

# Estimating Latent Processes on a Network From Indirect Measurements

Edoardo M. AIROLDI and Alexander W. BLOCKER

In a communication network, point-to-point traffic volumes over time are critical for designing protocols that route information efficiently and for maintaining security, whether at the scale of an Internet service provider or within a corporation. While technically feasible, the direct measurement of point-to-point traffic imposes a heavy burden on network performance and is typically not implemented. Instead, indirect aggregate traffic volumes are routinely collected. We consider the problem of estimating point-to-point traffic volumes,  $\mathbf{x}_t$ , from aggregate traffic volumes,  $\mathbf{y}_t$ , given information about the network routing protocol encoded in a matrix  $A$ . This estimation task can be reformulated as finding the solutions to a sequence of *ill-posed* linear inverse problems,  $\mathbf{y}_t = A \mathbf{x}_t$ , since the number of origin-destination routes of interest is higher than the number of aggregate measurements available.

Here, we introduce a novel multilevel state-space model (SSM) of aggregate traffic volumes with realistic features. We implement a naïve strategy for estimating unobserved point-to-point traffic volumes from indirect measurements of aggregate traffic, based on particle filtering. We then develop a more efficient two-stage inference strategy that relies on model-based regularization: a simple model is used to calibrate regularization parameters that lead to efficient/scalable inference in the multilevel SSM. We apply our methods to corporate and academic networks, where we show that the proposed inference strategy outperforms existing approaches and scales to larger networks. We also design a simulation study to explore the factors that influence the performance. Our results suggest that model-based regularization may be an efficient strategy for inference in other complex multilevel models. Supplementary materials for this article are available online.

**KEY WORDS:** Approximate inference; Ill-posed inverse problem; Multilevel state-space model; Multistage estimation; Network tomography; Origin-destination traffic matrix; Particle filtering; Polytope sampling; Stochastic dynamics.

## 1. INTRODUCTION

A pervasive challenge in multivariate time series analysis is the estimation of nonobservable time series of interest  $\{\mathbf{x}_t : t = 1, \dots, T\}$  from indirect noisy measurements  $\{\mathbf{y}_t : t = 1, \dots, T\}$ , typically obtained through an aggregation or mixing process,  $\mathbf{y}_t = \mathbf{a}(\mathbf{x}_t) \forall t$ . The inference problem that arises in this setting is often referred to as an *inverse*, or *deconvolution*, problem (e.g., Hansen 1998; Casella and Berger 2001; Meister 2009) in the statistics and computer science literatures, and qualified as *ill-posed* because of the lower dimensionality of the measurement vectors with respect to the nonobservable estimands of interest. Ill-posed inverse problems lie at the heart of a number of modern applications, including image super-resolution and positron emission tomography where we want to combine many two-dimensional images in a three-dimensional image consistent with two-dimensional constraints (Shepp and Kruskal 1978; Vardi, Shepp, and Kaufman 1985); blind source separation where there are more sound sources than sound tracks (i.e., the measurements) available (Liu and Chen 1995; Lee et al. 1999; Parra and Sajda 2003); and inference on cell values in contingency tables where two-way or multiway margins are prespecified (Bishop, Fienberg, and Holland 1975; Dobra, Tebaldi, and West 2006).

We consider a setting in which high-dimensional multivariate time series  $\mathbf{x}_{1:T}$  mix on a network. Individual time series correspond to traffic directed from a node to another. The aggregation

process encodes the routing protocol—whether deterministic or probabilistic—that determines the path traffic from any given source follows to reach its destination. This type of mixing can be specified as a linear aggregation process  $A$ . This problem setting leads to the following sequence of ill-posed linear inverse (or deconvolution) problems,

$$\mathbf{y}_t = A \mathbf{x}_t, \quad \text{s.t. } \mathbf{y}_t, \mathbf{x}_t \geq 0 \quad \text{for } t = 1, \dots, T, \quad (1)$$

since the observed aggregate traffic time series are low dimensional,  $\mathbf{y}_t \in \mathbb{R}^m$ , while the latent point-to-point traffic time series of interest are high dimensional,  $\mathbf{x}_t \in \mathbb{R}^n$ . Thus, the matrix  $A_{m \times n}$  is rank deficient,  $r(A) = m < n$ , in this general problem setting.

The application to communication networks that motivates our research is (*volume*) *network tomography*; an application originally introduced by Vardi (1996), which has quickly become a classic (see, e.g., Vanderbei and Iannone 1994; Tebaldi and West 1998; Cao et al. 2000; Coates et al. 2002; Medina et al. 2002; Liang and Yu 2003a; Zhang et al. 2003b; Airoldi and Faloutsos 2004; Castro et al. 2004; Lakhina et al. 2004; Lawrence et al. 2006b; Fang, Vardi, and Zhang 2007; Blocker and Airoldi 2011). An established engineering practice is at the root of the inference problem in network tomography. Briefly, the availability of point-to-point (or origin-destination (OD)) traffic volumes over time is critical for reliability analysis (e.g., predicting flows and failures), traffic engineering (e.g., minimizing congestion), capacity planning (e.g., forecasting requirements), and security management (e.g., detecting anomalous traffic patterns). While technically possible, however, the direct measurement of point-to-point traffic imposes a heavy burden on network performance and is never implemented, except for special purposes over short time periods. Instead,

Edoardo M. Airoldi is Assistant Professor, Department of Statistics, and Alfred P. Sloan Research Fellow, Harvard University, Cambridge, MA 02138 (E-mail: [airoldi@fas.harvard.edu](mailto:airoldi@fas.harvard.edu)). Alexander W. Blocker is Ph.D. candidate, Department of Statistics, Harvard University, Cambridge, MA 02138 (E-mail: [ablocker@fas.harvard.edu](mailto:ablocker@fas.harvard.edu)). The authors thank Jin Cao and Matthew Roughan for providing point-to-point traffic volumes, and the referees for providing valuable suggestions. This work was partially supported by the Army Research Office under grant 58153-MA-MUR, and by the National Science Foundation under grants IIS-1017967, DMS-1106980, and CAREER IIS-1149662.

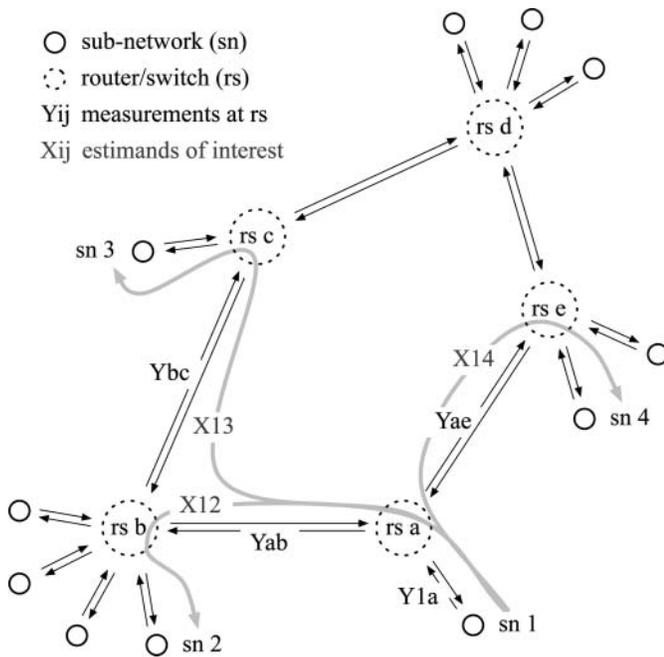


Figure 1. Illustration of the mathematical quantities in network tomography. Traffic  $x_{13}$  from sn 1 to sn 3 contributes to counter  $Y_{1a}$  into rs a, to counter  $Y_{ab}$  into rs b, to counter  $Y_{bc}$  out of rs b (which is the same as counter  $Y_{bc}$  into rs c), and to counter  $Y_{c3}$  out of rs c. Traffic volumes on these counters are recorded every few minutes. This routing information is summarized in the column of the routing matrix  $A$  that corresponds to the OD traffic volume  $x_{13}$ .

indirect aggregate traffic volumes are routinely collected. As a consequence, network engineers must solve the ill-posed linear inverse problems in Equation (1) to recover point-to-point traffic. We give pointers to a vast literature that spans statistics, computer science, and operations research in Section 1.1.

Figure 1 provides an illustration of a communication network and the key mathematical quantities in the application to network tomography. Dashed circles,  $rs\ a-e$ , represent routers and switches. Solid circles,  $sn\ 1-11$ , represent subnetworks. Intuitively, messages are sent from a subnetwork (origin) to another (destination) over the network. Routers and switches are special-purpose computers that quickly scan the messages and route them according to a prespecified routing protocol. Black arrows represent physical cables connecting routers and switches to subnetworks and indicate the direction in which traffic flows. On each router, a set of (software) counters measures aggregate traffic volumes,  $y_{ij}$ , corresponding to incoming and outgoing cables, routinely (every 5 min). The traffic recorded by each of these counters is the sum of a known subset<sup>1</sup> of nonobservable point-to-point traffic,  $x_{ij}$ , represented by gray arrows, over the same time window. For example, in Figure 1, traffic volumes  $x_{12}, x_{13}, x_{14}$  all contribute to counter  $y_{1a}$ , and traffic volumes  $x_{12}, x_{13}$  both contribute to counter  $y_{ab}$ . To establish a formal connection between measurements  $y_{ij}$  and estimands  $x_{ij}$ , it is convenient to simplify notation. Let us order all the (from-to) counter measurements collected over a 5-min window, into a vector  $y_t \in \mathbb{R}^m$ . We have  $m = 32$  measurements in Figure 1. Let

us also order<sup>2</sup> all the nonobservable point-to-point traffic volumes of interest over a 5-min window, into a single vector  $x_t$ . We have  $n = 11^2$  point-to-point traffic volumes in Figure 1. Using this more compact notation, we can write  $y_{it} = \sum_{j=1}^n A_{ij} x_{jt}$ , where  $t$  denotes time, and the matrix  $A_{m \times n}$  is built using information about the prespecified routing protocol. In particular,  $A_{ij} = 1$  if point-to-point traffic  $i$  contributes to counter  $j$  and  $A_{ij} = 0$  otherwise, in the case of deterministic routing. More complicated routing schemes, including probabilistic routing and dynamic load-balancing protocols that minimize the expected congestion, can also be formulated in terms of Equation (1), as discussed in Section 6.2.

From a statistical perspective, the application to communication networks we consider has additional features that make the inference task harder, and more interesting, than in a traditional deconvolution problem. First, we have low-dimensional observations,  $x_t$ , and high-dimensional estimands,  $y_t$ . In Figure 1, for example,  $m = 32$  and  $n = 121$ . In a general network with  $d$  subnetworks,  $m = O(d)$  is often orders of magnitude lower than  $n = O(d^2)$ , depending on the redundancy of some counters and on whether we are interested in traffic volumes on all possible OD pairs. Second, the space where the estimands live is highly constrained. We prove in Section 2.2 that the solution space is a convex polytope of dimension  $n - m$ . The dimensionality of this convex polytope gives the true complexity of the problem, in a computational sense. Working in a constrained solution space helps the inference to a point (see, e.g., Theorem 1). We gain additional information from modeling traffic dynamics explicitly. Sampling from such an extremely constrained solution space, however, proves to be a challenge. We approach this sampling problem by combining a random direction sampler (Smith 1984) with model-based regularization and a sequential sample-importance-resample-move (SIRM) particle filter (Gilks and Berzuini 2001).

In this article, we introduce a new dynamic multilevel model for aggregate traffic volumes, which posits two latent dynamic processes, in Section 2. The first is a heavy-tailed traffic process, in which the amount of traffic on each OD route is proportional to its variability up to a scaling factor shared by all OD routes. The second is an additional error process for better capturing near-zero traffic volumes. We carry out inference via a sequential SIRM particle filter, and we develop a novel two-stage strategy (inspired by Clogg et al. 1991), in Section 3. A transformation of the heavy-tailed layer of the multilevel model can be embedded into a Gaussian state-space formulation with identifiable parameters. We use the fit for such a reformulation to calibrate informative priors for key parameters of the multilevel model, and to develop an efficient particle filter that is statistically efficient, numerically stable, and scales to large problems. In Section 4, we show that the proposed methods are more accurate than published state-of-the-art solutions on two time series datasets. In Section 5, we then design experiments that combine real and simulated data to investigate comparative performance. In Section 6, we offer remarks on modeling, inferential and computational challenges with the proposed methods, and discuss limitations and extensions.

<sup>1</sup>This information is encoded by the routing protocol.

<sup>2</sup>Any such two orderings can be arbitrary and defined independently of each other. Different pairs of orderings will lead to different  $A$  matrices.

The R package `networkTomography` includes the two unpublished datasets we analyze, as well as robust code implementing all the seven methods we compare. It is available at the Comprehensive R Archive Network (CRAN) at <http://cran.r-project.org/>.

## 1.1 Related Work

Applied research related to the type of problems we consider can be traced back to literature on transportation and operations research (Bell 1991; Vanderbei and Iannone 1994). There the focus is on estimating a single set of OD traffic volumes,  $y$ , from integer-valued traffic counts over time,  $x_t$ . The line of research in statistics with application to communication networks is due to Vardi (1996) who coined the term *network tomography* by extending the approach to positron emission tomography by Shepp and Vardi (1982). In this latter setting, statistical approaches may be able to leverage knowledge about a physical process, explicitly specified by a model, to assist the inference task. In the network tomography setting, in contrast, we can only rely on knowledge about the routing matrix and statistics about traffic time series.

From a technical perspective, Vardi (1996) developed an estimating equation framework to estimate a single set of OD traffic volumes from time series data; the same data setting and estimation task considered in the transportation and operations research literature. Tebaldi and West (1998) developed a hierarchical Bayesian model that can be fit at each epoch independently, thus recovering time-varying OD traffic volumes. They pointed out that the hardness of the problem lies in having to sample from a very constrained solution space. Informative priors are advocated as a means to mitigate issues with nonidentifiability and multimodality that arise when making inference from aggregated traffic volumes at each point in time. In previous work (Airoldi and Faloutsos 2004), we extended their approach by explicitly modeling complex dynamics of the nonobservable time series. Cao et al. (2000) developed a local likelihood approach to attack the nonidentifiability issue. They developed a Gaussian model with a clever parameterization that leads to identifiability of the point-to-point traffic volumes, if they are assumed independent over a short time window—approximately 1 hr. Cao et al. (2001) extended this approach to inference on large networks by adopting a divide-and-conquer strategy. Zhang et al. (2003b) developed gravity models that can scale to large networks and used them to analyze point-to-point traffic on the AT&T backbone in North America. This approach was extended by Fang, Vardi, and Zhang (2007) and Zhang et al. (2009). Works in this area by Soule et al. (2005) and Erramilli, Crovella, and Taft (2006) provide slightly different approaches to this class of problems. Recent reviews of this literature are available (Castro et al. 2004; Lawrence et al. 2006b).

One of the key technical problems that we face during inference is that of sampling solutions from a convex polytope. In this sense, the problem of sampling a feasible set of OD traffic volumes given aggregate traffic is equivalent to that of sampling square tables given row and column totals, when the routing matrix corresponds to a star network topology. As we consider more complicated topologies, the equivalence still holds for more elaborate specifications of marginal totals. Airoldi and

Haas (2011) characterized such a correspondence using projective geometry and the Hermite normal form decomposition of the routing matrix  $A$ . Leveraging this equivalence, the iterative proportional fitting procedure (IPFP; Deming and Stephan 1940; Fienberg 1970) provides a baseline for the traffic matrix estimation at each epoch in Section 5.2. Other approaches to the problem of sampling tables given row and column margins include a sequential Markov chain Monte Carlo (MCMC) approach (Chen et al. 2005), a dynamic programming approach that is very efficient for matrices with a low maximum marginal total (Harrison *in press*; Miller and Harrison *in press*), and sampling strategies based on algebraic geometry (Diaconis and Sturmfels 1998; Dobra *in press*) or on an explicit characterization of the solution polytope (Airoldi and Haas 2011).

A related body of work on tomography focuses on the problem of *delay (network) tomography*, in which the times traffic reaches/leaves the routers are recorded at the router level, instead of the volumes (Presti et al. 2002; Liang and Yu 2003b; Lawrence, Michailidis, and Nair 2006a; Deng et al. 2012). However, inference in delay tomography has a different structure from inference in volume tomography, which we focus on in this article.

## 2. A MODEL OF MIXING TIME SERIES ON A NETWORK

Given  $m$  observed traffic counters over time,  $y_t = \{y_{it} : i = 1, \dots, m\}$ , the aggregate traffic loads, we want to make inference on  $n$  nonobservable point-to-point traffic time series,  $x_t = \{x_{jt} : j = 1, \dots, n\}$ . The routing scheme is parameterized by the routing matrix  $A$ , of size  $m \times n$ . Without loss of generality, we consider the case of a fixed routing scheme. In this case, the matrix  $A$  has binary entries; element  $A_{ij}$  specifies whether counter  $i$  includes the traffic on the point-to-point route  $j$ . Extensions to probabilistic routing and dynamic protocols for congestion management are discussed in Section 6.2.

The main observation that informs model elicitation is that the measured traffic volumes,  $y_t$ , are heavy tailed and sparse. For instance, peak traffic may be very high during certain hours of the day, and traffic is often zero during night hours. We develop a multilevel state-space model (SSM) to explain such a variability profile of the observed aggregate traffic volumes. The proposed multilevel model involves two latent processes:  $\{\lambda_t : t \geq 1\}$  at the top of the hierarchy and  $\{x_t : t \geq 1\}$  in the middle of the hierarchy. The observation process is at the bottom of the hierarchy. Intuitively, we posit a heavy-tailed  $\{\lambda_t\}$  process and a thin-tailed  $\{x_t | \lambda_t\}$  process, specifying  $x_t | \lambda_t$  as additive error around  $\lambda_t$ , constrained to be positive. The key point is that we need both temporal correlation and independent errors to induce positive density for near-zero traffic. In previous work, we assumed a heavy-tailed  $\{\lambda_t\}$  process and a heavy-tailed  $\{x_t\}$  process, specifying  $x_t | \lambda_t$  as independent log-normal variation conditionally on  $\lambda_t$  (Airoldi and Faloutsos 2004). This set of choices, however, leads to some computational instability during inference when actual point-to-point traffic is zero (or nearly zero), as the likelihood for  $x_t$  had zero density at  $x_{j,t} = 0$ .

In detail, we posit that each point-to-point traffic volume  $x_{j,t}$  has its own time-varying intensity  $\lambda_{j,t}$ . This underlying intensity

evolves through time according to a multiplicative process

$$\log \lambda_{j,t} = \rho \log \lambda_{j,t-1} + \varepsilon_{j,t},$$

where  $\varepsilon_{j,t} \sim N(\theta_{1,j,t}, \theta_{2,j,t})$ . Such a process leads to heavy-tailed traffic volumes that are not sparse. Moreover, small differences between low traffic volumes receive quite different probabilities under the log-normal model. Thus, conditional on the underlying intensity, we posit that the latent point-to-point traffic volumes  $x_{j,t}$  follow a truncated normal error model,

$$x_{j,t} | \lambda_{j,t}, \phi_t \sim \text{TruncN}_{(0,\infty)}(\lambda_{j,t}, \lambda_{j,t}^\tau (\exp(\phi_t) - 1)),$$

where  $\tau$  and  $\phi_t$  regulate temporally independent variation. The mean-variance structure of the error model is analogous to that of a log-normal distribution for  $\tau = 2$ ; in particular, if  $\log(z) \sim N(\mu, \sigma^2)$ ,  $\mathbb{E}(z) = \exp(\mu + \sigma^2/2)$ , and  $\text{var}(z) = \exp(2\mu + \sigma^2) \cdot (\exp(\sigma^2) - 1)$ . Thus,  $\lambda_{j,t}$  is analogous to  $\exp(\mu + \sigma^2/2)$ , and  $\phi_t$  is analogous to  $\sigma^2$ . The observed aggregate traffic is obtained by mixing point-to-point traffic according to the routing matrix,  $y_t = Ax_t$ . The model specification is complete by placing diffuse independent log-Normal priors on  $\lambda_{j,0}$ . We also place priors on  $\phi_t$  for stability, assuming  $\phi_t \sim \text{Gamma}(\alpha, \beta_t/\alpha)$ .

This multilevel structure provides a realistic model for the aggregate traffic volumes we measure, which are both heavy tailed and sparse. The error model induces sparsity while maintaining analytical tractability of the inference algorithms, detailed in Section 3, by decoupling sparsity control from the bursty dynamic behavior. The log-Normal layer provides heavy-tailed dynamics and replicates the intense traffic bursts observed empirically, whereas the truncated Normal layer allows for near-zero traffic levels with nonnegligible probability. By combining these two levels, we induce a posterior distribution on point-to-point traffic volumes, the estimands of interest in this problem, which can account for both extreme volumes and sparsity.

In summary, we developed a model for observed  $m$ -dimensional time series  $\{y_t\}$  mixing on a network according to a routing matrix  $A$ . The model involves two  $n$ -dimensional latent processes  $\{\lambda_t, x_t\}$ , a set of latent variables  $\{\phi_t\}$ , and constants  $\rho, \tau, \alpha, \{\beta_t\}$  and  $\{\theta_{1t}, \theta_{2t}\}$ . While the parameters  $\tau, \rho$ , and  $\alpha$  provide some flexibility to the model and can be calibrated through exploratory data analysis on the observed traffic time series, the parameters  $\{\theta_{1t}, \theta_{2t}, \beta_t\}$  are key to the inference. Strategies for parameter estimation and posterior inference are discussed in Section 3.

### 2.1 Qualitative Model Checking

As part of the model validation process, we looked at whether the simulated time series from the model in Section 2 possessed qualitative features of real time series; namely, sparse traffic localized in time, and heavy tails in distribution of the traffic loads.

We generated a number of time series using parameter values  $\tau = 2, \rho = 0.92, \theta_{1t} = 0, \theta_{2t} = \frac{2 \log(5)}{4}$  for all  $t$ . In addition, we set  $\phi_t = 0.25$ , rather than setting the constants  $\alpha, \beta_t$  underlying the distribution of  $\phi_t$ , for simplicity. These are realistic values for the constants; they were calibrated on the actual point-to-point traffic volumes from the Bell Labs dataset in Section 4.1. We used the empirical mean, standard deviation, and autocorrelation

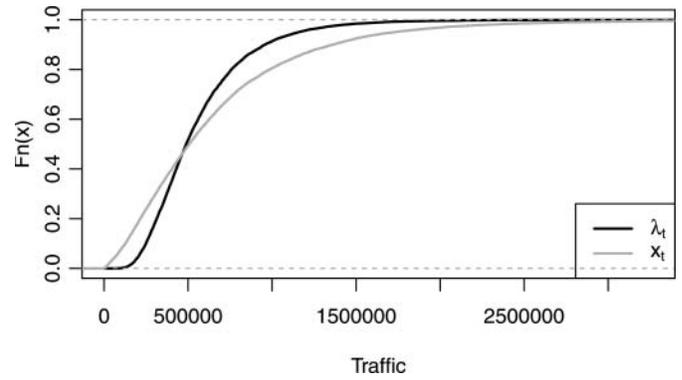


Figure 2. Comparison of CDFs for  $\lambda_{i,t}$  (solid black line) and  $x_{i,t}$  (solid gray line).

of  $\{\log x_{it}\}$  for each of these time series combined with the observed level of sparsity to create the results below.

Figure 2 shows the empirical CDF of the two latent processes  $\{\lambda_t\}$  and  $\{x_t\}$  for one simulated time series. The  $\{\lambda_t\}$  process places more mass in any  $\epsilon$  ball around zero relatively to the  $\{x_t\}$  process. The figure confirms our intuition about how the truncated Gaussian error operates.

Figure 3 shows an OD traffic time series, in the left panel, and a simulated  $\{x_t\}$  time series, in the right panel. Real point-to-point traffic volumes from the router/switch to the local sub-network were measured using special software installed on the routers, for validation purposes, courtesy of Bell Labs. The Bell Labs data are further discussed in Section 4.1. The simulated time series displays two key qualitative characteristic of the real point-to-point traffic time series. Specifically, we observe sudden traffic surges, typical for a heavy tail distribution of traffic volumes, and localized periods of low traffic, as expected from our (truncated) additive Gaussian correction.

The anecdotal findings above hold for most of the real point-to-point traffic volumes and simulated time series we considered. This suggests that the proposed model is capable of generating data that qualitatively resemble real traffic volumes. There are two important ways in which observed and simulated traffic differs, though. First, simulated point-to-point traffic peaks last longer than real traffic peaks. This is due

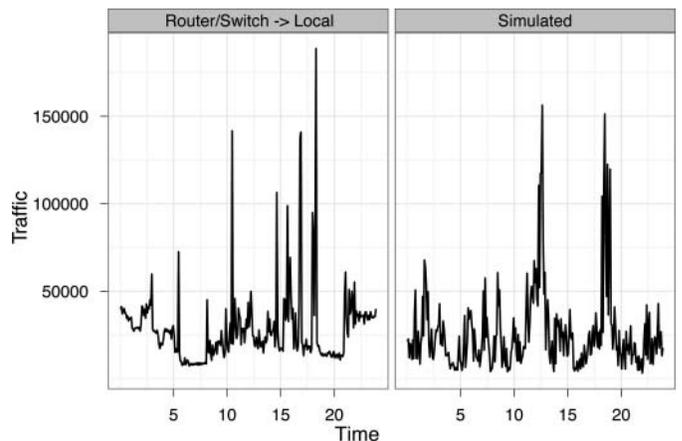


Figure 3. Actual (left panel) versus simulated (right panel) point-to-point traffic volumes.

to the autoregressive dependence in  $\lambda_{j,t}$ . Second, simulated point-to-point traffic volumes are more variable than real traffic volumes in low traffic regimes. This is due to the truncated Normal noise. Thus, the proposed model is not a perfect generative mechanism for point-to-point traffic. However, the structure the model captures is sufficient to provide useful posterior inference, as we demonstrate in Section 4.

## 2.2 Theory

Multimodality has been reported in the literature as an issue in network tomography. This issue has previously been illustrated only in a toy example; a small network with Poisson traffic (Vardi 1996). We investigated multimodality from a geometric perspective. Intuitively, our main result is that whenever a well-behaved distribution is used as a prior for the individual OD traffic volumes, however diffused, the posterior cannot have disconnected modes. The main result applies directly to the case of real-valued traffic volumes (e.g., Cao et al. 2000; Airoldi and Faloutsos 2004, and this article). For the case of integer-valued traffic volumes, analyzed by others (e.g., Vardi 1996; Tebaldi and West 1998), only a weaker condition is possible.

Consider the case of real-valued nonnegative traffic volumes. Feasible traffic volumes  $\mathbf{x}_t$  must be nonnegative and satisfy  $y_{it} \geq \sum_j A_{ij} x_{jt}$ . In other words, the space where  $\mathbf{x}_t$  lives can be characterized as the intersection between the positive orthant and  $m$  half-spaces in  $n - m$  dimensions. This is a convex polyhedron. Since both  $\mathbf{y}_t$  and  $A_{ij}$  are nonnegative, the polyhedron is bounded and the space of feasible solution is a convex polytope. The main result is a consequence of the fact that the space in which solution vectors  $\mathbf{x}_t$  to Equation (1) is a convex polytope.

*Theorem 1.* Assume  $f(\mathbf{x}_t)$  is quasi-concave. Let  $\mathbf{y}_t = A\mathbf{x}_t$ . Then,  $f(\mathbf{x}_t|\mathbf{y}_t)$  will also be quasi-concave, and will have no separated modes. The set  $\{\mathbf{z} : f(\mathbf{z}|\mathbf{y}_t) = \max_w f(\mathbf{w}|\mathbf{y}_t)\}$  is connected.

*Proof.*  $f(\mathbf{x}_t|\mathbf{y}_t) \propto I(y_t = A\mathbf{x}_t)f(\mathbf{x}_t)$ , so  $f(\mathbf{x}_t|\mathbf{y}_t)$  has support on only a bounded  $n - m$  dimensional subspace of  $\mathbb{R}^n$ , which forms a closed, bounded, convex polytope in the positive orthant. Denote this region  $B(\mathbf{y}_t)$ . Denote the mode of  $f(\mathbf{x}_t)$  as  $\hat{\mathbf{x}}_t$ . We now consider two cases.

*Case 1.*  $\hat{\mathbf{x}}_t \in B(\mathbf{y}_t)$ . Then, the mode of  $f(\mathbf{x}_t|\mathbf{y}_t)$  is also  $\hat{\mathbf{x}}_t$ .

*Case 2.*  $\hat{\mathbf{x}}_t \notin B(\mathbf{y}_t)$ . Then, we must do a little more work. Consider the level surfaces of  $f(\mathbf{x}_t|\mathbf{y}_t)$ , denoting  $C(z) = \{\mathbf{u} : f(\mathbf{u}|\mathbf{y}_t) = z\}$ . Define  $z^* = \max_{B(\mathbf{y}_t)} f(\mathbf{x}_t|\mathbf{y}_t)$ ; this is well defined and attained as  $B(\mathbf{y}_t)$  is closed. Now, denoting  $C_0(z) = \{\mathbf{u} : f(\mathbf{u}) = z\}$ , we have  $C(z) = C_0(z) \cap B(\mathbf{y}_t)$ . As  $f(\mathbf{x}_t)$  is quasi-concave, its superlevel sets  $U_0(z) = \{\mathbf{u} : f(\mathbf{u}) \geq z\}$  are convex. Thus, the superlevel sets of  $f(\mathbf{x}_t|\mathbf{y}_t)$ , denoted  $U(z) = U_0(z) \cap B(\mathbf{y}_t)$  analogously, are also convex. So, we have that the set  $U(z^*) = C(z^*)$  is convex and nonempty. Therefore, we have established that set of modes for  $f(\mathbf{x}_t|\mathbf{y}_t)$  is convex, hence connected.  $\square$

Next, consider the case of integer-valued nonnegative traffic volumes. To precisely state conditions under which pathological behavior is not possible in the integer-valued case, we need to introduce some concepts from integral geometry. A square integer matrix is defined to be *unimodular* if it is invertible as an integer matrix (so its determinant is  $\pm 1$ ). By extension, we

define a rectangular matrix to be *unimodular* if it has full rank and each square submatrix of maximal size is either unimodular or singular (its determinant is 0).

With integer-valued traffic, the inferential goal is to sample solutions to Equation (1), where  $A$  is a given unimodular  $m \times n$  matrix with  $\{0, 1\}$  entries and  $\mathbf{y}_t$  is a given integer positive vector. In the case of real-valued traffic, it was straightforward to show that the space of solutions to (1) is a convex polytope. In the case of integer-valued traffic we have

*Theorem 2* (Airoldi and Haas 2011). The space of real solutions  $x$  to equation  $y = Ax$ ,  $x \geq 0$ , is an integral polytope, whenever  $A$  is unimodular.

*Proof.* The vertices are the intersections of the affine solution space of  $Ax = y$  with the  $(n - m)$ -coordinate planes bordering the nonnegative orthant. So a vertex  $x$  has  $n - m$  zero coordinates. Let us gather the rest of the coordinates into a positive integer vector  $x'$  of dimension  $m$ . And let us gather the corresponding columns of  $A$  into a square matrix  $A_1$ ; so we get the equation  $A_1x' = y$ . If  $A_1$  was singular, the latter system would have either none or infinitely many solutions, which would contradict that  $x$  is a vertex. So  $A_1$  is unimodular and  $x' = A_1^{-1}y$ . And since  $y$  is integer,  $x'$  is also integer, and so is  $x$ .  $\square$

We can precisely characterize the space of feasible traffic volumes in the integer case, however, we cannot directly address multimodality. The concept of multiple modes and local maxima are not well defined in this setting. This result, however, provides insight into the connection between our results and the pathological case demonstrated by Vardi (1996).

The theory above helps us settle an important question about our model: how will posterior inference behave under dynamic updates? If dynamic updates were allowed to “grow” modes over time or exhibit other pathological behavior, the computation would be quite difficult and inference results would be less credible. Fortunately, this is not the case. In general, we have established that the quasi-concavity of a predictive distribution  $f(\mathbf{x}_t|\mathbf{y}_{t-1}, \dots)$  implies the quasi-concavity of the posterior  $f(\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}, \dots)$ ; thus, the set of maxima for  $f(\mathbf{x}_t|\mathbf{y}_t, \mathbf{y}_{t-1}, \dots)$  will form a convex set under the given condition. Since we initialize our model with a unimodal (quasi-concave) log-Normal distribution and impose log-Normal dynamics on the underlying intensities  $\lambda_t$ , Theorem 1 provides a useful limit on pathological behavior during inference with our model.

The situation is somewhat similar, but less constrained, in the case of integer traffic volumes, for unimodular routing matrices. While it is not known under what conditions a network routing scheme translates into a unimodular routing matrix  $A$ , the routing matrices in the cited literature are all unimodular. Thus, extreme forms of multimodality can be ruled out from the literature on dynamic network tomography in many cases. Our theory also suggests that models based upon real-valued traffic volumes will exhibit more predictable behavior under posterior updates than those based upon integer-valued volumes, making the former much more attractive for inference in cases where integer constraints provide little additional information.

### 3. PARAMETER ESTIMATION AND POSTERIOR INFERENCE

Here, we develop two inference strategies to produce estimates for the point-to-point traffic time series of interest, using the model described in the previous section. The first strategy is based on a variant of the sequential SIRM particle filter (Gilks and Berzuini 2001). This filter is simple to state and implement, but is computationally expensive due to the large number of particles needed to explore probable trajectories in high-dimensional polytopes. Details are given in Section 3.1. The second strategy combines the sequential SIRM filter with a model-based regularization step that leads to efficient particles. This strategy preferentially explores trajectories in regions of the solution polytopes with high posterior density. The model-based regularization step involves fitting a Gaussian SSM with an identifiable parameterization, which leads to informative priors for the multilevel model that are leveraged by a modified SIRM filter. Details are given in Section 3.2.

#### 3.1 An SIRM Filter for Multilevel State-Space Inference

Inference in the multilevel SSM is performed with a sequential sample-resample-move algorithm, akin to Gilks and Berzuini (2001). Its structure is outlined in Algorithm 1.

##### Sample-Importance-Resample-Move algorithm

```

for  $t \leftarrow 1$  to  $T$  do
  Sample step:
  for  $j \leftarrow 1$  to  $m$  do
    Draw a proposal
     $\log \lambda_{i,t}^{(j)*} \sim N(\theta_{1,i,t} + \log \lambda_{i,t-1}^{(j)}, \theta_{2,i,t})$ 
    Draw  $\phi_t^{(j)} \sim \text{Gamma}(\alpha, \beta_t / \alpha)$ 
    Draw  $\mathbf{x}_t^{(j)*}$  from a truncated Normal distribution
    with mean  $\boldsymbol{\mu}^* = \rho/m \sum_{j=1}^m \boldsymbol{\lambda}_{t-1}^{(j)}$  and covariance
    matrix  $\Sigma^* = (\exp(\beta_t) - 1)\text{diag}(\boldsymbol{\mu}^{*2})$  on the
    feasible region given by  $\mathbf{x}_t^{(j)*} \geq 0, y_t = A\mathbf{x}_t^{(j)*}$ 
    using Algorithm 2
  Resample our particles  $(\lambda_t^{(j)*}, \phi_t^{(j)*}, \mathbf{x}_t^{(j)*})$  with
  probabilities proportional to our weights  $w_t^{(j)}$ 
  Move each of our resampled particles  $(\lambda_t^{(j)}, \phi_t^{(j)}, \mathbf{x}_t^{(j)})$ 
  using a MCMC algorithm (Metropolis–Hastings within
  Gibbs, with proposal on  $\mathbf{x}_t$  given by Algorithm 2)
return  $(\lambda_t^{(j)}, \phi_t^{(j)}, \mathbf{x}_t^{(j)})$  for  $j \leftarrow 1$  to  $m, t \leftarrow 1$  to  $T$ 

```

**Algorithm 1:** SIRM algorithm for inference with multilevel SSM

In Algorithm 1, we use a random walk prior on the latent intensities  $\lambda_t$ . Thus, we fix  $\theta_{1,i,t} = 0$  for all  $i, t$ , and calibrate the constants  $\{\theta_{2,i,t}\}, \{\beta_t\}, \alpha, \tau$ , and  $\rho$  as discussed in Section 3.1.1.

Sampling particles that correspond to feasible trajectories, that is, to point-to-point traffic volumes  $\mathbf{x}_t$  in the convex polytope implied by  $\mathbf{y}_t = A\mathbf{x}_t$ , is nontrivial. The use of a random direction proposal on the region of feasible point-to-point traffic is a vital component of the SIRM filter.

We use the *random directions algorithm* (RDA; Smith 1984) to sample from the distributions of feasible traffic volumes on a constrained region, in the SIRM filter. This method constructs a random walk proposal on a convex region, such as the feasible

regions for  $\mathbf{x}_t$ , by first drawing a vector  $\mathbf{d}$  uniformly on the unit sphere. It then calculates the intersections of a line along this vector with the surface of the bounding region, and samples uniformly along the feasible segment of this line. Computing the feasible segment is facilitated by decomposing  $A$ . We decompose  $A$  as  $[A_1 | A_2]$  by permuting the columns of  $A$ , and the corresponding components of  $\mathbf{x}_t$ , so that  $A_1$  ( $r \times r$ ) is of full rank. Then, splitting the permuted vector  $\mathbf{x}_t = [\mathbf{x}_t^1, \mathbf{x}_t^2]$ , we obtain  $\mathbf{x}_t^1 = A_1^{-1}(\mathbf{y}_t - A_2\mathbf{x}_t^2)$ . This formulation can be used to construct an efficient RDA to propose feasible values of  $\mathbf{x}_t$ . We have included pseudocode for this algorithm in Algorithm 2.

#### Random Directions Algorithm

##### Initialization

###### begin

```

  Decompose  $A$  into  $[A_1 | A_2]$ ,  $A_1$  ( $r \times r$ ) full-rank
  Store  $B := A_1^{-1}$ ;  $C := A_1^{-1}A_2$ 

```

##### Metropolis step

###### given $\mathbf{x}_t$

###### begin

```

  Draw  $\mathbf{z} \sim N(0, I)$ ,  $\mathbf{z} \in \mathbb{R}^{c-r}$ 
  Set  $\mathbf{d} := \mathbf{z}/\|\mathbf{z}\|$ 
  Calculate  $\mathbf{w} := C \cdot \mathbf{d}$ 
  Set  $h_1 := \max\{\min_{k:w_k>0}(\mathbf{x}_{1,t})_k/w_k, 0\}$ 
  Set  $h_2 := \max\{\min_{k:d_k<0}-(\mathbf{x}_{2,t})_k/d_k, 0\}$ 
  Set  $h := \min\{h_1, h_2\}$ 
  Set  $l_1 := \max\{\max_{k:w_k<0}(\mathbf{x}_{1,t})_k/w_k, 0\}$ 
  Set  $l_2 := \max\{\max_{k:d_k>0}-(\mathbf{x}_{2,t})_k/d_k, 0\}$ 
  Set  $l := \max\{l_1, l_2\}$ 
  Draw  $u \sim \text{Unif}(l, h)$ 
  Set  $\mathbf{x}_{2,t}^* := \mathbf{x}_{2,t} + u \cdot \mathbf{d}$ ;  $\mathbf{x}_{1,t}^* = \mathbf{x}_{1,t} - u \cdot \mathbf{w}$ ;
   $\mathbf{x}_t^* = (\mathbf{x}_{1,t}^*, \mathbf{x}_{2,t}^*)$ 
  Set  $\mathbf{x}_t := \mathbf{x}_t^*$  with probability  $\min\{f(\mathbf{x}_t^*)/f(\mathbf{x}_t), 1\}$ 

```

###### return $\mathbf{x}_t$

**Algorithm 2:** RDA algorithm for sampling from  $f(\mathbf{x}_t)$ , truncated to the feasible region given by  $A \cdot \mathbf{x}_t = \mathbf{y}_t$

All draws from this proposal have positive posterior density, since they are feasible. This property allows our sampler to move away from problematic boundary regions of the extremely constrained solution polytope. In contrast, methods that use Gaussian random walk proposal rules, for instance, can perform quite poorly in these situations, requiring an extremely large number of draws to obtain feasible proposals. For example, with  $\mathbf{x}_t \in \mathbb{R}^{16}$ , it can sometimes require on the order of  $10^9$  draws to obtain a feasible particle, when using the conditional posterior from  $t - 1$  as proposal. This is a situation we encountered with alternative estimation methods described in Section 4.

**3.1.1 Setting the Constants.** To carry out inference, we must set values for the constants underlying the distributions at the top layer of the multilevel model;  $\{\theta_{1,j,t}\}, \{\theta_{2,j,t}\}, \{\beta_t\}, \alpha, \rho$ , and  $\tau$ . Choices can be evaluated using small sets of point-to-point traffic collected for diagnostic purposes, as in Section 2.1

The (fixed) autocorrelation parameter  $\rho$  drives the dynamics of  $\log \lambda_{j,t}$ . We typically set  $\rho = 0.9$ . A high value for  $\rho$  is a practically plausible assumption, as point-to-point traffic volumes tend to be highly autocorrelated in communication networks (Cao et al. 2002).

The parameter  $\tau$  controls the skew of point-to-point traffic volumes. The distribution of point-to-point traffic has been found to be extremely skewed empirically (Airoldi 2003), and this skew is comparable to the skew of the aggregate traffic volumes. Cao et al. (2000) found that the local variability of the aggregate traffic volumes is well described by  $\tau = 2$ . In our analyses, we fix  $\tau = 2$ . This assumption was checked on pilot data as in Cao et al. (2000).

The inference strategy based on the SIRM filter is amenable to a wide range of techniques for regularization. The simplest of these is a random walk prior on  $\log \lambda_t$ . For this, we fix  $\theta_{1,i,t} = 0$  for all  $t$  and set  $\{\theta_{2,i,t}\}$  by looking at the observed variability of  $\{y_t\}$ . On the datasets we consider,  $\theta_{2,i,t} = (2 \log 5)/4$  appeared reasonable based on the variability of the observed aggregate traffic. That is, we set  $\theta_2$  by rescaling the average variance of  $\log y_{j,t} - \rho \log y_{j,t-1}$  to correct for aggregation. This is a somewhat crude approach, but it provides a reasonable starting point.

The collection of constants  $\{\beta_t\}$  controls the common scale of variation in the point-to-point traffic. These constants were set by examining the observed marginal distribution of  $\{y_t\}$ . We selected  $\beta_t = 1.5$  as reasonable value based on the observed excess abundance of values near zero. Last, the constant  $\alpha$  is a fixed tuning parameter. We set it to  $n/2$  to provide a moderate degree of regularization for our inference, providing a weight equivalent to  $1/2$  of the observed data.

The random walk prior is a simple starting point for our inference and provides cues of computational issues. Its use is not recommended in practice. The inverse problem we confront in network tomography is too ill-posed for such a simplistic approach to regularization. A more refined, adaptive strategy is necessary to provide useful answers in realistic settings.

### 3.2 Two-Stage Inference

Here, we develop an inference strategy that improves the SIRM filter in Algorithm 1 by adding a regularization step that guides our inference, focusing our particle filter and sharing information across multiple classes of models. The idea is to leverage a first-stage estimation step to calibrate informative priors for key parameters in the multilevel model, in the spirit of empirical Bayes (Clogg et al. 1991). Different forms of model-based regularization are feasible (and useful) depending upon traffic dynamics and the topology of a given network. One approach is to use simple, well-established methods such as gravity-based methods (Zhang et al. 2003a,b; Fang, Vardi, and Zhang 2007). Another approach, developed below, uses a specific parameterization of a Gaussian SSM to approximate Poisson traffic. We find that these two approaches are useful in different situations (namely, local area networks and Internet backbone networks) as we discuss in Section 4.

**3.2.1 Model-Based Regularization.** Here, we describe a simple model used to calibrate key regularization parameters  $\{\theta_{1t}, \theta_{2t}, \beta_t\}$  of the multilevel SSM. We posit that  $\mathbf{x}_t$  follows a Gaussian autoregressive process,

$$\begin{cases} \mathbf{x}_t = F \cdot \mathbf{x}_{t-1} + Q \cdot \mathbf{1} + \mathbf{e}_t \\ \mathbf{y}_t = A \cdot \mathbf{x}_t + \boldsymbol{\epsilon}_t. \end{cases} \quad (2)$$

This model can be subsumed into a standard Gaussian state-space formulation, as detailed in Equation (3).

$$\begin{aligned} &= \begin{cases} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{1} \end{bmatrix} = \begin{bmatrix} F & Q \\ 0 & I \end{bmatrix} \begin{bmatrix} \mathbf{x}_{t-1} \\ \mathbf{1} \end{bmatrix} + \begin{bmatrix} \mathbf{e}_t \\ \mathbf{0} \end{bmatrix} \\ \mathbf{y}_t = [A | \mathbf{0}] \begin{bmatrix} \mathbf{x}_t \\ \mathbf{1} \end{bmatrix} + \boldsymbol{\epsilon}_t \end{cases} \\ &= \begin{cases} \tilde{\mathbf{x}}_t = \tilde{F} \cdot \tilde{\mathbf{x}}_{t-1} + \tilde{\mathbf{e}}_t \\ \mathbf{y}_t = \tilde{A} \cdot \tilde{\mathbf{x}}_t + \boldsymbol{\epsilon}_t. \end{cases} \end{aligned} \quad (3)$$

We estimate  $Q$  and cov  $\mathbf{e}_t$ , fixing the remaining parameters.  $F$  is fixed at  $\rho I$  for simplicity of estimation, with 0.1 a typical value for  $\rho$ . We also fix cov  $\boldsymbol{\epsilon}_t$  at  $\sigma^2 I$ , with 0.01 a typical value for  $\sigma^2$ . We assume  $Q$  to be a positive, diagonal matrix,  $Q = \text{diag}(\lambda_t)$ , and specify cov  $\mathbf{e}_t$  as  $\Sigma_t = \phi \text{diag}(\lambda_t)^\tau$ , where the power is taken entry-wise. We obtain inferences from this model via maximum likelihood on overlapping windows of a fixed length. We develop an inference strategy for this model in Section 4, and provide computational details in the Appendix.

The model in Equation (2) contains the local likelihood model of Cao et al. (2000) as a special case, when  $\rho = 0$ . The marginal likelihood for this model depends only upon the means and covariances of the data. A desirable property of this model is that its parameters are identifiable, under conditions analogous to those given in Cao et al. (2000), for a fixed value of  $\rho$ .

**3.2.2 Identifiability.** Identifiability in network tomography is a delicate and complex issue (Singhal and Michailidis 2007). For the proposed model, however, it suffices to consider the marginal distribution of the  $\mathbf{y}_t$ 's. Under the conditions on the routing matrix  $A$  analogous to those in Cao et al. (2000), the marginal mean and covariance of  $\mathbf{y}_t$  is an invertible function of the parameters  $\lambda$  and  $\phi$ . This argument is straightforward with the steady-state initialization discussed in the Appendix, but it extends to more general settings. In the case of steady-state initialization, the following result holds.

*Theorem 3.* Assume  $\mathbf{y}_1, \dots, \mathbf{y}_T$  are distributed according to the model given by Equation (3) and described above. Further assume that  $|\rho| < 1$  and the model is initialized from its steady state—that is,

$$\mathbf{x}_0 \sim N\left(\frac{1}{1-\rho}\lambda, \frac{\phi}{1-\rho^2}D\right),$$

where  $D = \text{diag}(\lambda_t)^\tau$ . Then,  $(\lambda, \phi, \rho)$  is identifiable under the same conditions required for the identifiability of the locally iid model of Cao et al. (2000).

*Proof.* The observations  $(\mathbf{y}_1, \dots, \mathbf{y}_T)$  are jointly normally distributed under the given model. Further, assuming  $|\rho| < 1$  and steady-state initialization,  $\mathbf{y}_t \sim N(\frac{1}{1-\rho}A\lambda, \frac{\phi}{1-\rho^2}ADA^\top + \sigma^2 I)$  marginally for  $t = 1, \dots, T$ . Define  $B$  as the matrix containing the rows of  $A$  and all distinct pairwise component-wise products of  $A$ 's rows. Fixing  $\rho$ , these marginal moments are invertible functions of  $(\lambda, \phi)$  if and only if the matrix  $B$  has full column rank by Theorem 1 of Cao et al. (2000). Further, as cov  $(\mathbf{y}_t, \mathbf{y}_{t+k}) = \frac{\phi \rho^{|k|}}{1-\rho^2} ADA^\top$ ,  $\rho$  is also identifiable from the component-wise autocorrelations of  $\mathbf{y}_t$ .  $\square$

A sufficient condition for  $B$  to have full column rank is for  $A$  to include aggregate incoming and outgoing (source and destination) traffic for each node, as discussed in Cao et al. (2000). This condition holds for all examples we consider and can be checked in practice with pilot studies; such aggregate traffic volumes are easily obtained via network management protocols such as Simple Network Management Protocol and are standard input for the widely used gravity method. Less restrictive conditions are possible based on the results of Singhal and Michailidis (2007); however, they are not needed for the situations we consider.

Next, we describe how this model is used to calibrate priors for  $\{\lambda_t\}$  and  $\{\phi_t\}$  in the multilevel SSM.

**3.2.3 Calibrating Key Regularization Parameters.** To calibrate priors for  $\{\lambda_t\}$  and  $\{\phi_t\}$  in the multilevel SSM, we follow a few steps. First obtain estimates from the Gaussian SSM in the previous section. We correct the estimates at each epoch through the IPFP to ensure positivity and validity with respect to our linear constraints. We then smooth the corrected estimates using a running median with a small window size (consisting of five observations) to obtain a final set of  $\hat{x}_t$  estimates. This smoothing step is important as it removes outlying estimates, which often originate from computational errors, from the prior calibration procedure. These outliers can otherwise degrade the effectiveness of the regularization. We have observed some sensitivity to the choice of window sizes—too broad and it smooths out bursts of traffic, too narrow and outlying estimates compromise our regularization. We selected 5 as the narrowest window that empirically removed outliers; we recommend this as a guideline for other settings. These final  $\hat{x}_t$  estimates are used to set the mean traffic intensity for  $\lambda_t$  as follows,

$$\theta_{1,j,t} = \log \hat{x}_{j,t} - \rho \log \hat{x}_{j,t-1}.$$

The variability of the traffic intensity  $\theta_{2,j,t}$  is set using the estimated variance of the final estimates  $\hat{x}_{j,t}$ . Denoting the estimated final variances with  $\hat{V}_{j,t}$ , we set  $\theta_{2,j,t}$  as follows,

$$\theta_{2,j,t} = (1 - \rho^2) \log (1 + \hat{V}_{j,t} / \hat{x}_{j,t}^2).$$

The estimated  $\{\hat{\phi}_t\}$  in the Gaussian SSM are used to calibrate the prior for the corresponding parameter  $\{\phi_t\}$  in the multilevel SSM. In particular, we set  $\beta_t = \log(1 + \hat{\phi}_t)$ . The form of this calibration is based on the log-Normal variance relationship described in Section 2. Remaining constants are calibrated as described in Section 3.1.1.

Alternative calibration approaches are possible, which use estimates from simple models to calibrate regularization parameters. For instance, in Section 4 we consider a simple gravity model in addition to the SSM described above. We take each gravity estimate to be  $\hat{x}_{j,t}$ , and we set each  $\theta_{1,j,t}$  as above. With simpler model we recommend using an empirical approach to setting  $\theta_2$ ; using the gravity model estimates, we set each  $\theta_{2,j,t}$  equal to the overall variance of  $\theta_{1,j,\dots}$ .

#### 4. EMPIRICAL ANALYSIS OF TRAFFIC DATA

Here, we present the analysis of three aggregate traffic datasets, for which OD traffic volumes were also collected with special software over a short time period. The first dataset involves traffic volumes on a local area network with 4 nodes (16 OD pairs) at Bell Labs, previously analyzed in Cao et al.

(2000). The second dataset involves traffic volumes on a local area network with 12 nodes (144 OD pairs) at Carnegie Mellon University (CMU), previously analyzed by Airoidi and Faloutsos (2004). The final dataset consists of traffic volumes from the Abilene network, an Internet2 backbone network with 12 nodes (144 OD pairs) previously analyzed in Fang, Vardi, and Zhang (2007). We use these three datasets to evaluate the proposed deconvolution methods. We compare the performance of our approach with that of several previously presented in the literature for this problem, focusing on accuracy, computational stability, and scalability.

We find that, of the seven methods we compare, the proposed methods consistently outperform all others both in terms of  $L_1$  and  $L_2$  estimation error. The empirical evaluation we provide below uses communication networks that are among the largest ever tried on this problem in the statistics and computer science literature. A quantitative evaluation is possible, since ground-truth OD traffic was laboriously collected for the three network we consider.

An R package that includes these three datasets and code to replicate the analyses below is available on CRAN, in the `networkTomography` package.

#### 4.1 Datasets

The first dataset was provided courtesy of Jin Cao of Bell Labs. We analyze the traffic volumes measured at *router1*, with four subnetworks organized as in Figure 4 (left panel). These yield eight observed aggregate traffic volumes (seven of them are independent, since the router does not send, nor receives traffic) and 16 OD traffic volumes (Cao et al. 2000). The aggregate traffic volumes are measured every 5 min over 1 day on the Bell Labs corporate network. This yields multivariate measurements at 287 points in time. The small size of this network allows us to focus on the fundamentals of the problem, avoiding scalability issues.

The second dataset was collected at the Information Networking Institute of CMU, courtesy of Russel Yount and Frank Kietzke. For the purpose of this article, given that the topology of Carnegie Mellon network is sensitive, we built a dataset of aggregate traffic volumes by mixing 2 days of OD traffic volumes on a slightly modified network topology. The network topology we use consists of 12 subnetworks, organized as in Figure 4 (right panel). These are connected by two routers, one with four of the nodes, the other with the remaining eight nodes. The

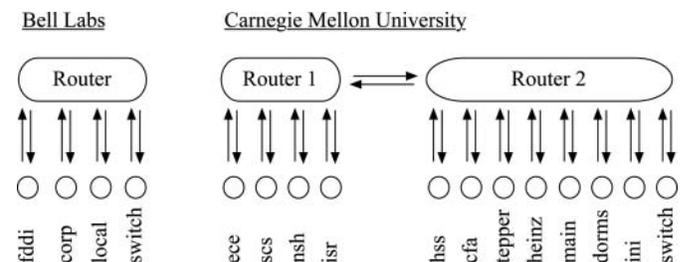


Figure 4. Topologies of the Bell Labs and Carnegie Mellon networks. The Abilene network has a more complex topology, an illustration of which is available in Fang, Vardi, and Zhang (2007). We also include these datasets in the `networkTomography` package.

routers are linked via a single connection. This configuration yields 26 observed aggregate traffic volumes and 144 OD traffic volumes, observed every 5 min at 473 points in time. This larger dataset allows us to compare network tomography techniques in a richer, more realistic setting. In combination with the *router1* data, it also allows us to explore the effect of dimensionality on performance and computational efficiency on real traffic data. Neither dataset contained any missing observations.

Our third dataset comes from the Abilene network, courtesy of Matthew Roughan. We use the *XI* dataset analyzed in Fang, Vardi, and Zhang (2007), which consists of aggregate and point-to-point traffic volumes measured every 5 min over a 3-day period. The underlying network has 12 nodes, yielding 144 point-to-point traffic time series. A total of 54 aggregate traffic volumes are observed at each time, consisting of 30 inner links and 24 edge links. Abilene is an Internet2 backbone network. Abilene’s traffic volumes, dynamics, and variability are quite different from those observed on local area networks such as Carnegie Mellon’s and Bell Labs’. Abilene’s topology is more complicated than simple star and dual-loop configurations. Thus, this dataset provides a different scenario for testing tomography methods.

We did not apply any seasonal adjustment or other more complex dynamic models to these datasets, given the short time period they span. We would recommend such an extension for time series spanning longer periods; indeed, even for data spanning only 2 days, usage patterns by time of day can be present. However, we endeavor to compare our deconvolution algorithms on equal footing—our focus is dynamic deconvolution. Thus, all methods are implemented with only local dynamics, without seasonal adjustment.

## 4.2 Competing Methods

We tested locally IID and smoothed Gaussian methods (Cao et al. 2000), a Bayesian MCMC approach (Tebaldi and West 1998), a simple gravity method (see, e.g., Fang, Vardi, and Zhang 2007), a tomogravity method (Zhang et al. 2003b), the Gaussian SSM developed for regularization in Section 3.2.1, the multilevel SSM with naïve regularization developed in Section 3.1, and the proposed two-stage estimation procedure with model-based regularization. All approaches were implemented in R with extensions in C to avoid computational bottlenecks (e.g., IPFP). For the methods of Cao et al. (2000) and the Gaussian SSM, which use windowed estimates, we selected a window width of 23 observations on the basis of prior work (Airoldi 2003). The final point-to-point traffic volume estimates were generally insensitive to a range of window sizes—this is largely attributable to the use of estimated  $x_t$ ’s instead of  $\lambda_t$ ’s for regularization. For the Abilene data, we considered the alternative model-based regularization procedure for the dynamic multilevel model, based on a simple gravity model as detailed at the end of Section 3.2.3.

For the approach of Tebaldi and West (1998), we tested both the original implementation and our own modification in which (following the authors’ original notation)  $\lambda_j$  and  $X_j$  are sampled with a joint Metropolis–Hastings step. The proposal distribution for this step is constructed by first proposing uniformly along the range of feasible values for  $X_j$  given all other values,

then drawing  $\lambda_j$  from its conditional posterior given the proposed  $X_j$ . This greatly improves the efficiency of the MCMC sampler, leading to improved convergence (we observed multivariate Gelman–Rubin diagnostics reduced by approximately an order of magnitude) and better predictions. These improvements allow us to compare the approach of Tebaldi and West (1998) on a more level playing field, focusing on the underlying model while mitigating computational issues.

For inference in the proposed dynamic model, we used 1000 particles and 10 MCMC iterations (in the move step of the SIRM filter) per time point in all experiments. We selected the former based on the number of effectively independent particles per time point, targeting a minimum of 10 in pilot runs. The number of MCMC iterations was chosen as a balance between computational burden and particle diversity. For the tomogravity method of Zhang et al. (2003b), we set  $\lambda$  to 0.01. Results were insensitive to the choice of  $\lambda$  across a wide range of values, as previously reported (Zhang et al. 2003b). These choices were kept consistent across experiments because they offered acceptable trade-offs and enabled a meaningful comparison of the competing methods.

## 4.3 Performance Comparison

We summarize performance of the methods described above on all three datasets in Tables 1–3. Each row corresponds to a method, and the columns provide mean  $L_1$  and  $L_2$  errors over time for the estimates of OD traffic in each dataset with corresponding standard errors. For the Bell Labs dataset, we provide errors in kilobytes; for the CMU and Abilene data, we provide errors in megabytes. We also provide Figure 5 below and Figures S1–S5 in the supplementary material as a visualization of our results on the Bell Labs dataset. We compare and discuss performance in terms of accuracy, computational stability, and scalability.

*Accuracy.* We obtain favorable performance for the two-stage approach (corresponding to the bottom rows of Tables 1–3) for all three datasets. For the Bell Labs data, mean (time-averaged)

Table 1. Performance comparison with Bell Labs data, all results in kB

Method	Bell Labs			
	$L_2$ error	SE	$L_1$ error	SE
Gravity	62.96	3.16	182.58	7.69
Tomogravity (Zhang et al. 2003b)	62.96	3.16	182.58	7.69
Locally IID model	104.59	5.54	160.24	6.53
Smoothed locally IID	104.25	5.52	157.87	6.48
Tebaldi and West (uniform prior)	76.60	4.91	173.94	7.49
Tebaldi and West (joint proposal)*	49.43	2.58	147.66	6.18
Gaussian state-space model	19.35	0.72	57.66	2.06
Dynamic multilevel model (naïve prior)	63.29	3.35	178.43	8.09
Dynamic multilevel model (SSM prior)	19.93	0.87	58.20	2.39

NOTE: \*Denotes our own improvement on the original algorithm by Tebaldi and West. Note that the performance of simple gravity and tomogravity is identical on this network due to its star topology.

Table 2. Performance comparison with CMU data, all results in MB

Method	CMU			
	$L_2$ error	SE	$L_1$ error	SE
Gravity	499.24	11.32	1521.66	30.09
Tomogravity (Zhang et al. 2003b)	310.61	5.95	1096.38	18.68
Locally IID model	592.49	9.91	1169.15	17.11
Smoothed locally IID	–	–	–	–
Tebaldi and West (uniform prior)	–	–	–	–
Tebaldi and West (joint proposal)*	167.94	4.42	712.37	14.68
Gaussian state-space model	110.47	6.19	389.14	16.72
Dynamic multilevel model (naïve prior)	311.21	6.25	1109.68	19.58
Dynamic multilevel model (SSM prior)	93.42	5.20	334.74	13.64

NOTE: \*Denotes our own improvement on the original algorithm by Tebaldi and West.

$L_1$  and  $L_2$  errors are statistically indistinguishable (within 1 SE) between the calibration procedure and dynamic multilevel model with model-based regularization. Both of these methods reduce average  $L_1$  and  $L_2$  errors by 60%–80% compared with the other approaches presented, representing a major gain in predictive accuracy. For the CMU data, we obtain a reduction of 53% in average  $L_2$  error and 44% in average  $L_1$  error from the algorithm of Cao et al. (2000) to our multilevel SSM; we observe 14%–15% reductions in average  $L_1$  and  $L_2$  errors from our Gaussian SSM to the multilevel state-space models. Furthermore, we observe large gains in filtering performance for both datasets compared with inference using naïve regularization with our multilevel SSM. Overall, our approach outperforms

Table 3. Performance comparison with Abilene data, all results in MB

Method	Abilene			
	$L_2$ error	SE	$L_1$ error	SE
Gravity	7.51	0.05	4.05	0.02
Tomogravity (Zhang et al. 2003b)	5.26	0.05	3.06	0.02
Locally IID Model	12.17	0.07	7.03	0.03
Tebaldi and West (joint proposal)*	12.74	0.07	7.44	0.04
Gaussian state-space model	15.48	0.09	8.42	0.05
Dynamic multilevel model (SSM prior)	14.89	0.08	8.09	0.05
Dynamic multilevel model (Gravity prior)	4.01	0.03	2.49	0.01

NOTE: \*Denotes our own improvement on the original algorithm by Tebaldi and West.

existing methods in accuracy, with greater gains from the Gaussian SSM to the multilevel SSM in our higher-dimensional setting. The multilevel SSM also outperforms gravity-based techniques substantially in these local area networks.

The comparative performance is somewhat different on the Abilene data. Gravity and tomogravity perform very well on this backbone network, while the locally iid and Poisson models (Tebaldi and West 1998; Cao et al. 2000) perform relatively poorly. The same is true for the Gaussian SSM, and the performance of the dynamic multilevel model suffers when using the SSM for regularization. Using a simple gravity model for regularization improves the performance of the dynamic multilevel model, leading to a reduction in mean  $L_1$  and  $L_2$  error of approximately 20%. The variability in traffic volumes requires

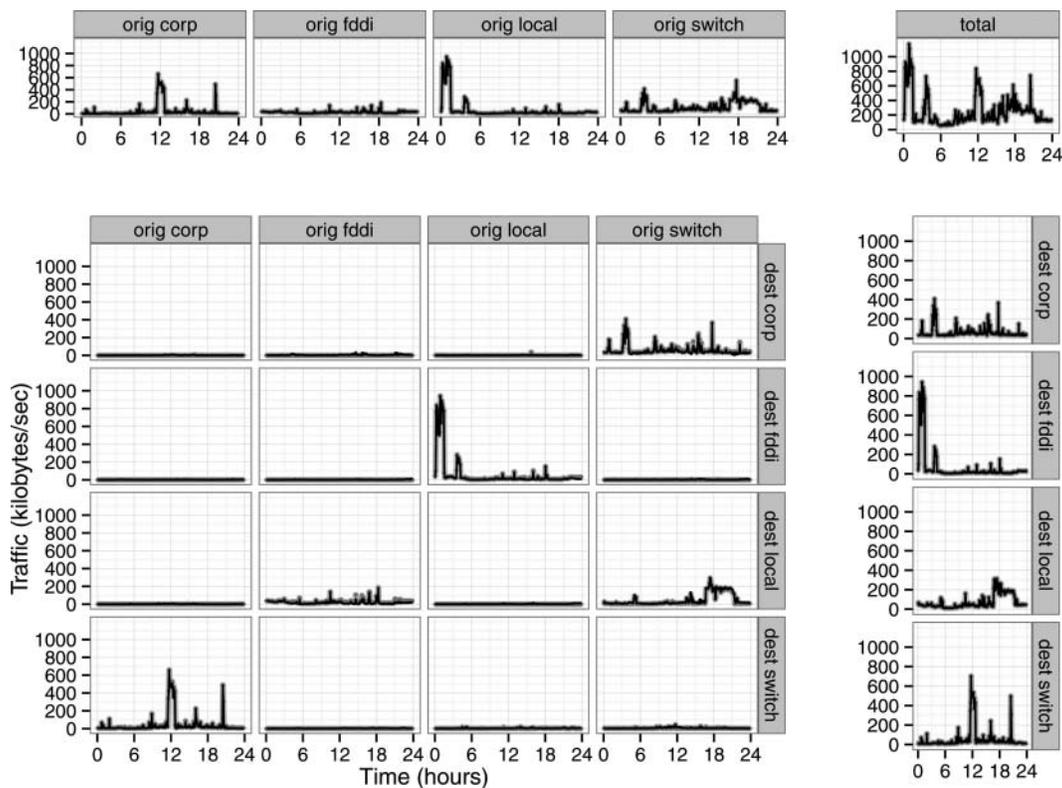


Figure 5. Actual and fitted (dynamic multilevel model with SSM prior) traffic for Bell Labs data. Actual traffic is in gray, fitted traffic is in black.

a smoothing based on instantaneous dynamics, rather than one based on a constant scaling fixed across point-to-point routes.

A combination of three factors can be used to understand the performance of our methods, listed on the bottom three rows in Tables 1–3: explicit dynamics, heavy tails, and regularization. The multilevel SSM fit with the naïve SIRM filter incorporates explicit dynamics and heavy tails, but its performance suffers because the distribution of  $\lambda$  is quite diffuse. The network tomography inverse problem requires more constraints and outside information to yield useful solutions. The Gaussian SSM used for calibration purposes is identifiable (as we show in Section 3.2.1) and incorporates explicit dynamics, but does not account for heavy tails. It performs well on the Bell Labs dataset, where the distributions of OD traffic are relatively symmetric (on the log scale), but suffers on the extremely heavy-tailed CMU traffic volumes. The multilevel SSM fit with the two-stage procedure overcomes the ill-posedness of the underlying inverse problem and accounts for heavy tails, attaining comparable performance on the Bell Labs data and outperforming considerably on the CMU data as a result. However, the Gaussian SSM is not realistic for every setting. Gravity methods offer a better source of regularization in backbone networks like Abilene. This last finding agrees with previous research (see, e.g. Zhang et al. 2003b, 2009).

*Computational stability.* We found a surprising amount of variation in computational stability among the methods evaluated. The local likelihood methods (Cao et al. 2000) remained stable across datasets, but the original method of Tebaldi and West (1998) encountered issues. On the Bell Labs data, it required a very large number of iterations to obtain convergence (as indicated by the Gelman–Rubin diagnostic); 150,000 iterations per time were used to provide the given estimates, 50,000 of which were discarded as burn-in. This method failed completely on several time points in the CMU data, becoming trapped in a corner of the feasible region. Our revised version of the original MCMC algorithm performed better, requiring far fewer iterations for convergence; 50,000 iterations were sufficient for all examined cases, although 150,000 iterations were used for the results presented for comparability.

Our calibration procedure proved computationally stable across all three datasets. The direct use of marginal likelihood, for maximum likelihood estimation, proved effective in both the low- and high-dimensional datasets. The multilevel SSM was also stable in both settings; however, it proved to be sensitive to some of the alternative specifications mentioned in Section 3. In particular, major problems arose in experiments using the posterior on  $x_t$  from the previous time as a proposal (as is common in applications of particle filtering); several time points in the Bell Labs data required more than 10 million proposals to obtain a single feasible particle. Additional care was needed with the “move” step due to similar issues. Furthermore, the use of a naïve, random walk regularization caused some computational difficulties, as the particles often became extremely diffuse in the feasible region. Overall, we found inference with the multilevel SSM computationally stable so long as sampling methods for highly constrained variables ( $x_t$  in particular) explicitly respected said constraints, proposing only valid values. Our RDA (detailed in Algorithm 2) handles this task well.

*Scalability.* All methods we evaluated fared well in scalability, including the computationally intensive, sequential SIRM inference we used for the multilevel SSM. On the Carnegie Mellon dataset, for each time point, the methods of Cao et al. (2000) required approximately 225 sec to obtain maximum likelihood estimates with a 23 observation window. Our modification of the sampler by Tebaldi and West (1998) required approximately 1500 sec to obtain 150,000 samples for a single time points—the original MCMC sampler required 2250 sec on average and often did not complete. In contrast, the simulation-based filtering method for the multilevel SSM required 270 sec per time point on average, on the Carnegie Mellon data, and 210 sec per time on average, on the Abilene data. Approximately 70% of this time was spent in the move step (MCMC) of the SIRM algorithm, with the vast majority of the remainder used for the random direction sampler. On the Bell Labs dataset, the filtering method required approximately 8 sec per time, whereas our modification of the sampler by Tebaldi and West (1998) required 150 sec per time—the original algorithm required approximately the same time.

These results are encouraging: the filtering algorithm is reasonably efficient (even written in R) and can run faster than real time with 144 point-to-point traffic volumes at 5-min sampling intervals. We note that the Abilene dataset required less computation per time than the Carnegie Mellon dataset, even though both involved 144 point-to-point traffic time series. This is because the effective dimensionality of the ill-posed linear inverse problem is substantially lower for Abilene; we observe 24 linearly independent aggregate traffic time series for Carnegie Mellon and 42 for Abilene, yielding 120 and 102 undetermined point-to-point traffic time series, respectively. The reduction in undetermined dimensions closely tracks the reduction in computation for the SIRM filter, as expected—the key computations of this sampler scale in complexity as the product of effective dimension and the number of point-to-point time series. The proposed method takes advantage of the geometric structure of each dataset to simplify sampling and guide inference.

Given more efficient implementation and parallelization, which are feasible for all sampling steps, the two-stage approach can scale to networks with many time more nodes. This is especially true given the sparsity of the traffic on many such point-to-point routes; the prevalence of zero (observable) aggregate traffic volumes in real world data further reduces the effective size of the deconvolution problem. By more efficient implementation, we refer to the actual code for the SIRM filter. In the `networkTomography` package, all parts of the SIRM filter itself are implemented in R. Moving this algorithm to a compiled language (C or Fortran) and eliminating many memory allocations promises an order of magnitude speedup, based on initial development.

## 5. SIMULATION STUDIES

Here, we further explore the relative performance of the methods we applied to the real data in Section 4.3 by designing two experiments that involve a mix of simulated and real data.

We sought to understand the source of the performance of the competing inference methods with two experiments. In the first experiment, we simulated data from the model and compared

the performance of the naïve random walk prior with the two-stage estimation strategy. The results show that the two-stage inference strategy leads to consistently better computational performance.

The second experiment involves a large-scale simulation study that compares the available methods by combining real OD traffic with simulated network topologies. We simulate a number of such scenarios by changing the degree of sparsity in the traffic and the complexity of the routing matrix, according to an experimental design that allows for an analysis of covariance (ANCOVA) analysis. The results show that significant error reduction can be expected by using the two-stage estimation strategy.

### 5.1 Evaluation of the Two-Stage Inference Strategy

In the first experiment, we sought to quantify whether the two-stage estimation strategy proposed in Section 3.2 leads to consistently lower  $L_1$  and  $L_2$  errors, on average, over time and OD routes.

We simulated OD traffic from our multilevel SSM under three network topologies: a three-node bidirectional chain, a three-node star topology, and a four-node star topology, corresponding to two-, four-, and nine-dimensional solution polytopes for inference on  $\mathbf{x}_t$ . For each of these cases, we produced 30 datasets consisting of 300 time points by drawing from the given multilevel model. We drew all initial OD traffic from log-Normal distributions with median 500 and geometric standard deviation 6. Subsequent evolution of these traffic volumes was simulated with  $\rho = 0.5$  and all other parameters as described in Section 3.1.1. We then computed the implied aggregate traffic volumes for each replicate and fit the multilevel model to these data using the two-stage estimation strategy. In addition to the two-stage approach outlined previously, we also performed filtering using our multilevel SSM with a naïve random walk regularization on the OD traffic, that is, we set  $\theta_{1,i,t} = 0 \quad \forall(i, t)$  and  $\theta_{2,i,t} = \log(5)/2$ . This allows us to directly evaluate the effects of regularization and the plausibility of our model.

The primary quantity of interest in our simulations are the relative mean  $L_2$  and  $L_1$  errors in estimated OD traffic for the naïve SIRM particle filter compared with our two-stage method. The distributions of these relative  $L_2$  errors is summarized in Figure 6. The magnitude of the errors is unchanged for the relative  $L_1$  errors.

We find that our two-stage method clearly outperforms the naïve SIRM particle filter when the dimensionality of the solution polytope is larger than two. Specifically, we have a mean relative error of  $1.09 \pm 0.49$  in two dimensions, increasing to  $1.57 \pm 0.45$  in four dimensions and  $1.40 \pm 0.26$  in nine dimensions.

Our experience with these simulations also highlighted the computational benefits of the proposed two-stage strategy. During iterations with the naïve SIRM filter, the *effective* number of particles rarely climbed above 2, whereas we typically obtained 10–50 with the two-stage approach, with an equivalent *actual* number of particles. With real data, we expect additional benefits from the two-stage estimation; in particular, we would expect it to have greater robustness to model misspecification. Essentially, we are using information from a simpler model to

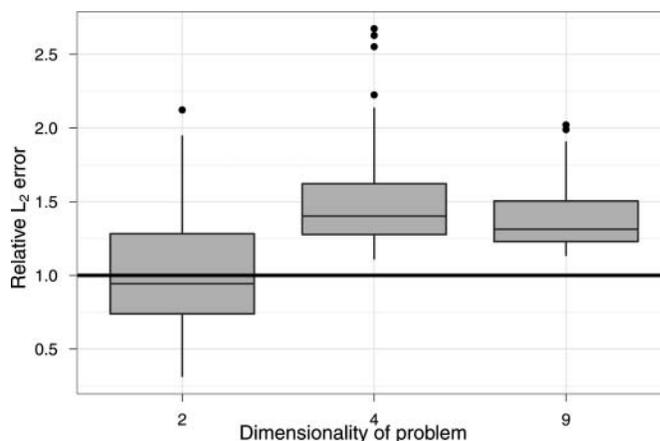


Figure 6. Relative  $L_2$  error for naïve versus two-stage method against dimensionality.

rein-in potential issues with the more delicate multilevel model. This strategy is expected to stabilize inference in the latter and limit problems of nonidentifiability. These intuitions are further explored in Section 5.2.

### 5.2 Quantifying the Factors That Affect Performance

In the second experiment, we sought to quantify the performance of the proposed method and existing methods relative to a the simple IPFP as baseline, explicitly controlling for the size for the network and the sparsity of the OD traffic volumes over time.

To get a better sense of the relative performances on real data, we designed experiments using real traffic time series. We selected a subset of the 1024 most active OD traffic volumes from the CMU dataset as the population of time series for designing these experiments. The remaining CMU OD pairs had negligible amounts of traffic. We used realistic but artificially designed routing matrices based on star topologies with three, four, five, and nine nodes, and we generated 10 datasets for each topology by randomly sampling from the population of time series.

This combination of active OD traffic volumes and star topologies defines a difficult scenario for inference underlying the network tomography problem. In cases with extremely sparse OD traffic time series, we can deterministically infer many of them from aggregate traffic, since a large portion of measurements will be zero. Hence, using high OD traffic volumes reduces the number of cases with deterministic solutions and puts the methods to a stringent test.

The star topology also creates a stringent test for our method. As noted by Singhal and Michailidis (2010), the star case is not worst with respect to, for instance, mean identifiability. However, it does provide the fewest measured aggregate traffic time series for a given number of OD routes. That is, given  $d$  nodes with  $n = d^2$  OD routes, and assuming all communications are bidirectional, any connected topology will have at least  $2d$  aggregate time series. Star topologies attain this lower bound, maximizing the dimension of the feasible region for OD traffic volumes given observed aggregate traffic. This dimension is the relevant measure of difficulty for inference in network tomography, so the star topology provides an appropriate benchmark.

Table 4. Log-linear ANCOVA model for simulation study;  $\log_{10}(L_1 \text{ errors})$  is outcome

	Estimate	Std. error	$t$ value	$\Pr(> t )$
(Intercept)	6.526	0.115	56.60	0.000
Dim = 4	0.058	0.111	0.53	0.599
Dim = 5	0.908	0.108	8.42	0.000
Dim = 9	2.103	0.108	19.47	0.000
Locally IID method	2.205	0.130	17.03	0.000
Tebaldi and West	-0.116	0.130	-0.89	0.373
Naïve prior	-0.021	0.130	-0.16	0.870
Calibration model	-0.241	0.130	-1.86	0.065
Two-stage inference	-0.244	0.130	-1.88	0.061
$\log_{10}$ sparsity	-8.237	3.179	-2.59	0.010

On each generated dataset, we ran five methods to estimate the OD traffic volumes: IPFP, the locally IID method of Cao et al. (2000), our implementation of the Poisson model from Tebaldi and West (1998), and the multilevel model with a naïve random walk prior, the proposed calibration procedure, and the proposed two-stage inference strategy.

The outcomes of interest are average  $L_1$  and  $L_2$  errors, over time and OD routes. We performed an ANCOVA analysis using the average errors from each of our experiments to understand what drives performance in this problem. The primary factor of interest for this analysis is the method used. We coded network size as a factor with three levels. We included sparsity as a covariate as well to capture the effect of having deterministically zero observed traffic as described above. Sparsity enters the ANCOVA analysis as  $\log_{10}$  (average proportion of measured traffic volumes that are not deterministically 0).

Table 4 summarizes the results of this analysis for  $L_1$  errors. We used a log-linear model for this analysis; initial diagnostics suggested that its variance structure is more appropriate for this experimental data than an untransformed model. We checked for interactions between dimensionality and method, but found no support in the data ( $p = 0.996$  for a standard  $F$  test). It appears that any such interactions would require larger networks to identify.

We find that the proposed methods significantly outperform the baseline IPFP approach, while the locally IID method performs significantly worse. The performance of the Bayesian model by Tebaldi and West (1998) and the multilevel methods fit with the naïve SIRM filter are inconclusively better than IPFP. The multilevel model consistently underperforms the model of Tebaldi and West (1998), as well as our other approaches, when used with naïve random walk priors. However, the calibration procedure alone performs quite well despite its simplicity. The performance of the proposed two-stage approach is not significantly higher than that for the calibration procedure in this setting, which suggests that our calibration procedure is the driver of our performance improvements at this scale. This agrees with our empirical findings with the Bell Labs dataset and is compatible with the observed increase in performance at larger scales, on the CMU dataset, in Table 2. These results are essentially unchanged (qualitatively and quantitatively) when we substitute  $L_2$  for  $L_1$  errors.

## 6. CONCLUDING REMARKS

In this article, we address the problem of (volume) network tomography in a dynamic filtering setting. For this application, we develop a novel approach to this problem by combining a new multilevel SSM that posits non-Gaussian marginals and non-linear probabilistic dynamics with a novel two-stage inference strategy. Our results and analyses substantiate several claims and suggest points for further discussion.

We analyzed three networks (Bell Labs, Carnegie Mellon, and Abilene) that span a wide range of dimensions, with different inference methods. The results demonstrate a clear improvement of the proposed methodology over previously published methods in estimating point-to-point traffic volumes. Comparison between Bell Labs and Carnegie Mellon results suggests that this gain increases with the dimensionality of the problem. Our results with the Abilene network highlight the differences between local area and Internet2 backbone networks. They differ in both topology and traffic dynamics, requiring different approaches to regularization.

### 6.1 Modeling Choices

Our model explicitly captures two critical features of point-to-point traffic—namely, skewness and temporal dependence. The substantial improvements in accuracy over existing methods can be attributed to these modeling improvements, to a large extent. The gains in computational efficiency are responsible for the improvements in accuracy only in part, as we discuss below. Previous modeling approaches have accounted for skewness (Tebaldi and West 1998), but never for explicit temporal dependence of the point-to-point traffic volumes. The intertemporal smoothing algorithm of Cao et al. (2000) includes elements of temporal dependence; however, the model assumes temporally independent time series and the dependence is imposed indirectly having observations within a time window that contribute to inference at any given time point. In summary, previous work has not accounted for the range of properties addressed by our model. The performance gains that stem from our modeling assumptions are clear on the three datasets tested; in particular, the gains from the model of Tebaldi and West (1998) to the Gaussian SSM and to the dynamic multilevel model for the CMU dataset reinforce the benefits of positing realistic dynamics in this problem.

We chose to increase the probability of near-zero traffic volumes using a truncated Gaussian error, rather than a log-Normal or Gamma distribution whose support is naturally on the non-negative reals. From a computational perspective, given that the particle filter involves a Metropolis step, the truncated Gaussian error is not particularly more tractable than log-Normal or Gamma errors. However, the truncated normal increases the probability assigned to any  $[0, \epsilon)$  interval, relatively to Gamma or log-Normal with same mean and scale. To obtain similar behavior with these nontruncated noise structures would either be impossible, for example, with a log-Normal distribution, or require ad hoc reparameterization linking the shape (intuitively, variability) and center (magnitude of point-to-point traffic volumes), for example, with a Gamma distribution. We believe that using a truncated distribution with the mode calibrated using

the underlying intensity process ( $\lambda_t$  in our notation) provides a cleaner solution. However, several alternatives are viable.

Fundamentally, we estimate the point-to-point traffic volumes by projecting aggregate traffic volumes onto the latent space point-to-point traffic inhabits, that is, we want to compute  $\mathbb{E}[x_t | y_t]$  under a given probabilistic structure. The relative variability of OD flows over time plays a large role in inference, as there is typically a strong relationship between the mean and variance of point-to-point traffic. Because of this, simple methods that do not model variability explicitly and realistically, including Moore–Penrose generalized inverse (Harville 2008) and independent component analysis (Hyvärinen, Karhunen, and Oja 2003), are of limited use in this context. Surprisingly, however, we found that the IPFP (Fienberg 1970) often produces reasonable solutions, likely because each such solution corresponds to a feasible set of OD flows. Liang, Taft, and Yu (2006) had recently capitalized on this finding. In contrast, our approach models this variability with a probabilistic structure, improving inference by using this additional information.

## 6.2 Applicability to More Complicated Routing Schemes

Routing schemes other than deterministic routing, such as probabilistic routing and dynamic load balancing, can be subsumed within the modeling framework we developed in Section 2.

A probabilistic routing scheme would be captured by a probabilistic  $A$  matrix. Column  $j$  of the  $A$  matrix specifies the  $m$  proportions of point-to-point traffic volume  $j$  that is measured at the  $m$  counters. In deterministic routing, each point-to-point traffic volume either contributes to a counter or does not,  $A_{ij} \in \{0, 1\}$ . In probabilistic routing, each point-to-point traffic volume contributes to multiple counters with different probabilities,  $A_{ij} \in [0, 1]$ .

Routing schemes that carry out dynamic load balancing to manage congestion would be captured by a time- and traffic-dependent matrix,  $A(t)$ . However, congestion can only be monitored at the router level, in practice, using the *observed* aggregate traffic counters. Thus, the counters can be considered as covariates, and the routing matrix can be modeled as a function of these covariates,

$$A(t) = \begin{cases} A^{\text{High}} & \text{if } y_t \text{ is high} \\ A^{\text{Low}} & \text{if } y_t \text{ is low.} \end{cases}$$

More nuanced specifications are possible. The key point is that implemented routing and switching protocols can only make use of the measurements collected at the router level,  $y_t$ .

## 6.3 Computational Challenges and Scalability

Our inference method is computationally efficient and scales to larger inference problems than have been previously addressed. The problem is fundamentally  $O(n - m)$  for each time point, so we cannot hope to do better than quadratic scaling in the number of nodes  $d$  in our network (excepting cases where a few aggregate traffic measurements are zero). Despite the sophistication of our dynamic multilevel model, the sequential Monte Carlo technique allows for inference in better than

real time for a network with  $144 = n \approx O(d^2)$  OD routes and  $26 = m \approx O(d)$  router level measurements. As SIRM is the portion of the inference algorithm that would be used in an on-line application, we have demonstrated a scalable technique for inference with a model of greater complexity and realism than has been previously found in the literature.

These gains in computational efficiency also reduce numerical instability and are ultimately responsible for additional gains in accuracy. Computational issues can be appreciated by considering the amount of effort needed to maintain  $\text{cov } e_t$  positive definite in the EM algorithm of Cao et al. (2000), especially when the traffic approaches zero. We can see some artifacts in the corresponding point-to-point traffic estimates in supplemental Figure S1 (green lines) due to this issue in low traffic OD routes, for example, see “orig local  $\rightarrow$  dest local.” We further quantified the effects of computational efficiency on inference in the original methods by Tebaldi and West (1998) in Table 1 by comparing the uniform prior and component-wise proposal to the joint proposal we developed. In addition to the gains in speed and convergence discussed in Section 4.3, we observe a large reduction in average error from the component-wise to joint proposal (35% in  $L_2$  error, 15% in  $L_1$ ), which correspond to no changes in priors nor to the underlying model.

The RDA plays an important role in the sequential Monte Carlo sampler. Without such an algorithm to sample directly from the feasible region of point-to-point traffic volumes, we would be forced to use a naïve proposal distribution. In our testing, such distributions proved extremely problematic (as discussed in Section 3), especially in higher-dimensional settings. In such cases, intelligent sampling techniques that fully use the geometric constraints implied by the data are necessary to obtain high accuracy and efficiency. This is particularly salient comparing the results presented here to our previous work (Airoldi and Faloutsos 2004); the method presented there suffered from computational instabilities. It was hampered both by inefficient sampling on the feasible space of solutions and by distributional assumptions that assigned low probability to low point-to-point traffic volumes.

Multimodality of the marginal posterior on point-to-point traffic volumes  $x_{it}$  appears negligible in our analyses. Our theoretical results suggest that multimodality in these problems is limited to that due to flat regions in the case of real-valued traffic volumes and models assuming independent traffic. We suspect that the issues with multimodality discussed in the literature are due mainly to the inefficiency of the samplers. This further reinforces the importance of efficient computation for inference in highly complex, poorly identified settings; even a simple model can falter on poor computation, and complex models require great care to obtain reliable inference results.

## 6.4 A Novel Two-Stage Inference Strategy in Dynamic Multilevel Models

As previously argued by Tebaldi and West (1998) in the static setting, informative priors are essential to identify the peak in the likelihood, which correspond to the true configuration of point-to-point traffic. This conclusion holds in the dynamic setting we consider, despite the additional information that temporal dependence makes relevant for the inference of point-to-point

traffic volumes. The technical choices at issue are (i) where to find such information—it is not obvious in the data; (ii) what parameters are most convenient to put priors on; and (iii) how do we translate the additional information into prior information for the chosen parameterization.

We use a simple identifiable model to find rough estimates of the point-to-point traffic volumes (in our first stage). These estimates provide some information about where the point-to-point traffic volumes live in the space of feasible solutions, enabling us to identify high-probability subsets of the feasible region at each time point before embarking on computationally intensive sampling. The expected benefits from this strategy are larger in higher dimensions, as the proportion of the feasible region's volume with high posterior density decreases rapidly with dimensionality (the classical curse of dimensionality). Practically, informative priors increase the efficiency of the particle filter by focusing the sampler on promising regions of the parameter space, avoiding wasted computation and improving inference.

To pass the first-stage information to the (non-Gaussian) dynamic multilevel model, we moved away from a standard linear state-space formulation with additive error to a nonlinear formulation with stochastic dynamics, which effectively provides a multiplicative error (second stage). The stochastic dynamics assumed for  $\lambda_t$  provide our parameters of choice for incorporating information obtained in the first stage of estimation. Prior calibration for the dynamics of  $\lambda_t$  guides inference without placing too tight of a constraint on the inferred point-to-point traffic volumes. In Section 3.2.3, we describe how we solve the problem of translating the first-stage estimates into priors for the parameters of the second-stage model. Essentially, we trade-off the need to pass as much information as possible from the first stage of estimation to the second with the controlled inaccuracy of the first-stage modeling assumptions.

Our two-stage procedure suggests a more general strategy for inference in dynamic hierarchical models with weak identifiability constraints. The use of simpler models to guide inference in more sophisticated, realistic models, and (if necessary) to provide regularization, can result in large performance gains, as we have demonstrated here. This strategy implements a principled approach to *cutting corners* (Meng 2010). Ultimately, the proposed two-stage approach has allowed us to use a sophisticated generative model for inference, leveraging the power of multilevel analysis, while maintaining efficiency for real time applications.

## APPENDIX: EFFICIENT INFERENCE FOR THE GAUSSIAN STATE-SPACE MODEL

Inference on the latent point-to-point time series  $x_t$  in the Gaussian SSM specified by Equation (3) can be carried out with standard Kalman filtering and smoothing. Estimating the constants underlying the model via maximum likelihood can be approached with two strategies: Expectation-Maximization (EM; Dempster, Laird, and Rubin 1977) and direct numerical optimization. The EM approach for unconstrained Gaussian SSMs requires Kalman smoothing for the E-step and maximization of the expected log-likelihood for the M-step (Ghahramani and Hinton 1996). While the E-step is straightforward and efficient to calculate using standard algorithms, the M-step requires expensive numerical optimization in our case. Due to the constraints on  $Q$  and cov  $e_t$  and the dependence of the observations, there is no

analytic form for the maximum of the expected log-likelihood. Since the EM requires numerical optimization, we decided to use direct numerical optimization on the marginal likelihood obtained from the Kalman smoother. This amounts to maximizing

$$\ell(Y | \theta) = - \sum_t \log |\hat{\Sigma}_t| - \frac{1}{2} \sum_t (y_t - \hat{y}_t)' \hat{\Sigma}_t^{-1} (y_t - \hat{y}_t),$$

where  $\hat{y}_t$  and  $\hat{\Sigma}_t$  are the estimated mean and covariance matrices from the Kalman smoother. With a fast (Fortran) implementation of the Kalman iterations, this approach yields favorable run times and stable results.

Such efficient computation is, however, sensitive to certain modeling decisions. Enforcing a steady-state starting value within each window is particularly useful. Formally, suppose that we index each window of width  $w$  with  $t = 1, \dots, w$ . We must specify a starting value  $x_0$  for each window. By linking  $x_0$  to  $\lambda$ , we can simplify our computation. For a given choice of  $\lambda$ , the steady-state mean of the process specified in Equation (3) is  $\frac{1}{1-\rho}\lambda$ . Fixing  $x_0 = \frac{1}{1-\rho}\lambda$  allows us to reduce the dimensionality of Equation (3). Formally, we can rewrite it as

$$\begin{aligned} x_t - \frac{1}{1-\rho}\lambda &= F \left( x_{t-1} - \frac{1}{1-\rho}\lambda \right) + e_t \\ y_t &= A \left( x_t + \frac{1}{1-\rho}\lambda \right) + \epsilon_t. \end{aligned} \quad (\text{A.1})$$

This reduces the dimensionality of our state variable by a factor of 2, greatly accelerating all Kalman filter and smoother calculations. As said calculations have complexity quadratic in the problem's dimensionality, this reduces the computational load by approximately a factor of 4.

## SUPPLEMENTARY MATERIALS

The supplement contains plots of the actual and fitted OD flows for the methods presented previously. We plot all OD flows for the Bell Labs data and the 12 most variable OD flows for CMU. Ground truth is always in black, with estimated values in color. Figures S1 through S5 cover the Bell Labs data, and Figures S6 through S10 cover the CMU data.

[Received September 2011. Revised November 2012.]

## REFERENCES

- Airoldi, E. M. (2003), "Advances in Network Tomography," Technical Report CMU-CALD-03-101, Carnegie Mellon University. [154,157]
- Airoldi, E. M., and Faloutsos, C. (2004), "Recovering Latent Time-Series From Their Observed Sums: Network Tomography With Particle Filters," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 10, pp. 30–39. [149,151,153,156,162]
- Airoldi, E. M., and Haas, B. (2011), "Polytope Samplers for Inference in Ill-Posed Inverse Problems," in *International Conference on Artificial Intelligence and Statistics*, 15, 110–118. [151,153]
- Bell, M. G. H. (1991), "The Estimation of Origin-Destination Matrices by Constrained Generalized Least Squares," *Transportation Research, Series B*, 25B, 13–22. [151]
- Bishop, Y., Fienberg, S. E., and Holland, P. (1975), *Discrete Multivariate Analysis: Theory and Practice*, Cambridge, MA: The MIT Press. [149]
- Blocker, A. W., and Airoldi, E. M. (2011), "Deconvolution of Mixing Time Series on a Graph," in *Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence (UAI)*, 51–60. [149]
- Cao, J., Cleveland, W. S., Lin, D., and Sun, D. X. (2002), "The Effect of Statistical Multiplexing on the Long Range Dependence of Internet Packet Traffic," Technical Report, Bell Labs. [154]
- Cao, J., Davis, D., Van Der Vriel, S., and Yu, B. (2000), "Time-Varying Network Tomography: Router Link Data," *Journal of the American Statistical Association*, 95, 1063–1075. [149,151,153,154,155,156,157,158,159,161,162]

- Cao, J., Davis, D., Van Der Viel, S., Yu, B., and Zu, Z. (2001). "A Scalable Method for Estimating Network Traffic Matrices From Link Counts," Technical Report, Bell Labs. [151]
- Casella, G., and Berger, R. L. (2001). *Statistical Inference*, Pacific Grove, CA: Duxbury Press. [149]
- Castro, R., Coates, M., Liang, G., Nowak, R., and Yu, B. (2004). "Network Tomography: Recent Developments," *Statistical Science*, 19, 499–517. [149,151]
- Chen, Y., Diaconis, P., Holmes, S., and Liu, J. S. (2005). "Sequential Monte Carlo Methods for Statistical Analysis of Tables," *Journal of the American Statistical Association*, 100, 109–120. [151]
- Clogg, C. C., Rubin, D. B., Schenker, N., Schultz, B., and Weidman, L. (1991). "Multiple Imputation of Industry and Occupation Codes in Census Public-Use Samples Using Bayesian Logistic Regression," *Journal of the American Statistical Association*, 86, 68–78. [150,155]
- Coates, A., Hero, A. O., III, Nowak, R., and Yu, B. (2002). "Internet Tomography," *Signal Processing Magazine, IEEE*, 19, 47–65. [149]
- Deming, W. E., and Stephan, F. F. (1940). "On a Least Squares Adjustment of a Sampled Frequency Table When the Expected Marginal Totals are Known," *Annals of Mathematical Statistics*, 11, 427–444. [151]
- Dempster, A., Laird, N., and Rubin, D. (1977). "Maximum Likelihood From Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society, Series B*, 39, 1–38. [163]
- Deng, K., Li, Y., Zhu, W., Geng, Z., and Liu, J. S. (2012). "On Delay Tomography: Fast Algorithms and Spatially Dependent Models," *IEEE Transactions on Signal Processing*, 60, 5685–5697. [151]
- Diaconis, P., and Sturmfels, B. (1998). "Algebraic Algorithms for Sampling From Conditional Distributions," *The Annals of Statistics*, 26, 363–397. [151]
- Dobra, A. (2012). "Dynamic Markov Bases," *Journal of Computational and Graphical Statistics*, 21, 496–517. [151]
- Dobra, A., Tebaldi, C., and West, M. (2006). "Data Augmentation in Multi-Way Contingency Tables With Fixed Marginal Totals," *Journal of Statistical Planning and Inference*, 136, 355–372. [149]
- Erramilli, V., Crovella, M., and Taft, N. (2006). "An Independent-Connection Model for Traffic Matrices," in *ACM SIGCOMM Internet Measurement Conference (IMC06)*, New York: ACM, pp. 251–256. [151]
- Fang, J., Vardi, Y., and Zhang, C.-H. (2007). "An Iterative Tomography Algorithm for the Estimation of Network Traffic," in *Complex Datasets and Inverse Problems: Tomography, Networks and Beyond* (Vol. 54 of Lecture Notes–Monograph Series), eds. R. Liu, W. Strawderman, and C.-H. Zhang, IMS, pp. 12–23. [149,151,155,156,157]
- Fienberg, S. E. (1970). "An Iterative Procedure for Estimation in Contingency Tables," *The Annals of Mathematical Statistics*, 41, 907–917. [151,162]
- Ghahramani, Z., and Hinton, G. E. (1996). "Parameter Estimation for Linear Dynamical Systems," Technical Report CRG-TR-96-2, Department of Computer Science, University of Toronto. [163]
- Gilks, W. R., and Berzuini, C. (2001). "Following a Moving Target–Monte Carlo Inference for Dynamic Bayesian Models," *Journal of the Royal Statistical Society, Series B*, 63, 127–146. [150,154]
- Hansen, P. C. (1998). *Rank-Deficient and Discrete Ill-Posed Problems: Numerical Aspects of Linear Inversion*, Philadelphia, PA: SIAM. [149]
- Harrison, M. T. (2009). "A Dynamic Programming Approach for Approximate Uniform Generation of Binary Matrices With Specified Margins," arXiv:0906.1004. [151]
- Harville, D. A. (2008). *Matrix Algebra From a Statistician's Perspective*, New York: Springer. [162]
- Hyvärinen, A., Karhunen, J., and Oja, E. (2003). *Independent Component Analysis*, New York: Wiley. [162]
- Lakhina, A., Papagiannaki, K., Crovella, M., Diot, C., Kolaczyk, E. D., and Taft, N. (2004). "Structural Analysis of Network Traffic Flows," *SIGMETRICS Performance Evaluation Review*, 32, 61–72. [149]
- Lawrence, E., Michailidis, G., and Nair, V. (2006a). "Network Delay Tomography Using Flexicast Experiments," *Journal of the Royal Statistical Society, Series B*, 68, 785–813. [151]
- Lawrence, E., Michailidis, G., Nair, V., and Xi, B. (2006b). "Network Tomography: A Review and Recent Developments," in *Frontiers in Statistics*, eds. J. Fan and H. L. Koul, London, UK: Imperial College Press, pp. 365–368. [149,151]
- Lee, T.-W., Lewicki, M. S., Girolami, M., and Sejnowski, T. J. (1999). "Blind Source Separation of More Sources Than Mixtures Using Overcomplete Representations," *IEEE Signal Processing Letters*, 6, 87–90. [149]
- Liang, G., Taft, N., and Yu, B. (2006). "A Fast Lightweight Approach to Origin-Destination IP Traffic Estimation Using Partial Measurements," *IEEE/ACM Transactions on Networking*, 14, 2634–2648. [162]
- Liang, G., and Yu, B. (2003a). "Pseudo-Likelihood Estimations in Network Tomography," in *Proceedings of IEEE INFOCOM*, pp. 2101–2111. [149]
- (2003b). "Maximum Pseudo Likelihood Estimation in Network Tomography," *IEEE Transactions on Signal Processing*, 51, 2043–2053. [151]
- Liu, J. S., and Chen, R. (1995). "Blind Deconvolution via Sequential Imputations," *Journal of the American Statistical Association*, 90, 567–576. [149]
- Medina, A., Taft, N., Salamatian, K., Bhattacharyya, S., and Diot, C. (2002). "Traffic Matrix Estimation: Existing Techniques and New Directions," *SIGCOMM Computer Communication Review*, 32, 161–174. [149]
- Meister, A. (2009). *Deconvolution Problems in Nonparametric Statistics: Lecture Notes in Statistics*, New York: Springer. [149]
- Meng, X. L. (2010). "Machine Learning With Human Intelligence: Principled Corner Cutting ( $pc^2$ )," in *Plenary Invited Talk, Annual Conference on Neural Information Processing Systems (NIPS)*. [163]
- Miller, J. W., and Harrison, M. T. (2011). "Exact Enumeration and Sampling of Matrices With Specified Margins," arXiv preprint arXiv:1104.0323. [151]
- Parra, L., and Sajda, P. (2003). "Blind Source Separation via Generalized Eigenvalue Decomposition," *Journal of Machine Learning Research*, 4, 1261–1269. [149]
- Presti, F. L., Duffield, N. G., Horwitz, J., and Towsley, D. (2002). "Multicast-based Inference of Network-Internal Delay Distribution," *IEEE/ACM Transactions on Networking*, 6, 761–775. [151]
- Shepp, L. A., and Kruskal, J. B. (1978). "Computerized Tomography: The New Medical x-Ray Technology," *The American Mathematical Monthly*, 85, 420–439. [149]
- Shepp, L. A., and Vardi, Y. (1982). "Maximum Likelihood Reconstruction for Emission Tomography," *IEEE Transactions on Medical Imaging*, 1, 113–122. [151]
- Singhal, H., and Michailidis, G. (2007). "Identifiability of Flow Distributions From Link Measurements With Applications to Computer Networks," *Inverse Problems*, 23, 1821–1849. [155]
- (2010). "Optimal Experiment Design in a Filtering Context With Application to Sampled Network Data," *The Annals of Applied Statistics*, 4, 78–93. [160]
- Smith, R. L. (1984). "Efficient Monte Carlo Procedures for Generating Points Uniformly Distributed Over Bounded Regions," *Operations Research*, 32, 1296–1308. [150,154]
- Soule, A., Lakhina, A., Taft, N., Papagiannaki, K., Salamatian, K., Nucci, A., Crovella, M., and Diot, C. (2005). "Traffic Matrices: Balancing Measurements, Inference and Modeling," in *ACM Sigmetrics*, pp. 362–373. [151]
- Tebaldi, C., and West, M. (1998). "Bayesian Inference on Network Traffic Using Link Count Data," *Journal of the American Statistical Association*, 93, 557–573. [149,151,153,157,158,159,161,162]
- Vanderbei, R. J., and Iannone, J. (1994). "An EM Approach to OD Matrix Estimation," Technical Report SOR 94-04, Princeton University. [149,151]
- Vardi, Y. (1996). "Network Tomography: Estimating Source-Destination Traffic Intensities From Link Data," *Journal of the American Statistical Association*, 91, 365–377. [149,151,153]
- Vardi, Y., Shepp, L. A., and Kaufman, L. (1985). "A Statistical Model for Positron Emission Tomography," *Journal of the American Statistical Association*, 80, 8–20. [149]
- Zhang, Y., Roughan, M., Duffield, N., and Greenberg, A. (2003a). "Fast Accurate Computation of Large-Scale IP Traffic Matrices From Link Loads," in *Proceedings of SIGMETRICS*, pp. 206–217. [155]
- Zhang, Y., Roughan, M., Lund, C., and Donoho, D. (2003b). "An Information-Theoretic Approach to Traffic Matrix Estimation," in *Proceedings of SIGCOMM*, 301–312. [149,151,155,157,159]
- Zhang, Y., Roughan, M., Willinger, W., and Qui, L. (2009). "Spatio-Temporal Compressive Sensing and Internet Traffic Matrices," in *Proceedings of ACM SIGCOMM, Barcelona, Spain*, pp. 267–278. [151,159]