



# EXTRACTING LATENT ATTRIBUTES FROM VIDEO SCENES USING TEXT AS BACKGROUND KNOWLEDGE

Anh Tran, Mihai Surdeanu, Paul Cohen  
*University of Arizona*

*Third Joint Conference on Lexical and Computational Semantics (\*SEM 2014)*

## Motivation: Surveillance Application

- Two chase scenarios: Should the police be notified?



dark gun angry / fear / distress



light no-gun joyful / happy

- A surveillance system could automatically dispatch the police if it can identify that someone is in a **state of distress**.
  - Need **latent attributes** (e.g., mental state information) about the scene

## Problem Definition

---

*Identifying latent attributes from video scenes, with a focus on the mental states of activity participants, given some contextual information about the scenes.*

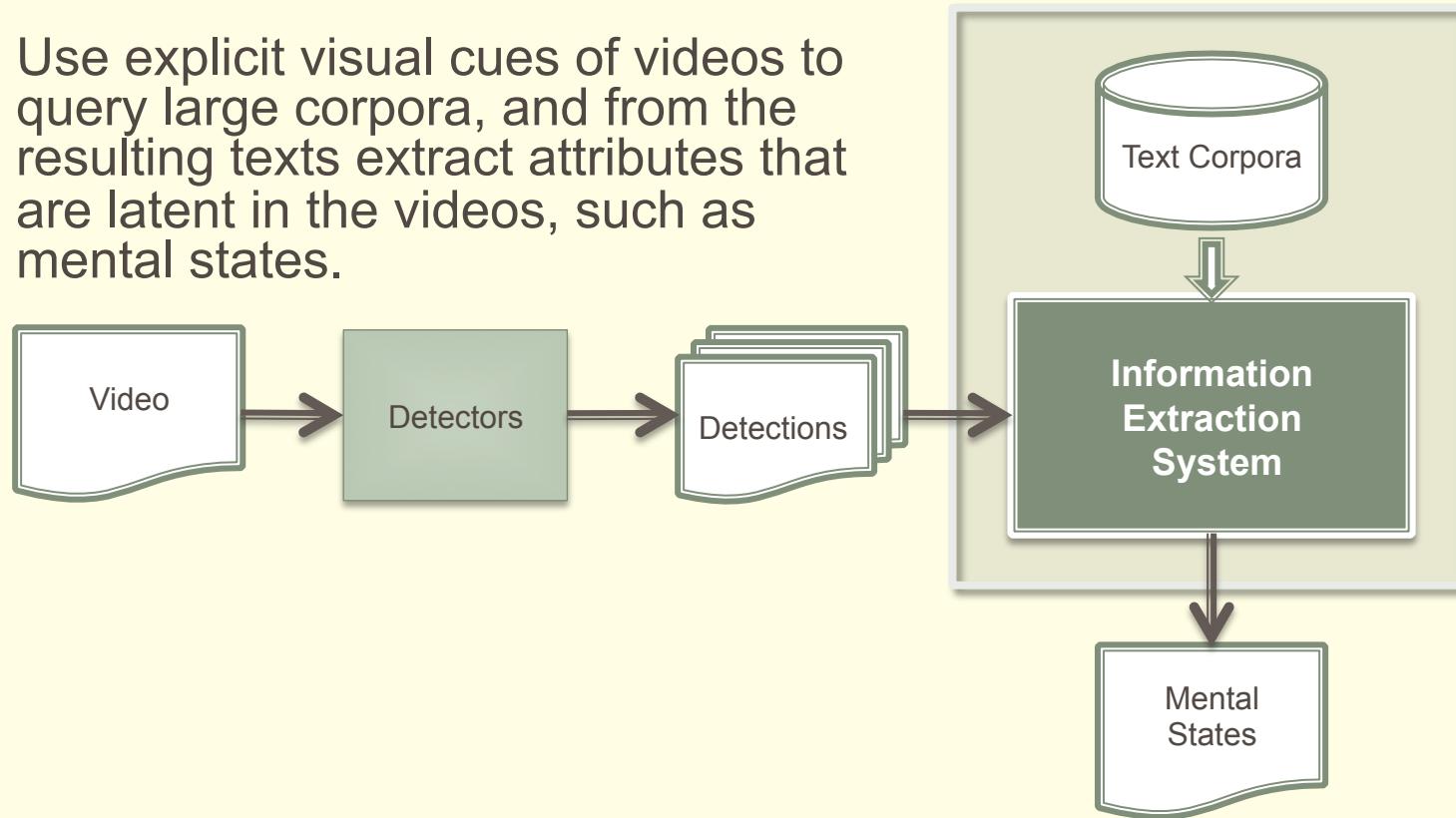
- **Latent attributes:** unobservable elements about the scene.
  - Mental states, motives, intentions, etc.
- **Contextual knowledge:** any observable elements (or cues).
  - Activity, object, actor type (child vs. policeman)
- Automatic identification of latent attributes is a challenging task.
  - No access to the same background knowledge that humans possess.
  - Machines can only “detect” explicit contents (e.g., observable cues).
  - How to identify latent information using these cues as the contextual knowledge?

# The Approach

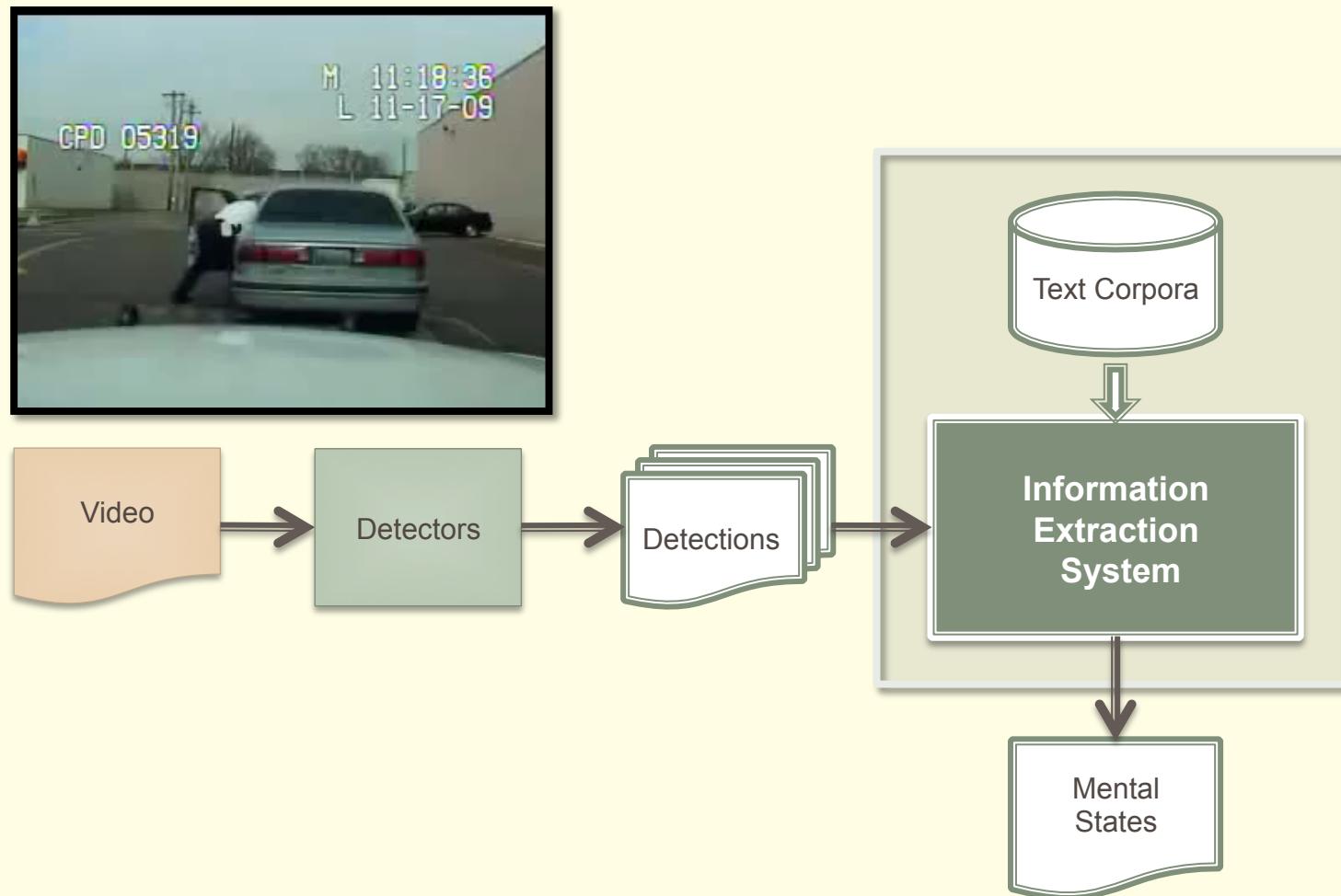
---

*Attributes that are latent in videos are often explicit in text.*

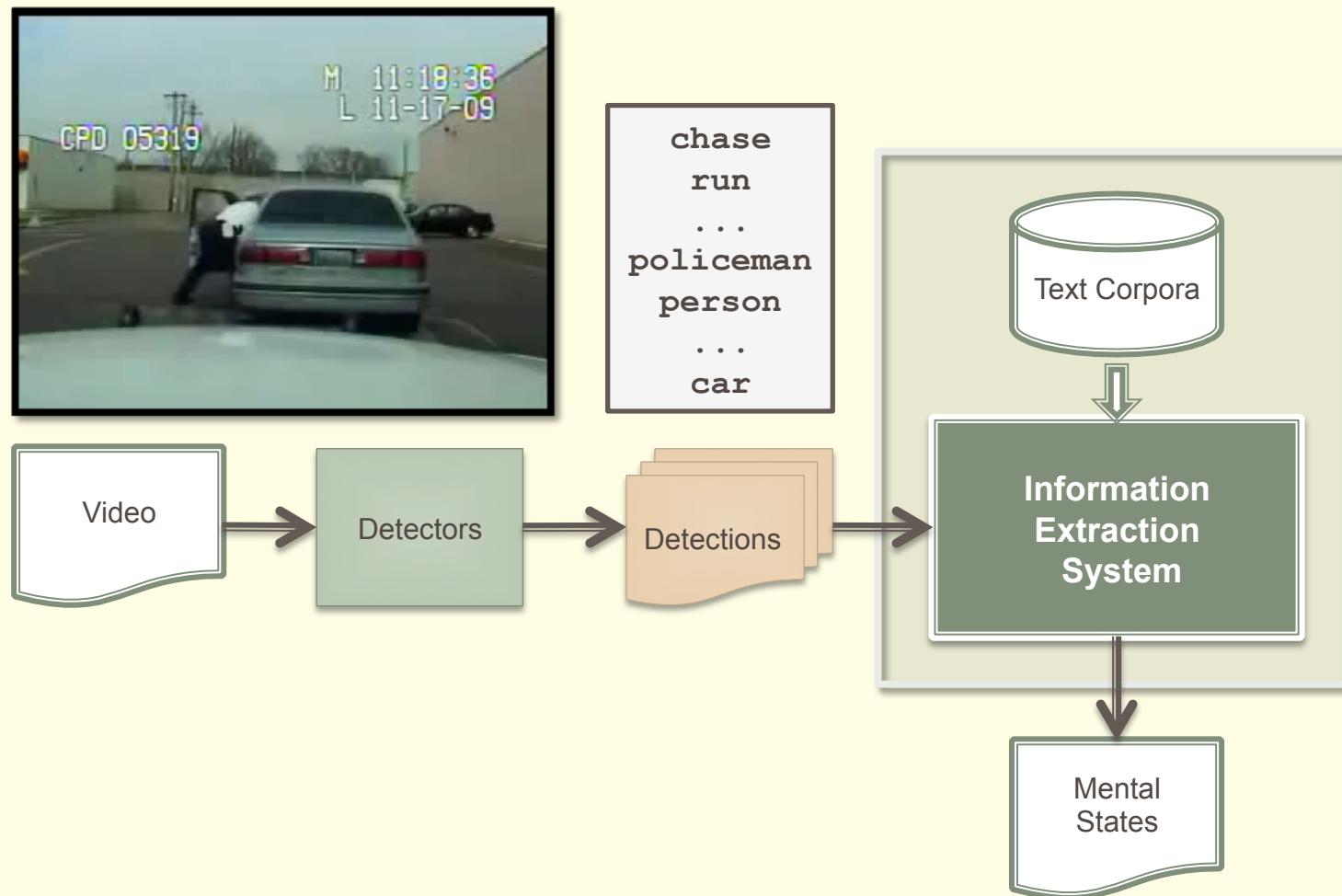
- Use explicit visual cues of videos to query large corpora, and from the resulting texts extract attributes that are latent in the videos, such as mental states.



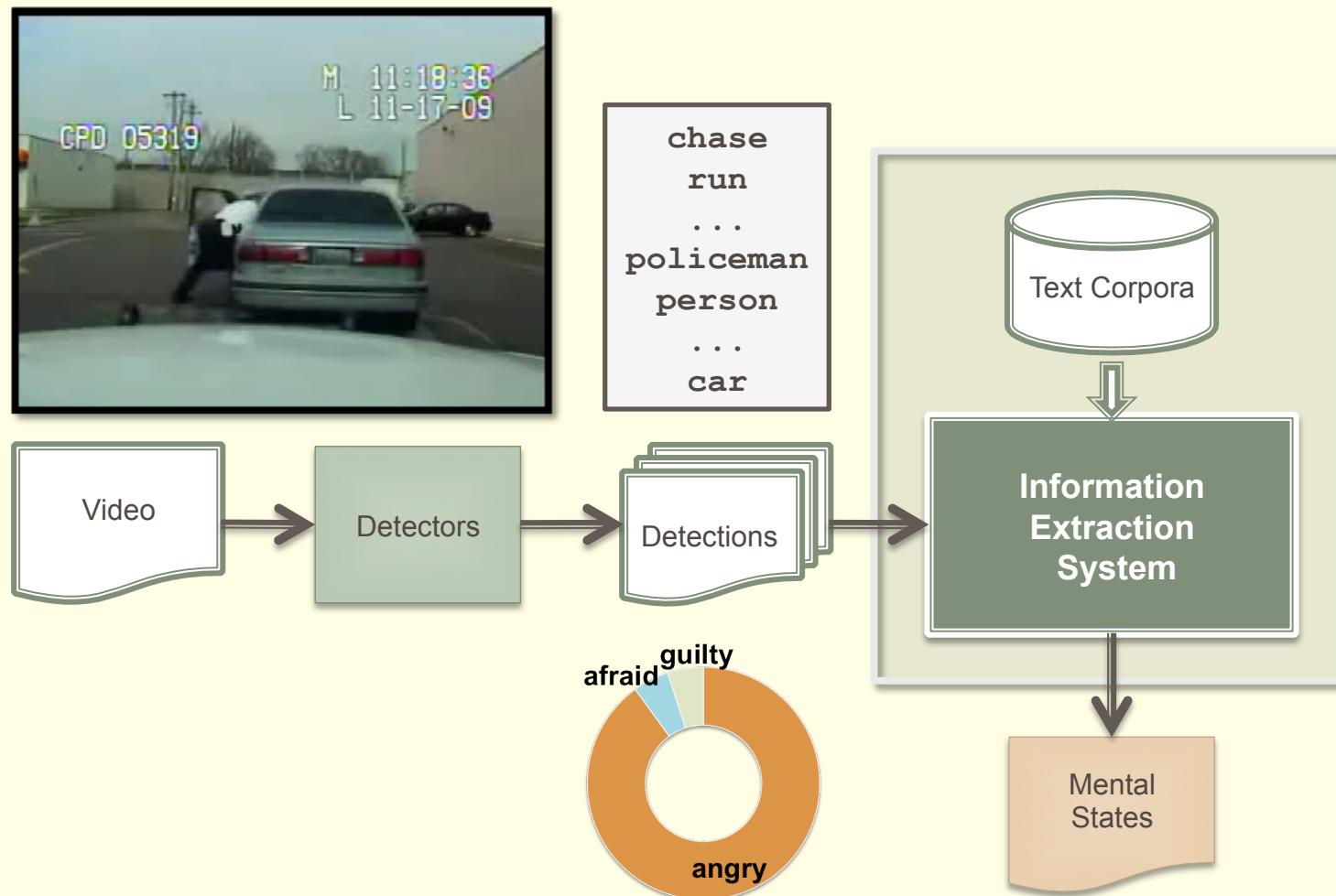
## Run-through Example – Video



## Run-through Example – Detection Labels

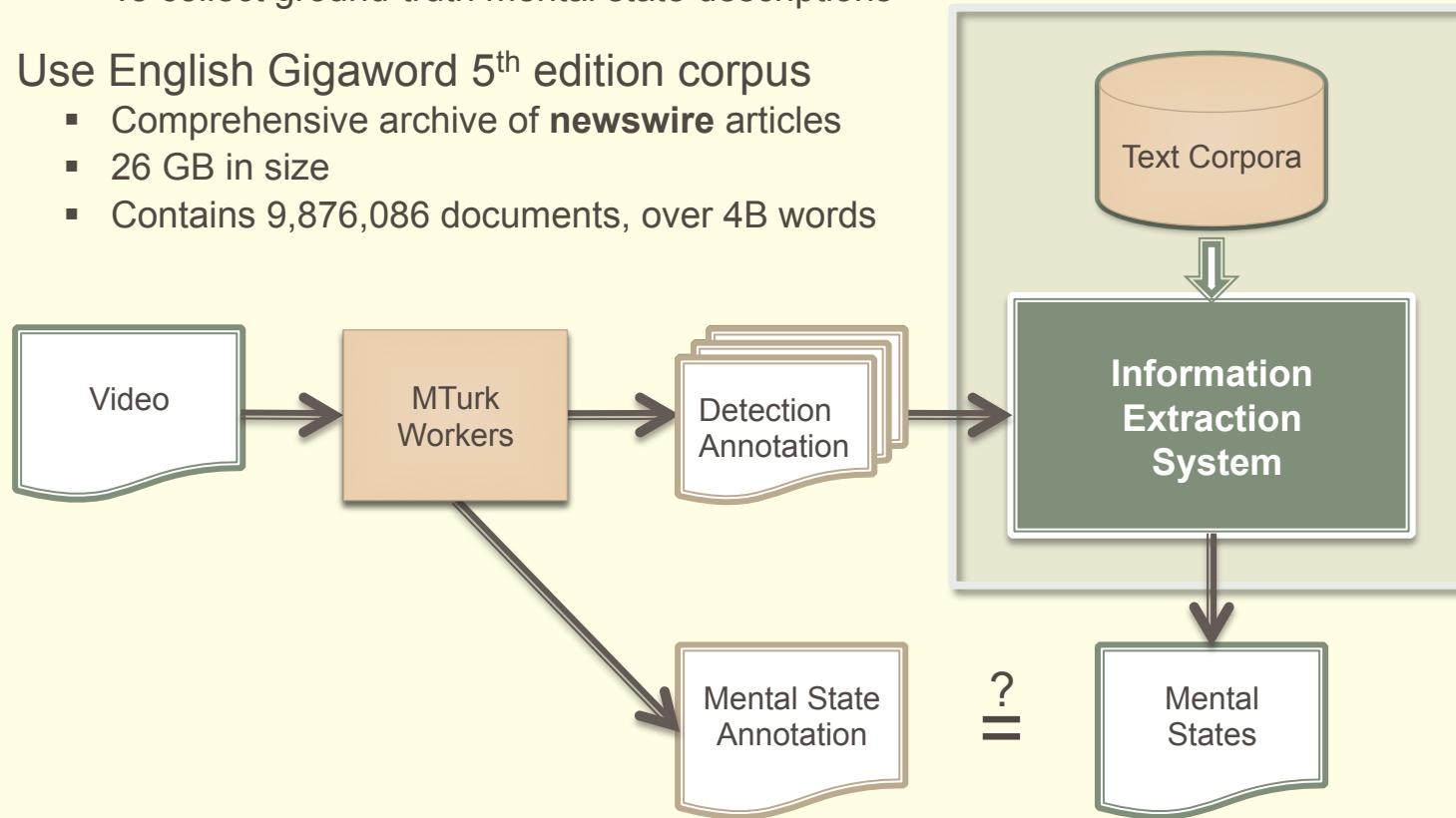


## Run-through Example – Mental States



# Data Sources

- Use Amazon Mechanical Turk workers
  - As proxy for automatic detection system
  - To collect ground-truth mental state descriptions
- Use English Gigaword 5<sup>th</sup> edition corpus
  - Comprehensive archive of **newswire** articles
  - 26 GB in size
  - Contains 9,876,086 documents, over 4B words



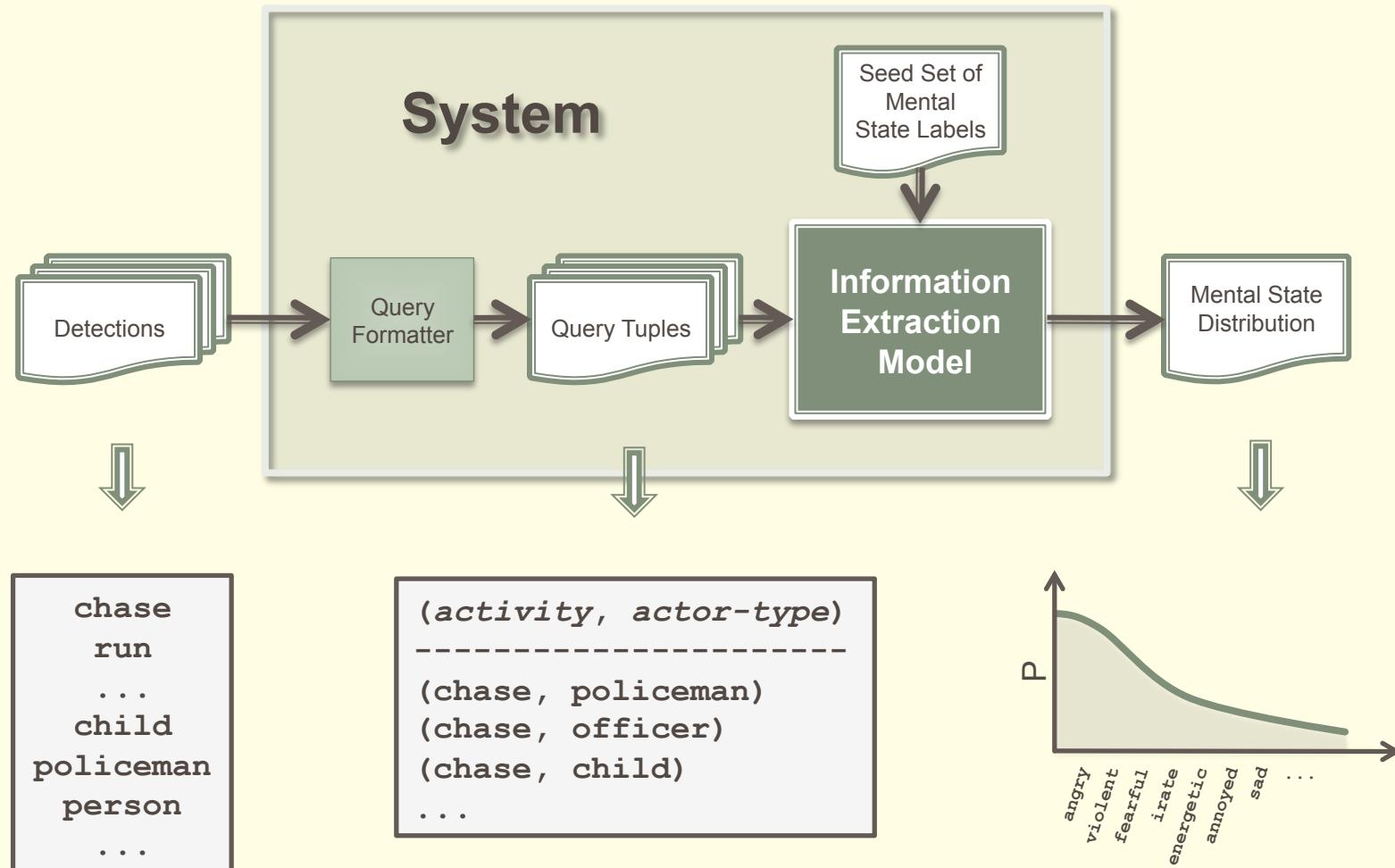
## Video Dataset

---

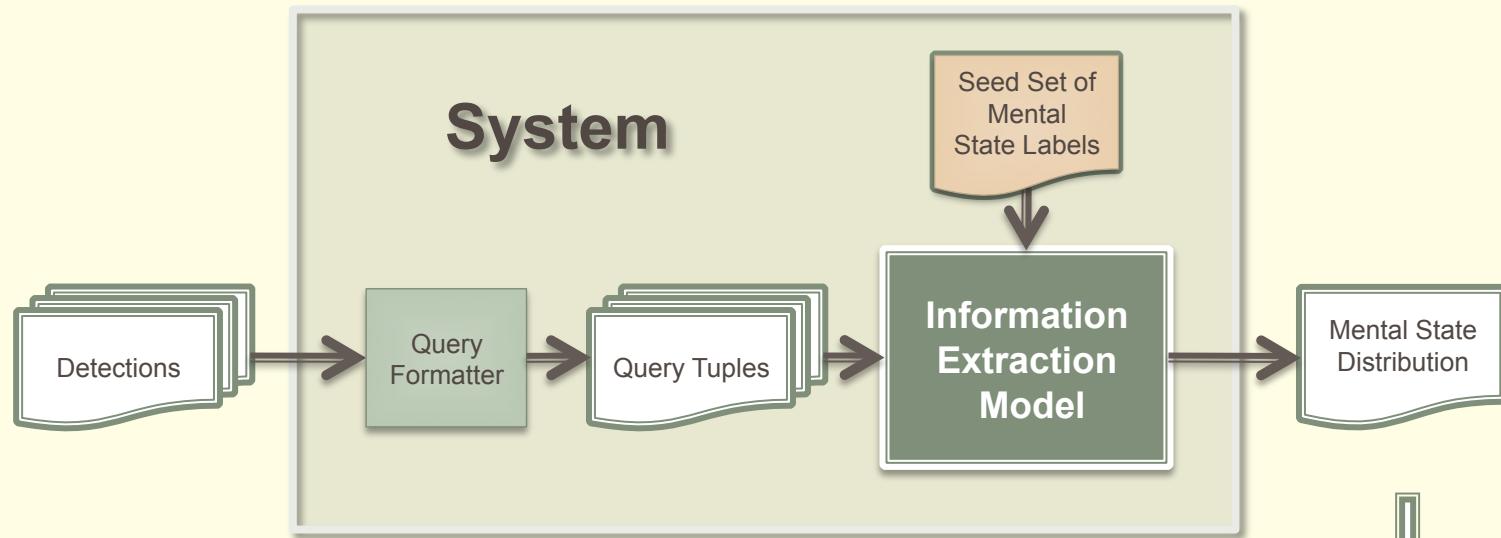
- Generated a dataset of 26 chase videos:
  - Mixture of police (5), children (7), sport-related (4), and other (12) chases.



# Detailed System Overview

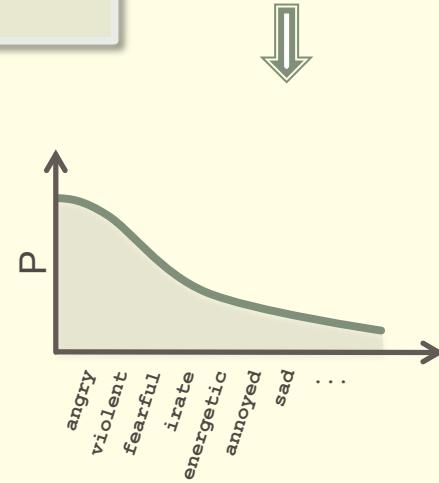


# Seed Set of Mental State Labels

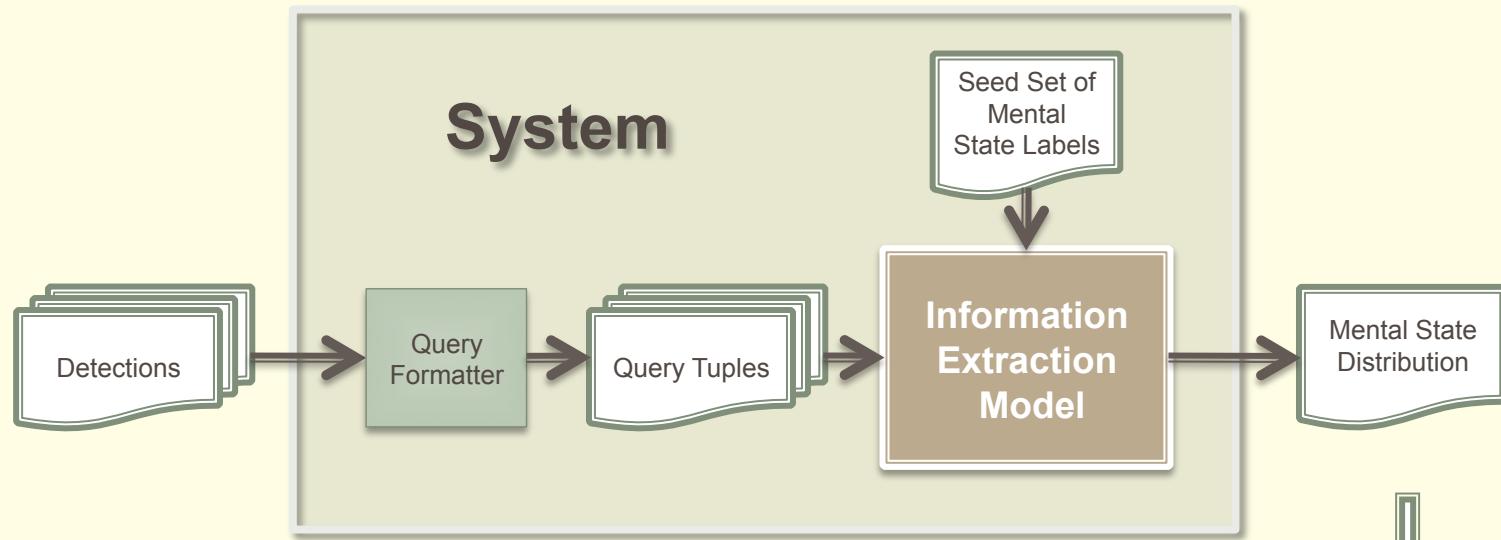


Source	Example Mental State Labels
POMS	alert, annoyed, energetic, exhausted, helpful, sad, terrified, unworthy, weary, etc.
Plutchik	angry, disgusted, fearful, joyful/joyous, sad, surprised, trusting, etc.
Others	agitated, competitive, cynical, disappointed, excited, giddy, happy, inebriated, violent, etc.

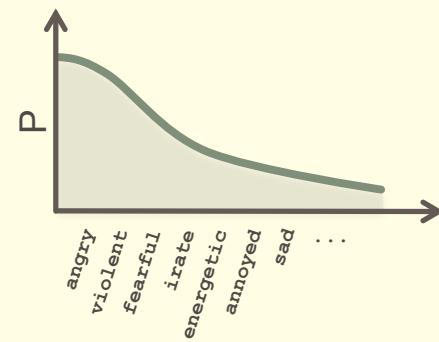
160 mental state labels total



# Neighborhood Models

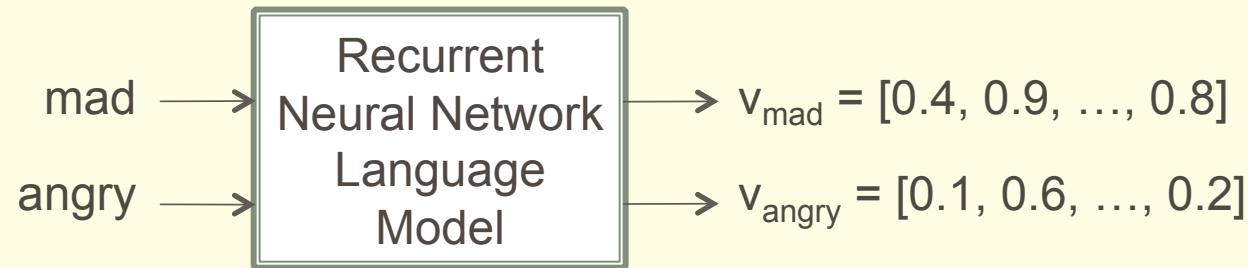


**Three largely unsupervised information extraction models.**

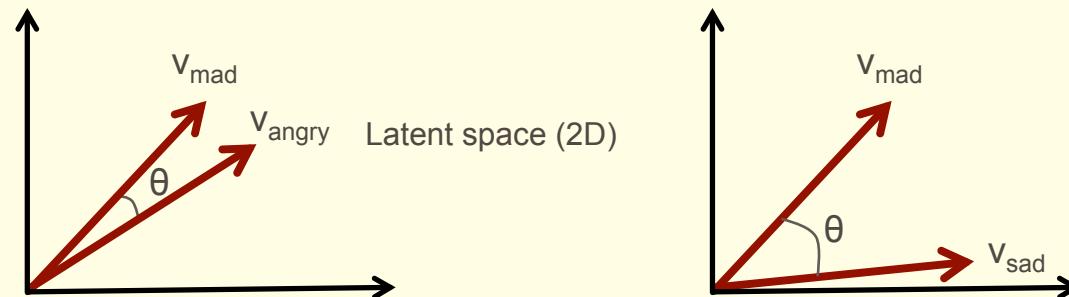


# Vector Space with Back-off Linear Interpolation

- Idea: Project mental state labels and search context into latent conceptual space produced by a RNNLM (Mikolov et al., 2013a).

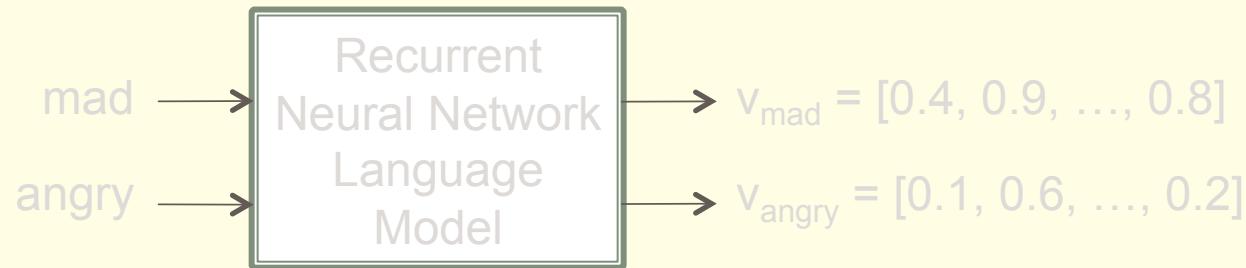


- Compare in latent space using the angle between the vectors.

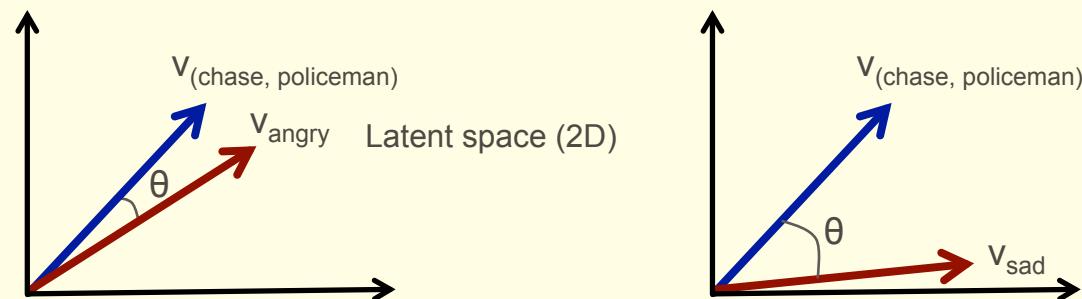


# Vector Space with Back-off Linear Interpolation

- Idea: Project mental state labels and search context into latent conceptual space produced by a RNNLM (Mikolov et al., 2013a).



- Compare mental state labels to query tuple in latent space.



# Vector Space with Back-off Linear Interpolation

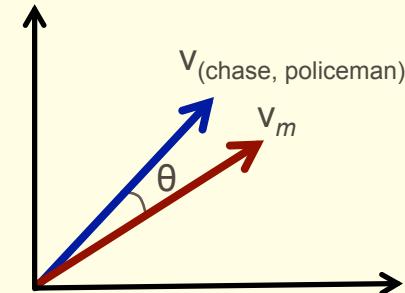
- Compute context-vector for query tuple:

$$\text{vec}(\text{chase}, \text{policeman}) = \text{vec}(\text{chase}) + \text{vec}(\text{policeman})$$

- Compute similarity to each mental state  $m$ :

$$\cos(\Theta_m) = \frac{\text{vec}(m) \cdot \text{vec}(\text{context tuple})}{\|\text{vec}(m)\| \|\text{vec}(\text{context tuple})\|}$$

- → 160 scores per context (or query) tuple.
- Normalize scores to generate a distribution per tuple, average across tuples to create one distribution, and prune to yield final response distribution.
- Improve robustness with back-off model (see paper for details)



## Sentence Co-occurrence with Deleted Interpolation

- **Idea:** Words in the same sentence are likely to be related.
- Rank mental state labels based on the likelihood that they appear in sentences cued by query tuples.
- Interested in the conditional probability:

$$P(m|activity, actor-type) = \frac{f(m, activity, actor-type)}{f(activity, actor-type)}$$

- Normally, we could compute this probability based on relative frequencies.
- However, estimation is unreliable due to sparse data!

## Sentence Co-occurrence with Deleted Interpolation

---

- Cannot estimate probability of trigrams reliably from the corpus, so we estimate probability as linear interpolation of unigrams, bigrams, trigrams.
- Define maximum likelihood probabilities  $\hat{P}$  based on relative frequencies:

$$\text{Unigram: } \hat{P}(m) = \frac{f(m)}{N}$$

$$\text{Bigram: } \hat{P}(m|activity) = \frac{f(m, activity)}{f(activity)}$$

$$\text{Trigram: } \hat{P}(m|activity, actor-type) = \frac{f(m, activity, actor-type)}{f(activity, actor-type)}$$

- $N$  = total number of tokens in the corpus
- $f(m, activity)$  = number of sentences containing both  $m$  as an adjective and  $activity$  as a verb

$$P(m|activity, actor-type) = \lambda_1 \hat{P}(m) + \lambda_2 \hat{P}(m|activity) + \lambda_3 \hat{P}(m|activity, actor-type)$$

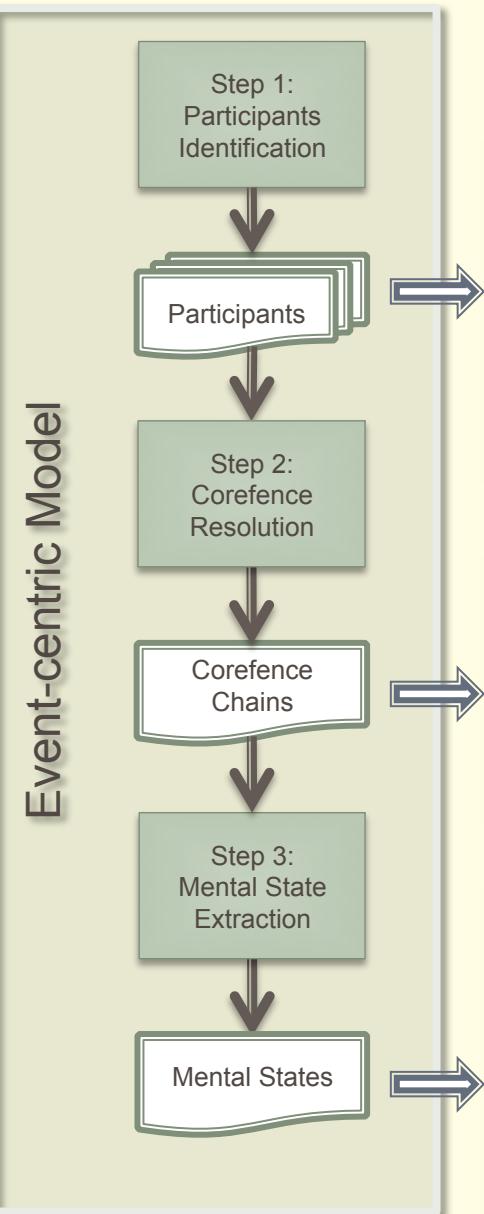
- Use deleted interpolation to estimate lambdas.
- 160 trigram probabilities for each query tuple, average across all query tuples and prune to yield final response distribution.



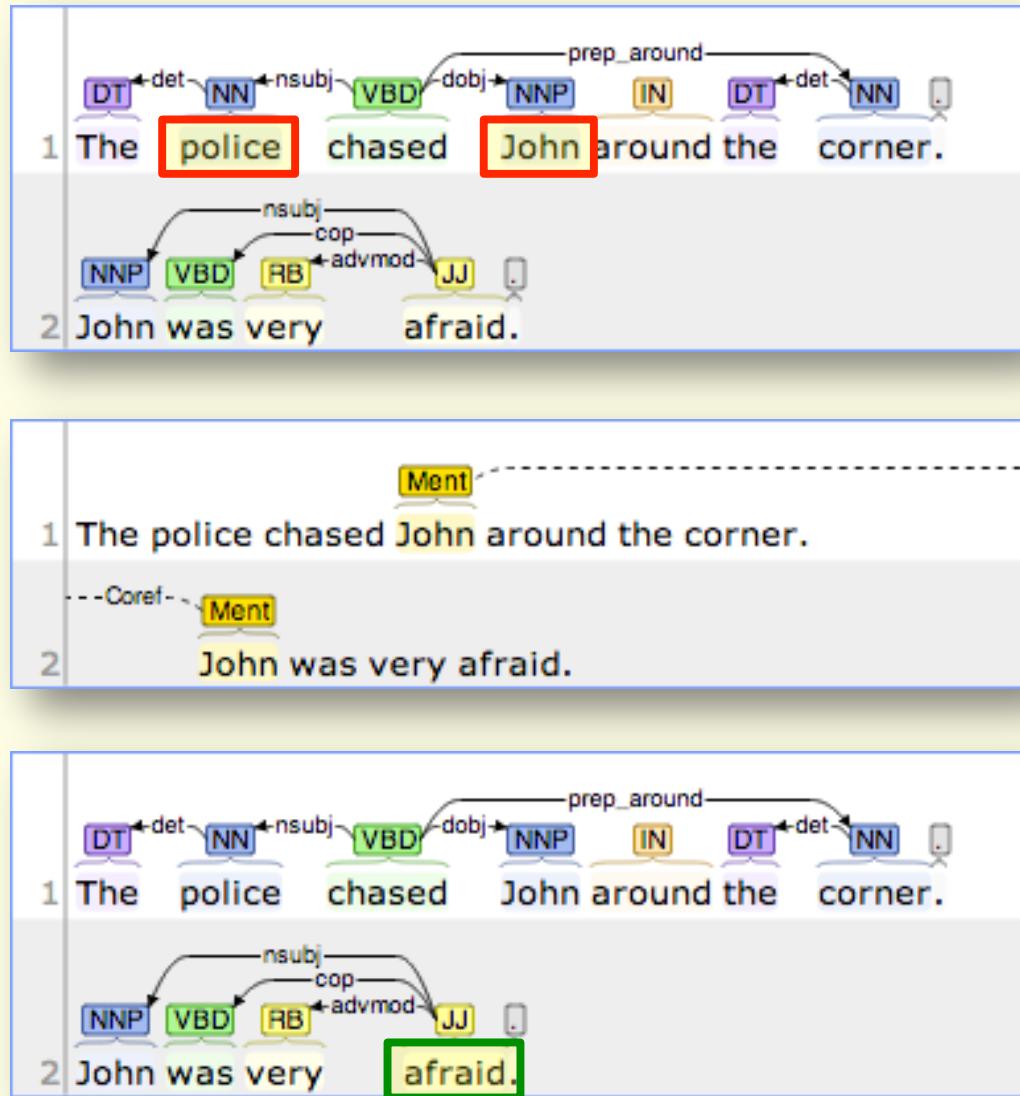
## Event-centric with Deleted Interpolation

---

- **Idea:** Identify the event + its participants in the relevant sentences and focus only on the mental states of event participants.
- A smarter, more robust, way to find collocating mental states for joint frequency estimation.
  - Go beyond sentence boundary.
  - Focus on mental states of participants.



## Example Output



# Evaluation Measure

---



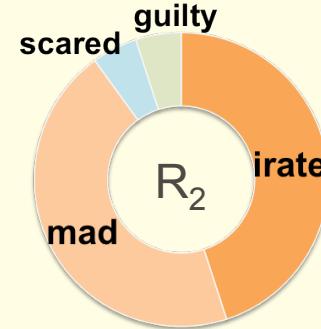
---

- New task → No standard performance measure
- Need to compare two normalized distributions over mental state labels.
- Similarity of **distribution shapes (weights)**
  - A good measure must account for the similarity between the shapes of the two distributions (i.e., ratios between weights)
- Semantic similarity of **distribution elements (synonyms)**
  - A good measure must allow for semantic comparisons at the level of distribution elements (i.e., recognize that irate and angry are similar)

Gold $G$	(angry, 0.9), (afraid, 0.05), (guilty, 0.05)
Response $R_1$	(angry, 0.1), (afraid, 0.2), (guilty, 0.7)
Response $R_2$	(irate, 0.45), (mad, 0.45), (scared, 0.05), (guilty, 0.05)

## Evaluation Toy Example

---




---

Gold  $G$        $(\text{angry}, 0.9), (\text{afraid}, 0.05), (\text{guilty}, 0.05)$

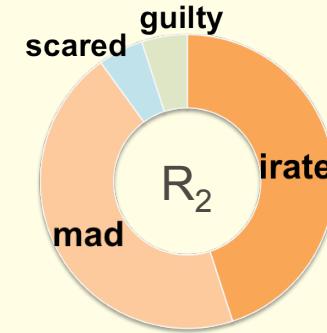
---

Response  $R_1$        $(\text{angry}, 0.1), (\text{afraid}, 0.2), (\text{guilty}, 0.7)$

---

Response  $R_2$        $(\text{irate}, 0.45), (\text{mad}, 0.45), (\text{scared}, 0.05), (\text{guilty}, 0.05)$

# Constrained Weighted Similarity-Aligned F1



	$F_1$			CWSA- $F_1$		
	p	r	$f_1$	p	r	$f_1$
$R_1$	1	1	1	0.2	0.2	0.2
$R_2$	0.25	0.33	0.29	1	1	1

red = non-intuitive score  
green = intuitive score

# Mental State Identification in Chase Videos

---

	$F_1$			CWSA- $F_1$		
	p	r	$f_1$	p	r	$f_1$
baseline	.107	<b>.750</b>	.187	.284	.289	.286
<i>sentence</i>	.194	.293	.227	.366	.376	.368
<i>vector</i>	.226	.145	.175	.399	.392	.393
<i>event</i>	.231	.303	.256	.446	.488	.463
<i>event+vector</i>	<b>.259</b>	.296	<b>.274</b>	<b>.488</b>	<b>.517</b>	<b>.500</b>

The average evaluation performance across 26 different chase videos are shown against the baseline scores for our neighborhood information extraction models. Bold font indicates the best score in a given column.

\* All average improvements over the baseline responses are significant ( $p < 0.01$ ). All significance tests were one-tailed and were based on nonparametric bootstrap resampling with 10,000 iterations.

# Mental State Identification in Chase Videos

	$F_1$			CWSA- $F_1$		
	p	r	$f_1$	p	r	$f_1$
baseline	.107	<b>.750</b>	.187	.284	.289	.286
<i>sentence</i>	.194	.293	.227	.366	.376	.368
<i>vector</i>	.226	.145	.175	.399	.392	.393
<i>event</i>	.231	.303	.256	.446	.488	.463
<i>event+vector</i>	<b>.259</b>	.296	<b>.274</b>	<b>.488</b>	<b>.517</b>	<b>.500</b>

- *event+vector* outperformed baseline by almost 75%.
- Ensemble outperformed individual components.
  - Operating in latent space (*vector*) and operating on text (*event*) yield complementary information.
- Incremental improvements due to NLP.
  - *sentence* is similar to current state-of-the-art (e.g., de Marneffe et al. 2010)
  - *sentence < vector < event*

## Additional Results\* – Actor-specific

---

	CWSA- $F_1$		
	p	r	$f_1$
baseline	.196	.195	.195
<i>vector</i>	.351	.338	.342
<i>event</i>	.353	.340	.340
<i>event+vector</i>	<b>.395</b>	<b>.400</b>	<b>.396</b>

The average evaluation performance for the mental state of the **subject** across 26 different chase videos.

	CWSA- $F_1$		
	p	r	$f_1$
baseline	.191	.181	.185
<i>vector</i>	.358	.374	.363
<i>event</i>	.383	.407	.391
<i>event+vector</i>	<b>.389</b>	<b>.415</b>	<b>.399</b>

The average evaluation performance for the mental state of the **object** across 26 different chase videos.

\* Results not included in paper

## Additional Results\* – Hug Dataset

	CWSA- $F_1$		
	p	r	f <sub>1</sub>
baseline	.226	.210	.217
<i>vector</i>	.347	.334	.339
<i>sentence</i>	.388	.378	.382
<i>event</i>	.406	.384	.394
<i>event+vector</i>	<b>.443</b>	<b>.437</b>	<b>.439</b>



- Performance average across 45 hug videos
- *event+vector* outperformed baseline by over 100%
- Consistent behavior as in chase videos:
  - Ensemble outperformed individual components.
  - Incremental improvement with each NLP module.

\* Results not included in paper

## Additional Results\* – Noisy Detections

---

- Introduce noise into detections based on published precision rates of actual state-of-the-art detectors in computer vision
  - False-positive: randomly insert new detection label
  - False-negative: randomly withhold an annotated detection label
- Average errors introduced per movie

Stat Type	No. Occurrences per Movie		
True Positives	3.17		
False Negatives	2.83		
True Negatives	9.75		
False Positives	7.25		

- Performance average across 20 different simulations (26 chase videos per simulation)

	$F_1$			CWSA- $F_1$		
	p	r	$f_1$	p	r	$f_1$
<i>event+vector</i>	.231	.255	.239	.443	.437	.439
baseline	.107	<b>.750</b>	.187	.284	.289	.286

\* Results not included in paper

# Conclusions

---

- Summary

- **Problem:** Identifying latent attributes in videos, given some context.
- **Data:** Videos from web, annotations via crowd sourcing
- **Solution:** Largely unsupervised information extraction models
  - Lexical semantic in vector space (*vector*)
  - Sentence co-occurrence in text (*sentence*)
  - Event-centric in text (*event*)
- **Evaluation:** CWSA-F<sub>1</sub> score to compare mental state distributions

- Findings

- First to show how to identify latent information from videos using text collections as the sole background knowledge
- More NLP → better performance
- Robust models: work on different datasets, tolerate noisy detections

- All code (Scala) + data (annotations & videos) at
  - <https://trananh.github.io/vlsa/>



# THE END



# BACKUP SLIDES

admiring	cranky	fearful	instinctive	pleased	sympathetic
afraid	crazy	focused	interested	protective	tense
aggressive	curious	forgetful	irate	raging	terrified
agitated	cynical	frantic	irritated	rebellious	terrifying
alarmed	demented	friendly	jealous	refreshed	thankful
alert	depressed	frightened	joyful	relaxed	threatening
ambitious	desperate	frustrated	joyous	relieved	tired
amazed	determined	fun	lively	reluctant	trustful
amused	devious	furious	loathsome	remorseful	trusting
angry	disappointed	giddy	lonely	resentful	uncomfortab
annoyed	discontent	glamorous	loved	restless	uneasy
anxious	discouraged	gleeful	mad	revengeful	unhappy
apprehensive	disgusted	grateful	mellow	romantic	unworthy
ashamed	distracted	grumpy	merciless	sad	upset
assertive	drunken	guilty	mischievous	satisfied	urgent
bitter	eager	happy	miserable	scared	vengeful
bored	ecstatic	helpful	motivated	selfish	vigilant
calm	encouraged	helpless	naughty	selfless	vigorous
carefree	energetic	homicidal	nervous	serious	violent
cautious	energized	hopeful	numb	shaky	wary
cheerful	enraged	hopeless	optimistic	shocked	weary
competitive	enthusiastic	hostile	panicked	sickened	weird
complacent	envious	hurried	panicky	spiteful	welcoming
concerned	excited	impressed	peaceful	stressed	worried
confused	exhausted	indifferent	peeved	submissive	worthless
considerate	exhilarating	inebriated	pessimistic	surprised	
content	fatigued	infuriated	playful	suspenseful	

```
participants-nlp-3
1 ===== DOCUMENT =====
2 The police chased John around the corner . John was very afraid .
3
4
5 ==> Step 1: PARTICIPANTS IDENTIFICATION
6
7 Target Sentence(s):
8 "The police chased John around the corner ."
9
10 Identified subject (CoreNLP): "police"
11 Identified object (CoreNLP): "John"
12
13
14 ==> Step 2: COREFERENCE RESOLUTION
15
16 Coreference chain(s) for SUBJECTS:
17 One chain found containing the following mentions:
18     sentence (0): [The police] chased John around the corner .
19
20 Coreference chain(s) for OBJECTS:
21 One chain found containing the following mentions:
22     sentence (1): [John] was very afraid .
23     sentence (0): The police chased [John] around the corner .
24
25
26 ==> Step 3: COMPLEMENTS EXTRACTION
27
28 Complement(s) for SUBJECTS:
29     sentence (0): NONE
30
31 Complement(s) for OBJECT:
32     sentence (0): NONE
33     sentence (1): afraid
```

# Constrained Weighted Similarity-Aligned F1

- We start with the standard  $F_1$  measure.

$$precision = \frac{|R \cap G|}{|R|}, recall = \frac{|R \cap G|}{|G|}, F_1 = 2 \cdot \frac{precision \cdot recall}{precision + recall} \quad F_1$$

- Generalize the formulas to address our criteria.

$$\begin{aligned} precision &= \frac{1}{|R|} \sum_{r \in R} \max_{g \in G} \sigma(r, g) & \sigma(r, g) &= \begin{cases} 1, & \text{if } r = g \\ 0, & \text{otherwise} \end{cases} & SA-F_1 \\ &= \sum_{r \in R} R(r) \cdot \max_{g \in G} \sigma(r, g) & R(r) &= \frac{1}{|R|} \\ &= \sum_{r \in R} R(r) \cdot \sigma_G^*(r) & & & WSA-F_1 \end{aligned}$$

- Address greedy problem of WSA-F<sub>1</sub>

$$M_S(\ell) = \{e \mid \sigma(\ell, e) = \sigma_S^*(\ell), \forall e \in S\}$$

$$precision = \sum_{r \in R} \min \left( R(r), \sum_{e \in M_G(r)} G(e) \right) \cdot \sigma_G^*(r) \quad CWSA-F_1$$

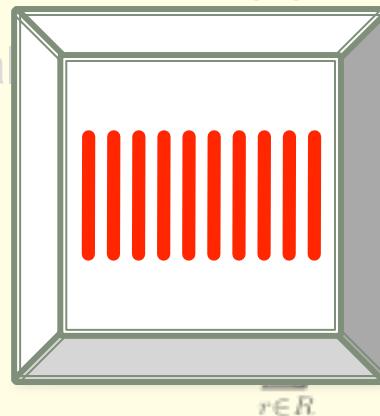
# Constrained Weighted Similarity-Aligned F<sub>1</sub>

- We start with the standard F<sub>1</sub> measure.

$$\text{precision} = \frac{|R \cap G|}{|R|}, \text{recall} = \frac{|R \cap G|}{|G|}, F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

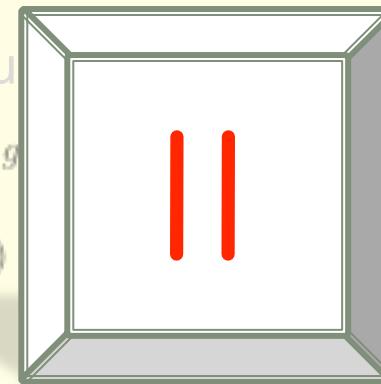
F<sub>1</sub>

- Generalize to address our problem



to address our problem

$$\begin{aligned} & \max_{g \in G} \sigma(r, g) && \sigma(r, g) \\ & \cdot \max_{g \in G} \sigma(r, g) && R(r) \\ & \cdot \sigma_G^*(r) && \end{aligned}$$



SA-F<sub>1</sub>

WSA-F<sub>1</sub>

- Address greedy problem of WSA-F<sub>1</sub>

$$M_S(\ell) = \{e \mid \sigma(\ell, e) = \sigma_S^*(\ell), \forall e \in S\}$$

$$\text{precision} = \sum_{r \in R} \min \left( R(r), \sum_{e \in M_G(r)} G(e) \right) \cdot \sigma_G^*(r)$$

CWSA-F<sub>1</sub>

## Limitation: Biases in Underlying Data

Categories	Baseline	<i>event+vector</i>	Change
children	0.2082	0.3599	+0.1517
police	<b>0.3313</b>	<b>0.6006</b>	<b>+0.2693</b>
sports	0.2318	0.4126	+0.1808
others	0.3157	0.5457	+0.2300

The average CWSA-F<sub>1</sub> scores for the ensemble model *event+vector* are shown in comparison to the baseline performance, categorized by video scenarios.

*children* = video contains a child participant

*police* = video contains a policeman participant

*sports* = video is sports-related

*other* = video does not fit in the first categories (e.g., civilian adults)

## Limitations: Biases in Underlying Data

---

Categories	Baseline	<i>event+vector</i>	Change
children	0.2082	0.3599	+0.1517
police	<b>0.3313</b>	<b>0.6006</b>	<b>+0.2693</b>
sports	0.2318	0.4126	+0.1808
others	0.3157	0.5457	+0.2300

- Baseline did worse on children and sports-related videos than police related videos.
- Baseline uses all 160 mental states with uniform probability.
- **Initial seed set is more fit to describe police chases.**
- See biggest improvement over baseline on police videos, least improvement on children videos.
- Gigaword corpus = newswire articles
- **Underlying corpus is biased towards police chases (i.e., news-worthy events).**

# Effectiveness of Coreference Resolution

Models	CWSA-F1	Versus <i>coref</i>	<i>p</i> -value
<i>win-0</i>	0.388682	−0.027512	0.0067
<i>win-1</i>	0.415328	−0.000866	0.4629
<i>win-2</i>	0.399777	−0.016417	0.0311
<i>win-3</i>	0.392832	−0.023362	0.0029

Comparing the average CWSA-F1 scores of a naïve windowing model, under different window sizes, to the performance of the *coref* model. The *p*-values, based on the average differences, were obtained using one-tailed nonparametric bootstrap resampling with 10,000 iterations.

*win-n* extends the single sentence boundary of *sentence* to also include the *n* preceding and *n* following sentences, while also piecing all relevant sentences of a document together to generate 1 neighborhood per document.

# Effectiveness of Coreference Resolution

---

Models	CWSA-F1	Versus <i>coref</i>	<i>p</i> -value
<i>win-0</i>	0.388682	−0.027512	0.0067
<i>win-1</i>	0.415328	−0.000866	0.4629
<i>win-2</i>	0.399777	−0.016417	0.0311
<i>win-3</i>	0.392832	−0.023362	0.0029

- *coref* outperformed all tested windowing configurations.
- Improvement over *win-1* is not significant.
  - *coref* and *win-1* generate very similar neighborhoods (extracted roughly the same number of sentences relevant to *chase*).
- ***coref does not do worse + provides references to participants for downstream process.***

Models	Total Sentences
<i>win-0</i>	90,399
<i>win-1</i>	260,423
<i>coref</i>	281,666
<i>win-2</i>	418,827
<i>win-3</i>	567,706

# Ensemble Models

---



---

	$F_1$			CWSA- $F_1$		
	p	r	$f_1$	p	r	$f_1$
<i>vector</i>	.226	.145	.175	.399	.392	.393
<i>sentence</i>	.194	.293	.227	.366	.376	.368
<i>sentence+vector</i>	.192	.377	.250	.434	.444	.438
<i>coref</i>	.264	.251	.253	.382	.461	.416
<i>coref+vector</i>	.231	.337	.271	.448	.481	.462
<i>event</i>	.231	.303	.256	.446	.488	.463
<i>event+vector</i>	.259	.296	.274	.488	.517	.500

- Combine a deleted interpolation model (text space) with vector model (latent space) creates an ensemble model.
- Every ensemble model outperformed its respective individual components.
- **Information gained from operating on text and operating in the latent vector space are highly complementary.**
  - Improvement to each will improve the resulting ensemble model.