated networks in which road links are described by relatively simple speed-flow relationships connected, if at all, by simple queuing.

Yet the relationships spilling across links are crucial to understanding the development of highly congested systems, where queues can quickly spread and perhaps can also form spontaneously when flow approaches a saturation level. There is considerable uncertainty about the nature of flow under such conditions. It is known that on a single link, a given flow may occur at two different speeds, one relatively high and the other much lower and less stable. We shall use the terminology, common in economics, of “congestion” for the former case and “hypercongestion” for the latter.<sup>1</sup> But the exact process of transition from one to the other is much debated and seems to depend critically on how one link interacts with another.<sup>2</sup>

One common way to model severe congestion is through deterministic queuing at a bottleneck, perhaps including the spillback of queues from one link to another. Such analysis almost invariably makes the simplification that the bottleneck capacity is constant. Yet it is well known that flow tends to be unstable when it is near its maximum, and in fact capacity is often defined as the largest flow sustainable over a moderate time period rather than the maximum

---

<sup>1</sup> In much of the engineering literature, the corresponding terms are “free-flow” and “congested flow”.

<sup>2</sup> Small and Verhoef (2007, sect 3.4.1) provide a review of these arguments. See for example the difference in opinion about the spontaneous onset of hypercongestion as a type of phase transition, reflected in Kerner and Rehborn (1997) and Daganzo, Cassidy, and Bertini (1999). Verhoef (2001) and Small and Chu (2003) argue that hypercongestion does not exist in a stable steady-state equilibrium, but rather is generated dynamically when queues form behind bottlenecks. McDonald, d’Ouille, and Liu (1999), however, claim to observe stable hypercongestion on Chicago area expressways.

possible to achieve.<sup>3</sup> Furthermore, discharge rates from bottlenecks tend to rise to a temporarily high value, then fall as a queue forms, and then partially recover; a well-documented example is demonstrated by Cassidy and Bertini (1999), leading them to “view the [lower] long-run queue discharge flow as the bottleneck capacity” (p. 40). With this definition, flow can exceed capacity for short periods of uncertain duration, resulting in considerable stochastic variability in the travel times experienced and the marginal effects of an additional vehicle.

Another approach, used for city street networks, is to model average flows and speeds throughout an area. Both simulation and aerial photography have suggested that such average flows and speeds can be related by an aggregate speed-flow function that has both congested and hypercongested regimes (May, Shepherd, and Bates 2000, Ardekani and Herman 1987). Small and Chu (2003) develop a dynamic aggregate model based on such a relationship that can be used to measure the marginal cost of a vehicle entering the area, but it cannot describe heterogeneity of conditions within the area. At the opposite extreme, one can model the behavior of traffic at individual signalized intersections within street networks; but this analysis becomes extremely complex when queues at one link obstruct flow on another, a situation typically requiring dynamic computer simulations with individual vehicles.

Another difficulty in modeling dynamic congestion arises in the process of empirical estimation. Such estimation requires data on the traffic flow and speed (or either of these quantities along with density) at each of many locations and times. The most common source of such data is magnetic loop detectors placed in roadways. However, the resulting data contain serious errors due to periodically non-functioning equipment and uncertain assumptions about vehicle sizes and flow homogeneity needed to convert the observed timing and spacing of axle passages into vehicle flows and speeds (Steimetz and Brownstone, 2007). Furthermore, the causal relationship between aggregate traffic flow and speed is ambiguous.

This paper provides an empirical description of congestion formation throughout a freeway network covering a part of Denmark. We are able to solve many of the problems just described by taking advantage of an unusually detailed data set containing reliable speed measurements on each link at five-minute intervals over the entire day. Because the data are extensive, we can model congestion on these links using flexible dynamic functions.

---

<sup>3</sup> See, for example, Institute of Traffic Engineers (1982), p. 471, and the *Highway Capacity Manual* published regularly by the Transportation Research Board in the United States.

Specifically, we allow a dynamic relationship explaining travel time on a link in terms of present and past conditions on the link itself and also in terms of conditions downstream. The functional specification allows for spontaneous hypercongestion as well as hypercongestion caused by spillbacks from downstream congestion. We use the resulting model to simulate the pattern of marginal external costs associated with adding a vehicle to the traffic flow.

The results show that dynamic effects are quite important, causing perturbations in flow to persist for well over five minutes in many cases. They also show that marginal external costs arise both from the link itself, through the usual speed-flow relationship, and from the downstream link when it is congested. However, our results are quite sensitive to details of model specification, leading us to suspect that our approximate solution to a full reduced-form model of travel flow by link does not capture all the interactions that occur under heavy congestion.

The layout of the paper is as follows. The general model is formulated in Section 2, while Section 3 describes the data. The empirical model specification and estimation results are contained in section 4. Section 5 applies these results to calculate the marginal external cost associated with adding additional vehicles to traffic flow. Section 6 concludes.

## 2. Model Specification

The links on our network, indexed by  $n$ , are defined as sections of roadway between two intersections. The time periods, indexed by  $t$ , are 5 minutes in duration.

Let  $T_t^n$  be the link travel time (in minutes per kilometer) observed for vehicles exiting link  $n$ , with physical length  $L^n$  (in km), during time interval  $t$  of duration  $\Delta t=5$  min. Vehicles exit the link at rate  $F_t^n$  (in vehicles per lane per minute). These are the quantities on which we have direct measurements. Nearly all of our links have identical numbers of lanes, so expressing flows in per-lane units is mainly a convenience.

We assume exit flow  $F_t^n$  is the smaller of the potential flow reaching the end of the section and the capacity of the section, the latter being reduced by blockages from the downstream link. This potential exit flow rate is equal to the link's entering rate  $E_t^n$  plus the discharge over time interval  $\Delta t$  of any accumulated internal queue,  $q_t^n$ . Downstream blockage

depends in some unknown way on downstream density  $D_t^{n+1}$  (measured in vehicles per lane-kilometer). Thus:

$$F_t^n = \text{Min}\left\{\left[E_t^n + \left(q_t^n / \Delta t\right)\right], g_1(D_t^{n+1})\right\} \quad (1)$$

where  $g_1(\cdot)$  is a strongly nonlinearly decreasing function. The size of the internal queue is an accumulation of past excesses of entry flows over exit flows:

$$q_t^n = \text{Max}\left\{\sum_{t'=t_0^n}^{t-1} \Delta q_{t'}^n, 0\right\}; \quad \Delta q_{t'}^n = (E_{t'}^n - F_{t'}^n) \cdot \Delta t \quad (2)$$

where  $t_0^n$  is the most recent time period  $t'$  for which  $q_{t'}^n = 0$ .

Travel time follows a speed-density relationship:

$$T_t^n = h_1(D_t^n) \quad (3)$$

where density  $D$  is given by flow  $F$  divided by speed  $S=1/T$ :

$$D_t^n = T_t^n \cdot F_t^n . \quad (4)$$

Finally, we approximate entry flow during interval  $t$  based on what we know of the exit flows from the upstream link at previous times. (This will be inexact because we lack data on entry and exit ramps.) Due to the lengths of our sections and the five-minute duration of our time interval, we need the upstream flow for the current and up to two previous time periods. The result is

$$E_t^n = w_0 F_t^{n-1} + w_1 F_{t-1}^{n-1} + w_2 F_{t-2}^{n-1} \quad (5)$$

with weights summing to one and determined from the link length as described in the Appendix.

Equations (1)–(5) form a simultaneous system in various flows and travel times. We would like to solve them for the flows and travel times as functions of other variables. These endogenous variables affect each other in several highly nonlinear and interconnected ways. First,  $T$  appears on both sides of (3), since it is part of definition (4) of density. Second, current values are highly nonlinear functions of lagged values through queue formation as described in (2). Third, values of flow and travel time for section  $n$  are functions of their values for downstream sections through the term  $g_1(\cdot)$  in (1), representing blockage from downstream congestion; and they depend on upstream sections through the last term on the right-hand side of (5), representing how entry to section  $n$  depends on exit from section  $n-1$ . Of course, the same

equations apply to these upstream and downstream sections, so that congestion effects on a given section can propagate in both directions.

For these reasons, we find it intractable to estimate equations (1)–(5) structurally or to solve them explicitly for the endogenous variables. Instead, we suggest the following heuristic approximation of a solution for travel time  $T_t^n$ . It is motivated by our assessment of the most important sources of simultaneity. First, the solution will imply a strong dependence of current travel time on entry flow, which for convenience we represent as a flexible function of the natural logarithm of entry flow,  $f(\log E_t^n)$ . Second, the impact of recent flow imbalances via current queue length,  $q_t^n$ , will be closely related to recent past values of travel times; we therefore approximate it by including in our reduced-form equation two lagged values of travel time,  $T_{t-1}^n$  and  $T_{t-2}^n$ . It is important for our later simulations to recognize that these lagged travel times represent congestion dynamics and therefore play a significant role in the response of the system to any perturbation of entry flow. Third, the impact of the queue will depend strongly on recent past flow differences, which we proxy in some specifications by including the variable:

$$Q_t^n = \text{Max}\{E_{t-1}^n - F_{t-1}^n, 0\}. \quad (6)$$

Fourth, downstream blockage will affect travel time through the term  $g_1(\cdot)$  in (1); we approximate this effect by including a flexible function  $g(\log D_t^{n+1})$  in the reduced-form equation for  $T_t^n$ . Finally, we include link-specific constants and two control variables,  $W$  and  $H$ , as explained in Section 4.1.

We represent most variables by logarithms, except for  $Q$  which is often zero. The result is the following empirical equation:

$$\begin{aligned} \log T_t^n &= \beta_0^n + \beta_1 \log T_{t-1}^n + \beta_2 \log T_{t-2}^n + f(\log E_t^n) + g(\log D_t^{n+1}) + h(Q_t^n) \\ &+ \beta_W^n W_t^n + \beta_H^n H_t^n + \varepsilon_t^n \end{aligned} \quad (7)$$

with  $E_t^n$  measured by (5). As discussed later, some of the right-hand-side variables, namely  $E_t^n$ ,  $D_t^{n+1}$ , and  $Q_t^n$ , are endogenous.

The implied steady-state speed-density relationship is seen by substituting  $T = T_{-1} = T_{-2} \equiv \bar{T}$  and  $\varepsilon=0$  into (7), and solving with other variables held steady at values  $\bar{D}^{n+1}$ ,  $\bar{Q}^n$ , and  $\bar{W}^n$ . The result is:

$$\log \bar{T}^n = \frac{\beta_0^n + f(\log \bar{E}^n) + g(\log \bar{D}^{n+1}) + h(\bar{Q}^n) + \beta_w \bar{W}^n + \beta_H \bar{H}^n}{1 - \beta_1 - \beta_2} \quad (8)$$

provided  $-1 < \beta_1 + \beta_2 < 1$ , a condition that is necessary for dynamic stability and which we find true empirically in every case. Thus the effect of an exogenous shock to steady-state entry flow is determined from:

$$\frac{\partial \log \bar{T}^n}{\partial \log \bar{E}^n} = \frac{f'(\log \bar{E}^n)}{1 - \beta_1 - \beta_2}.$$

It is worth noting that by distinguishing between entry flow and exit flow, our formulation solves one of the dilemmas of empirical specification of speed-flow functions. Engineering realism suggests a functional form with a maximum possible flow, such as a backward-bending speed-flow curve. But such a function cannot tell us what happens when quantity demanded exceeds capacity; furthermore, it leads to unstable and nonsensical apparent equilibria when interacted with certain demand curves. This is because flow is typically treated as a single variable, depicting both the flow that determines congestion (a supply relationship) and the quantity of travel chosen at a given level of congestion (a demand relationship). But then the backward-bending part of the speed-flow relationship makes the supply curve downward-sloping, as though one could improve conditions by adding more cars to the link. In our formulation, we can think of entry flow as quantity demanded; it can exceed exit capacity without contradiction because there are entrances and exits along the link and queue lengths can change so as to absorb imbalances between entry flow and exit capacity. We hope that the net effect of these factors is captured by the dynamics in (7) and of the terms involving queuing variable  $Q$ .

### 3. Data

The data are collected through the period January 16 – May 8, 2007 on the freeway network in South-East Denmark.<sup>4</sup> The 91.1 km network links the cities of Odense in the east to Vejle in the north and Kolding in the west, as shown in Figure 1. It includes the Lillebælt Bridge, over which flows all road traffic between Copenhagen (east of Odense) and continental Denmark and

---

<sup>4</sup> We are grateful to the Danish Road Directorate for providing these data.

Germany. Cameras are placed near each intersection, dividing the network into 15 pieces, with data recorded separately for the two directions giving observations for each of 30 one-way links. The links range from 1.7 to 11.9 kilometers in length and two to three lanes in width. Data are recorded for five-minute intervals.

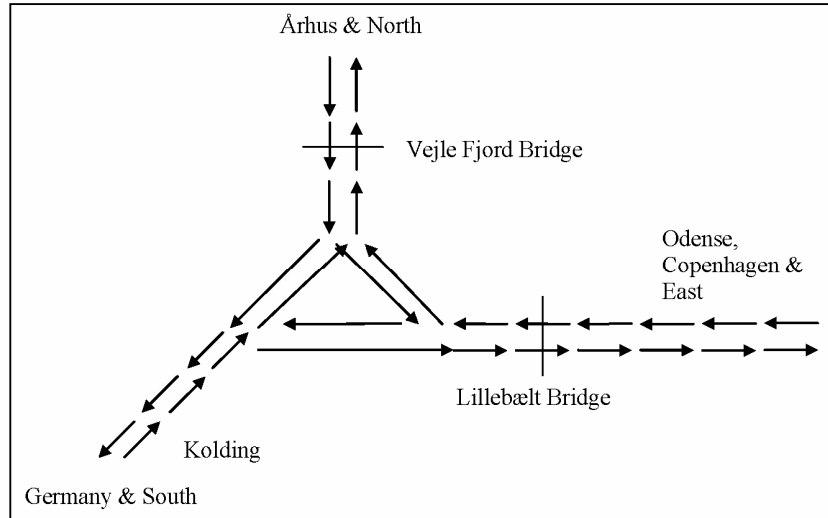


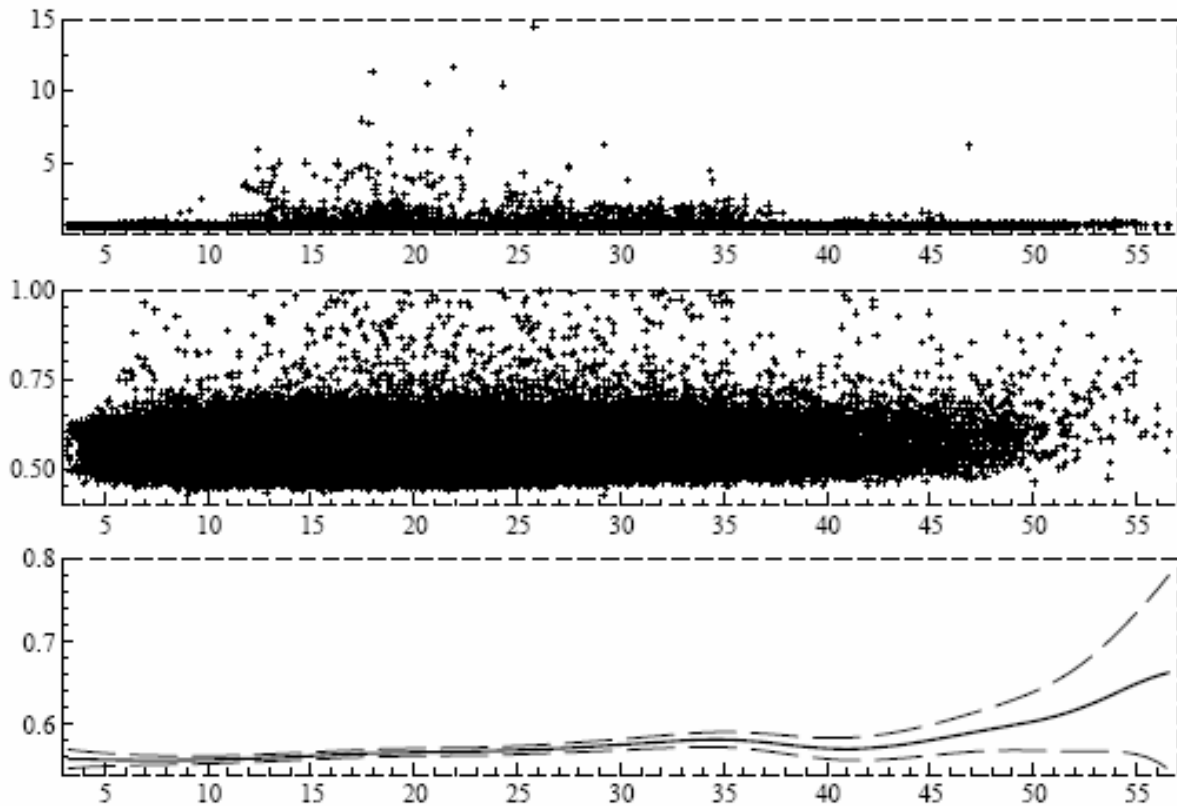
Figure 1 Network layout

We use data for all links for which we have observations on the link itself and on the first and second links upstream and downstream of it. We also require that the upstream link is a single link, in order to ensure that entry flow is unambiguous and to ease the computation of the relevant variables. This yields 246,230 observations from nine one-way links. For every five-minute observation period, the data record the exit flows and average travel times for both light and heavy vehicles, the distinction between vehicle types being approximate as it is based on the license plate. An observation is omitted when the exit flow is less than 10 vehicles per five minutes. We compute traffic flow in passenger car equivalents (pce) using a conversion factor of 2.25 pce per truck. Travel times have been divided by distance and are expressed in minutes per kilometer, while flows are divided by number of lanes and expressed in pce per lane per minute.

Figure 2 plots the observations of travel time against entry flow, with the latter averaged over a one-hour period. This and later plots of the same type show, in the upper panel, a scatter plot of the data and, in the lower panel, a kernel smooth including the mean and 95 percent pointwise confidence band. Using a normal density kernel, the bandwidth for the smooth here and later has been set to 5 percent of the range of the independent variable, which in Figure 2 is flow. The smoothed mean indicates that average travel time mostly increases slowly with flow

up to a flow of about 40 pce/lane/min, after which it rises more steeply. The overall average travel time in the sample is 0.57 min/km, corresponding to a speed of 105 km/h. Although most observations are in the lower-flow region, we also have many observations of larger flows, which of course are important for measuring congestion effects.

The scatter plot reveals that there is a very large dispersion of travel times: most observations are near the average but a considerable number are much larger. We believe these observations with high travel times are real and therefore we include them in the analysis; most of them occur at low entry flows, probably indicating conditions where entry flow is blocked by queues forming at bottlenecks within or downstream of the link in question. In order to reveal more detail in the region with most data, Figure 2 and later similar figures includes a middle panel showing data within a restricted vertical range.



**Figure 2** Travel time against hourly average entry flow

Figure 3 plots the entry flow against time of day. There are morning and afternoon peaks even though the data include both weekdays and weekends. Data are mostly missing during the hours 1:00-5:00 a.m. when there is too little traffic for reliable measurement.

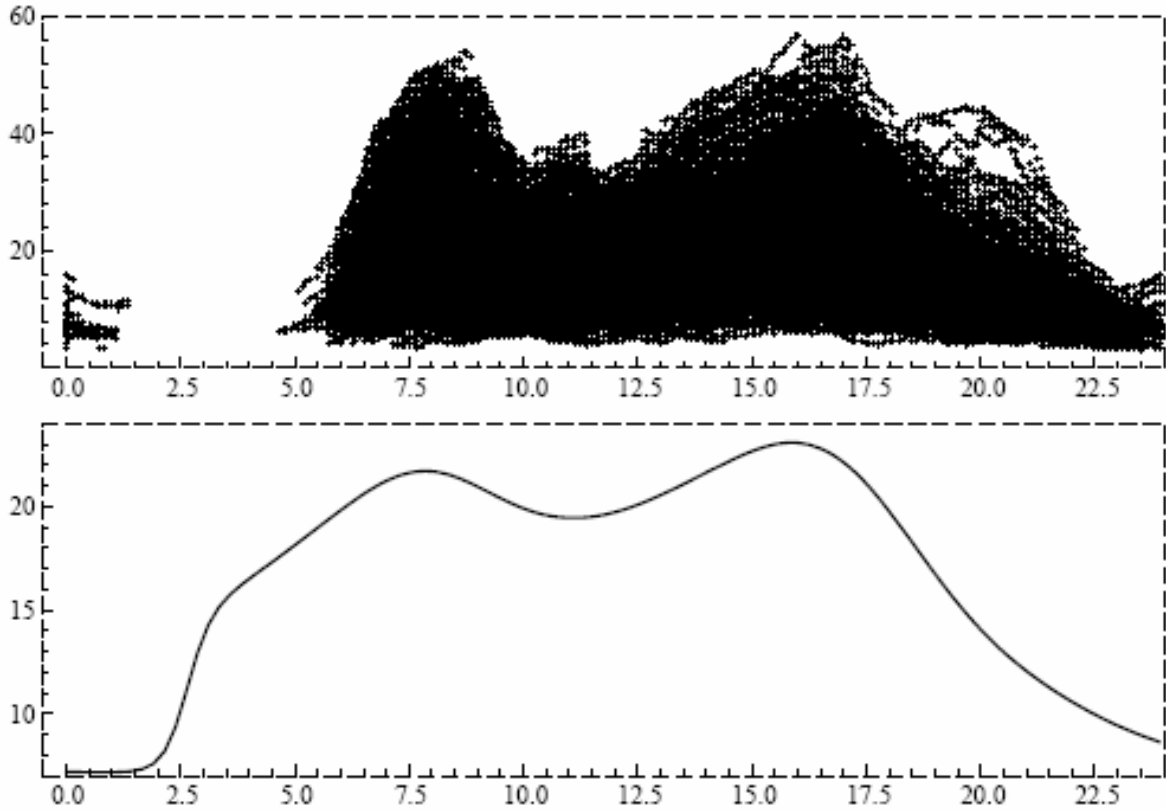


Figure 3 Hourly average flow against time of day

## 4. Empirical results

### 4.1 Model specification

We need to specify the model in (7), describing the current travel time on a link as a function of lagged travel times on the same link, entry flow, queue length, downstream blockage, and controls  $W$  and  $H$ .

The function  $f(\cdot)$  is expected to rise slowly at low entry flows, then steeply at some value approximating the capacity of an expressway lane. After some experimentation, we find a simple piecewise linear function with one breakpoint works well. Similarly, we use one breakpoint for  $g(\cdot)$  and two for  $h(\cdot)$ . We also experimented with cubic functions for  $f$ ,  $g$ , and  $h$ ; but those models do not fit as well, and also the cubic functions are overly sensitive to our numerous low-congestion observations and in fact display regions with wrong-sign derivatives. We did, however, use the estimated cubic functions visually to choose our breakpoints.

We use two types of controls. First, travel time on a link is affected by weather and other conditions not related to flow. We utilize the log travel time experienced concurrently on the same roadway in the opposite direction to control for these effects and label this variable  $W$ . This variable is averaged over three periods around the current time interval. Second, the speed for trucks is more restricted than it is for passenger cars; hence we include a variable  $H$  equal to the share of heavy vehicles measured in pce in the current exit flow.

We estimate link-specific fixed effects as well as link-specific parameters for control variables  $W$ . All other parameters are common across links.

Table 1 presents some descriptive statistics for the variables in the estimated equations.

**Table 1 Descriptive statistics**

	$\log T_{t-1}^n$	$(\log D_t^{n+1})$	$Q_t^n$	$\log D_t^{n+1}$	$W_t^n$	$H_t^n$
Mean	-0.5774	2.902	6.160	2.082	0.5200	0.5813
Median	-0.5913	2.961	0.000	2.081	0.5431	0.5508
Maximum	2.842	4.226	81.41	4.588	0.9183	19.40
Minimum	-0.8519	0.8310	0.000	.08948	0.1106	0.4208
Std. Dev.	0.1374	0.4997	10.08	0.5477	0.1428	0.3071
Skewness	5.836	-0.4590	1.991	.01343	-0.3641	22.45
Kurtosis	78.01	2.871	7.275	2.972	2.409	723.4

Finally, we turn to endogeneity of variables. According to the discussion in Section 2, we must regard entry flow, queue, and downstream blockage as endogenous since these variables are all affected by current congestion. We therefore use an instrumental variables (IV) estimator, which requires us to specify instrumental variables that are correlated with the endogenous variables but uncorrelated with the current residual in (7). We use the following two variables as instruments: the flow two links upstream of the current link, and the density two links downstream. The rationale for these variables is that they influence entry flow, queue size, and upstream density directly, but they are unlikely to be correlated with the residual in (7) because blockages seldom if ever are observed to extend across more than two links. We include also lags and some powers of these two instruments in order to gain as much power as possible in

explaining the endogenous variables, while testing to avoid weak instruments as explained in the next section.

## 4.2 Estimation results

We present estimates from three models based on piecewise linear specifications of functions  $f$ ,  $g$  and  $h$ , all using instrumental variables unless otherwise noted. The function  $f$  involving entry flow has a breakpoint at 40 pce/lane/min, which corresponds to the point in Figure 2 where travel time begins to rise and which just slightly exceeds the Danish design standard for lane capacity.<sup>5</sup> The function  $g$  involving downstream density is zero until a density of 50 pce/lane/km and linear from there. To interpret this breakpoint value, note that it corresponds to a point where downstream flow divided by downstream speed equals 50: for example, to a flow at capacity of 50 pce/lane/min and a speed of 1 km/min (60 km/h) which is roughly half free-flow speed.

The specification of the function  $h(\cdot)$  varies across models. In model M1,  $h$  is a piecewise linear function (with two pieces) in  $Q$ , defined as the positive part of the sum of two lagged differences between entry and exit flow — a natural extension of (6). Model M2 replaces  $Q$  by its first constituents, namely the first lagged values of entry and exit flows, entered as logarithms, estimating a separate coefficient for each. Finally, model M3 omits the  $Q$  variable altogether.

Results are shown in Table 2. The three models are estimated in EViews by two stage least squares (TSLS). They yield an adjusted R-square of about 0.5 and a Durbin-Watson statistic close to 2, indicating little autocorrelation of the residuals. Table 2 furthermore shows the result of estimating model M3 using OLS.

All three models portray stable and statistically significant dynamics. The coefficients for first and second lags of travel time are positive, very significant, and sum to less than one (about 0.7). These values imply that the remaining coefficients should be multiplied by about  $1/(1-0.7) \approx 3.3$  to get the values that apply when the model is solved in steady state, as shown in equation (8).

---

<sup>5</sup> The Danish design standard states an ideal capacity of 2300 pce/lane/h, or 38.3 pce/lane/min ([www.vejregler.dk](http://www.vejregler.dk)).

**Table 2 Estimation results**

Dependent variable: natural logarithm of travel time per km (min/km)								
Model:	M1		M2		M3		M3 OLS	
Variable	Coeff.	t-Stat.	Coeff.	t-Stat.	Coeff.	t-Stat.	Coeff.	t-Stat.
Const.	-0.261	-34.2	-0.253	-32.0	-0.264	-33.6	-	-
T <sub>-1</sub>	0.379	75.5	0.375	74.8	0.381	76.8	0.395	108.4
T <sub>-2</sub>	0.309	62.6	0.308	61.2	0.310	63.5	0.323	89.8
lnE	0.007	3.8	0.075	5.4	0.008	4.2	0.005	5.4
(lnE-ln(40))*1 <sub>{E&gt;40}</sub>	0.232	4.0	0.178	3.3	0.144	2.7	0.054	4.0
(Q-20)*1 <sub>{Q&gt;20}</sub>	-0.001	-4.0						
(Q-40)*1 <sub>{Q&gt;40}</sub>	-0.001	-1.7						
lnE <sub>-1</sub>			-0.070	-5.6				
lnF <sub>-1</sub>			-0.001	-0.7				
(lnD <sup>n+1</sup> -ln(50))*1 <sub>{D<sup>n+1</sup>&gt;50}</sub>	3.583	5.0	4.449	6.2	3.679	5.1	0.849	13.5
H	0.045	13.5	0.042	12.5	0.044	13.6	0.037	13.4
link-specific constants	yes		yes		yes		yes	
link-specific const's * W	yes		yes		yes		yes	
Number of observations	66902		67352		67777		67777	
Adjusted R-squared	0.503		0.495		0.509		0.524	
Sum squared resid (SSR)	645.214		663.418		652.246		632.680	
Durbin-Watson stat	2.083		2.049		2.089		2.160	
Second-stage SSR	627.699		631.086		633.877			

Note: 1<sub>{ }</sub> denotes the indicator function for the event in the curly brackets.

Other control variables are also stable across models. The coefficient for  $H$ , the share of heavy vehicles in the exit flow, indicates that the travel time of heavy vehicles is 13–15 percent larger than for light vehicles in the same traffic stream.<sup>6</sup> With nine links included, there are nine link-specific effects of control variable  $W$ . The latter control variables almost all have

<sup>6</sup> That is, for given values of other right-hand-side variables, the travel times for truck and car,  $T_T$  and  $T_C$ , are related by  $T_T/T_C = \exp[\beta_H/(1-\beta_1-\beta_2)] \approx 1.13-1.15$ . This may be somewhat too small, as the speed limit for trucks is 80 km/h compared to 110 or 130 km/h for cars. Actual speeds tend to be higher. We tried including interactions between the share of heavy vehicles and the functions for entry flow and downstream density, but these interactions were jointly insignificant. A more complicated alternative would be to develop a model with travel time for cars and trucks as separate dependent variables.

statistically significant effects, typically in the range of 0.02–0.15, indicating that travel time on the opposing link is mildly correlated with travel time on the link in question.<sup>7</sup>

Turning to the variables of main interest, consider first the role of our queuing proxy,  $Q$ , in explaining travel time. Model M1 represents the effect of  $Q$  as two linear pieces, one for  $Q$  between 20 and 40 and one for  $Q$  above 40. We expect a positive relationship because an internal queue should cause delay; however, the estimated coefficients are both negative. In model M2, we replace  $Q$  by the entry and exit flow of the last period; we see that the lagged exit flow becomes insignificant and that the lagged entry flow receives a large negative parameter, while the parameter for current entry flow (already positive in model M1) becomes much larger. We conclude that this variable does a poor job of capturing the effect of internal queuing, which is not altogether surprising given the discussion in section 2.1. In particular, the entry flow is not measured but is approximated in (5) as a weighted average of past exit flows from the upstream link; and we have no information on how many vehicles enter and exit the freeway along the way. There is also the possibility that the effect of internal queuing is just not very strong in our dataset.

In model M3, we therefore discard the internal queuing variable. The results are reassuringly similar to model M1 except that the effect of downstream density is greater. We therefore consider this our most reliable model. Model M3 shows steady-state elasticities exceeding the one-period elasticities by a factor of  $(1-0.381-0.310)^{-1}=3.2$ . Entry flow has a significant positive effect on travel time, with a steady-state elasticity about  $0.008*3.2\approx 0.026$  at entry flows less than 40 pce/lane/min and a much larger elasticity of  $0.152*3.2\approx 0.49$  for larger flows. The coefficient for the downstream density implies a large steady-state elasticity of  $3.679*3.2=11.9$  when density is greater than the breakpoint. Our results thus confirm that queue spillbacks can be an important contribution to congestion, as has long been assumed throughout the engineering literature (e.g. May 1990).

Just how high a degree of statistical significance should we expect from this model? We note that the effective sample for the congestion variables is much smaller than the full sample, because most observations show no congestion (Fig. 1). Therefore, the moderate asymptotic t-

---

<sup>7</sup> In addition the estimation procedure effectively estimates a fixed effect model, but without explicitly estimating the fixed-effect coefficients, by subtracting from each independent variable its mean value (across time) for a given link.

statistics we find for the associated parameters, typically between four and six, seem satisfactory. Furthermore, we recognize that the specification with two lagged values of travel time is only one of many types of dynamics that could be present. If we re-estimate the same model but allowing for first-order autocorrelation, we find it difficult to achieve convergence; but it appears that these key coefficients are not stable and the estimated autocorrelation, although small ( $\sim 0.1$ ), is statistically significant. Thus, our subsequent calculations of external costs are sensitive to the assumed time-series properties of the model.

The last columns of Table 2 show for comparison the results of estimating M3 using OLS; that is, without taking endogeneity of entry flow and downstream density into account. While most parameter estimates are largely unaffected, those corresponding to the endogenous variables change markedly as is expected when endogeneity is present and important.

The instruments were chosen by removing instruments, one by one, from the full list of potential instruments described earlier, until the Sargan test indicated acceptance of over-identifying restrictions. The corresponding regression of residuals against instruments yielded a significance level of 0.59, indicating that residuals are not seriously correlated with instruments. This is in agreement with the maintained hypothesis of model M3.

To check the strength of the instruments, we furthermore carried out the first-stage regressions in a separate procedure, regressing the endogenous variables on all exogenous variables and testing the joint significance of the instruments (i.e., the significance of the exogenous variables that do not enter directly in the model). For each of the endogenous variables we found that the instruments were very significant in explaining the endogenous variables.<sup>8</sup> The first-stage estimates are shown in the Appendix.

## 5. Calculation of marginal external costs

This section presents the calculation of marginal external cost (*mec*). We begin by observing that with entry flow  $E$  and travel time  $T$ , the internal cost (ignoring monetary costs) is  $T$  and the total cost of all users is  $TE$ . Then we may find the *mec* by differentiating total cost with respect to entry flow and subtracting internal cost. This yields

---

<sup>8</sup> The F-statistics are 1693.7 for variable  $(\ln E)$ ; 222.1 for variable  $((\ln E - \ln(40)) * 1_{\{E > 40\}})$ ; and 11.5 for variable  $((\ln D^{n+1} - \ln(50)) * 1_{\{D_{n+1} > 50\}})$ . Corresponding significance levels in the F-distribution are virtually zero.

$$mec = \frac{\partial(TE)}{\partial E} - T = \frac{\partial T}{\partial E} E. \quad (9)$$

We now turn to our preferred empirical model. It is possible in principle to simulate with the model in order to compute an estimate of the marginal external cost associated with adding a vehicle to one link during some time interval. There are, however, numerous complications associated with doing this and the computational cost is high. We shall therefore use a simplified approach that allows for easier computation.

We take the model to be representative of a generic link and hence omit superscripts relating to a specific link. The generic link will act as current as well as downstream link (the latter denoted by superscript +1). Solving for steady-state travel time as a function of other steady-state variables (and ignoring constants and the error term), we can write the model as follows:

$$\log T = \gamma_1 \log E + \gamma_2 \log\left(\frac{E}{40}\right) \cdot 1_{\{E^n > 40\}} + \gamma_3 \log\left(\frac{D^{+1}}{50}\right) \cdot 1_{\{D^{+1} > 50\}} \quad (10)$$

where  $\gamma$  represents a coefficient from Table 2 divided by  $(1-\beta_1-\beta_2)$  and where the notation  $1_{\{\}}$  denotes the indicator function for the event in the curly brackets.

We note from (10) that a given flow affects travel time on two different sections. First, flow in the current section affects that section's current travel time, through the terms involving  $\gamma_1$  and  $\gamma_2$ . Second, by applying (10) to the next section upstream, we see that flow in the current section affects upstream travel time, via the term involving  $\gamma_3$ , if current-section density is above 50 vehicles per lane-km (or is raised by the flow in question to that level). The latter effect may involve a lag, so for simplicity we consider steady states. Thus  $mec$  in (9) (the total effect of adding to flow on the current link) is the sum of two components:  $mec_E$ , concerning the effect on travel time on the current link, and  $mec_D$ , concerning the effect on travel time on the upstream link. We ignore further linkages between the travel time and flow on links.<sup>9</sup>

---

<sup>9</sup> This decomposition is analogous to those of Yang and Huang (1998), who also explicitly consider upstream and downstream links, and to that of Mun (1999), who divides the link being analyzed into a queued portion and a portion subject to normal congestion.

We calculate both components by differentiating the relevant terms of (10) with respect to  $E$  and using (9), keeping track of when we are considering the effect on the current or the upstream link from the link where we change  $E$ . The first component is:

$$mec_E = (\gamma_1 + \gamma_2 \cdot 1_{\{E>40\}})T. \quad (11)$$

The second component is:

$$\begin{aligned} mec_D &= \frac{\partial T}{\partial D^{+1}} \cdot \frac{\partial D^{+1}}{\partial E^{+1}} \cdot E^{+1} \\ &= \gamma_3 \cdot 1_{\{D^{+1}>50\}} \cdot \frac{T \cdot E^{+1}}{D^{+1}} \cdot \frac{\partial D^{+1}}{\partial E^{+1}}. \end{aligned} \quad (12)$$

The last derivative in (12) is:

$$\frac{\partial D^{+1}}{\partial E^{+1}} = \frac{\partial(T^{+1}E^{+1})}{\partial E^{+1}} = \left( \frac{\partial T^{+1}}{\partial E^{+1}} E^{+1} + T^{+1} \right). \quad (13)$$

Since the current and the downstream are the same generic link, we may omit the superscripts (+1) to obtain the following simplified expression for  $mec_D$  to be

$$mec_D = \gamma_3 (mec_E + T) \cdot 1_{\{D>50\}} \quad (14)$$

Combining (11) and (14),

$$\begin{aligned} mec &= mec_E + mec_D \\ &= (\gamma_1 + \gamma_2 \cdot 1_{\{E>40\}})T + \gamma_3 (mec_E + T) \cdot 1_{\{D>50\}}. \end{aligned} \quad (15)$$

We compute both components of (15) for every observation in our sample, thus depicting for each value of  $n$  and  $t$  what the  $mec$  would be if the flow, travel time, and downstream density for that observation were maintained for several periods. An advantage of using (15) in this way is that the formula uses realisations of  $T$  including error terms. This matters for the result as the dependency of  $T$  on error terms in (7) is nonlinear and so the distribution of error terms is important. The present formula preserves this information in a way that is easy to handle.

Figure 4 presents scatter plots and a smoothed mean of the  $mec$  against the entry flow, where the entry flow is expressed as an hourly average. Each data point on the scatter corresponds to a five minute interval on a section in the network. The scatter in the vertical direction is therefore due to both variation in entry flow within this one-hour average, and (more importantly) to the random term in travel time, which appears in equation (15) for  $mec$ . The discontinuity of the derivative of the fitted model is visible as a vertical gap between a cloud of

low *mec*'s and a cloud of larger *mec*'s. (The positive skewness of the distribution of marginal external cost reflects the positive skewness of the distribution of travel times.) The smooth of the *mec* may be considered to be an estimate of the expected *mec* conditional on the hourly average entry flow.

The *mec* is initially small and rises slowly until an entry flow of about 30 pce/lane/min. At this point the *mec* of a vehicle is about 0.15 min/km, which corresponds to about 25 percent of the average travel time (and a much larger proportion of the private congestion delay). From this point, the *mec* rises more steeply and at a flow of 40, *mec* has reached 0.3 min/km or about a little more than half the average travel time. At 50 pce/lane/min, *mec* has risen to about 0.7 min/km which is more than the (increased) travel time. This result confirms the view, expressed in many economic models of congestion, that external cost rises slowly at first, then rapidly as the entry flow approaches and then exceeds the capacity of a highway. This rising *mec* confirms common perceptions and is quite important for the welfare effects of pricing policies.

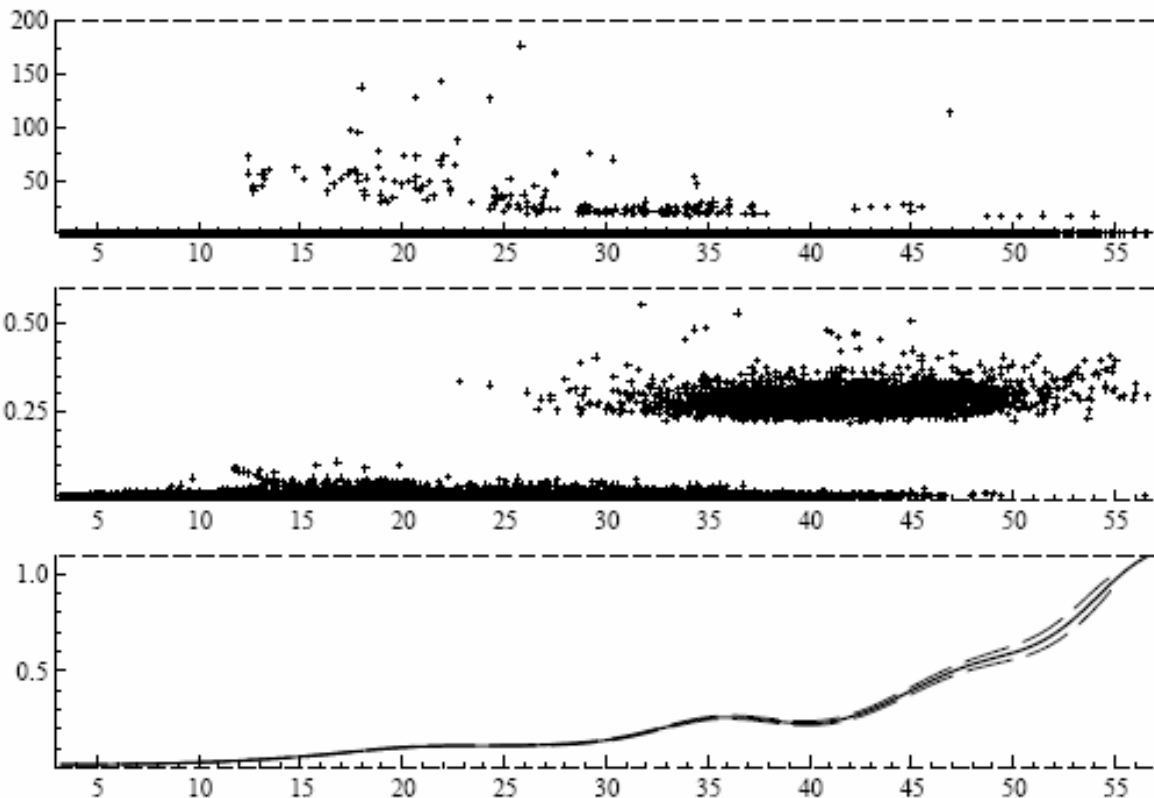


Figure 4 *mec* against hourly average entry flow

Figures 5 and 6 present the two components of the *mec*. It seems the component reflecting entry flow congestion,  $mec_E$ , is extremely variable and its average dominates until the component reflecting downstream congestion,  $mec_D$ , becomes the larger.

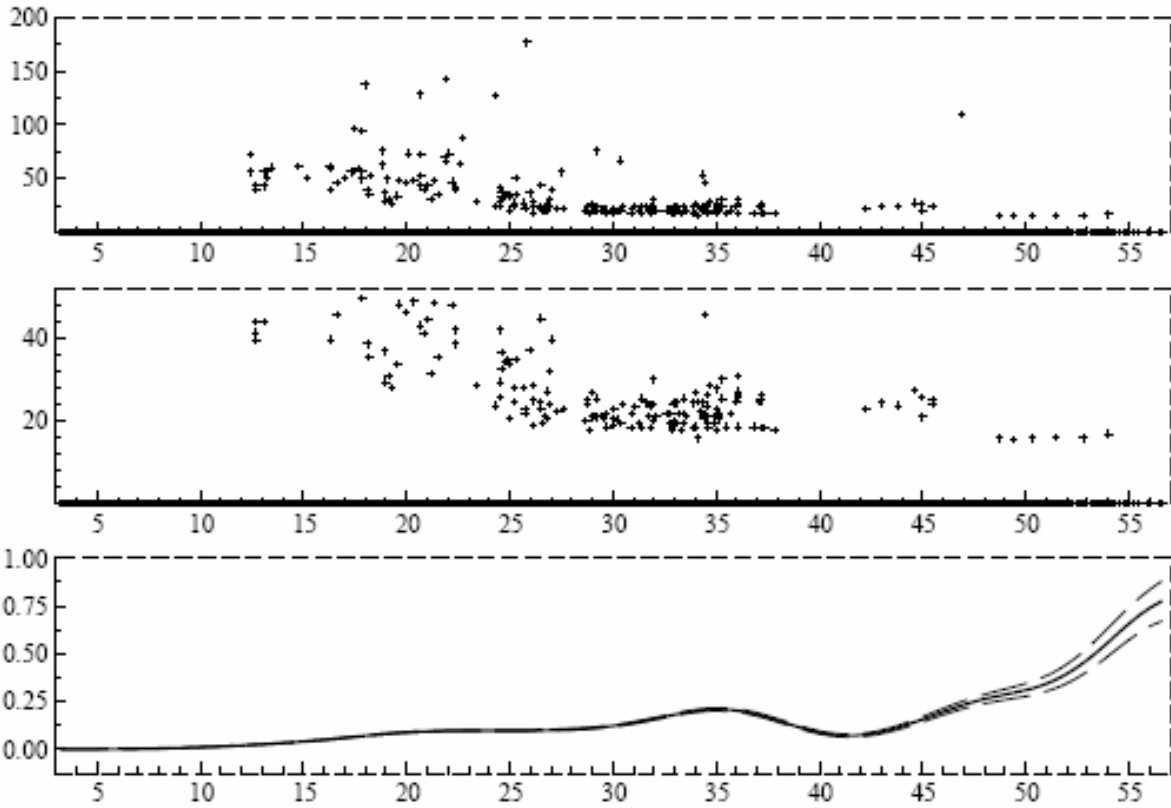


Figure 5  $mec_E$  against hourly average entry flow

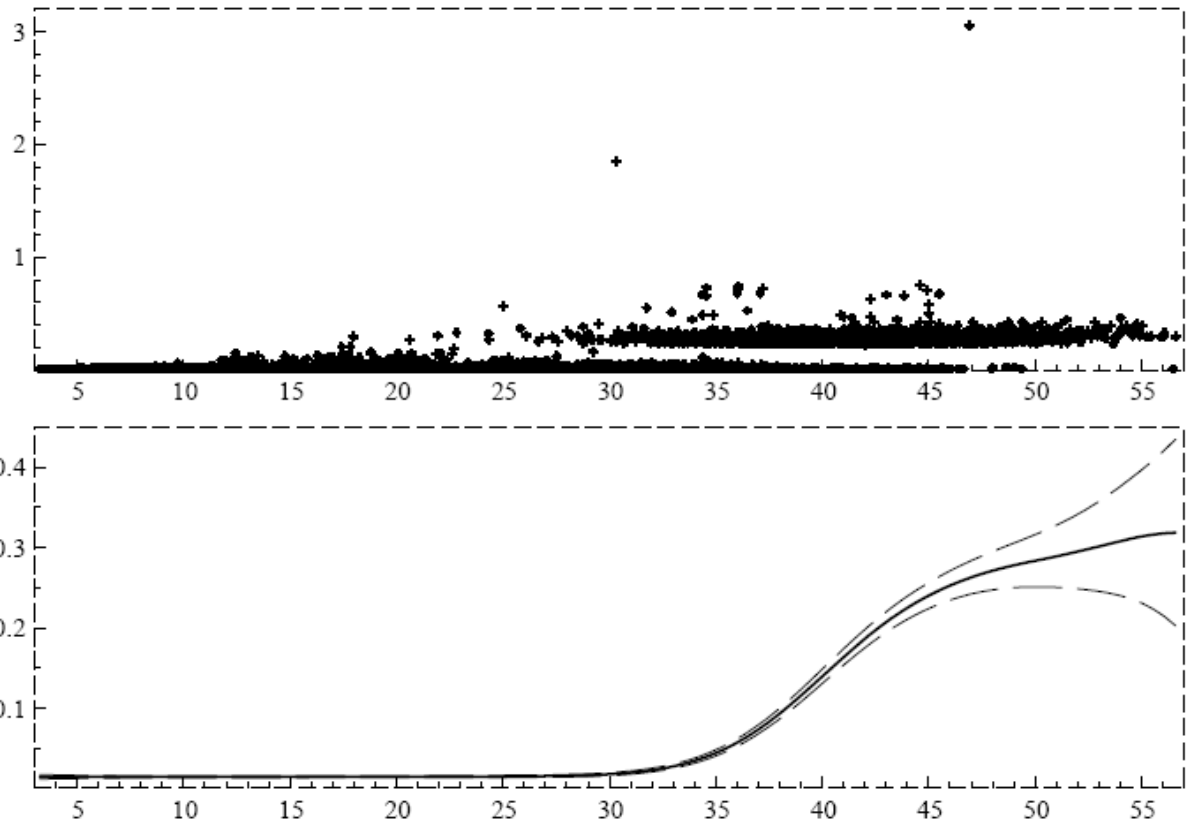
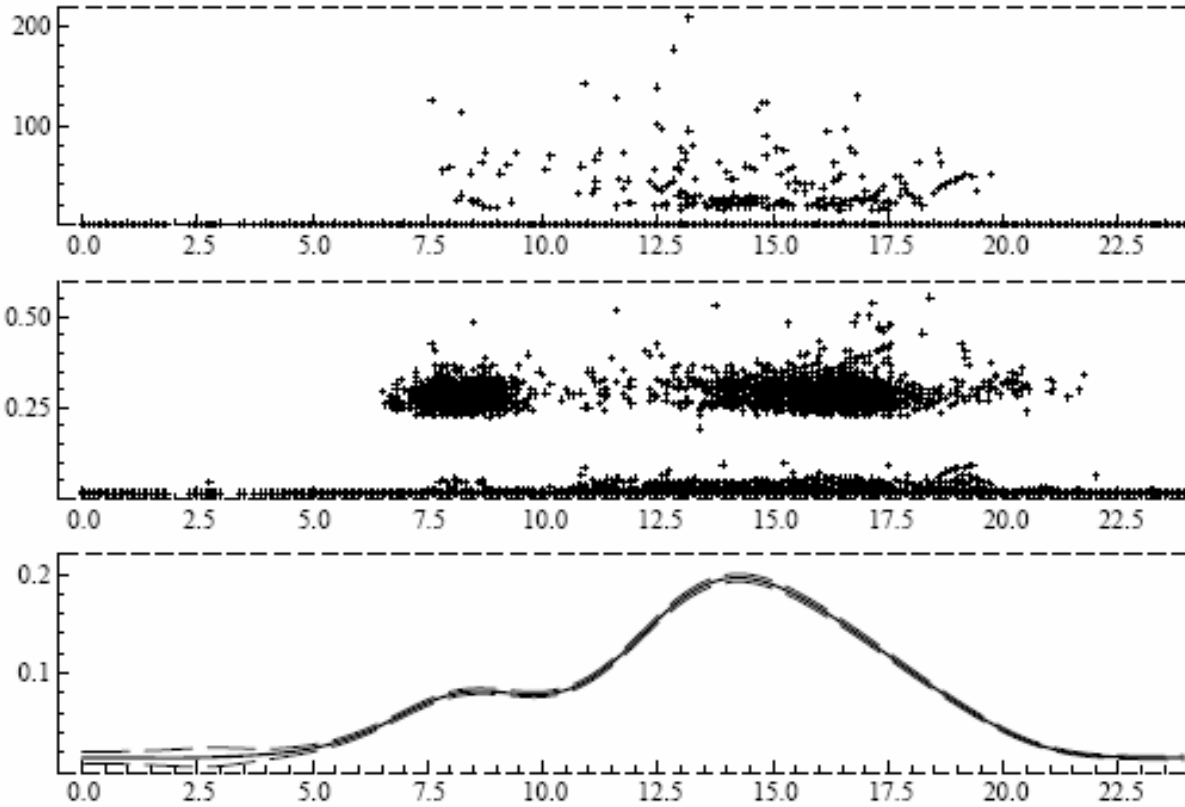


Figure 6  $mec_D$  against hourly average entry flow

Finally, Figure 7 shows the *mec* against the time of day. Its average follows the peaks in traffic and seems to be highest at about 0.2 min/km at around 3 p.m. Evidently, on this network the lower-flow situations are most common even at the peaks, causing the average *mec* to be well below the values shown in the cloud of calculated points at high average flows. Of course, there are many individual data points where *mec* is much higher than this, a reminder that marginal external cost can vary a lot due to randomness in conditions.



**Figure 7** *mec* against time of day

We noted in Section 4 that our empirical results are apparently sensitive to the specification of autocorrelation in the error terms. Furthermore, while the lagged dependent variables in (7) are intended to capture dynamic effects involving internal queues operating in a regime of hypercongestion, they could also be proxies for autocorrelation resulting from persistent influences not captured by our control variables  $W$  and  $H$ . In that case, the time-series correlations in the data would not be explained by congestion dynamics, and thus would not raise the external costs. Our best assessment is that lagged travel times have a robust effect and therefore the dynamics are real, but that the magnitude of the direct effect of entry flow and downstream density are uncertain due to specification uncertainties.

## 6. Conclusions

This paper has contributed to the measurement of the marginal cost of freeway congestion in several ways. Simultaneity of speed and traffic flow is likely to be important when there is congestion. An indication of this issue is found in Figure 2, which seems to show that average travel time may even decrease as flow increases. This would imply a negative marginal external cost of adding vehicles to flow which is blatantly nonsensical if one considers only travel time.

Using our structural model, we have argued for the use of observations from other links in the same network as instruments for flow in an equation describing travel time as a function of flow. Thus we believe we are able to go some way in tackling the simultaneity issue. Our results using these instruments indicate a positive and increasing marginal external cost in accordance with the a priori expectation. But we acknowledge the sensitivity of these results with respect to specification of the time-series properties of the residuals.

Another feature of our model is the effect of downstream congestion on the travel time on the current link. This effect is found to be empirically significant and important for the marginal external cost.

The resulting model seems plausible. So do the estimates of the marginal external cost at various levels of flow and during the day. From the structural model we derive an expression for the marginal external cost that is easy to compute for each observation in the sample.

## References

- Ardekani, Siamak, and Robert Herman (1987), "Urban Network-Wide Traffic Variables and Their Relations," *Transportation Science*, 21: 1-16.
- Cassidy, Michael J. and Robert L. Bertini (1999), "Some traffic features at freeway bottlenecks," *Transportation Research Part B*, 33: 25-42.
- Daganzo, Carlos F., Michael J. Cassidy and Robert L. Bertini (1999), "Possible explanations of phase transitions in highway traffic," *Transportation Research Part A*, 33: 365-379.
- Institute of Traffic Engineers (1982), *Transportation and Traffic Engineering Handbook*, Englewood Cliffs, New Jersey: Prentice-Hall.
- Kerner, B.S. and H. Rehborn (1997), "Experimental properties of phase transitions in traffic flow," *Physical Review Letters* 79: 4030-4033.
- May, Adolf D. (1990), *Traffic Flow Fundamentals*, Upper Saddle River, NJ: Prentice-Hall.

- May, Anthony D., S.P. Shepherd, and J.J. Bates (2000), "Supply Curves for Urban Road Networks," *Journal of Transport Economics and Policy*, 34: 261-290.
- McDonald, John.F., Edward d'Ouille, and Louis Nan Liu (1999) *Economics of Urban Highway Congestion and Pricing*, Kluwer.
- Mun, Se-il (1999) "Peak-load pricing of a bottleneck with traffic jam," *Journal of Urban Economics* 46: 323-349.
- Small, Kenneth A. and Xuehao Chu (2003), "Hypercongestion," *Journal of Transport Economics and Policy* 37: 319-352.
- Small, Kenneth A. and Erik T. Verhoef (2007), *The Economics of Urban Transportation*, London and New York: Routledge.
- Steimetz, Seiji S.C., and David Brownstone (2007), "Estimating commuters' 'value of time' with noisy data: a multiple imputation approach," *Transportation Research Part B*, 39: 865-889.
- Verhoef, Erik T. (2001), "An integrated dynamic model of road traffic congestion based on simple car-following theory: Exploring hypercongestion," *Journal of Urban Economics* 49: 505-542.
- Yang, Hai, and Hai-Jun Huang (1998), "Principle of marginal-cost pricing: How does it work in a general road network?," *Transportation Research Part A* 32(1): 45-54.

## Notation

$F$  = exit flow (pce/min per lane)

$E$  = entry flow (pce/min per lane)

$T$  = travel time (minutes/km)

$D$  = density (pce/lane-km)

$q$  = size of internal queue on link (pce/lane)

$Q$  = proxy for size of internal queue on link (pce/lane)

$W$  = travel time on control section (proxy for weather, etc.)

$H$  = the share of heavy vehicles in the exit flow

$L$  = length of link (km)

$\Delta t$  = width of time interval, min (=5 in our data)

$t$  = time period (integer)

$n$  = section (larger numbers are downstream)

## Appendix A: Approximating entering flow from observed upstream flows

The relevant upstream flows are those during the current and previous time periods: just one previous time period in the case of short sections (those that take less than five minutes to traverse), and two previous time periods in the case of longer sections. (No section takes longer than two time periods to traverse, so we need not consider three lags.) Thus we construct a flow variable equal to a weighted average of those three observed upstream flows, with the weights equal to the proportions of vehicles that could be expected to have been observed during the current and previous time period, respectively:

$$F_t^{*n} = w_t F_t^{n-1} + w_{t-1} F_{t-1}^{n-1} + w_{t-2} F_{t-2}^{n-1}.$$

In order not to introduce endogeneity into the flow variable, we compute these weights using the average speed on the entire network,  $S^*$ , expressed in km/min. Consider the vehicles exiting link  $n$  during the five-minute time interval  $t$ , which we take to begin at time 0 and end at time 5. For a link with length  $L \leq 5S^*$ , all the vehicles exiting before time  $L/(5S^*)$  entered the section during interval  $t-1$ , while the rest entered during interval  $t$ ; so  $w_{t-1} = L/(5S^*)$  and  $w_t = 1 - w_{t-1}$ . For a longer link, all the vehicles exiting before time  $-5 + L/(5S^*)$  entered during interval  $t-2$ , the rest during interval  $t-1$ ; so  $w_{t-2} = L/(5S^*) - 1$  and  $w_{t-1} = 1 - w_{t-2}$ . We can summarize for both cases as follows:

$$\begin{aligned} w_t &= \text{Max}\{0, [1 - L/(5S^*)]\} \\ w_{t-2} &= \text{Max}\{0, [L/(5S^*) - 1]\} \\ w_{t-1} &= 1 - w_t - w_{t-2} \end{aligned}$$

## Appendix B: First-stage regression for Model M3

Table B1 gives results for the first-stage regressions, *i.e.* those explaining the endogenous right-hand-side variables in (7), for the two-stage least squares estimation of Model M3 of Table 2.

**Table B1. First-stage regression results for Model M3**

Dependent variable:	lnE		$(\ln E - \ln(40)) * 1_{\{E > 40\}}$		$(\ln D^{n+1} - \ln(50)) * 1_{\{D^{n+1} > 50\}}$	
Variable	Coeff.	t-Stat.	Coeff.	t-Stat.	Coeff.	t-Stat.
C	1.316	93.8	-0.026	-18.4	0.007	22.1
T <sub>-1</sub>	0.002	0.2	0.008	7.8	0.004	19.5
T <sub>-2</sub>	-0.047	-4.7	0.000	-0.5	0.004	19.7
$(\ln F^{n-2} - \ln 10) * 1_{\{F^{n-2} > 10\}}$	0.263	38.9	-0.006	-8.6	0.000	-0.8
$(\ln F^{n-2} - \ln 20) * 1_{\{F^{n-2} > 20\}}$	0.009	0.6	0.023	15.9	0.000	0.9
$(\ln F^{n-2} - \ln 30) * 1_{\{F^{n-2} > 30\}}$	0.050	1.5	0.070	21.0	-0.001	-1.6
$(\ln F^{n-2} - \ln 40) * 1_{\{F^{n-2} > 40\}}$	-0.130	-1.6	0.035	4.4	0.001	0.7
$(\ln F^{n-2} - \ln 50) * 1_{\{F^{n-2} > 50\}}$	-0.241	-1.2	-0.188	-9.4	-0.002	-0.4
$\ln(F^{n-2})_{-2}$	0.356	82.4	0.008	17.7	0.000	-0.8
$\ln(D^{n+2})_{-1}$	0.128	23.8	0.001	1.9	0.000	-0.3
$\ln(D^{n+2})_{-2}$	0.086	16.3	0.002	3.7	0.000	0.9
$(\ln(D^{n+2}) - \ln 15) * 1_{\{D^{n+2} > 15\}}$	0.026	2.5	0.011	10.1	0.000	0.2
$(\ln(D^{n+2}) - \ln 30) * 1_{\{D^{n+2} > 30\}}$	-0.416	-9.7	-0.050	-11.7	0.016	16.6
$(\ln(D^{n+2}) - \ln 45) * 1_{\{D^{n+2} > 45\}}$	0.464	7.2	0.044	6.9	-0.019	-13.3
H	-0.037	-4.8	0.004	4.6	-0.002	-12.1
link-specific constants	yes		yes		yes	
link-specific const's * W	yes		yes		yes	
Number of observations	67777		67777		67777	
Adjusted R-squared	0.702		0.162		0.043	