# Detecting Crime Types using Twitter

**Deepak Venugopal**
**University of Memphis**

**Motivation** Predict fine-grained crime occurrences that could help in efficient allocation of resources

**Approach** Use Twitter to provide real-time predictions of crime events specific to a Geo-Location

## Data
- Collect tweets from Twitter API for a specific geolocation
- Filter tweet text with a vocabulary of common. Words suggesting crimes
- Annotate each tweet with one of following categories of crime (Violent Crime, Narcotics, Racism, Fraud)
- Data collected from co-ordinates corresponding to TN
- Challenging to obtain sufficient data from specific geo-locations (e.g. Memphis)

### Twitter-based real-time prediction
- Built on top of Amazon Cloud Services
- Collect and extract language features from tweets using Spark Streaming (cluster computing framework)
- Apply machine learning methods to categorize tweets based on language features
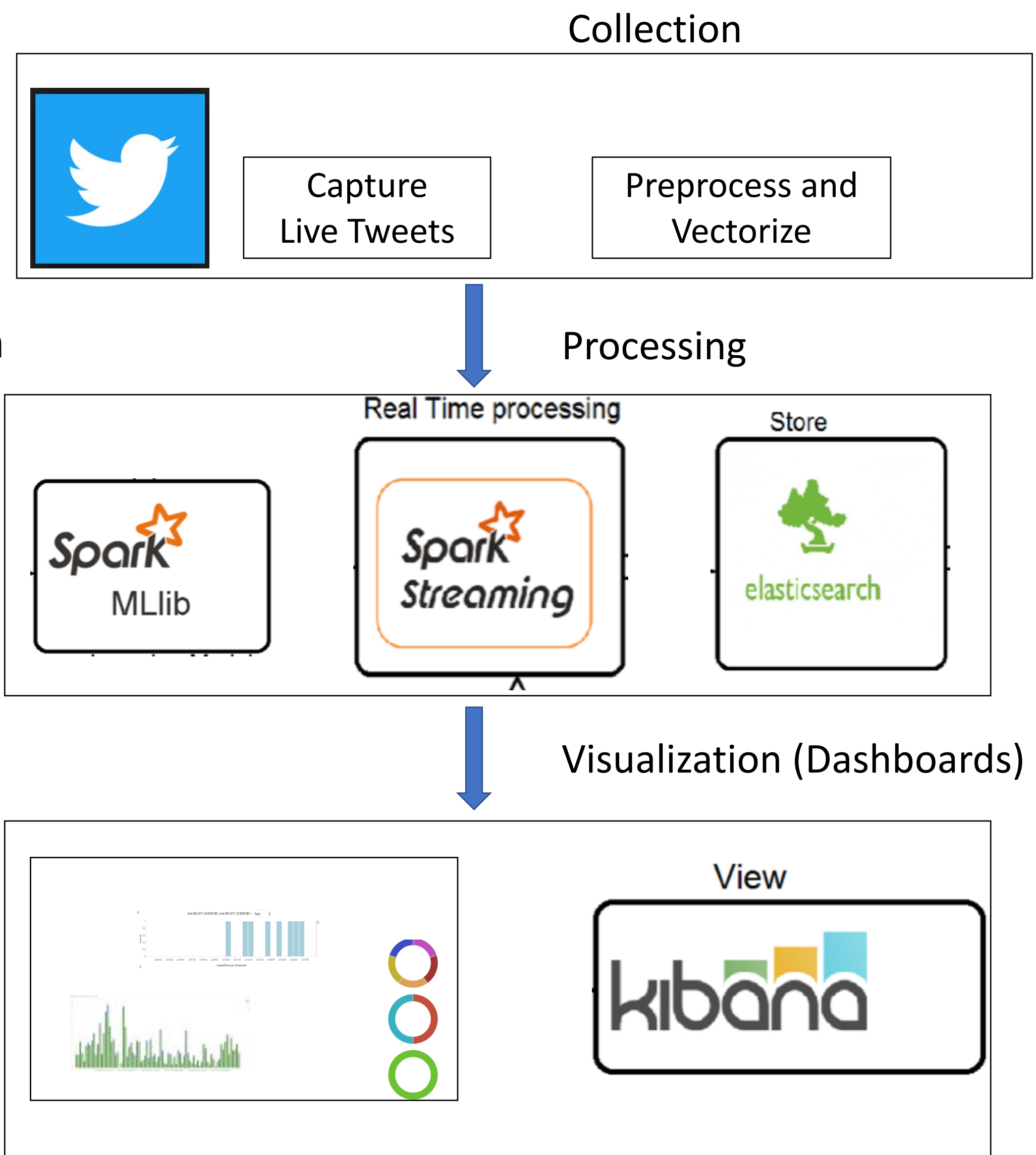- Rich visual interface to interpret prediction results through Amazon cloud-based Kibana

### HMM model
- Twitter data has the advantage that it is real-time
- However, twitter data can be very noisy and it is quite difficult to get reliable signals for crime prediction from twitter text alone
- Augment with real crime data from **BlueCrush**
- **BlueCrush** stores real crime events in Memphis, location and type of crime
- Develop a Hidden Markov Model (HMM) to predict future crime types based on recorded crimes
- Model latent crime states
- From observed incidents learn to transition between crime states
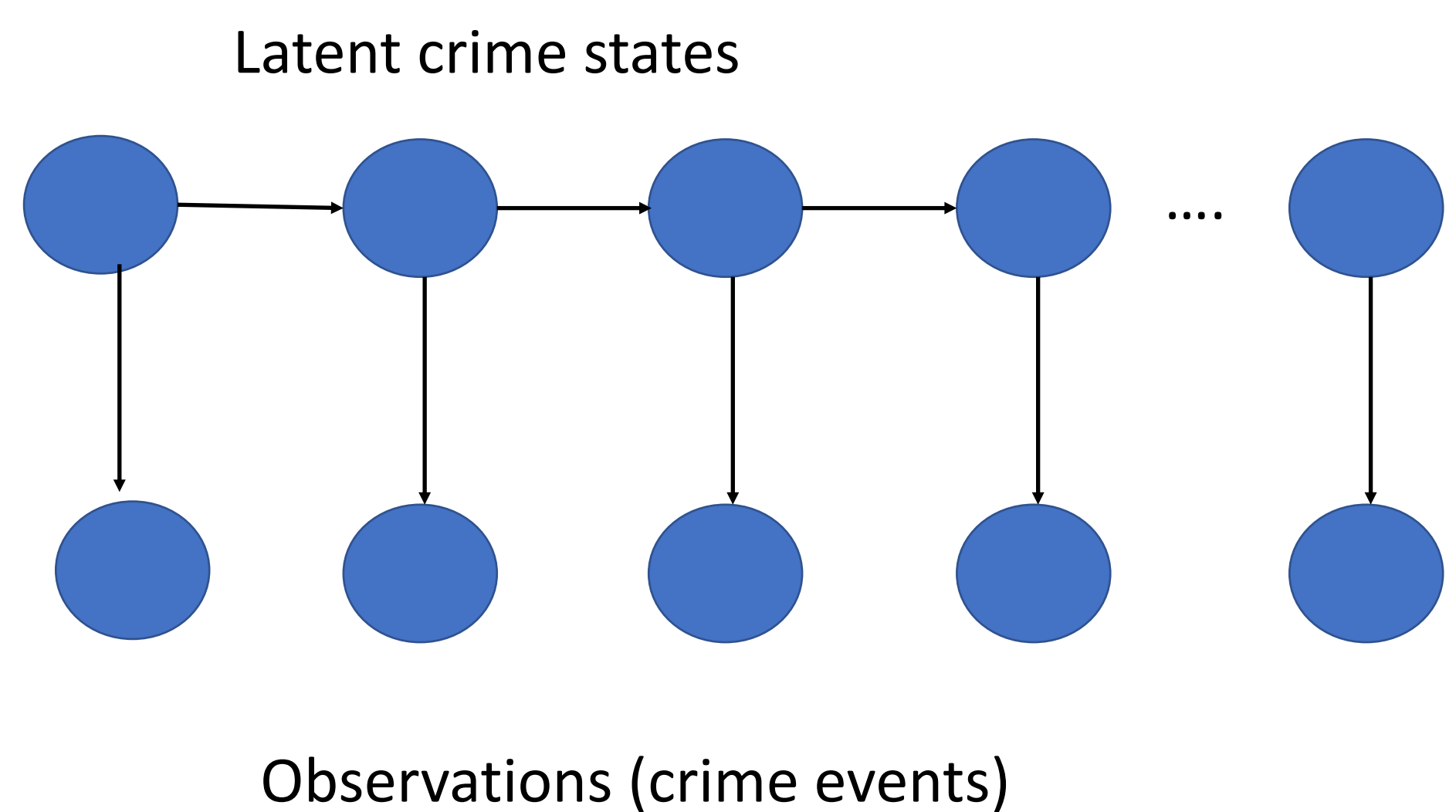
## Results
- 5000 Tweet Corpus Manually annotated for training

## Twitter-based real-time predictor
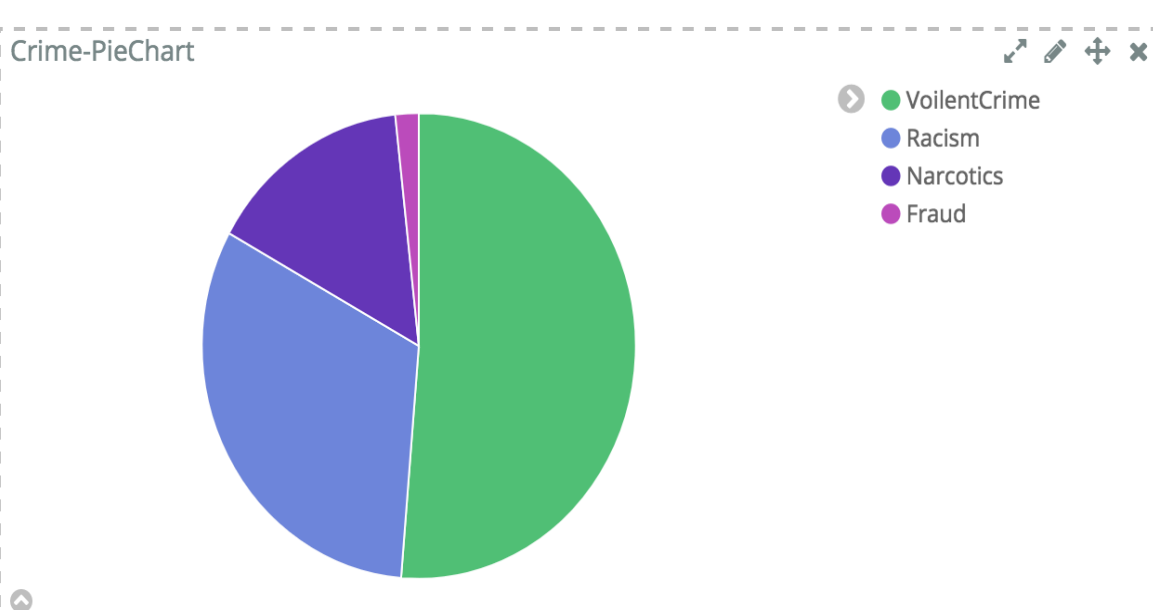


## HMM Model



## Overall System (Multi-source predictions)
- Real-time predictions using Twitter text
- Sample from the HMM model to obtain a distribution of crime states based real crime data
- Stack the predictions (ongoing work)

## Future Work
- Combine predictors from Twitter and HMM in a principled manner
- Derive advanced linguistic features using neural embeddings by taking advantage of unlabeled data
- Come up with detailed evaluation measures



Crime categories predicted on 5 days of tweets (thousands of tweets processed)

| Classifier | Fold-1 | Fold-2 | Fold-3 | Fold-4 | Fold-5 |
|---|---|---|---|---|---|
| Multinomia lNB | 0.79 | 0.8 | 0.8 | 0.8 | 0.83 |
| Logistic Regression | 0.83 | 0.86 | 0.83 | 0.86 | 0.86 |
| SVM (RBF) | 0.5 | 0.6 | 0.48 | 0.51 | 0.5 |
| **SVM (linear)** | **0.86** | **0.88** | **0.88** | **0.86** | **0.86** |

| States=3 | States=5 | States=10 |
|---|---|---|
| 0.15 | 0.12 | 0.17 |

HMM predictions using a corpus of incidents from 12/1/2018 – 2/16/2018 (Error for the predicted values in the last 10 days in the corpus)

5-fold Cross Validation on annotated twitter corpus (weighted F1-scores)