

Statistical Modeling of Large-Scale Sensor Count Data

Padhraic Smyth
Department of Computer Science
University of California, Irvine

Acknowledgements

- Collaborators:

Jon Hutchins



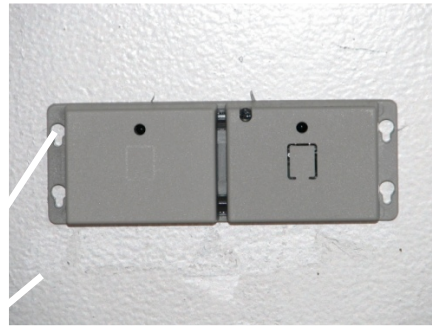
Alex Ihler



- Funding:

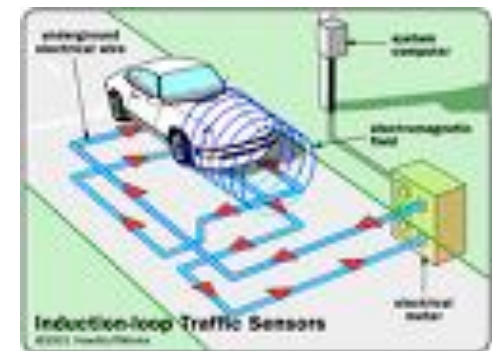
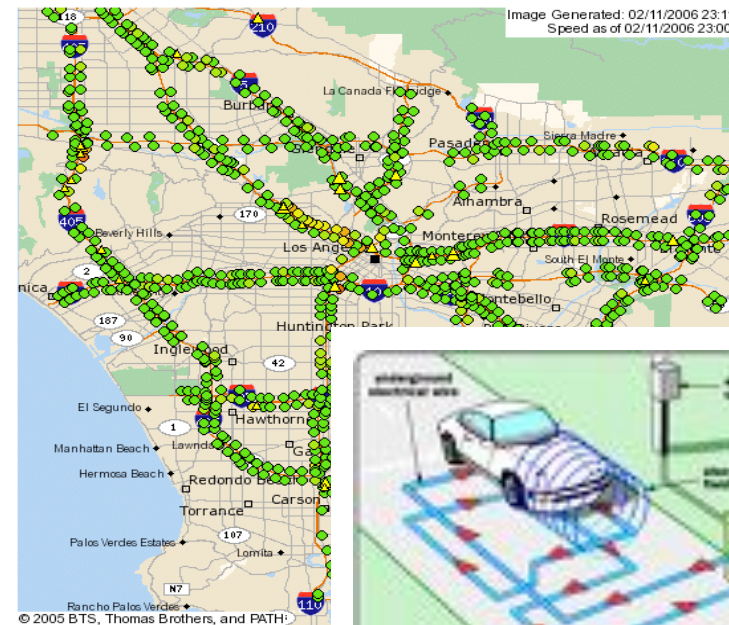
- US National Science Foundation
- California Department of Transportation

Sensors Measuring Human Activity

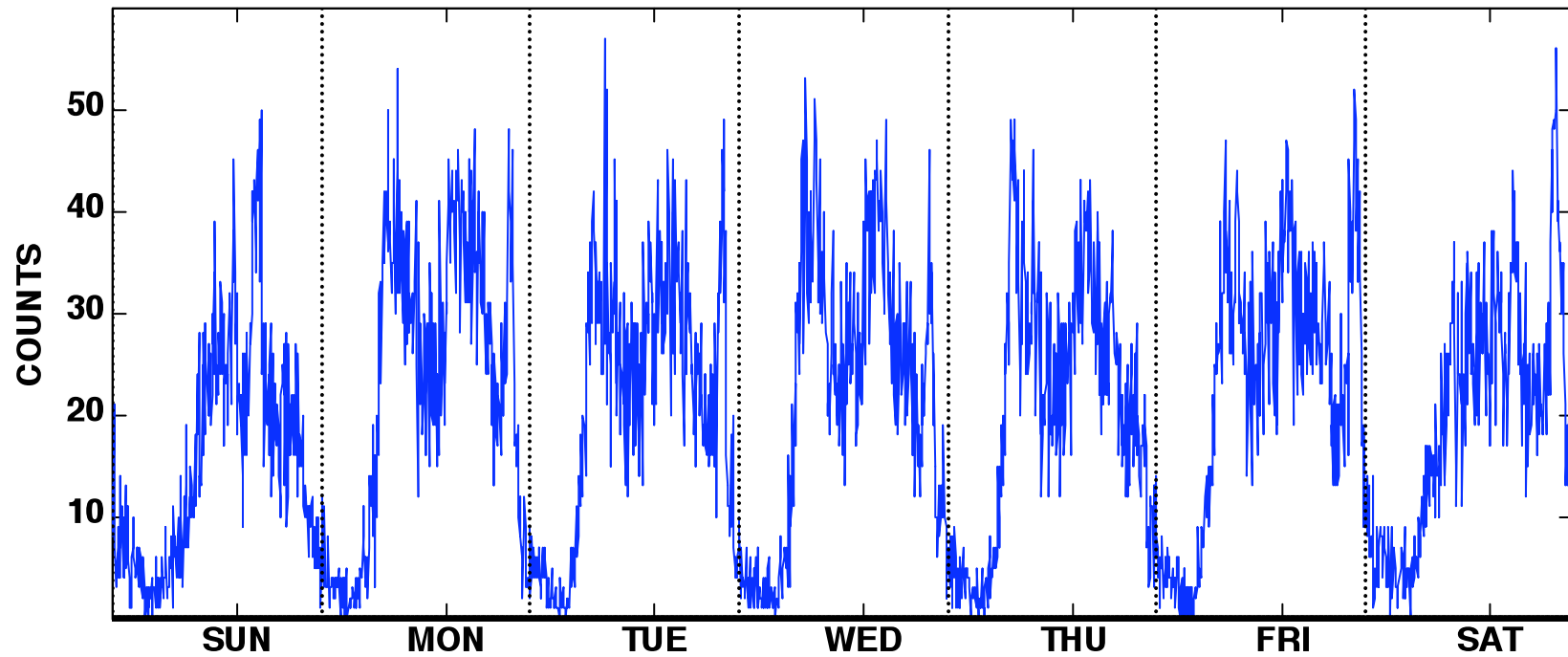


Optical people counter at a building entrance on campus

Loop sensors on Southern California freeways

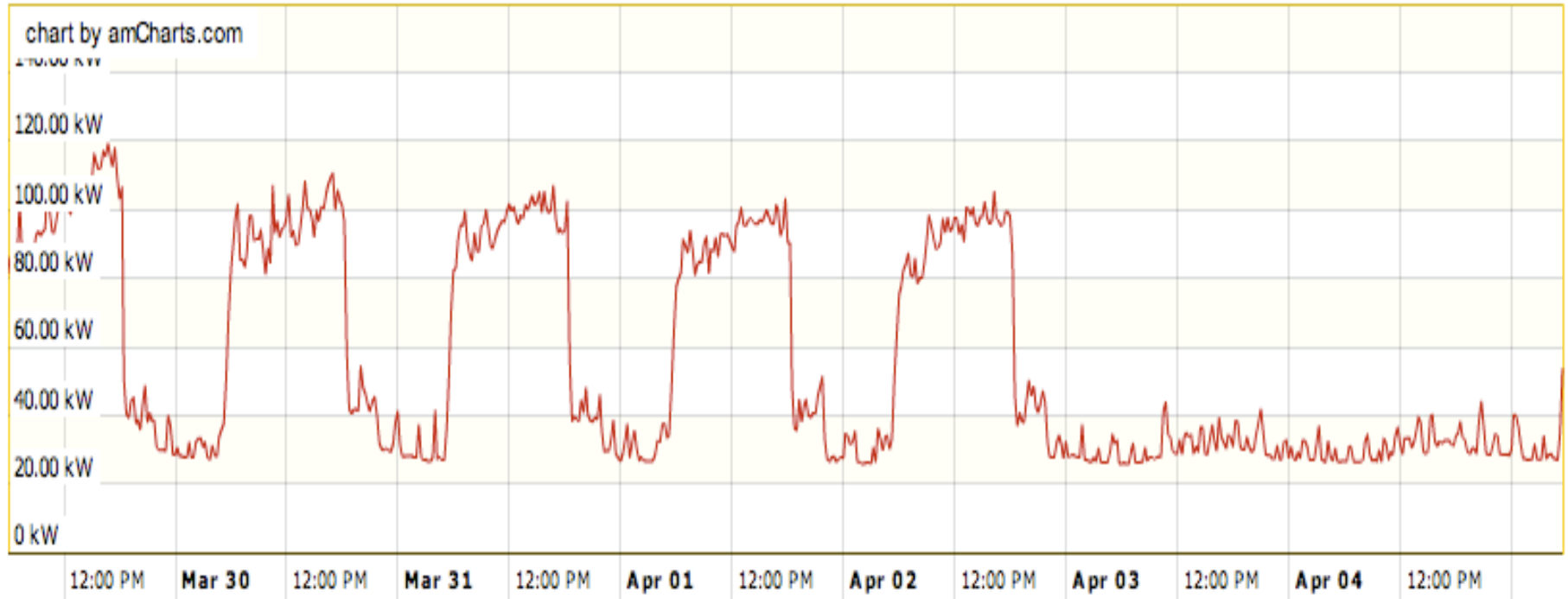


One Week of Traffic Data



● 1st Average kW in highlighted period: 56.39 kW

Mar 29, 2010 6:00:00 AM - Apr 05, 2010 5:30:00 AM



UCSanDiego | Energy Dashboard

[Login](#)
[Home](#) | [Individual Devices](#) | [Campus Meters](#) | [Research](#) | [About](#)

CSE Building / EBU3B > Campus Meter

[Meter Graph](#) | [Time Comparison](#) | [Add to compare list](#)

Device Information

Name: CSE Building Mechanical Load

Description: Combined power consumption for mechanical and elevator loads in the CSE building. Includes things like HVAC.

Overall Energy Statistics

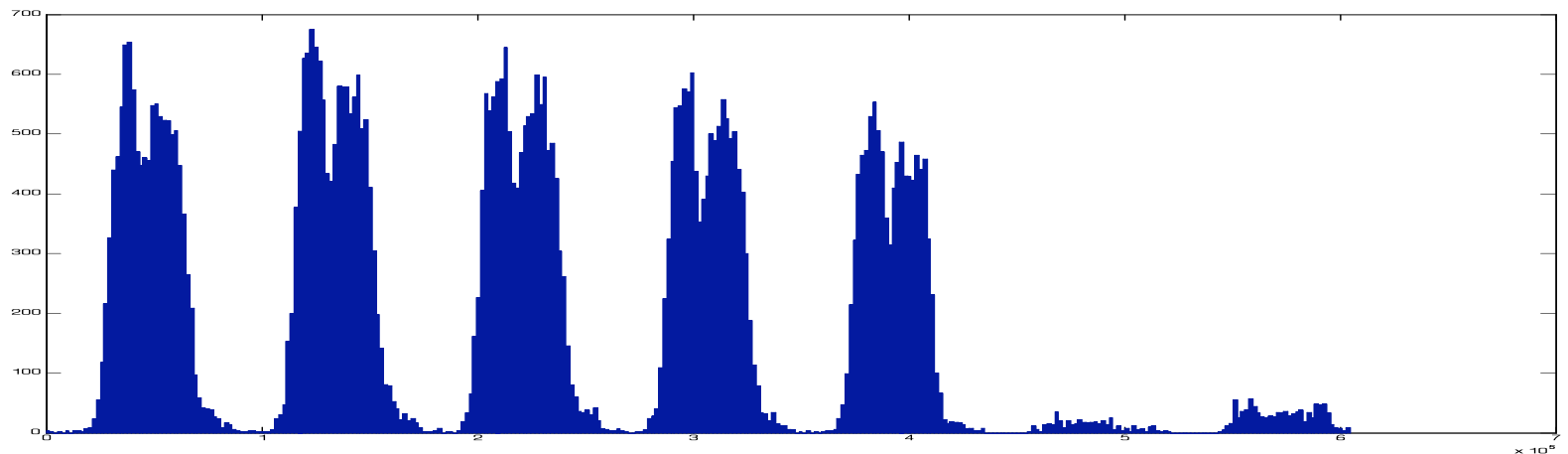
kW-Hours: 9473.4 kW-H

Avg kW: 56.39 kW

Energy costs: \$1231.542

Email Activity over 1 Week

number
of emails
sent



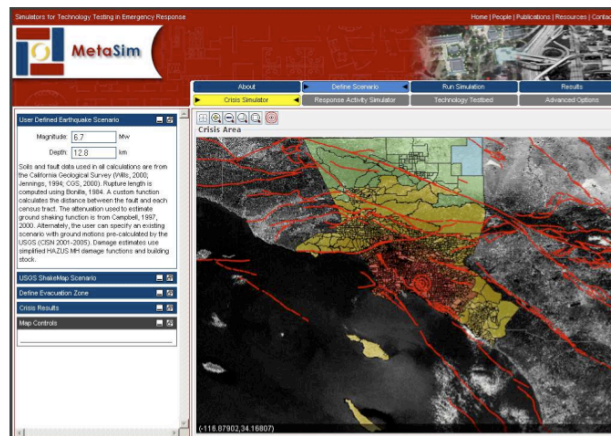
Time

Research Problems

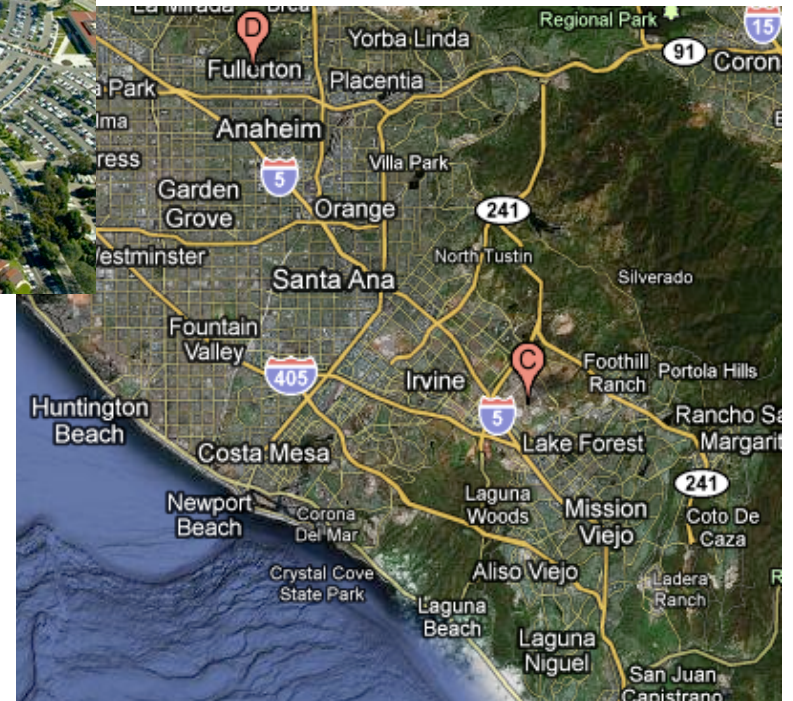
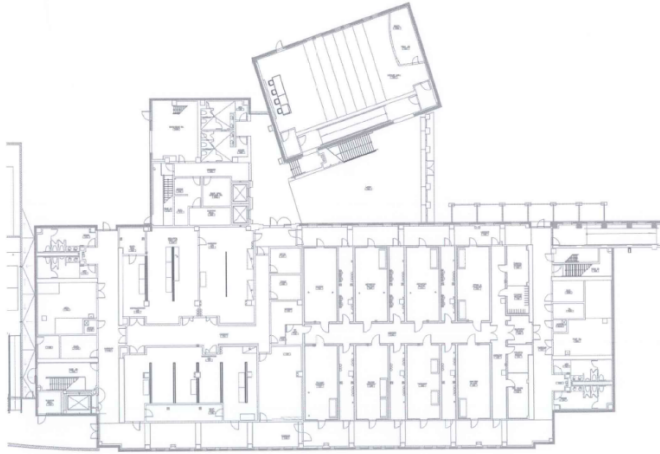
- Detection of anomalous events
- Inference of “typical behavior”
- Occupancy and flow estimation
- Prediction
- Missing data inference
-and more

Applications

- Traffic engineering
- Urban planning
- Energy monitoring
- Disaster response



Spatial Scale



Urban Sensing

- Combining multiple data sources
 - Satellite and aerial images (e.g., MODIS)
 - Energy usage (e.g., electricity)
 - Road sensors
 - Network traffic
 - Search engine data
 - Cell-phone data
 - Census data
 -

Outline

- Graphical models for count time-series
- Experimental results
 - Detecting events in traffic data and building sensor data
 - Estimating building occupancy
- Ongoing work and conclusions

General Context

- Aggregated counts at discrete times (people, vehicles, emails, etc)
- Data has underlying “signal”
 - e.g., daily/weekly rhythms of human behavior
 - occasional “bursts” of counts related to events
- Large amounts of data (e.g., years)
- ...but little prior knowledge
 - detailed time patterns of human behavior are unknown
 - algorithm has no knowledge of scheduled events

Related Work

- Markov-modulated Poisson processes
 - Heffes and Lucantoni (1994)
 - Scott (1998)
 - Scott and Smyth (2003)
- Segmentation with Markov/Poisson models
 - Kleinberg (2002)
 - Salmenkivi and Mannila (2005)
- Statistical modeling of traffic
 - Bickel et al, *Statistical Science*, 2007 + related work

Points to make re: Vehicle Data

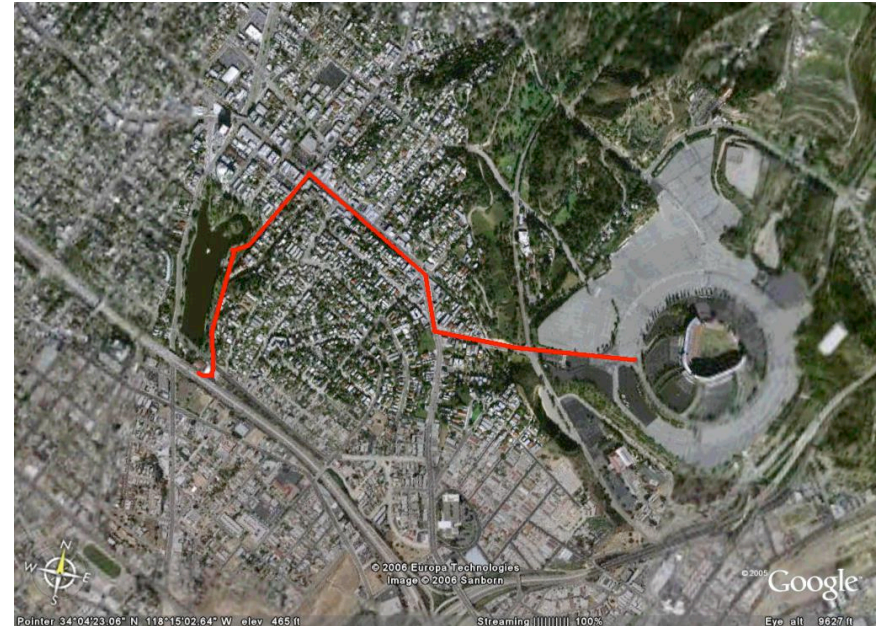
- Volume of data (# sensors x # measurements)
- Basic challenges and difficulties (missing, bad data)
- Why is this interesting to machine learning?

More points to make

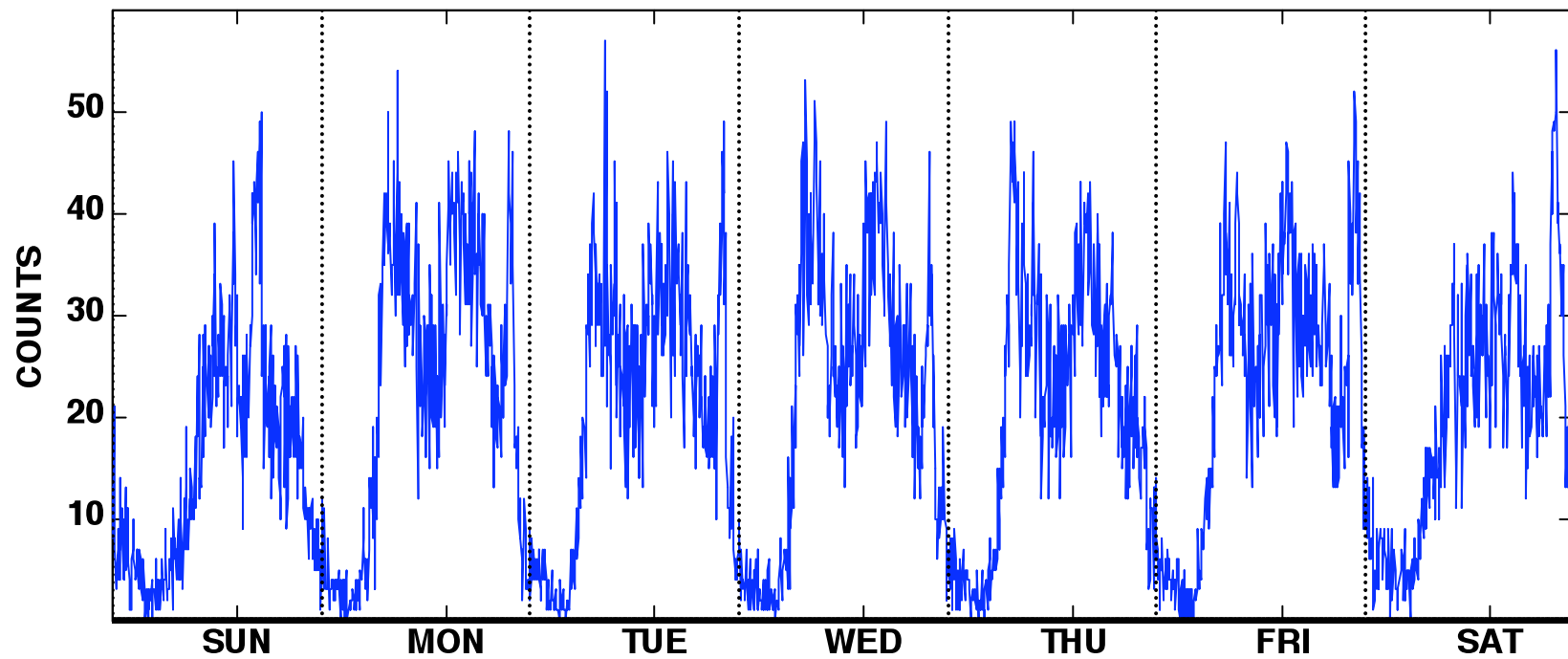
- Why traffic data is important
 - Air quality and the environment
 - Transportation and planning
 - Real-time prediction
 - Census correlations
 - Economics

Freeway Traffic Data

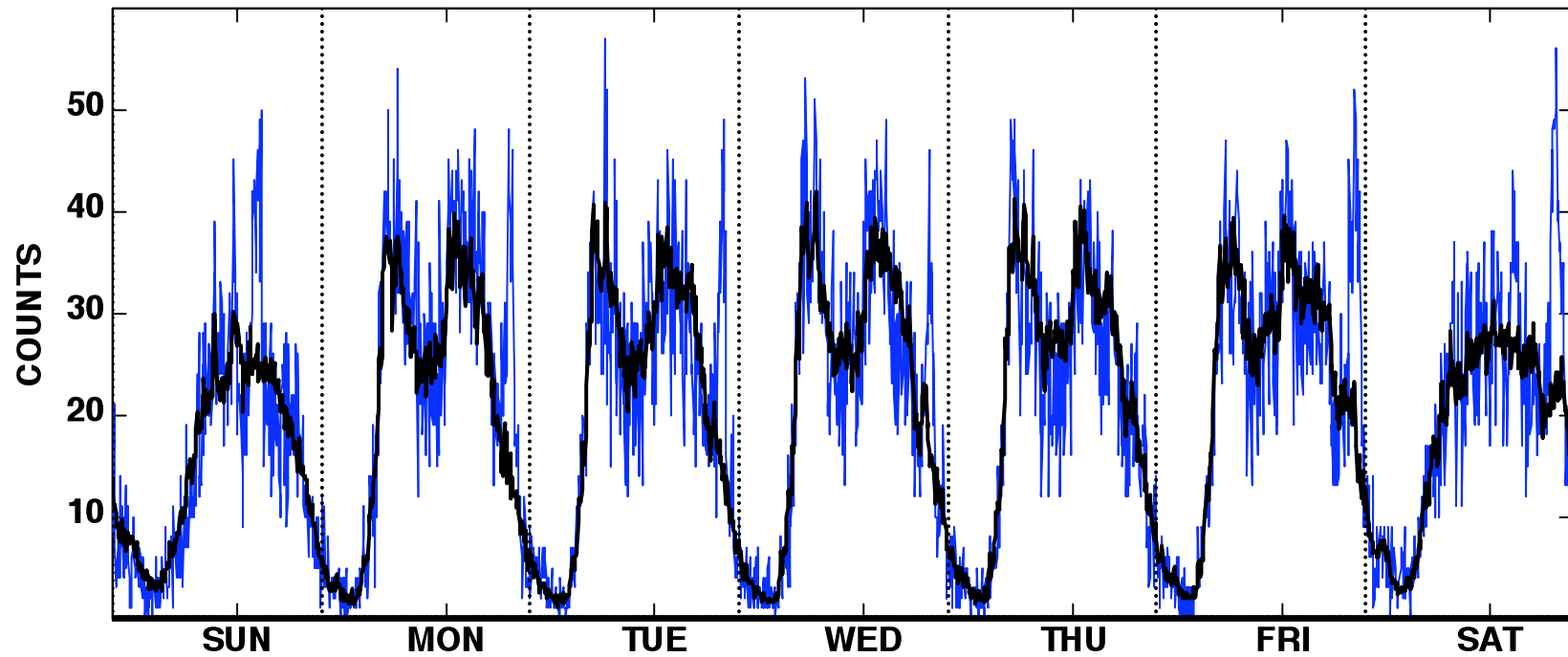
- Onramp sensor near Dodger baseball stadium
- Aggregated vehicle counts every 5 minutes
- 2 years of data
- Partial ground truth for event data = baseball games



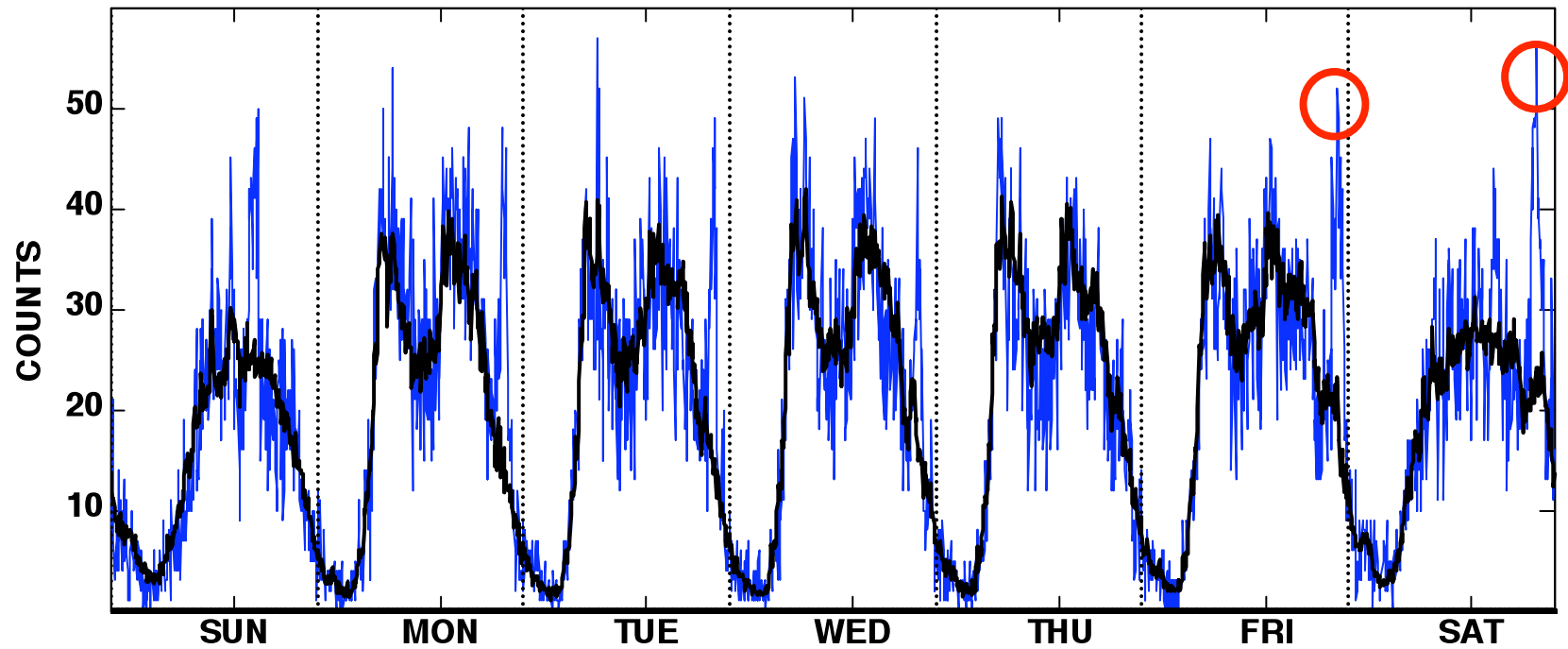
One Week of Traffic Data



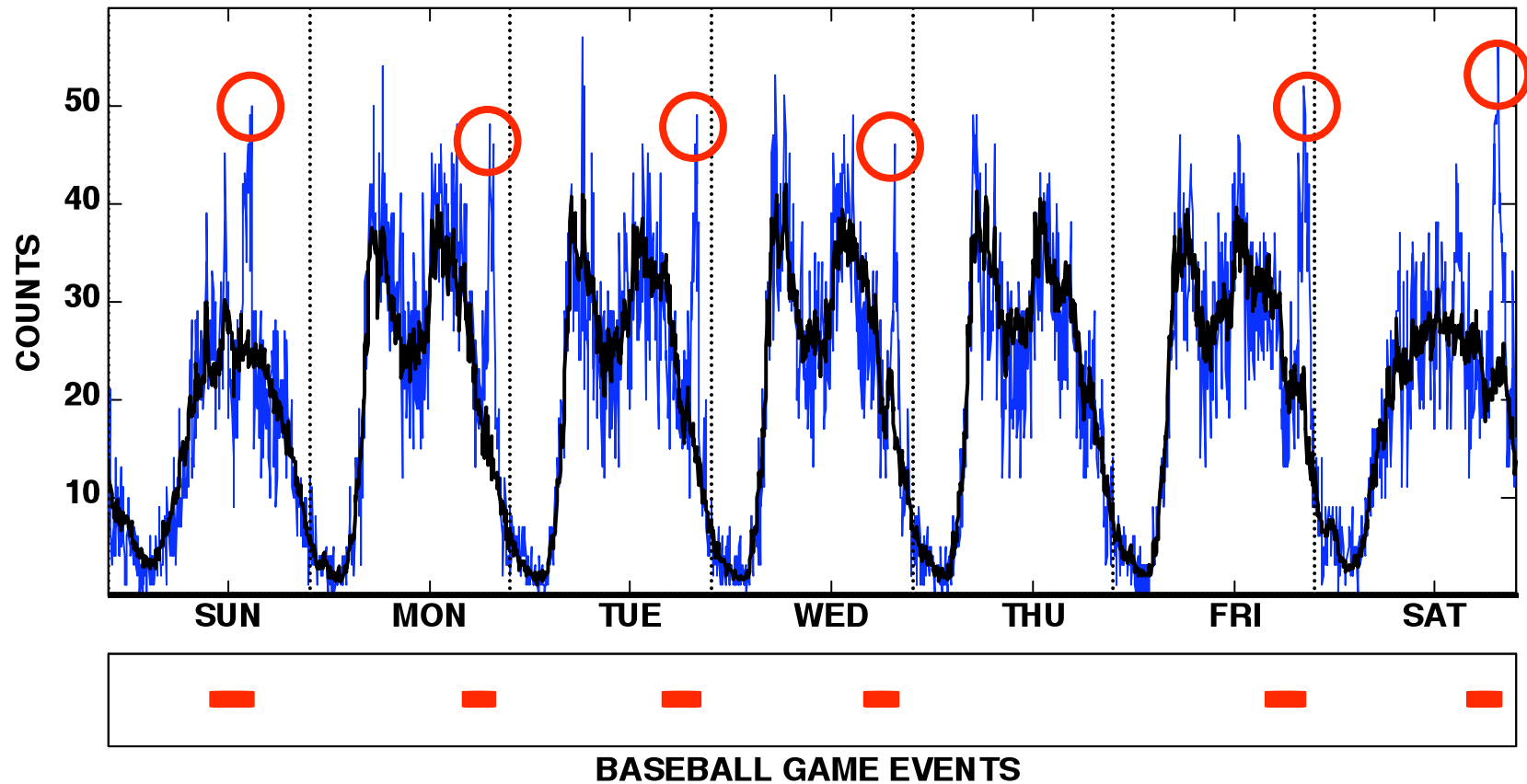
One Week of Traffic Data



One Week of Traffic Data



One Week of Traffic Data



“Event” = large scale activity,
unusual relative to normal patterns

Building Entry Data

Optical sensors at door entrances of Calit2 building on UC Irvine campus

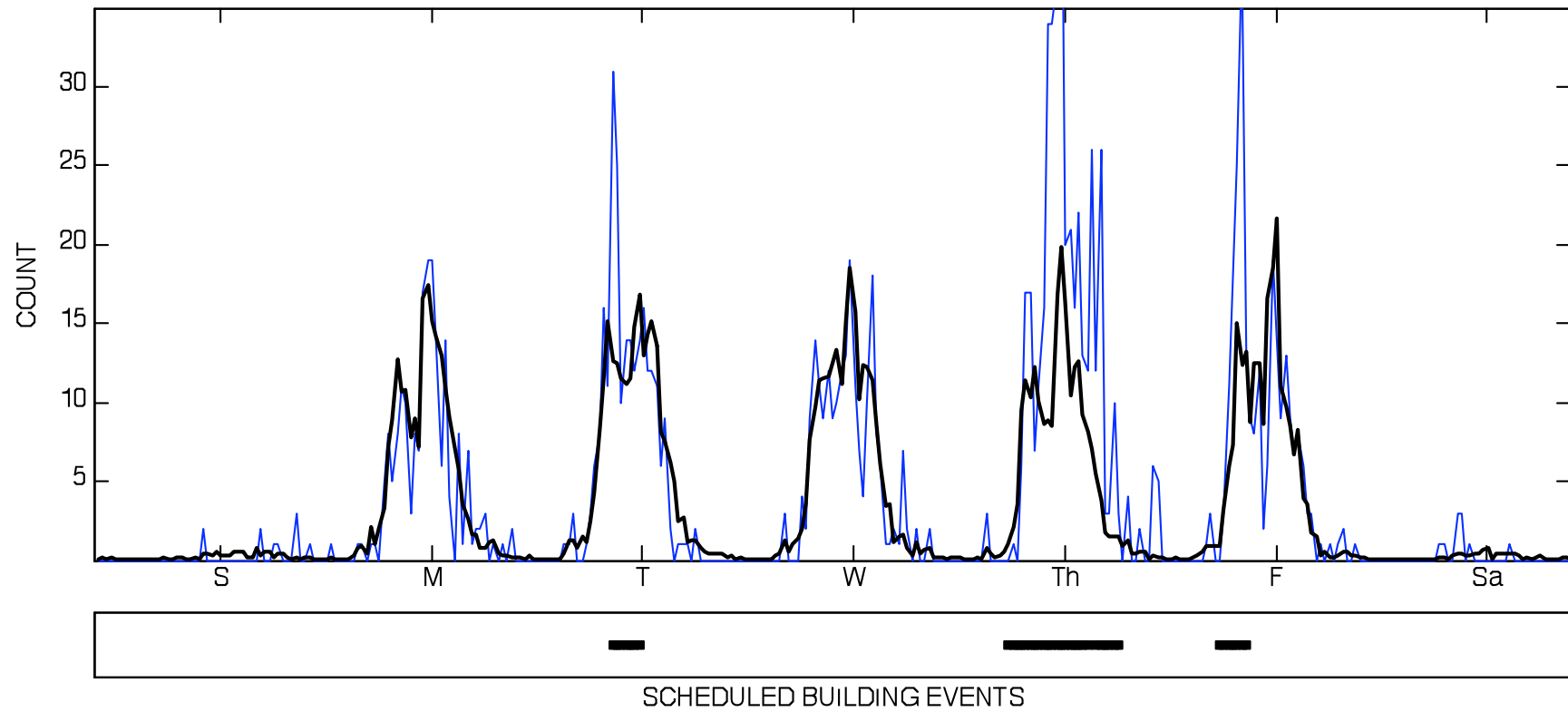
Aggregated counts every 30 minutes

6 months of data

Partial ground truth for events = building calendar



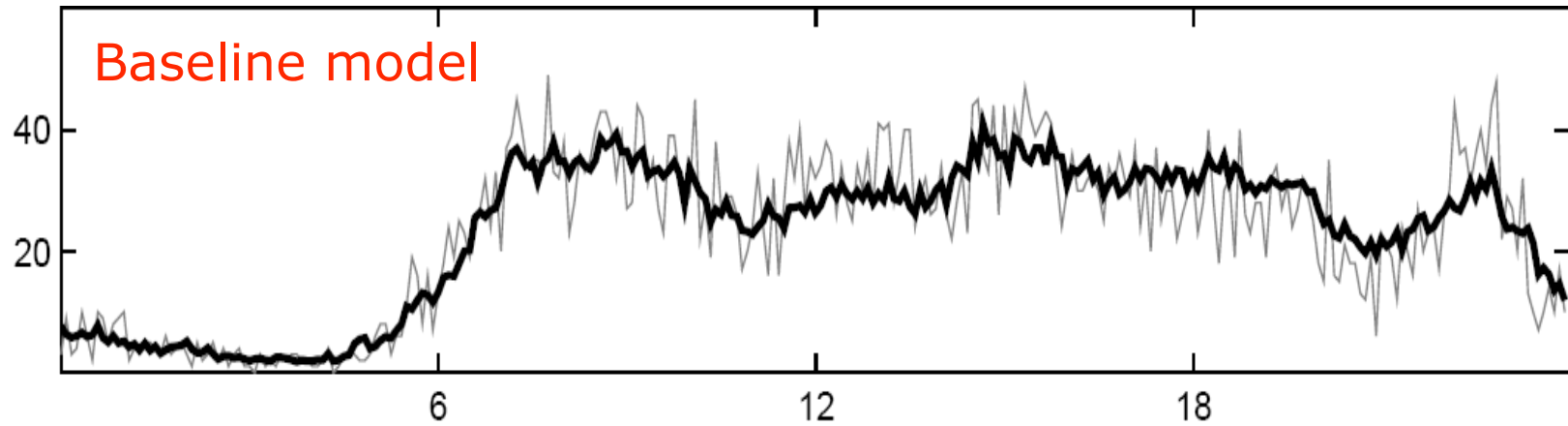
1 Week of Building Data



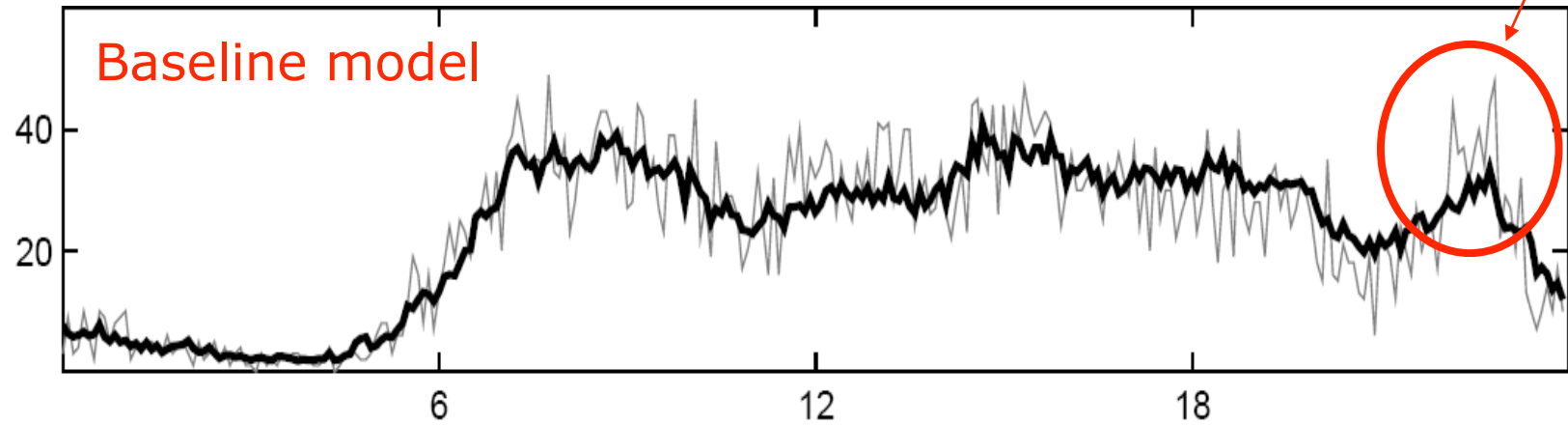
Problem Statement

- Given:
 - Data in the form of count time-series (of people, vehicles, etc)
 - Regular patterns of behavior + bursty events
 - Unlabeled
- Learn:
 - How to model normal recurrent behavior
 - How to detect abnormal events

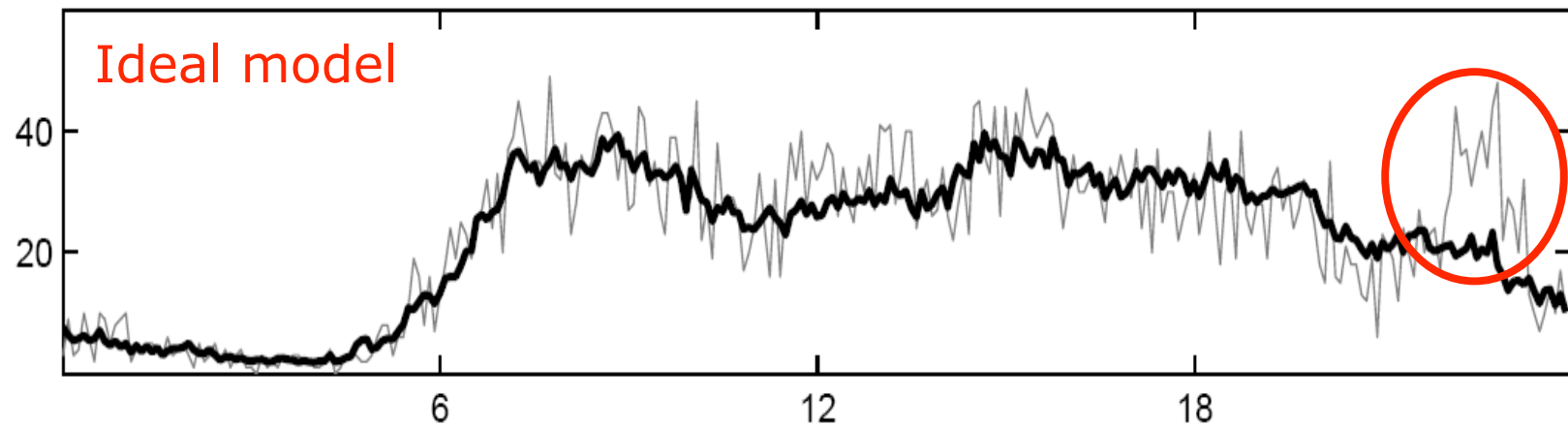
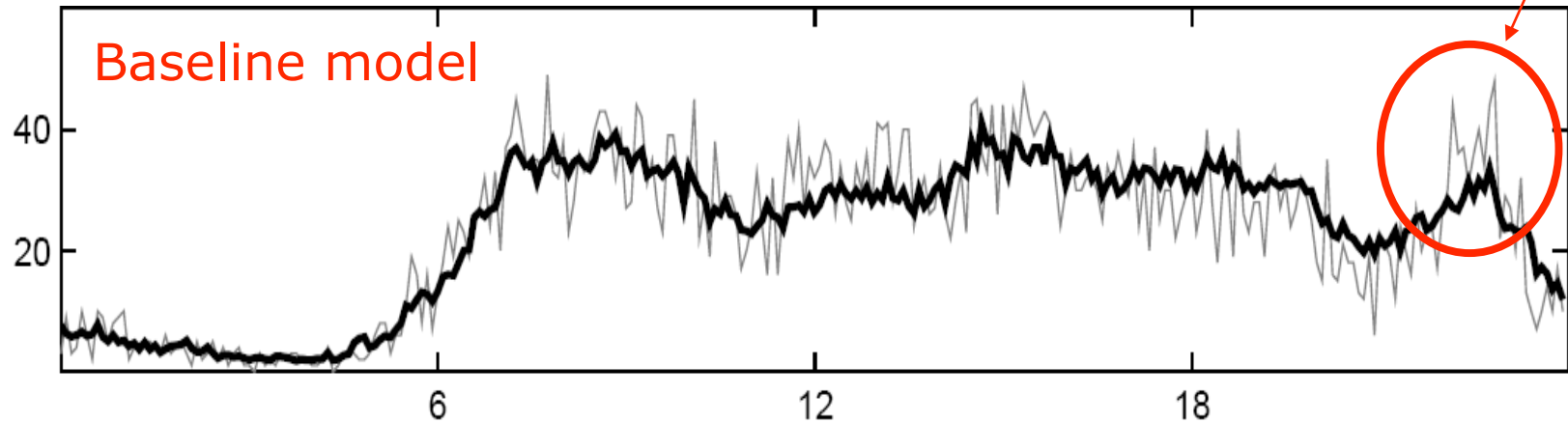
The Chicken and Egg Problem



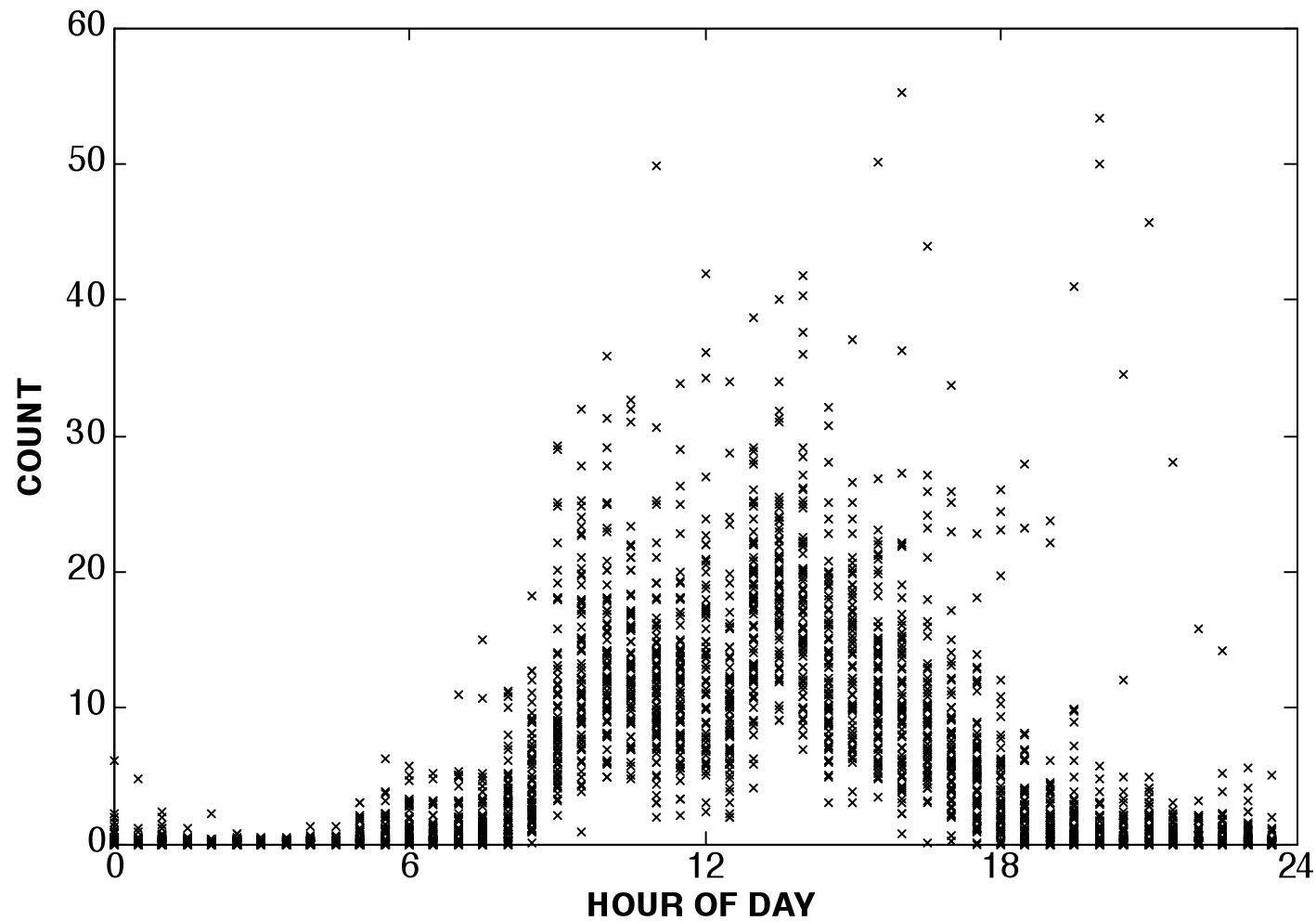
The Chicken and Egg Problem



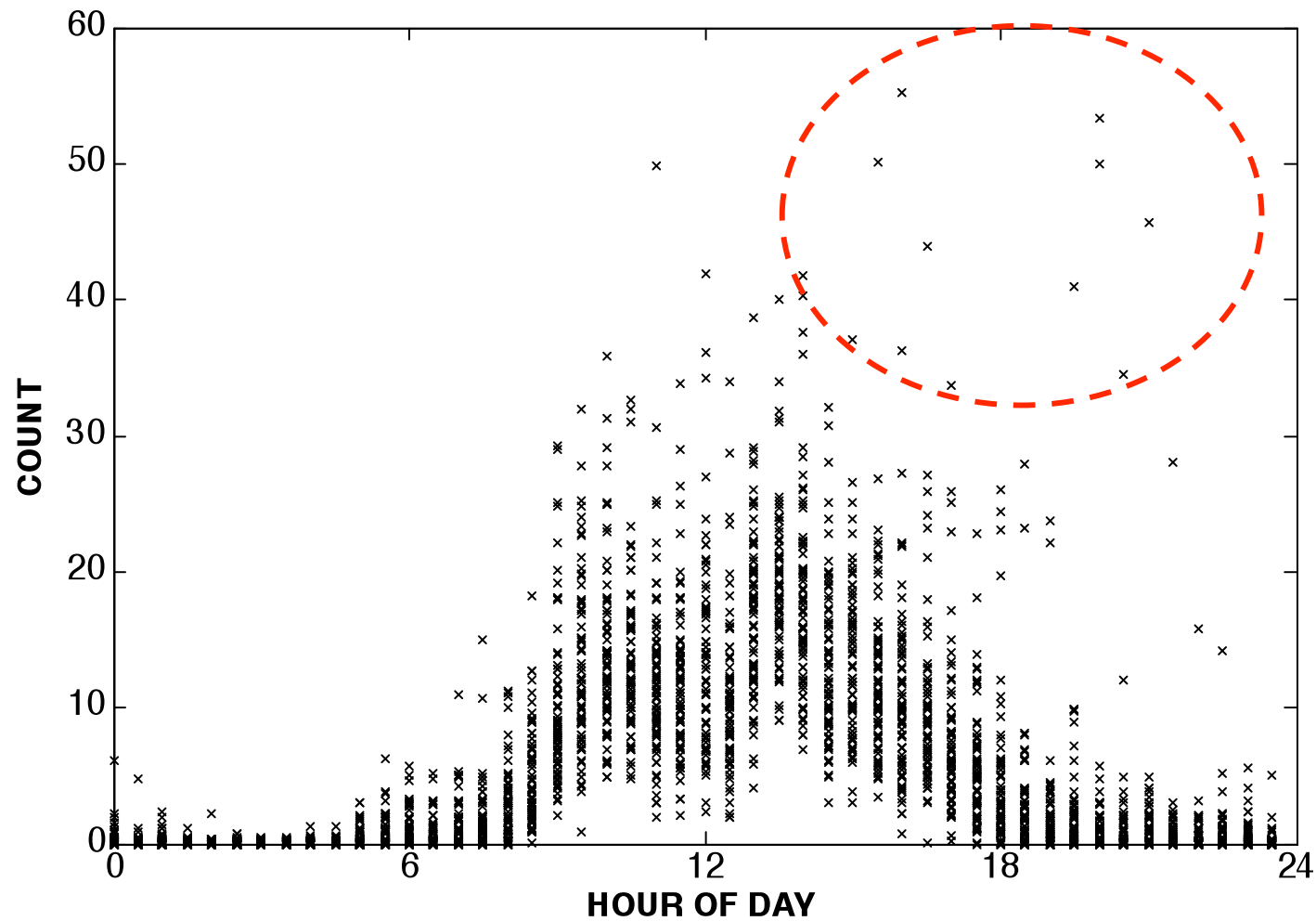
The Chicken and Egg Problem



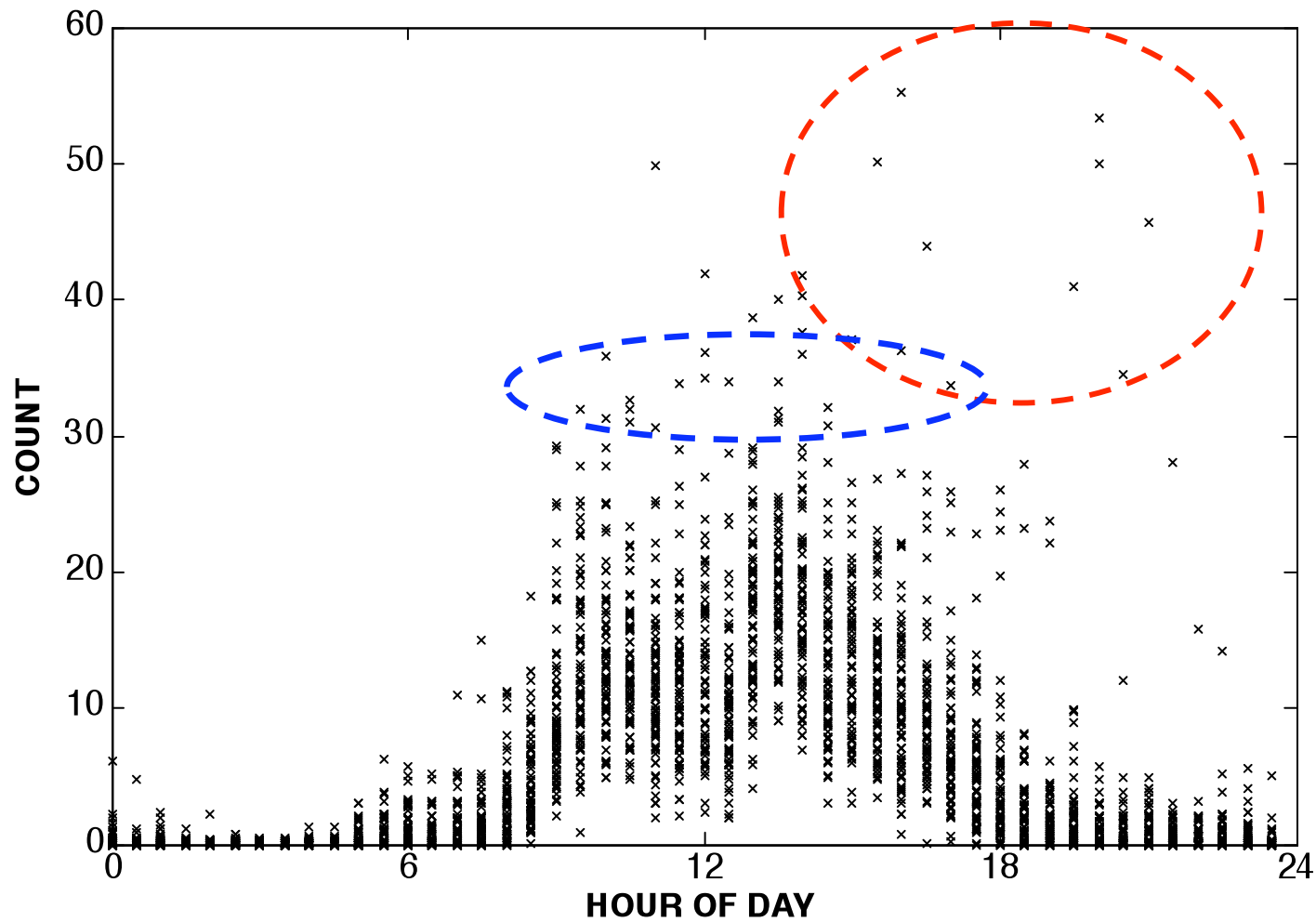
Why this problem is hard....



Why this problem is hard....



Why this problem is hard....



Proposed Model

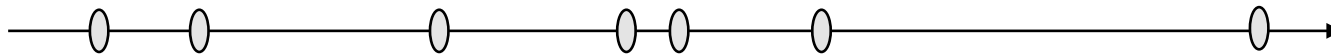
$$N(t) = N_0(t) + N_E(t)$$

OBSERVED
COUNT

NORMAL
COUNT
(UNOBSERVED)

EVENT
COUNT
(UNOBSERVED)

Poisson Process

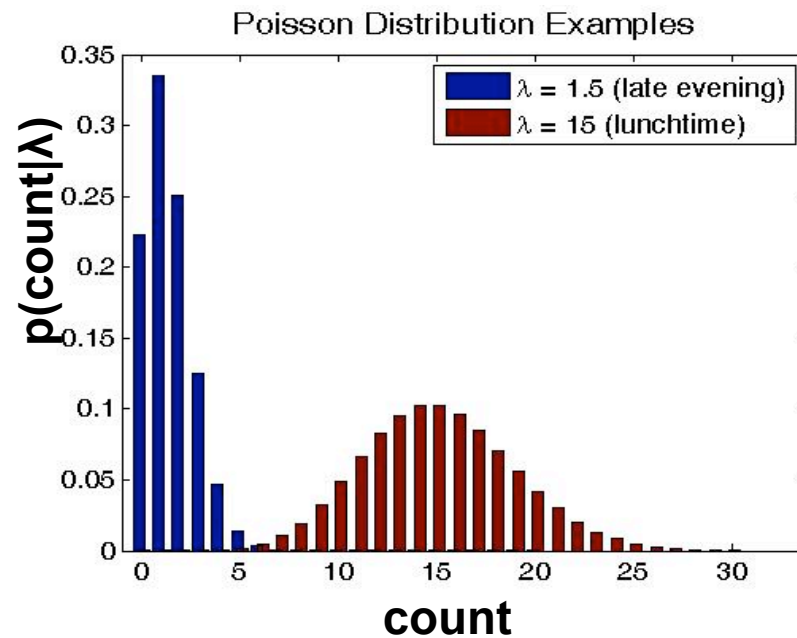


- Simple stochastic model for count data
 - Single parameter λ = “rate”

λ = expected number of events per unit time

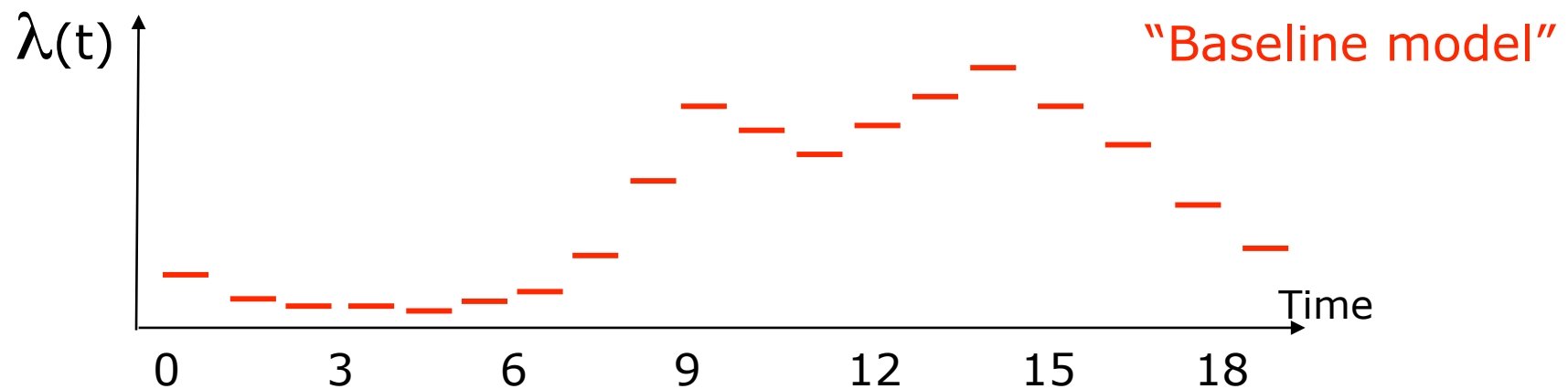
Time-Varying Poisson Process

- The rate $\lambda(t)$ is a function of time (“inhomogeneous”)
- Simple parametrization
 - Different λ 's for different discrete times of day, day of week

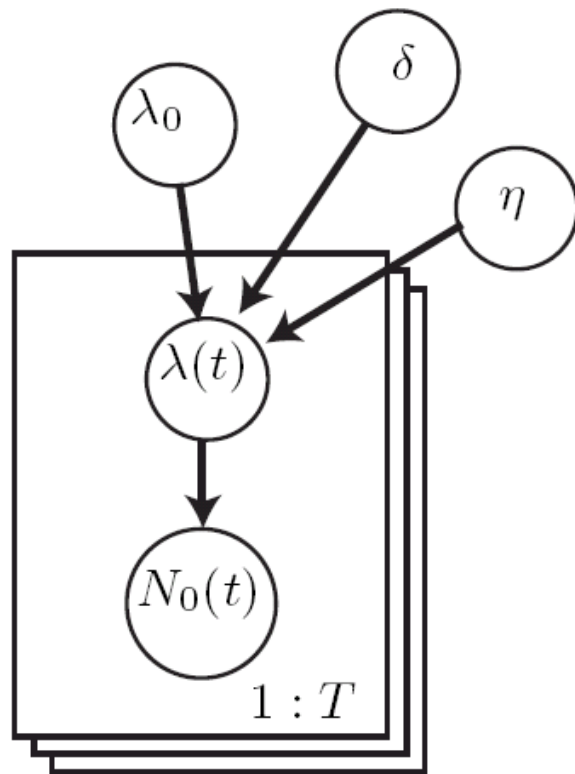


Inhomogeneous Poisson Process

- The rate $\lambda(t)$ is a function of time (“inhomogeneous”)
- Piecewise constant parametrization
 - Different λ 's for different discrete times of day, day of week

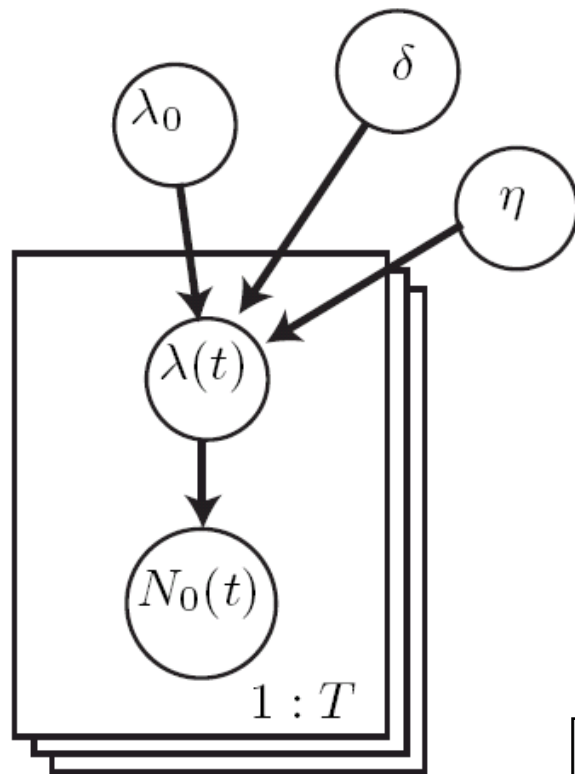


Model for Normal Counts



$$\lambda(t) = \lambda_0 \delta_{d(t)} \eta_{d(t),h(t)}$$

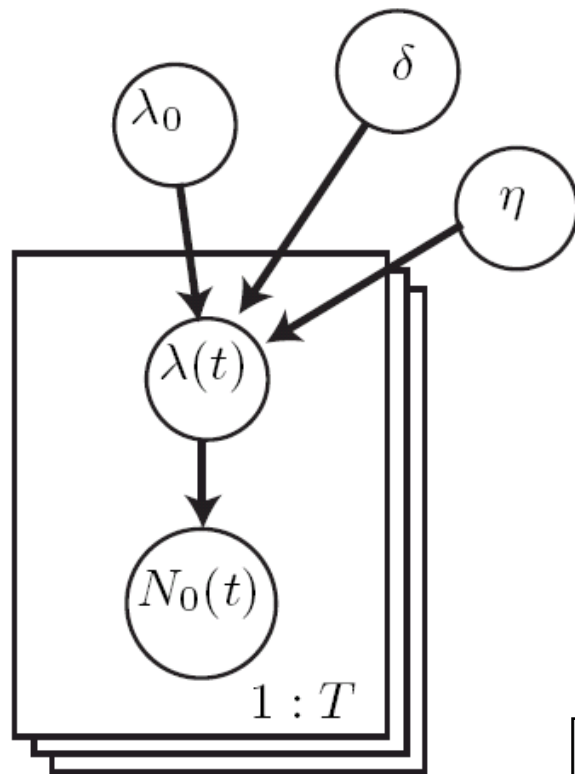
Model for Normal Counts



$$\lambda(t) = \lambda_0 \delta_{d(t)} \eta_{d(t),h(t)}$$

Overall mean rate

Model for Normal Counts

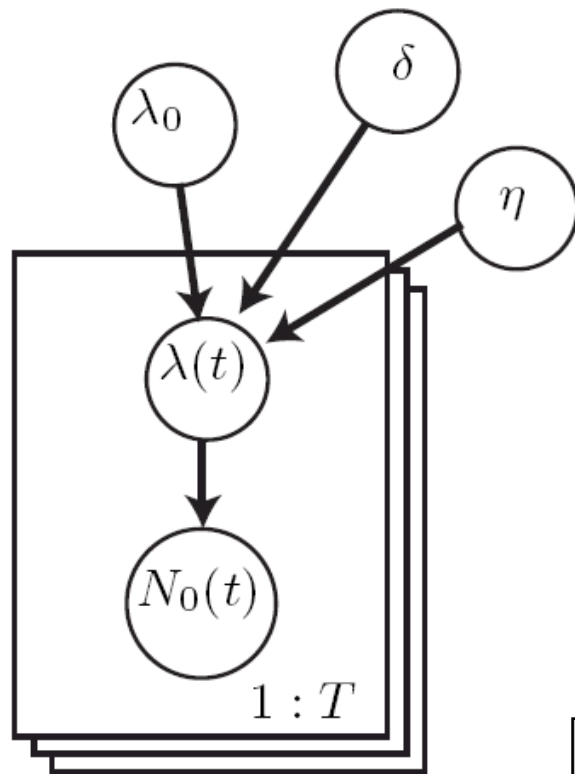


$$\lambda(t) = \lambda_0 \delta_{d(t)} \eta_{d(t),h(t)}$$

Daily rate component

Overall mean rate

Model for Normal Counts



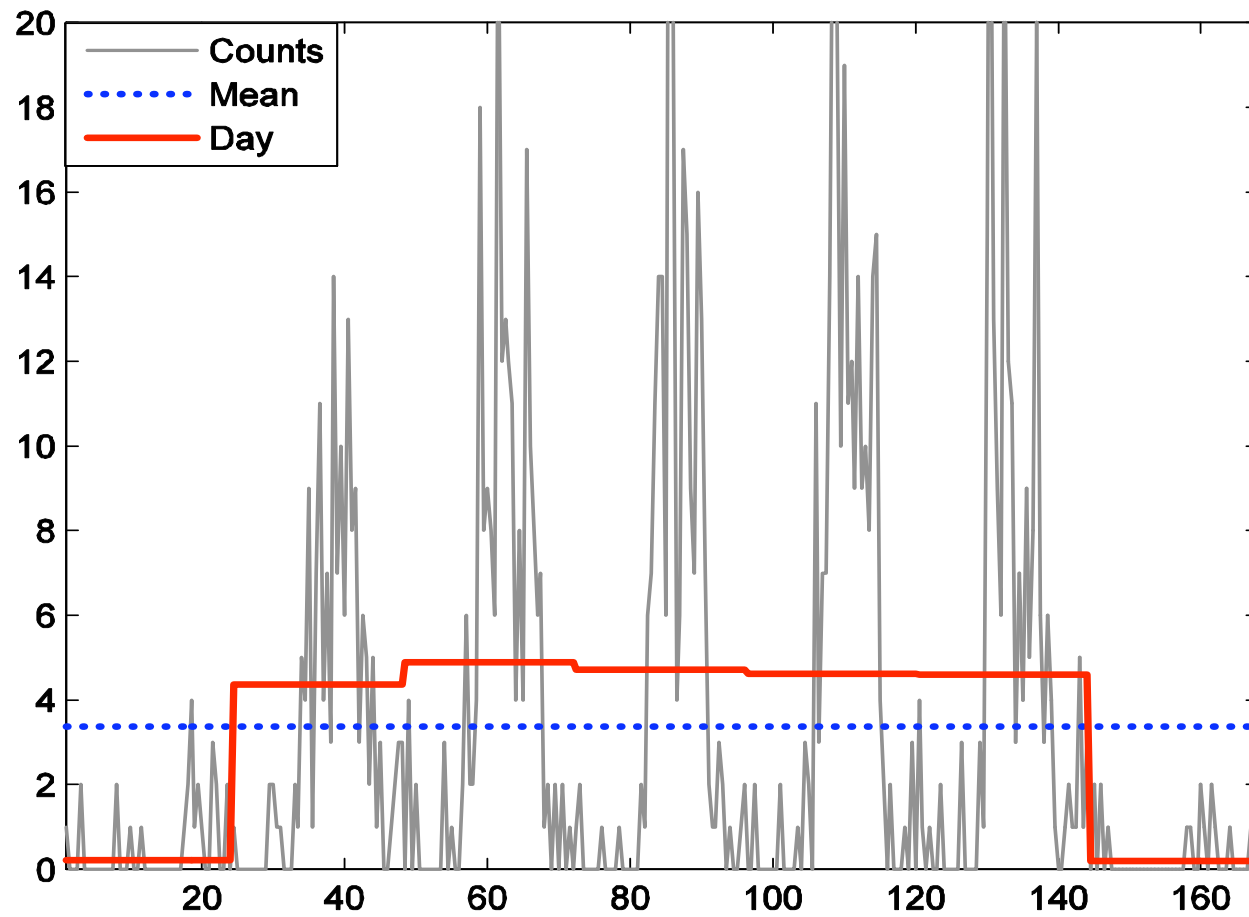
$$\lambda(t) = \lambda_0 \delta_{d(t)} \eta_{d(t),h(t)}$$

Daily rate component

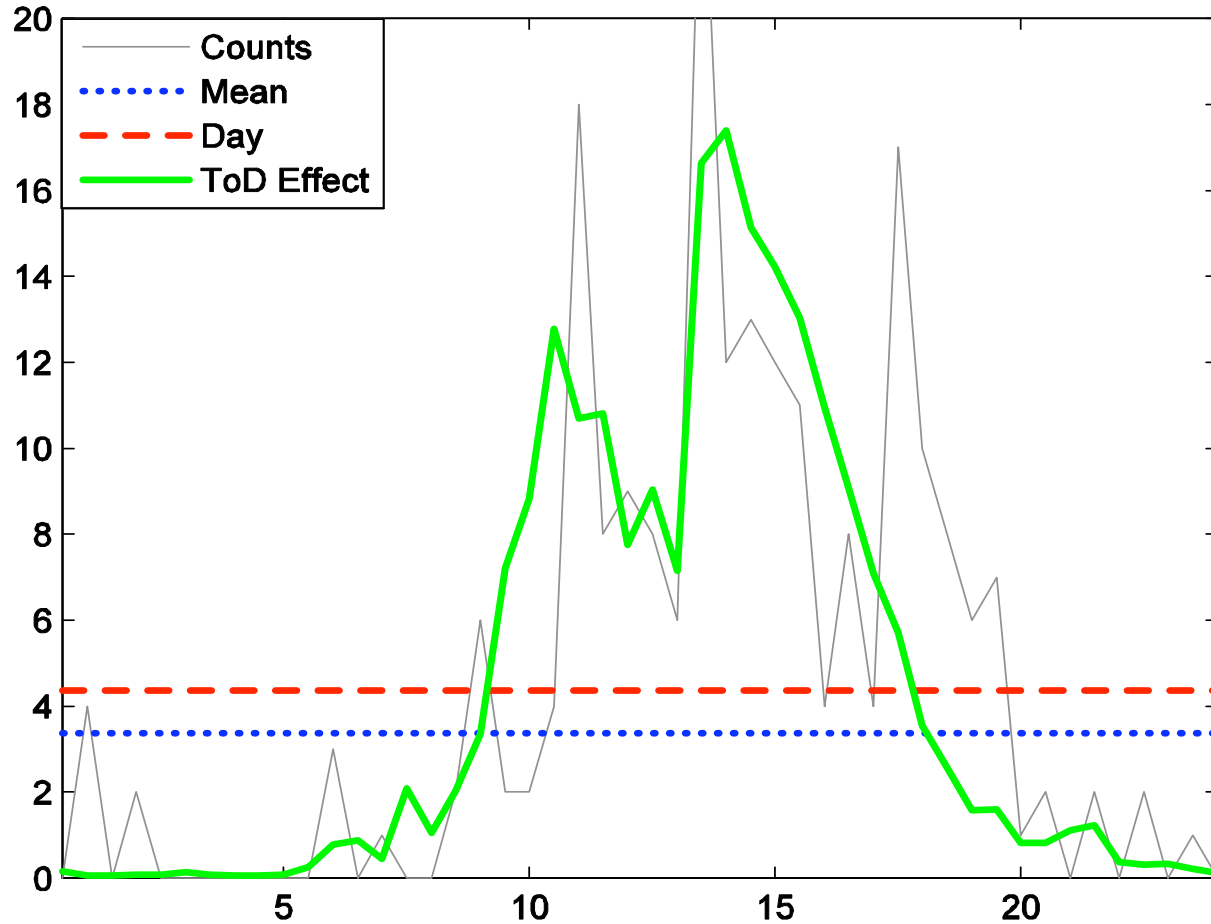
Overall mean rate

Local time bin component, given day and time

Daily Rate Component



Time of Day Rate Component



Proposed Model

$$N(t) = N_0(t) + N_E(t)$$

OBSERVED
COUNT

The diagram illustrates the proposed model equation $N(t) = N_0(t) + N_E(t)$. Three callout boxes are positioned below the equation. The first box, labeled 'OBSERVED COUNT', points to the entire equation. The second box, labeled 'NORMAL COUNT (UNOBSERVED)', points to the $N_0(t)$ term. The third box, labeled 'EVENT COUNT (UNOBSERVED)', points to the $N_E(t)$ term.

NORMAL
COUNT
(UNOBSERVED)

EVENT
COUNT
(UNOBSERVED)

Proposed Model

$$N(t) = N_0(t) + N_E(t)$$

OBSERVED
COUNT

NORMAL
COUNT
(UNOBSERVED)

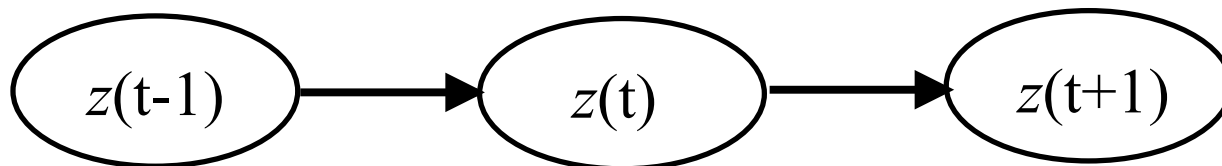
EVENT
COUNT
(UNOBSERVED)

Time-varying
Poisson

Markov with
Poisson counts

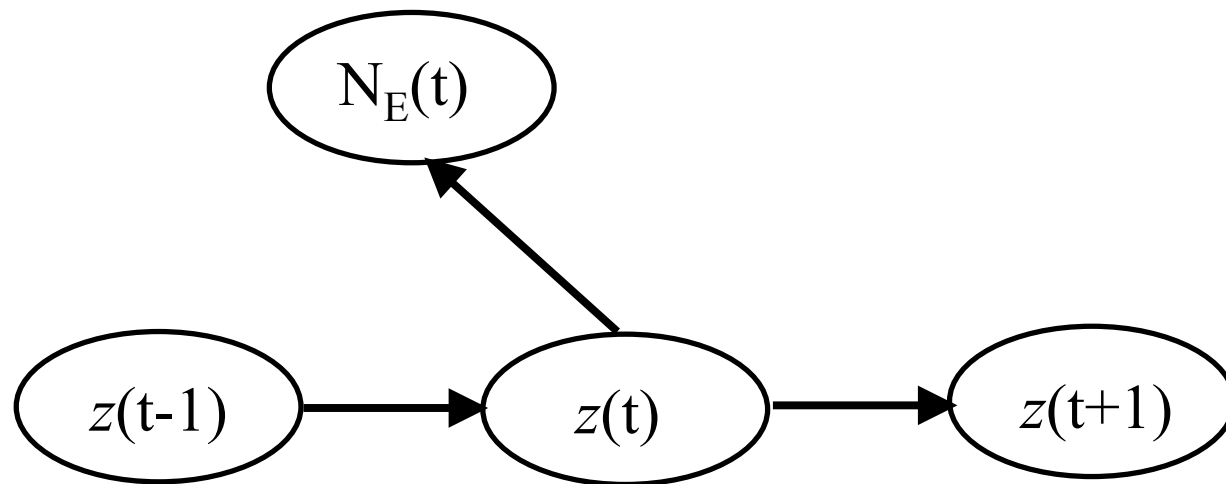
Adding a Hidden Event Process

$$z(t) = \begin{cases} 0 & \text{if there is no event at time } t \\ +1 & \text{if there is a positive event} \\ -1 & \text{if there is a negative event} \end{cases}$$

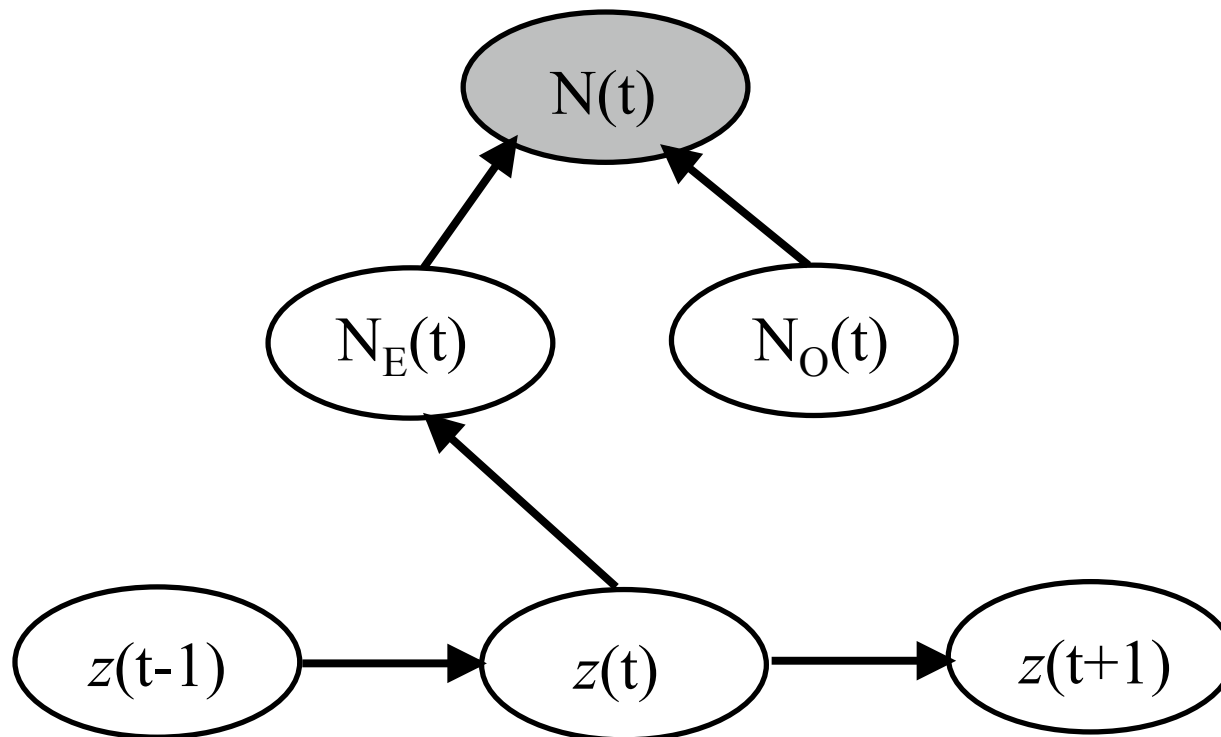


Adding a Hidden Event Process

$$N_E(t) \sim \begin{cases} 0 & z(t)=0 \\ \text{Poisson}(N; \gamma) & z(t)=1 \text{ or } -1 \end{cases}$$

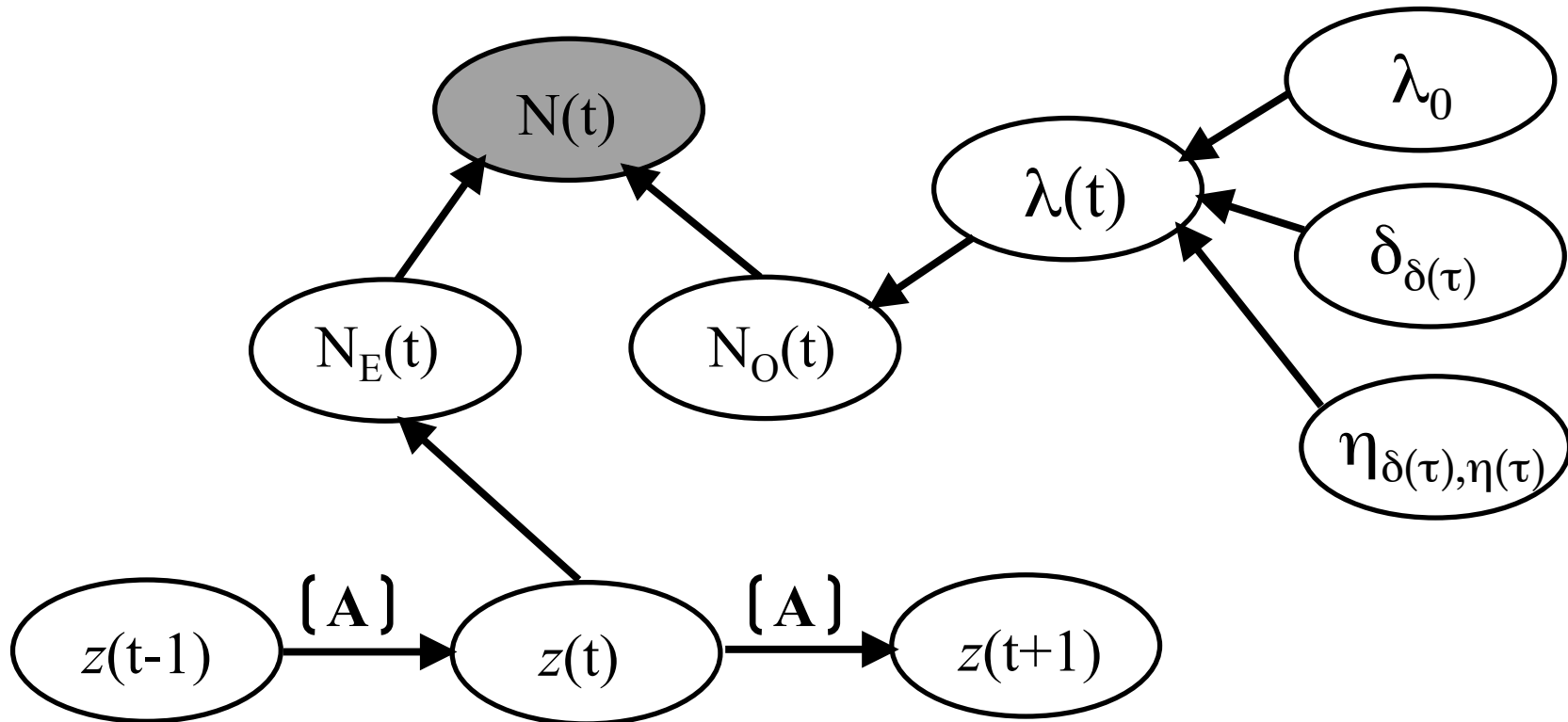


Adding a Hidden Event Process



Markov Modulated Poisson Process (MMPP) model
e.g., see Heffes and Lucantoni (1994), Scott (1998)

Graphical Model



Bayesian Learning

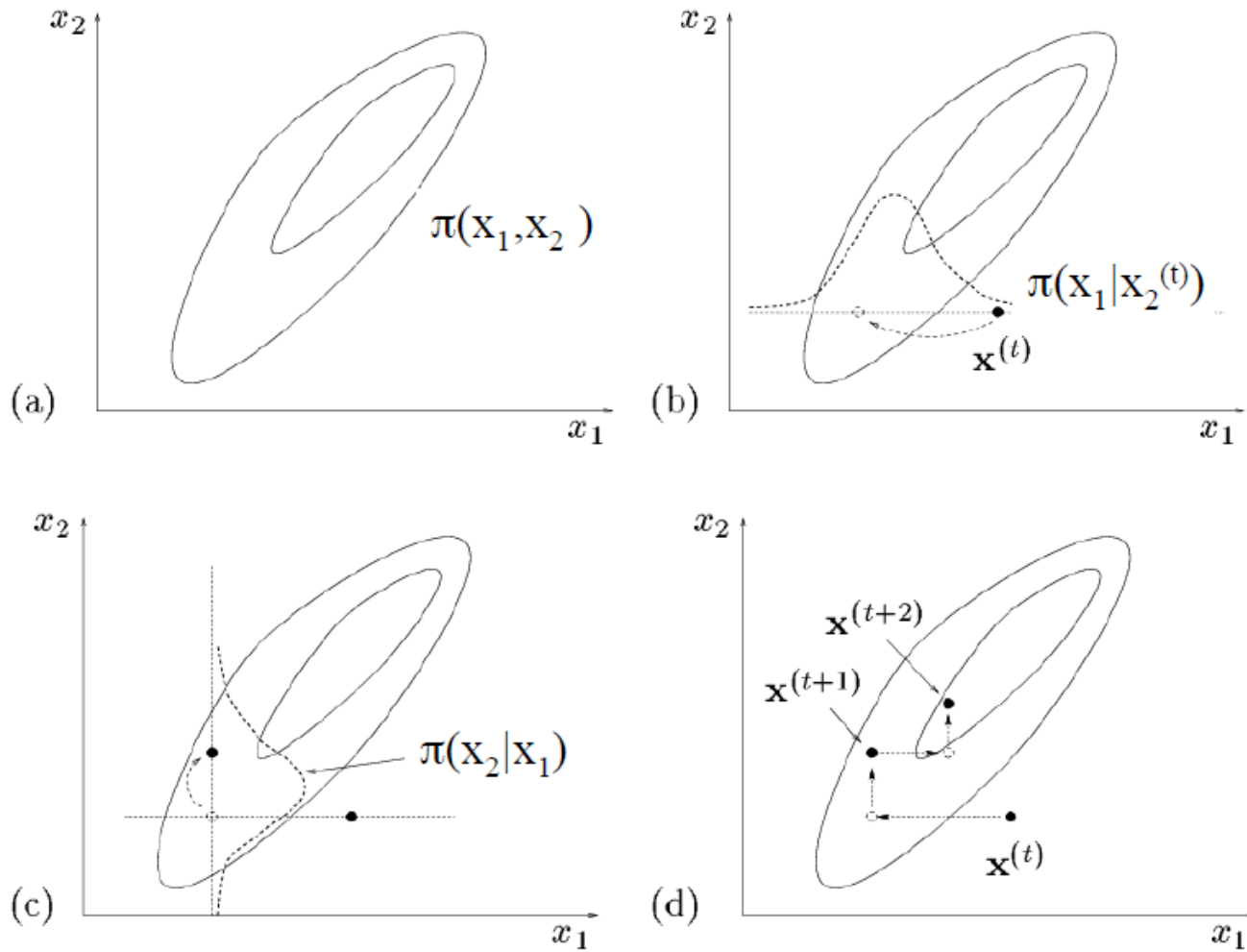
- Unknowns:
 - Poisson parameters: $\lambda, \delta(t), \eta(t)$
 - Markov parameters: $A, N_E(t)$
 - Hidden state sequence: $z(t)$

- Use uninformative priors in general
 - E.g., Gamma priors on Poisson rates
 - Informative priors on event transition probabilities

- Use Gibbs sampling to get a sample-based estimate of the posterior distribution given the counts

- Time complexity is linear in T , number of time points
 - Typically 10 iterations for burn-in, and then 50-100 samples after that

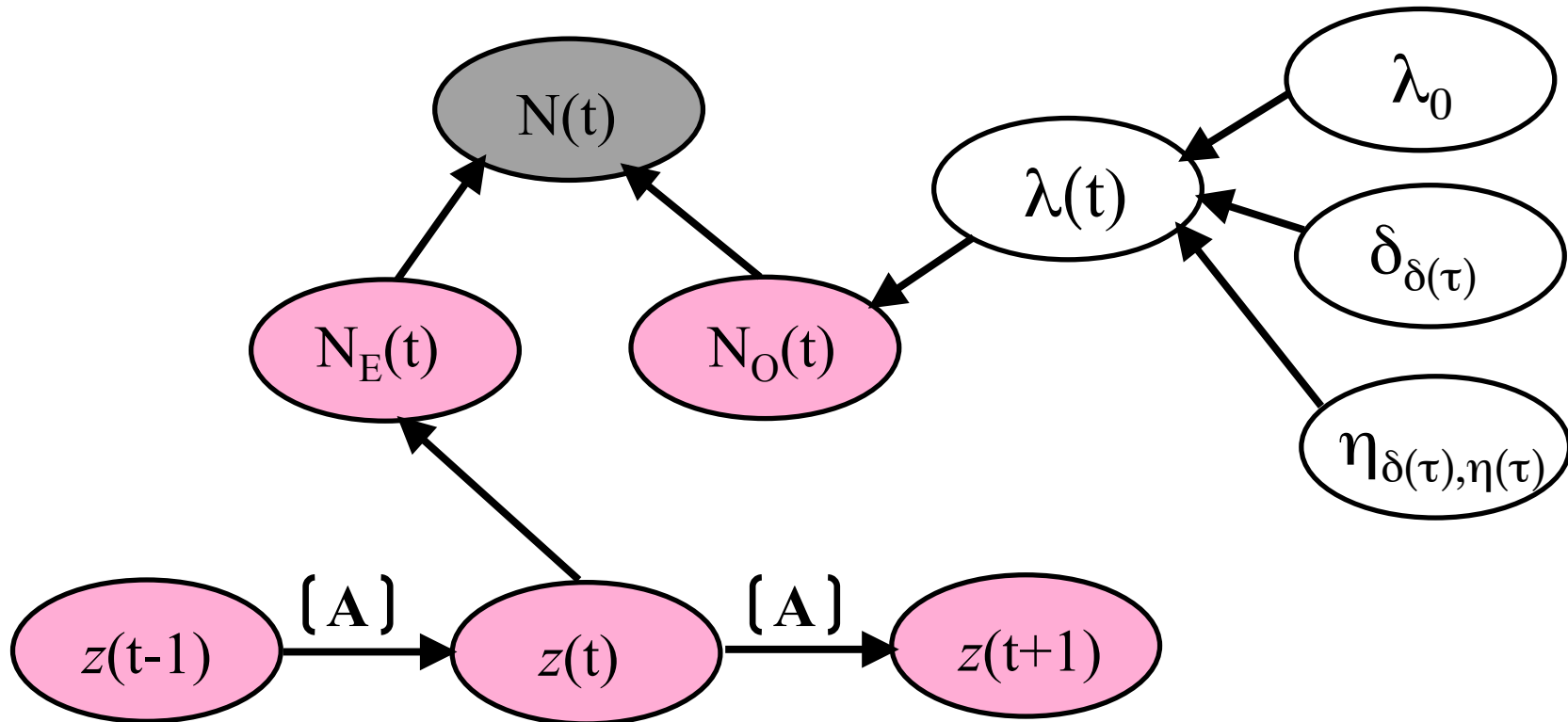
Gibbs Sampling



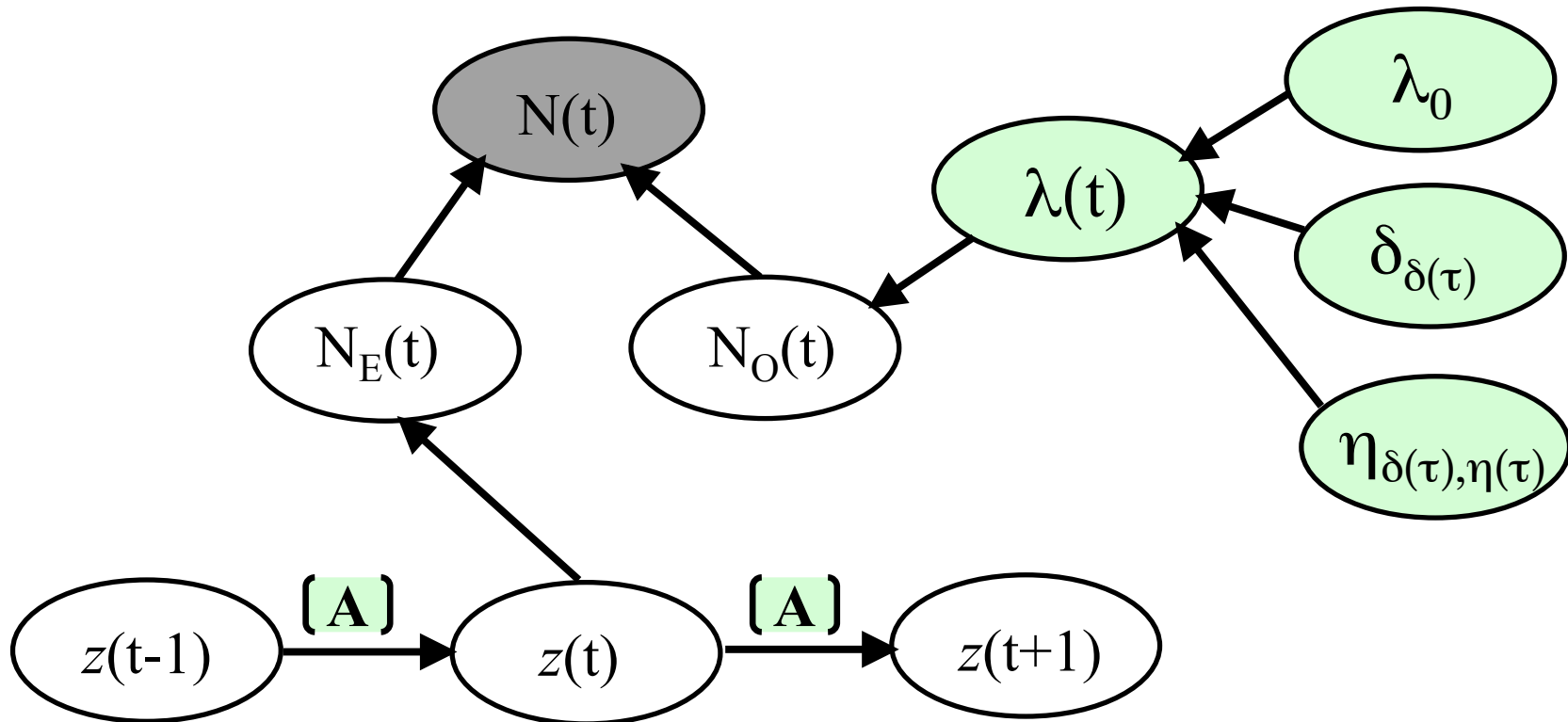
From online MCMC tutorial notes by Frank Dellaert, Georgia Tech

Sampling the Hidden Variables

S. Scott, JASA, 2002



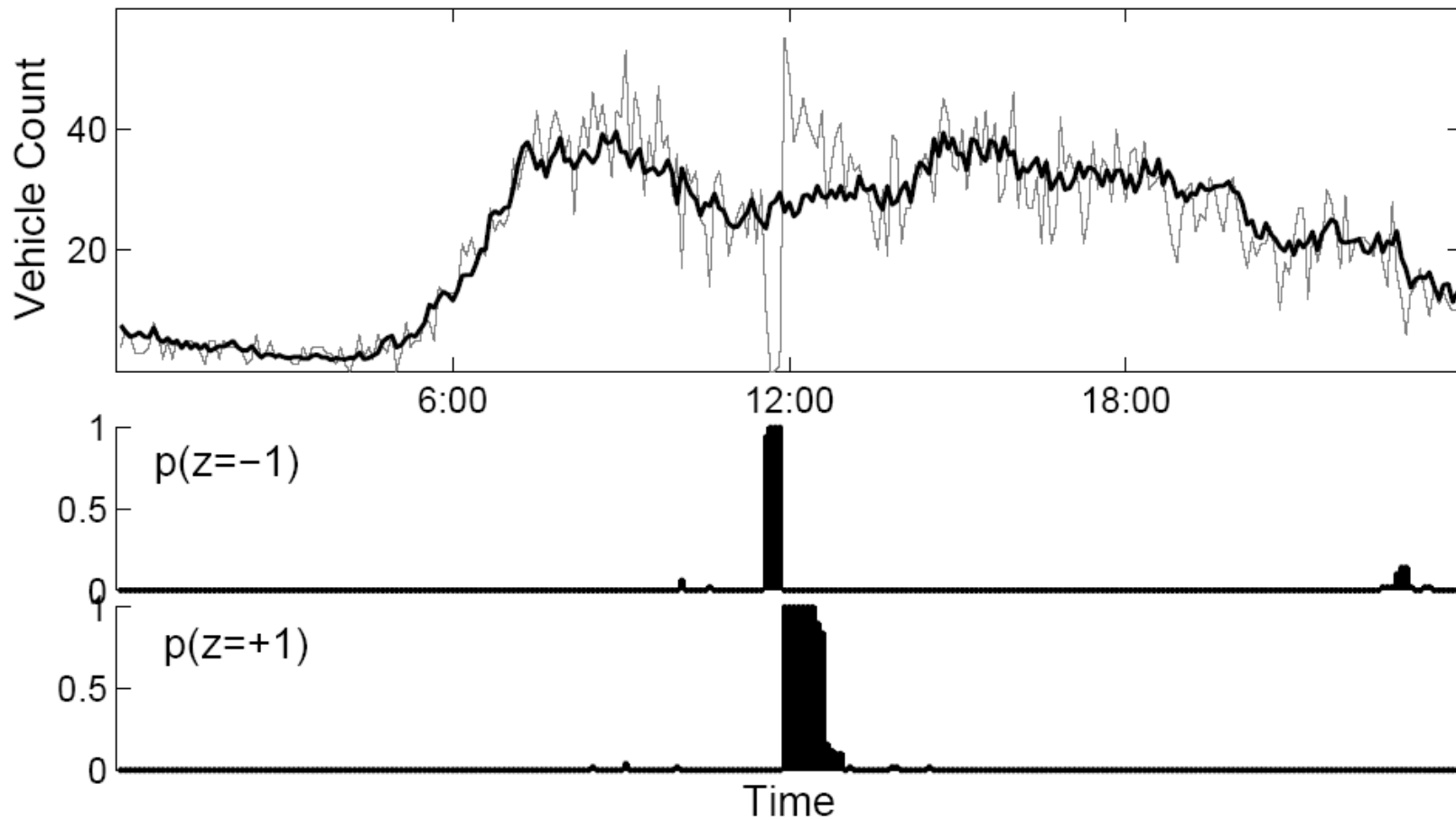
Sampling the Parameters



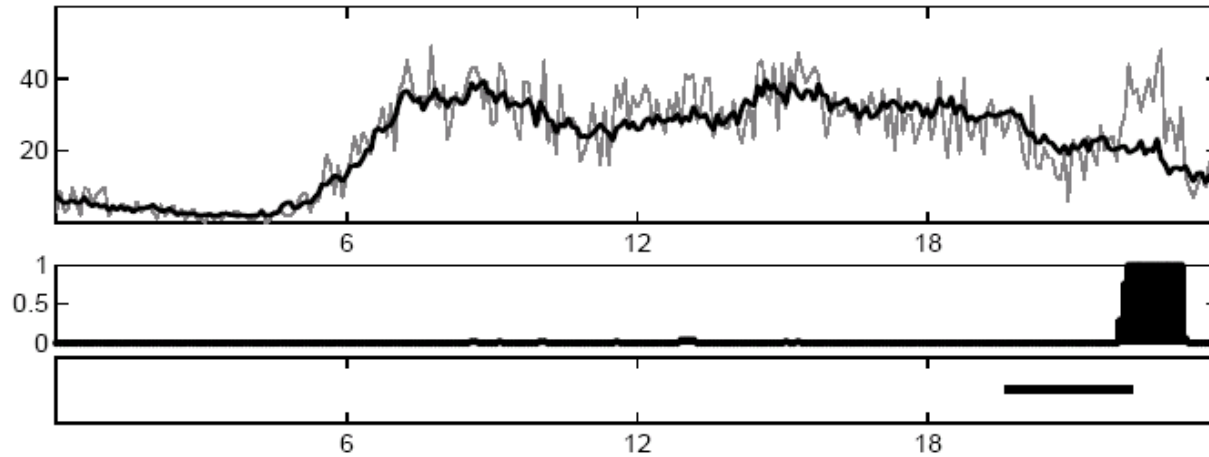
Experimental Results

- Fit our probabilistic model to sensor streams:
 - UC Irvine building people counts
 - Freeway on-ramp near Dodger Stadium
- Perform inference on several months of data
 - Order of several minutes using Gibbs sampling
- Evaluate models by comparing with (hidden) ground truth
 - Model detects events when $p(z = 0) < 0.5$

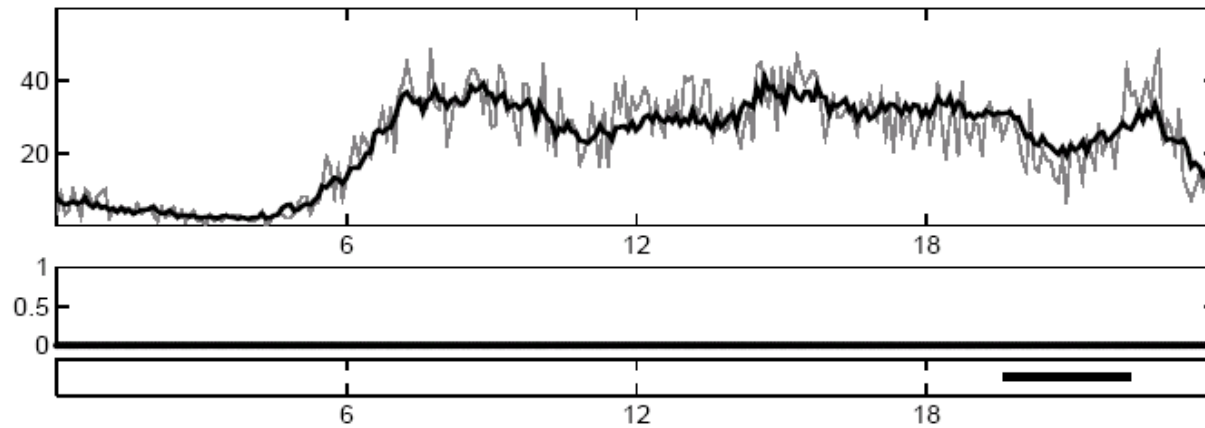
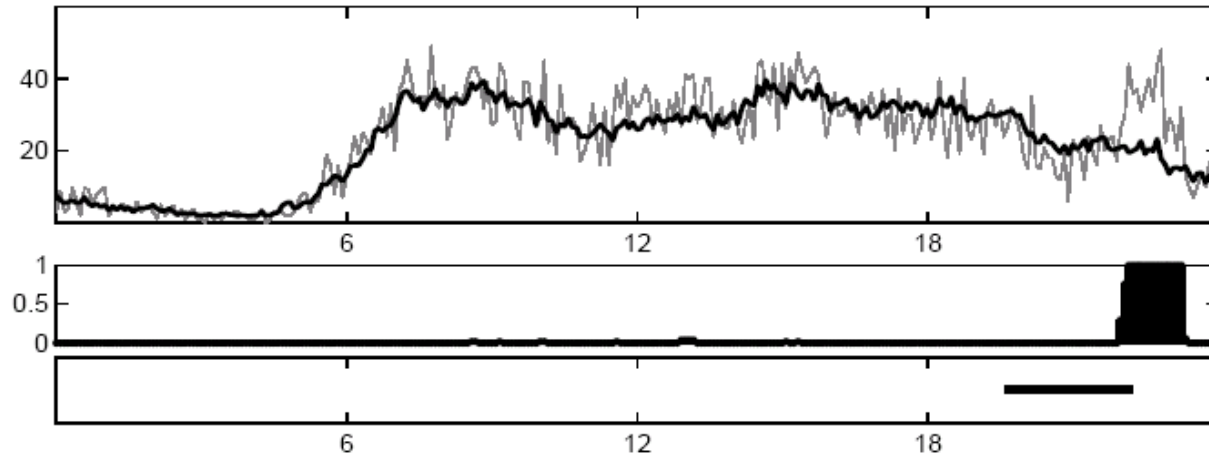
Example: Traffic Data



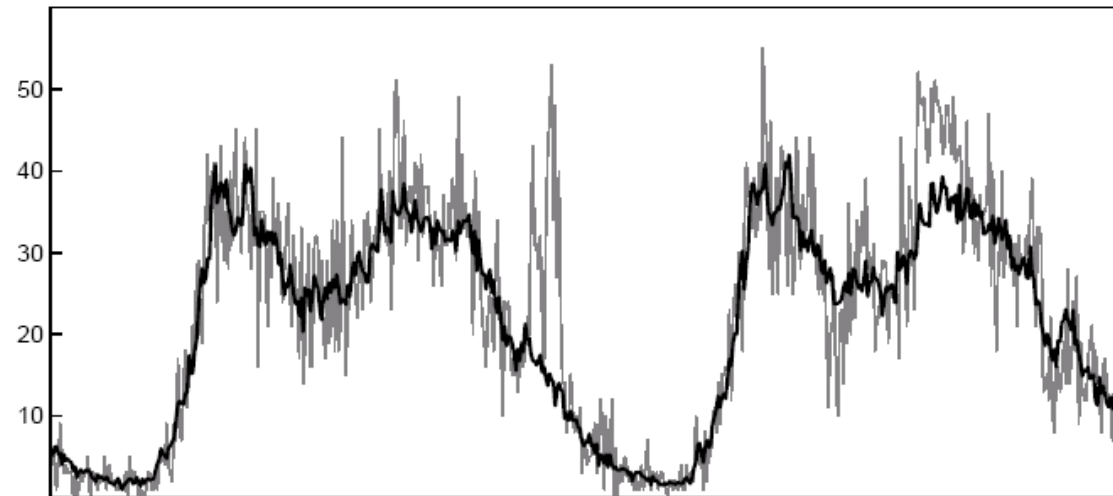
Detection of Friday Night Baseball Game



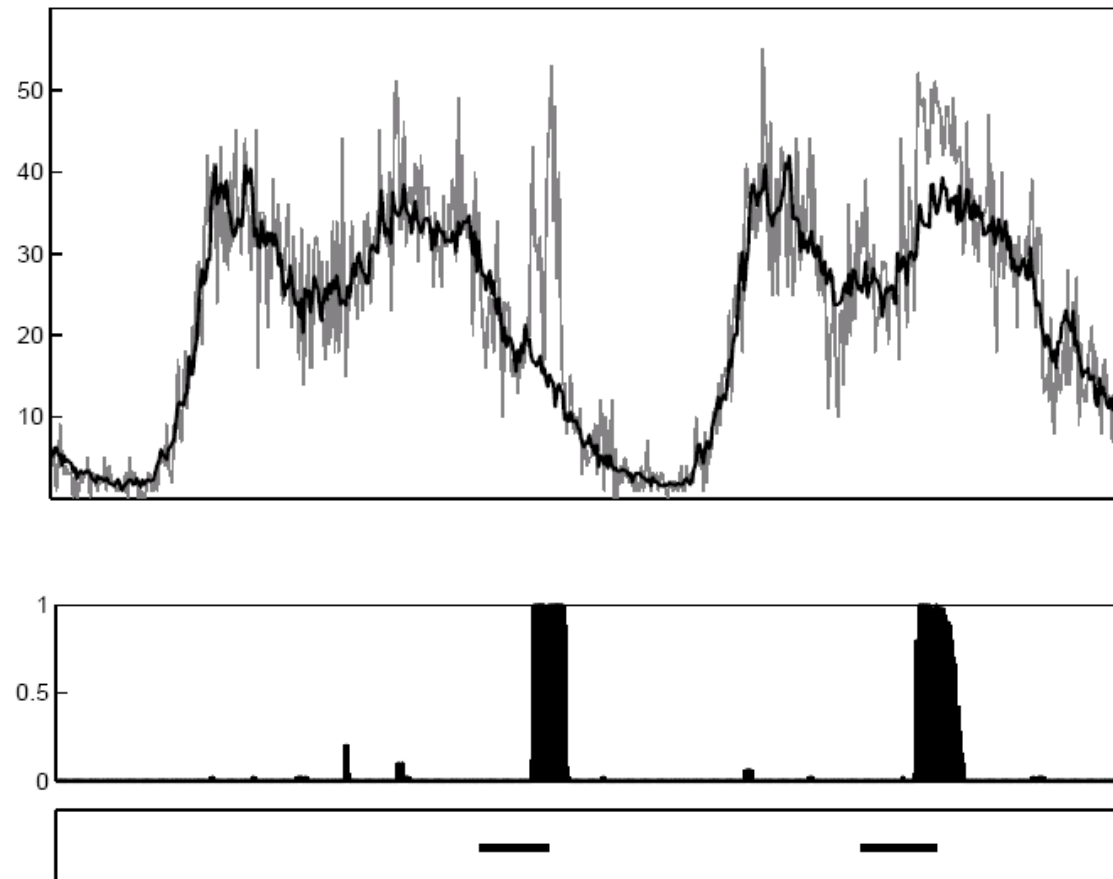
Detection of Friday Night Baseball Game



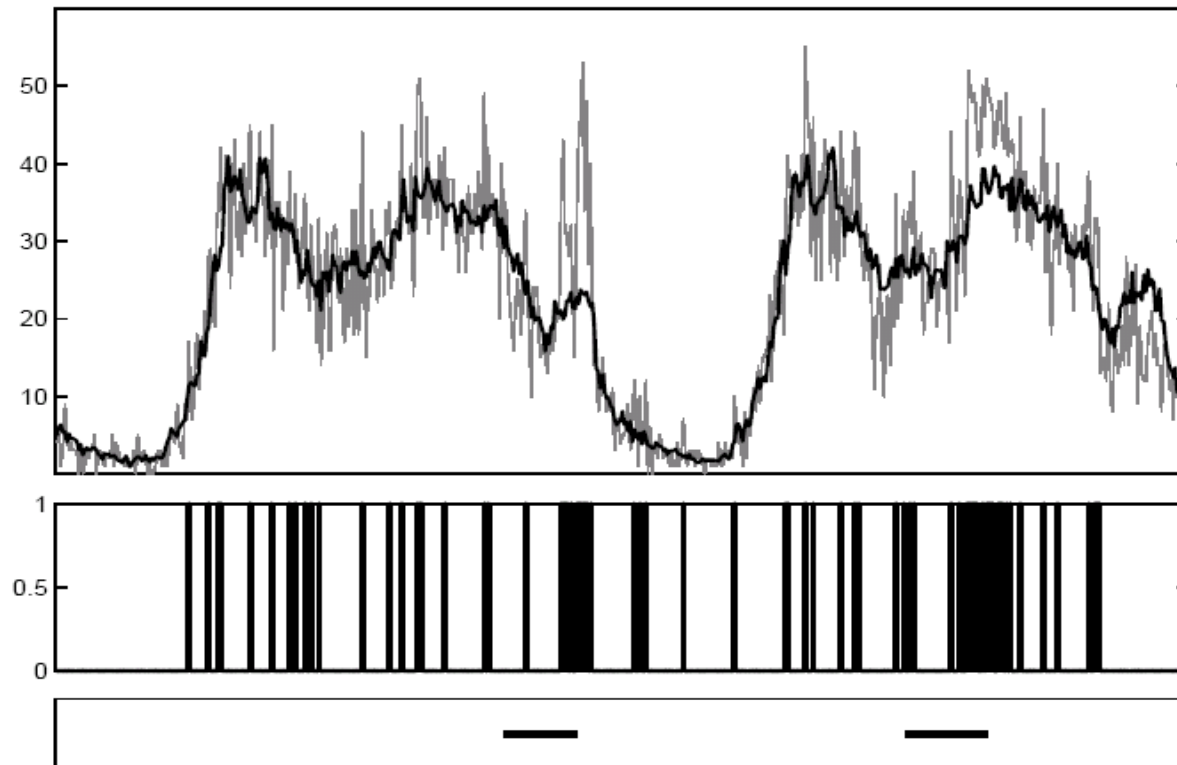
Detection of Baseball Games

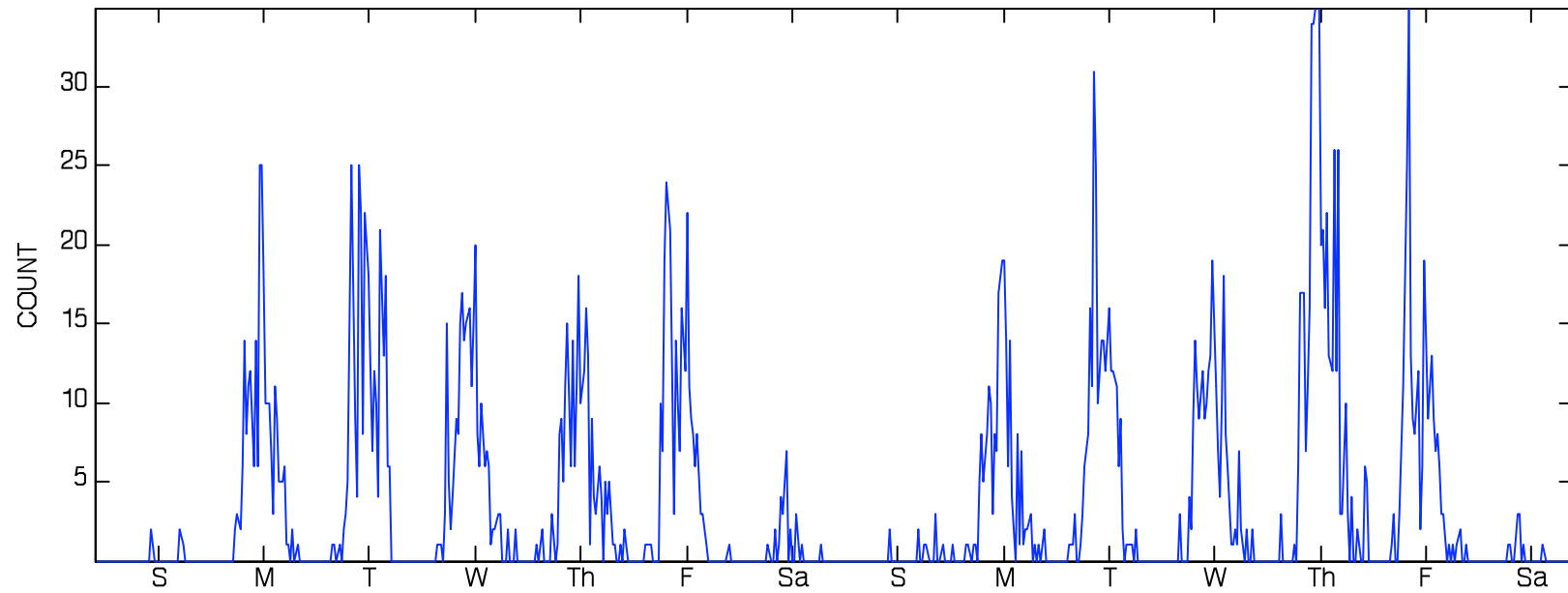


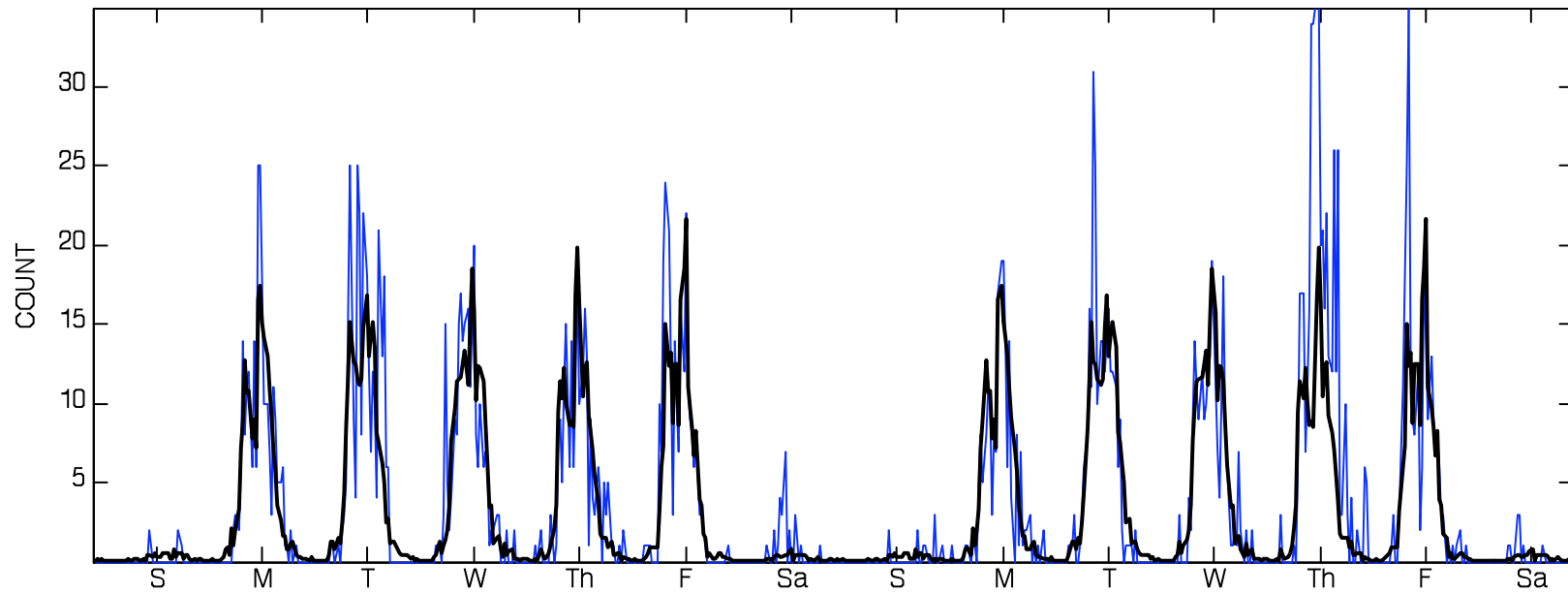
Detection of Baseball Games

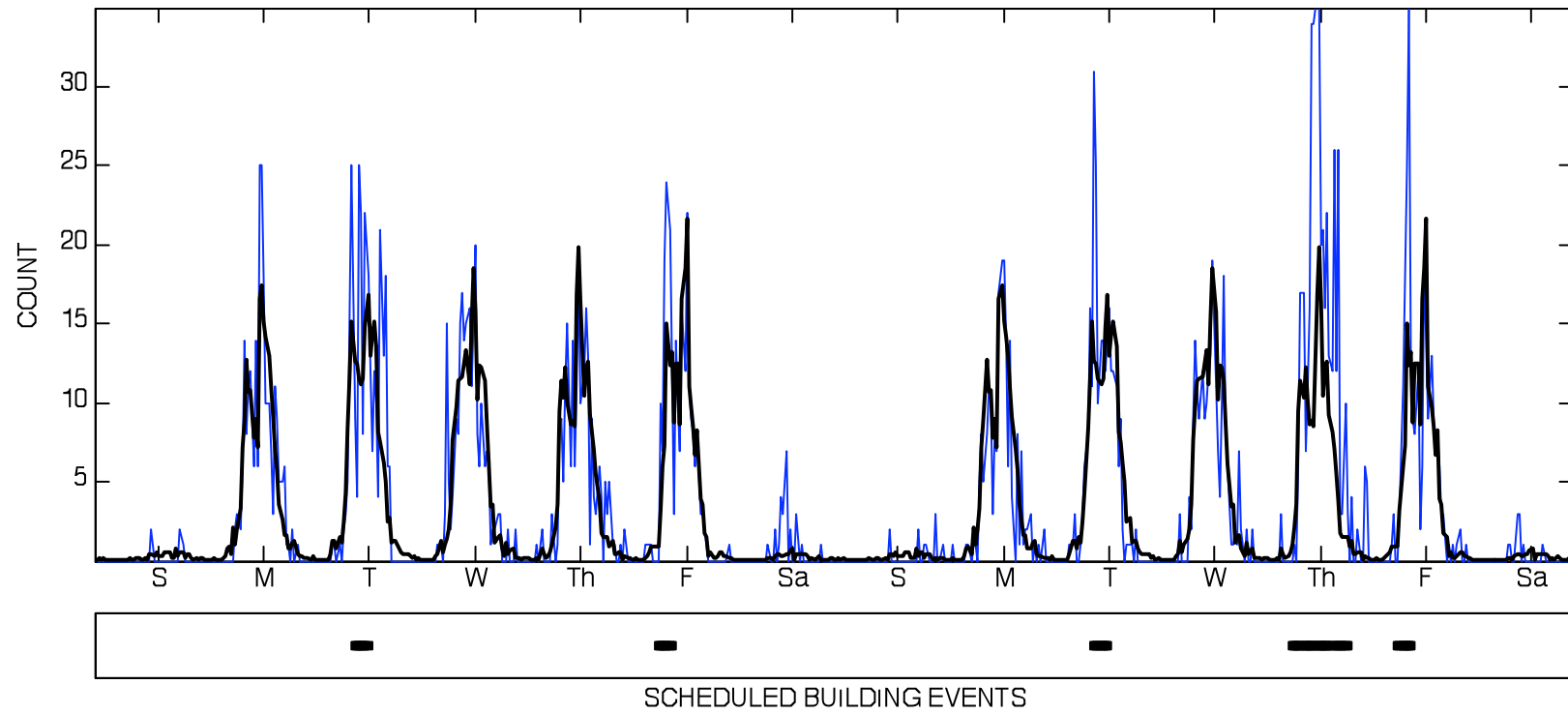


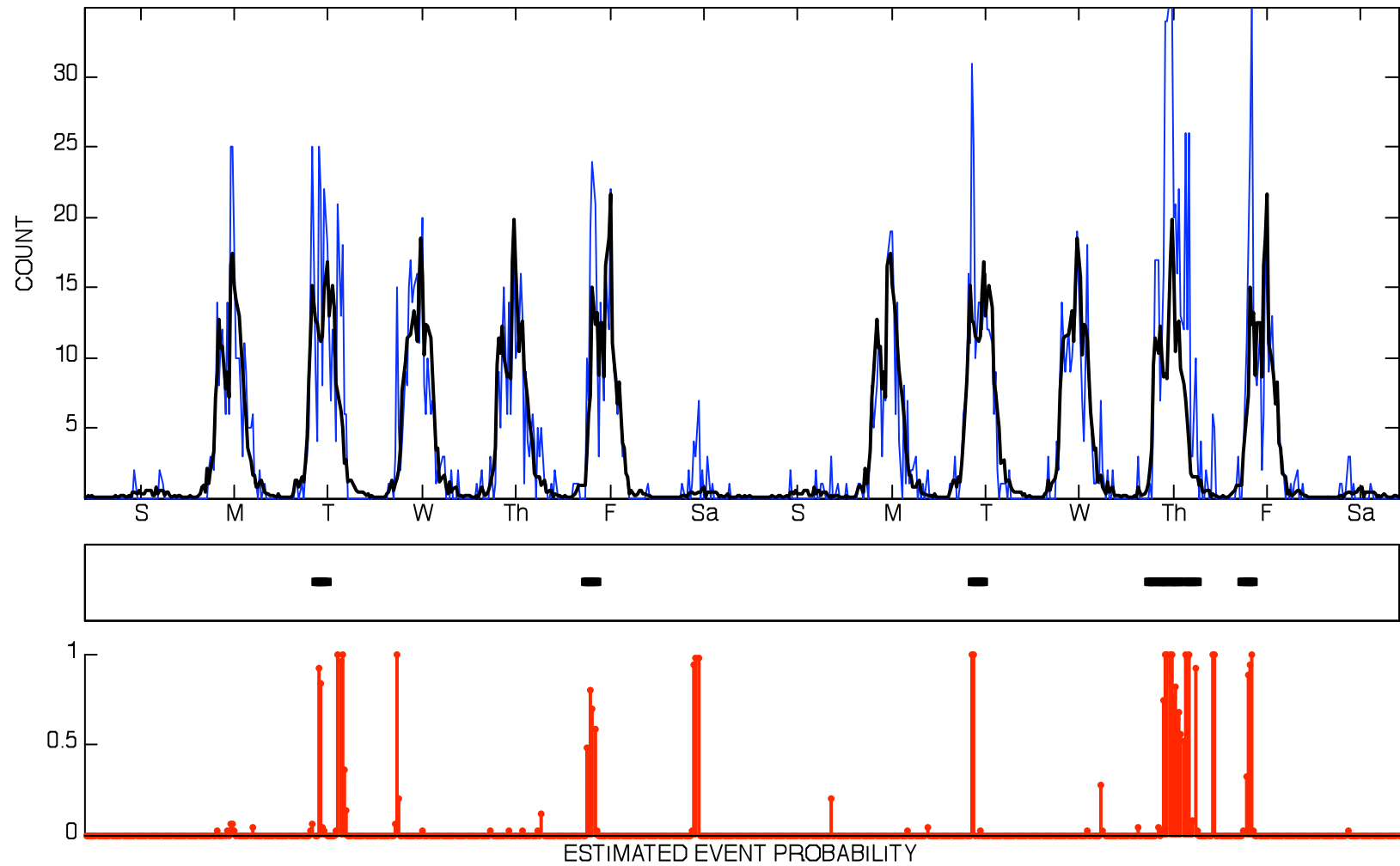
Baseline Poisson Model without Events











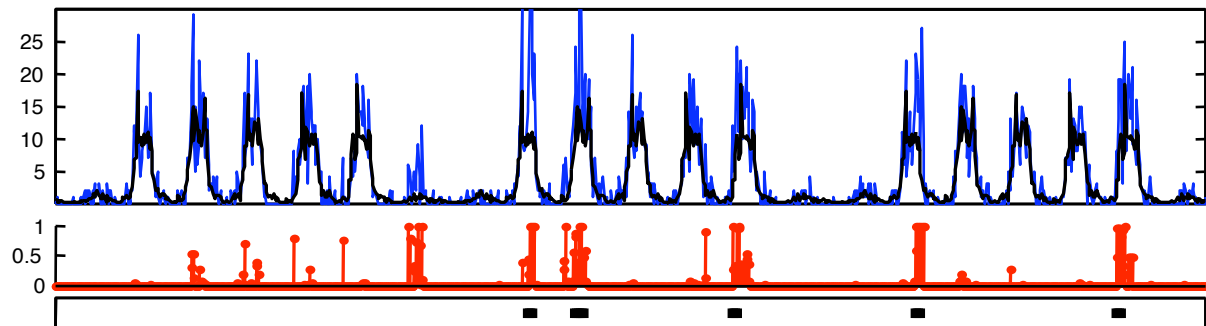
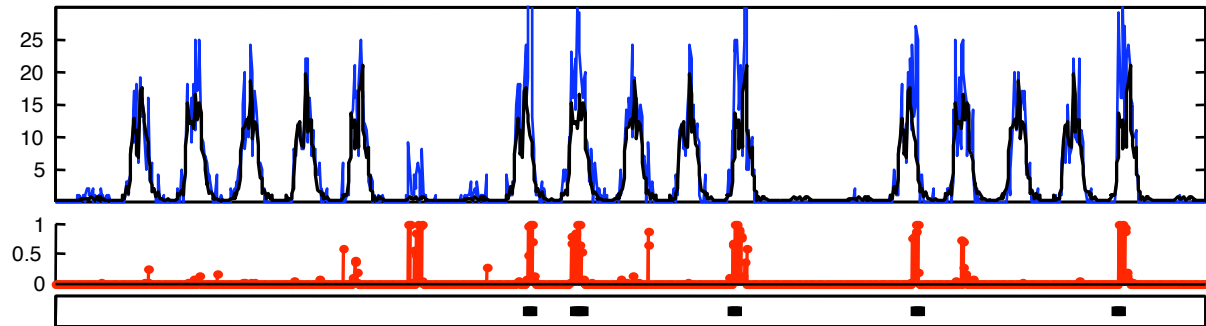
6 Weeks of Building Data

- **Actual data**

- **Mean profile $\lambda(t)$**

- **Event probability $p(z)$**

- **Actual events**



Evaluation

- Compare predictions with ground truth
 - Building data: official schedule of events
 - Freeway: times/days of Dodger baseball games
- Ground truth is partial
 - Model can detect additional “unknown” events
- Important
 - Known events were hidden from model during training, training is completely unsupervised

Fraction of known events that are detected:

Building Data

MMPP Model	Threshold Model
100.0%	89.7%
86.2%	82.8%
79.3%	65.5%

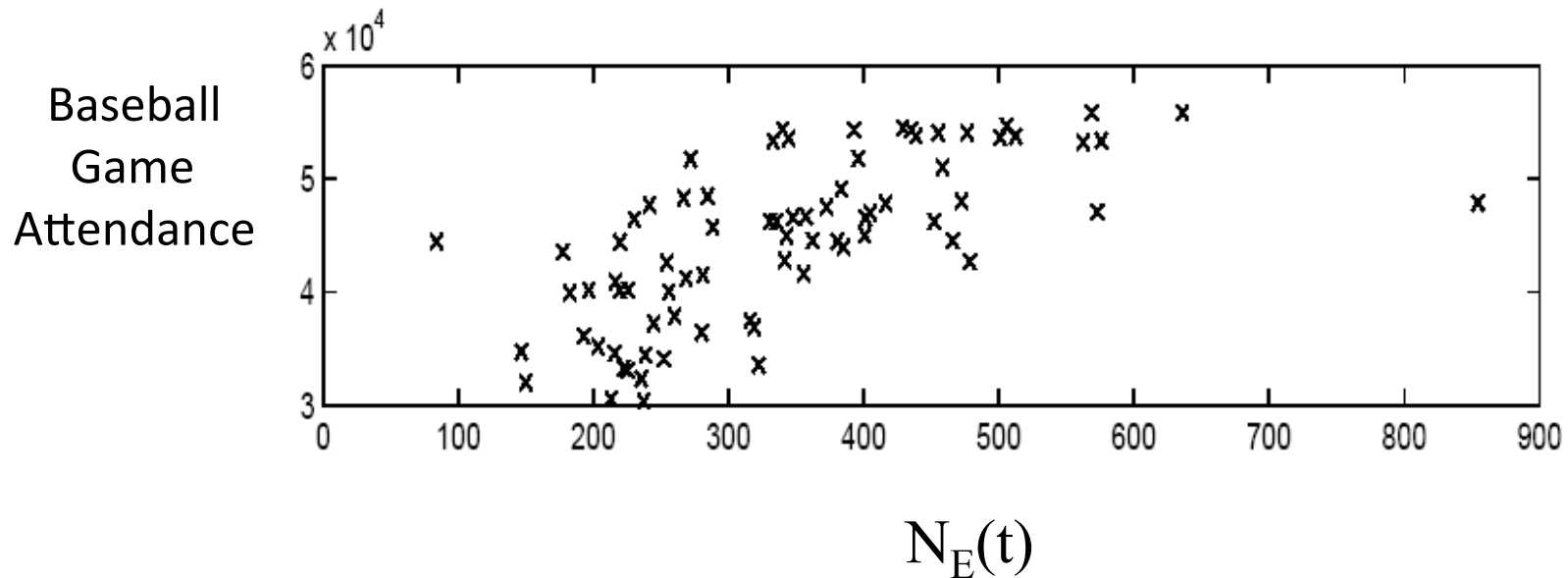
Traffic Data

MMPP Model	Threshold Model
100.0%	85.9%
100.0%	82.1%
100.0%	66.7%
97.4%	55.1%

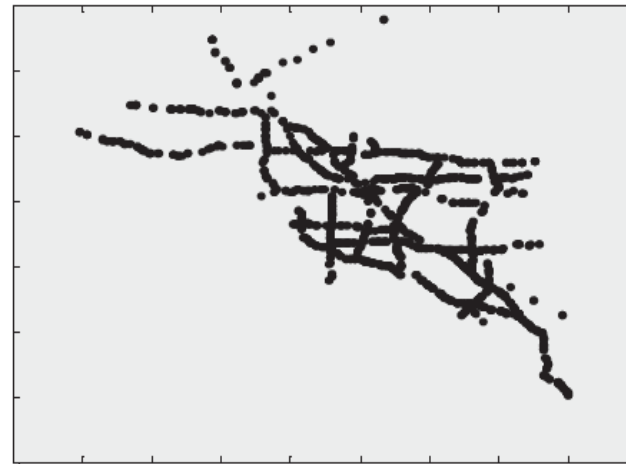
Adjust priors on transition probabilities and on thresholds to vary number of events predicted

Estimation of Event Size

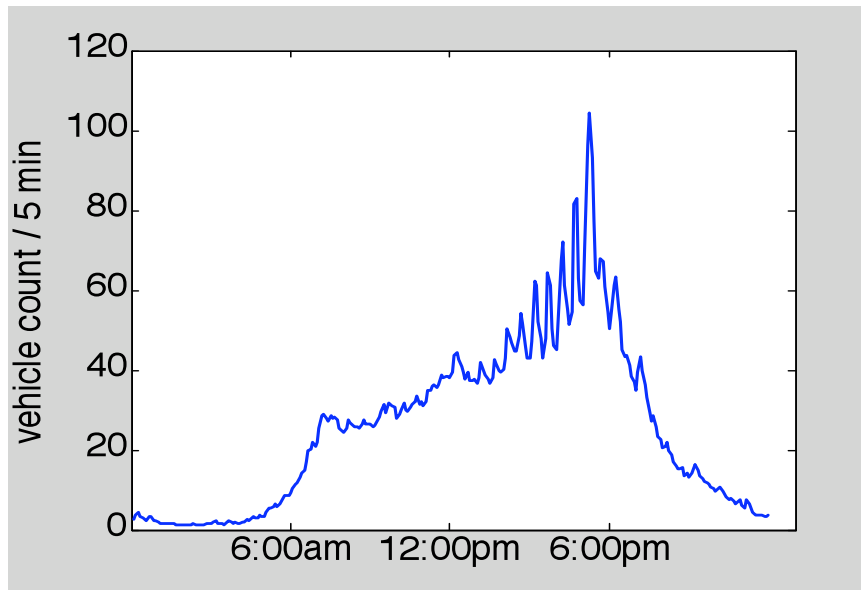
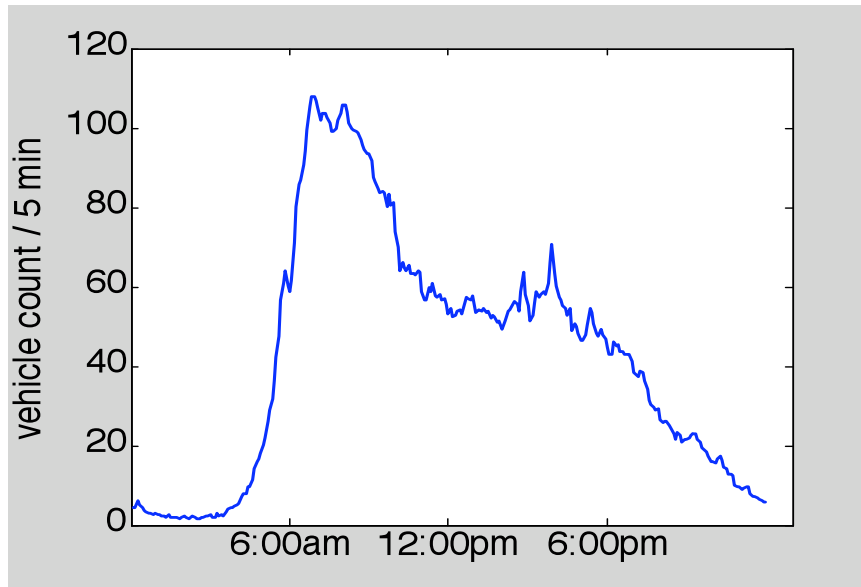
$N_E(t)$ estimates number of additional counts per event...



A Large-Scale Study

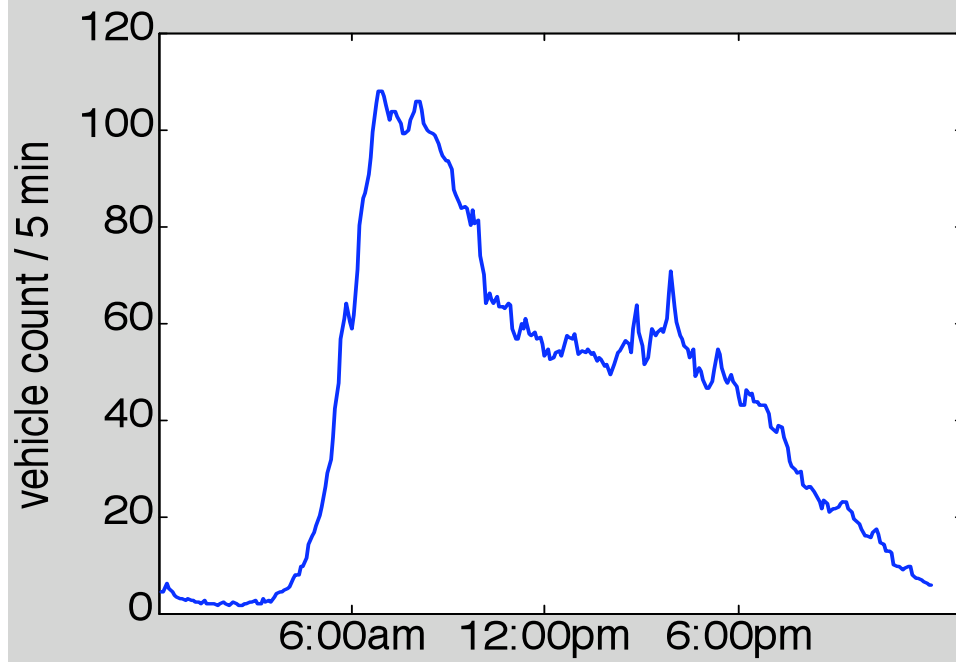


- All entrance and exit ramps in LA and OC for 7 months
- 1716 traffic sensors with data
 - 300 million measurements
 - 29% missing measurements per sensor on average

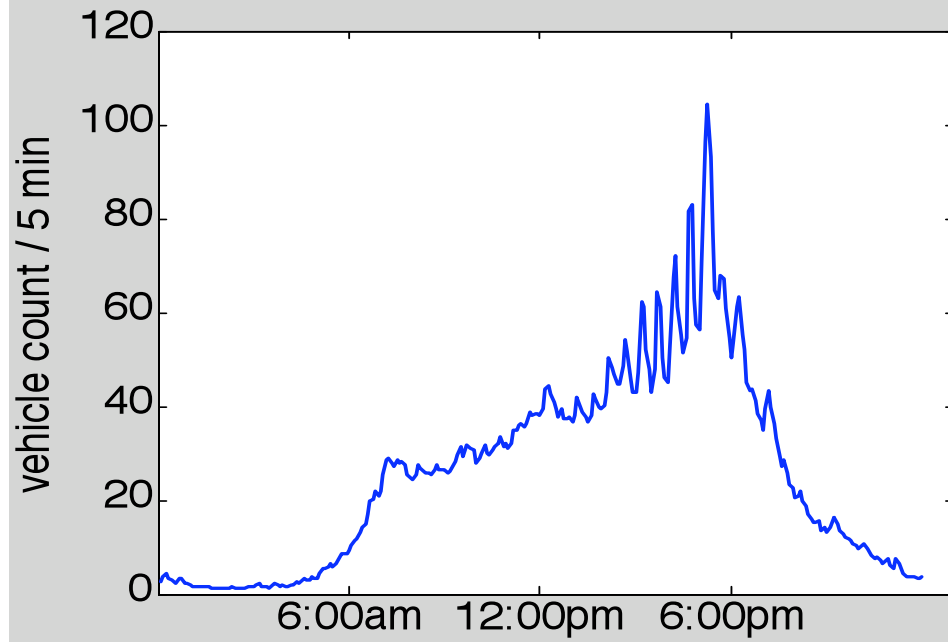


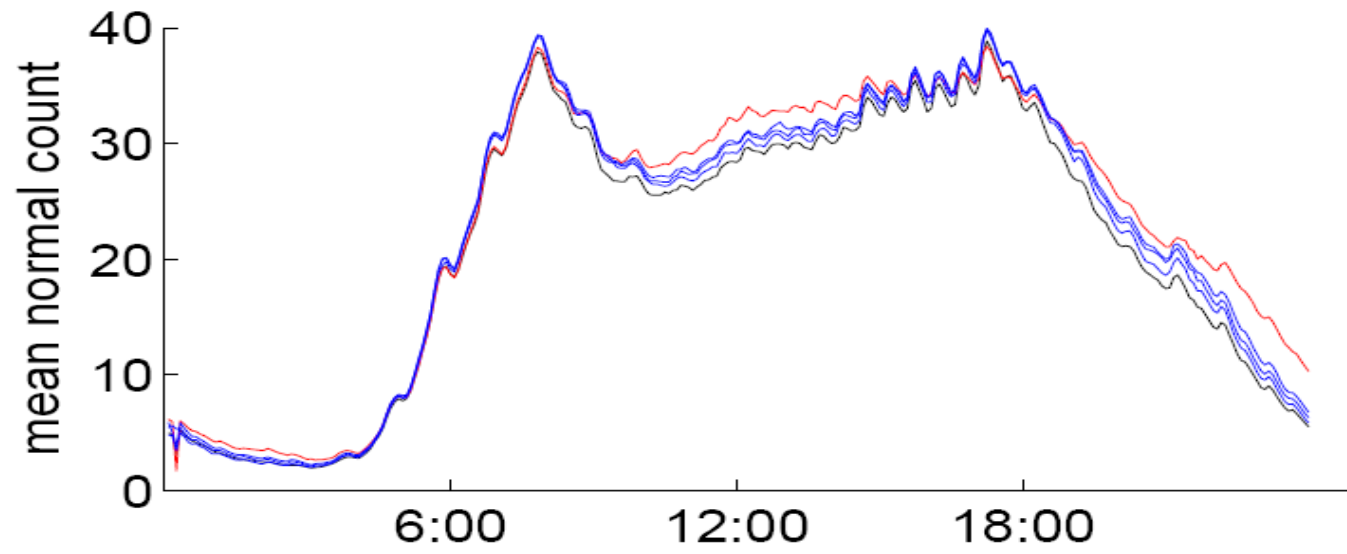
Profile Shapes (on-ramps)

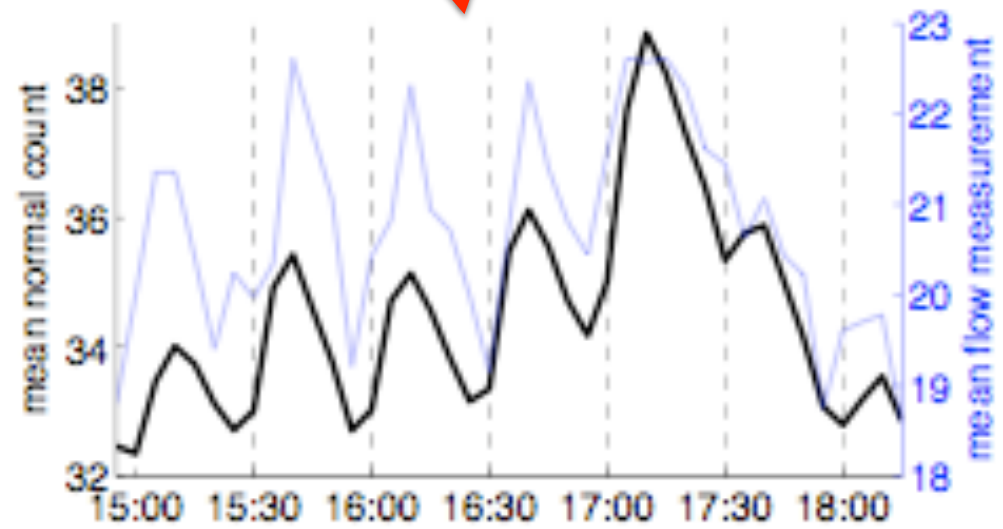
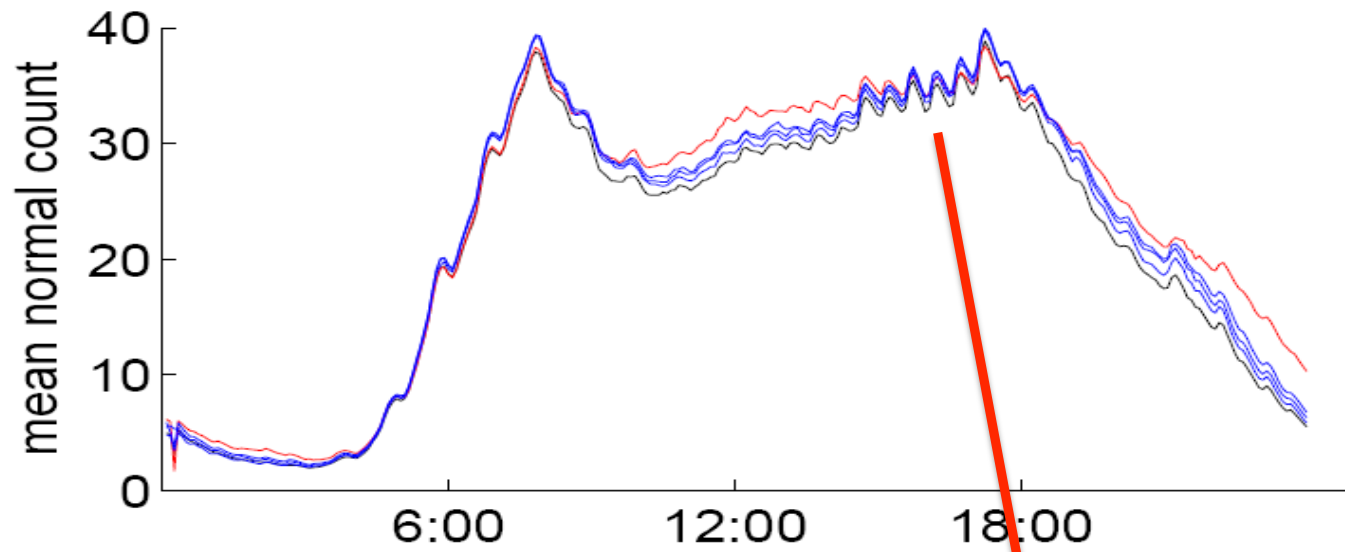
Suburban Area

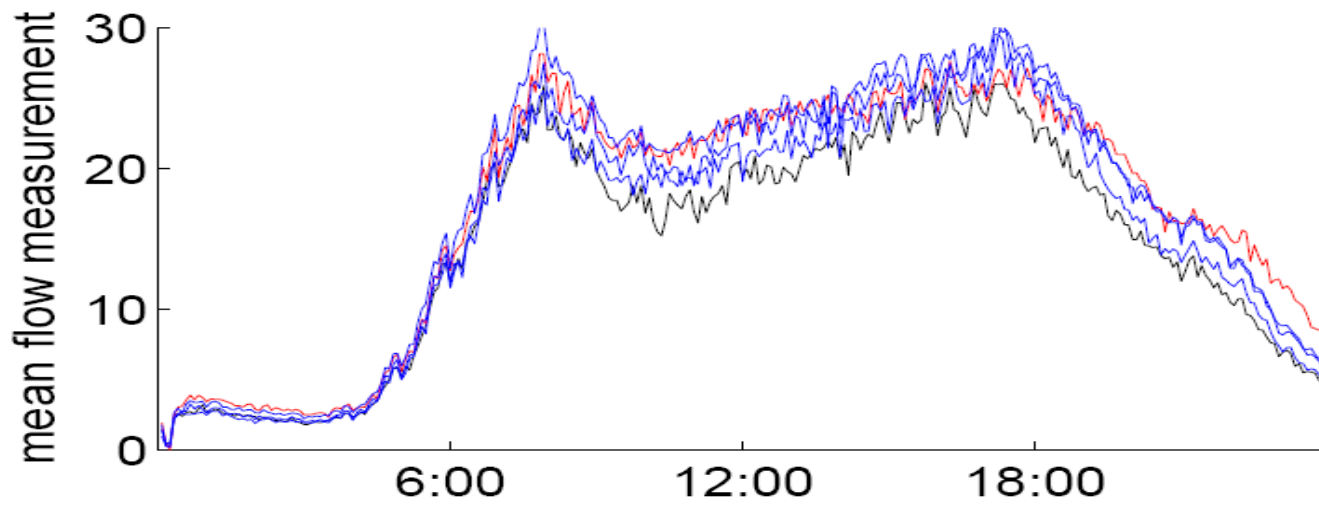
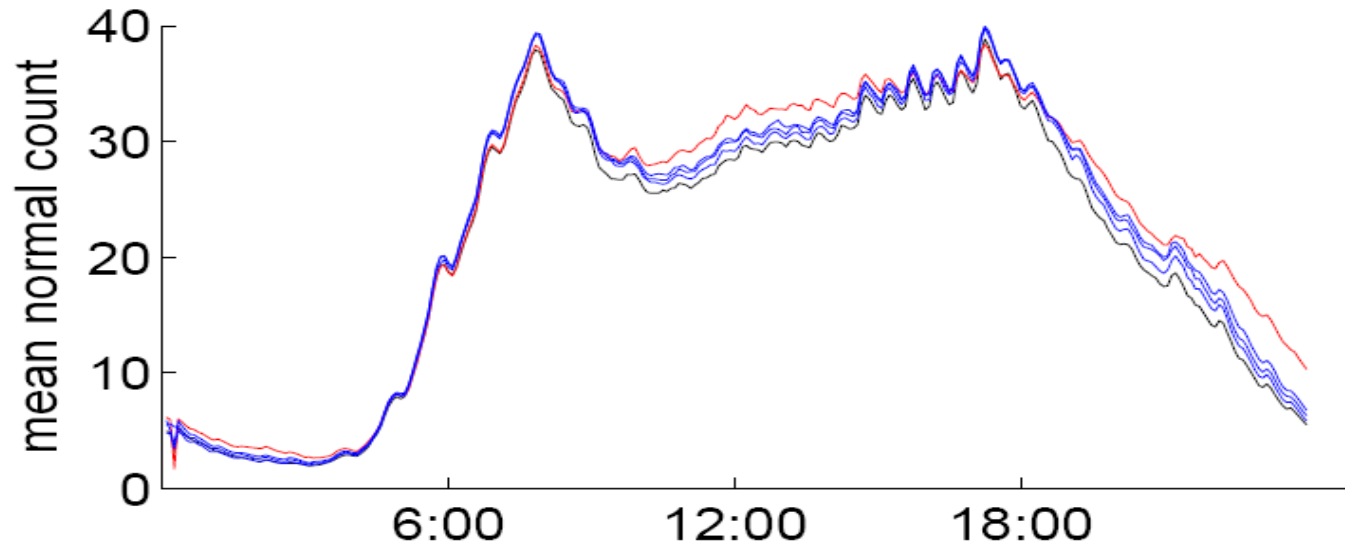


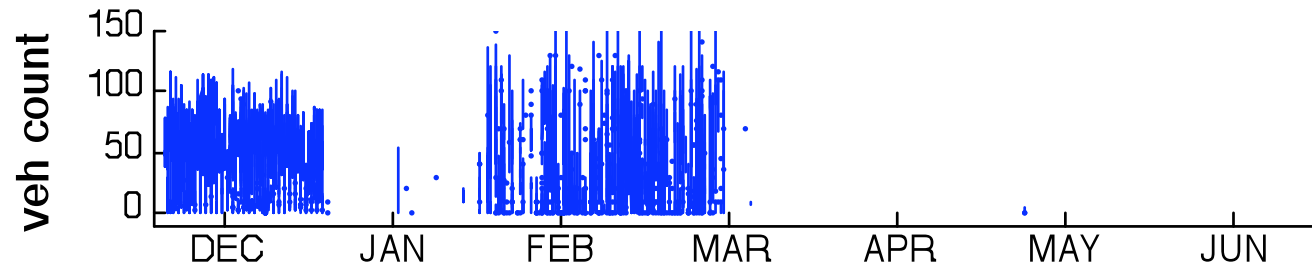
Industrial Area

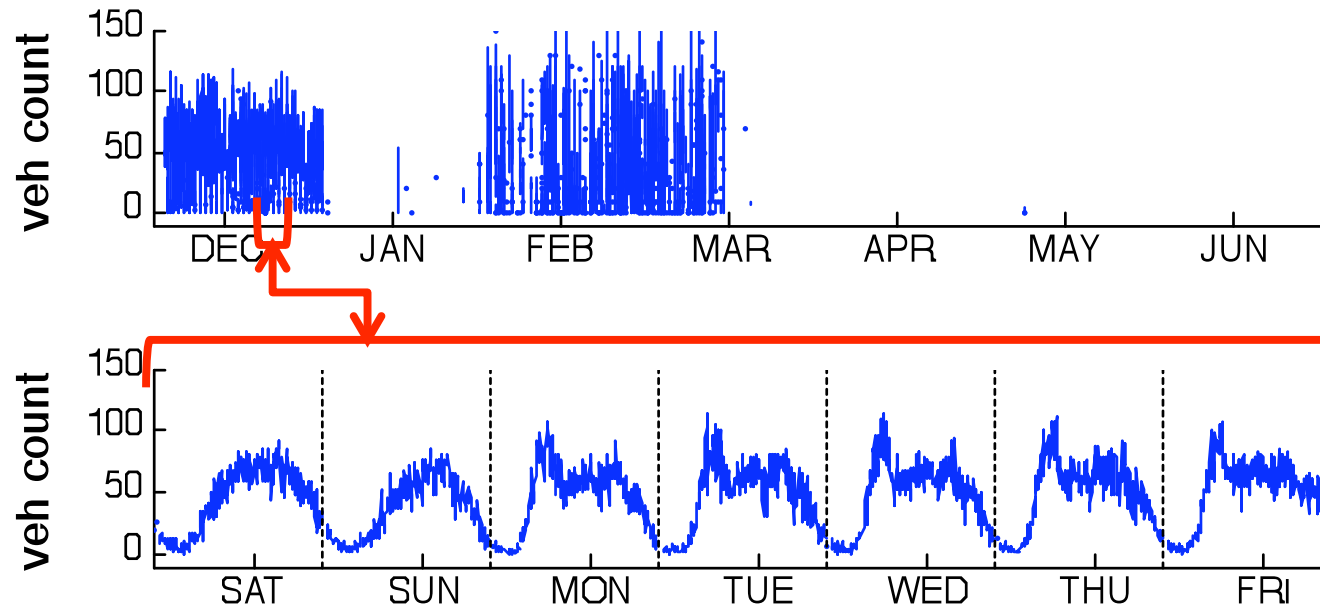


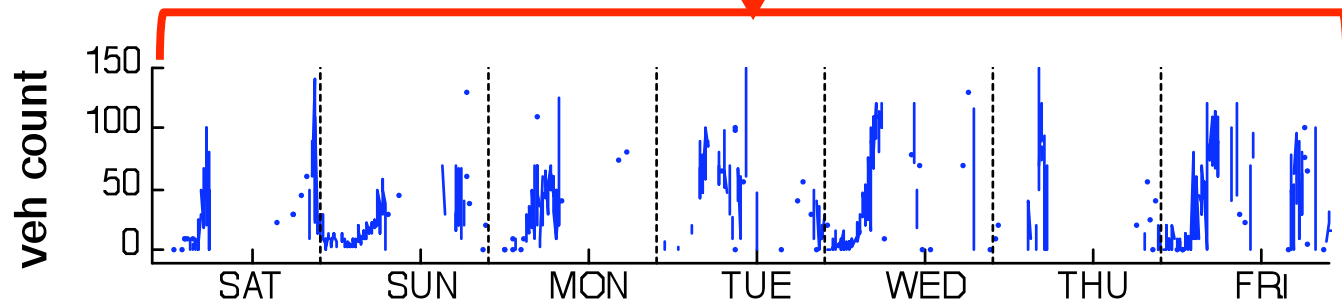
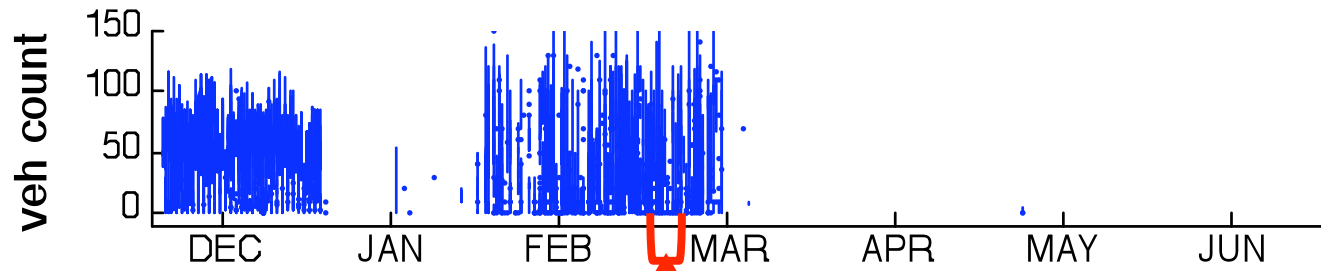


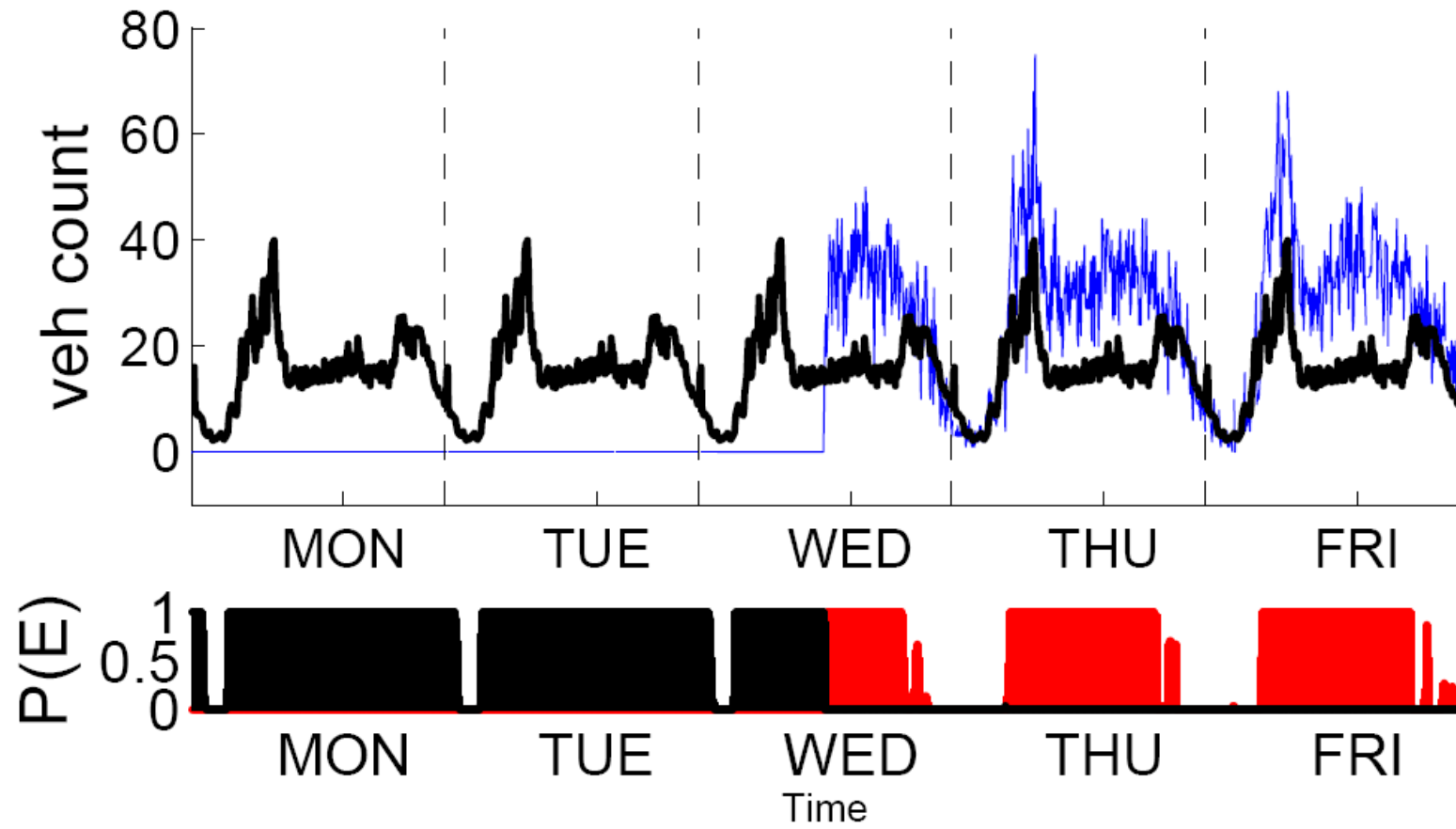


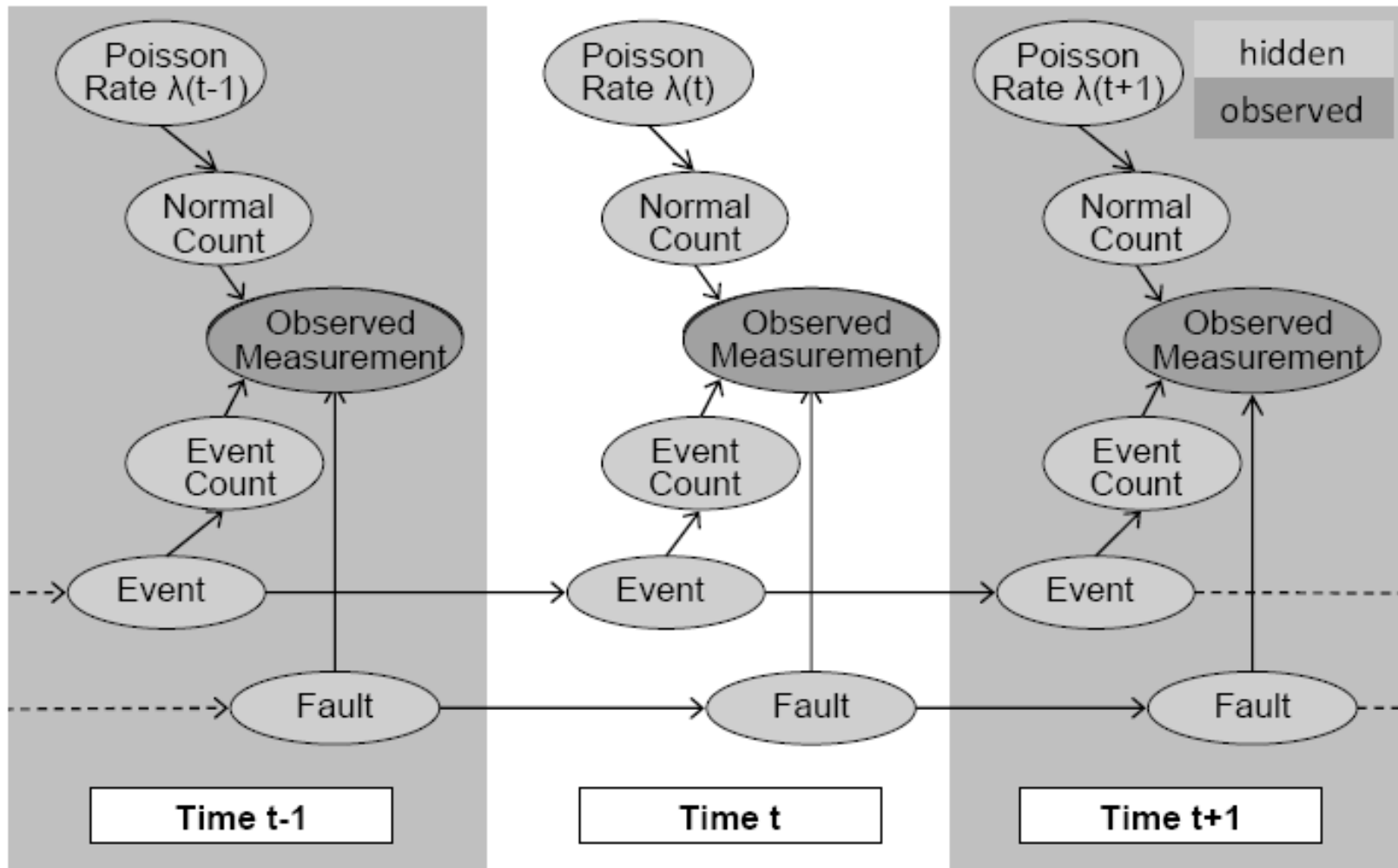


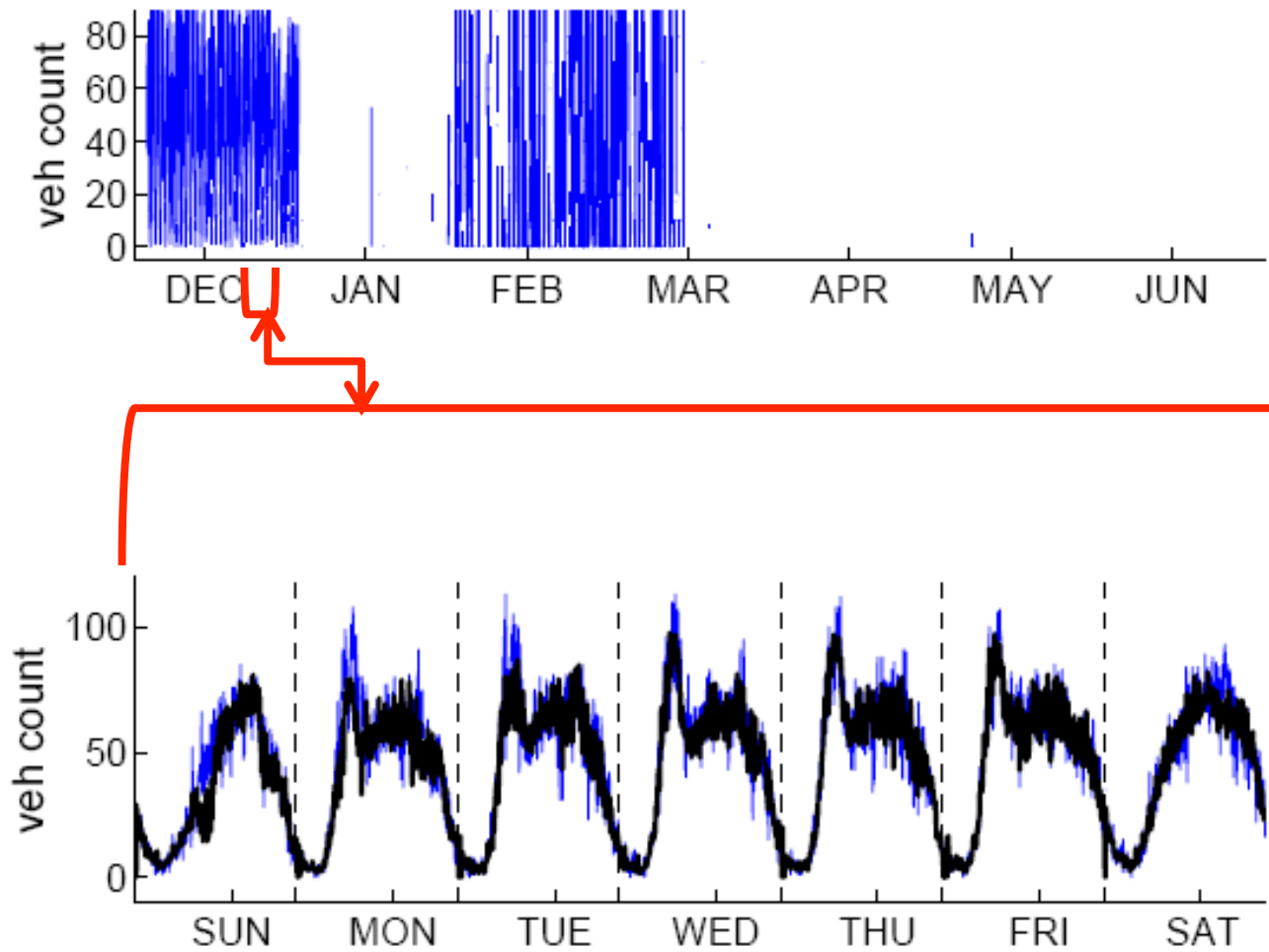








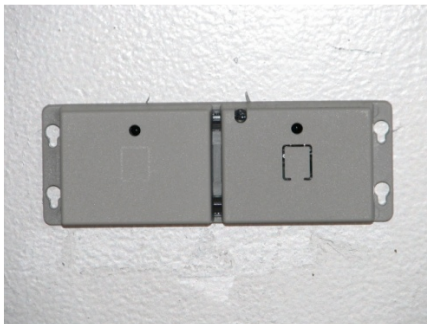
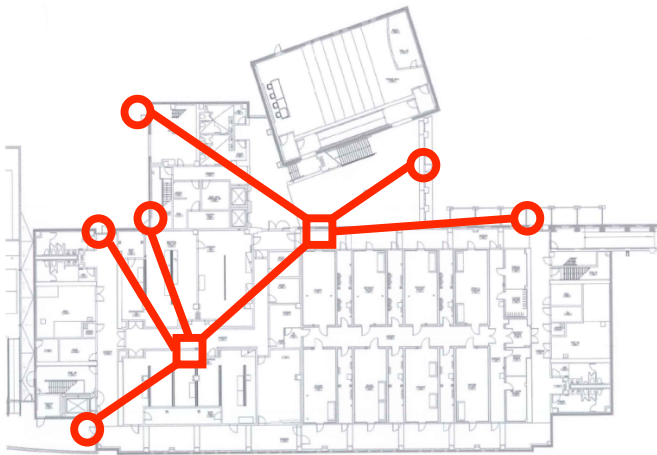




Animation showing spatio-temporal
signature of traffic event

Estimating building occupancy....

Multiple Door Sensors



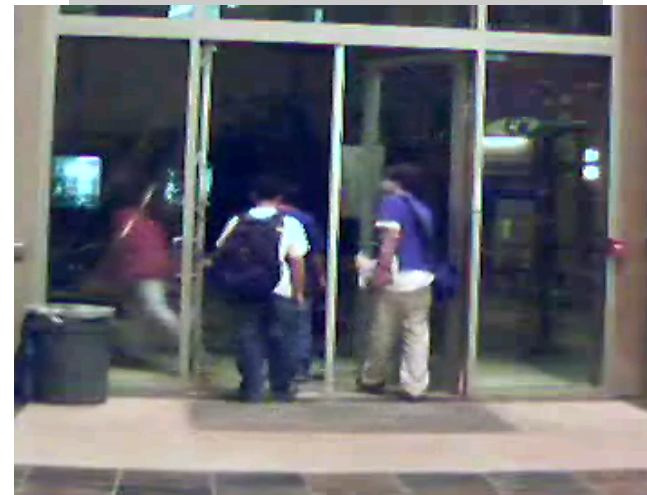
Callt2 Building, UC Irvine campus

Sensor Noise

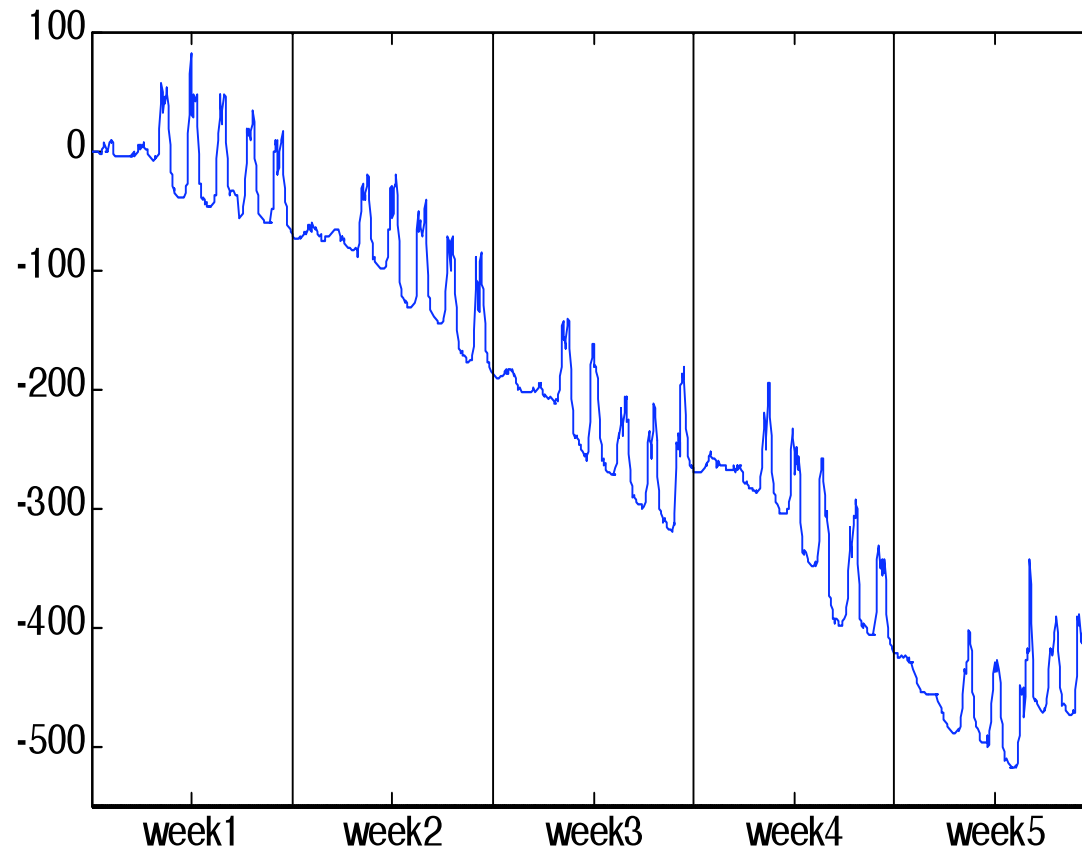
Over-counting



Under-counting

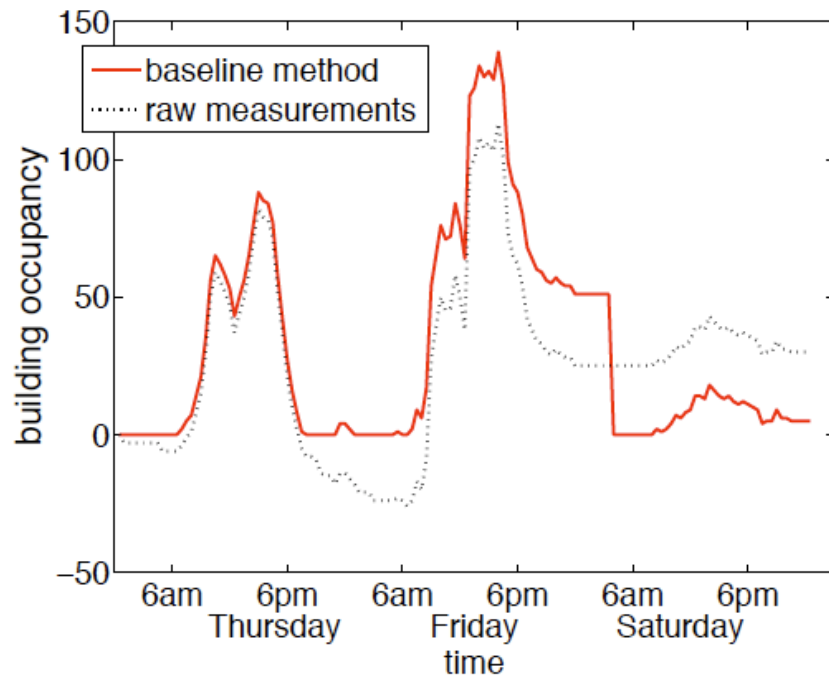


Difference of In and Out Counts

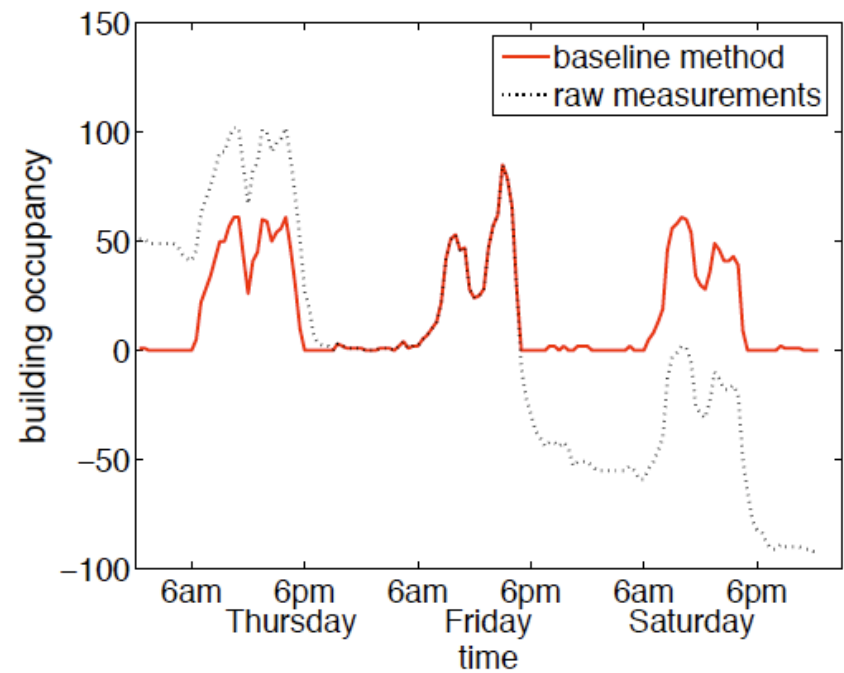


Systematic bias in raw counts

Extra In Counts



Extra Out Counts



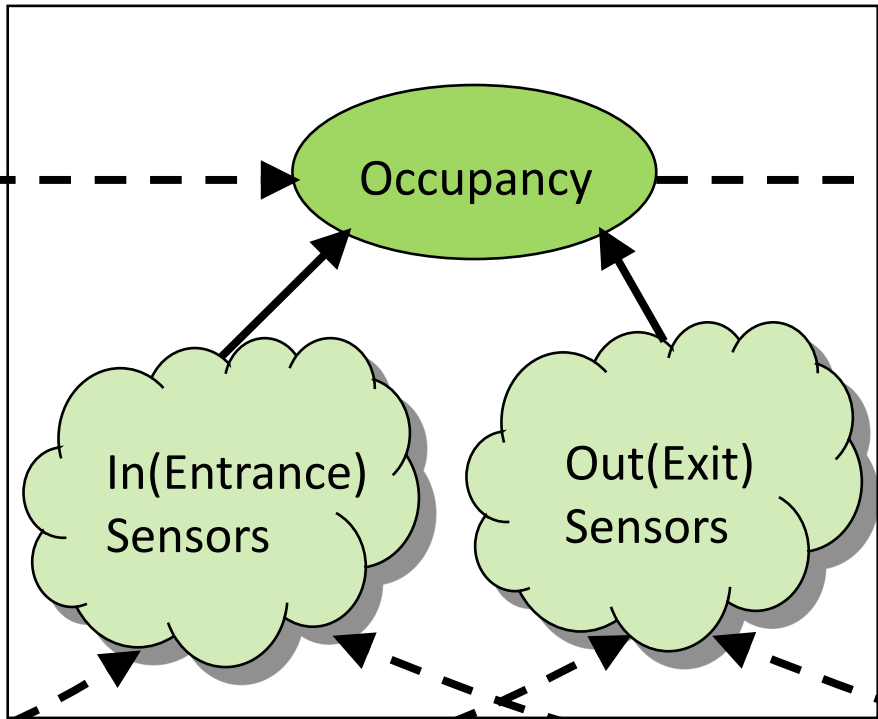
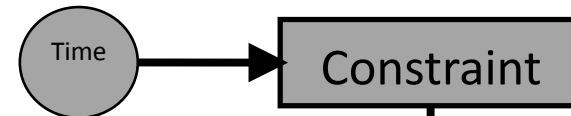
Probabilistic Model

- For each sensor:
 - Same MPMM model as before, learn normal Poisson behavior, subtracting out events
 - Add a sensor noise model:

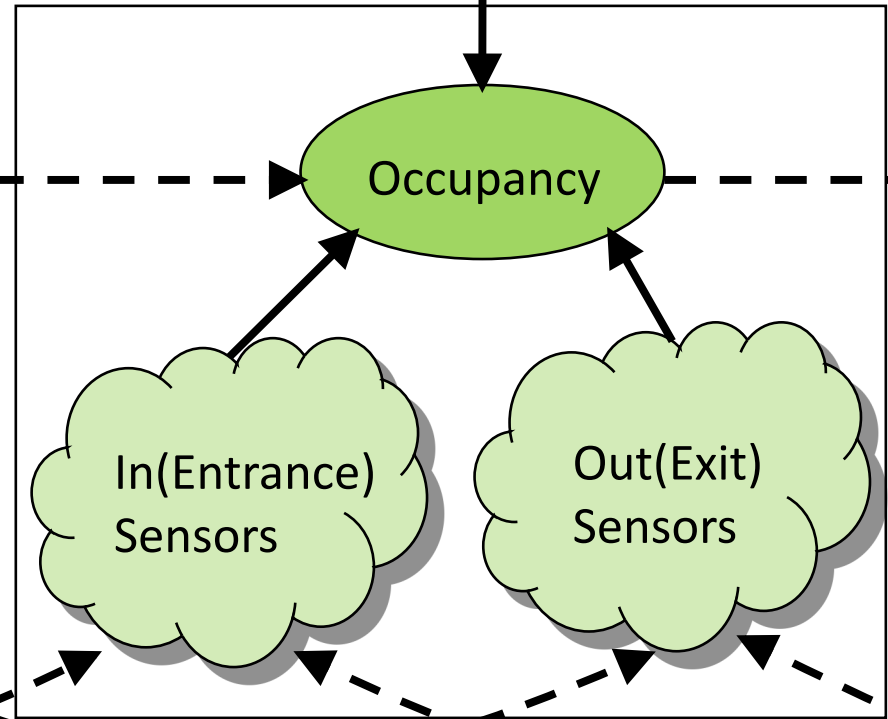
$$o_t = a_t + \Upsilon_t^O - \Upsilon_t^U$$

- Occupancy(t)
 - = sums/differences of unobserved sensor counts + Occupancy(t-1)
- Add a prior on late night occupancy
 - Geometric(0.9) prior on occupancy at 3am - mean \sim 1.1 person

Sketch of Graphical Model

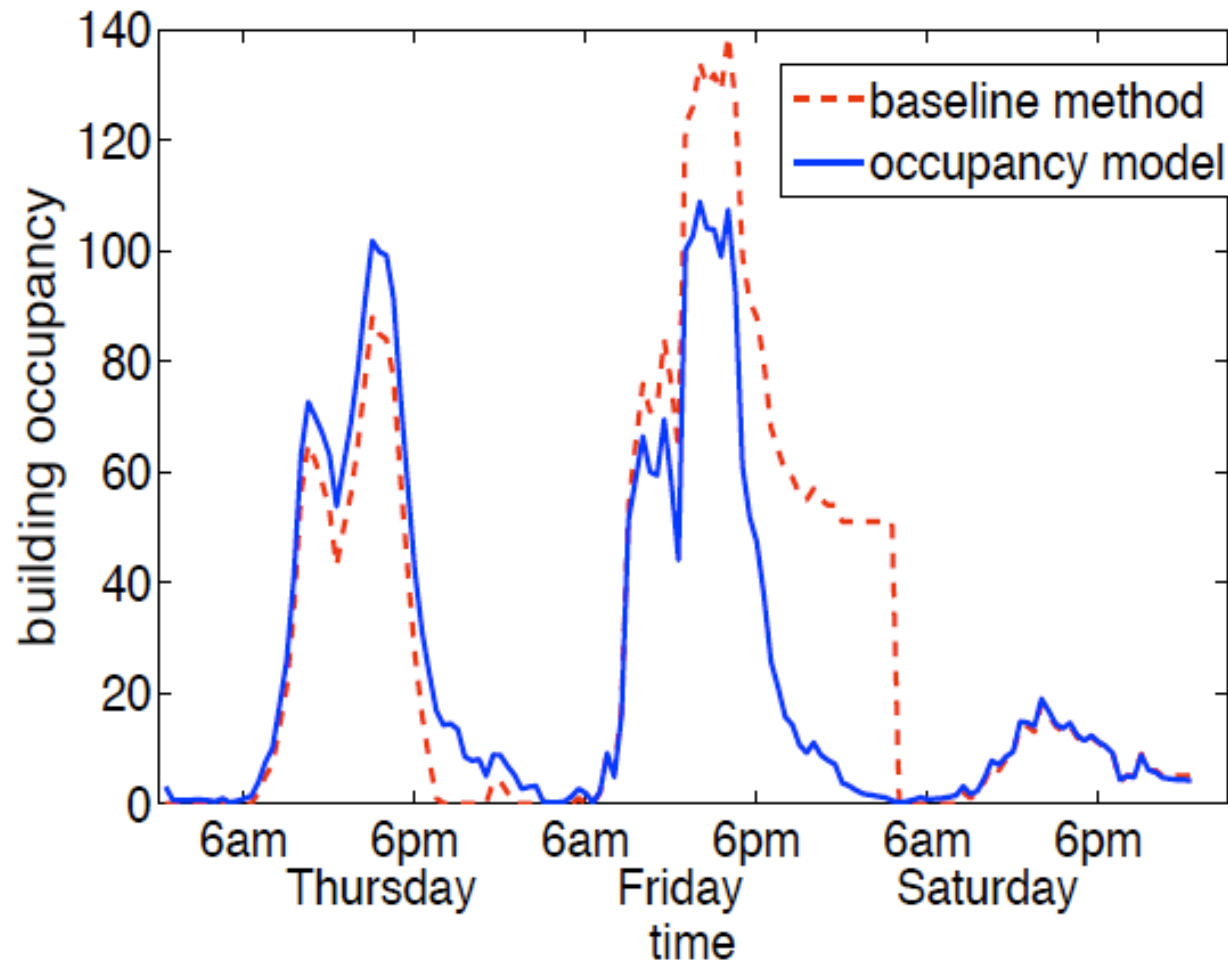


Time t

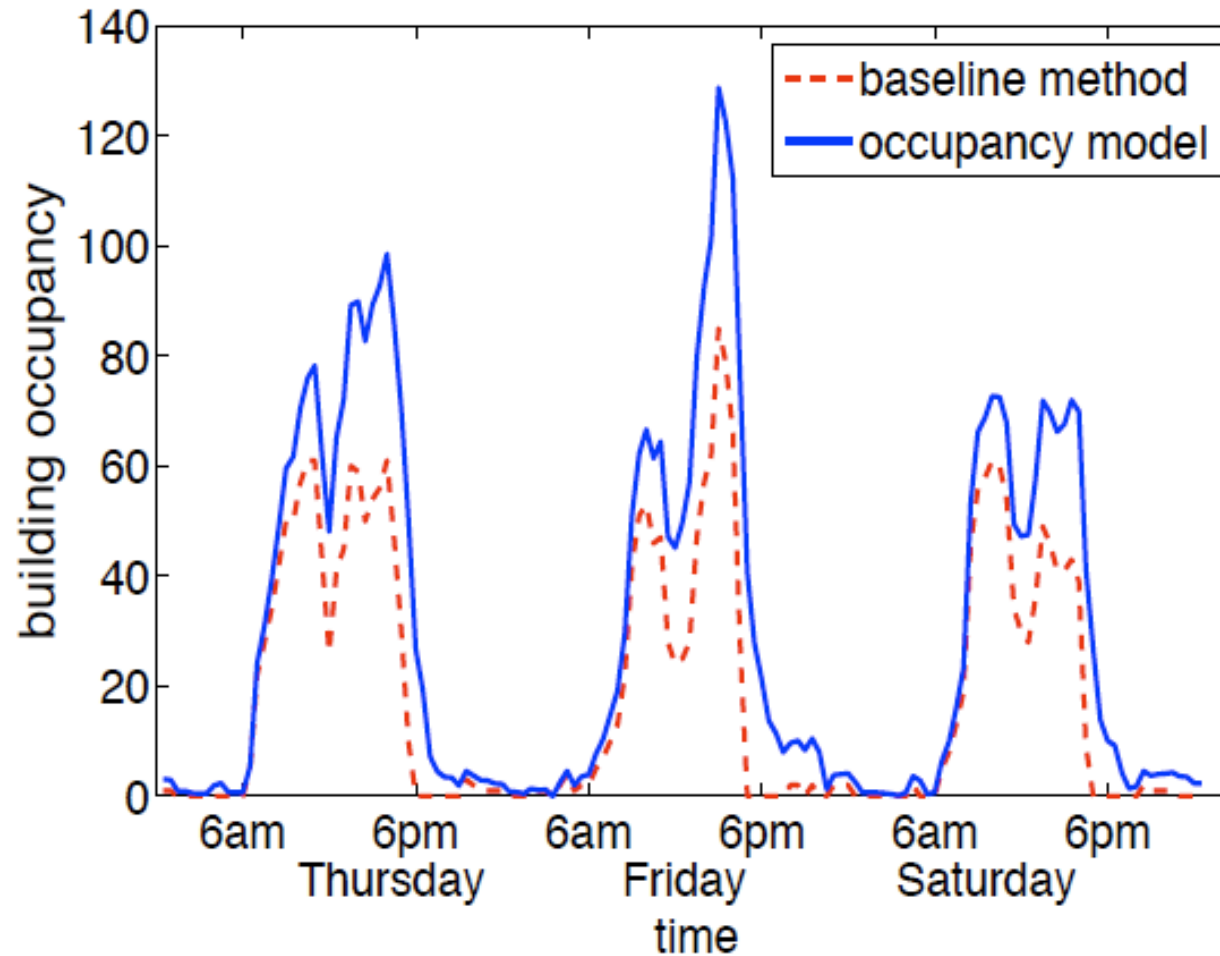


Time $t+1$

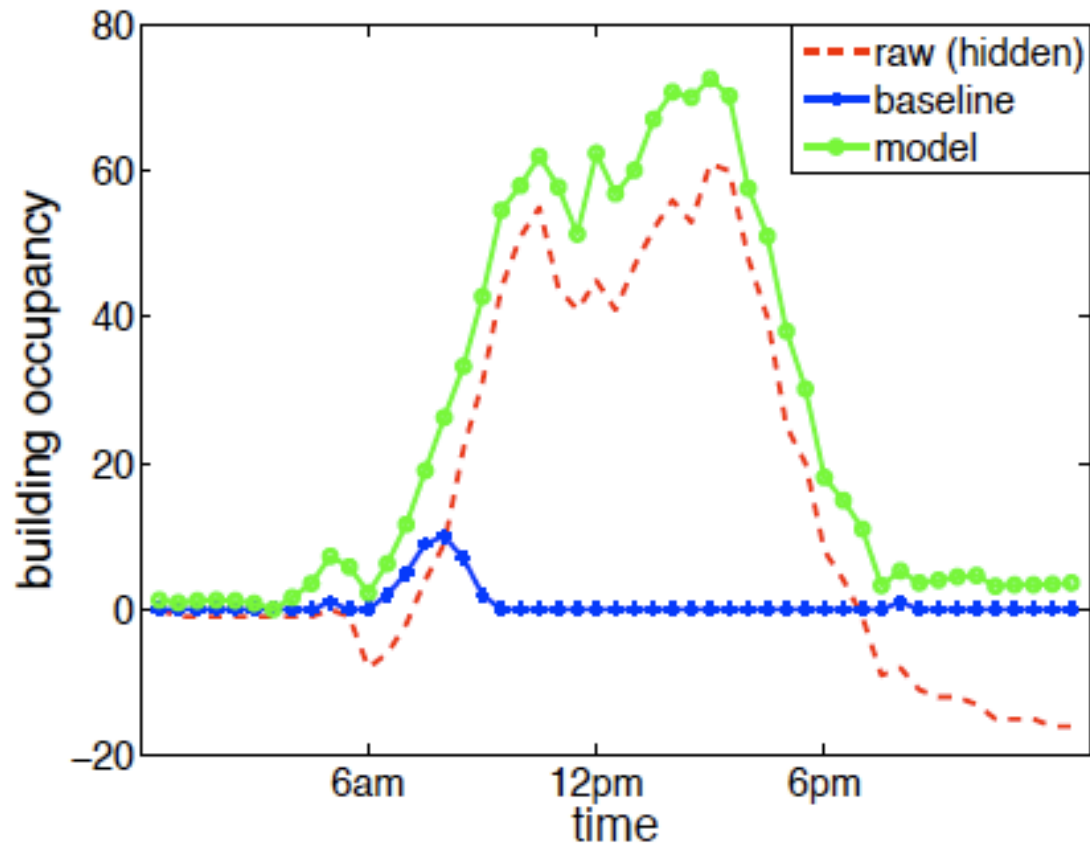
Estimation with Extra In Counts



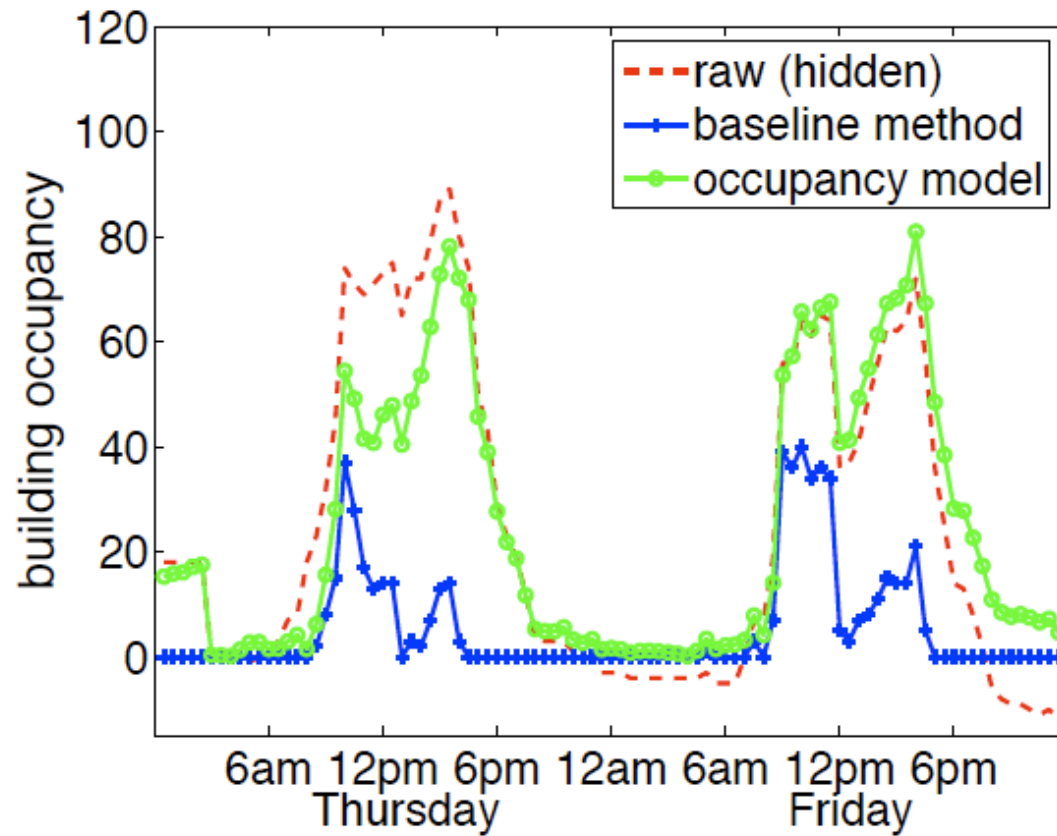
Estimation with Extra Out Counts



Missing Data



Corrupted Data Experiment



Summary

- Time-series data of human behavior
 - Increasingly available
 - Interesting research problems and applications
 - Strong temporal regularities + bursts of anomalous events
 - Hidden Markov Poisson models are useful and scalable
 - Same models can be applied to multiple data types and applications
- Ongoing and future directions
 - Accident detection and characterization
 - Spatial models and missing data, e.g., sparse Markov networks
 - Linking traffic sensors to census and image data
 - Modeling trends, changes, random effects...