# Events, Relationships, and Script Learning

Presenter: Ed Morris

# Events, Relationships, and Script Learning Summary

**Problem:** Lack of automated methods that reliably recognize *actors*, *activities*, and *objects* comprising an *event* or *sequences of events* (i.e., *scripts*) within textual data sources

## FY17 Activities

1. Extract events from unstructured text sentences

2. DTRA BMIP and Cornerstone Support
   - BMIP: provides access to biological material information for situational awareness
   - Cornerstone: automated ID of dangerous biological material holdings worldwide

3. Extract structured summaries from document text
   - Use multiple clues within a document to improve accuracy
   - Complete pre-defined event templates

4. Extract events from social media (tweets) (Dr. Alan Ritter, OSU)
   - Use minimally supervised approaches for training NLP algorithms
   - Use redundancy (multiple tweets) to improve accuracy

Demonstration available at: http://kb1.cse.ohio-state.edu:8123/events/shooting

# Unstructured Text – Sentence Level Event Extraction Process
## Kevin Pitstick



Trigger: the word that most clearly expresses the Event's occurrence

**Example Event:**

# Unstructured Text: Trigger Identification
## Joint Event and Entity Model*

Train linear chain Conditional Random Field (CRF) models to identify candidate triggers and entities

- Features – word, part-of-speech tag, context words, word type, gazetteer-based, pre-trained word embedding

Uses a joint inference approach to find globally-optimal assignments for all trigger, argument, and entity variables

| Results | P | R | F1 |
|---|---|---|---|
| Reported | 77.6 | 65.4 | 71.0 |
| Replicated | 61.4 | 71.0 | 65.9 |

Significant false positives and negatives

**Precision**: $\dfrac{\text{true positives}}{\text{true positives + false positives}}$

**Recall**: $\dfrac{\text{true positives}}{\text{true positives + false negatives}}$

**F1**: weighted average of P and R

Differences between reported, replicated results could be due to slightly different models, or different dev, test datasets. We used replicated numbers to ensure consistency.

*Described in "Joint Extraction of Events and Entities within a Document Context" (Yang & Mitchell, 2016)*

# Unstructured Text: Trigger Identification
## Bi-directional LSTM + CRF

Concatenate character embeddings with pre-trained word embeddings to get vectors representing each word

Run a bi-LSTM over the sequence of word vectors to obtain the two hidden states

Use a CRF to find the sequence with the highest probability

| Results | P | R | F1 |
|---|---|---|---|
| JointEventEntity | 61.4 | 71.0 | 65.9 |
| Bi-LSTM + CRF | 66.7 | 66.0 | **66.4** |

Significant  false positives and negatives

Bottom line: Neither technique performs accurate trigger identification

# Unstructured Text: Trigger Classification
## Support Vector Machine (SVM)

Train SVMs to label triggers as belonging to one of 33 subtypes (e.g. attack, sentence, convict, etc.)

| Results | P | R | F1 |
|---|---|---|---|
| JointEventEntity (reported) | 75.1 | 63.3 | 68.7 |
| JointEventEntity (replicated) | 61.5 | 71.2 | 66.0 |
| SVM | 81.6 | 54.5 | 65.3 |

Relatively few false positives

Many false negatives

However, with perfect trigger identification, F1 for trigger classification is 80.0.

Bottom line: If we can find a better way to identify triggers, classification works fairly well

# Redirected Effort: Support for DTRA BMIP and Cornerstone
# Javier Vazquez-Trejo

Biological Materials Information Program: BMIP is a dynamic compendium of information concerning potentially dangerous biological material holdings worldwide and their security status



**BMIP Current Approach**

- Manual, time consuming data entry (1 week/facility)

**Cornerstone Objective**

- Automate ingest of data

**Cornerstone Approach**

- Mine PubMed by facility, collecting pathogen, equipment, personnel data.
- Incorporate analytic algorithms from U.S. and allies

SEI Role:
- Prototype NLP algorithms to extract information from abstract/body of PubMed articles
- Strategy to monitor the performance and behavior of Cornerstone algorithms (separately funded)

# Document Level Macro-Event Extraction Process
## Andrew Hsi, Daegun Won, Petar Stojanov, Dr. Jaime Carbonell

Use multiple clues within a document to improve accuracy of extraction

Complete pre-defined templates



| ATTACK Macro-Event | |
| --- | --- |
| Perpetrator | Michael Dunn |
| Victim – Dead | Jordan Davis |
| Victim – Injured | None |
| Time | November 23, 2012 |
| Location | Jacksonville |

| ARREST Macro-Event | |
| --- | --- |
| Arrestee | Michael Dunn |
| Time | (Unknown) |
| Location | (Unknown) |

| TRIAL Macro-Event | |
| --- | --- |
| Defendant | Michael Dunn |
| Crime | Murder |
| Verdict | Guilty |
| Sentence | Prison |
| Time | Friday |
| Location | Florida |

# Document Level Macro-Event Extraction Algorithms
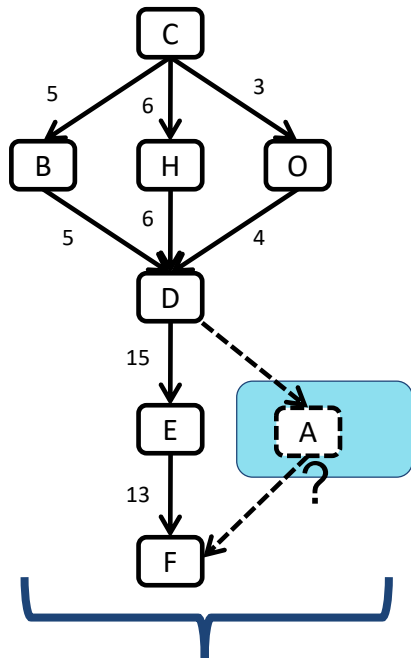
Novel ML-based algorithms for solving this problem

1. A structured prediction model based on Learning to Search

   - Reframes structured prediction as reinforcement learning problem with rewards for correct answers based on current state

   - Goal is to maximize the rewards for the system

   - Advantage: the entire document is considered without prohibitive expensive

2. A deep neural network based on machine comprehension with no reliance on target domain training data

   - Allows efficient retargeting to different domains

Preliminary results on the *attack* and *elections* domains show significantly improved performance against baseline methods

Currently gathering annotated data via Mechanical Turk

# Document Level Macro-Event Extraction
# Future Work: Relating Events to Scripts (CMU)

Better representation of the world

- Finer-grain event representation (actor, location, etc.)
- Probability distributions over the possible arguments

More robust script manipulation

- Better script addition (avoiding adding rare instances, etc.)
- Splitting / pruning existing scripts

Use of macro-event knowledge for inferring *constraints*

| ATTACK Macro-Event | |
|---|---|
| Perpetrator | Michael Dunn |
| Victim – Dead | Jordan Davis |
| Victim – Injured | None |
| Time | November 23, 2012 |
| Location | Jacksonville |

| TRIAL Macro-Event | |
|---|---|
| Defendant | Michael Dunn |
| Crime | Murder |
| Verdict | Guilty |
| Sentence | Prison |
| Time | Friday |
| Location | Florida |

Scripts are stereotypical sequences of related events

# Conclusion

**Summary**

- *Problem:* Lack of automated methods that reliably recognize *actors*, *activities*, and *objects* comprising an *event* or *sequences of events* (i.e., *scripts*) within textual data sources

- *FY17 Goal:* Develop event recognition strategies for sentences, documents, and social media

- *Results:*
  - Slight (but insufficient) improvement for event extraction from sentences, redirection of effort to support DTRA BMIP/Cornerstone
  - Good preliminary results for macro-event extraction from documents
  - Prototype for event extraction from social media

**Future Work**

- SEI Support for DTRA BMIP/Cornerstone
- CMU dissertation proposal for macro-event extraction from documents
- CMU pursuing development of scripts from macro-events

# Contact Information

**Presenter**

Ed Morris

MTS – Senior Engineer

Email:  ejm@sei.cmu.edu

**SEI Team**

- Kevin Pitstick
- Javier Vazquez-Trejo

**CMU Collaborators**

- Dr. Jaime Carbonell, Head, Language Technology Institute
- Andrew Hsi
- Petar Stojanov
- Daegun Won

**Ohio State University Collaborator**

- Dr. Alan Ritter