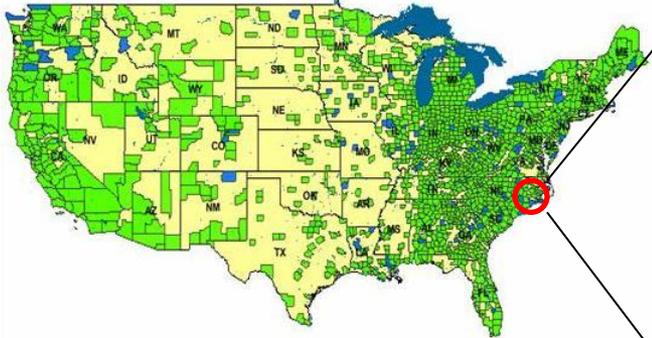


# Fast Subset Sums for Multivariate Bayesian Scan Statistics

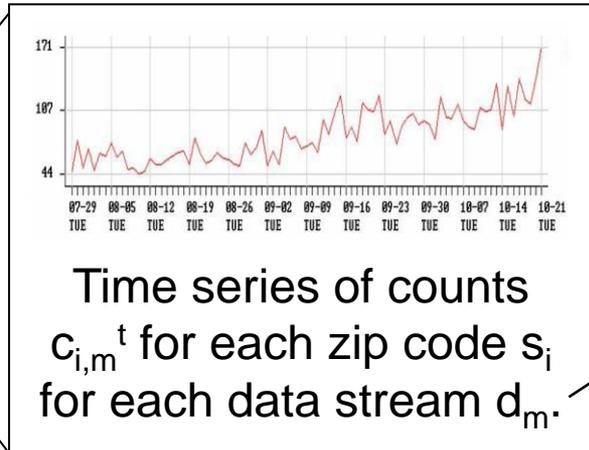
**Daniel B. Neill**  
**Carnegie Mellon University**  
**H.J. Heinz III College**  
**E-mail: [neill@cs.cmu.edu](mailto:neill@cs.cmu.edu)**

This work was partially supported by NSF grants IIS-0916345, IIS-0911032, and IIS-0325581.

# Multivariate event detection



Daily health data from thousands of hospitals and pharmacies nationwide.



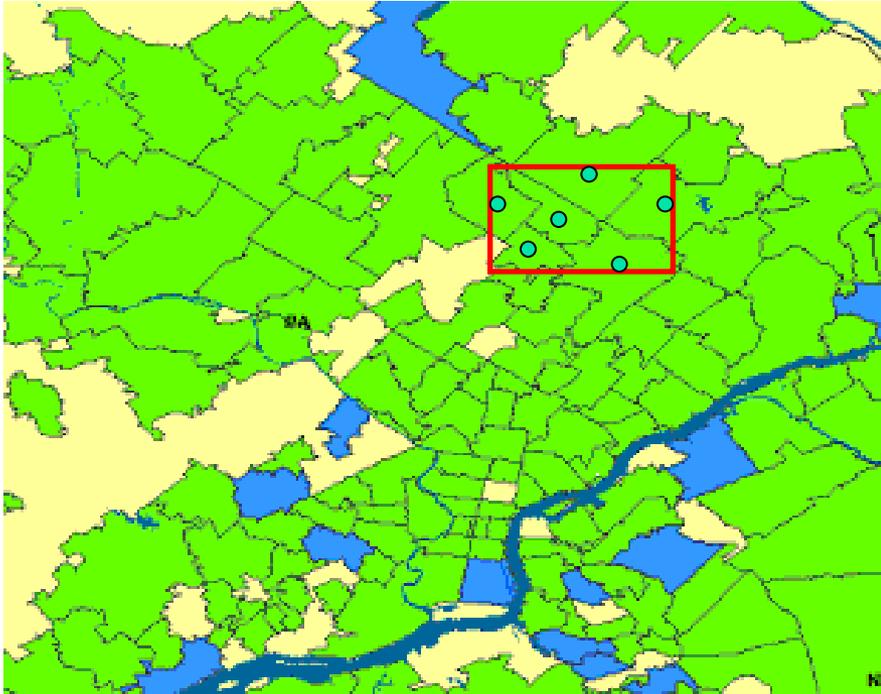
$d_1$  = respiratory ED  
 $d_2$  = constitutional ED  
 $d_3$  = OTC cough/cold  
 $d_4$  = OTC anti-fever  
etc.

Given all of this nationwide health data on a daily basis, we want to obtain a complete situational awareness by integrating information from the multiple data streams.

More precisely, we have three main goals: to detect any emerging events (i.e. outbreaks of disease), characterize the type of event, and pinpoint the affected areas.

# Expectation-based scan statistics

(Kulldorff, 1997; Neill and Moore, 2005)

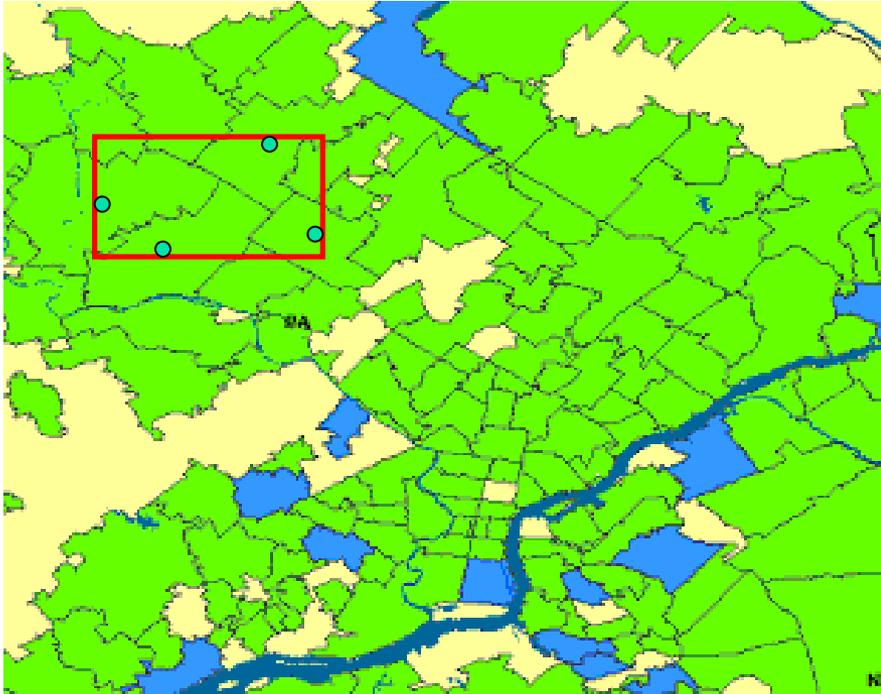


To detect and localize events, we can search for space-time regions where the number of cases is higher than expected.

Imagine moving a window around the scan area, allowing the window size, shape, and temporal duration to vary.

# Expectation-based scan statistics

(Kulldorff, 1997; Neill and Moore, 2005)

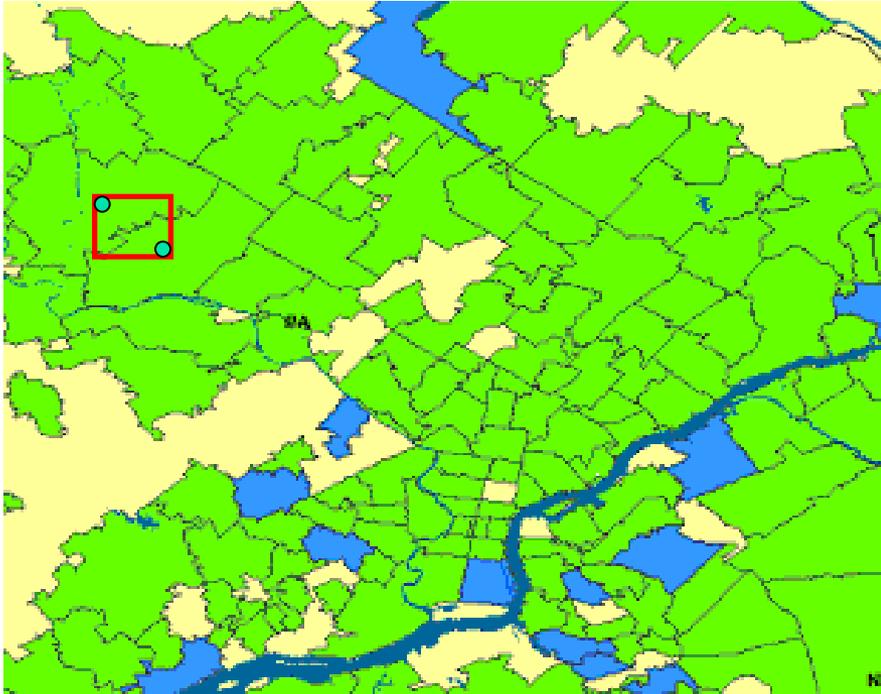


To detect and localize events, we can search for space-time regions where the number of cases is higher than expected.

Imagine moving a window around the scan area, allowing the window size, shape, and temporal duration to vary.

# Expectation-based scan statistics

(Kulldorff, 1997; Neill and Moore, 2005)

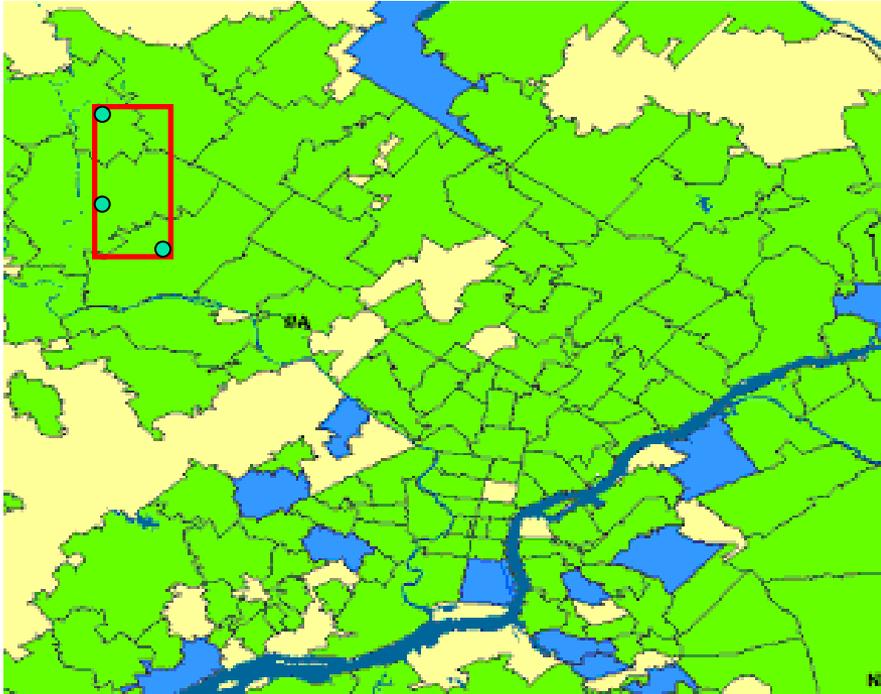


To detect and localize events, we can search for space-time regions where the number of cases is higher than expected.

Imagine moving a window around the scan area, allowing the window size, shape, and temporal duration to vary.

# Expectation-based scan statistics

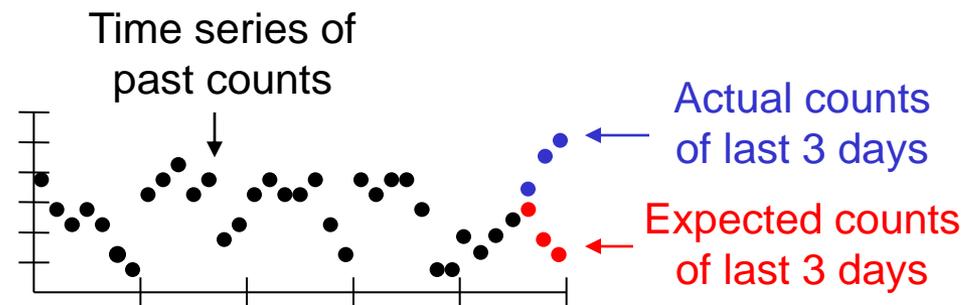
(Kulldorff, 1997; Neill and Moore, 2005)



To detect and localize events, we can search for space-time regions where the number of cases is higher than expected.

Imagine moving a window around the scan area, allowing the window size, shape, and temporal duration to vary.

For each subset of locations, we examine the aggregated time series, and compare actual to expected counts.



# Overview of the MBSS method

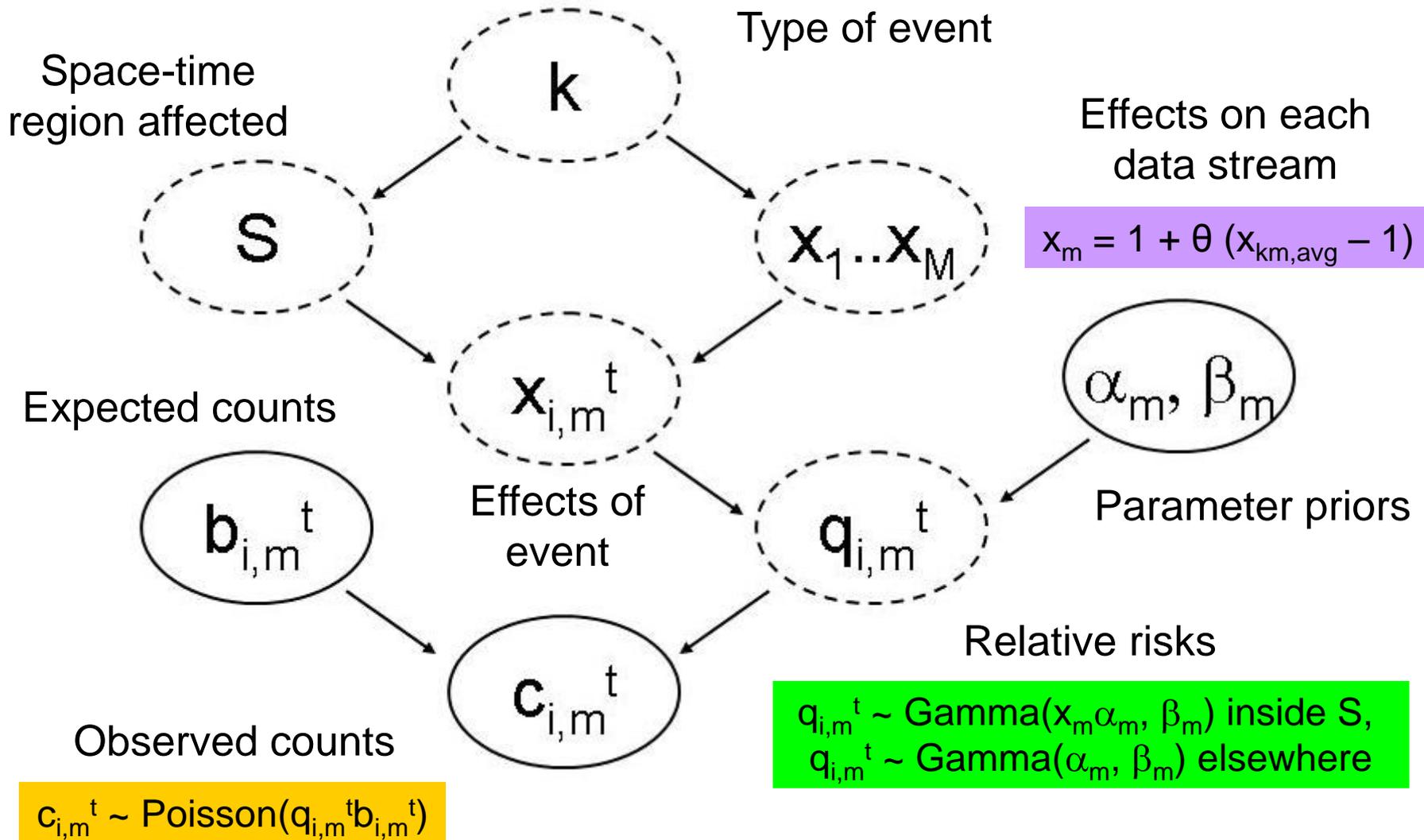


Given a set of event types  $E_k$ , a set of space-time regions  $S$ , and the multivariate dataset  $D$ , MBSS outputs the posterior probability  $\Pr(H_1(S, E_k) | D)$  of each type of event in each region, as well as the probability of no event,  $\Pr(H_0 | D)$ .

We must provide the prior probability  $\Pr(H_1(S, E_k))$  of each event type  $E_k$  in each region  $S$ , as well as the prior probability of no event,  $\Pr(H_0)$ .

MBSS uses Bayes' Theorem to combine the data likelihood given each hypothesis with the prior probability of that hypothesis:  $\Pr(H | D) = \Pr(D | H) \Pr(H) / \Pr(D)$ .

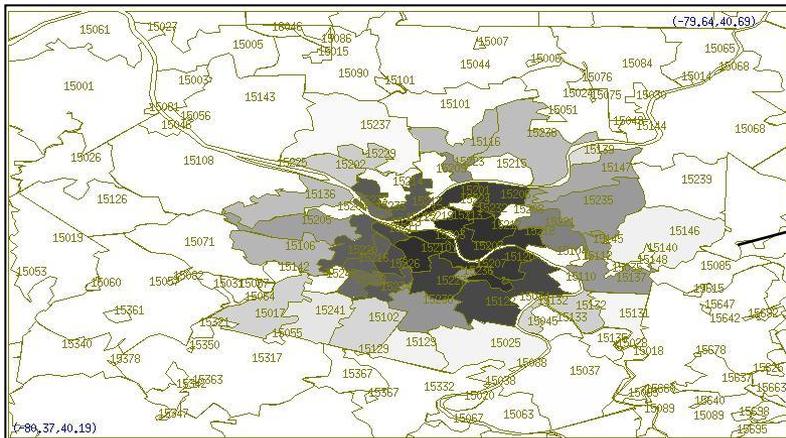
# The Bayesian hierarchical model



# Interpretation and visualization

MBSS gives the total posterior probability of each event type  $E_k$ , and the distribution of this probability over space-time regions  $S$ .

Visualization:  $\Pr(H_1(s_i, E_k)) = \sum \Pr(H_1(S, E_k))$   
for all regions  $S$  containing location  $s_i$ .



## Posterior probability map

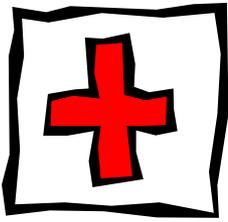
Total posterior probability of a respiratory outbreak in each Allegheny County zip code.

Darker shading = higher probability.

# MBSS: advantages and limitations

MBSS can detect faster and more accurately by integrating multiple data streams.

MBSS can model and differentiate between multiple potential causes of an event.



MBSS assumes a uniform prior for circular regions and zero prior for non-circular regions, resulting in low power for **elongated** or **irregular** clusters.

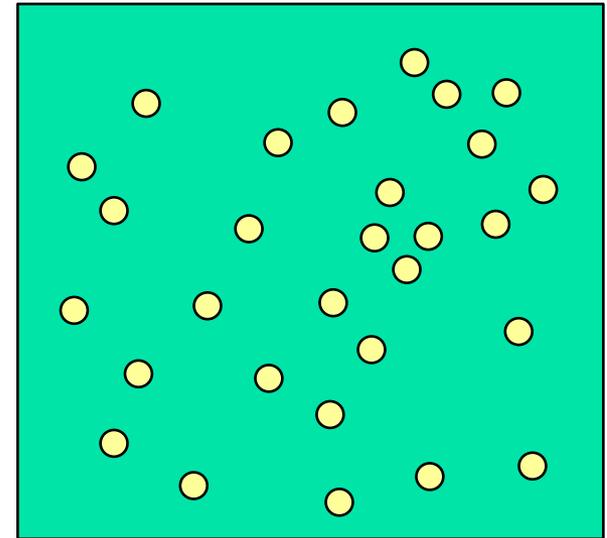
There are too many subsets of the data ( $2^N$ ) to compute likelihoods for all of them!

How can we extend MBSS to **efficiently** detect irregular clusters?

# Hierarchical prior distribution

We define a non-uniform prior  $\Pr(H_1(S, E_k))$  over all  $2^N$  subsets of the data.

This prior has hierarchical structure:

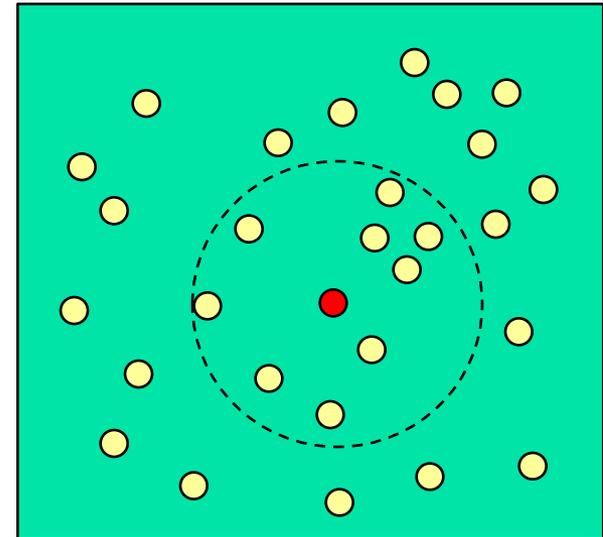


# Hierarchical prior distribution

We define a non-uniform prior  $\Pr(H_1(S, E_k))$  over all  $2^N$  subsets of the data.

This prior has hierarchical structure:

1. Choose the **center location**  $\mathbf{s}_c$  uniformly at random from  $\{s_1 \dots s_N\}$ .
2. Choose the **neighborhood size**  $n$  uniformly at random from  $\{1 \dots n_{\max}\}$ .

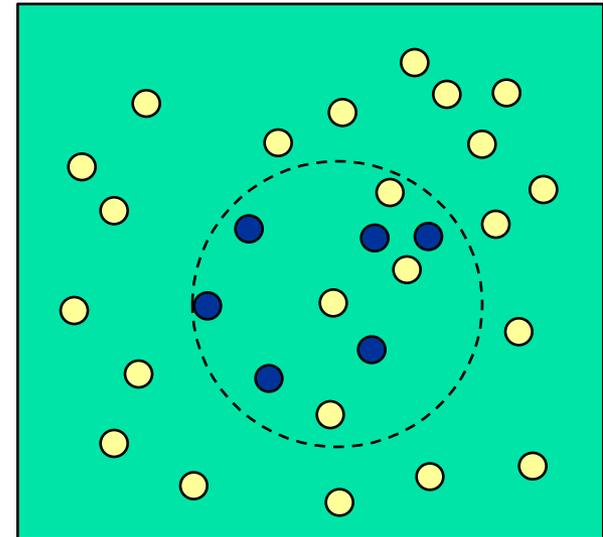


# Hierarchical prior distribution

We define a non-uniform prior  $\Pr(H_1(S, E_k))$  over all  $2^N$  subsets of the data.

This prior has hierarchical structure:

1. Choose the **center location**  $\mathbf{s}_c$  uniformly at random from  $\{s_1 \dots s_N\}$ .
2. Choose the **neighborhood size**  $n$  uniformly at random from  $\{1 \dots n_{\max}\}$ .
3. Choose **region**  $\mathbf{S}$  uniformly at random from the  $2^n$  subsets of  $S_{cn} = \{s_c \text{ and its } n - 1 \text{ nearest neighbors}\}$ .



This prior distribution has non-zero prior probabilities for any given subset  $S$ , but more compact clusters have larger priors.

# Fast Subset Sums (FSS)

Naïve computation of posterior probabilities using this prior requires summing over an exponential number of regions, which is infeasible.

However, the total posterior probability of an outbreak,  $\Pr(H_1(E_k) \mid D)$ , and the posterior probability map,  $\Pr(H_1(s_i, E_k) \mid D)$ , can be calculated efficiently **without** computing the probability of each region  $S$ .

In the original MBSS method, the **likelihood ratio** of spatial region  $S$  for a given event type  $E_k$  and event severity  $\theta$  can be found by multiplying the likelihood ratios  $LR(s_i \mid E_k, \theta)$  for all locations  $s_i$  in  $S$ .

In FSS, the **average likelihood ratio** of the  $2^n$  subsets for a given center  $s_c$  and neighborhood size  $n$  can be found by multiplying the quantities  $((1 + LR(s_i \mid E_k, \theta)) / 2)$  for all locations  $s_i$  in  $S$ .

Since the prior is uniform for a given center and neighborhood, we can compute the posteriors for each  $s_c$  and  $n$ , and marginalize over them.

# Evaluation

- We injected simulated disease outbreaks into two streams of Emergency Department data (cough, nausea) from 97 Allegheny County zip codes.
- Results were computed for ten different outbreak shapes, including compact, elongated, and irregularly-shaped, with 200 injects of each type.
- We compared FSS to the original MBSS method (searching over circles) in terms of run time, timeliness of detection, proportion of outbreaks detected, and spatial accuracy.

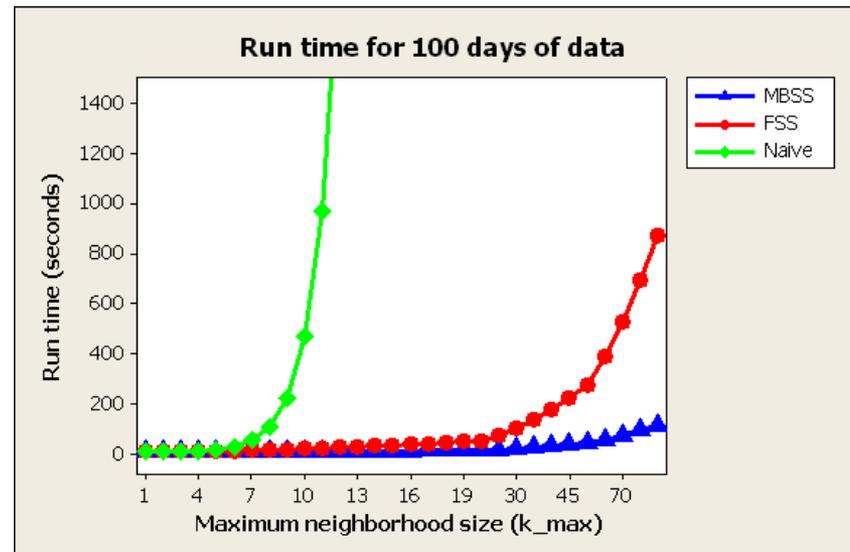
# Computation time

We compared the run times of MBSS, FSS, and a naïve subset sums implementation as a function of the maximum neighborhood size  $n_{\max}$ .

Run time of MBSS increased gradually with increasing  $n_{\max}$ , up to 1.2 seconds per day of data.

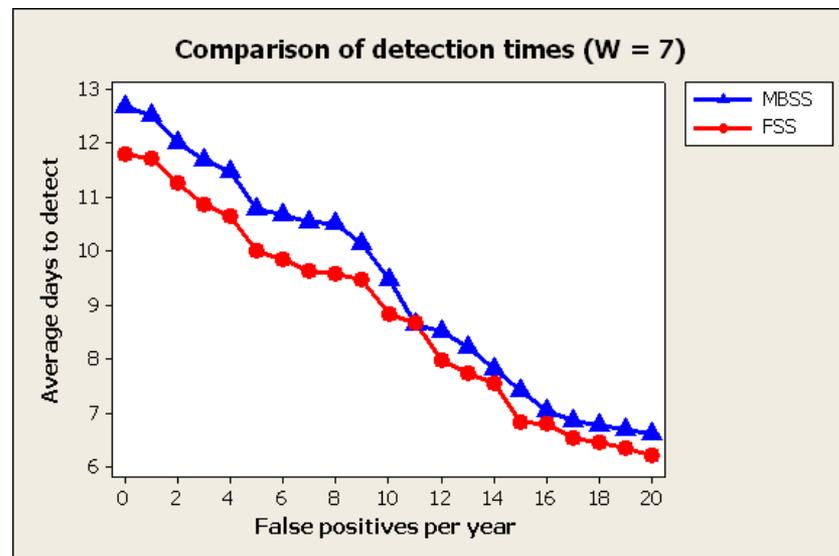
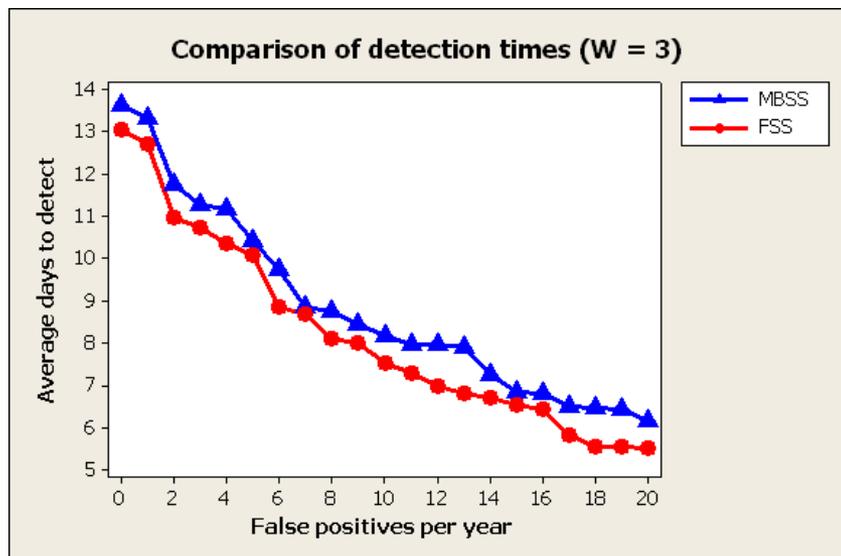
Run time of Naïve Subset Sums increased exponentially, making it infeasible for  $n_{\max} \geq 25$ .

Run time of FSS scaled quadratically with  $n_{\max}$ , up to 8.8 seconds per day of data.



Thus, while FSS is approximately 7.5x slower than the original MBSS method, it is still extremely fast, computing the posterior probability map for each day of data in under nine seconds.

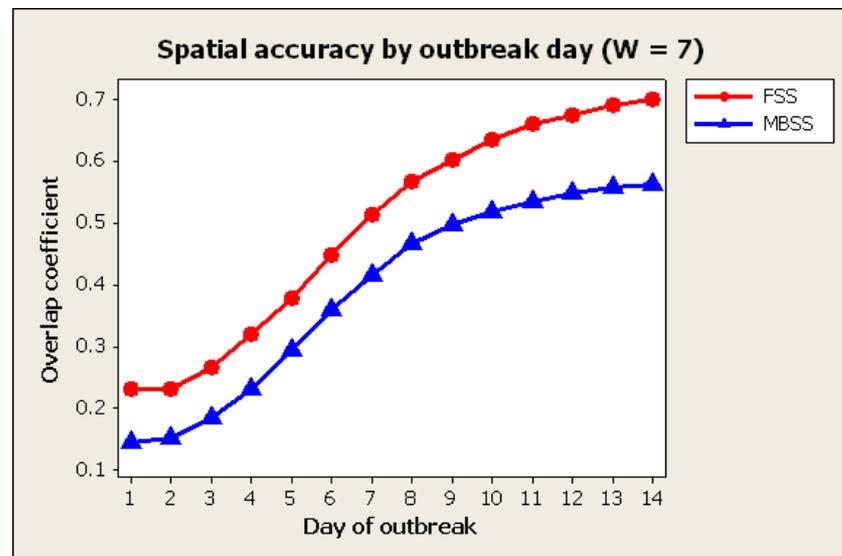
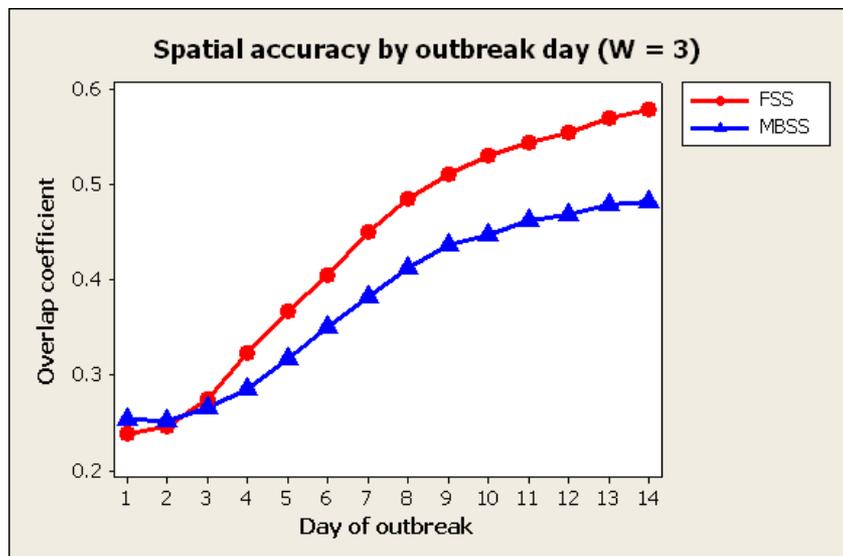
# Timeliness of detection



FSS detected an average of **one day earlier** than MBSS for maximum temporal window  $W = 3$ , and **0.54 days earlier** for  $W = 7$ , with less than half as many missed outbreaks.

Both methods achieve similar detection times for compact outbreak regions. For highly elongated outbreaks, FSS detects 1.3 to 2.2 days earlier, and for irregular regions, FSS detects 0.3 to 1.2 days earlier.

# Spatial accuracy



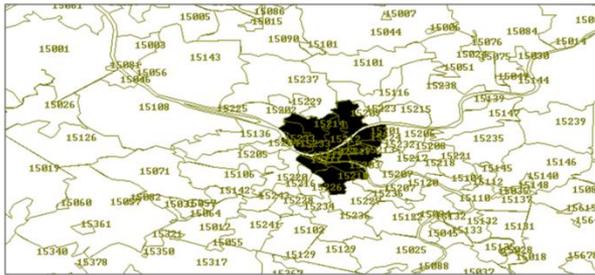
As measured by the average overlap coefficient between true and detected clusters, FSS outperformed MBSS by 10-15%.

For elongated and irregular clusters, FSS had much higher precision and recall. For compact clusters, FSS had higher precision, and MBSS had higher recall.

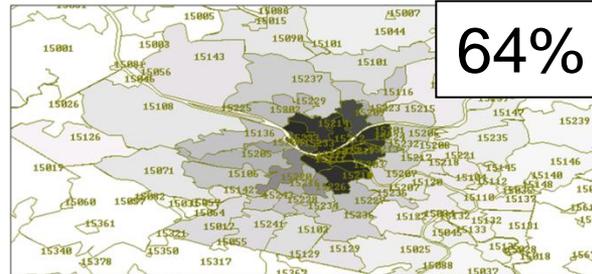
# Posterior probability maps

Spatial accuracy of FSS was similar to MBSS for compact clusters.

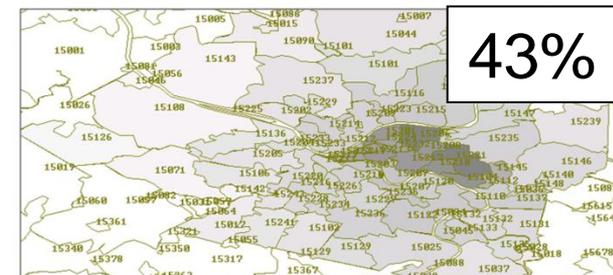
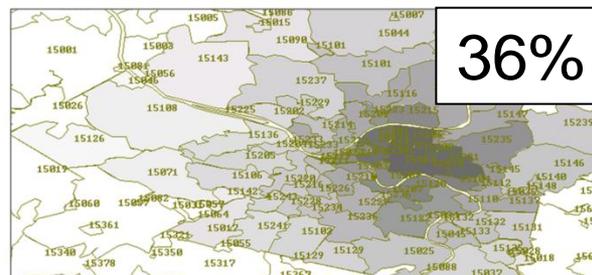
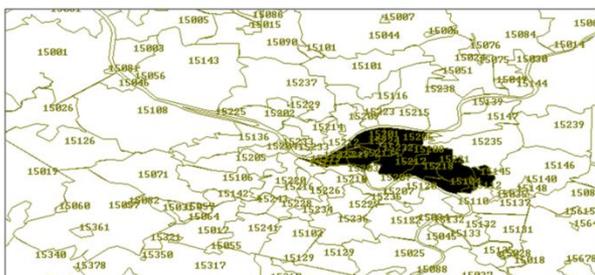
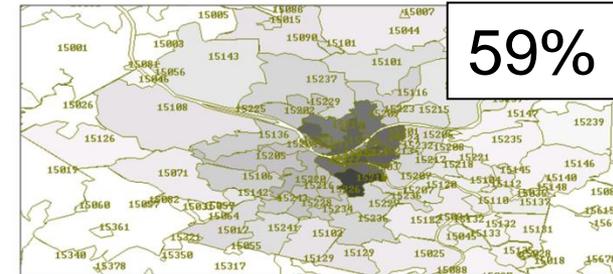
True outbreak region



MBSS



Fast Subset Sums







# Conclusions

FSS shares the essential advantages of MBSS: it can integrate information from **multiple data streams**, and can accurately distinguish between **multiple outbreak types**.

As compared to the original MBSS method, FSS substantially improves **accuracy** and **timeliness** of detection for elongated or irregular clusters, with similar performance for compact clusters.

While a naïve computation over the exponentially many subsets of the data is computationally infeasible, FSS can **efficiently** and **exactly** compute the posterior probability map.

Future work includes generalizing the hierarchical prior for FSS while maintaining efficient computation.

We can also **learn** the prior distribution from a small amount of labeled training data, as in Makatchev and Neill (2008).

# References

- **D.B. Neill. Fast Bayesian scan statistics for multivariate event detection and visualization. Submitted for publication.**
- D.B. Neill and G.F. Cooper. A multivariate Bayesian scan statistic for early event detection and characterization. *Machine Learning*, 2009, in press.
- M. Makatchev and D.B. Neill. Learning outbreak regions in Bayesian spatial scan statistics. *Proc. ICML Workshop on Machine Learning in Health Care Applications*, 2008.
- D.B. Neill. Incorporating learning into disease surveillance systems. *Advances in Disease Surveillance* 4: 107, 2007.
- D.B. Neill, A.W. Moore, and G.F. Cooper. A multivariate Bayesian scan statistic. *Advances in Disease Surveillance* 2: 60, 2007.
- D.B. Neill, A.W. Moore, and G.F. Cooper. A Bayesian spatial scan statistic. In *Advances in Neural Information Processing Systems* 18, 2006.