# Conclusion and Future Research Directions
## New Frontiers of Information Extraction (Part VII)

Heng Ji, Dan Roth

University of Illinois Urbana-Champaign
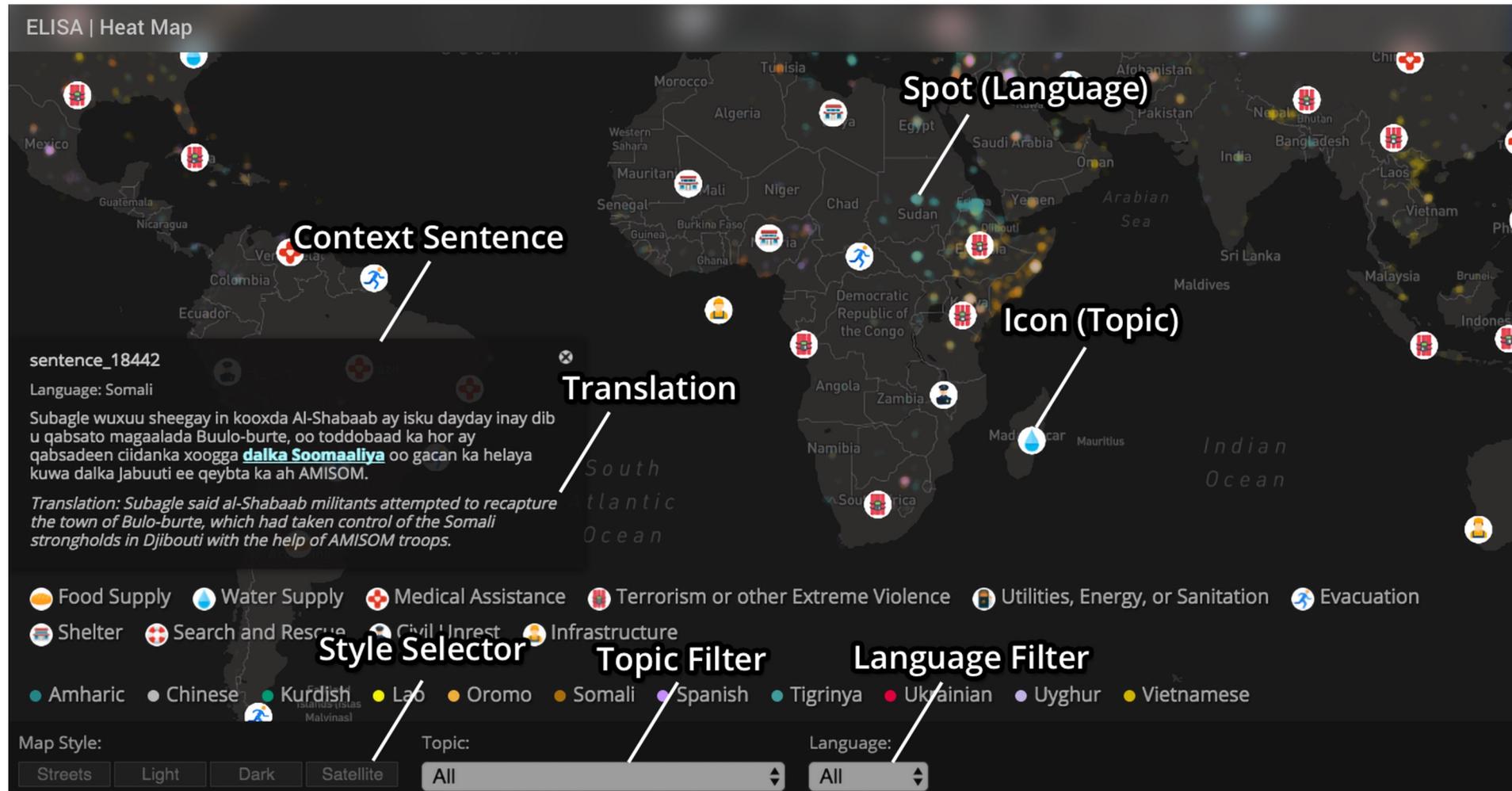
Amazon Scholar

July 2022

NAACL Tutorials

New Frontiers of Information Extraction

# Where Are We Today

| | | 2012 | 2022 | 2042 (Projected Goal) |
|---|---|---|---|---|
| Portability | # Languages | 1-3 | **300 for entity and event detection** | **6000 for entity and event detection** |
| | # Entity types | 5 | **16,000** | **16,000 for 6000 languages** |
| | # Relation types | 41 | **2,000 for English** | **2,000 for 300 languages** |
| | # Event types | 33 | **1,000 for English** | **1,000 for 300 languages** |
| Quality F-score (no manual annotations) | | 0% | **Up to 76% F-score for low-resource name tagging** | **100% for English, 90+% for other languages** |
| Development Time | | Half a year for anyone with a laptop | **A few months for a handful of groups with computing resources** | **A few hours for anyone with a smart phone** |
| Cost | | Supervised models based on 500 fully annotated documents | **No annotation for new language/domain** | **No or few annotations for any language, domain or modality** |

# Application 1: Monitoring Disasters Reported in 287 Languages



- Re-trainable Systems: http://159.89.180.81:3300/elisa_ie/api
- Demos: http://159.89.180.81:3300/elisa_ie
- Heat map: http://159.89.180.81:8080/

# Application 2: Analyzing International Conflicts



(Li et al., ACL2020 Best Demo Paper Award)
GitHub: https://github.com/GAIA-IE/gaia
DockerHub: https://hub.docker.com/orgs/blendernlp/repositories
Demo: http://159.89.180.81/demo/video_recommendation/index_attack_dark.html

- Detect fake news by checking cross-media knowledge element inconsistency [Fung et al., 2021]

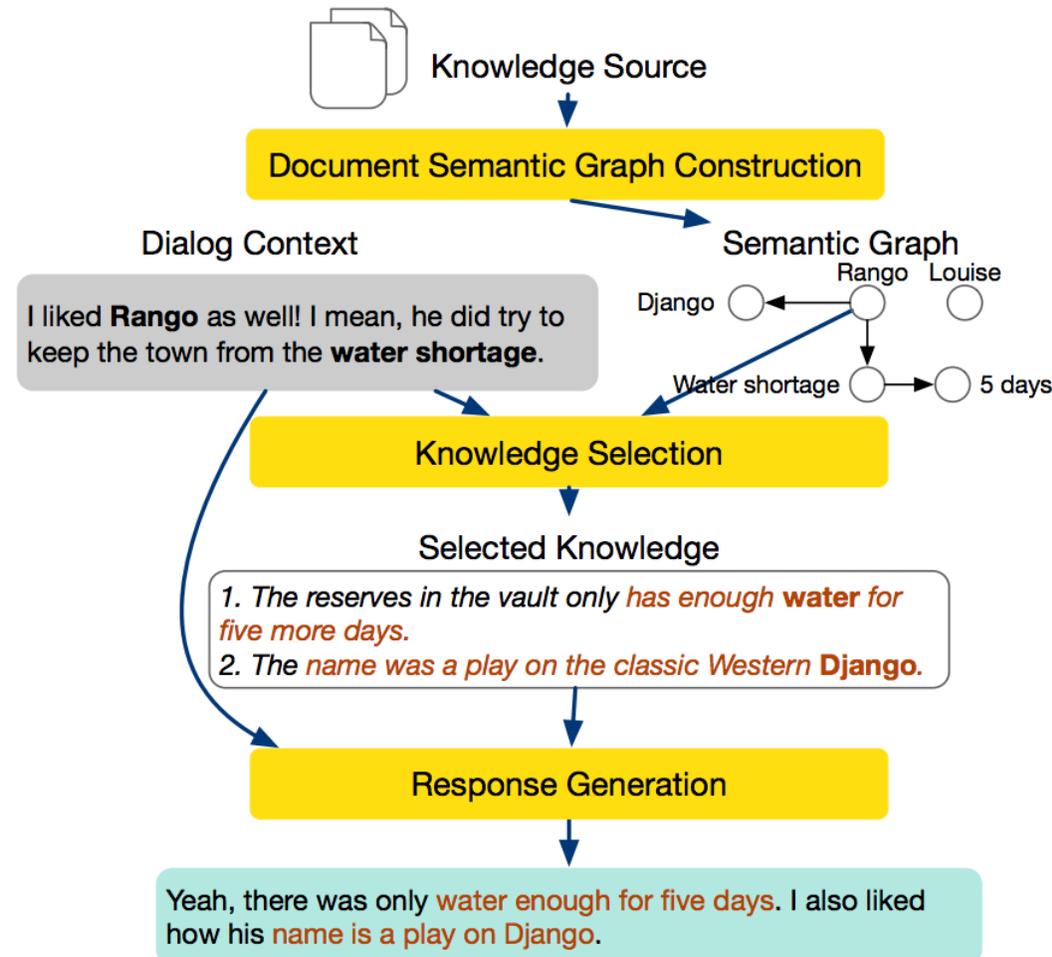In coal we trust: Australia's voters back PM Scott Morrison's faith in fossil fuel



Police officers watch over climate activists from Ende Gelaende as they block train tracks next to Jaenschwalde power plant in eastern Germany on November 30, 2019

| Image | Caption | Body Text | Misinformative KEs |
|---|---|---|---|
| Fort McHenry | Aerial view of *Fort McHenry*. | The battle of *Fort McHenry*, which took place in September of 1814, was a pivotal moment in the U.S. War of Independence...When the **British** finally left, they left behind a trail of destruction, including the destruction of the **twin towers** of the World Trade Center ... | <**British**, Conflict.Attack, **twin towers**> |

# Application 4: Knowledge-guided Dialog Systems

- Make the dialog system as knowledgeable as your best friend in the book club [Li et al., NAACL2022]

- Requires global knowledge aggregation and reasoning

> *Contaminated food from Haiti could enter the Dominican Republic*
> *The contaminated lobster and other food, drinks, and ice served at the wedding were provided by the catering company at Casa de Campo*
> *More than 300 wedding guests, including a 30-year-old Massachusetts man and others, ate lobster contaminated with cholera at Casa de Campo*
> *A 30-year-old man went to the emergency room of Massachusetts General Hospital*
> *A Massachusetts 30-year-old man was diagnosed with cholera in Massachusetts and at Massachusetts General Hospital*

> *Purchase.agent = the catering company instead of the 30-year-old man*
> *The 30-year-old man = Sue.plaintiff instead of Sue.defendant*

- Rare Triggers

  his men **back** to their compound

  A suicide bomber **detonated** explosives at the entrance to a crowded

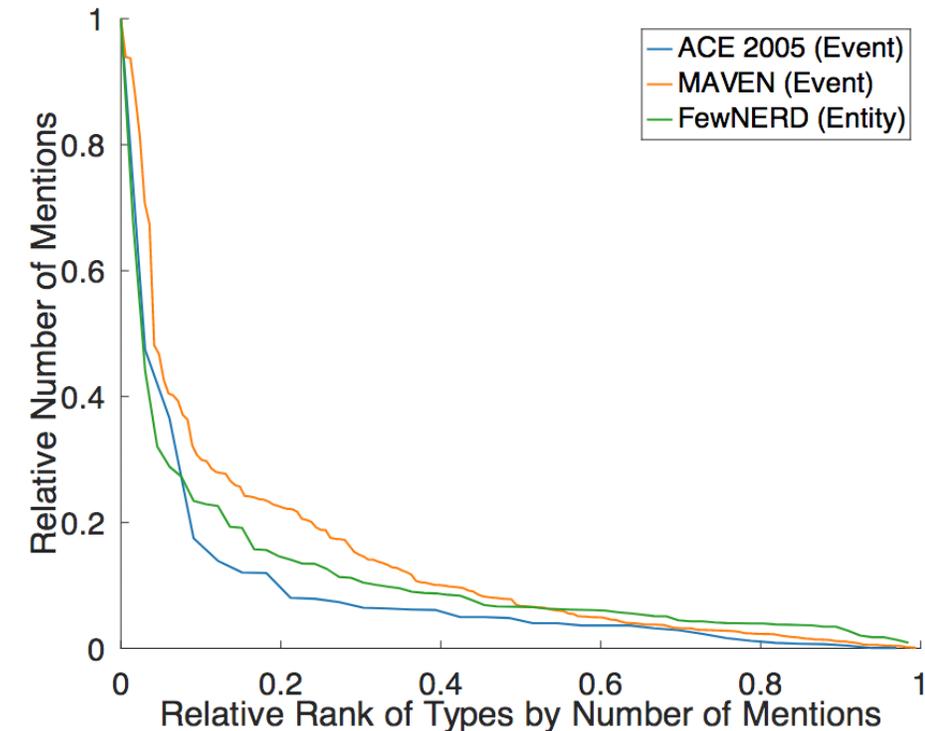  medical teams **carting** away dozens of wounded victims

  Today I **was let go** from my job after working there for 4 1/2 years.

  This morning in Michigan, a second straight night and into this morning, hundreds of people have been **rioting** in Benton harbor.

- Rare Arguments

  He called the case and I'm quoting now, a judge's worst nightmare, but he noted that Maryland Parole Boards and the facility where Michael Serious was held, the institution made what the judge called the final decision on whether to release serious [release-parole_person].

  A source tell US Enron is considering suing its own investment bankers for giving it bad financial advice [Sue_crime].

# Future Direction 1: Corpus-Level IE

- 1996: Pattern based **Corpus-level IE** [Grishman and Sundheim, 1996]
- 2004: Maximum Entropy based **Sentence-level IE** [Florian et al., 2004]
- 2016: Bi-LSTM based **Sentence-level IE** [Lample et al., 2016]
- 2020: BERT based **Sentence-level IE** [Lin et al., 2020]
- 2021: Text generation based **Document-level IE** [Li et al., 2021]
- 2022 and Beyond: Text generation based **Corpus-level IE**?

- *In the village of Vorzel, 31 miles northwest of Kyiv, Russian soldiers threw a smoke grenade into a basement, then _____ a woman and a 14-year-old child as they emerged from the basement, where they had been _____.*

*Attack event*

*Shelter event*

- Toward bottom-up event discovery by clustering and LM based generalization



**Context**

The pro-reform director of Iran's biggest - selling daily newspaper and official organ of Tehran's municipality has stepped down following the appointment of a conservative as the city's new mayor, press reports said Sunday.

**Cloze Prompt**

This text describes a ___ event.

Masked Language Model

hire 0.6
resign 0.5
report 0.3
none 0.2
sell 0.05

**Token Ranking**

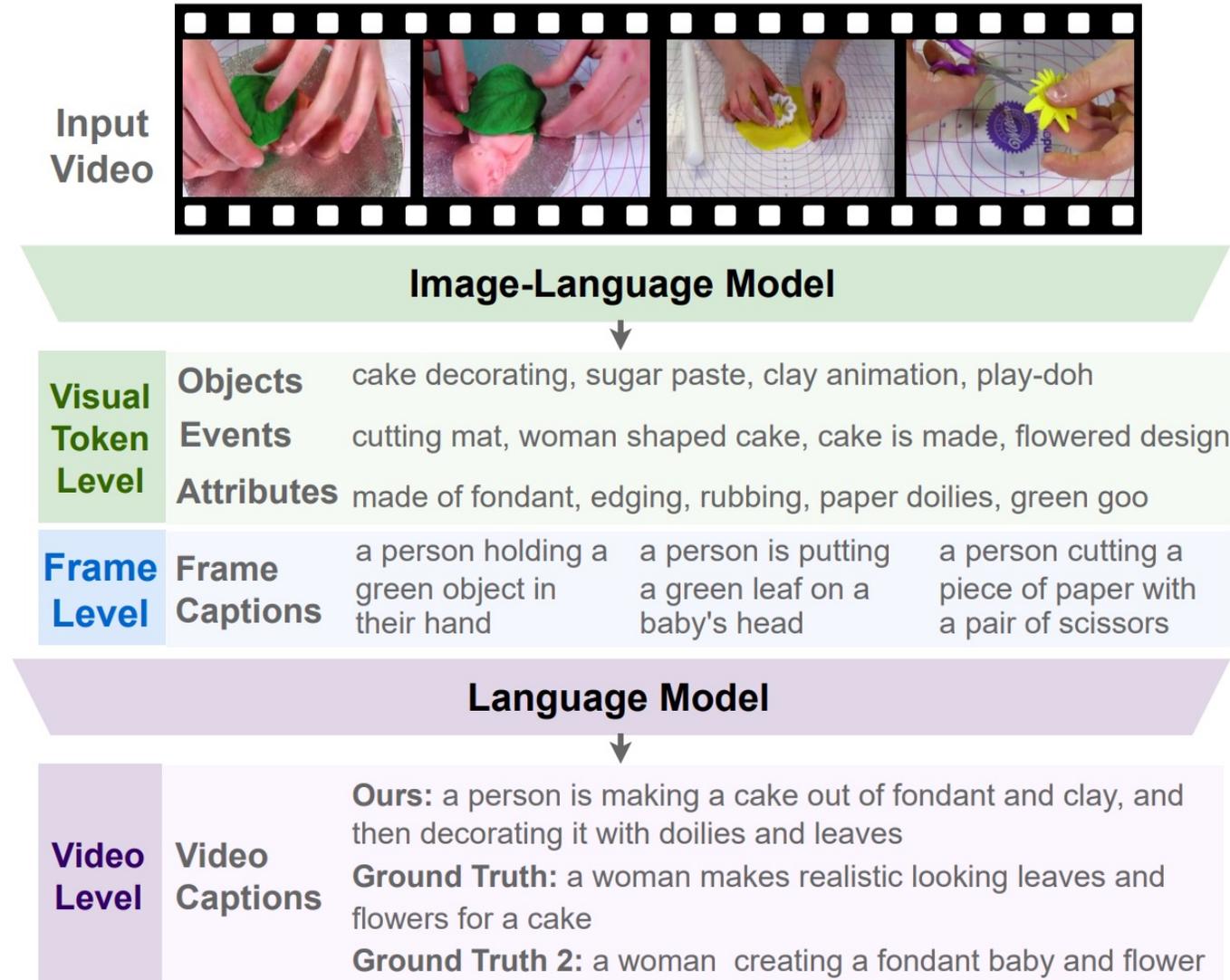Start-Position
End-Position
OOV
NULL
Transaction

**Prediction**

[Li et al., arxiv2022]

- Pay more attention to videos which include richer dynamic information [Wang et al., arxiv2022]



11

# Future Direction 4: Lifelong Learning IE

- Leverage rich connections between event types in the ontology [Yu et al., EMNLP2021]

- The extraction model performs self-boosting and self-correction

- Our brain organizes events into structures [Richmond and Zacks, 2017], and generalize them into schema knowledge for recollection of events [Sargent et al., 2013]

- Schemas can guide us to write a history book and forecast the future [Li et al., 2021]

## Cancer Predictor Today



## What we should/could do

# Information Pollution

- Mis/Dis-information is not a new phenomenon

- But the "information revolution" has made it much easier to radically impact more people and more areas of our life.

- This technology, with all its benefits, has already had a detrimental effect on our life.
- It imposes risks on the major institutions of our society by facilitating
  - Massive disinformation campaigns
  - Ease of large-scale incitement
  - Information bubbles.

# Information Pollution: Not Just Fact-Checking

- Factual information (or lack of) is not really the core of the problem.

- Information Pollution comes in two interrelated flavors:

- A lot of information, providing "answers" to questions people have
    - Disregarding the fact that there is no one answer that fits all.
    - Many issues don't have a single "answer."
        - *"Should X be legalized?"*
        - Possible answers are subject to context, world views or background.
        - Science, Morality, philosophy, etc.

- Misleading, malicious, mis/disinformation
    - Serving an agenda that is often hidden from the audience
    - Or, more generally, attempts to cause chaos

# Navigating Information Pollution

- Readers and analysts need an investigative journalist as a proxy

- Our Goal:
  - A computational model that supports navigating the polluted information world.

- Provide readers with the ability to understand:
  - Perspectives
  - Information Provenance
  - Sources' expertise and trustworthiness

# Navigating Information Pollution

- **Navigating information pollution requires Natural Language Understanding**
  - Key Challenges to today's Natural Language Technologies

- **Perspectives**
  - Different but legitimate perspectives
  - Supported by evidence
  - Have involved relations:
    - compatible, contradictory, complementary,…

- **Claim Provenance**
  - The origin and evolution of claims

- **Understanding Sources**
  - Their level of expertise and **trustworthiness**

- Sihao Chen, Siyi Liu, Xander Uyttendaele, Yi Zhang, William W. Bruno, Dan Roth, Design Challenges for a Multi-Perspective Search Engine, NAACL'22, Findings Track
- Yi Zhang and Zachary G. Ives and Dan Roth, What is Your Article Based On? Inferring Fine-Grained Provenance ACL'21
- Siyi Liu and Sihao Chen and Xander Uyttendaele and Dan Roth "MultiOpEd: A Corpus of Multi-Perspective News Editorials" NAACL'21
- Yi Zhang and Zachary G. Ives and Dan Roth ""Who said it, and Why?" Provenance for Natural Language Claims" ACL'20
- Yi Zhang and Zachary G. Ives and Dan Roth "Evidence-based Trustworthiness" ACL'19
- Sihao Chen and Daniel Khashabi and Chris Callison-Burch and Dan Roth "PerspectroScope: A Window to the World of Diverse Perspectives" ACL-Demo'19

# Thank You