

### ACE Event Extraction

- The Automatic Content Extraction (ACE) evaluations have included several information extraction tasks, including **entity recognition**, **value recognition**, **time recognition**, **relation extraction**, and **event extraction**.
- The **2005 ACE event extraction task** included 8 general event types and 33 event subtypes, with 35 possible event roles.
- An event mention is limited in scope to one sentence. The role fillers for each event mention must be identified within the event sentence. (This is different from template-based event extraction, which extracts event information from an entire document.)
- Each event has **event role fillers**, which include participants, objects, and other arguments like date and time. Importantly, an event can have multiple fillers (arguments) for a single role.

1

### Event Roles vs. Semantic Roles

- Semantic roles** (akin to thematic roles) capture arguments to a target word and represent the roles played in the action or concept directly expressed by the target.
- Example: *John was arrested by police for the murder of George.*

TARGET = **arrested** → THEME = **John** AGENT = **police**  
 TARGET = **murder** → THEME = **George**

- Event roles** capture arguments associated with a trigger word (in ACE) describing an action or concept associated with a higher-level event.
- Example: *John was arrested by police for the murder of George.*

EVENT = **DIE** →  
 PERPETRATOR = **John**  
 VICTIM = **George**

2

### ACE Terminology

- Entity**: an object that belongs to a semantic category.
- Entity mention**: a reference to an entity.
- Timex**: a time expression (e.g., day, year, date)
- Event mention**: a phrase or sentence that describes the occurrence of an event.
- Event trigger**: the main word that most clearly expresses an occurrence of a relevant event.
- Event mention arguments (event role fillers)**: entity mentions that are involved in an event and their relation to the event.

3

### ACE 2005 Entity Types and Subtypes

Type	Subtype
Life	Bi-Born, Mary, Divorce, Injure, Die
Movement	Transport
Transaction	Transfer Ownership, Transfer Money
Business	Start-Up, Merge/Org, Decline/Bankruptcy, Invest
Conflict	Attack, Demonstrate
Contact	Meet, Phone-Write
Personnel	Start-Position, End-Position, Nominate, Firing
Justice	Arrest-Jail, Release-Parole, Trial-Hearing, Charge-Indict, Sue, Convict, Sentence, Fine, Execute, Extradite, Acquit, Appeal, Pardon

An entity mention can be a proper name, nominal, or pronoun.

4

### ACE 2005 Event Types and Subtypes

Table 7: ACE05 Event Types and Subtypes	
Types	Subtype
Life	Bi-Born, Mary, Divorce, Injure, Die
Movement	Transport
Transaction	Transfer Ownership, Transfer Money
Business	Start-Up, Merge/Org, Decline/Bankruptcy, Invest
Conflict	Attack, Demonstrate
Contact	Meet, Phone-Write
Personnel	Start-Position, End-Position, Nominate, Firing
Justice	Arrest-Jail, Release-Parole, Trial-Hearing, Charge-Indict, Sue, Convict, Sentence, Fine, Execute, Extradite, Acquit, Appeal, Pardon

5

### ACE Event Extraction Example

Barry Diller on Wednesday quit as chief of Vivendi Universal Entertainment.

Trigger = **quit** (**RESIGNATION**)

Arguments (Roles)

Person: **Barry Diller**

Organization: **Vivendi Universal Entertainment**

Position: **chief**

Time-within: **Wednesday**

7

### ACE 2005 Event Roles (from Liao's dissertation)

Person	Place	Buyer	Seller
Beneficiary	Price	Artifact	Origin
Destination	Giver	Recipient	Money
Org	Agent	Victim	Instrument
Entity	Attacker	Target	Defendant
Adjudicator	Prosecutor	Plaintiff	Crime
Position	Sentence	Vehicle	Time-After
Time-Before	Time-At-Beginning	End	Time-Starting
Time-Ending	Time-Holds	Time-Within	

Table 1.2 - 35 Argument roles defined by ACE 2005

6

### ACE Event Extraction Example

Barry Diller on Wednesday quit as chief of Vivendi Universal Entertainment.

Trigger = **quit** [for End-Position event]

Arguments (Roles)

Person: **Barry Diller**

Organization: **Vivendi Universal Entertainment**

Position: **chief**

Time-within: **Wednesday**

8

### Multiple Triggers Example

Three murders occurred in France today, including the senseless slaying of Bob Cole and the assassination of Joe Westbrook.

Event	Trigger	Place	Victim	Time
DIE	<b>murder</b>	France	--	today
DIE	<b>slay</b>	France	<b>Bob Cole</b>	today
DIE	<b>assassinate</b>	France	<b>Joe Westbrook</b>	today

9

### Multiple Triggers Example

Three **murders** occurred in France today, including the senseless **slaying** of **Bob Cole** and the **assassination** of **Joe Westbrook**.

Event	Trigger	Place	Victim	Time
DIE	<b>murder</b>	France	--	today
DIE	<b>slay</b>	France	<b>Bob Cole</b>	today
DIE	<b>assassinate</b>	France	<b>Joe Westbrook</b>	today

10

### Sentence-level Event Extraction System

- [Grishman, Westbrook, & Meyers, 2005] developed a sentence-level event extraction system for ACE 2005.
- Using training data, event patterns are learned from sequences of constituent heads that separate a trigger word and its arguments.
  - A Trigger Labeler uses the patterns to distinguish event mentions from non-event mentions and classify each mention by type.
  - An Argument Detector is trained with MaxEnt classifier to distinguish arguments from non-arguments with respect to a trigger word.
  - A Role Classifier is trained with MaxEnt to label detected arguments with respect to event roles.
  - A Reportable-Event Classifier is trained with MaxEnt to determine whether a potential trigger, event type, and set of arguments are describing a true event occurrence.

11

### Sentence-level Extraction Pipeline

- Each test document is searched for instances of trigger words that occurred in the training documents.
  - For each trigger word, the patterns learned for that trigger are applied to identify arguments (with role labels) of the trigger.
  - The argument detector is applied to the remaining entity mentions in the sentence to look for more arguments.
  - If new arguments are found, the role classifier is applied to assign a role to each one.
  - Finally, the reportable-event classifier is applied to the entire context to decide whether it is truly an event.
- [Ji & Grishman, 2008] and [Liao & Grishman, 2010] use this system as a component in their document-level event extraction pipelines.

12

### Sentence-Level vs. Document-Level

#### Sentence-Level Event Extraction:

Traditionally, most systems have extracted information about an event from an isolated sentence. Each sentence in a document is processed independently of the others.

#### Document-Level Event Extraction:

Recently, researchers have begun to incorporate discourse properties and information about associations across sentences in a document to improve event extraction performance.

13

### Motivating Observations

- "Within a document, there is a strong trigger consistency: if one instance of a word triggers an event, other instances of the same word will trigger events of the same type."

True > 99.4% of the time in the ACE corpus.

- "Normally one entity, if it appears as an argument of multiple events of the same type in a single document, is assigned the same role each time."

True > 97% of the time in the ACE corpus.

15

### Enforcing Consistency

[Gale, Church, & Yarowsky, 1992] identified the widely recognized **One Sense Per Discourse** heuristic: within a discourse, instances of the same word have a strong tendency to share the same sense.

[Ji & Grishman] made a related observation that strong sense and event role consistency exists across related documents.

**One Trigger Sense Per Cluster:** In topically-related documents, event trigger words have a strong tendency to share the same sense.

**One Argument Role Per Cluster:** In topically-related documents, an entity has a strong tendency to participate in the same event role (argument).

14

### Shared Trigger Sense Example

#### Test Sentence:

*Most US army commanders believe it is critical to pause the breakneck **advance towards** Baghdad to secure the supply lines and make sure weapons are operable and troops resupplied ...*

#### Related Document:

*British and US forces report gains in the **advance on** Baghdad and take control of Umm Qasr, despite a fierce sandstorm which slows another flank.*

"**Advance towards**" wasn't in the training data, but "**advance on**" was. Identifying "**advance**" as a Movement\_Transport event trigger in a related document suggests the same sense for in the new document.

16

### Correcting Trigger Senses

Test Sentence:  
But few at the Kremlin forum suggested that Putin's own standing among voters will be hurt by Russia's apparent diplomacy failures.

Related Document:  
Putin boosted ties with the United States by throwing his support behind its war on terrorism after the Sept. 11 attacks, but the Iraq war has hurt the relationship.

Originally, hurt was mistakenly identified as a Life\_Injure event in the test sentence because it is a common trigger word for that event type.  
But ... hurt is never a trigger for Life\_Injure events in the typically related documents, so that trigger label can be discarded.

17

### Shared Argument Role Example

Test Sentence:  
Vivendi earlier this week confirmed months of press speculation that it planned to sell its entertainment assets by the end of the year.

Related Documents:  
Vivendi has been trying to sell assets to pay off huge debt, estimated at the end of last month at more than \$13 billion.

Under the reported plans, Blackstone group would buy Vivendi's theme park division, including Universal Studios Hollywood, ...

Originally, "Vivendi" was not recognized as a seller in the test document.  
But it was extracted as a seller in several typically related documents, which suggests it is likely to be a seller in the new documents too.

18

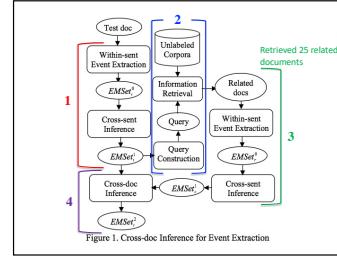
### Cross-Sentence and Cross-Document Extraction

[Ji & Grishman, 2008] created a system that improved event extraction output by making inferences that enforce consistency across within-document sentences as well as related documents.

A pipeline architecture gradually refines the output:

1. sentence-level event extraction
2. document-level (cross-sentence) inferences: rules enforce consistency across sentences in the same document.
3. cross-document inferences: an IR system retrieves related documents and rules enforce consistency with these documents.

19



20

2/2:

### Empirical Evidence for Cluster Heuristics

Candidate Triggers	Event Type	Per./Freq. as trigger	Per./Freq. as target	Per./Freq. as argument	Per./Freq. as related mentions
Correct Event Triggers	Movement Transport	31% of 2	50% of 2	88.9% of 27	
	Personnel End Position	7% of 2	100% of 2	100% of 2	
	Center Area	14% of 2	100% of 2	100% of 2	
	Replace	5% of 20	100% of 1	83.3% of 6	
	Form	Business Assn/Org	1% of 2	100% of 2	
	Contract/Mkt	59% of 24	100% of 2	100% of 26	
	Hurt	24% of 33	0% of 24	0% of 7	
Incorrect Event Triggers	Life_Inj	12% of 8	0% of 4	4% of 24	

21

### Rules for Cluster-Wide Consistency

**Rule (1): Remove Triggers and Arguments with Low Cluster-wide Confidence**  
If XDoc-Trigger-Preg (trigger, arg) <  $\delta_1$ , then delete EM.

**Rule (2): Adjust Trigger Classification to Achieve Cluster-wide Consistency**  
If XDoc-Trigger-Merge(trigger) >  $\delta_2$ , then propagate the most frequent trigger to all event mentions with trigger in the document, and propagate the trigger for the event type.

**Rule (3): Adjust Trigger Identification to Achieve Document-wide Consistency**  
If XConf-Trigger-Arg (trigger, arg) >  $\delta_3$ , then propagate trigger to all other triggers propagate the event.

**Rule (4): Adjust Argument Identification to Achieve Document-wide Consistency**  
If XConf-Arg-EM (trigger, arg) >  $\delta_4$ , then in the document, for each sentence containing an event mention EM with trigger, add any unlabeled mention in that sentence with the same head as arg as an argument of EM with role.

23

### Rules for Within-Document Consistency

**Rule (1): Remove Triggers and Arguments with Low Local Confidence**  
If Doc-Trigger-Preg (trigger, arg) <  $\delta_1$ , then delete arg.

**Rule (2): Adjust Trigger Classification to Achieve Document-wide Consistency**  
If XConf-Trigger-Arg (trigger, arg) >  $\delta_2$ , then propagate the most frequent trigger to all event mentions with trigger in the document, and propagate the trigger for the event type.

**Rule (3): Adjust Trigger Identification to Achieve Document-wide Consistency**  
If XConf-Trigger-Arg (trigger, arg) >  $\delta_3$ , then propagate trigger to all other triggers propagate the event.

**Rule (4): Adjust Argument Identification to Achieve Document-wide Consistency**  
If XConf-Arg-EM (trigger, arg) >  $\delta_4$ , then in the document, for each sentence containing an event mention EM with trigger, add any unlabeled mention in that sentence with the same head as arg as an argument of EM with role.

**Rule 1:** removes questionable triggers  
**Rule 2:** fixes event types and roles (hopefully)  
**Rules 3 & 4:** discover more triggers and arguments

22

Performance System/Human	Trigger Identification +Classification			Argument Identification			Argument Classification Accuracy			Argument Identification +Classification		
	P	R	F	P	R	F	P	R	F	P	R	F
Within-Sentence (E with Rule (1)) (Baseline)	53.5	47.8	38.3	42.5	36.0	41.2	32.9	36.6	36.6	86.0	81.2	77.7
Cross-Sentence (E with Rule (1))	53.5	47.8	38.3	42.5	36.0	41.2	32.9	36.6	36.6	86.0	81.2	77.7
Cross-doc (E with Rule (1))	64.7	65.4	55.7	59.4	46.2	54.1	51.3	56.7	56.7	45.6	45.6	45.6
Human Annotator1	59.2	59.4	59.3	60.0	69.4	64.4	64.8	64.8	64.8	54.1	54.1	54.1
Human Annotator2	69.2	75.0	72.0	63.7	85.4	73.3	86.3	86.3	86.3	73.7	73.7	73.7
Inter-Annotator Agreement	41.9	38.8	40.3	55.2	46.7	50.6	91.7	90.6	90.6	42.9	46.4	46.4

The gold standard was adjudicated by two people (H1 & H2).  
The last three rows are:  
H1 vs. System , H2 vs. System, H1 vs. H2

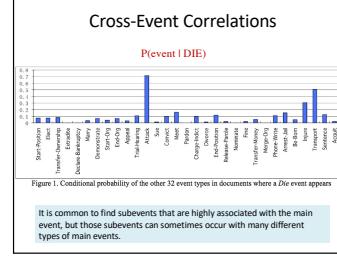
24

### Cross-Event Inference for Event Extraction

[Liao & Grishman, 2010] observed that certain types of events frequently co-occur, so they incorporated **cross-event information** into their event extraction system.

For example:  
S1: He left the company.  
S2: He planned to go shopping before heading home.  
left → TRANSPORT EVENT  
S2: His colleagues threw a retirement party for him.  
left → END-POSITION EVENT

25



26

2/2:

### Two-Pass Approach

- Liao & Grishman adopt a two-pass approach that first identifies the "easy cases" and then uses that knowledge to help identify the harder cases.

*The pro-reform director of Iran's biggest-selling daily newspaper and official organ of Tehran's municipality has died again following the appointment of a conservative ... it was founded a decade ago ... but a conservative city council was elected in the February 28 municipal polls ... Mohammad Ahmadinejad, reported to be a hardliner among conservatives, was appointed mayor on Saturday ... Founded by former mayor Gholamhosse Karbaschi Hamshahri ...*

*British officials say they believe Isis was a blindfolded woman seen being shot in the head by a beheaded militant on video obtained but not yet shown to the public by Britain's Al-Jazeera channel. She would be the first foreign woman to die in the wave of kidnappings in Iraq — she's been killed by (men in jinamas) ...*

29

### Confident Event Table

- The sentence-level event extraction system is applied to a document to identify high confidence predictions of event triggers and arguments.
- Event triggers and arguments (roles) that are labeled with high confidence are stored in a **Confident Event Table**.
- A word is assumed to be a trigger for only one type of event, and an entity is assumed to belong to just one role for an event trigger.

If multiple labels are assigned to the same word/entity, then the highest scoring label is chosen if the difference between scores is large. If the difference between scores is small, then there is a conflict so the information is recorded in a separate **Conflict Table**.

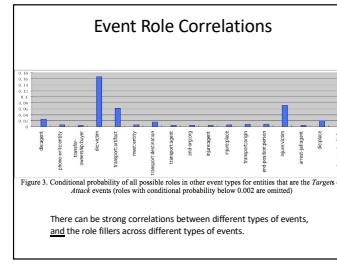
30

### Events that Co-occur with DIE

Event	Cond. Prob.
Attack	0.714
Transport	0.507
Injure	0.306
Meet	0.164
Arrest/Jail	0.124
Sentence	0.126
Phone-Write	0.111
End-Position	0.116
Trial-Hearing	0.105
Convict	0.100

Table 3. Events co-occurring with a Die events with conditional probability > 10%

27



28

Confident table		
Event type table		
Trigger	Event Type	Role
Met	Meet	Time Within
Exploded	Attack	Place
Went	Transport	Agent
Injured	Injure	Target
Attacked	Attack	Place
Died	Die	Agent

Argument role table		
Entity ID	Event type	Role
0004-2	Die	Time Within
0004-6	Die	Place
0004-8	Die	Agent
0004-7	Die	Target
0004-11	Attack	Place
0004-73	Attack	Agent
0004-9	Attack	Target
0004-10	Attack	Attacker

Conflict table		
Entity ID	Event type	Role
0004-8	Die	Victim, Agent

Table 4. Example of document-level confident-event table (event type and argument role entries) and conflict table.

31

### Document-Level Trigger Classifier

- A MaxEnt classifier is trained to predict whether a word is the trigger of an event, and if so, what type.
- The information in the confident event table is used to create features representing the other event types that have been found in the document.
- Each feature is the conjunction of:
  - the base form of the word
  - for each of the 33 event types, a binary value indicating whether this event type is present elsewhere in the document.

32

### Document-Level Argument (Role) Classifier

- A MaxEnt classifier is trained to predict whether a given mention is an argument of a given event, and if so, what role it plays.
- The information in the confident event table is used to create features for other event roles associated with this entity in the document.
- Each feature is the conjunction of:
  - the type of the given event
  - for each of the other 32 event types, the role of the given entity with respect to that event type (if one exists) or else *null*.

33

### Putting it All Together

- First, the sentence-level event extraction system is applied and high-confidence triggers and arguments are labeled.
  - Next, the document-level trigger classifier is applied to all words that do not already have a label. This will often identify some additional event triggers.
  - Finally, the document-level argument tagger is applied to all event triggers. Only entity mentions in the same sentence that have not already been assigned a role are considered.
- This tagger can identify arguments for the newly identified triggers as well as new arguments for the previously identified triggers.

34

### Conclusions

- Event extraction is a difficult problem – recall and precision are still only mediocre.
- But a variety of recent systems have shown that considering the entire document, as well as related documents, can be beneficial.
- These systems are still relatively shallow in their understanding of event descriptions. More explicit, richer event representations are probably needed to push performance to a higher level.

37

### Experiments

- The 2005 ACE data set was used for evaluation: 549 training texts, 10 tuning texts, 40 test texts.
- Two baseline systems were evaluated:
  - the sentence-level event extraction component by itself
  - the [Ji & Grishman, 2008]'s approach for cross-sentence and cross-document ("within-event-type") inference rules.
- They also looked at the performance of two human annotators on 28 documents in the test set.

35

### Evaluation Results

Performance system/human	Trigger classification			Argument classification			Role classification		
	P	R	F	P	R	F	P	R	F
Sentence-level baseline system	67.56	53.54	59.74	46.45	37.15	41.29	41.02	32.81	36.46
Within-event-type rules	63.03	59.90	61.43	48.59	46.16	47.35	45.33	41.16	42.21
Cross-event statistical model	68.71	68.87	68.79	50.85	49.72	50.28	45.06	44.05	44.55
Human annotation1	59.2	69.4	59.3	60.0	69.4	64.4	51.6	59.5	55.3
Human annotation2	69.2	75.0	72.0	62.7	85.4	72.3	54.1	73.7	62.4

Table 5. Overall performance on blind test data

36

3/:

### Neural Models for Event Extraction

- A variety of neural models have been developed for sentence-level event extraction. As with many other tasks, NNS began with simple architectures and have gotten increasingly complex.
  - Some models improve F scores, but it's not always easy to understand why or whether they are truly better in general.
- A common theme in IE is to jointly tackle multiple tasks, for example performing entity recognition and event extraction at the same time. This idea is appealing because:
  - In a pipeline architecture, errors that occur early are propagated to later models and can't be undone.
  - Intuitively, it makes sense that different tasks can inform each other!

1

### A Joint Neural Model for Information Extraction with Global Features

- [Lin et al., ACL 2020]
- This paper describes a recent neural network architecture called **OneIE** that performs all of the subtasks required for sentence-level event extraction in a single framework.
  - A key aspect of this model is that it creates an **information network** that captures information extracted across the entire sentence.
    - This network allows the model to consider **cross-task** and **cross-instance** interactions.
  - The overall model is refreshingly straightforward: relatively simple individual task classifiers plus a decoding step at the end to find a globally optimal information network (graph) with learned features.
  - This model is also language independent, so can be used for multiple natural languages.

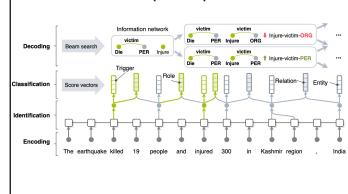
2

### The Four Stages of OneIE

1. Sentence Encoding: a contextualized embedding representation is produced for the input sentence.
2. Identification of entity mentions and event triggers:
  - these will become nodes in the information network
3. Local Classification of Nodes and Edges:
  - Labeling the nodes: entities with entity types, and the event triggers with event types.
  - Labeling edges with relations and event roles.
4. Finding the globally optimal information graph with a beam decoder.

5

### Joint Information Extraction Network (OneIE)



6

### The IE Tasks

- Entity Extraction:** identifying and labeling entity mentions, which can be names, nominals, or pronouns.  
Ex: finding all location phrases, such as *Kashmir region*
- Relation Extraction:** assigning a relation type to an ordered pair of entity mentions.  
Ex: PART-WHOLE(*Kashmir region*, *India*)
- Event Extraction:** identifying event triggers and their arguments  
Ex: *'assassination of a government official'* should identify *assassination* as a DIE event trigger and *government official* as the Victim of the DIE event.

3

### Motivating Example

Prime Minister Abdullah Gul resigned earlier Tuesday to make way for Erdogan, who won a parliamentary seat in by-elections Sunday.

The graph below shows system output with an extra person edge.

We'd like an IE system to recognize that ELECT events should not have 2 person arguments! Plus, Gul already has a role.

4

### Step 1: Sentence Encoding

- Input:** a sentence of L words
- The sentence is passed to a pre-trained BERT encoder (transformer model), which produces a contextualized word embedding for each word in the sentence.
- BERT sometimes splits words into pieces, so they average the embedding vectors of each word's pieces.
- Ex: *Mondrian* → *Mon #dr #ian*
- These word embedding vectors are the input to the next stage of processing.
- They use a multilingual BERT model for the Chinese & Spanish data.

7

### Multilingual BERT

- BERT can be trained to handle multiple natural languages!
- The multilingual BERT model (on HuggingFace) is a language model that has been pre-trained with the 104 languages that have the most Wikipedia pages.
- The texts are lower-cased and tokenized. This produced a shared vocabulary size of 110,000.
- For languages like Chinese, Japanese Kanji and Korean Hanja that don't have space, a CJK Unicode block is added around every character.
- The languages with a larger Wikipedia are under-sampled and the ones with lower resources are oversampled.

8

### Step 2: Identifying Entities and Event Triggers

- Two feed-forward neural networks (FFNs) are trained: one for **entity recognition** and one for **event trigger identification**.
- Each FFN produces a score vector for every word over the label set.
  - They use BIO tagging for both tasks.
  - The possible labels for entities are the entity types.
  - The possible labels for event triggers are the event types.
- A conditional random field (CRF) layer is added on top of each FFN to handle dependencies when predicting the sequence of labels.
- The predicted types are not actually used though, to avoid propagating errors. The decoding stage at the end will ultimately decide on the entity and event trigger types.

9

### Cross-Subtask Interactions

Interactions across different subtasks (recognizing entities, relations, events) can be valuable constraints! For example:

*A civilian aid worker from San Francisco was killed in an attack in Afghanistan.*



(a) Cross-subtask Interaction

*San Francisco* may be predicted to be a **Victim** because it preceded **was killed**. But a **GPE** is unlikely to be a **Victim**.

11

### Step 3: Local Classification of Nodes & Edges

- For each identified node (entity or event trigger), if it contains multiple words then an embedding is produced for the node by averaging the vectors for each word.
- Feed-forward neural networks are trained for each node-labeling task (entity types and event trigger types) to produce scores over each label set.
- Feed-forward neural networks are trained for each edge-labeling task (relations and event roles).
- An edge  $(v_i, v_j)$  is represented as the concatenation of the  $v_i$  and  $v_j$  span embeddings. Each edge has a start, end, and label type.
- For each task, the label with the highest scores is predicted to produce a locally best graph.

10

### Cross-instance Interactions

Interactions across different events or relations in the same sentence can be valuable constraints.

*South Carolina boy, 9, dies during hunting trip after his father accidentally shot him on Thanksgiving Day.*



(b) Cross-instance Interactions

Recognising that **boy** is the ATTACK Victim (trigger = **shot**) is tough because of the long distance between the words. But a DIE Victim is also likely to be the ATTACK Victim.

12

3/

### Global Feature Templates

#### Event Role Features

- The number of entities that act as  $\langle \text{role} \rangle$  > arguments at the same time.
- The number of event types, > events with  $\langle \text{entity}, \text{role} \rangle$  > arguments.
- The number of occurrences of  $\langle \text{event.type}, \text{role} \rangle$ , and  $\langle \text{entity.type}, \text{role} \rangle$  combination.
- The number of events that have multiple  $\langle \text{role} \rangle$  > arguments.
- The number of entities that act as  $\langle \text{role} \rangle$ , argument of an  $\langle \text{event.type} \rangle$  > event and a  $\langle \text{role} \rangle$ , argument of an  $\langle \text{event.type} \rangle$  > event at the same time.

#### Relation Features

- The number of occurrences of  $\langle \text{entity.type}_1, \text{entity.type}_2 \rangle$ , and  $\langle \text{relation.type}_1 \rangle$  combination.
- The number of occurrences of  $\langle \text{entity.type}_1, \text{entity.type}_2 \rangle$ , and  $\langle \text{relation.type}_2 \rangle$  combination.
- The number of occurrences of a  $\langle \text{relation.type}_1 \rangle$ , relation between a  $\langle \text{role} \rangle$ , argument and a  $\langle \text{role} \rangle$ , argument of the same event.
- The number of entities that have a  $\langle \text{relation.type}_1 \rangle$  relation with multiple entities.
- The number of entities involving in  $\langle \text{relation.type}_1 \rangle$  and  $\langle \text{relation.type}_2 \rangle$  relations simultaneously.

#### Trigger Features

- Whether a graph contains more than one  $\langle \text{event.type} \rangle$  > event.

The local score for G is:  $s(G) = s'(G) + u f_G$

### Using the Feature Templates

- Each feature template is instantiated with all possible values to generate a large feature set.
- Given a graph G, each function returns a scalar value. For example:  $f_1(G) = \begin{cases} 1, & G \text{ has multiple ATTACK events} \\ 0, & \text{otherwise.} \end{cases}$
- The neural network is tasked with learning a weight vector  $u$  for the features during training.
- The global score of G is the sum of its local score and global feature score:  $s(G) = s'(G) + u f_G$

14

### The Decoding Process

- The beam is initialized with 1 candidate, an empty graph:  $B = \{\emptyset\}$
- At each step, each candidate in B is expanded with a node and edge.
- Node Step:** a node  $v_i$  is selected (presumably with the highest-scoring label). Given a hyper-parameter  $\beta_v$ , the  $\beta_v$  best labels are used. The beam is updated with a copy of each candidate graph that has  $v_i$  added to it with one of the selected labels.
- Edge Step:** edges are added between  $v_i$  and previous nodes. Given a hyper-parameter  $\beta_e$ , the edges  $(v_i, v_j)$  with the highest label scores are selected. The beam is updated with a copy of each candidate graph that has one of the selected edges added to it.
- If  $|B| >$  beam width k, then the global score for each candidate is computed and only the top k candidates are kept.

16

### Step 4: Decoding the Globally Best Graph

- OneIE makes joint decisions for all nodes and edges to obtain a globally optimal graph.
- Ideally, we'd like to generate every possible candidate graph, calculate its global score and pick the best one. But exhaustive search is not feasible.
- SOLUTION:** decoding is done via **beam search**.
  - Beam search is a greedy algorithm for heuristic search.
  - When exploring a graph, the k best partial solutions are expanded at each step of the search process. The value k is called the **beam width**.

15

3/

### Decoding Illustration

*He also brought a check from Campbell to pay the fine and fees.*



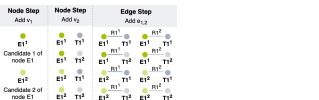
**Campbell** (E1) is selected as the first node to add, with 2 possible label candidates (E1<sup>1</sup> and E1<sup>2</sup>).

No edges are added because there are no previous nodes to connect with E1.

17

### Decoding Illustration

*He also brought a check from Campbell to pay the fine and fees.*



Edges are selected to link T1 with E1.

You can see here how the combinatorics can quickly create a large number of candidate graphs!

19

### Decoding Illustration

*He also brought a check from Campbell to pay the fine and fees.*

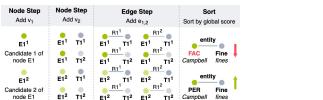


**fine** (T1) is selected as the second node to add, with 2 possible label candidates (T1<sup>1</sup> and T1<sup>2</sup>).

18

### Decoding Illustration

*He also brought a check from Campbell to pay the fine and fees.*



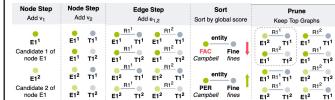
If the set of candidate graphs is greater than the beam width k, then the graphs are all scored and only the top k are kept in the beam.

20

3/

### Decoding Illustration

*He also brought a check from Campbell to pay the fine and fees.*



And then the process repeats, until all of the nodes have been added.

21

### Data Set Statistics

Dataset	Split	#Sentences	#Entities	#Roles	#Events
ACE05-R	Train	10,081	26,972	3,362	1,131
	Dev	1,424	5,362	1,131	-
	Test	2,050	5,476	1,131	-
ACE05-E	Train	17,172	29,006	4,664	4,202
	Dev	935	2,451	569	-
	Test	1,332	3,017	403	-
ACE05-CN	Train	6,841	29,657	7,934	2,926
	Dev	526	2,250	598	217
	Test	670	1,190	357	-
ACE05-E'	Train	19,340	47,525	7,152	4,419
	Dev	902	3,422	728	468
	Test	676	3,673	802	424
ER-E-EN	Train	14,219	38,864	5,045	6,419
	Dev	1,067	3,811	1,020	1,020
	Test	1,129	3,291	477	559
ER-E-S	Train	7,067	11,839	1,698	3,272
	Dev	556	886	128	210
	Test	546	811	108	209

23

### Data Sets

- ACE 2005: entities, relations, and events for English, Chinese, and Arabic.
  - 7 entity types, 6 (coarse) relation types, 33 event types, and 22 event roles.
  - ACE05-R only includes named entities and relations.
  - ACE05-E includes all entities, relations, and events.
  - ACE05-E+ includes ordered relation arguments, pronouns, and multi-token event triggers, which had been largely ignored in recent work!
- ERE-EN: Entities, Relations and Events (ERE) from a DEFT Program.
  - 7 entity types, 5 relation types, 38 event types, and 20 event roles.
  - 458 documents; 16,516 sentences.
  - They also created a Spanish version (ERE-ES)

22

### Evaluations

- Entity:** offsets and type must match the gold.
- Relation:** offsets and type must match the gold.
- Event Trigger:**
  - Trigger-I: offsets must match the gold (*identification only*)
  - Trigger-C: offsets and event type must match the gold
- Argument (Event Role):**
  - Arg-I: offsets must match the gold (*identification only*)
  - Arg-C: offsets and role type must match the gold

24

## Experimental Results

Dataset	Task	DyGE++	Baseline	OneIE
ACE05-R	Entity Relation	88.6	-	88.8
		63.4	63.9	67.5
		89.7	90.2	89.3
		76.6	78.2	78.2
ACE05-E	Trig-C	69.7	73.5	74.7
	Arg-I	53.8	56.8	59.2
	Arg-C	48.8	53.9	56.8

Scores

- DyGE++ is a previous state-of-the-art model.
- Baseline is their system but without the globally optimal graph decoding layer.

25

## Results for Chinese and Spanish

Task	Entity	Trig-I	Trig-C	Arg-I	Arg-C	Relation
ACE05-B		89.6	75.6	72.8	57.3	54.8
ERB-EN		87.0	68.4	57.0	50.1	46.5

Task	Entity	Trig-I	Trig-C	Arg-I	Arg-C	Relation
ACE05-CN		85.3	82.4	78.7	67.7	53.2
ERB-ES	ES	81.8	52.9	59.1	42.3	

Note that adding English training data to the Chinese and Spanish models improved their performance!

26

## Positive Learned Global Features

Positive Feature	Weight
1 A TRANSPORT event has only one DESTINATION argument	2.61
2 An ATTACK event has only one PLACE argument	2.31
3 A TRANSPORT event has only one ORIGIN argument	2.01
4 An END-POSITION event has only one PERSON argument	1.51
5 A PER-SOC relation exists between two PER entities	1.08
6 A GEN-AFF relation exists between ORG and LOC	0.96
7 A BERSERKARY argument is a PER entity	0.93
8 A GEN-AFF relation exists between ORG and GPE entities	0.90

27

## Negative Learned Global Features

Negative Feature	Weight
9 An entity has an ORG-AFF relation with multiple entities	-3.21
10 An entity has an PART-WHOLE relation with multiple entities	-2.49
11 An event has two PLACE arguments	-2.47
12 A TRANSPORT event has multiple DESTINATION arguments	-2.25
13 An entity has a GEN-AFF relation with multiple entities	-2.02
14 An ATTACK event has multiple PLACE arguments	-1.86
15 An entity has a PHYS relation with multiple entities	-1.69
16 An event has multiple VICTIM arguments	-1.61

28

## Feature Analysis #5

- #5: He also brought a check from **Campbell** to **pay** the fines and fees.
- Global feature category: 3
- Analysis: As “Campbell” is likely to be an ENTITY argument of a FINE event, the model corrects its entity type from FAC to PER.

Baseline



+ Global Features



33

## Error Analysis

They manually analyzed 75 errors to better understand where the remaining challenges lie:



34

## Error Category Examples

- Background Knowledge:** “And Putin’s media aide, Sergei Yastrzhembsky, told Kommersant Russia would not forgive the Iraqi debt”. *Kommersant* was labeled as a PERSON, but is an ORG. (It is a newspaper.)
- Rare Words:** the long-tail problem is a well-known challenge for NLP. For example, “*corretaker*” is a PERSON. Event triggers can be multi-word expressions that are rare, and even adverbs occasionally.
- Multiple Types per Trigger:** “*named*” can refer to both NOMINATE and START-POSITION events, and “*killed*” can refer to both ATTACK and DIE events (but usually only one is annotated).

35

## Error Category Examples

- Complex Syntactic Structure:** “As well as previously holding senior positions at *Barclays Bank*, *BZW* and *Kleinwort Benson*, McCarthy was formerly a top civil servant at the Department of Trade and Industry.” Their model missed all 3 employers: *Barclays Bank*, *BZW*, and *Kleinwort Benson*.
- Uncertain Events:** future planned events can be mistaken for events that have already happened, for example: “*The statement did not give any reason for the move, but said Lahoud would begin consultations Wednesday aimed at the formation of a new government*.”
- Metaphor:** “*Russia hints ‘peace camp’ alliance with Germany and France is dying*.”

36

## Feature Analysis #1

#1: Russia’s foreign minister expressed outrage at suggestions from a top Washington official last week that Moscow should forgive the eight billion dollars in Soviet-era debt that Baghdad owes it, as a gesture of good will.

- Global feature category: 8
- Analysis: It is unlikely for a person to have an ORG-AFF relation with multiple entities.

Baseline



+ Global Features



29

## Feature Analysis #2

#2: They also deployed along the border with Israel.

- Global feature category: 9
- Analysis: It is uncommon that an ORIGIN argument and a DESTINATION argument have a PART-WHOLE relation.

Baseline



+ Global Features



30

## Feature Analysis #3

#3: Prime Minister Abdullah Gul resigned earlier Tuesday to make way for Erdogan , who won a parliamentary seat in by-elections Sunday.

- Global feature categories: 2 and 5
- Analysis: 1. An ELECT usually has only one PERSON argument;
- 2. An entity is unlikely to act as a PERSON argument for END-POSITION and ELECT events at the same time.

Baseline



+ Global Features



31

## Feature Analysis #4

#4: Diller will continue to play a critical role in the future of Vivendi ’s entertainment arm.

- Global feature category: 6
- Analysis: A PART-WHOLE relation should not exist between PER and ORG entities.

Baseline



+ Global Features

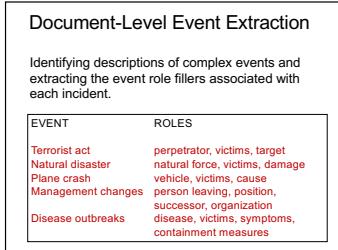


32

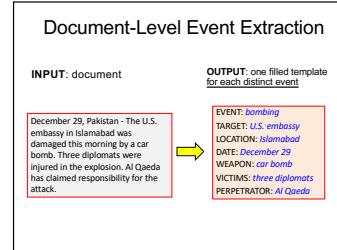
## Summary

- The OneIE model is a framework that simultaneously learns multiple IE subtasks and incorporates global features to recognize subtask and instance interactions that are desirable or undesirable.
  - A nice combination of feature engineering with neural nets.
  - These types of features are clearly beneficial!
- However, this work is still only tackling *sentence-level event extraction*! And the results are still far from perfect.
- The challenges grow even greater when extracting event information from an entire document, where the information can be scattered about and a coherent event representation must be produced.

37



1



2

### Event Template for Disease Outbreaks

```

Story: <document id>
ID: <template id>
Date: <date>
Event: OUTBREAK
Status: <set fill>
Containment: <set fill>
Country: <set fill>
Victims: <string list>
Disease: <string>

```

5

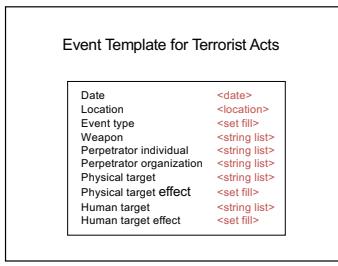
### Filled Event Template for Disease Outbreaks

```

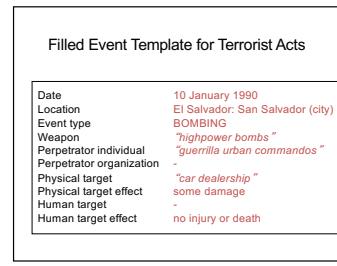
Story: 20020714.4756
ID: 1
Date: August 14, 2002
Event: OUTBREAK
Status: confirmed
Containment: none
Country: Switzerland
Victims: the 27 reported cases
Disease: Creutzfeldt-Jakob Disease / [sporadic] Creutzfeldt-Jakob disease (CJD) / CJD / Sporadic CJD / hereditary dominant CJD / Swiss CJD / sporadic Creutzfeldt-Jakob disease NOTE: disjunctive options!

```

6



3



4

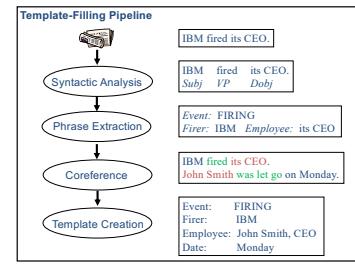
### Patterns/Rules vs. Sequence Tagging

Two general approaches to event extraction:

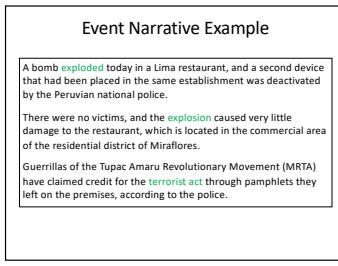
**Pattern-based systems** use patterns or rules which identify phrases that should be extracted for each event role.

**Machine learning classifiers** label individual tokens indicating whether they should be extracted, and if so, what role they play.

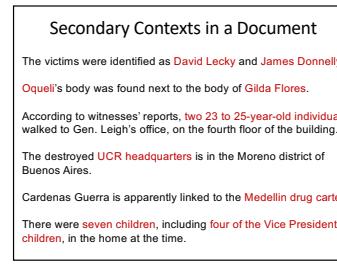
7



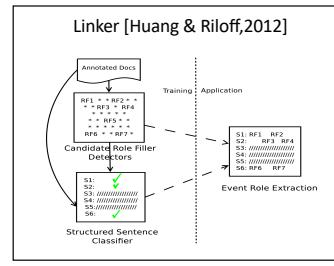
8



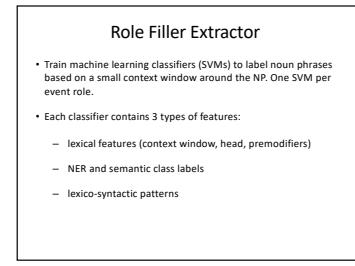
9



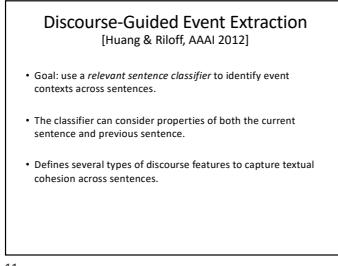
10



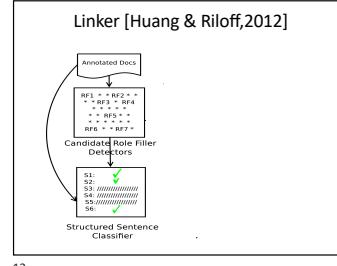
13



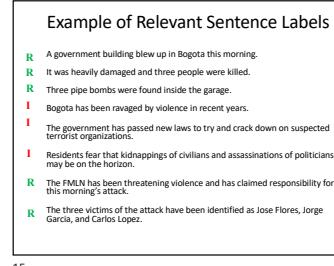
14



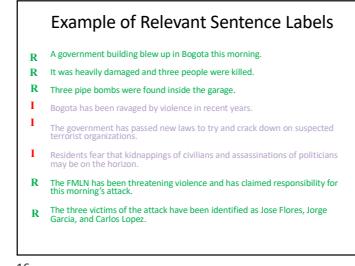
11



12

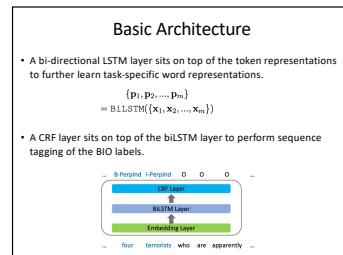


15

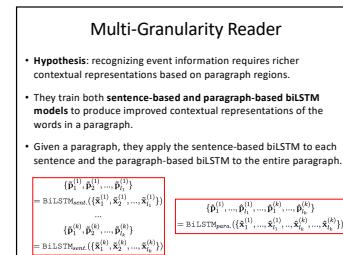


16

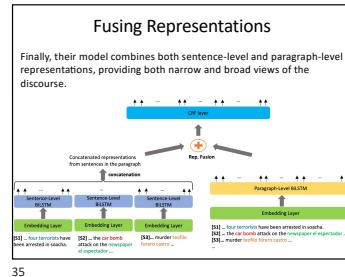




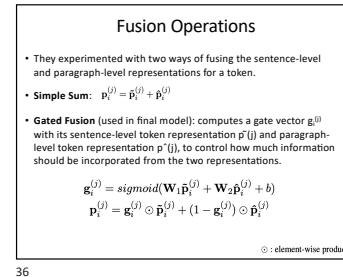
33



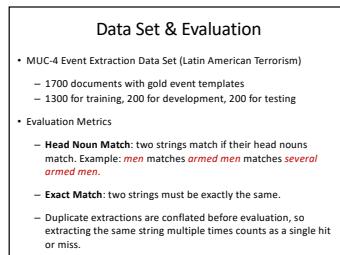
34



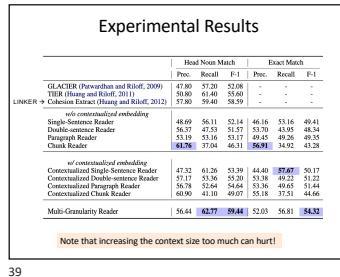
35



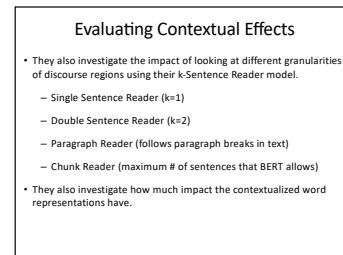
36



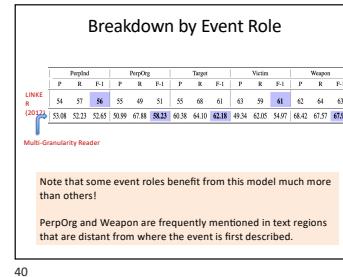
37



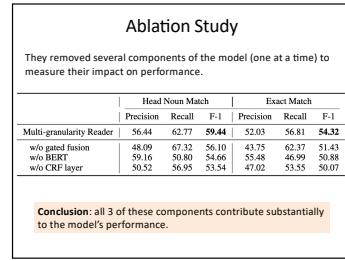
39



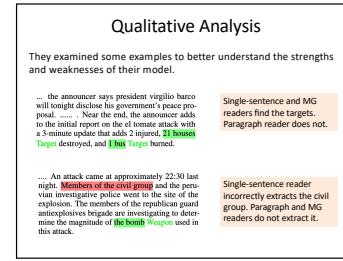
38



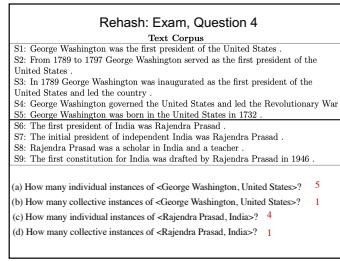
40



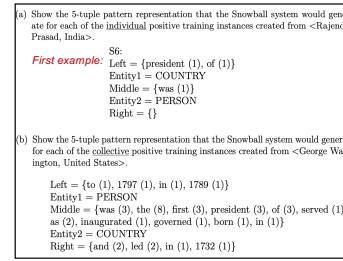
41



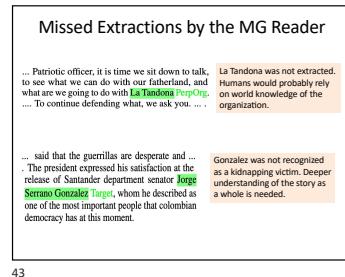
42



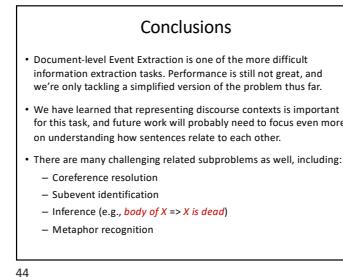
1



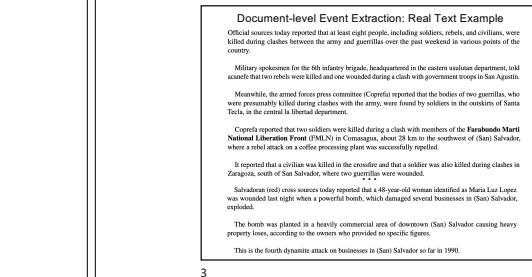
2



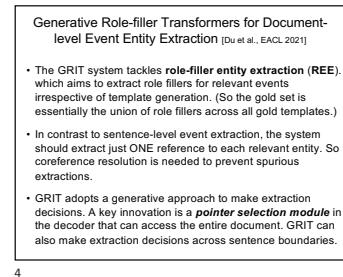
43



44



3



4

## Event Extraction Example

Role		Role-filler Entities	
Perpetrator Individual		two men, two men wearing sports clothes, Shining Path members	
Perpetrator Organization		Shining Path	
Physical Target		water pipes, water pipes	
Weapon		Pimal telephone company building, telephone company building, telephone company offices	
Victim		public telephone booth	
		125 to 150 grams of Tnt	
		Water pipes	

A bomb exploded in a Pimal city destroying some [water pipes]. According to unconfirmed reports, the bomb contained [125 to 150 grams of Tnt] and was set on the back of the [Pimal telephone company building]. The explosion occurred at 2350 on 16 January, causing panic but no casualties. The explosion caused damage to the [telephone company building] and [telephone company offices] and [water pipes]. Witnesses reported that the bomb was planted by [two men] wearing sports clothes, who escaped into the night. They were later identified as [Shining Path] members; Their focus is on entity-level decisions!

5

## A Sequence-to-Sequence Task Model

- The source sequence consists of the tokens from the input document, with a [CLS] token prepended and a [SEP] token appended.
- The target sequence is the concatenation of the extractions for each event role, with a [SEP] token between each one.
- Each role filler entity is represented by the beginning (b) and end (e) tokens for its first mention. Multiple fillers for the same slot seem to be appended without a separator between them.
- The roles are always presented in a fixed order, and the model learns to generate them in this same order.

7

## The GRIT Model

- GRIT is built upon the pre-trained transformer model BERT.
- Role filler extraction is treated as a sequence-to-sequence task, with an encoder-decoder framework.
  - The encoder generates a representation for the input ("source" sequence).
  - The decoder produces role filler extractions as its output ("target" sequence).
- A single pre-trained transformer model is used for both, without fine-tuning.

6

## Notation Examples

Formatting for target sequence:

$$\langle\!\langle S \rangle\!\rangle = e_1^{(1)}, e_2^{(1)}, \dots, [SEP] \\ e_1^{(2)}, e_2^{(2)}, \dots, [SEP] \\ e_1^{(3)}, e_2^{(3)}, \dots, e_3^{(3)}, [SEP]$$

Perpetrator Individual	two men, two men wearing sports clothes, Shining Path members
Perpetrator Organization	Shining Path
Physical Target	water pipes, water pipes
Weapon	Pimal telephone company building, telephone company building, telephone company offices
Victim	public telephone booth

125 to 150 grams of Tnt [SEP]

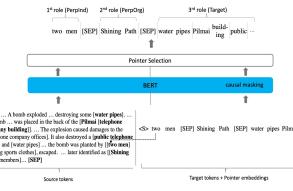
8

## Pointer Embeddings

- GRIT captures "pointers" that allow the decoder to refer back to places in the source document where the extractions came from.
- As the input to BERT, they use the sum of token, position, and segment embeddings ("segment" is essentially the sentence number).
- The decoder generates a representation for the input ("source" sequence).
- For the target representation, the position for a word corresponds to its position in the source document.
- BERT is given the source document as sequence A and the target information as sequence B.

9

## GRIT Example



11

## Causal Mask for Target Embedding

- Since they use a single BERT model for both encoding and decoding, they incorporate causal masks to ensure realistic self-attention.
- When embedding the source document, the target tokens are masked.
- When embedding the target information, the mask only allows self-attention to prior tokens (to make the model autoregressive).

10

## Pointer Decoding

- The next layer is *pointer decoding*, which selects "pointers" to the source document rather than predicting words.
- When generating the target information at time t, they compute the dot product between the target information at time t-1 and each position in the source document.
- Softmax is then applied to get the probability of pointing to each source document term. Greedy decoding selects the source token with the highest probability.
- The predicted token is added to the target information for the next time step.
- Decoding stops when the last [SEP] token (representing the last event role) is generated!

12

## Comparative Results

GRIT produced better F1 scores than previous models, primarily due to much higher precision.

Models	P	R	F1
LINKER	58.38	39.53	47.14
CohesionExtract (Huang and Riloff, 2012)	56.82	48.92	52.58
NST (Du and Cardie, 2020)	57.04	46.77	51.40
DyGIE++ (Wadden et al., 2019)	64.19**	47.36	54.50*

13

## Performance on documents with nested role fillers across Perplnd & PerpOrg roles

PerpOrg (all docs)	PerpOrg (33/200)	
	P / R / F1	P / R / F1
NST	56.00 / 34.35 / 42.42	80.00 / 44.44 / 57.14
DyGIE++	60.00 / 43.00 / 50.70	61.54 / 35.56 / 45.07
GRIT	66.04 / 42.65 / 51.85	80.77 / 46.67 / 59.15

GRIT seems to help when there are dependencies across roles.

Gold Role-filler Entities	Perpetrator		Physical Target	
	• guerrillas	• FARC	• guerrillas	• FARC
NST & DyGIE++	• guerrillas, guerrillas and FARC	• FARC	-	-
GRIT	• guerrillas	• FARC	• guerrillas	• FARC

15

## A Deeper Dive into the Results

Results are Precision/Recall/F1 using an entity-based metric that they defined, called CEAFF-REE.

	PerpOrg	PerpOrg	Target	Victim	Weapon
K=1	1 < K < 2.5	3.25 < K < 1.5	1.5 < K < 1.75	K > 1.75	
NST	48.39 / 37.01 / 38.38	60.00 / 44.96 / 50.70	54.96 / 34.51 / 53.88	62.50 / 31.16 / 62.82	61.67 / 45.67 / 61.67
DyGIE++	59.49 / 34.09 / 43.21	56.00 / 34.35 / 42.42	53.49 / 30.74 / 32.08	60.00 / 46.52 / 60.00	57.14 / 53.33 / 55.17
GRIT	65.85 / 39.80 / 49.53	66.04 / 42.68 / 51.85	59.05 / 44.12 / 48.91	76.32 / 61.05 / 87.41	61.82 / 56.67 / 79.15

They also examined subsets of documents with increasing numbers of co-referent mentions per event role.

K=1

1 < K < 2.5

3.25 < K < 1.5

1.5 < K < 1.75

K > 1.75

NST

DyGIE++

GRIT

72.59 / 50.00 / 59.18

70.00 / 40.00 / 50.81

60.48 / 48.39 / 53.78

52.94 / 38.57 / 44.61

66.96 / 48.73 / 56.41

65.85 / 46.55 / 54.55

74.42 / 45.71 / 56.64

73.20 / 55.81 / 56.35

67.44 / 41.43 / 51.35

66.73 / 52.53 / 59.95

## GRIT is also more efficient

- Previous models require an additional classification layer on top of BERT.
- This means that many more parameters must be learned, which substantially increases training time.

	additional params	training cost
DyGIE++	2H(#roles + 1)	~20h
NST	H(2#roles + 1)	~1h
GRIT	0	<40min

17

## Template Filling Example

Event 1 Template	Attack
Perpetrator Indx.	• Zarate armed forces
Perpetrator Org.	• Zarate armed forces
Perpetrator Indx.	• Zarate armed forces
Physical Target	• Zarate armed forces
Weapon	• Zarate armed forces
Victim	-

Event 2 Template	Bombing
Perpetrator Indx.	• Zarate armed forces
Perpetrator Org.	• Zarate armed forces
Perpetrator Indx.	• Zarate armed forces
Physical Target	• Zarate armed forces
Weapon	-
Victim	-

Event 3 Template	Arson
Perpetrator Indx.	• Zarate armed forces
Perpetrator Org.	• Zarate armed forces
Perpetrator Indx.	• Zarate armed forces
Physical Target	• old church
Weapon	-
Victim	-

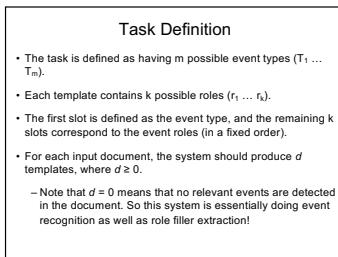
19

## Template Filling with Generative Transformers

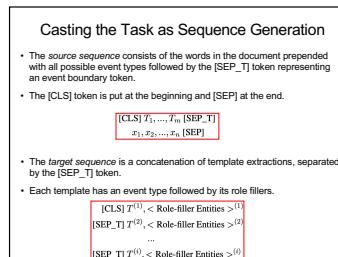
[Du et al., NAACL 2021]

- Du et al. recently revisited the challenge of automatically creating and filling distinct event templates.
- Their model (GTT) tackles the problem of event recognition and role filler extraction at the same time using generative transformers.
- This model is built upon their earlier GRIT system.
  - The GRIT model was extended to also predict event types.
  - They modify the decoder to attend to the event types.

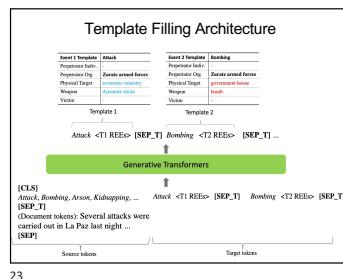
20



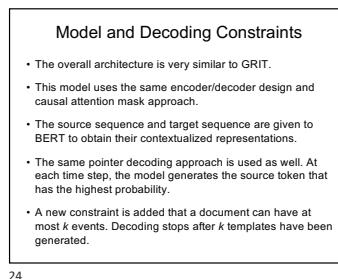
21



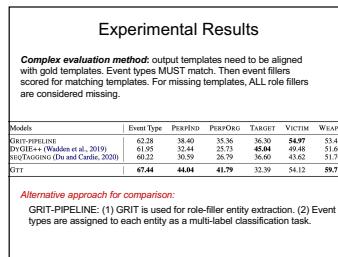
22



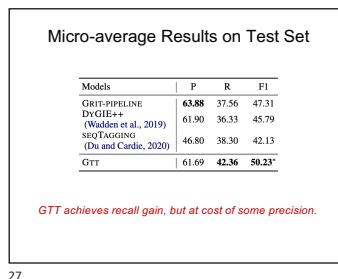
23



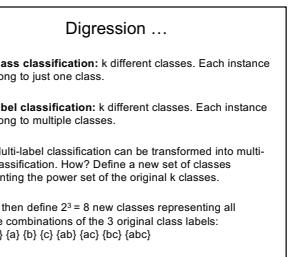
24



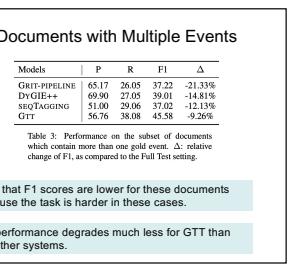
25



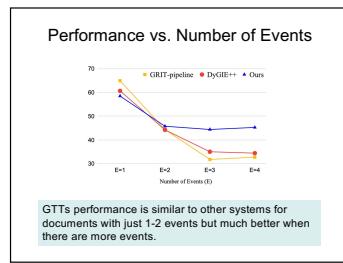
27



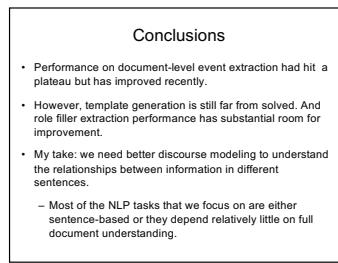
26



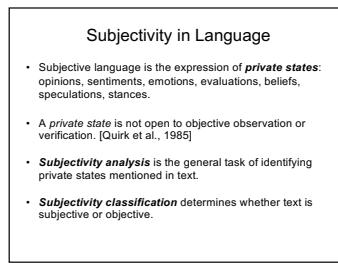
28



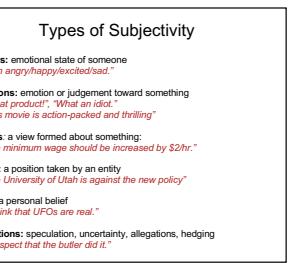
29



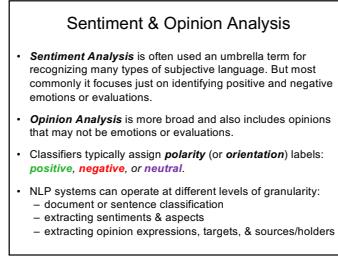
30



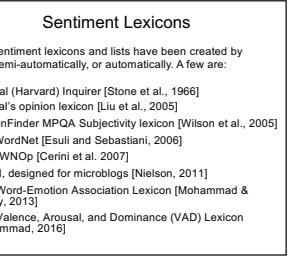
1



2



3



4

### Semi-Supervised Induction of Sentiment Lexicons

- Sentiment lexicons are often produced with semi-supervised learning methods that begin with "seed" words and then apply heuristics to learn the polarity of new words.
- One of the earliest works by Hatzivassiloglou and McKeown (1997) exploited the observation that conjoined adjectives usually have the same polarity.
- Ex: "*happy and excited*", "*sad and depressed*"
- A general class of methods often used for this problem are **label propagation** algorithms. Terms are encoded in a graph structure with edges capturing their similarity and sentiment values are iteratively refined.

5

### Examples

Domain	Positive seed words	Negative seed words
Standard English	good, lovely, excellent, awesome, pleasant, successful, kind, love, happy	bad, horrid, terrible, annoying, unkind, bad, bad, bad, bad, bad, bad, unhappy
Finance	success, excellent, profit, beneficial, improving, improved, success, gains, positive	negligent, loss, volatile, wrong, losses, damages, bad, litigation, failure, down, negative
Twitter	love, loved, loves, awesome, nice, amazing, best, fantastic, correct, happy	hat, hated, hates, terrible, nasty, awful, worse, horrible, wrong, sad

a. Run random walks from seed words.  
b. Assign polarity scores based on frequency of nouns with links.

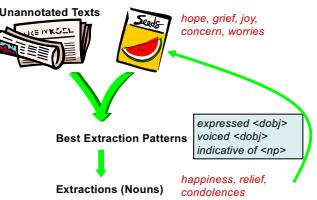
7

### Sentiment Propagation Algorithm

1. Build a graph with vertices representing the words in your corpus.
2. For each vertex  $v$ , add an edge between  $v$  and each of its  $k$  nearest neighbors, based on cosine similarity of the word embedding vectors. The cosine similarity score is also used as the edge weight.
3. Perform "random walks" starting with the seed word nodes. The next node  $w$  is chosen with probability proportional to the edge weights from  $v$  to  $w$ .
4. Positive and negative scores can then be computed for each node based on the likelihood of reaching it from a positive or negative seed word.
5. To be more robust over different seed sets, confidence estimates were produced by repeated experiments with different subsets of the seed words.

6

### Bootstrapped Learning of Subjective Nouns and Patterns



8

### Learning Subjective Expressions

- |                    |  |
|--------------------|--|
| expressed <obj>    | condolences, hope, grief, views, worries |
| indicative of <np> | compromise, desire, thinking             |
| inject <obj>       | vitality, hatred                         |
| reaffirmed <obj>   | resolve, position, commitment            |
| voiced <obj>       | outrage, support, skepticism,            |
| show of <np>       | opposition, gratitude, indignation       |
| <subj> was shared  | support, strength, goodwill, solidarity  |

9

### Examples of Strong Subjective Nouns

aberration	eyebrows	resistant
allusion	failures	risk
apprehensions	inclination	sincerity
assault	intrigue	slump
beneficiary	liability	spirit
benefit	likelihood	success
blood	peaceful	tolerance
controversy	persistent	trick
credo	play	trust
distortion	pressure	unity
drama	promise	
eternity	rejection	

11

### Examples of Strong Subjective Nouns

anguish	exploitation	pariah
antagonism	evil	repudiation
apologist	fallacies	revenge
atrocities	genius	rogue
barbarian	goodwill	sanctimonious
belligerence	humiliation	scum
bully	ill-treatment	smokescreen
condemnation	injustice	sympathy
denunciation	innuendo	tyranny
devil	insinuation	venom
diatribe	ilar	
exaggeration	mockery	

10

### Contextual Polarization

- Sentiment lexicons capture the **prior polarity** of words and phrases.
- However, the polarity of a word often depends on context due to polysemy, negation, polarity shifters, scoping, expressions, etc.

Example from [Wilson, Wiebe, &amp; Hoffmann 2005]:

Philip Clapp, president of the National Environment Trust, sums up well the general thrust of the reaction of environmental movements: "There is no reason at all to believe that the polluters are suddenly going to become reasonable."

12

### Negation and Polarity Shifters

- People often try simple approaches to sentiment classification that simply count the number of positive vs. negative words found in a lexicon.
- But this approach is often inaccurate because it ignores negation and polarity shifters, which TOGGLE polarity.

```
There was a lot of damage. → NEG
There was no damage. → POS
There was little damage. → POS

She showed much empathy. → POS
She never showed empathy. → NEG
She showed a lack of empathy. → NEG
```

13

### Opinion Extraction

Opinion extraction systems typically aim to recognize three components of opinions:

Opinion Expression: phrase that describes an attitude toward or evaluation of something

Opinion Holder (Source): the entity whose opinion is being expressed

Opinion Target: the entity, object, or concept that the opinion is about (toward)

According to UN officials, the human rights record in Syria is horrendous.

14

### Opinion Holders and Targets

- Opinion holders are typically either the *speaker/writer* (implicitly) or a Person or Organization entity (but could also be something like a report).
- Opinion targets can be almost anything! So identifying the boundaries is difficult. For example:
  - I dislike John.
  - I dislike the uniforms of the Utah Jazz.
  - I dislike the recent actions of the government that raised tariffs on tens of thousands of German products, which resulted in ....

15

### Extracting Opinion Propositions and Holders

[Bethard et al., 2004] developed one of the earliest systems to identify propositional opinions and the opinion holders (sources).

• **Opinion:** answer to the question "How does X feel about Y"

• **Propositional Opinion:** an opinion localized in an argument of a verb, generally a sentential complement.

• **Opinion Holder:** the entity who holds the opinion

For example:

- She believes [you have to use the system to change it].
- Still, Vista officials realize [they're relatively fortunate].
- I'd be destroying myself] replies Mr. Korotich.

16

### Extracting Opinion Propositions and Holders

[Bethard et al., 2004] developed one of the earliest systems to identify propositional opinions and the opinion holders (sources).

• **Opinion:** answer to the question "How does X feel about Y"

• **Propositional Opinion:** an opinion localized in an argument of a verb, generally a sentential complement.

• **Opinion Holder:** the entity who holds the opinion

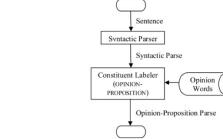
For example:

- She believes [you have to use the system to change it].
- Still, Vista officials realize [they're relatively fortunate].
- I'd be destroying myself] replies Mr. Korotich.

17

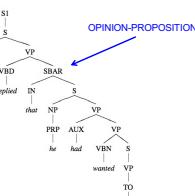
### One-Tiered Architecture

The first system is an SVM classifier that labels syntactic constituents in a parse tree as either OPINION-PROPOSITION or NULL.



19

### Example



20

### Opinion-Proposition Classifier

They followed the same design as a semantic role labeling classifier by [Pradhan et al., 2003] with 8 syntactic features:

1. the verb
2. verb's cluster
3. subcategorization type of the verb
4. syntactic phrase type of the potential argument
5. head word of the potential argument
6. before/after position of the argument relative to the verb
7. parse tree path between verb and potential argument
8. voice (active/passive) of the verb

This feature set was later augmented with features derived from an acquired set of opinion words.

21

### Opinion Word Features

Given a constituent to classify, the following features captured opinion word information:

- Counts:** the number of opinion words in the constituent.
- Score Sum:** the sum of the opinion scores for each opinion word in the constituent, sometimes with a minimum score threshold.

- ADJP:** a binary feature indicating whether the constituent contains a complex adjective phrase. (Simple adjectives produce many false hits.) For example:

excessively affluent  
more bureaucratic

[Note: I've observed "ADV ADJ" to be a useful pattern tool!]

22

### Results for Two-Tiered Architecture

The first component that labels PROPOSITION constituents achieved 62% recall with 82% precision. (This was a 10% precision gain over the more general semantic role classifier.)

The results for the 3 models to determine which PROPOSITION constituents are opinions are shown below:

Train on	Predict on	Metric	Weights	Ngrams	Features	Prec	Recall	Orientatio
Sentence	Sentence	Recall	33%	29.9%	100%	62.7%	67.4%	72%
Sentence	Sentence	Precision	67.84%	61.13%	62.50%	65.55%	67.47%	
Sentence	Proposition	Recall	37.48%	37.32%	37.79%	36.03%	28.81%	
Sentence	Proposition	Precision	53.95%	59.00%	57.91%	55.07%	68.41%	
Proposition	Proposition	Recall	39.73%	37.84%	37.52%	37.53%	37.53%	
Proposition	Proposition	Precision	59.56%	61.63%	60.43%	58.77%	61.66%	

Table 5: Two-tiered Approach Results for Opinion Propositions.

25

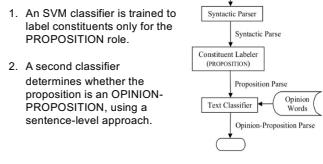
### Results and Conclusions

- This work focused on one type of opinion recognition, **propositional opinions**, and identified the opinion holders.
- This approach is very syntactically-oriented, requiring an alignment between the propositions/holders and syntactic constituents.
- This approach cannot identify cases where a proposition spans multiple sentences, or the holder is in a different sentence than the proposition.
- The two architectures exhibited a recall/precision trade-off:  
51% R with 58% P for 1 Tiered  
43% R with 68% P for 2 Tiered

26

### Two-Tiered Architecture

The second system performs two steps:



23

### Labeling Propositions as Opinions

Three Naive Bayes classifiers were trained to determine whether a proposition is an OPINION-PROPOSITION.

1. The first model is trained using approximate sentence labels from fact-heavy vs. opinion-heavy texts.  
Sentences in editorials and letters to the editor are assumed to contain opinions.  
Sentences in news and business articles are assumed to be factual.  
The sentence containing each proposition is classified and the proposition is assigned the label of its sentence.
2. The second model is trained at the sentence-level but predictions are based only on the text of the proposition.
3. For the third model, both training and testing use only the text of the propositions (with the same approximate labeling during training).

24

### Opinion Role Extraction

- [Wiegand & Ruppenhofer, CoNLL 2015] observed that **opinion roles** (holders and targets) of **opinion verbs** are often aligned with the verb's semantic roles, but the mapping depends on the verb.

Peter	criticized	Mary	Holder	Target
Peter	disappoints	Mary	Target	Holder

- Also, many opinion verbs implicitly express the sentiment of the speaker, so the opinion holder does not appear. For example:

At my work, they are constantly gossiping.

27

### Verb Categorization for Opinion Roles

- [Wiegand & Ruppenhofer, CoNLL 2015] propose that verbs can be grouped into 3 categories that share **opinion role subcategorization** frames.
- Thus, categorizing a verb into one of these 3 classes essentially dictates how to extract its opinion roles!
- There are often multiple ways to interpret an opinion, partly due to polysemy, but for the sake of simplicity they only seek to learn the most common interpretation for a verb.

28

### Verbs with Agent View (AG)

Verbs with an Agent View convey the opinion of its agent, so the agent is the opinion holder and the "patient" is the opinion target.

[Here, "patient" refers to a semantic role capturing the object of the verb's action, such as what was criticized or argued. Similar to its theme.]

Peter criticized Mary.  
They like the exam.  
The guests complained about the noise.  
They argue that this plan is infeasible.

Note that the patient of the verb can appear in different types of syntactic constituents, including NPs, PPs, and infinitive or complement phrases!

29

### Verbs with Patient View (PT)

Verbs with a Patient View are the opposite, and have the agent as the opinion target and the patient as the opinion holder.

Peter disappoints Mary.  
The noise irritated the guests.  
The gift pleased her very much.  
The policy raising tariffs on German products upset Tom.

30

### Verbs with Speaker View (SP)

Verbs with a Speaker View express an opinion that comes from the speaker/writer (implicitly), rather than from an entity involved in the action. The agent is usually the target.

At my work, they are constantly gossiping.  
They cheated in the exam.  
The young baseball player has improved.  
He besmirched the King's name.

Occasionally, the patient of Speaker View verbs can also be a target, as for "besmirched" above. But since this is less common, they only assign the agent to be a target.

31

### Opinion Verbs & Gold Standard

- This work aims to classify an existing set of opinion verbs to extract the holders and targets from their arguments.
- They use the 1,175 opinion verb lemmas in the Subjectivity Lexicon [Wilson et al., 2005].
- They annotated their semantic roles to be consistent with PropBank and assigned holder, target, and speaker labels.
- They measured inter-annotator agreement on 400 verbs:  
opinion holders:  $\kappa = 60.8$   
opinion targets:  $\kappa = 62.3$   
speaker views:  $\kappa = 59.9$

32

### Distribution of Verb Types in Gold Standard Data

Agent (AG)	Patent (PT)	Speaker (SP)			
Freq	Percent	Freq	Percent	Freq	Percent
450	38.3	188	16.0	537	45.7

33

### Pattern-based Seed Initialization

- To induce seed AG verbs: identify verbs that co-occur with **prototypical opinion holders** (e.g., opponents and critics). For example:  
*Opponents claim these arguments miss the point.*  
*Critics argue that the proposed limits were unconstitutional.*
- To induce seed PT verbs: identify **morphologically related adjectives**, which then reveal a PT verb. Specifically, they identify verbs in past-participle form that are identical to a predicate adjective. For example:  
*He has upset me.*    *I am upset.*
- To induce seed SP verbs: extract verbs using 3 patterns: (1) *accused of X<sub>agg</sub>*, (2) *blamed for X<sub>agg</sub>*, and (3) *help to X<sub>agg</sub>*. For example:  
*He was accused of falsifying the documents.*  
*The UN was blamed for misinterpreting climate data.*

34

### Extracted Verb Seeds

They used the North American News Text Corpus for seed extraction and comparison of verb similarities.

The top 12 extracted verb seeds for each category were:

AG	argue, contend, speculate, fear, doubt, complain, consider, praise, recommend, view, acknowledge, hope
PT	interest, surprise, please, excite, disappoint, delight, impress, shock, trouble, embarrass, annoy, distress
SP	murder, plot, incite, blaspheme, bewitch, bungle, despoil, plagiarize, privatize, instigate, molest, conspire

Ultimately they used the top 40 AG verbs, 30 PT verbs, and 50 SP verbs.

35

### Similarity Metrics

1. **Word Embeddings:** get embedding vectors for opinion words using word2vec, and compare with cosine similarity.
  2. **WordNet:** used a similarity metric based on WordNet's graph structure (WordNet:Similarity).
  3. **Coordination:** coordination typically preserves semantic coherence, so they apply a dependency parser to detect verb coordination (e.g., *They criticize and hate him*).
- The similarity score is just the frequency of two verbs appearing in a conjunction.

36

### Digression: Pointwise Mutual Information

**Pointwise mutual information (PMI)** measures the degree to which two words are statistically dependent.

$$\text{PMI}(w_1, w_2) = \log_2 \frac{P(w_1 \& w_2)}{P(w_1) * P(w_2)}$$

Mutual information can be used to measure the strength of association between a relation and its slot fillers:

$$\text{MI}(o, \text{Slot}, b) = \log_2 \frac{|\{o, \text{Slot}, b \mid X \mid *, \text{Slot}, *\}|}{|\{o, \text{Slot}, * \mid X \mid *, \text{Slot}, b\}|}$$

37

### MI Example for Relation & Filler

$$\text{MI}(\text{will buy}, \text{SlotX}, \text{Smith}) =$$

$$\log_2 \left[ \frac{|\{\text{will buy}, \text{SlotX}, \text{Smith} \mid X \mid *, \text{SlotX}, *\}|}{|\{\text{will buy}, \text{SlotX}, * \mid X \mid *, \text{SlotX}, \text{Smith}\}|} \right]$$

where:

$$\begin{aligned} |\{\text{will buy}, \text{SlotX}, \text{Smith} \mid X \mid *, \text{SlotX}, *\}| &= \# \text{ times } \text{Smith} \text{ fills SlotX for path will buy} \\ |\{X\mid, \text{SlotX}, *\}| &= \text{total # fillers for SlotX in all paths} \\ |\{\text{will buy}, \text{SlotX}, *\}| &= \# \text{ fillers in SlotX for path will buy} \\ |\{X\mid, \text{SlotX}, \text{Smith}\}| &= \# \text{ times } \text{Smith} \text{ fills SlotX in all paths} \end{aligned}$$

38

### Dependency-based Similarity

**4. Distributional similarity based on dependency relations:**  $(x, r, y)$  where  $x$  and  $y$  are words and  $r$  is a relation.

Example: (argue-V, nsubj, critics-N)

$$\text{sim}(v_1, v_2) = \frac{\sum_{(r,w) \in T(v_1) \cap T(v_2)} \text{MI}(v_1, r, w) + \text{MI}(v_2, r, w)}{\sum_{(r,w) \in T(v_1)} \text{MI}(v_1, r, w) + \sum_{(r,w) \in T(v_2)} \text{MI}(v_2, r, w)}$$

where  $T(x)$  is the set of pairs  $(r, y)$  such that

$$\log_2 \left[ \frac{|\{x, r, y \mid X \mid *, r, *\}|}{|\{x, r, * \mid X \mid *, r, y\}|} \right] > 0$$

From DIRT = [Lin & Pantel, 2001]

39

### Semi-Supervised Learning: Propagation Methods

They tried 2 methods to propagate labels from the seed verbs to new verbs.

1. **K nearest neighbor classification:** kNN methods identify the  $k$  closest (most similar/related) neighbors to input  $X$  and typically assign the predominant label to  $X$ .
2. **Label Propagation:** label propagation is a class of semi-supervised learning methods that iteratively propagate labels across nodes in a graph based on edges with weights that capture similarity/relatedness.

40

### Experimental Results

	Acc	Prec	Rec	F1
<b>Baselines</b>				
Majority Class Only Seeds	45.7	14.2	33.3	20.9
<b>Coordination</b>				
KNN graph	8.9	9.7	9.3	17.6
42.7	68.7	39.7	50.4	
<b>WordNet</b>				
kNN graph	52.8	51.5	50.7	51.1
51.1	51.9	51.5	51.5	
<b>Embedding</b>				
kNN graph	59.3	58.4	61.0	59.7
64.0	70.5	59.4	64.5	
<b>Dependency</b>				
kNN graph	65.7	63.8	65.4	64.5
70.3	72.0	68.0	70.6	

For kNN models,  $k=3$ .

41

### Multilingual Results

They also replicated their system to handle German! This illustrates the generality of their approach.

Overall, the same conclusions held.

	Major.	Coordin.	WordNet	Embedd.	Depend.
English	20.9	53.4	50.4	51.1	51.5
German	22.9	43.8	48.9	53.2	59.9

43

### Inspecting the Similarity Metrics

Inspecting the most similar words identified by each metric reveals that distributional similarity based on dependency relations indeed looked the best.

The 12 verbs most similar to *outrage*, which is PT. The underlined words are not PT verbs.

<b>Coordin.</b>	argue, believe, bring, vow, want, offend, shock, help, experience, strengthen, nominate, distract
<b>WordNet</b>	spoil, scandalize, anger, rage, sicken, temper, hate, fear, love, alarm, dread, tingle
<b>Embedd.</b>	anger, disgust, disgust, protest, alarm, enraged, shock, regret, concern, horrify, appal, sorrow
<b>Depend.</b>	anger, infuriate, alarm, shock, stun, enraged, incense, dismay, upset, appal, offend, disappoint

Coordination was probably poor due to sparsity.

42

### In-Context Classification

- Next, they built a classifier to use their induced verb knowledge for extracting opinion information in sentence contexts.
- They sampled ~1100 sentences from the N. American News Corpus in which their opinion verbs occurred.
- They manually annotated the opinion holders, targets, and targets evoked by speaker views. The data has: 753 holders, 745 targets, and 499 speaker view targets.
- They trained 3 SVM classifiers, one for each type of opinion information.

44

### Feature Set for Contextual Classifier

Features	Descriptions
cond\_is\_root	head lemma of candidate (phrase)
cond\_pos	part-of-speech tag of head of candidate phrase
cond\_base	candidate base form
cond\_person	is candidate a person
verb\_lemma	verb lemma
verb\_tense	verb tense
verb\_subj	verb subject
verb\_obj	verb object
word	bag of words within the sentence
pos	part-of-speech sequence between cond and verb
distance	distance between cond and verb
const	path from sentence parse tree from cond to verb
relat	vectorization of frame of verb
obj\_framegraph	semantic role annotations between cond and verb (semantic roles based on PropBank)
brown	Brown clusters of cond word/verb
obj\_parent	Brown cluster of word/verb's parent
frame\_lex	frame element lexicon (and the frame name to which frame element belongs)
fine\_grain\_lex	is fine-grained lexicon
course\_grain\_lex	is candidate holder/target/participant according to course-grain lexicon
inde\_graph	is candidate holder/target/participant according to the course-grained lexicon automatically induced with graph clustering (and induced word (41))

45

### Contextual Classification Results

Evaluation results using 10-fold cross-validation:

Features	Holder	Target	Target_Speaker
standard	63.59	54.18	40.00
+std <sub>framenet</sub>	65.44*	55.70*	42.16
+induc <sub>graph</sub>	68.06**	59.61**	46.66**
+std <sub>framenet</sub> +induc <sub>graph</sub>	69.70**	60.47**	47.33**
+induc <sub>graph</sub> .lex	68.56**	59.89**	54.31**
+std <sub>framenet</sub> .lex	69.70**	60.68**	54.06**
+induc <sub>graph</sub> .lex	69.83**	62.89**	56.71**
+std <sub>framenet</sub> .+fine-grain.lex	70.80**	63.72**	56.64**

46

### Conclusions

- Opinion extraction can be decomposed into subproblems: identifying opinion expressions, and their holders & targets.
- Aligning opinion targets of verbs with semantic roles makes sense linguistically and constrains the boundaries. This work offered a new insight that different classes of verbs behave differently in terms of their opinion roles.
- However, the results are still far from perfect. Opinions vary greatly in scope.
- And not all opinions are expressed as verbs! For other cases, opinion target boundaries can be elusive.

47

### Opinion Roles and Semantic Roles

Prior research observed that opinion holders and opinion targets often align with the arguments of verbs (or verb phrases) that express an opinion. For example:

(1) Australia said [it]<sub>H</sub> **fear[ed]** [violence]<sub>T</sub> if voters thought the election had been stolen.

Nouns can also take syntactic arguments, and opinion role fillers can also be found in their arguments too.

*The disgust of citizens toward the policy was apparent at the rally.*

holder	target
--------	--------

Research question: can semantic role labeling help with opinion extraction?

1

### PropBank Definitions

PropBank provides *Frames files* which defines a set of roles (*roleset*) for verb senses from VerbNet. There are two types of roles: numbered arguments and adjuncts.

#### Numbered Arguments: A0-A5

- Arg0 usually refers to the verb's agent.
- Arg1 usually refers to the verb's patient/theme (if it has one)
- All other arguments vary from verb to verb.

Adjuncts: optional, general arguments that any verb can take

AM-ADP	general purpose	AM-MOD	modif. verb
AM-CAU	cause	AM-NEG	negation marker
AM-DEP	deposition	AM-PER	perception
AM-DIS	descriptive marker	AM-REC	reciprocal
AM-EXT	extent	AM-LOC	location
AM-GEN	generality	AM-MNR	manner
AM-HDR	header	AM-NOM	nominal

3

### PropBank-style Semantic Roles

- The **PropBank** project produced semantic role annotations on the Wall Street Journal portion of the Penn Treebank (which already had parse tree annotations).
- PropBank defines predicate-argument structures for verbs with semantic role assignments for each verb's arguments.
- The predicate is labeled as REL (for relation) and is either a verb or a verb + particle (e.g., "keep up").
- PropBank's semantic role arguments are not named, but indicated as Arg0, Arg1, Arg2, etc. The meaning is specific to one verb sense! They do not have the same meaning for different verbs or different senses of the same verb.

2

### PropBank Frame File Example

Frame File for the verb 'expect':

Roles:

Arg0: expecter

Arg1: thing expected

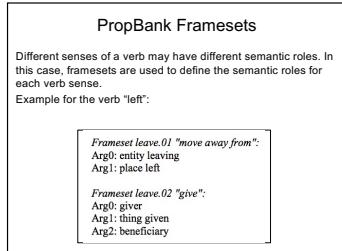
Example: Transitive, active:  
*Portfolio managers expect further declines in interest rates.*

Arg0: Portfolio managers

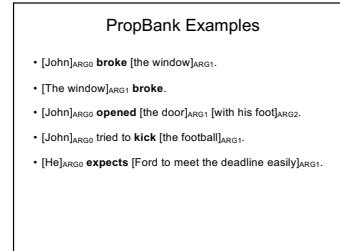
REL: expect

Arg1: further declines in interest rates

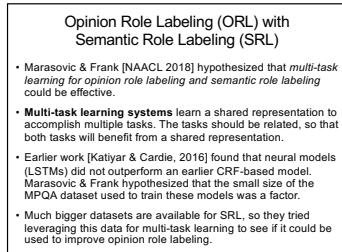
4



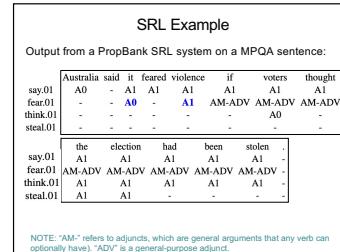
5



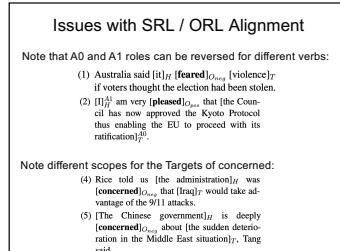
6



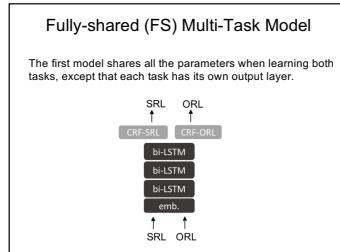
7



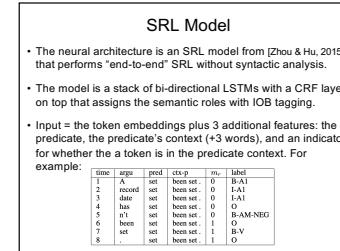
8



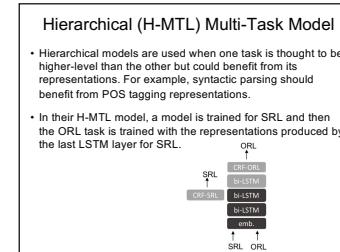
9



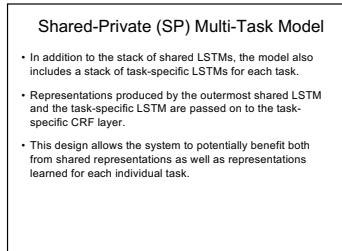
11



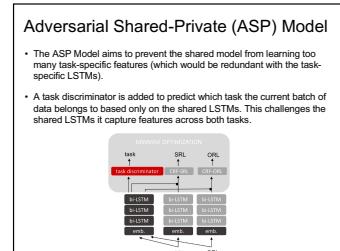
10



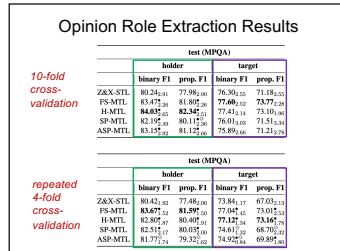
12



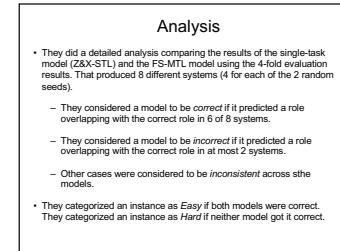
13



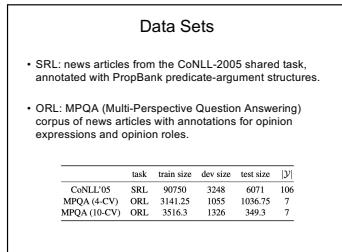
14



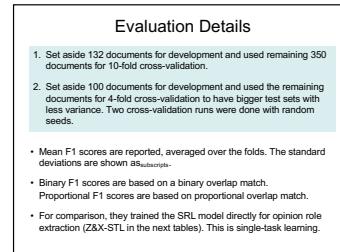
17



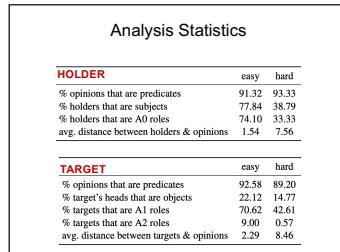
18



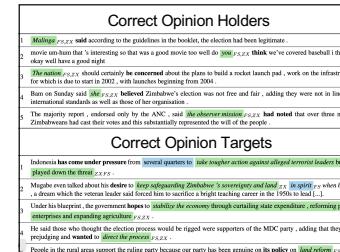
15



16



19



20

### Incorrect Opinion Holders

1. It would be entirely improper if, *in its defense* of *Israel*, the United States continues to exert pressure on [...].  
 2. **Indonesia** has come under pressure from several quarters to take tougher action against alleged terrorist leaders but has refused to do so.  
 3. The Geneva Convention and the *General Principle of International Law*, which states that all nations in the world must first respect and promote the humanitarian interests and progress of all humankind.  
 4. The government said that it will cost \$500 for an HIV/AIDS patient per year at this rate, and the following years this cost will increase by 10%.  
 5. The Organization of African Unity (OAU) *also* backed Zimbabwe President Robert Mugabe's revolution, *with* its members *saying* that the conflict was "irreconcilable".  
 6. Regarding the American proposed Anti-Missile Defense System too, neither Russia, China, Japan, nor even the European Union has been consulted about the project.  
 7. The president received his pledge to throw *terrorism* *over* *heads*, who want to "turn up" with regimes hoping to acquire weapons of mass destruction and said "nations will come with us" if the US-led war on terrorism is extended.

The correct answer is *italicized*.  
 Blue = FS-MTL Yellow = Z&X-STL Green = Both

21

### Incorrect Opinion Targets

1. State sanctioned land invasions, *several times declared* illegal by Zimbabwe's courts, as well as a drought have disrupted Zimbabwe's food production and famine is already looming in much of the country.  
 2. But he told *Interviewer*, *z,x*, that in spite of self-government in the agreement referred from previous West countries' concern over Zimbabwe's human rights record, he would not be able to implement policies in such a manner.  
 3. If the Europeans wish to influence Israel in the political arena – is a direction that many in Israel would support whole heartedly – they can not be able to pursue their policies in such a manner.  
 4. They are not going to be able to do anything, *including* *negotiations*, *in* *order* to bring about peace during a press conference [...].  
 5. And he has *listened* *just* *complained* – I don't want *to* *turn* *up* *with* *the* *Kyoto* *Protocol* *z,r,x*.  
 6. During *President Bush's* *peace*, I thought of *hockey* *z,z*: What are you going to do with the Kyoto Protocol *z,r,x*?  
 7. After I didn't want *to* *apply* *for* *z,r,x*. Not the principal called me during the summer months and said, "Sandra the time is running out, you need to *apply*".

The correct answer is *italicized*.  
 Blue = FS-MTL Yellow = Z&X-STL Green = Both

22

### Holders: FS-MTL is Correct, Z&X-STL Wrong

1. Yoshihisa Maruwa, a management consultant for Booz & Allen Hamilton Japan Inc., *said* *his* *team* *z,r,x* *will* *likely* *be* *recommending* *acquisitions* *of* *Japanese* *companies* *more*, *z,x* *often* *to* *foreign* *clients* *in* *the* *future*.  
 2. The *source* *r,z,x*, interviewed by *Forbes* in *Cosby*, *expressed* *concern* *that* *the* *command* *of* *the* *Russian* *forces* *is* *not* *adequate* *to* *defend* *the* *country*.  
 3. The Commonwealth man earlier this week *z,x*, *said* that "the conditions in Zimbabwe did not adequately allow the free and fair expression of will by the electorate".  
 4. Publishing such stand reports *will* *only* *create* *misrat* *among* *nations* *z,x* *regarding* *the* *objectives* *and* *independence* *of* *UN* *and* *its* *members*.  
 5. The Islamic Freedom Party – Democratic Alliance, New National Party, African Christian Democratic Party, the Pak Atlantic Congress and the United Christian Democratic Party *z,x* had disagreed with *the ANC*, *r,z,x*, *concluded*.  
 6. The Nigerian leader, President Olusegun Obasanjo *z,x*, had urged the *minister*, *r,z,x* *not* *to* *attack* *Algeria* *frontier* *outposts*.  
 7. US diplomats *z,x* *say* *Bush*, *r,z,x* *will* *seek* *to* *support* *Kim*'s *Nobel* *Pri* *winning* *policy* *by* *offering* *new* *talks* *with* *the* *North*, while *remaining* *firm* *with* *North* *Korea's* *missile* *tests* *and* *its* *forsen* *chemical* *and* *biological* *weapon* *programmes*.

The correct answer is *italicized*.  
 Blue = FS-MTL Yellow = Z&X-STL Green = Both

23

### Targets: FS-MTL is Correct, Z&X-STL Wrong

1. In most cases *described* *the* *legal* *punishment*, *z,x*, *like* *floggings* *and* *executions* *of* *mercenaries* *and* *major* *drug* *traffickers* *who* *are* *applied* *based* *on* *Sharia*, *Islamic* *law* *or* *human* *rights* *violations*.  
 2. In another verbal attack *Khartoum* accused *the* *United* *States*, *r*, *of* *wanting* *to* *exercise* "world dictatorship" since the "hostile" *US* *has* *been* *invading* *countries* *such* *as* *Iraq* *and* *Afghanistan*.  
 3. He said those who brought the election process would be regarded *as* *supporters* *of* *the* *MDC* *party*, *saying* *that* *they* *were* *prejudiced* *and* *wanted* *to* *direct* *the* *process*, *z,x*.  
 4. These factors *are* *not* *the* *only* *reason* *to* *have* *a* *more* *balanced* *view* *of* *the* *conflict*, *z,x* *is* *not* *the* *only* *reason* *to* *do* *that* *and* *small* *decisions*, *z,x*, *in* *fact*, *he* *admitted*.  
 5. But *his* *rough* *road* *to* *power*, *z,x*, *has* *provoked* *concerns* *in* *Seoul*, *z,x*, *where* *President* *Kim* *Tae-chung*, *who* *is* *in* *the* *last* *year* *of* *his* *five-year* *term*, *has* *been* *trying* *to* *prize* *the* *hermit* *state* *out* *of* *isolation*.

The correct answer is *italicized*.  
 Blue = FS-MTL Yellow = Z&X-STL Green = Both

24

### Temporal Information Extraction

- Time expressions and temporal relations are important to identify for many (most!) applications and domains.
- Temporal IE is especially relevant for event-related applications, such as event detection, event extraction, and event timeline construction.
- Most current work focuses on events in the past. Identifying future references poses additional types of challenges.
- Temporal recognition is also a basic element of language understanding!

1

### Temporal Understanding Examples

**Question:** When did airlines as a group last raise fares?

Last week, Delta boosted thousands of fares by \$10 per round trip, and most big network rivals immediately matched the increase.

A fare increase initiated *last week* by United Airlines was matched by competitors over the *weekend*, marking the second fare increase in two weeks.

On *Monday*, American Airlines raised fares on all domestic flights in the United States. Other airlines were slow to follow suit, but *two weeks later* other airlines have not yet matched the fare increases.

2

### Types of Temporal Expressions

Absolute	Relative	Durations
11/02/17 2020 May 4 May 11:00pm 11 o'clock 23:59 MDT Noon Christmas Christmas eve Pioneer Day 9/11	yesterday tonight in the morning last month 10:00-11:00 2 days ago next quarter in an hour before noon my birthday night when Xmas is over after the election soon	5 days 2 months 10 minutes 10:00-11:00 until midnight a few years for an hour in the winter dawn to dusk January-April before 2020 before this semester
yesterday tonight in the morning last month 10:00-11:00 2 days ago next quarter in an hour before noon my birthday night when Xmas is over after the election soon		
5 days 2 months 10 minutes 10:00-11:00 until midnight a few years for an hour in the winter dawn to dusk January-April before 2020 before this semester		

3

### Lexical Triggers for Time Expressions

**Nominal Nouns:** *morning, noon, night, winter, dusk, dawn, eve, hour, minute, sunrise, sunset*

**Proper Nouns:** *January, Monday, New Year's Eve, Labor Day, Easter, Passover, Ramadan*

**Adjectives/Prefixes:** *recent, last, next, annual, early, late, mid*

**Adverbs:** *hourly, daily, monthly, yearly, annually*

Most time expression recognizers identify trigger words and look at context around them:

Numeric "shape" patterns are needed too (e.g., 3/14/2017)

4

### Conclusions

- Multi-task learning can be an effective strategy for improving performance for a task when small amounts of task-specific gold data are available but there exists a related task with substantially more gold data.
- Their analysis found that long-distance dependencies remain challenging for opinion extraction.
- Reminder: these models were given an opinion expression and then extracted its role fillers. But finding the opinion expressions is also a key part of the task!
- Opinion extraction is a challenging IE task. Boundaries are tough (for humans & NLP systems!), and people bring a lot of knowledge to bear when identifying opinions.

25

### General Approaches

**Rule-based Systems:** finite-state automata can be built to recognize patterns for temporal expressions. Cascaded finite-state machines typically start with simple expressions and build patterns of greater complexity.

**Sequential Classifiers:** given annotated training data, sequential taggers can label time expressions using BIO tags.

**Constituent-based Classifiers:** given annotated training data, a constituent-based classifier can learn to label syntactic constituents (e.g., chunks or parse tree nodes).

**Pro:** boundary issues are separate syntactic problem  
**Con:** time expressions must align with syntactic constituents.

5

### Common Classifier Features for Time Expressions

- Lexical token
- Tokens in local context window
- Part-of-speech tags for target and window tokens
- Syntactic chunk/phrase type for target and window tokens
- Temporal keywords ("lexical triggers"), including days of the week, months, holidays, and general time words (e.g., "morning")
- Character-based "shape" features, for example:  
 #####  
 ##-##-##  
 ##-##  
 ##:#(am,pm)  
 ##:#(am,pm)(time-zone) 03/14/2017  
 ##:#(am,pm)(time-zone) 03-14-2017  
 ##:#(am,pm)(time-zone) 11:50  
 ##:#(am,pm)(time-zone) 11:50am  
 ##:#(am,pm)(time-zone) 11:50am (EDT)

6

### False Hits

Even strong time keywords can produce false hits for a variety of reasons, such as:

- Film, TV, book, and song titles, such as:  
 - *Any Given Sunday* (Oliver Stone film), *48 Hours* (TV series), *1984* (Orwell), *Tuesday Morning* (Mamas & Papas)
- Organization names, such as:  
 - *Black September* (terrorist group), *Tuesday Morning* (discount store)
- People names, such as:  
 - *April, May, June* are common female names
- Word sense (and part-of-speech) ambiguity  
 - *March, May*

7

The Weather Channel said dress for the mid 70s today.



8

## Temporal Normalization

- For real applications, time expressions need to be anchored to real dates (e.g., *Tuesday* → 3/29/2022) and mapped into a standardized format.
- Often, this requires extracting the dateline of a document (e.g., news article) or identifying its publication date.
- The normalization process can be complex and typically relies on hand-coded rules. Some examples:
  - last Tuesday*
  - this weekend*
  - 3 weeks ago*
  - in a week*
  - on Thanksgiving* (note: the date changes each year!)

9

## Temporal Relations

For many event-related applications, we want to know the temporal relationship between events (i.e., the ordering of events relative to each other).

Time expressions are usually present for only some of the events. Recognizing temporal relations between events allows for a relative timeline to be constructed.

The **TimeBank** corpus is a widely used resource that contains annotated temporal expressions, event mentions, and *Temporal Links* (TLINKs) between events and temporal expressions.

**TimeML** is a mark-up language for temporal information.

11

## Time Normalization Format Examples

Some normalization formats for fully specified dates:

Title	Pattern	Sample Value
Fully specified date	YYYY-MM-DD	1991-09-28
Weeks	YYYY-Www	2007-W27
Weekends	PwW	P1WE
24-hour clock times	HHMM:SS	11:13:45
Dates and times	YYYY-MM-DDTHH:MM:SS	1991-09-28T11:00:00
Financial quarters	Qn	1999-Q3

And ... dates are not always fully specified, so standards need to include partially specified dates too. For example:

- partial date **1995**
- open-ended ranges, such as before/after a specific date  
- **08 JAN 90**      **08 JAN 90**

10

## TimeBank 1.2

- 183 news articles annotated with time expressions, event mentions, and temporal links (and some other things too)
- Thirteen types of temporal links between events:
  - after/before**
  - includes/is\_included**
  - during/simultaneous** (for events/states that persist)
  - latter/before** (immediately after or before)
  - identity** (coreference)
  - begins/ends** (one event is beginning/end of another)
  - begun\_by/ended\_by** (inverse relation to begins/ends)

12

## (Simplified) TimeBank Annotation Example

The Russian airline Aeroflot has been **<E1: hit>** with a writ for loss and damages, **<E2: fled>** in Hong Kong by the families of seven passengers **<E3: killed>** in an air **<E4: crash>**.

All 75 people **<E7: state-on board>** the Aeroflot Airbus **<E5: died>** when it **<E6: ploughed>** into a Siberian mountain in **<T1: March 1994>**.

TLINKS:  
 (is included **<E6:ploughed T1:March 1994>**  
 (before **E3:killed E2:fled**)  
 (is included **E7:state-on board E4:crash**)  
 (includes **E7:on board E5:died**)  
 (after **E5:died E6:ploughed**)  
 (identiy **E6:ploughed E4:crash**)

13

## Task A: Recognizing Temporal Expressions

Task A is to recognize time expressions as defined by the TimeML TIMEX3 tag, as well as "type" and "value" attributes.

types = {DATE, TIME, DURATION, SET}

The possible values depend on the type.

<TIME3 id="1" type="DATE" value="1999-SU">  
*the summer of 1999* </TIME3>  
<TIME3 id="12" type="TIME" value="T24:00">  
*twelve o'clock midnight* </TIME3>  
<TIME3 id="13" type="DURATION" value="P2D">  
*two entire days* </TIME3>  
<TIME3 id="14" type="SET" value="XXXX-10">  
*every October* </TIME3>

15

## TIPSem Temporal IE System

- TIPSem [Llorens et al., 2010] is a set of temporal IE systems that performed well in TempEval-2 tasks. (TIPSem = Temporal Information Processing based on Semantic information)
- They used similar sequential tagging (CRF) models across six different temporal tasks.
- An unusual emphasis of their work is the incorporation of semantic information.
- They created systems for both English and Spanish, demonstrating the generality of their approach.

14

## Task B: Recognizing Events

- Task B is to recognize and classify events, as defined by the TimeML EVENT tag.

Events can be expressed as verbs, nominalizations, adjectives, predicative clauses, or PPs (but only the heads are annotated). Events also have 4 types of information:

- Polarity** captures negative (e.g., for events that did not happen)
- Tense** captures temporal verb forms (e.g., past, present and future)
- Aspect** captures verbal information about how events & states extend over time. In English: neutral, progressive, perfect, progressive perfect, and (in past tense) habitual. For example: *"I lose* vs. *"I am losing* vs. *"I have been losing"*.
- Modality** captures ability, possibility, permission or obligation as indicated by modal verbs (e.g., *could*, *should*, *must*).

16

## Example of Event Annotations

Five days after he **<EVENT eid="e1" class="OCCURRENCE> came**  
**<EVENT> back ...**

A major **<EVENT eid="e2" class="OCCURRENCE> earthquake**  
**<EVENT> in Indonesia ...**

After many months of renewed **<EVENT eid="e3" class="STATE> hostility**  
**<EVENT> ...**

The **<EVENT eid="e4" class="OCCURRENCE> attack** **<EVENT> was not** **<EVENT eid="e5" class="STATE> expected** **<EVENT> at all ...**

The full annotations would include polarity, tense, aspect, and modality as well.

17

## Tasks C-F: Temporal Relation Links

**Task C:** Determine the temporal relation between an event and a time expression in the same sentence, where the event syntactically dominates the time expression or they occur in the same NP.

**Task D:** Determine the temporal relation between an event and the *document creation time* (DCT).

**Task E:** Determine the temporal relation between two main events in consecutive sentences.

**Task F:** Determine the temporal relation between two events where one event syntactically dominates the other event. For example: *"she heard an explosion"* or *"he said they postponed the meeting"*.

18

## Types of Relation Links

Temporal relations have 6 relation types: **Before**, **After**, **Overlap**, **Before-or-Overlap**, **Overlap-or-After**, or **Vague**.

### Examples:

*Mary taught<sub>e1</sub> on Tuesday<sub>n1</sub>* → Overlap(e1, n1)  
*The country defaulted<sub>e2</sub> on debts for that entire year* → Before(e2, DCT)  
*The students heard<sub>e1</sub> a fire alarm<sub>e2</sub>* → Overlap(e1, e2)  
*He said<sub>e1</sub>, they had postponed<sub>e2</sub> the meeting* → After(e1, e2)

19

## Classification Models

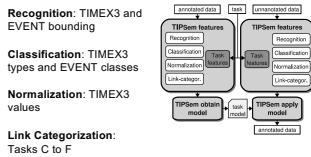
- TIPSem systems are CRF models for sequential tagging trained with supervised learning.
- For Tasks A and B, the input is a word sequence and the output is BIO labels on the words.
- For Tasks C-F, the input is instances of the classes (e.g., TIMEX3 and EVENT instances) and the output is relation links.
- In addition to traditional features, they use information from a semantic role labeling (SRL) system. For example:

*John<sub>AGENT</sub> sold<sub>VERB</sub> Mary<sub>RECIPIENT</sub> a car<sub>THEME</sub>*

20

## TipSem Architecture

They grouped the tasks into 4 categories, which each use different types of feature sets:



21

## General Features

- Morphological:** lemmas and POS tags, for context windows of size +/- 2.
- Syntactic:** Phrase level syntactic information from parse trees produced by constituency parsers.
- Polarity, tense, and aspect:** hand-crafted rules that use POS tags

For example: *will + VERB → FUTURE tense*

22

## Semantic Features

- Role:** the semantic role based on the verb it depends on in the parse. The CCG SRL tool was used for English, and AnCora for Spanish.
- Governing Verb:** the verb that the token depends on in the parse.
- Role+Verb combination:** governing verb paired with role.
- Role configuration:** for verbs that head a sentence or nested sentence, the set of roles that depend on the verb are captured.
- Lexical semantics:** the top 4 semantic classes for a word, based on WordNet for English and EuroWordNet for Spanish.

23

## Link Categorization Features

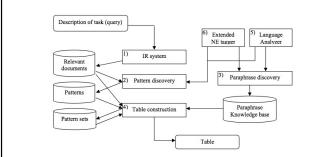
- Head preposition:** if TIMEX or EVENT in a PP.
  - Syntactic relation of the TIMEX and EVENT:** same sentence, subsentence, or subphrase.
  - Time position:** if EVENT is not directly linked to a TIMEX, then it is before, after, or overlapping with the TIMEX.
  - Interval:** if there is an interval indicator near the TIMEX type.
  - Semantic roles:** if the EVENT or TIMEX are labeled with a temporal role. For example: *After<sub>verb</sub> he left<sub>1</sub> home ...*
- (NOTE: Different tasks used different combinations of the features.)

24



### Prior Work: On-Demand Information Extraction

[Sekine 2006] proposed “on-demand” information extraction, where a user would provide a query for a desired relation and the system would automatically learn paraphrases and build a table of extracted information.



5

### Preemptive Information Extraction

- [Shinyama and Sekine, 2006] explored the idea of **preemptive information extraction** and proposed:
  - “a technique called Unrestricted Relation Discovery that discovers all possible relations from texts and presents them as tables.”*
- Their system used clustering, pattern learning, and meta-clustering to build a set of tables filled with information extracted for different relations, without training data.
- This work was among the earliest research similar to Open IE, and a preliminary system was built. But the effort did not continue on a large scale. However, similar efforts were undertaken by other research groups on a much larger scale...

6

### UW’s Open IE Research Efforts

Following KnowItAll, UW embarked on a long-term research effort focused on Open IE from the Web and it continued to evolve.

- ReVerb:** identifies and extracts unspecified binary relations.
- RESOLVER:** a probabilistic relational model for determining whether two relation expressions are “synonymous”.
- TextRunner** generates labeled examples using heuristics and trains a classifier for unrestricted relation extraction. RESOLVER is incorporated to identify synonymous relation phrases.
- SHERLOCK:** learns first-order Horn Clauses as inference rules. For example: `Contains(Food, Chemical) :- isMadeFrom(Food, Ingredient) A Contains(Ingredient, Chemical);`

9

### KnowItAll Rule Examples [Etzioni et al., 2004]

```
Extraction Rule:
NP1 ["is"] "such as" NP2
NP1 ["is"] "and other" NP2
NP1 ["is"] "including" NPList2
NP1 ["is"] NP2
NP1 ["is the"] NP2 "of" NP3
"the" NP1 "of" NP2 "is" NP3

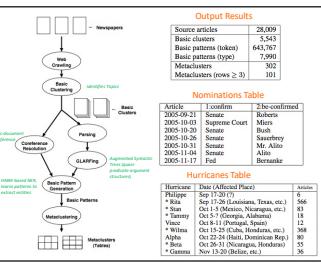
Extraction Rule for a Binary Relation:
NP1 ["plays for" NP2
& properNoun(head(NP1))
& head(NP2)="Seattle Mariners"
instanceOf(Athlete,head(NP1))
instanceOf(SportTeam,head(NP2))
& playsFor(head(NP1),head(NP2))
keywords="plays for", "Seattle Mariners"
```

10

### Why is Open IE more difficult?

- Open IE has to identify **both** the set of entities that participate in a relation as well as the textual clues that reveal the relation.
- A relation-independent process means that relation-specific features cannot be exploited.
- Many relations can have a wide variety of argument types, so anchoring the relation with named entities is not sufficient. Even the general types of arguments are not constrained in any way.

12



7

### KnowItAll

- A research group at the University of Washington began an OpenIE research project called **KnowItAll**, which produced a steady stream of research results related to open information extraction for many years.
- The emphasis of this project has been massive Web-scale IE, with an emphasis on speed and extracting large volumes of information.
- Consequently, many of the methods use very shallow pattern matching and few NLP tools.
- The original KnowItAll system used Hearst’s hyponym patterns to identify relation instances in an iterative learning framework.

**Hyponym Patterns:** <> such as <y>, <x> including <y>, etc.

8

### Tradeoffs Between Open and Traditional Relation Extraction [Banko & Etzioni, ACL 2008]

- “Open IE is a *relation-independent* extraction paradigm that is tailored to massive and heterogeneous corpora such as the Web.”
- “An Open IE system extracts a diverse set of relational tuples from text without *any* relation-specific input.”
- In contrast, traditional relation extraction systems begin with training examples for a specific type of relation and learn to identify instances of that type of relation.
  - Each relation requires its own training process.

11

### Common Syntactic Patterns

- 500 randomly sampled sentences were reviewed to manually identify the types of constructions that captured a relation expression.
- 95% of the identified patterns could be grouped into 8 lexico-syntactic categories.
- While these patterns are not sufficient to identify a relation, these results suggest that most relation expressions can be captured by a small set of pattern templates.

13

### Seed Labeled Data

Relation-independent heuristics were applied to the Penn Treebank to obtain labeled relation instances, which were designed to approximate syntactic dependencies and semantic roles.

NO parsing or semantic analysis! Just NP chunking and POS tags.

For example:

Class: + Heuristic: Subject, Verb, Object (SVO) Triple

Example: "*Einstein* received *the Nobel Prize*"

Class: - Heuristic: ADVP crossing

Example: "*He studied Einstein's work* when visiting *Germany*"

3/3

### Common Lexico-syntactic Patterns

95% of the 500 sampled sentences had relation expressions matching one of these patterns.

Relative Frequency	Category	Simplified Lexico-Syntactic Pattern
37.8	Verb	E, Verb E <sub>2</sub> X established Y
22.8	Noun+Prep	E, NP Prep E <sub>2</sub> X mentioned with Y
16.0	Verb+Prep	E, Verb NP Prep E <sub>2</sub> Y mentioned with X
9.4	Infinitive	E, to Verb E <sub>2</sub> X plans to acquire Y
5.2	Modifier	E, Verb E <sub>2</sub> Noun E <sub>3</sub> X modify Y
1.8	Coordinate	E <sub>1</sub> (and,;) E <sub>2</sub> NP X deal
1.0	Coordinate	E <sub>1</sub> (and,;) E <sub>2</sub> Verb X, Y merge
0.8	Affiliative	E, NP (G,J) E <sub>2</sub> X affiliates : Y

14

### O-CRF

- Labeled instances for training are generated heuristically.
- A sequential tagging model (CRF) is trained to label tokens that express a binary relation using IOB tags.
- A noun phrase chunker is applied and all NPs pairs within a certain distance from each other are candidates for a relation instance.
- The feature set includes POS tags, regular expressions to detect things like capitalization and punctuation, context words, and conjunctions of features for adjacent positions in a context window of size +/- 6 words.
- Context words are only captured for closed class words and **not** for open class words! Presumably for improved generality.

16

### Relation Extraction as Sequence Labeling

Each relation must be anchored by two noun phrases, which are called “entities” (ENT).

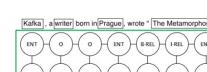


Figure 1: Relation Extraction as Sequence Labeling: A CRF is used to identify the relationship, *wrote in*, between *Kafka* and *Prague*.

17

### Evaluation of Open Relation Extraction

The first evaluation compares the performance of O-CRF with the TextRunner Open IE system (O-NB). TextRunner had extracted 7.5 million tuples from 9 million Web pages.

Both systems were tested on 500 sentences.

Category	O-CRF			O-NB		
	P	R	F1	P	R	F1
Verb	93.9	65.1	76.9	100	38.6	55.7
Noun+Prep	89.1	36.0	51.3	100	9.7	55.7
Verb+Prep	95.2	59.5	65.6	95.2	25.3	40.0
Infinitive	95.7	46.8	62.9	100	25.5	40.6
Modifier	0	0	0	0	0	0
Coordinate	0	0	0	0	0	0
All	<b>88.3</b>	<b>45.2</b>	<b>59.8</b>	86.6	23.2	36.6

19

### O-CRF’s Limitations

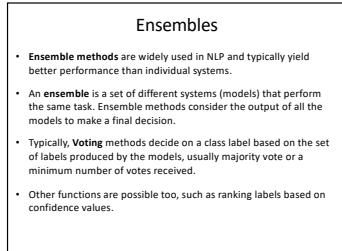
- Relations can only be identified if they are explicitly mentioned in a text.
- Relations can only be identified through lexical context. Document style features are not considered.
- Relations can only be identified between NPs within the same sentence.
- O-CRF does not cluster/normalize relations.
  - Relation “synonyms” (paraphrases) were identified by a different system called RESOLVER [Yates and Etzioni, 2007].

18

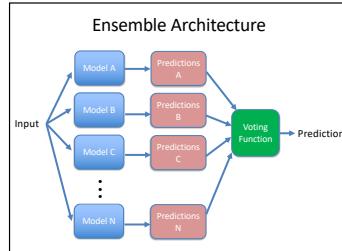
### Relation-Specific Extraction

- For comparison, a traditional relation extraction system was trained with a CRF model, which they called R1-CRF.
- R1-CRF is identical to O-CRF except:
  - R1-CRF was trained from manually labeled positive and negative instances of a specific relation R.
  - R1-CRF used both closed-class and open-class words as features. (O-CRF could only use closed-class words.)
  - No additional steps are needed to identify the relation type, since it is trained to identify only instances of relation R.

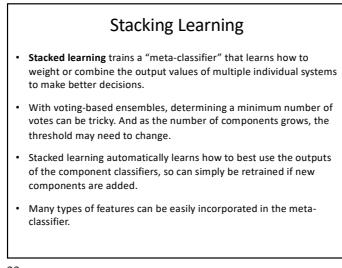
20



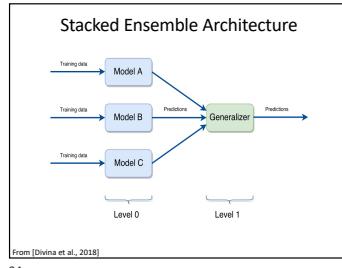
21



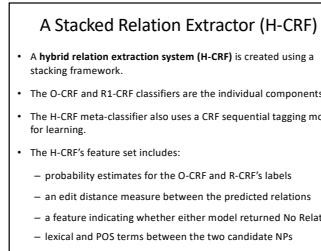
22



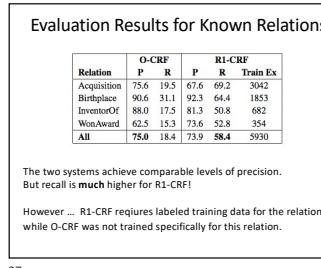
23



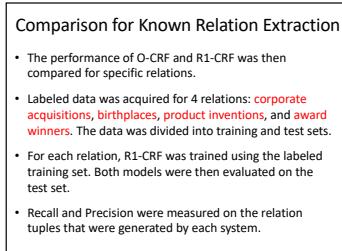
24



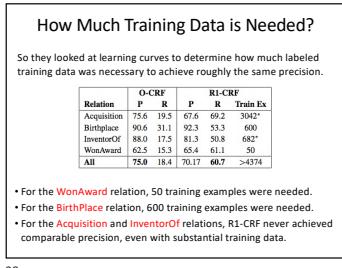
25



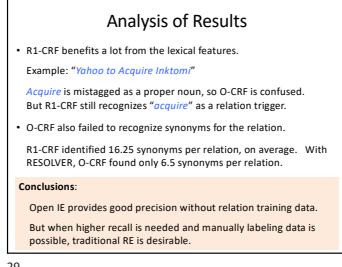
27



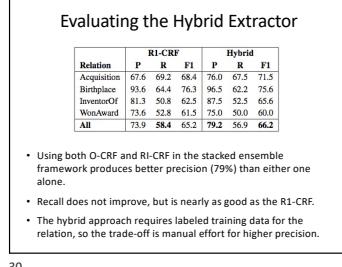
26



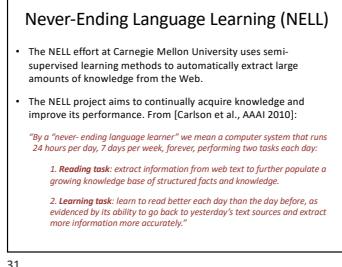
28



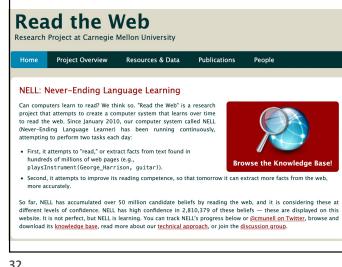
29



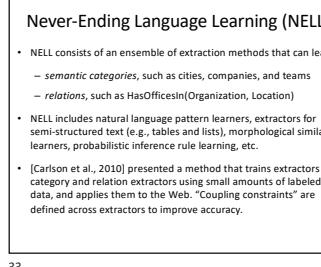
30



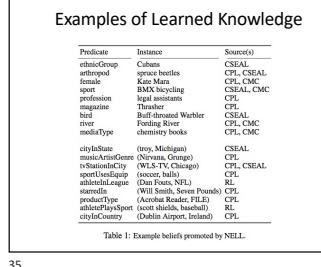
31



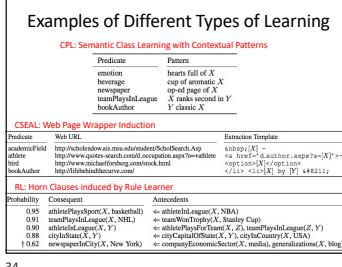
32



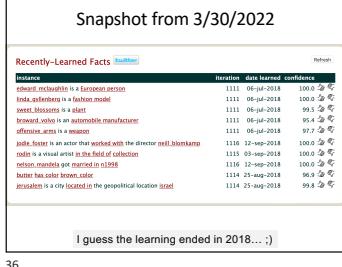
33



35



34



36

**NELL** (@nellcmu)  
I am a machine reading research project at Carnegie Mellon, periodically tweeting facts I need. Please follow me, and reply with corrections so I can improve!  
Pittsburgh PA / [rw.ml.cmu.edu](http://rw.ml.cmu.edu) Joined March 2010  
581 Following 3,070 Followers

NEIL (@nellcmu) · Feb 22, 2019  
True or False? "school golf" is a #Sport ([bit.ly/2SUS82n](https://bit.ly/2SUS82n))  
6 4 1 1

NEIL (@nellcmu) · Feb 22, 2019  
True or False? "prescription pseudoephedrine" is a #Music ([bit.ly/2E1DQG0](https://bit.ly/2E1DQG0))  
3 1 1 1

3/3

### Never-Ending Image Learner (NEIL)

Interesting Aside: NELL also inspired a follow-on effort at CMU called NEIL to continually extract visual knowledge!

**NEIL: Never Ending Image Learner**

I crawl, I see, I learn.

**OBJECTS**  
How does a computer know what our books look like? How does a sheep know what a computer looks like? Can a computer learn all the objects in the world? NEIL is a computer program that tries to do just that! It is an effort to build the world's longest visual knowledge base. NEIL has learned to identify over 100,000 different objects in over 100,000 images. It is learning to identify many common objects and 40 others. See current statistics about how much NEIL knows about our world!

**SCENES**  
TO SHOW THE VISUAL KNOWLEDGE FULLY

**ATTRIBUTES**  
To show the knowledge learned, you can browse the knowledge base by clicking on categories in the left-hand panel. Simply type the search term in the top right. Each page shows the visual examples and the common sense facts about a category.

**TRAIN A CONCEPT**

38

### Conclusions

- Open information Extraction holds great promise for automatically constructing large and rich knowledge bases.
- These efforts have advanced the state-of-the-art for robustly and efficiently extracting large volumes of diverse knowledge from unstructured, often unwieldy Web text.
- However, there is ample room for improvement in the accuracy, organization, and richness of the learned knowledge.  
– organizing the learned knowledge is a key challenge!
- Open IE learners tend to learn the most prevalent facts and relations, and are less able to learn less common knowledge or acquire specialized concepts with domain-specific idiosyncrasies.

39

### Common Sense Knowledge for NLP

- People rely on a great deal of common sense knowledge when understanding language.
- Most NLP systems are trained for a specific task using the words in the input and (sometimes) additional features. The features typically represent relatively shallow or general information (e.g., proximity, orthography, syntax, or general semantic knowledge).
- Our NLP models need to know the same common sense knowledge that people do in order to:
  - make more intelligent decisions
  - make the same inferences that people do

1

### The Winograd Challenge (Coreference Resolution)

The city councilmen refused the demonstrators a permit because they feared violence. → they = councilmen  
they advocated violence. → they = demonstrators

The lawyer asked the witness a question, but he was reluctant to answer it.  
→ he = lawyer

The man couldn't lift his son because he was so weak.  
→ he = man

I poured water from the bottle into the cup until it was full.  
→ it = cup  
empty.

2

### Machine Translation

Machine translation systems often do quite well with word sense disambiguation by relying on the words in the surrounding context.

For example:

The electrician is working. (working → Labor)  
The telephone is working. (working → Functional)

But these systems often get confused when words associated with different senses are interspersed. For example:

English original	Google translation
The electrician is working.	Der Elektriker arbeitet.
The electrician that came to fix the telephone is working.	Der Elektriker die auf das Telefon zu beben kann funktioniert.
The telephone is working.	Das Telefon funktioniert.
The telephone on the desk is working.	Das Telefon auf dem Schreibtisch arbeitet.

3

### Temporal Knowledge

The duration of events is a form of temporal knowledge that relies on common sense knowledge.

Julie dropped her iPhone on the floor. (seconds)

Julie made a sandwich for lunch. (minutes)

Julie watched a movie. (hours)

Julie went on vacation to Seattle. (days/weeks)

Julie took a class on Natural Language Processing. (months)

Julie got her computer science degree at the Univ. of Utah. (years)

4

4/

### Story Understanding

- Mary went to a restaurant.
- George went to the dentist.
- Julie finished the watermelon.
- Julie finished the book.
- Tom ordered a pizza.
- Tom ordered a taxi.
- The boy had a bone stuck in his throat.
- Max needed money to get his car fixed. He called his sister.
- The lion spotted an antelope on the hill.

5

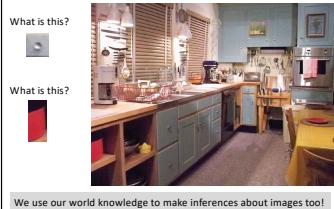
### A Short Story

John got up one morning and discovered his power was out. Unable to shave, he called his next door neighbor and asked if he could come over to borrow the bathroom. But everyone on the street was out. So John called FG&E and drove to work hoping no one would see him before he found a bathroom with hot water. Unfortunately, he ran into his boss on the elevator. He explained his predicament, but did not feel reassured by Mr. Carver's silence. John stumbled through the rest of the week half-expecting to find a pink slip in his mailbox.

We make numerous commonsense inferences when we read!

6

### Common Sense for Visual Understanding



7

### Plausible Reasoning

- When we communicate, we make inferences about things that are very likely to be true, but not guaranteed (defeasible inferences).
  - We make assumptions about the most common or typical situations.
  - We usually assume that we are being told the truth and that we are being given complete information (Gricean maxims).
  - We often reason based on similar situations that we know about ("case-based reasoning").
- Generating a knowledge base of common sense information can capture the types of default ("prior") knowledge that we assume to be true in the absence of information to the contrary.

8

### Transitivity across Relations?

- Commonsense knowledge is often captured in ISA (hypernym/hyponym) hierarchies, and we often assume transitivity for hierarchical relations. For example:  
(dog ISA mammal) & (mammal ISA animal) → (dog ISA animal)
- But this is not always true ... the world is complicated.

chair IS furniture	I fit inside my clothes
car seat IS chair	My clothes fit inside a drawer.
car seat is NOT furniture!	I DO NOT fit inside a drawer.

9

### CYC

- CYC was an ambitious project started in 1984 aimed at manually compiling a massive repository of common sense knowledge. It was originally envisioned as a 10-yr project, but still exists today!
- As of 2012, the public version OpenCyc 4.0 contained 239,000 concepts and over 2 million facts, organized in a taxonomy.
- A larger ResearchCyc is available with a license, and contains 500,000 concepts and 5 million facts.
- Some people have reported using CYC for Web query expansion, question answering, and intelligence analysis.
- As with many large, manually curated KBs, people have complained that it is organized poorly and unevenly. How to represent knowledge and organize knowledge is a critical but often overlooked problem!

10

### ConceptNet

ConceptNet [Speer et al., AAAI 2017] is a large knowledge graph that connects words with phrases that capture a variety of relations.

- Facts acquired from Open Mind Common Sense (OMCS) (Speer et al. 2009) and sister projects in other languages (Anand et al. 2009)
- Information extracted from parsing Wiktionary, in multiple languages, with a custom parser ("Wikisense")
- Commonsense knowledge from the Penn Treebank (Baker et al. 2011) (Koehn et al. 2009)
- Open Multilingual WordNet (Miller et al. 1998) and its English synonym network (Nakabayashi et al. 2006)
- IMDb (Borgman 2004), a movie database and movie dictionary
- OpenCyc (Carroll 1998), a system that represents common sense knowledge in a formal logic
- OpenCyc, a hierarchy of hypernyms provided by Cyc (Leutgeb and Giela 1998), a system that represents common sense knowledge in a formal logic
- A subset of DBpedia (Auer et al. 2007), a network of facts extracted from Wikipedia infoboxes

11

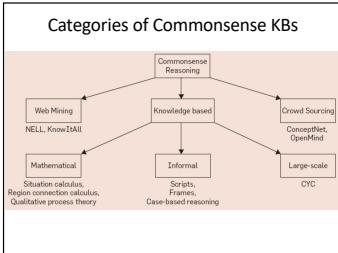
### ConceptNet's Content

- ConceptNet contains over 21 million edges and 8 million nodes. Its English vocabulary contains 1.5 million nodes.
- ConceptNet includes 36 core relations.
  - Symmetric relations:** *Anonym, DistinctFrom, Etyymology, IsPartOf, LocatedNear, RelatedTo, SimilarTo, and Synonym*
  - Asymmetric relations:** *AtLocation, CapableOf, Causes, CausallyRelatedTo, LocatedNear, RelatedTo, SimilarTo, and Synonym*
  - Transitive relations:** *Entails, Entailment, ForAll, HasBody, HasHome, HasProperty, InstanceOf, IsA, MadeOf, MannerOf, MotivationOf, ObstructedBy, PartOf, ReceivesAction, SenseOf, SymbolOf, and UsedFor*

12

Sample of ConceptNet's "knife" Info		
knife is used for...	knife is capable of...	Location of knife
stabbing ->	spread peanut butter ->	the kitchen ->
cutting ->	spread peanut butter ->	the drawer ->
cutting food ->	spreading butter bread ->	the kitchen drawer ->
carving wood ->	cut ->	a pocket ->
slicing ->	cut that apple ->	your back ->
boning ->	hurt a dog ->	a backpack ->
cut meat with ->	hurt ->	a drawer in a kitchen ->
cut string ->	butter brie ->	a fishing boat ->
cutting steak ->	cut that cake ->	a knife block ->
killing ->	cut cheese ->	a knife-holder ->
pare an apple ->	cut leather ->	a knife store ->
scratching ->	cut a man's hand ->	a plate ->
slicing bread ->		a sheath ->

13



14

- Reporting Bias**
- A dream for NLP is to conquer the *knowledge acquisition bottleneck* by automatically extracting common sense knowledge from texts.
  - We often assume that the more often we read something, the more likely it is to be true. Sometimes this is the case, but not always!
    - For example, some body parts are mentioned in text much more often than others. The Knes system found > 1 million instances of "*people may have eyes*" but < 1,500 instances of "*people may have a spleen*". But all body parts are equally likely in people!
  - The discrepancy between reality and coverage in text has been called "*reporting bias*" [Gordon & Van Durme].
    - The problem is especially acute for common sense knowledge because these facts are so obvious to people that they are rarely mentioned!

15

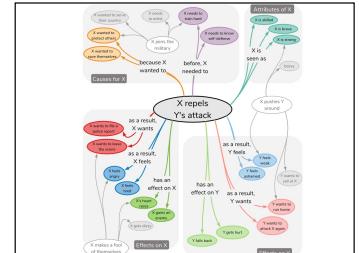
- Examples of Reporting Bias Issues**
- Events like murders and robberies are reported far more often than events like sleeping and breathing, which we all do every day.
  - Airplane crashes are mentioned far more often than motorcycle crashes, but the latter are much more common.
  - There are far more reports of people winning races than losing races, but there are many more losers than winners.
  - In a forest fire story, it's common to mention that homes were destroyed and people killed, but rare to mention that deer, raccoons, and squirrels were killed.
  - We're more likely to mention someone's hair color if it's purple or red, than if it's brown or black.
  - If we mention a grocery store trip, we rarely mention details that almost certainly happened, such as: grabbing a shopping cart, walking down the aisles, putting items in the cart, standing in the checkout line, paying for the groceries, etc.

16

## ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning

- [Sap et al., 2019] create a large common sense knowledge base called ATOMIC that contains 877k textual descriptions of inferential knowledge.
- ATOMIC was built semi-automatically by extracting a large set of events from text corpora and crowd-sourcing information about event relationships.
- The goal is to capture common world knowledge and enable inferences about everyday events, their causes, and their effects.

17



18

## Types of Knowledge

- If-Event-Then-Mental-State:** captures mental pre- and post-conditions of an event: <X, action, Y>
- X Intent:** likely reason why the agent performed the action
- X Reaction:** emotional reaction of the agent
- Other Reaction:** emotional reaction of others

For example, given event: X compliments Y

- X wants to be nice
- X feels good
- Y feels flattered

19

## Types of Knowledge

- If-Event-Then-Event:** probable preceding and following events.
  - X Need:** pre-condition for the event for the agent
  - Effect on X:** voluntary post-conditions for the agent
  - Effect on Other:** voluntary post-conditions for other
  - X Want:** involuntary post-conditions for the agent
  - Other Want:** involuntary post-conditions for other

For example, given event: X makes Y's coffee

- X needs to put coffee in the filter
- X adds cream and sugar
- Y drinks coffee
- Y gets thanked by X

20

## Types of Knowledge

- 3. If-Event-Then-Persona:** static relation that captures how the agent is described or perceived.

- a) **X Attribute:** resulting perceived attributes for the agent

For example, given the event: X calls the police

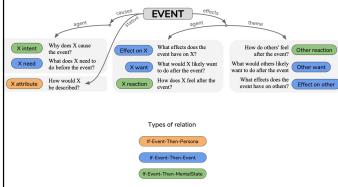
- a) X is lawful

For example, given the event: X returns wallet containing money

- a) X is honest

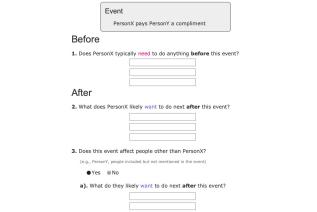
21

## 9 Types of Knowledge in All



22

## Example of Crowd-Sourcing Interface



25

## Examples of If-Event-Then-X rules in ATOMIC

Event	Type of relation	Inference examples	Inference dim
PersonX wants to be nice	If-Event-Then-Mental-State	PersonX will be appreciative	Agent
PersonX will feel flattered	If-Event-Then-Mental-State	PersonX will feel flattered	Agent
PersonX will want to chat with PersonY	If-Event-Then-Event	PersonX will talk to PersonY	Want
PersonX will be helpful	If-Event-Then-Event	PersonX will help	Want
PersonX will compliment PersonX back	If-Event-Then-Event	PersonX will compliment PersonX back	Want
PersonX is flattering	If-Event-Then-Event	PersonX is flattery	Want
PersonX is helpful	If-Event-Then-Event	PersonX is helpful	Want
PersonX is kind	If-Event-Then-Event	PersonX is kind	Want
PersonX is nice	If-Event-Then-Event	PersonX is nice	Want
PersonX is thoughtful	If-Event-Then-Event	PersonX is thoughtful	Want
PersonX is friendly	If-Event-Then-Event	PersonX is friendly	Want
PersonX is helpful	If-Event-Then-Mental-State	PersonX will be helpful	Agent
PersonX will be appreciated	If-Event-Then-Mental-State	PersonX will be appreciated	Agent
PersonX needs to be nice	If-Event-Then-Mental-State	PersonX needs to be nice	Agent
PersonX will feel flattered	If-Event-Then-Mental-State	PersonX will feel flattered	Agent
PersonX will want to chat with PersonY	If-Event-Then-Event	PersonX will talk to PersonY	Want
PersonX will be helpful	If-Event-Then-Event	PersonX will help	Want
PersonX will compliment PersonX back	If-Event-Then-Event	PersonX will compliment PersonX back	Want
PersonX is flattering	If-Event-Then-Event	PersonX is flattery	Want
PersonX is helpful	If-Event-Then-Event	PersonX is helpful	Want
PersonX is kind	If-Event-Then-Event	PersonX is kind	Want
PersonX is nice	If-Event-Then-Event	PersonX is nice	Want
PersonX is thoughtful	If-Event-Then-Event	PersonX is thoughtful	Want
PersonX is friendly	If-Event-Then-Event	PersonX is friendly	Want
PersonX needs to be nice	If-Event-Then-Mental-State	PersonX will be appreciated	Agent
PersonX will feel flattered	If-Event-Then-Mental-State	PersonX will feel flattered	Agent
PersonX will want to chat with PersonY	If-Event-Then-Event	PersonX will talk to PersonY	Want
PersonX will be helpful	If-Event-Then-Event	PersonX will help	Want
PersonX will compliment PersonX back	If-Event-Then-Event	PersonX will compliment PersonX back	Want
PersonX is flattering	If-Event-Then-Event	PersonX is flattery	Want
PersonX is helpful	If-Event-Then-Event	PersonX is helpful	Want
PersonX is kind	If-Event-Then-Event	PersonX is kind	Want
PersonX is nice	If-Event-Then-Event	PersonX is nice	Want
PersonX is thoughtful	If-Event-Then-Event	PersonX is thoughtful	Want
PersonX is friendly	If-Event-Then-Event	PersonX is friendly	Want
PersonX needs to be nice	If-Event-Then-Mental-State	PersonX will be appreciated	Agent
PersonX will feel flattered	If-Event-Then-Mental-State	PersonX will feel flattered	Agent
PersonX will want to chat with PersonY	If-Event-Then-Event	PersonX will talk to PersonY	Want
PersonX will be helpful	If-Event-Then-Event	PersonX will help	Want
PersonX will compliment PersonX back	If-Event-Then-Event	PersonX will compliment PersonX back	Want
PersonX is flattering	If-Event-Then-Event	PersonX is flattery	Want
PersonX is helpful	If-Event-Then-Event	PersonX is helpful	Want
PersonX is kind	If-Event-Then-Event	PersonX is kind	Want
PersonX is nice	If-Event-Then-Event	PersonX is nice	Want
PersonX is thoughtful	If-Event-Then-Event	PersonX is thoughtful	Want
PersonX is friendly	If-Event-Then-Event	PersonX is friendly	Want
PersonX needs to be nice	If-Event-Then-Mental-State	PersonX will be appreciated	Agent
PersonX will feel flattered	If-Event-Then-Mental-State	PersonX will feel flattered	Agent
PersonX will want to chat with PersonY	If-Event-Then-Event	PersonX will talk to PersonY	Want
PersonX will be helpful	If-Event-Then-Event	PersonX will help	Want
PersonX will compliment PersonX back	If-Event-Then-Event	PersonX will compliment PersonX back	Want
PersonX is flattering	If-Event-Then-Event	PersonX is flattery	Want
PersonX is helpful	If-Event-Then-Event	PersonX is helpful	Want
PersonX is kind	If-Event-Then-Event	PersonX is kind	Want
PersonX is nice	If-Event-Then-Event	PersonX is nice	Want
PersonX is thoughtful	If-Event-Then-Event	PersonX is thoughtful	Want
PersonX is friendly	If-Event-Then-Event	PersonX is friendly	Want

26

## Statistics for ATOMIC's Content

The resulting knowledge graph contains over 300k nodes for the original 24k events. (A node is a short phrase, 2.7 words on avg.)

Each triple is of the form <event, relation, event>

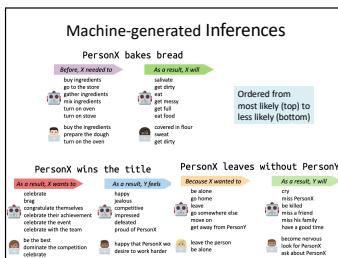
	Count	words
# triples: If-Event-Then-X	877,108	-
- Mental-State	212,598	-
- Event	521,334	-
- Person	143,706	-
# nodes: If-Event-Then-X	209,515	2.7
- Mental-State	51,928	2.1
- Event	245,905	3.3
- Person	11,490	1.0
Base events	24,313	4.6
# nodes appearing > 1	47,356	-

27

## Learning to Perform Commonsense Inference

- They framed the inference task as conditional sequence generation: given an event phrase **e** and an inference dimension **c**, train a model to generate the target (inference) : **f(e,c)**
- The architecture is an **encoding-decoding** framework:
  - First, they create an embedding representation for the event phrase beginning with pre-trained GloVe vectors concatenated with pre-trained ELMo contextualized embeddings and then further encoded with a bidirectional GRU.
  - Next, they train a decoder (unidirectional GRU) to generate an output string given the event phrase's encoding.

28



29

**Experiments**

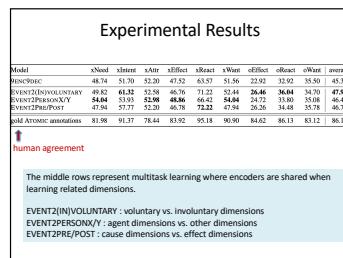
- They split the seed events into 80% training, 10% validation, and 10% test.
- Given a test event, the model generated phrases for each of the 9 dimensions of If-Then inferences.
- Automatic Evaluation: BLEU scores (n-gram matching) for the system's top 10 predictions compared with the answers from the crowd workers. (Hard to interpret results, see the paper for more.)
- Human Evaluation: randomly selected 100 events and generated 10 most likely system predictions. 5 crowd workers judged how many of them were valid.
- Results presented as **Precision@10**: average number of correct predictions in the top 10.

4/

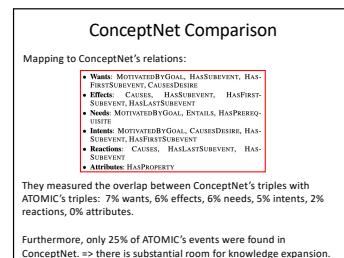
**Conclusions**

- ATOMIC was among the first large-scale efforts to build a knowledge graph related to events and common sense inferences associated with events.
- In principle, the dimensions capture many different aspects of common sense reasoning associated with events (causes, effects, emotional reactions, etc.)
- They then used ATOMIC to show how a neural network can be trained to generate common sense inferences for previously unseen events.
- This effort represents a major push toward trying to capture common sense knowledge, but still relied on crowd-sourcing to populate the information in the knowledge graph.
- Nevertheless, this is a proof-of-concept showing that a model can be trained to generate such inferences for new events.
- Results are still coarsely evaluated (precision@10), although it is an open question of how to properly evaluate common sense knowledge!

33



31



32

4/

**Using Commonsense Knowledge to Assess Semantic Plausibility**

- [Wang et al., NAACL 2018] investigated the problem of recognizing the *semantic plausibility* of novel events. For example:  
*the boy swallowed a bottlecap* → semantically plausible  
*the boy swallowed a desk* → NOT semantically plausible
- Many unusual events are perfectly understandable to people, even though they are unlikely to appear in a text corpus so our models may conclude that they are impossible.
- NLP challenge: can we develop systems that can distinguish between an event that is impossible from an event that is entirely possible but hasn't been seen before!

1

**Semantic Plausibility ≠ Selectional Preferences**

- There has been prior research on learning the selectional preferences for verbal arguments.
- Refresher:** selectional preferences characterize the types of semantic entities expected by a verb's argument. For example, consider *eat*: its agent is typically ANIMATE and its object is typically FOOD.
- Semantic plausibility goes beyond selectional preferences, including things that are possible even if atypical.

man-swallow*	PREFERRED?	PLAUSIBLE?
candy	✓	✓
paintball	✗	✓
desk	✗	✗

2

**Baselines for Physical Properties**

- They experimented with an existing neural network model for identifying selectional preferences to see how well it would work for this new task.
- As input, they concatenated the 300D GloVe embeddings of the 3 words in the S-V-O triple.
- The model achieved 68% accuracy.
- The data has a roughly 50/50 split of plausible vs. implausible, so this is better than baseline but still underwhelming.
- CONCLUSION: distributional similarity alone is not sufficient.
- NEXT AGENDA: if we had commonsense knowledge of physical properties, would it help?

5

**Knowledge about Physical Attributes**

- Wang et al. argued that NLP systems need knowledge about physical properties associated with objects for this task. In particular: *sentence*, *mass-count*, *phase*, *size*, *weight*, and *rigidity*.
- Why? Because these properties often characterize the types of objects that can be sensibly involved in an action. For example:  
*throw X* → X must be solid  
*X eats* → X must be sentient  
*put X in Y* → X must be smaller than Y

3

**Gold Plausibility Annotations**

- They used Amazon's Mechanical Turk to crowd-source gold labels.
- The authors collected 150 concrete verbs and 450 concrete nouns based on Brysbaert et al.'s (2014) word list with concreteness ratings. Five turkers then labeled S-V-O (Subject-Verb-Object) triples:
  - (a) Have Turk select plausible or implausible S-V and V-O selections;
  - (b) Randomly generate S-V-O triples from collected S-V and V-O pairs;
  - (c) Send resulting S-V-O triples to Turkers to filter for ones with high agreement (by majority vote).

The final gold data set consisted of 3,062 S-V-O triples for which at least 3 Turkers agreed on the label.

4

**Gold Physical Properties**

- To generate "gold" knowledge about physical objects, 5 Turkers labeled the 450 concrete nouns used to create the S-V-O triples.
- Given a noun and a physical property, the annotator has to decide which of the landmarks for the property the noun is closest to.
  - For example: the landmarks for SIZE are: twatch, book, cat, person, jeep, stadium.
  - Given dog, the landmark cat should be selected.
  - Given shed, the landmark jeep should be selected.

7

4/

**Using Landmarks to Represent Properties of Physical Objects**

Instead of comparing object pairs with respect to physical properties, they proposed defining **landmark categories** and binning objects. This approach is natural for people and avoids arbitrary numeric ranges or pairwise comparisons.

SENTENCE: rock, tree, ant, cat, chimp, man.
MASS-COUNT: milk, sand, pebbles, car.
PHASE: smoke, milk, wood.
SIZE: watch, book, cat, person, jeep, stadium.
WEIGHT: watch, book, dumbbell, man, jeep, stadium.
RIGIDITY: water, skin, leather, wood, metal.

For example, *dog* would be assigned to the *wood* landmark for PHASE and the *cat* landmark for SIZE.

6

**Injecting World Knowledge**

They created a neural network model for the semantic plausibility task and showed that adding features that captured the physical property knowledge improved its performance.

The diagram shows a neural network architecture for generating gold physical properties. It consists of two parallel paths: one for nouns (NN) and one for verbs (VV). Both paths have an "embeddings" layer. The NN path also includes a "binning" layer. The outputs from both paths are combined via a "concatenate" layer. The final output is passed through a "softmax" layer to produce the landmark probabilities.

8

### The WK Neural Net

- The WK (World Knowledge) NN encodes the physical properties associated with the Subject and Object. They experimented with 2 ways of encoding the knowledge for each property.
- 3-LEVEL** represents whether the Subject is  $<$ ,  $=$ , or  $>$  the Object.
- $f_{\text{BIN}}(\text{SIZE}(s), \text{SIZE}(o)) \in \{-1, 0, 1\}$
- BIN-DIFF** captures the relative distance between the landmark categories for the Subject (S) and Object (O):  

$$f_{\text{BIN}}(\text{SIZE}(s), \text{SIZE}(o)) = \text{BIN}(s) - \text{BIN}(o)$$

Example: consider  $\text{SIZE}(\text{watch}, \text{book}, \text{cat}, \text{person}, \text{jeep}, \text{stadium})$   
ant would be labeled as  $\text{watch}$ ; man would be labeled as  $\text{person}$   
 $\rightarrow \text{BIN-DIFF}(\text{ant}, \text{man}) = 1 - 4 = -3$

9

### Experimental Results

- They present results for 10-fold cross-validation for the semantic plausibility task.
- For comparison, they show a Random model, Logistic Regression (LR), and the Selectional Preference NN alone.
- The BIN-DIFF encoding outperformed 3-Level.

MODELS	ACCURACY
Random	0.50
LR baseline	0.64
NN (Van de Crux, 2014)	0.68
NN + WK-GOLD	<b>0.76</b>

*Physical object knowledge clearly improves performance!*

10

### Commonsense Knowledge of Quantities

- How much does a housecat weigh?
- How tall is a housecat?
- When do people typically eat breakfast?
- When do people typically sleep?
- How long are basketball games?
- What is the typical temperature on Christmas?
- How big is a ball? A house?
- How much does a ball cost? A house?

11

### How Large are Lions?

- [Elazar et al., ACL 2019] tackled the problem of learning “quantitative attributes” for words in their paper “How Large Are Lions? Inducing Distributions over Quantitative Attributes”.
- They developed an unsupervised method to extract quantitative information from large text corpora and coalesce the results into reliable data.
- This paper includes valuable discussions about general challenges in acquiring commonsense knowledge from text corpora:
  - dealing with noisy data because IE from text is never perfect
  - challenges with reporting bias

12

### Task

- Acquire quantitative distributions over 10 dimensions: **TIME, CURRENCY, LENGTH, AREA, VOLUME, MASS, TEMPERATURE, DURATION, SPEED, and VOLTAGE**.
- Distributions are learned for:
  - Nouns** For example: elephant, airplane, NBA game
  - Adjectives** For example: cold, hot, lukewarm
  - Verbs** For example: eating, walking, running
- The resulting resource is called **Distributions over Quantities (DoQ)** and contains over 350k entries (triples) that were each observed at least 1,000 times.

13

### Extracting Measurements

- They wrote a context-free grammar to identify measurement expressions with a parser.
- They also created a mapping table between units and dimensions.
- Examples: **inch**  $\rightarrow$  LENGTH   **acre foot**  $\rightarrow$  VOLUME
- The mapping table also defines each unit in terms of a standard unit for normalization purposes. For example, all TEMPERATURE mentions are normalized in terms of degrees Kelvin and SPEED mentions are normalized in terms of meters per second.
- Examples:
  - inch** = 0.02524 meters
  - acre foot** = 1233.48 cubic meters

15

### Extracting Measurement Information

- They use a rule-based approach to detect and extract measurement information from a large text corpus.
- First, they extract all measurement mentions that they can detect. Units serve as key anchors!
- Second, they associate the measurements with nearby objects and aggregate the results.
- They intentionally aimed for a simple approach that requires only shallow resources so that this method can be applied to different languages.

14

### Extracting Objects

- All nouns, adjectives, and verbs are extracted as the objects of measurements. (NOTE: “object” in this paper means the target of the measurement, not syntactic object.)
- For each 1-word object, they also look to create a multi-word object by extracting its syntactic head. If its head is also a targeted POS, then a 2-word phrase is created from both words.
- Example: “**the fast car was driving at 100 mph**.”
- 1) **“fast”**, **“car”**, and **“driving”** are each extracted as objects.
- 2) **“fast car”** is also extracted because the head of fast is car (noun).

NOTE: there are no details about the syntactic heads in the paper, but I’d guess that the head needs to be a noun or maybe adjective.

16

### Generating Measurement/Object Pairs

- They collected **billions** of English web pages and set up a framework that allows for parallel processing.
- Extract and normalize measurement phrases.
- Apply a POS tagger and dependency parser.
- Extract object phrases that occur within close proximity (same sentence or within a distance threshold) to a measurement.
- Discard all measurements that occur in the same sentence as a negation word, because determining the precise scope of negation is a non-trivial problem.
- Aggregate all instances with the same object head and measurement unit to obtain a distribution over values.

17

### Dataset Statistics

The table below shows the number of distinct object/measurement tuples that were extracted:

Filter/Type	Nouns	Adjectives	Verbs
none	117,953,900	2,513,033	2,121,448
5	16,188,215	598,563	603,799
100	1,497,753	130,534	160,060
1000	266,655	40,518	51,625

has at least this number of occurrences

18

### Mass Distributions for Animals

These are “violin” plots. The white dot is the median. Wider areas represent values with a higher probability and skinnier areas represent lower probabilities

19

### Distributions over Hours of the Day

20

### Speed Distributions for Car Modifiers

21

### Comparison with [Forbes & Choi, 2017]

- They discovered problems with the original [Forbes & Choi, 2017] set, so they re-labeled the data via a new crowd-sourcing process.
  - Their version of the data set is substantially smaller (because they tossed problematic cases), but overall they claim it has higher quality labels.
- Applying their DoQ to this data produces better results than the previous state-of-the-art system [Yang et al., 2018].

Model/Dataset	F&C Clean		New Data	
	Dev	Test	Dev	Test
Majority	0.54	0.57	0.51	0.50
Yang et al. (PCE LSTM)	0.86	0.88	0.85	0.87
DoQ	0.78	0.77	<b>0.62</b>	<b>0.62</b>
DoQ + 10-distance	0.78	0.77	<b>0.62</b>	<b>0.62</b>
DoQ + 3-distance	0.81	0.80	<b>0.62</b>	0.61

23

### Evaluation

- They evaluated the quality and utility of the DoQ by applying it to several existing datasets and also with an intrinsic evaluation.
- The first evaluation focuses on the task of identifying the relative physical relationships between object pairs (as in [Forbes & Choi, 2017]). Formally:
 

Given  $O_1$  and  $O_2$ , predict the relation  $\{<, =, >\}$
- To apply the DoQ, they look up the distributions for the designated property for  $O_1$  and  $O_2$  and compare their medians.

22

### Similar Results on Size Data

They also evaluated their DoQ on a data set produced by [Bagherinezhad et al., 2016] of 486 object pairs labeled with respect to physical size.

Model	Accuracy
Chance	0.5
Bagherinezhad et al.	0.835
Yang et al. (Transfer)	0.858
DoQ	0.872
DoQ + 10-distance	<b>0.877</b>
DoQ + 3-distance	0.858

24

### Adjective Evaluation

- Prior work has focused on learning the relative intensities of adjectives, so they exploited some of that existing data to evaluate the quality of their adjective results.
- They collected adjectives that had comparative intensity labels and manually assigned the appropriate type of dimension. For example:
  - Hot** and **Cold** are not comparable
  - Cold < Frigid** : TEMPERATURE dimension
  - Tiny < small** : SIZE dimension
- They also adopted a different procedure for applying the DoQ, because adjectives can lead to very different inferences depending on what they are modifying. For example, consider:
  - small dog** vs. **small car**

25

### Comparing Adjective Distributions

**Algorithm 1** Adjectives Comparison Inference

```

Input: adjectives  $x, z$ , dimension  $d$  and
object distributions  $H$ 
Output: comparison label
Procedure:
Initialize  $\hat{y}$ , the predictions per head
 $intersect \leftarrow \text{findHeadIntersection}(H, x, z, d)$ 
 $\triangleright$  the intersecting heads of  $x$  and  $z$ 
for  $a_i, b_i \in intersect$  do
   $\hat{y}_i \leftarrow \text{compare}(a_i, b_i, d)$ 
end for
Return majority( $\hat{y}$ )

```

26

### Results on Scalar Adjective Dataset

The table below shows results for the adjectives evaluation across 3 different data sets.

Model	deMelo	Wilks-intense	Wilks-all
Global Ranking	<b>0.642</b>	0.818	-
Cocos et al.	0.600	<b>0.841</b>	-
DoQ	0.617	0.700	0.870
DoQ + 10-distance	0.608	0.750	<b>0.891</b>
DoQ + 3-distance	0.567	0.500	0.761

27

### Intrinsic Evaluation

A sample of the DoQ was manually annotated to directly assess the quality of its entries.

Mass	Length	Speed	Currency	All
.61	.79	.77	.58	.69

The annotators were from India, and the authors wondered if the low Currency results could be from cultural differences in perceived prices. So they had U.S. annotators re-label Currency. Their agreement = 76%!

**Example:** Indian annotators said that a suit could not cost between \$1K-\$10K, while U.S.-based annotators reported it was possible.

We tend to forget about cultural differences when we crowd-source data but it can be a significant issue!

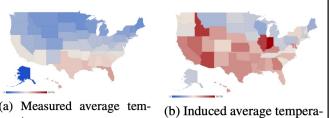
28

### Reporting Bias Discussion

- They thought that their technique would be relatively robust to reporting bias issues because they used a massive text collection and because they focused on numeric measurement data.
- However they observed that:
  - People tend to discuss objects when they are exceptional (e.g., *I saw an extremely tiny horse*).
  - People sometimes exaggerate measurements for rhetorical effect. For example, they found that people tend to exaggerate hot temperatures more than cold temperatures.

29

### Reporting Bias Example

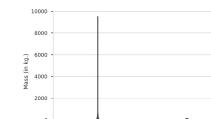


(a) Measured average temperature  
(b) Induced average temperature

Figure (b) shows that people tend to talk about or exaggerate hot temperatures more than cold ones!

30

### Reporting Bias Example



These temperature patterns reflect a reporting bias for the Northern Hemisphere!

31

Alfalfa is typically mentioned in farming contexts, so harvesting large amounts (tons). In contrast, watermelon is typically mentioned in terms of individual items so small units (grams).

32

### Conclusions

- [Elazar et al., 2019] showed that a simple, unsupervised approach for extracting measurements from a large corpus can produce a resource that compares favorably against methods that require more resources and offer less coverage.
  - Extracting massive volumes of data in a simple way and taking care to clean it properly can be remarkably effective sometimes!
- However, the quality of the data is still mixed if you look through the resource. The notion of “object” probably should be refined.
- Also, the representation of ranges is still relatively crude.
- Even for this relatively straightforward kind of data, reporting bias is still an issue! This is an on-going challenge for nearly all problems related to harvesting commonsense knowledge.

33

### Commonsense Knowledge Base Completion

- In recent years, some researchers have tackled the problem of **commonsense knowledge base completion (CKBC)**, which is the task of automatically expanding an existing common sense KB.
- These models typically use supervised learning and are trained with an existing KB such as ConceptNet. A held-out sample of the KB is reserved as a test set.
- These methods have produced evaluation scores that can look quite impressive.
- However, a deeper dive into the results has shown that they are primarily rephrasing the relations found in the training set. **Training set leakage** is a common problem.

This begs the question: is supervised learning really the answer?

1

### Training Set Leakage Problem

- Training set leakage refers to situations where the training and test sets accidentally share information.
  - Information is revealed to the model that gives it an unrealistic advantage to make better predictions
- As a result, performance on the test set is artificially high! Applying the system to a truly blind test set yields much lower results.
- Some examples of how this can happen:
  - A single document is split across the training and test sets.
  - Duplicates may exist (this is very common with Twitter data).
  - New articles from same time period.
  - Very similar cases exist (e.g., paraphrases). This case is hard to avoid, and likely a common reason why different data sources usually perform worse.
  - In general, train/test splits share more vocabulary than new data sets.

2

### Types of Neural Language Models

Causal language models process an input sequence from left-to-right and predict the next word in the sequence.

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n p(w_i | w_1, \dots, w_{i-1})$$

GPT (Generative Pre-trained Transformer) models are a family of neural LMs trained on massive text collections.

Masked Language Models (MLMs) are a new breed of neural LMs that are bidirectional and trained to predict a small number of words that have been “masked”:
 
$$p(w_i | w'_1, w'_2, \dots, w'_{i-1}, w'_{i+1}, w'_n)$$

$$w' \in V \setminus \{w\}$$
 where  $w$  is a special Mask token.

- The MLM returns a probability distribution over words in the Mask position.
- BERT (Bidirectional Encoder Representations from Transformers) is a MLM that is widely used and trained on massive text corpora.

3

### Recap: Masked Language Modeling

- Masked Language Models (MLMs) are trained by “masking” a position of the input.
- Causal transformer models are trained with a **causal mask** to predict the next word:
  - Peanut butter and \_\_\_\_\_
- Bidirectional transformer models are trained with a **cloze test**, which is a “fill in the blank” mask:
  - Peanut butter and \_\_\_\_\_ sandwich
- To train with a MLM, many input tokens are sampled. For BERT, 80% are then masked, 10% are randomly replaced, and 10% are left unchanged.

4

**Harvesting Commonsense Knowledge from Language Models**

- Pre-trained language models (LMs) have become the backbone of many NLP models recently, to produce contextual embeddings that are used as rich semantic representations for downstream tasks.
- Researchers have also begun to explore whether large language models can be useful for recognizing commonsense knowledge.
- A key observation is that a LM can estimate the probability of a sentence based on its training corpus. **The likelihood of the sentence can serve as a proxy for the truth of the knowledge.**

Dogs are mammals. ✓ (high prob)  
Dogs can run. ✓ (high prob)  
Dogs can fly. ✗ (low prob)

5

**Exploiting Masked Language Models**

- [Feldman et al., EMNLP 2019] explored the idea of using Masked Language Models (MLMs) to validate commonsense knowledge facts.
- Given a relational triple  $\langle h, r, t \rangle$ , they present an approach for using a MLM to estimate the likelihood that the relation is a valid sentence (and therefore true knowledge).
- This approach could be used to validate any relational triple, but their focus is on evaluating commonsense knowledge relations.
- Similar ideas have been explored for assessing the plausibility of a relation, for tasks like textual entailment.

6

**Representing Facts**

- Each candidate fact is represented as a **head-relation-tail triple**:  $x = \langle h, r, t \rangle$
- They assume that there is a fixed set of known relations  $R$ , so  $r \in R$ .
- The head and tail can be multi-word phrases, where each word comes from a known vocabulary  $V$ .  
 $h = (h_1, h_2, \dots, h_n) \quad t = (t_1, t_2, \dots, t_m)$
- The goal is to learn a function  $f(x) = y$  such that  $y$  reflects the model's confidence that  $x$  represents true knowledge.

7

**Generating Sentences**

Language models expect sentences as input, so each triple needs to be converted into a grammatical sentence. But this is not as simple as it sounds! For example, consider:

(ferret, AtLocation, pet store)

Possible sentences:  
A ferret is at a pet store.  
The ferret is at the pet store.  
The ferret is in the pet store.  
The ferret is located in the pet store.

Possible sentences:  
(water, AtLocation, Mars)  
A water is at Mars.  
The water is at the Mars.  
The water is in the Mars.  
The water is located in the Mars.

8

**Generating Sentences from Triples**

- They apply grammatical transformations to each head and tail:
  - If the first word is a N or ADJ, or the first word is a V and the second word is a N or ADJ, then prepend an indefinite or definite article.
  - If the first word is an infinitive form, make it a gerund (e.g., *jump* → *jumping*)
  - If the first word is a number, pluralize the next word (e.g., *two leg* → *two legs*)
- For each relation, they manually define a set of templates.

For example, the templates below are used for  $\langle X, CapableOf, Y \rangle$ :

X can Y      An activity X can do is Y  
X often Y      X sometimes Y

9

**Coherency Ranking Example**

LM scores for several candidate sentences generated from the triple: *(musician, CapableOf, play musical instrument)*

Candidate Sentence $S_i$	$\log p(S_i)$
"musician can playing musical instrument"	-5.7
"musician can be play musical instrument"	-4.9
"musician often play musical instrument"	-5.5
"a musician can play a musical instrument"	<b>-2.9</b>

11

**Coherency Ranking**

- All possible combinations of candidate sentences are generated.
- Then they select the candidate sentence with the highest log-likelihood score from a unidirectional LM:  $P_{\text{coh}}$ .

$$S^* = \arg \max_{S \in S} [\log P_{\text{coh}}(S)]$$

The expectation is that natural, grammatical sentences will have higher LM scores than unnatural or ungrammatical sentence.

In practice, they found that this approach produced "significantly higher quality" results than using deterministic rules alone.

10

**Scoring Sentences**

- Once the best sentence has been identified, it needs to be scored to assess its validity.
- A sentence is scored based on the PMI of the head and tail conditioned on the relation.

$$\text{PMI}(t, h|r) = \log p(t|h, r) - \log p(t|r)$$

Ultimately they use a weighted PMI measure, where the lambda ( $\lambda$ ) value is a tuned hyperparameter.

$$\text{PMI}_\lambda(t, h|r) = \lambda \log p(t|h, r) - \log p(t|r)$$

PMI is symmetric in theory, but this approximation is not. So they average  $\text{PMI}_\lambda(t, h|r)$  and  $\text{PMI}_\lambda(h, t|r)$ .

12

**PMI Redux**

$$\text{PMI}(X, Y) = \log_2 \left[ \frac{P(X, Y)}{P(X) * P(Y)} \right]$$

$$\text{PMI}(t, h | R) = \log_2 \left[ \frac{P(t, h | R)}{P(t|R) * P(h|R)} \right]$$

$$\log_2 \left[ \frac{P(t, h | R)}{P(t|R) * P(h|R)} \right]$$

$$\log_2 \left[ \frac{P(t | h, R)}{P(t|R)} \right]$$

$$= \log_2 (P(t | h, R)) - \log_2 (P(t | R))$$

13

**Estimating Probabilities from a MLM**

- We can estimate the probability of the tail with a single mask position in a bidirectional MLM model ( $P_{\text{cmp}}$ )

$$P(t | h, R) = P_{\text{cmp}}(w_t = t | w_{1:i-1}, w_{i+1:m})$$

- If the tail is a multi-word term, then they use a greedy approach to estimate a probability for the phrase. For a phrase with  $j$  terms:
  - Mask each word, one at a time, and compute the probability for each.
  - Select the word with the highest probability and insert it in the phrase
  - Repeat this process  $j$  times.
- Finally, the probability for the phrase is the product of the probabilities found for the individual terms:

$$p(t|h, r) = \prod_{k=1}^j p_k$$

14

**Task 1: Commonsense Knowledge Base Completion**

- Use test set from (Li et al., 2016), which contains 2400 triples that contain an equal number of **Valid/Invalid triples**.
- The Valid triples come from the crowd-sourced Open Mind Common Sense (OMCS) entries in the ConceptNet 5 dataset.
- The Invalid triples were produced by replacing one of the elements in a valid triple with a randomly selected item.
- The Coherency Rank model is applied to produce a score for each triple and the triples are grouped into two clusters based on their scores. All of the triples in the cluster with the highest mean PMI are labeled as **Valid**.

17

**Task 2: Mining Wikipedia**

- (Li et al., 2016) curated a data set mined from Wikipedia primarily using part-of-speech patterns: 1.7M triples across 10 relations.
- For this work, 300 triples were sampled for each relation, producing a test set of 3,000 relation triples.
- The Coherency Rank method scored each triple, and the top-scoring 100 triples were manually reviewed by 2 human annotators on a scale from 0 to 4:
  - 0: doesn't make sense
  - 1: not true
  - 2: opinion / don't know
  - 3: sometimes true
  - 4: generally true

NOTE: kappa agreement was only .23! But if only two buckets are used, then disagreement drops by 50%.

18

**Estimating Probabilities from a MLM**

Computing  $P(t | R)$  is similar but it needs to be estimated over all possible heads. So the head must be masked throughout.

For example, consider the sentence:  
*You are likely to find a ferret in the pet store.*

Initially, both the head and tail are masked:  
*You are likely to find a h<sub>i</sub> in the t<sub>j</sub>.*

The tail terms are gradually inserted based on their probabilities:  
*You are likely to find a h<sub>i</sub> in the t<sub>j</sub> store.* (→  $P_{\text{store}} > P_{\text{pet}}$ )  
*You are likely to find a h<sub>i</sub> in the pet store.* (→  $P_{\text{pet}}$ )  
The final probability is then  $P_{\text{store}} * P_{\text{pet}}$

15

**Experimental Set-up**

- For sentence ranking, they use the GPT-2 Language Model. For the masked language model, they use BERT (large model).
- As a baseline for sentence generation, they simply split the words in the relation name and **concatenate** the head and tail.
- For example:  
*(ferret, AtLocation, pet store)* → *ferret at location pet store*
- They also evaluate against a Commonsense Knowledge Base Completion (CKBC) approach by instantiating a single template with the head and tail.
- For example:  
*(ferret, AtLocation, pet store)* →  
*"you are likely to find ferret in pet store"*

16

**Results**

Model	Task 1	Task 2
Unsupervised		
CONCATENATION	68.8	2.05 ± 0.11
TEMPLATE	72.2	2.98 ± 0.11
TEMPL+GRAMMAR	74.4	2.56 ± 0.13
COHERENCY RANK	78.8	<b>3.00 ± 0.12</b>
Supervised		
DNN	<b>89.2</b>	2.50
FACTORIZED	89.0	2.61
PROTOTYPICAL	79.4	2.55

For Task 1: The Coherency Rank approach performed better than the other unsupervised approaches, but not as well as the supervised CKBC systems.

For Task 2: The Coherency Rank approach outperformed all the other models, including the supervised CKBC systems.

19

They analyzed their results based on whether the sentence had a grammatical error or misrepresented the relation's meaning.

For example, *(golf, HasProperty, good)* → *"golf is a good"*  
Is grammatical but captures a different meaning.

Task 1	N (100)	F1 Score
GRAMMATICAL	75	70.1
UNGRAMMATICAL	25	66.7
CORRECT MEANING	91	77.6
WRONG MEANING	9	66.7

Task 2	N	Quality
GRAMMATICAL	83	3.01
UNGRAMMATICAL	17	2.88
CORRECT MEANING	88	3.22
WRONG MEANING	12	1.18

20

### Error Analysis

**Most Confident Mistakes:** the following triples were in the top 100 predictions but received scores < 3 by the human annotators.

- (atomic number, HasProperty, positive)
- (the harbor, HasA, island)
- (the substance, HasA, drug)
- (plurality voting, HasPrerequisite, majority)
- (function, ReceivesAction, element of a)
- (minister, ReceivesAction, member of parliament)
- (bombing, IsA, war crime)
- (prime minister, ReceivesAction, head of state)
- (film, Causes, silent version)
- (cause, ReceivesAction, element of s)
- (monarchy, ReceivesAction, form of government)
- (law, ReceivesAction, cause of action)
- (weather, UsedFor, heavy rain)
- (example, UsedFor, word processing)

21

### Conclusions

- The Coherency Rank method is a creative approach for exploiting massive pre-trained LMs to confirm or disconfirm whether a fact is likely true or not.
- This is a powerful idea, since these LMs are essentially encoding massive amounts of textual data.
  - Instead of explicitly harvesting facts from a large corpus, can we essentially probe a LM to discover what it knows?
- This approach is also exciting because it does not use supervised learning so no manually annotated data is needed.
- We've just begun to scratch the surface of how these large LMs may be used! This direction seems very promising.

22

### Prompting Methods

- Prompting methods** have become a hot area in NLP that exploit pre-trained LMs as an alternative to supervised learning.
- Large pre-trained language models can be viewed not only as a repository of language but also as a *repository of knowledge*. **Prompt-based methods** aim to extract knowledge from pre-trained LMs with carefully designed inputs that essentially query the LM for missing information.
- Typically, a *template* is combined with an input value to create a *prompt string* that has one or more unfilled positions. The language model can then return the words or phrases that are most likely to occur in those positions with probability estimates.
- The trick is to design good templates for your task!

23

### Example: Sentiment Analysis

Suppose you want to get a sense of the popularity of a movie, such as "Arrival". The typical approach would be to collect a corpus of reviews for the movie and apply a sentiment classifier, ideally one trained on movie reviews.

Alternatively, you could use prompting with the template:

"*<MOVIE> was a <MASK> movie*"

The input to the LM would be "*Arrival was a [X] movie*", and the LM would then return a probability distribution over words that could fill the X position. For example:

great (.10), terrific (.09), boring (.07), sci-fi (.06), good (.06) ...

You could then use a sentiment lexicon to assess the overall polarity.

24

### Prompt Engineering

- Designing an effective prompt for your task is key. Language models can be sensitive to the specific words used in the prompt, even punctuation can matter.
- There are different types of prompts, such as:
  - Prefix prompts:** the LM completes the input
  - Cloze prompts:** the LM fills a blank in the middle
  - Zero-shot prompts:** no examples are given
    - Example: "The capital of Utah is [X]"
  - Few shot prompts:** one or more examples are given.
    - Example: "A list of cities in Utah. 1. Salt Lake City"

25

### Prompt Ensembling

Since any one prompt may not be ideal, multiple prompts can be used in a sort of "ensemble". The results are then pooled and can be ranked based on voting or more complex ranking methods.

See "Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing" [Liu et al., 2021] for a comprehensive survey of prompting methods.

27

### Challenges with Prompting

Prompting methods can sometimes work very well and their potential is intriguing, but there are also major challenges.

- Multi-word answers can be difficult to extract with a mask.
- Hard to control syntactic expectations (many types of answers are possible due to complex syntactic constructions).
- The most common answers dominate. It can be difficult, perhaps impossible, to extract uncommon/rare information.
- Complex prompts can be challenging because the LM often focuses on the local context around the mask and may ignore the full context of the prompt.

26

### Also ... LMs alone may not be enough

[Lin et al., EMNLP 2020] recently showed that masked language models are not very good at learning numerical common sense knowledge.

BERT-Large	1st:fly (79.5%)
Birds can [MASK].	Masked Word Prediction 2nd: sing (9.1%)
However, for Numerical Commonsense Knowledge :	
A bird usually has [MASK] legs.	1st:four (44.8%)
	2nd:two (18.7%)
A car usually has [MASK] wheels.	1st:four (53.7%)
	2nd:two (20.5%)
A car usually has [MASK] round wheels.	1st:two (37.1%)
	2nd:four (28.2%)

Even when the predictions are correct, they can be brittle.

28

### LMs can have strong biases

The table below shows the top 3 predictions when the [x] variable is instantiated with 1k random words.

Template:	a [x] usually has [MASK] legs.
BERT-L	four 28.3%, two 9.5%, three 10.1%
RoBERTa-L	four 28.3%, two 9.5%, three 10.1%
Template:	most [x] have [MASK] wheels.
BERT-L	four 25.3%, two 14.1%, three 5.1%
RoBERTa-L	four 22.3%, two 7.8%, three 4.6%
Template:	all [x] have [MASK] sides.
BERT-L	two 28.3%, three 12.9%, four 12.9%
RoBERTa-L	two 16.6%, no 2.9%, three 2.3%

There is often one dominant value that seems to emerge based on the most common cases.

29

### Summary: Commonsense Knowledge

- We have a long way to go to be able to **accurately** acquire the **vast** amounts of commonsense knowledge that are needed.
- We need tons of knowledge, but it also **must** be organized and represented in a useful way.
- However, there are many promising avenues for acquiring specific types of knowledge, automatically and semi-automatically.
  - Researchers nearly always focus on fully automatic techniques, but in the real-world semi-automatic techniques can be extremely valuable!
    - A small amount of human effort to "curate" automatically extracted information can insure high-integrity data and often produce substantially more data than manual efforts would.

30

### Summary: Information Extraction

- IE is a rich area of NLP that covers a wide variety of problems!
- Some topics correspond to fundamental aspects of text understanding, such as Named Entity Recognition, Semantic Class Learning, and Temporal IE.
- Other topics are closely tied to real-world applications, such as Relation Extraction, Event Extraction, Opinion Extraction.
- For most sequence labeling tasks, NLP models can be trained to perform reasonably well for a specific domain given sufficient domain-specific training.
- IE tasks tend to be more challenging than non-sequential classification tasks, and it is difficult to achieve high recall and high precision at the same time.
- Many challenges for the future, including cross-sentence (document-level) models, learning with small amounts of labeled data, and IE methods to harvest domain-specific knowledge!

New IE problems are always around the corner! You'll see some in the projects. ☺

31