

Multimedia IE

New Frontiers of Information Extraction (Part V)

Manling Li

Department of Computer Science

University of Illinois Urbana-Champaign

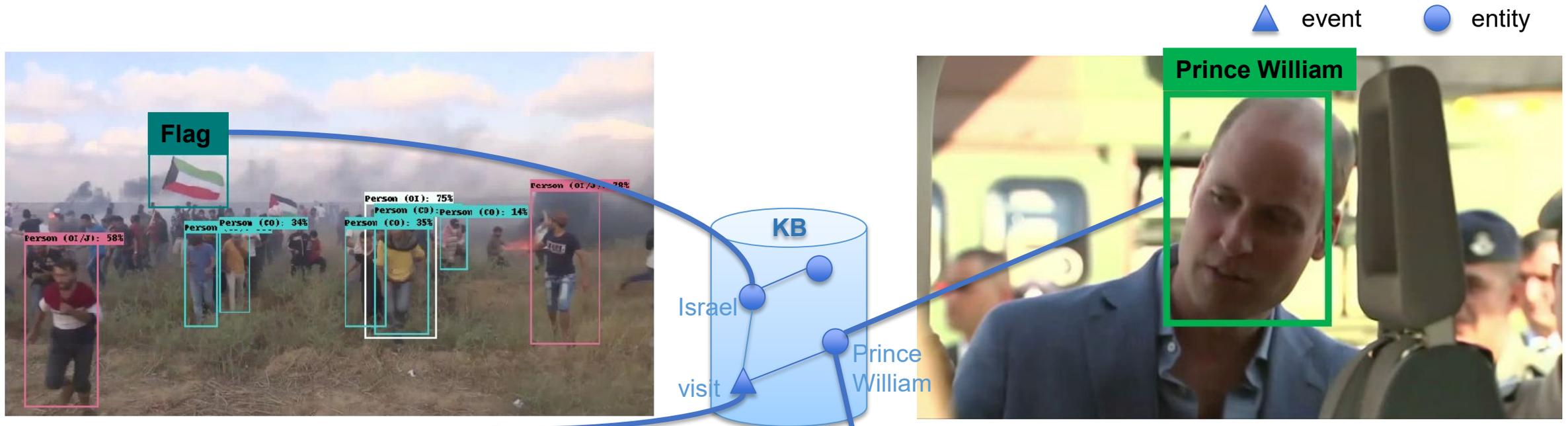
July 2022

NAACL Tutorials

New Frontiers of Information Extraction



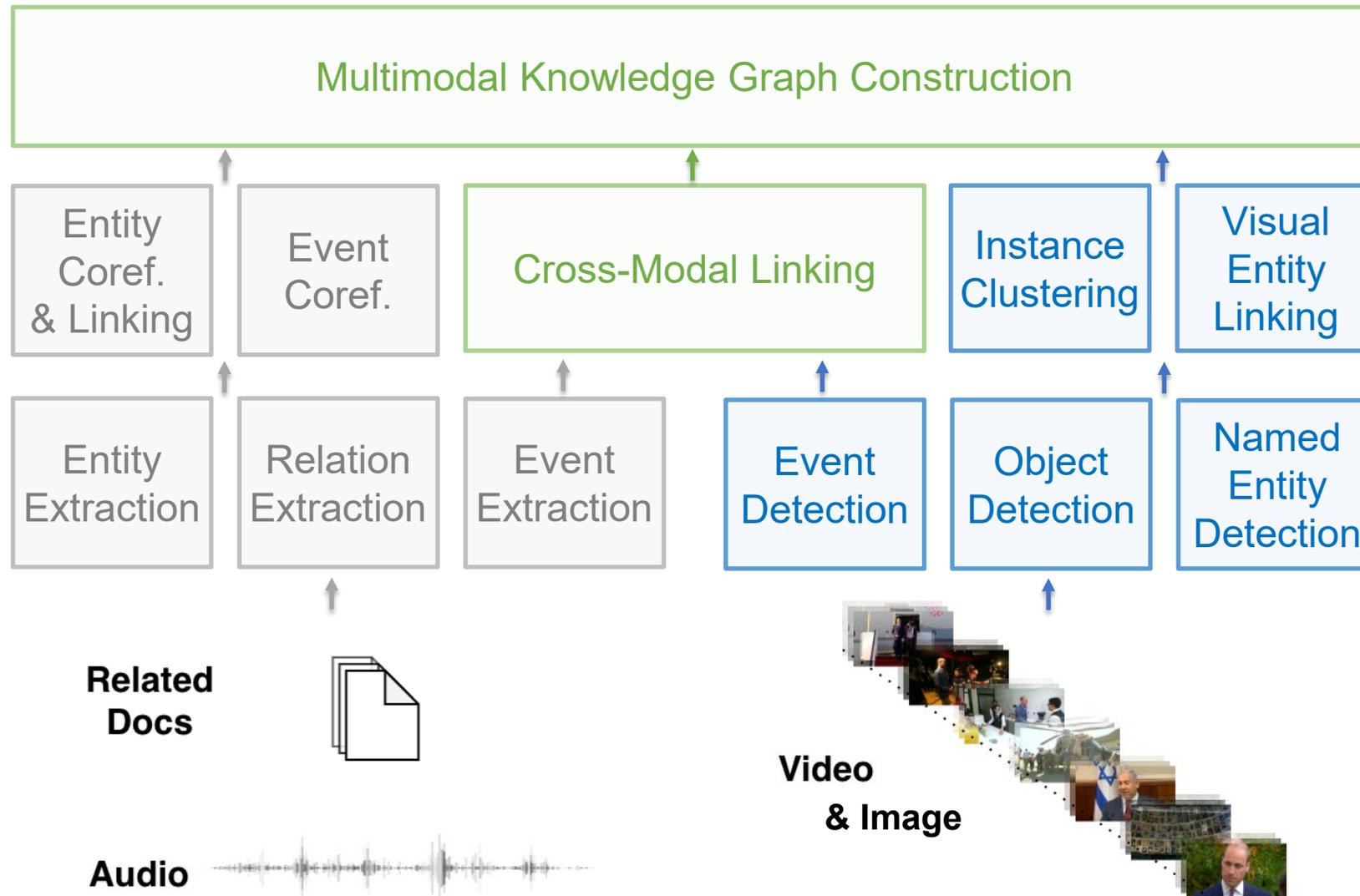
- Multimedia Knowledge Base with entities, relations and events.



The first-ever official **visit** by a British royal to **Israel** is underway. **Prince William** the 36 year-old Duke of **Cambridge** and second in line to the throne will **meet** with both **Israeli** and **Palestinian** leaders over the next three days.

Contact.Meet_Participant

Multimedia Information Extraction



- Treat Image/Video as a foreign language

Text	Image / Video Frame
Sentence	Image
Word	Image Region
Entity	Visual Object
Relation	Visual Relation
Entity-Relation Graph	Visual Scene Graph
Event Trigger	Visual Activity
Event Structure	Image Event Graph

O
b
j
e
c
t

Car



E
v
e
n
t

Car

Event

Bombing

Item

Car



A
r
g
u
m
e
n
t



Event

Attacking



Attacker

protesters

Target

police

Event

Attacking

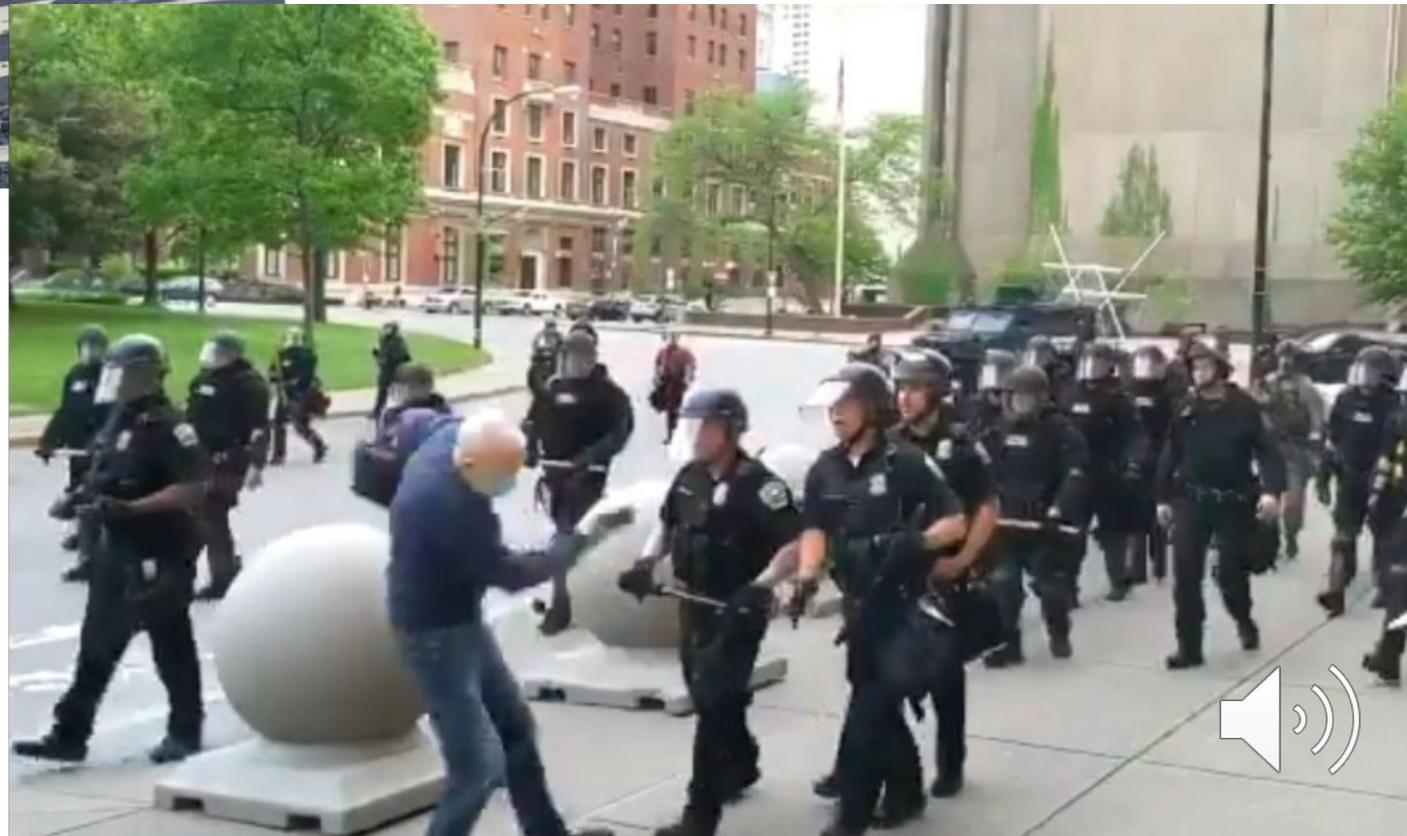
Attacker

police



Target

protester



Event **Wearing**

Item mask

Agent person



Event **Treatment**

Agent doctor

Target patient



Event **Researching**

Agent researcher

Target dropper



Event **Sanitizing**

Agent person

Tool sprayer



Event **Testing**

Agent woman

place car



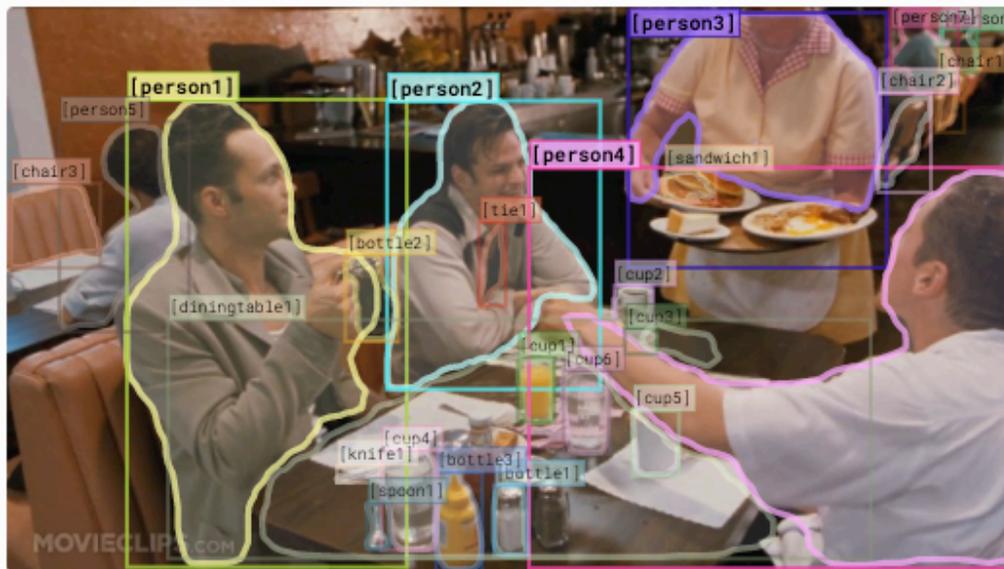
Event **Vaccination**

Agent woman

Target girl



- Real-world multimedia applications requires image-language models to understand multiple levels of alignments such as **events**, **objects**, as well as **semantic structures**.



hide all show all [person1] [person2] [person3] [person4]

[person5] [person6] [person7] [tie1] [bottle1]

[bottle2] [bottle3] [cup1] [cup2] [cup3] [cup4]

[cup5] [cup6] [knife1] [spoon1] [sandwich1] [chair1]

[chair2] [chair3] [diningtable1]

Why is [person4] pointing at [person1]?

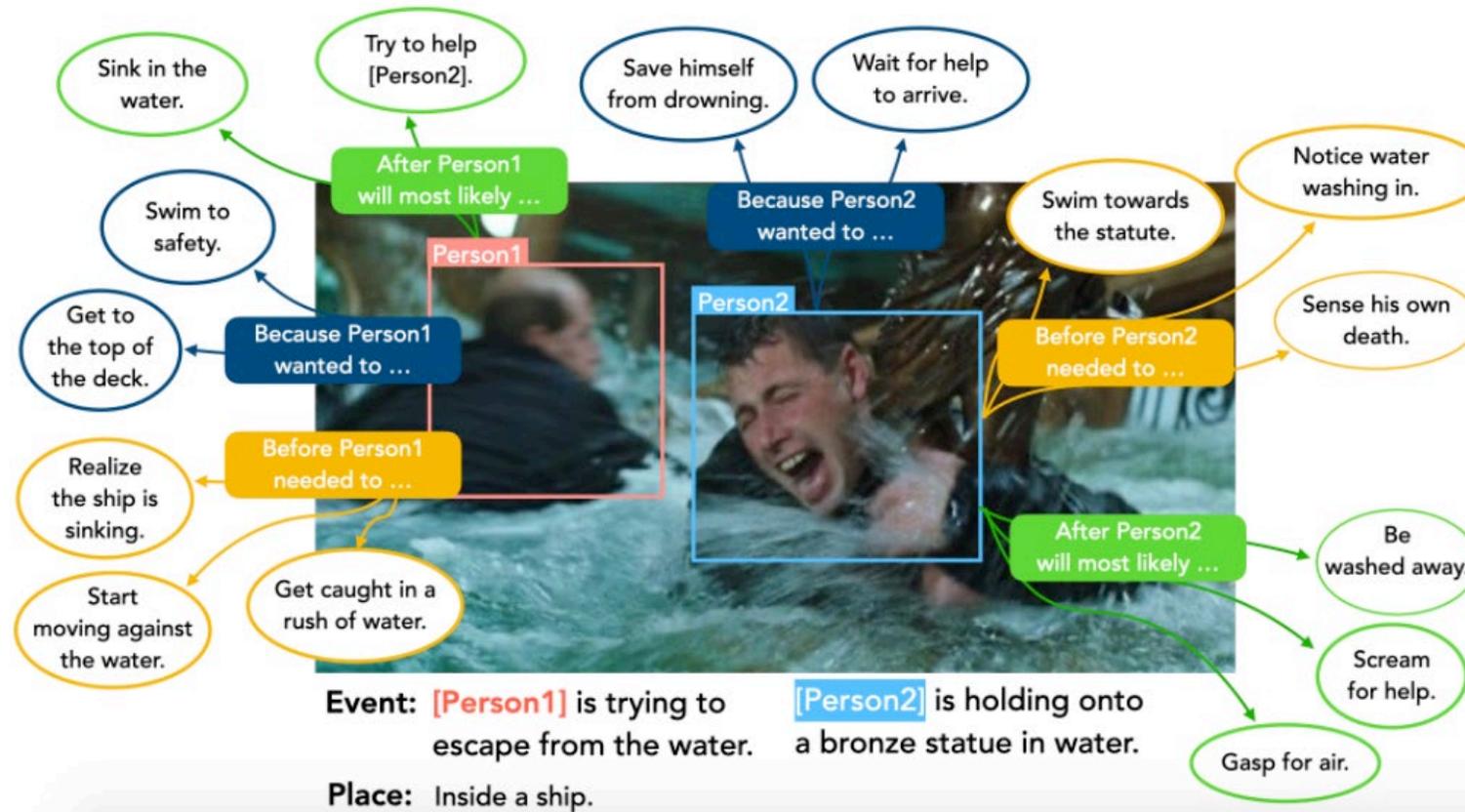
- a) He is telling [person3] that [person1] ordered the pancakes.
- b) He just told a joke.
- c) He is feeling accusatory towards [person1].
- d) He is giving [person1] directions.

Rationale: I think so because...

- a) [person1] has the pancakes in front of him.
- b) [person4] is taking everyone's order and asked for clarification.
- c) [person3] is looking at the pancakes both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.

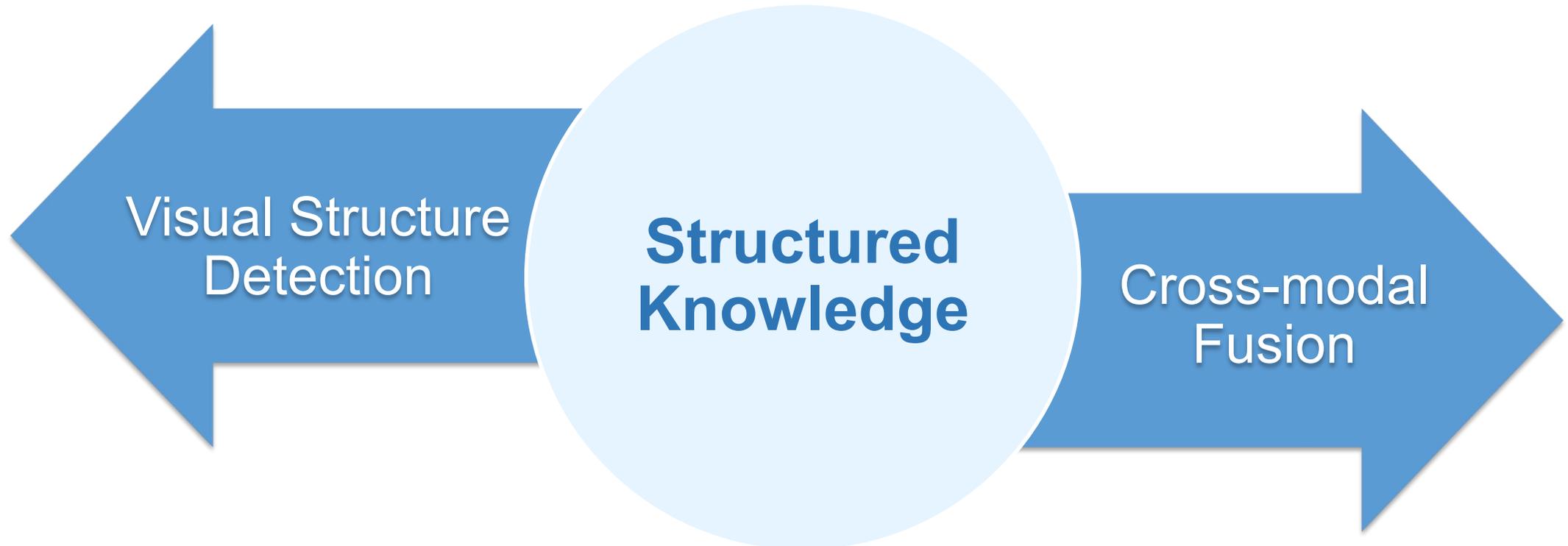
Visual Commonsense Reasoning

- Real-world multimedia applications requires image-language models to understand multiple levels of alignments such as **events**, **objects**, as well as **semantic structures**.

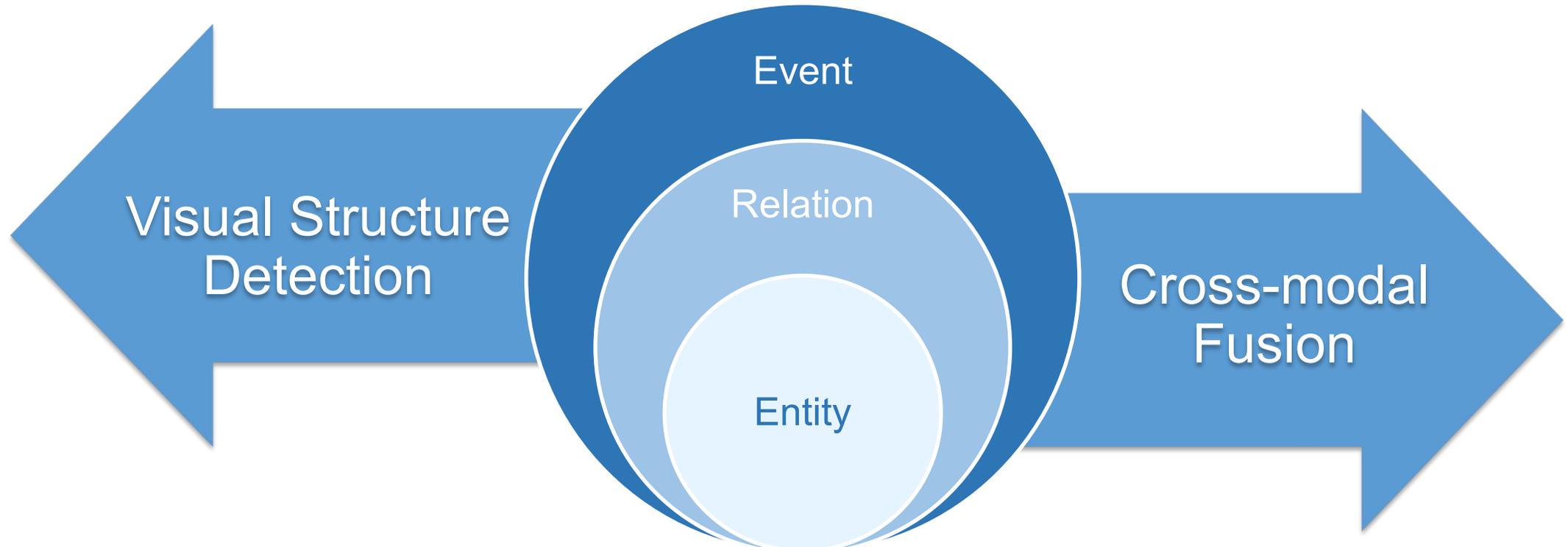


VisualCOMET (Visual Commonsense Reasoning in Time)

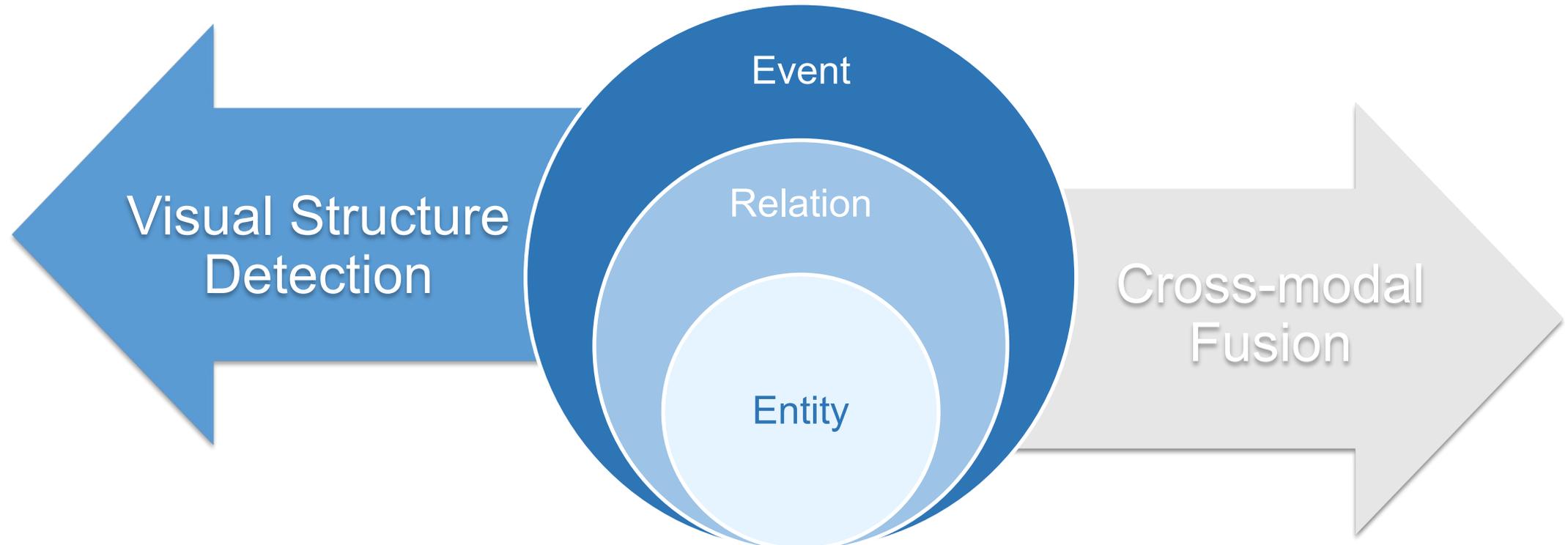
- 1. How do we find structured knowledge in vision data?
- 2. How do we align structured across vision and text?

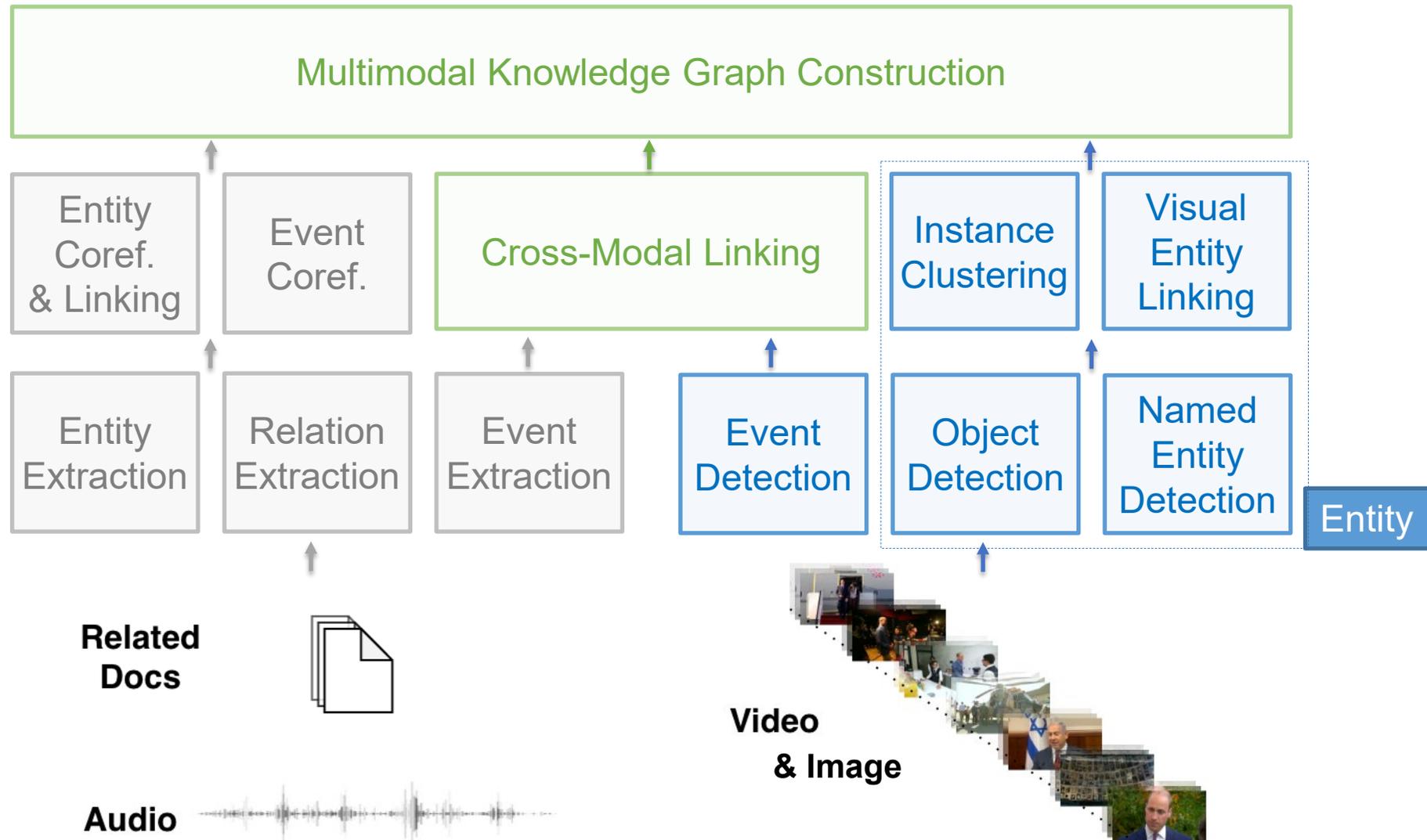


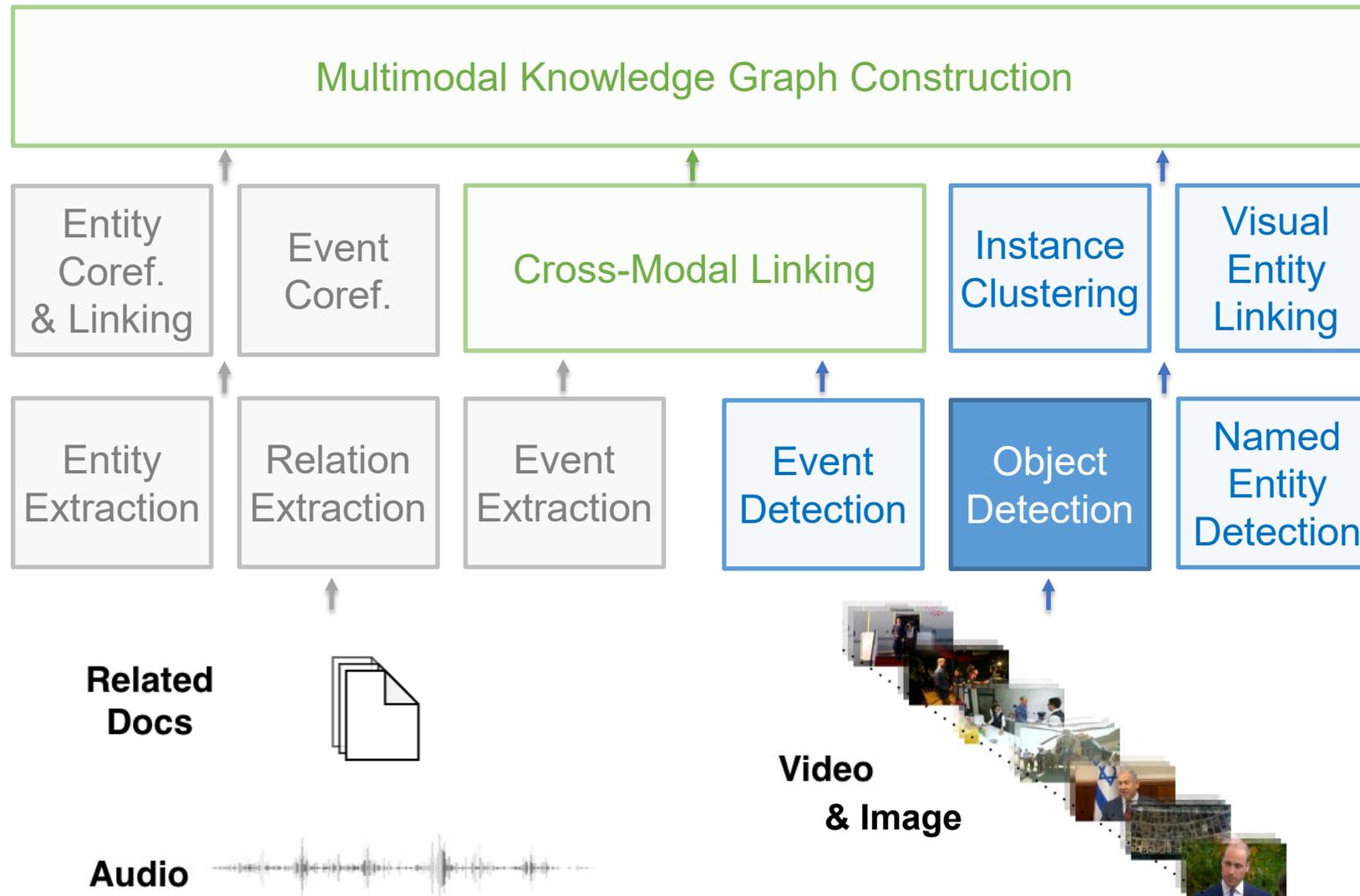
- 1. How do we find structured knowledge in vision data?
- 2. How do we align structured across vision and text?



- 1. How do we find structured knowledge in vision data?
- 2. How do we align structured across vision and text?





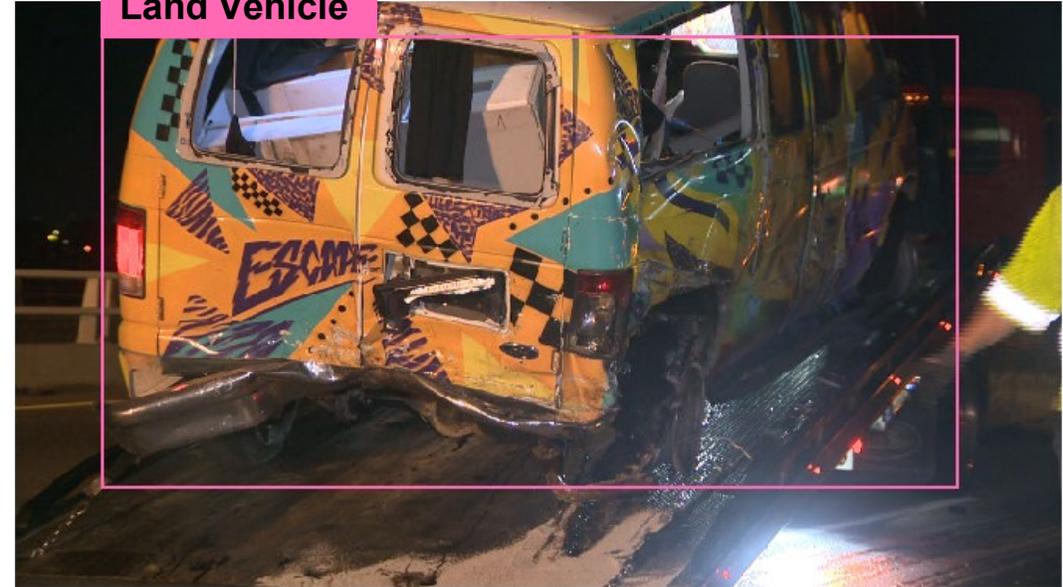


Visual Entity Extraction: Object Detection

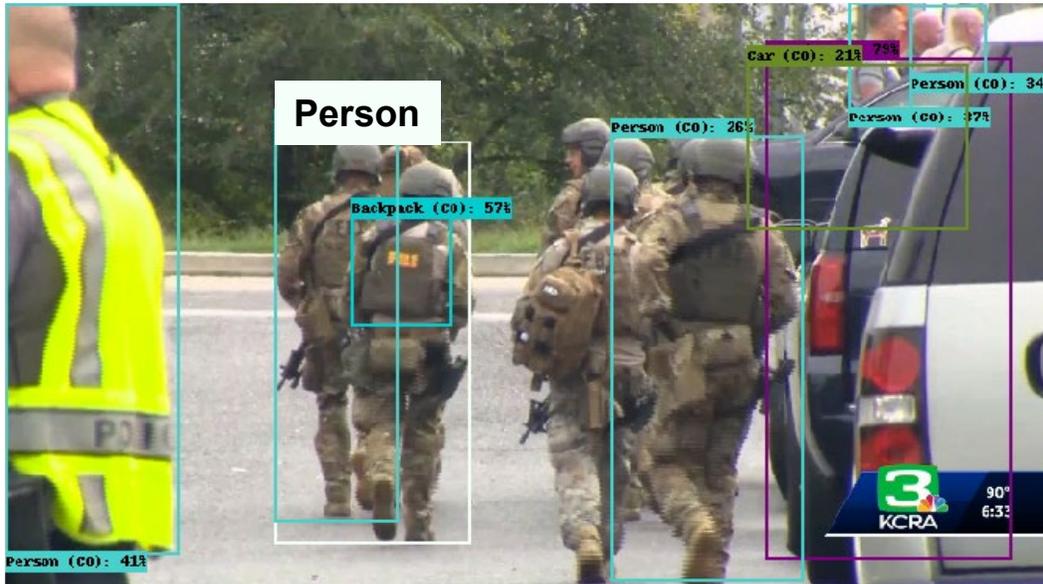
Helicopter



Land Vehicle



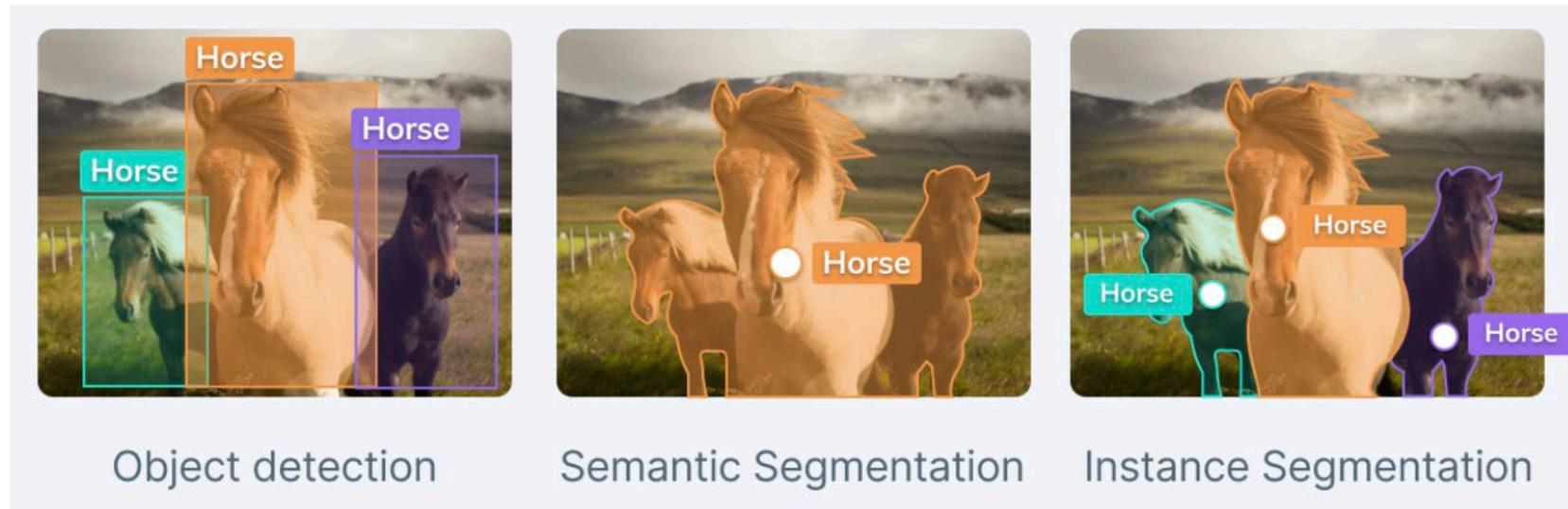
Person

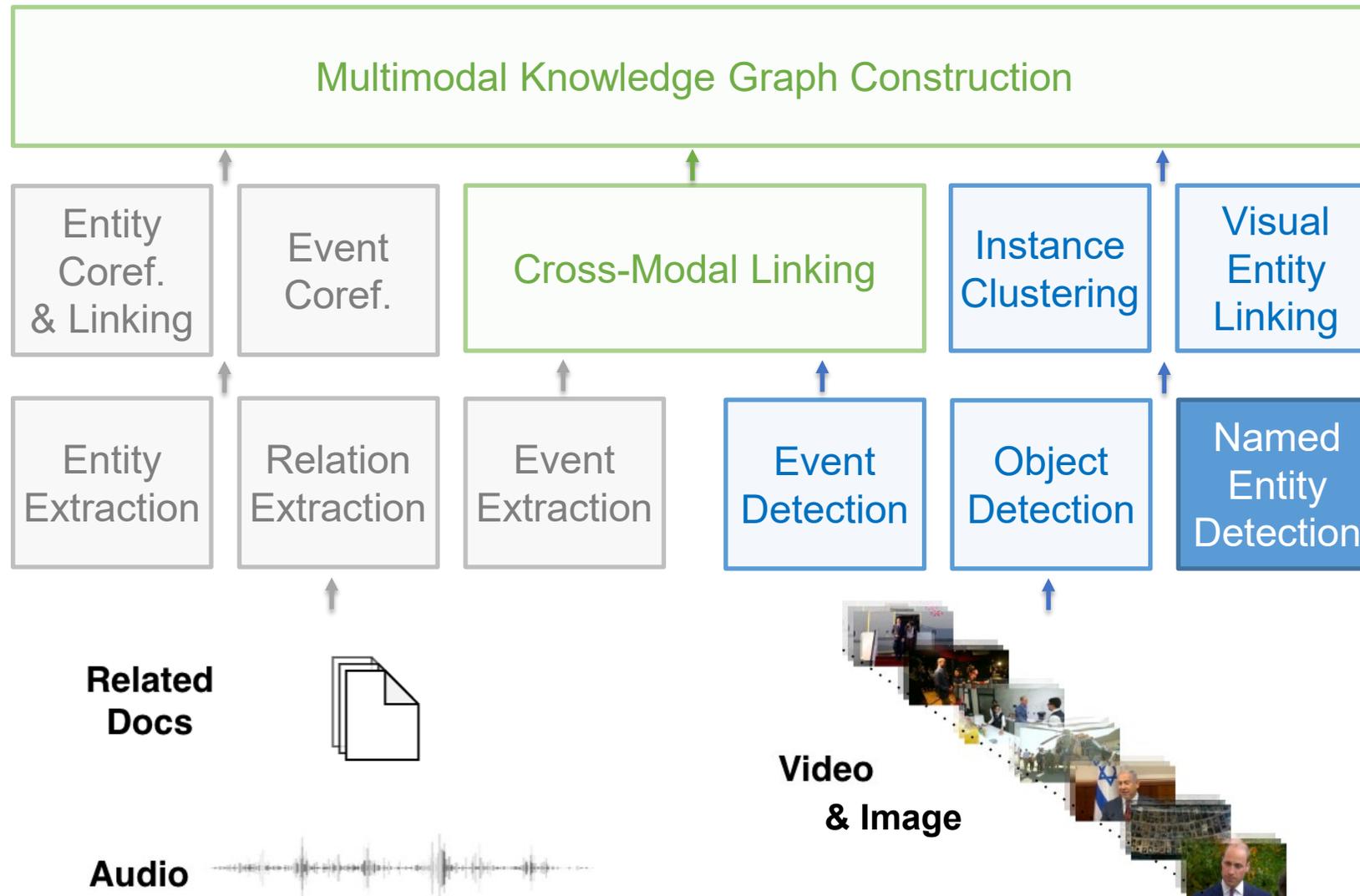


Car



- Object Detection: Object instances at the bounding box level
- Semantic Segmentation: Object class at the pixel level
- Instance Segmentation: Object instances at the pixel level





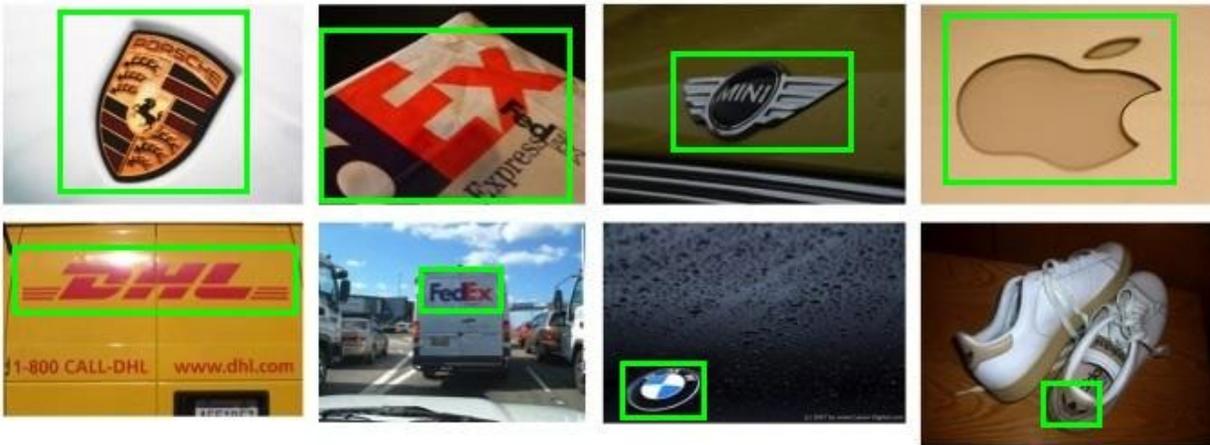
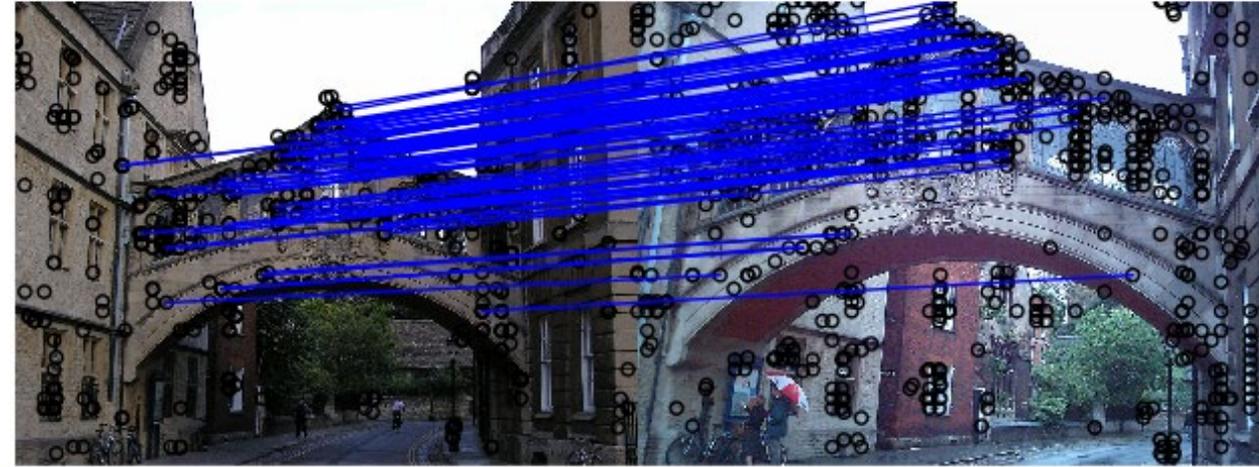
Visual Entity Linking: Named Entity Recognition



Named Entity Recognition

- Face Recognition
- Landmark Recognition
- Flag Recognition
- Logo Recognition
- ...

DELF correspondences



Schroff, Florian, Dmitry Kalenichenko, and James Philbin. "Facenet: A unified embedding for face recognition and clustering." *CVPR* 2015.

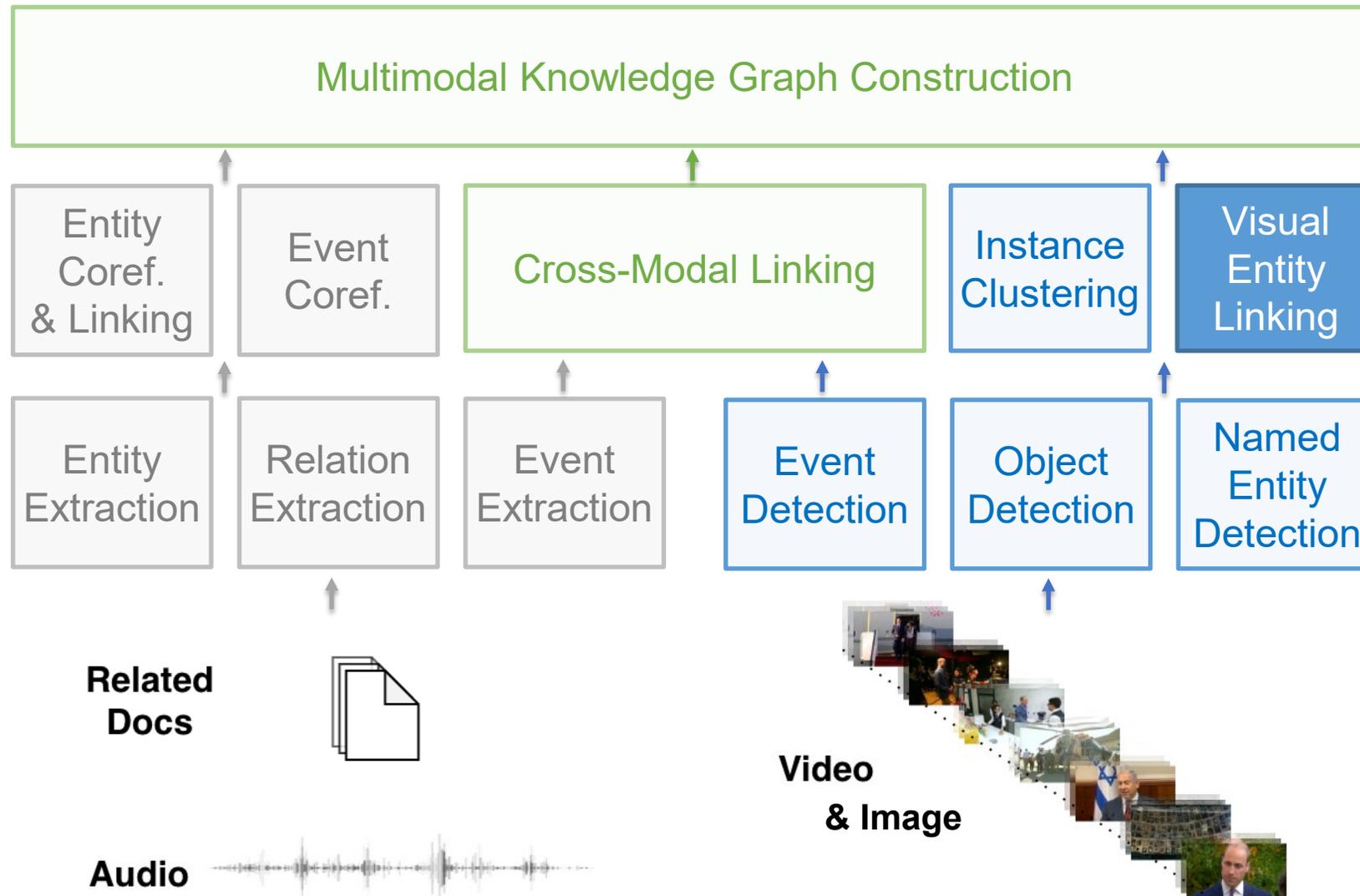
Weyand, Tobias, et al. "Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval." *CVPR* 2020.

Wu, Shou-Fang, et al. "FlagDetSeg: Multi-Nation Flag Detection and Segmentation in the Wild." *AVSS* 2021.

Bianco, Simone, et al. "Deep learning for logo recognition." *Neurocomputing* 245 (2017): 23-30.



Multimedia Information Extraction



Visual Entity Linking and Coreference

- Landmark recognition



Visual Entity Linking and Coreference

- Flag recognition

US Flag



Missed detection



Different shape and angle
Partial observe

Ukraine Flag



Low resolution

Russia Flag (X)

US Flag

UK Flag



Visual Entity Linking and Coreference

- Face detection
- Entity linking and coreference

  Unknown people



Kirstjen Nielsen



Joe Biden

Vladimir Putin

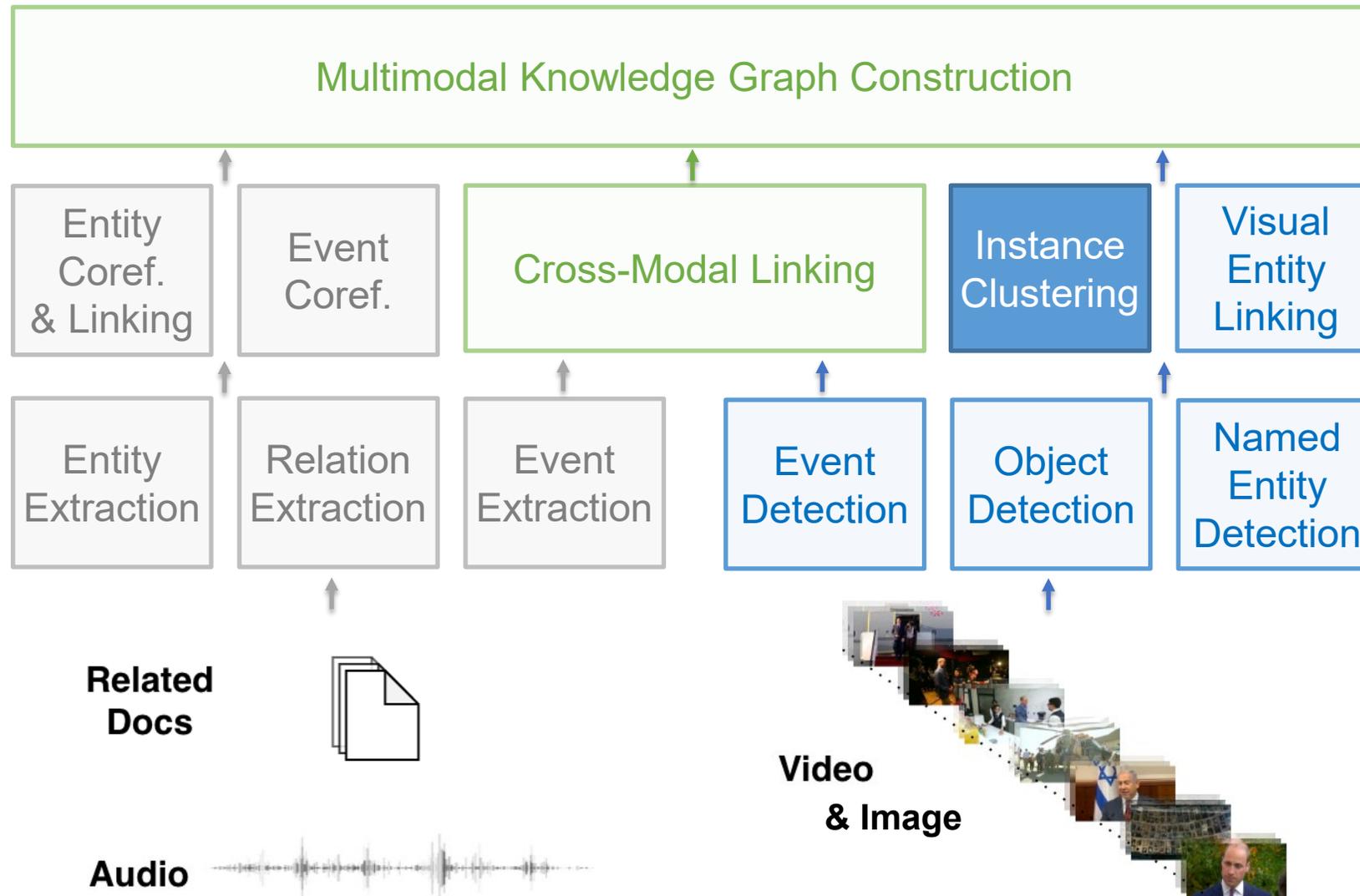


Joe Biden



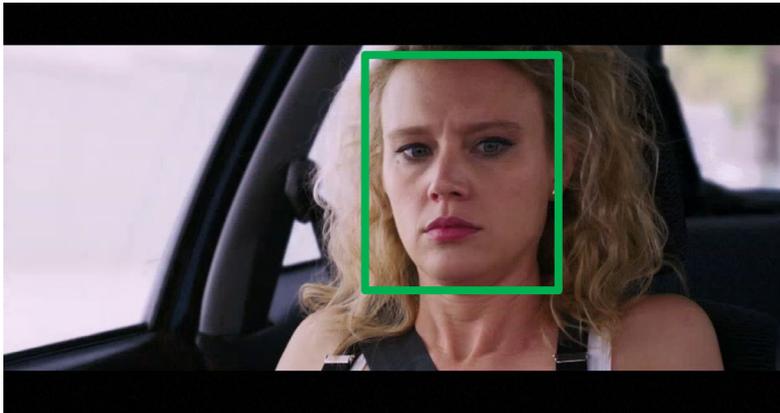
Rudy Giuliani

Multimedia Information Extraction

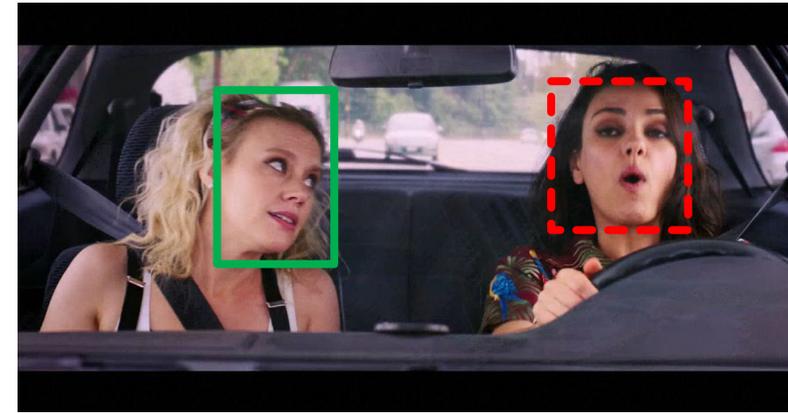


Visual Entity Linking and Coreference

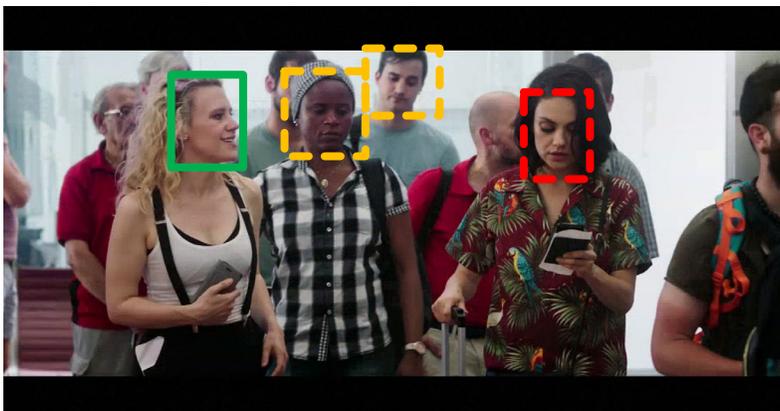
- Face detection
- Entity linking and coreference



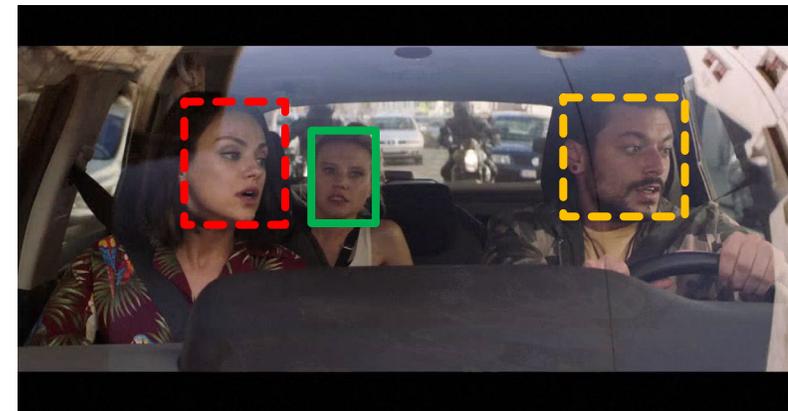
Kirstjen Nielsen



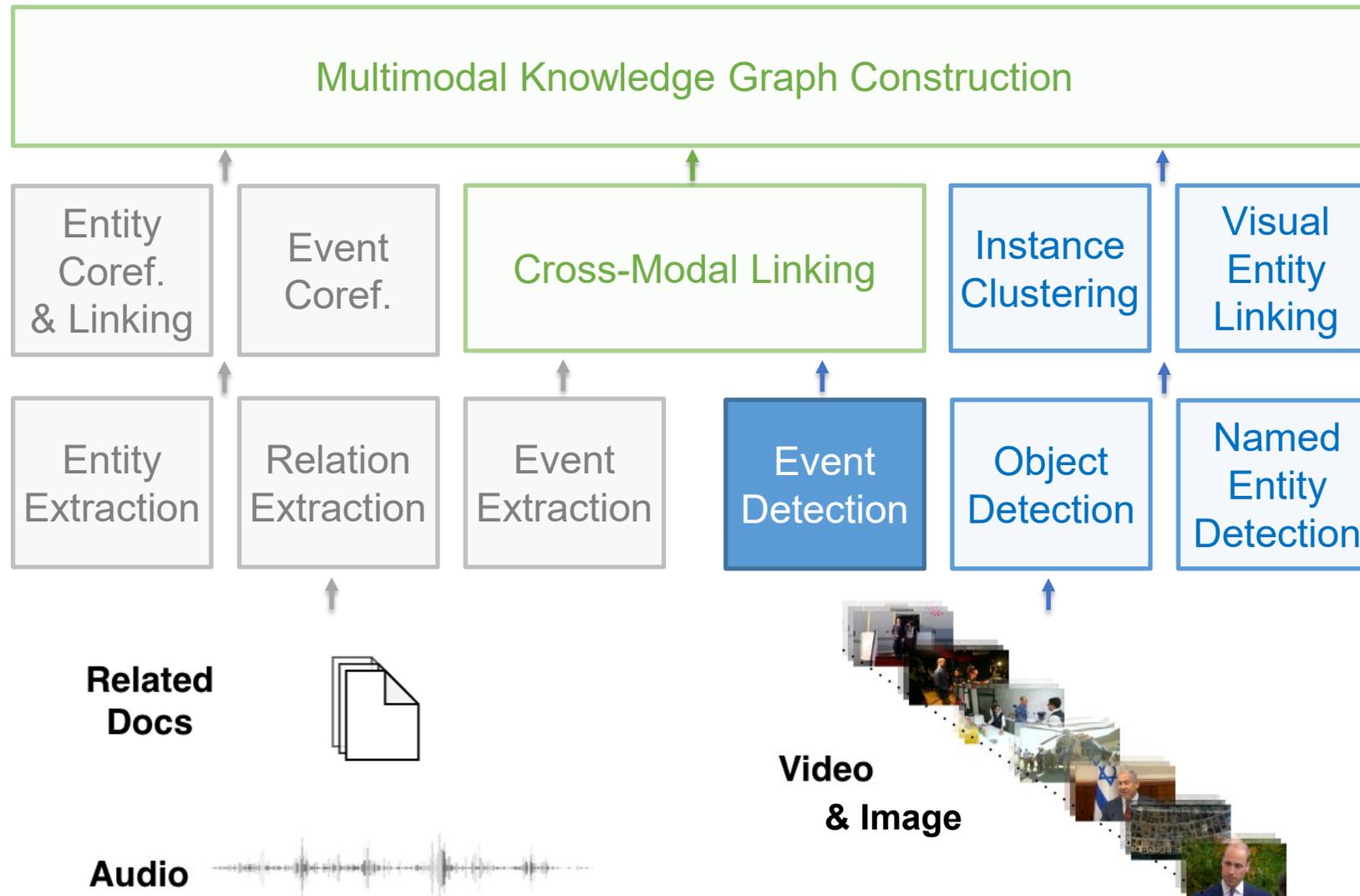
Kirstjen Nielsen



Kirstjen Nielsen



Kirstjen Nielsen



- Situation Recognition (image \rightarrow text + boundingbox)

RIDING							
ROLE	AGENT	VEHICLE	PLACE	ROLE	AGENT	VEHICLE	PLACE
VALUE	MAN	HORSE	OUTSIDE	VALUE	DOG	SKATEBOARD	SIDEWALK

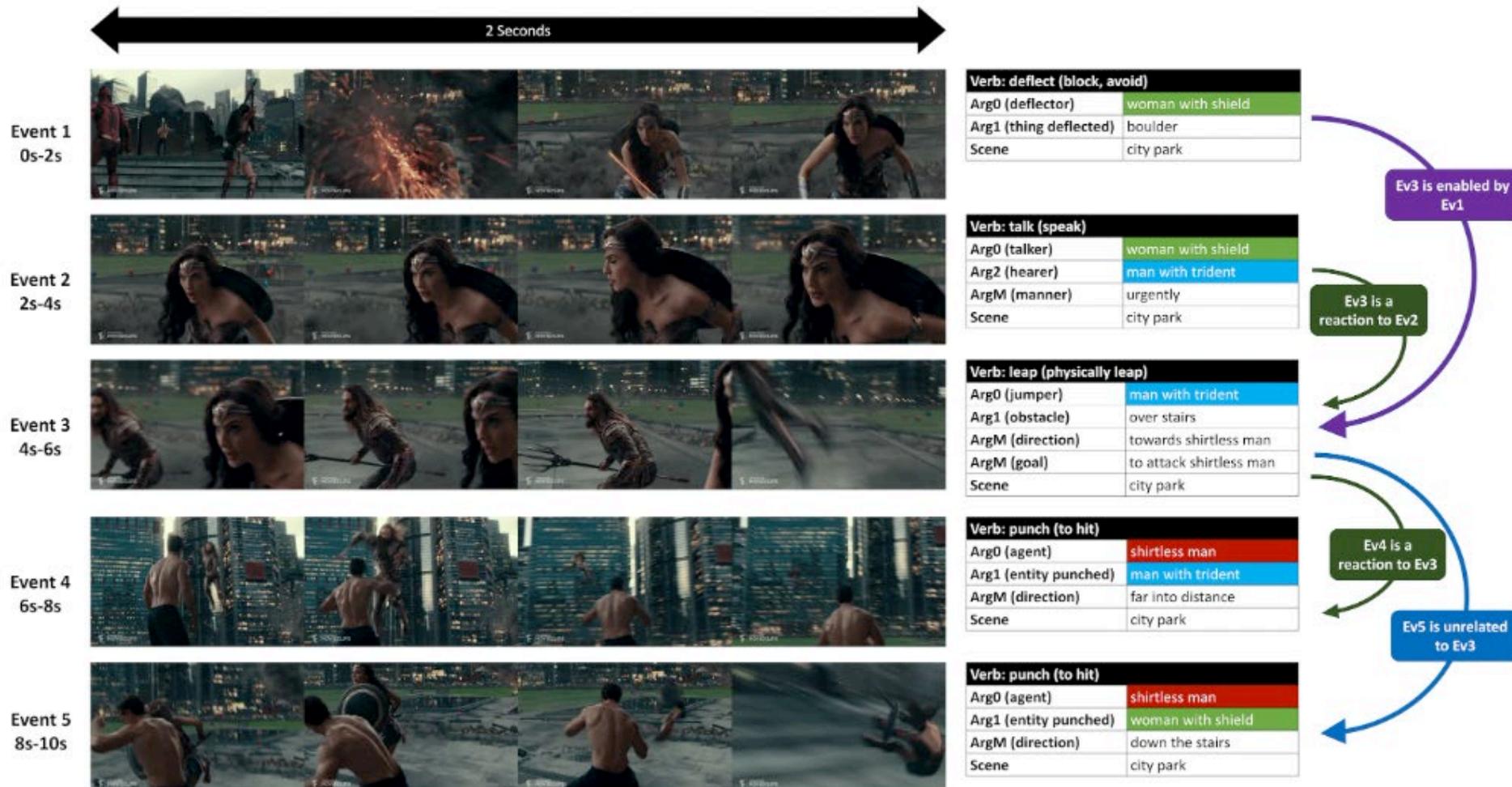
Situation Recognition [Yatskar et al. 2016]

Hitting				
Agent	Tool	Victim	Victim Part	Place
Ballplayer	Bat	Baseball	Ø	Field

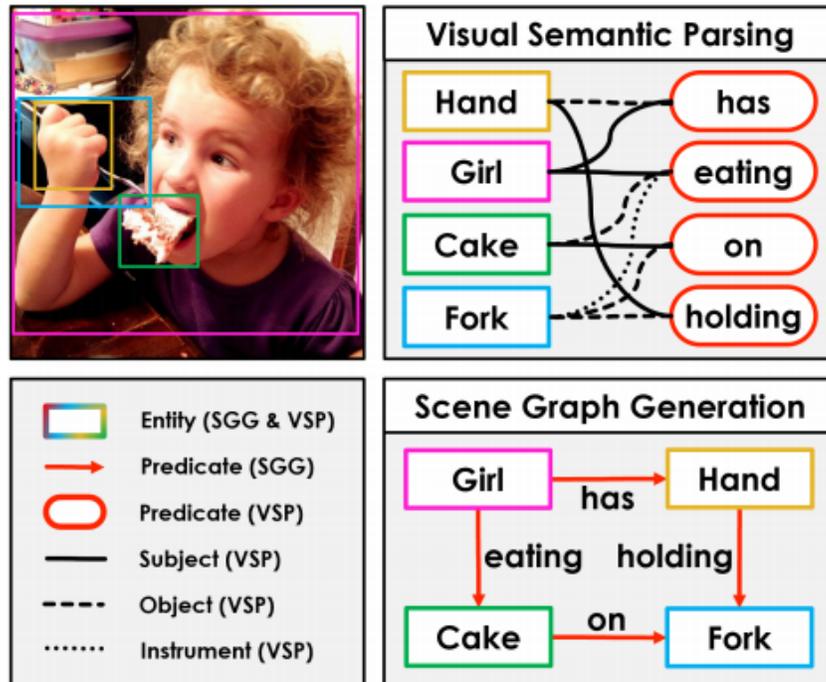
Catching			
Agent	Caught Item	Tool	Place
Bear	Fish	Mouth	River

Grounded Situation Recognition [Pratt et al, 2020]

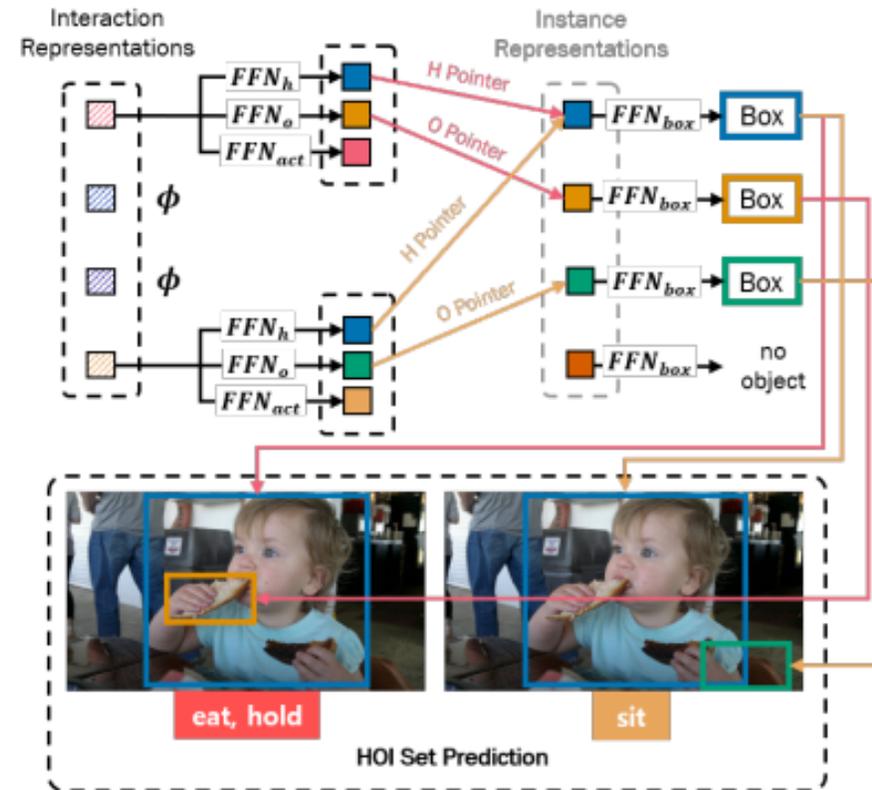
- Visual Features can provide more details for event interactions



- Scene Graph Based (image \rightarrow text + boundingbox)
 - Trained and evaluated on Scene Graph annotation (Visual Genome)
 - Added bounding box alignment

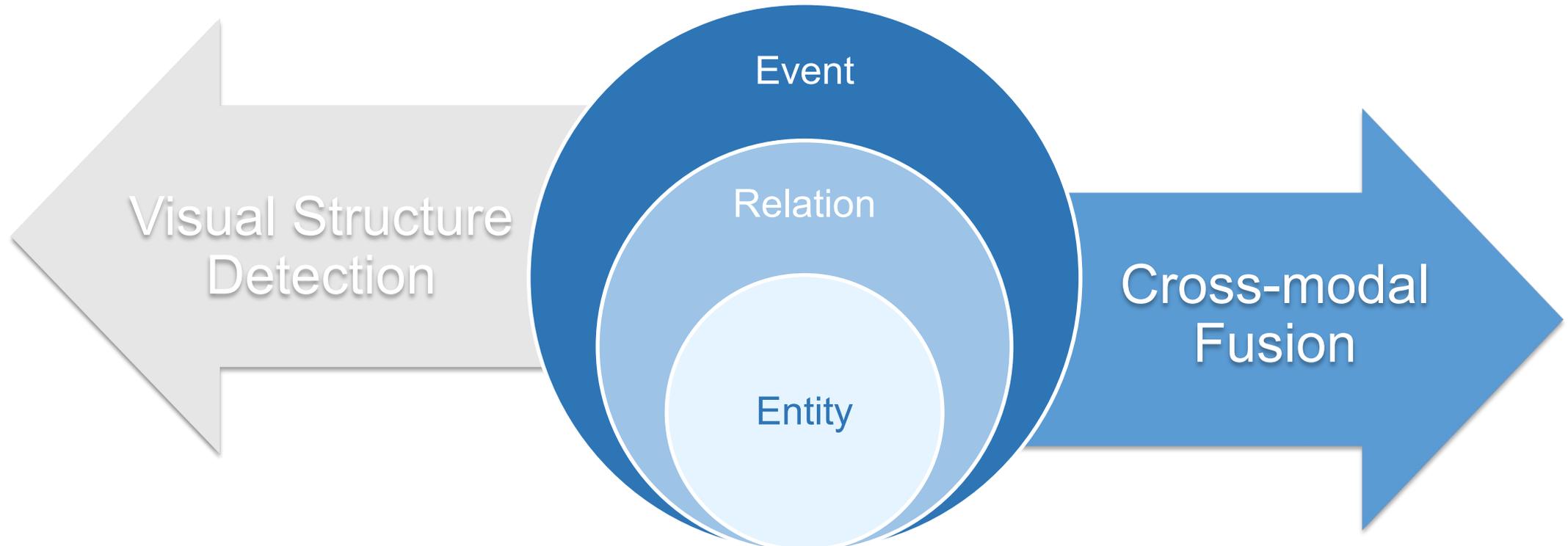


Visual Semantic Parsing [Zareian et al, 2020]

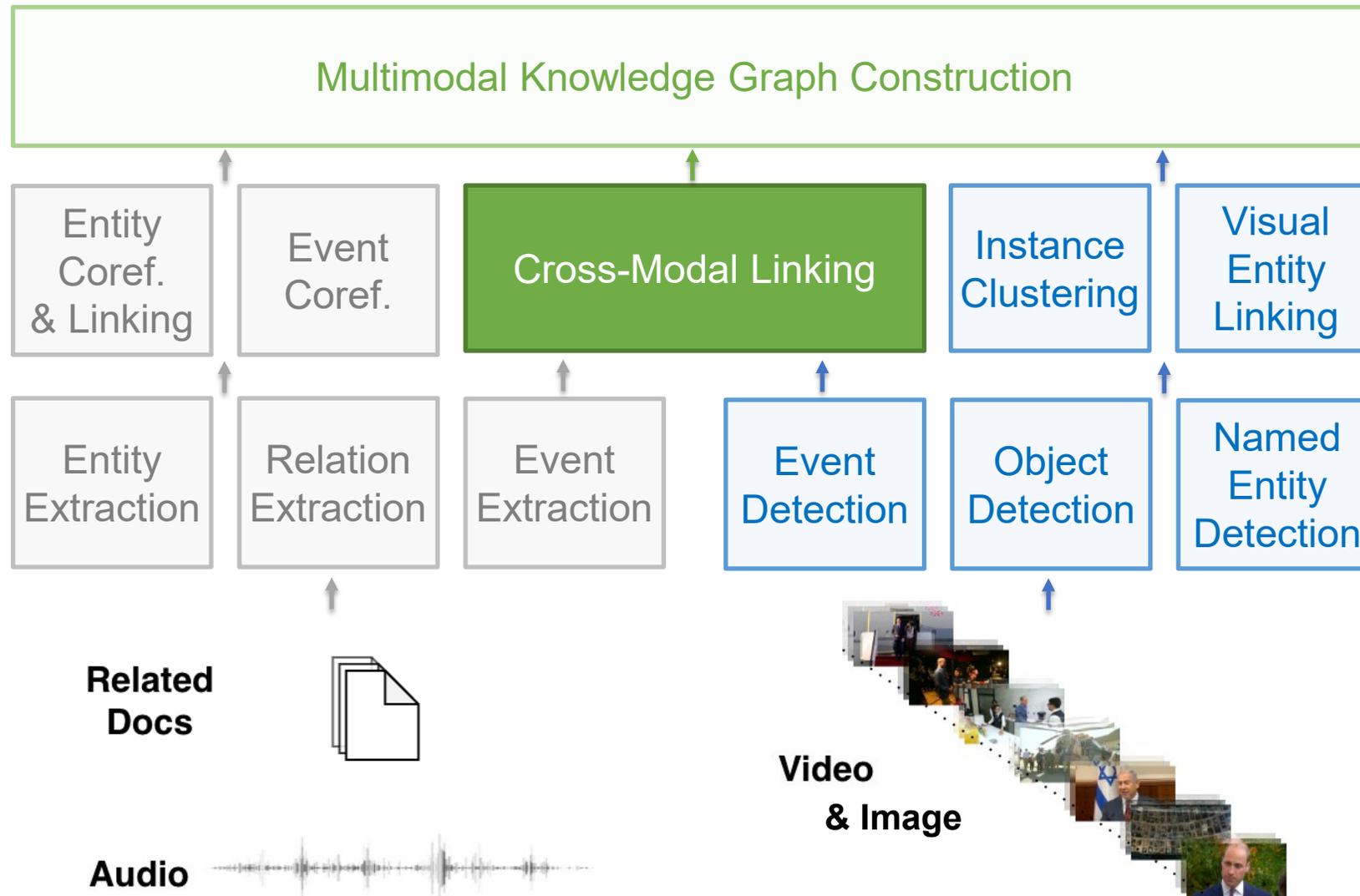


Human-Object Interaction [Kim et al, 2021]

- 1. How do we find structured knowledge in vision data?
- 2. How do we align structured across vision and text?



Multimedia Information Extraction



Cross-Media Fusion: Outline

Linking based

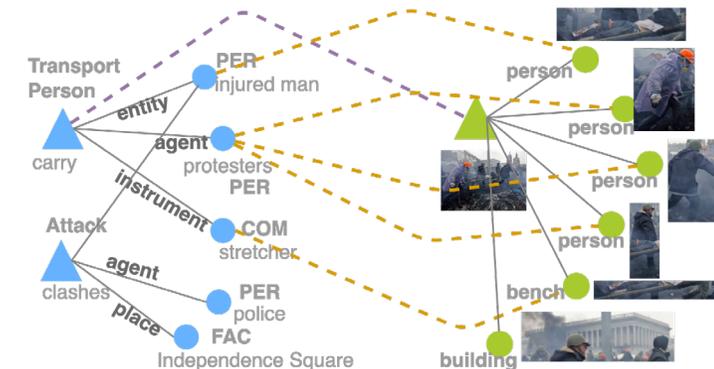


Grounding based

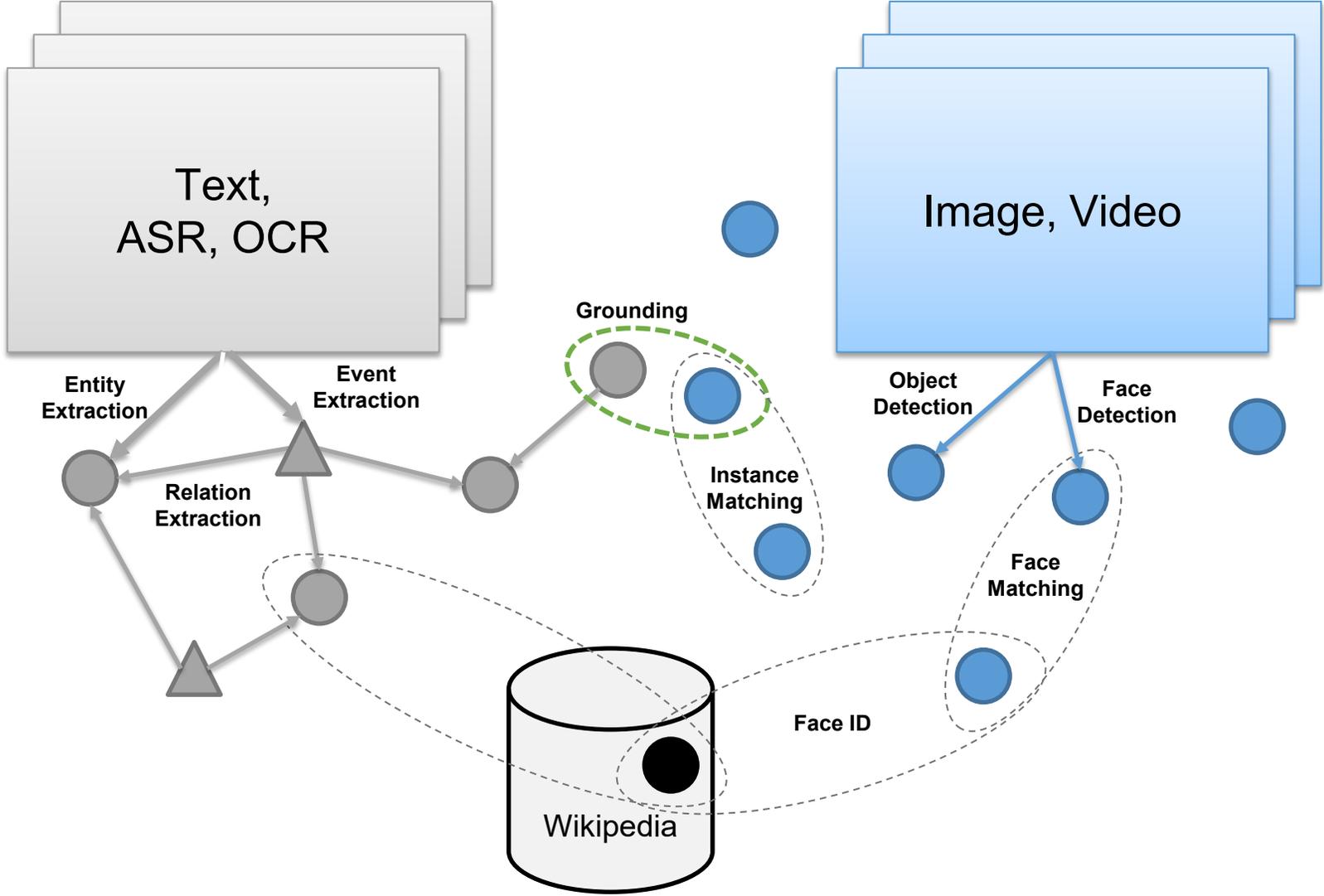
A crowd of **onlookers** on a **tractor ride** watch a farmer hard at work in the field



Structure based



Cross-Media Fusion: Linking-based



Cross-Media Fusion: Grounding-based

A man in red pushes his motocross bike up a rock



A crowd of onlookers on a tractor ride watch a farmer hard at work in the field

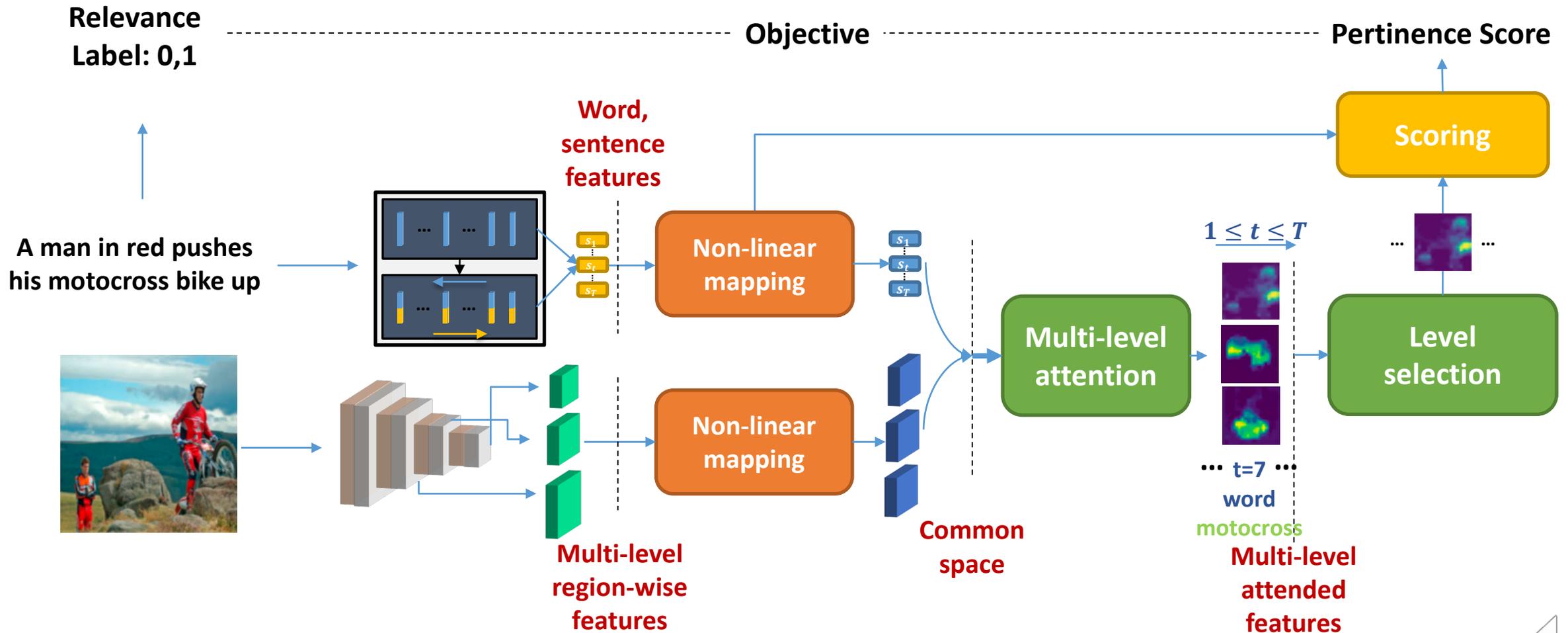


A cute young boy waving an american flag outside

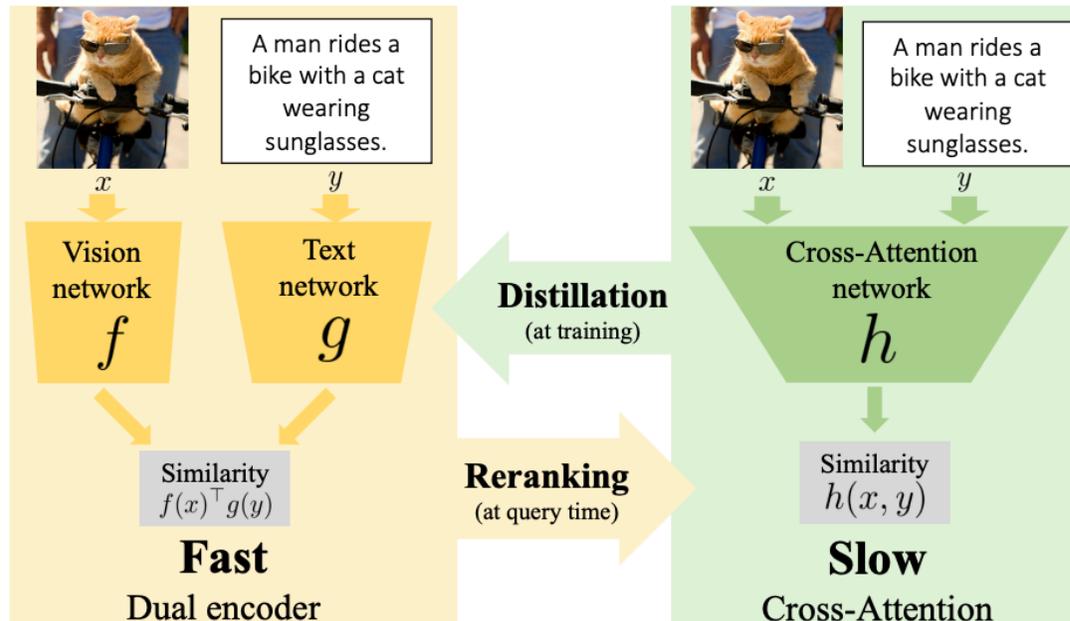


Cross-Media Fusion: Grounding-based

Cross-Attention between Word and Regions:

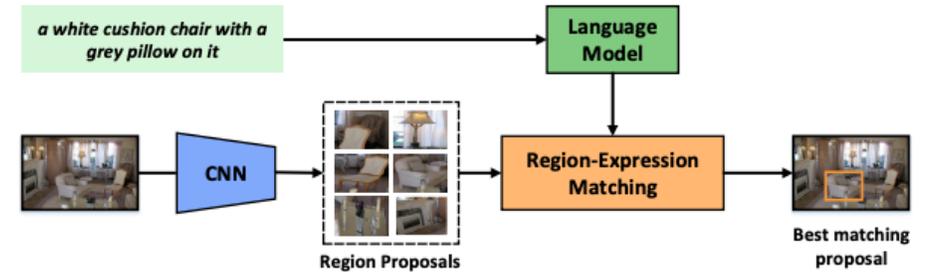


■ Cross-attention in Transformers

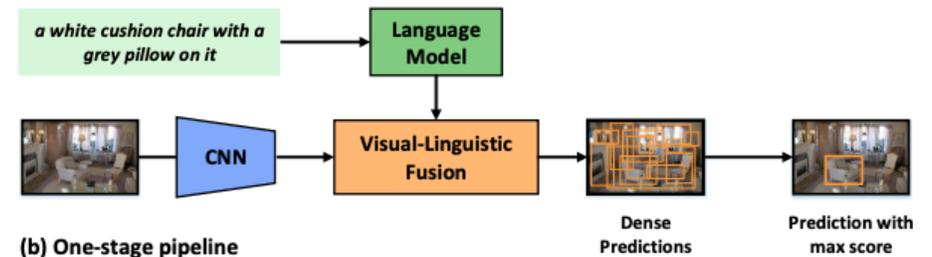


Fast and Slow

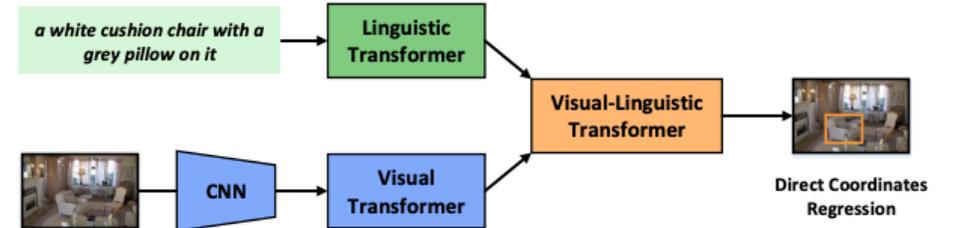
Miech, Antoine, et al. "Thinking fast and slow: Efficient text-to-visual retrieval with transformers." *CVPR* 2021.



(a) Two-stage pipeline



(b) One-stage pipeline



(c) TransVG (ours)

Deng, Jiajun, et al. "Transvg: End-to-end visual grounding with transformers." *CVPR* 2021.



Cross-Media Fusion: Outline

Linking based

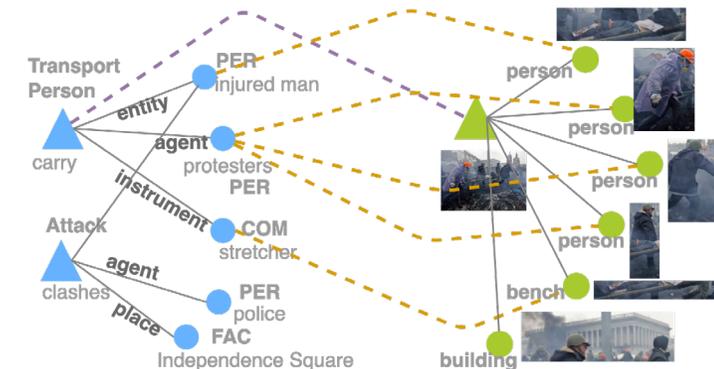


Grounding based

A crowd of **onlookers** on a **tractor ride** watch a farmer hard at work in the field

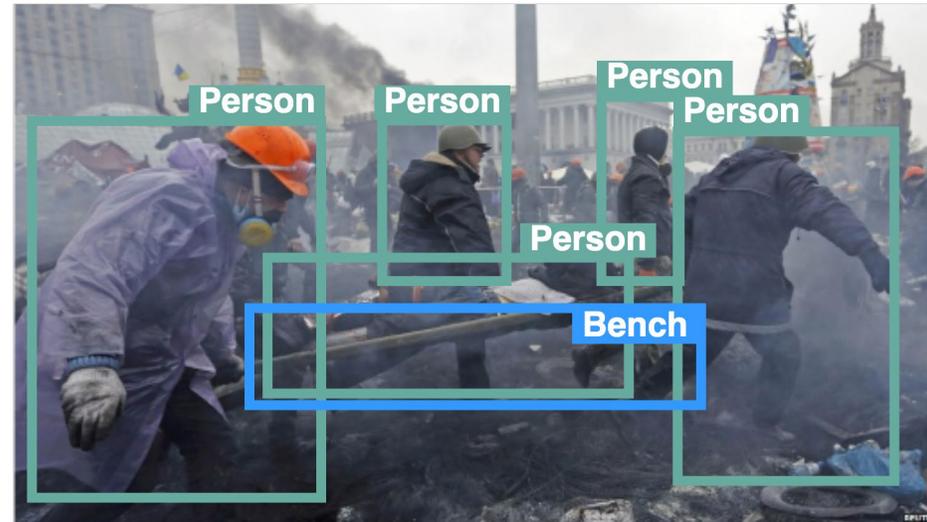


Structure based



Cross-Media Fusion: Structure-based

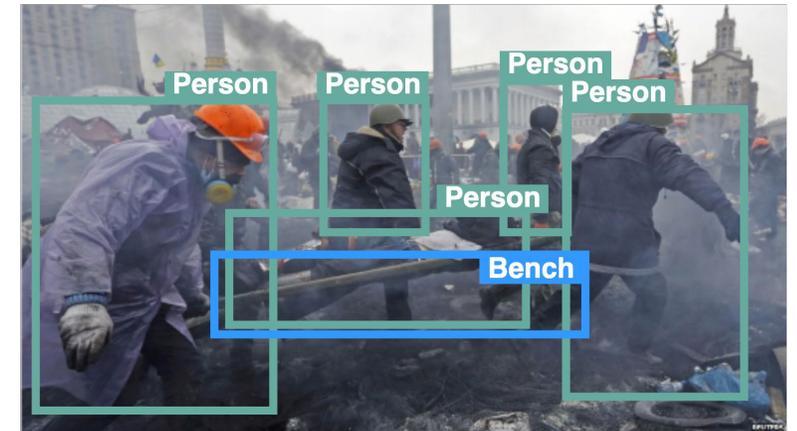
Antigovernment protesters **carry** an injured man on a **stretcher** after **clashes** with riot police on Independence Square in Kyiv on February 20, 2014.



Cross-Media Fusion: Structure-based

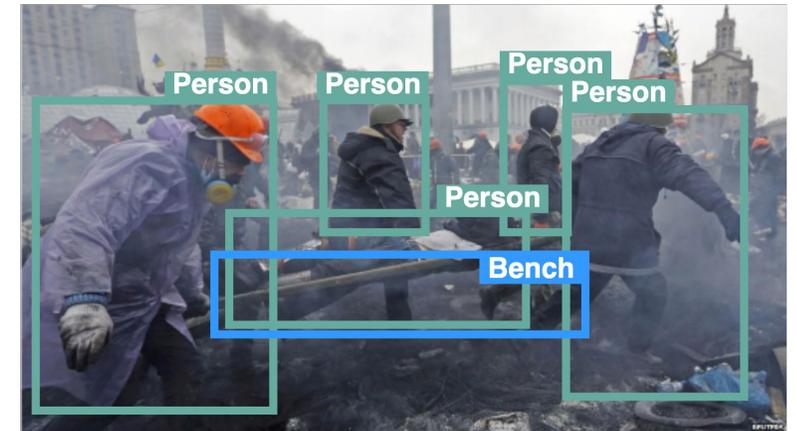
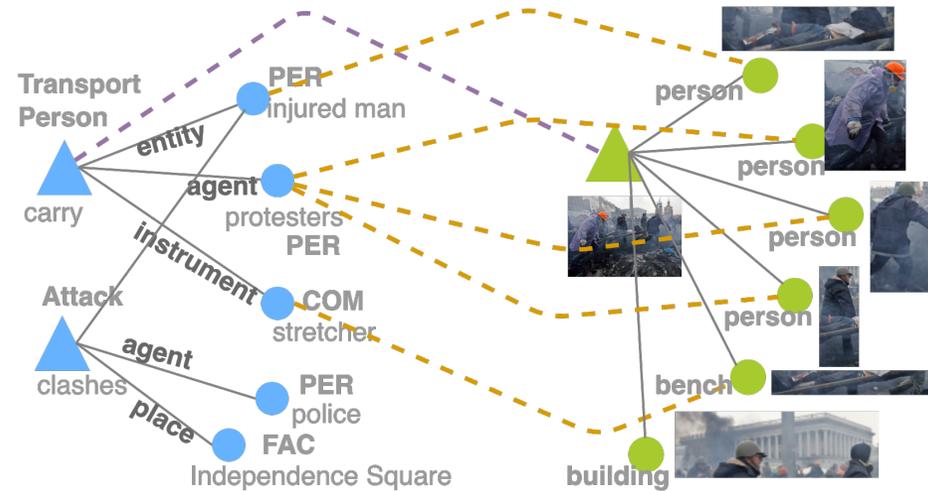
Antigovernment protesters **carry** an **injured man** on a **stretcher** after **clashes** with riot police on Independence Square in Kyiv on February 20, 2014.

	Event	Transport
Agent	protesters	
Entity	injured man	
Instrument	stretcher	



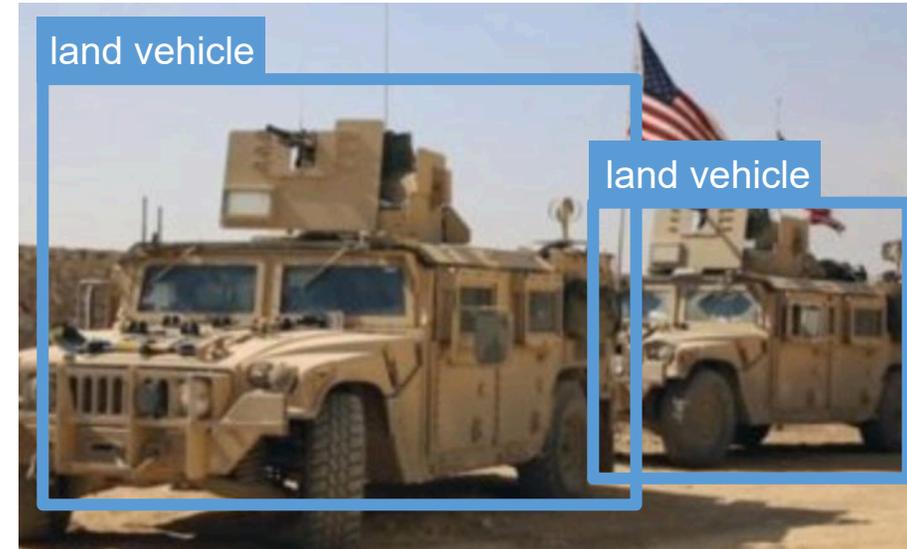
Cross-Media Fusion: Structure-based

Antigovernment protesters **carry** an **injured man** on a **stretcher** after **clashes** with riot police on Independence Square in Kyiv on February 20, 2014.



Input: News Article Text and Image

Last week , U.S . Secretary of State Rex Tillerson visited Ankara, the first senior administration official to visit Turkey, to try to seal a deal about the battle for Raqqa and to overcome President Recep Tayyip Erdogan's strong objections to Washington's backing of the Kurdish Democratic Union Party (PYD) militias. Turkish forces have attacked SDF forces in the past around Manbij, west of Raqqa, forcing the **United States** to **deploy** dozens of **soldiers** on the **outskirts** of the town in a mission to prevent a repeat of clashes, which risk derailing an assault on Raqqa.



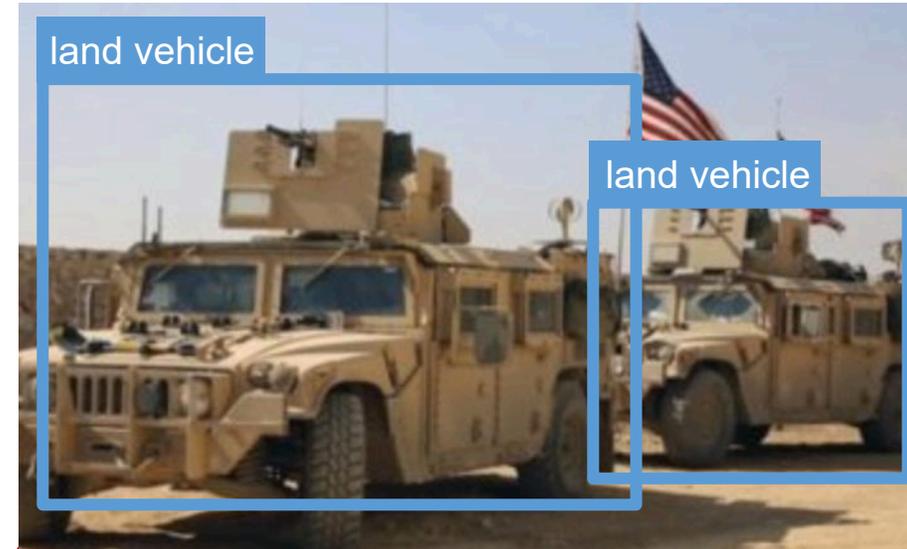
Output: Events & Argument Roles

Event Type	Movement.Transport
Text Trigger	deploy
Event	
Image	

Arguments	Agent	United States
	Destination	outskirts
	Artifact	soldiers
	Vehicle	
	Vehicle	

Input: News Article Text and Image

Last week , U.S . Secretary of State Rex Tillerson visited Ankara, the first senior administration official to visit Turkey, to try to seal a deal about the battle for Raqqa and to overcome President Recep Tayyip Erdogan's strong objections to Washington's backing of the Kurdish Democratic Union Party (PYD) militias. Turkish forces have attacked SDF forces in the past around Manbij, west of Raqqa, forcing the **United States** to **deploy** dozens of **soldiers** on the **outskirts** of the town in a mission to prevent a repeat of clashes, which risk derailing an assault on Raqqa.



Output: Multimedia Events & Argument Roles

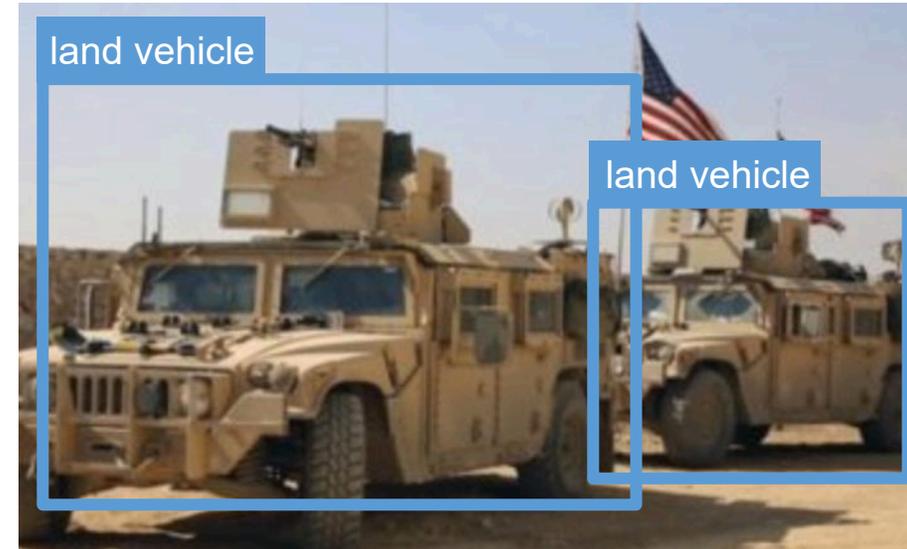
Event Type	Movement.Transport
Text Trigger	deploy
Event	
Image	

Arguments	Agent	United States
	Destination	outskirts
	Artifact	soldiers
	Vehicle	
	Vehicle	

Multimedia Event Extraction (M²E²)

Input: News Article Text and Image

Last week , U.S . Secretary of State Rex Tillerson visited Ankara, the first senior administration official to visit Turkey, to try to seal a deal about the battle for Raqqa and to overcome President Recep Tayyip Erdogan's strong objections to Washington's backing of the Kurdish Democratic Union Party (PYD) militias. Turkish forces have attacked SDF forces in the past around Manbij, west of Raqqa, forcing the **United States** to **deploy** dozens of **soldiers** on the **outskirts** of the town in a mission to prevent a repeat of clashes, which risk derailing an assault on Raqqa.

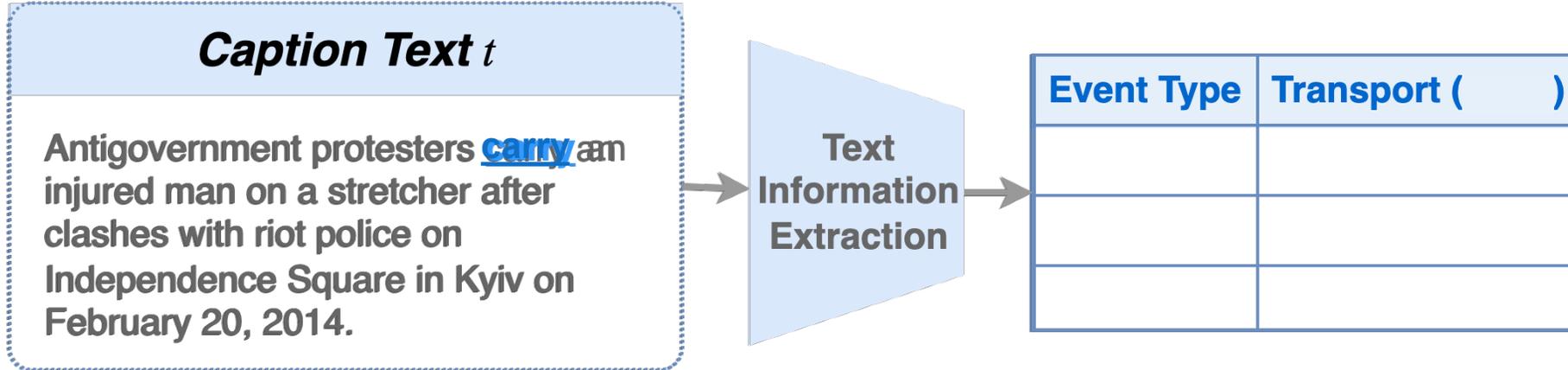


Output: Multimedia Events & Argument Roles

Event Type	Movement.Transport
Text Trigger	deploy
Event	
Image	

Arguments	Agent	United States
	Destination	outskirts
	Artifact	soldiers
	Vehicle	
	Vehicle	

Transfer Structured Knowledge Across Modalities

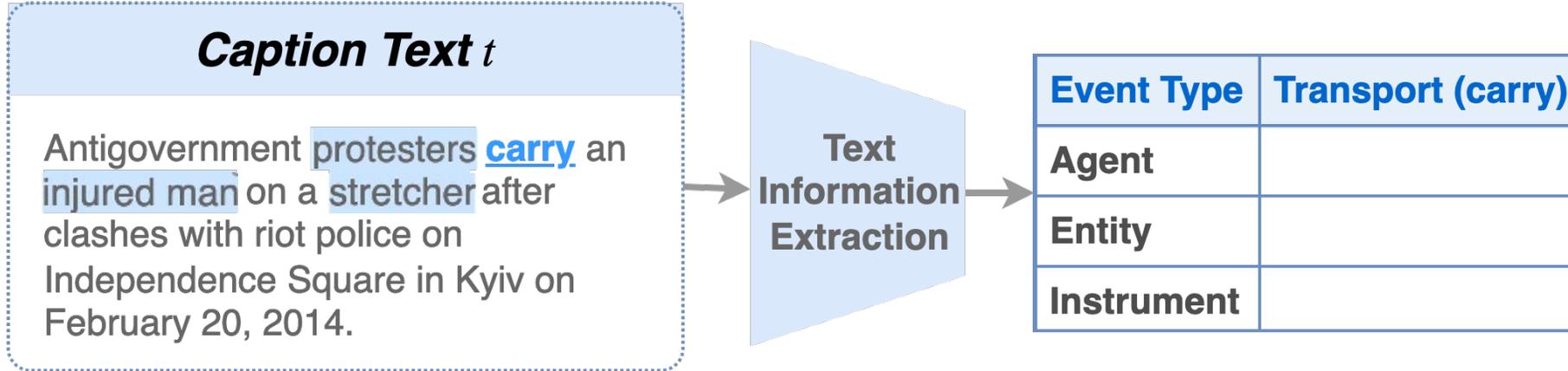


Natural Language Processing

Event Extraction { Trigger Word



Transfer Structured Knowledge Across Modalities



Natural Language Processing

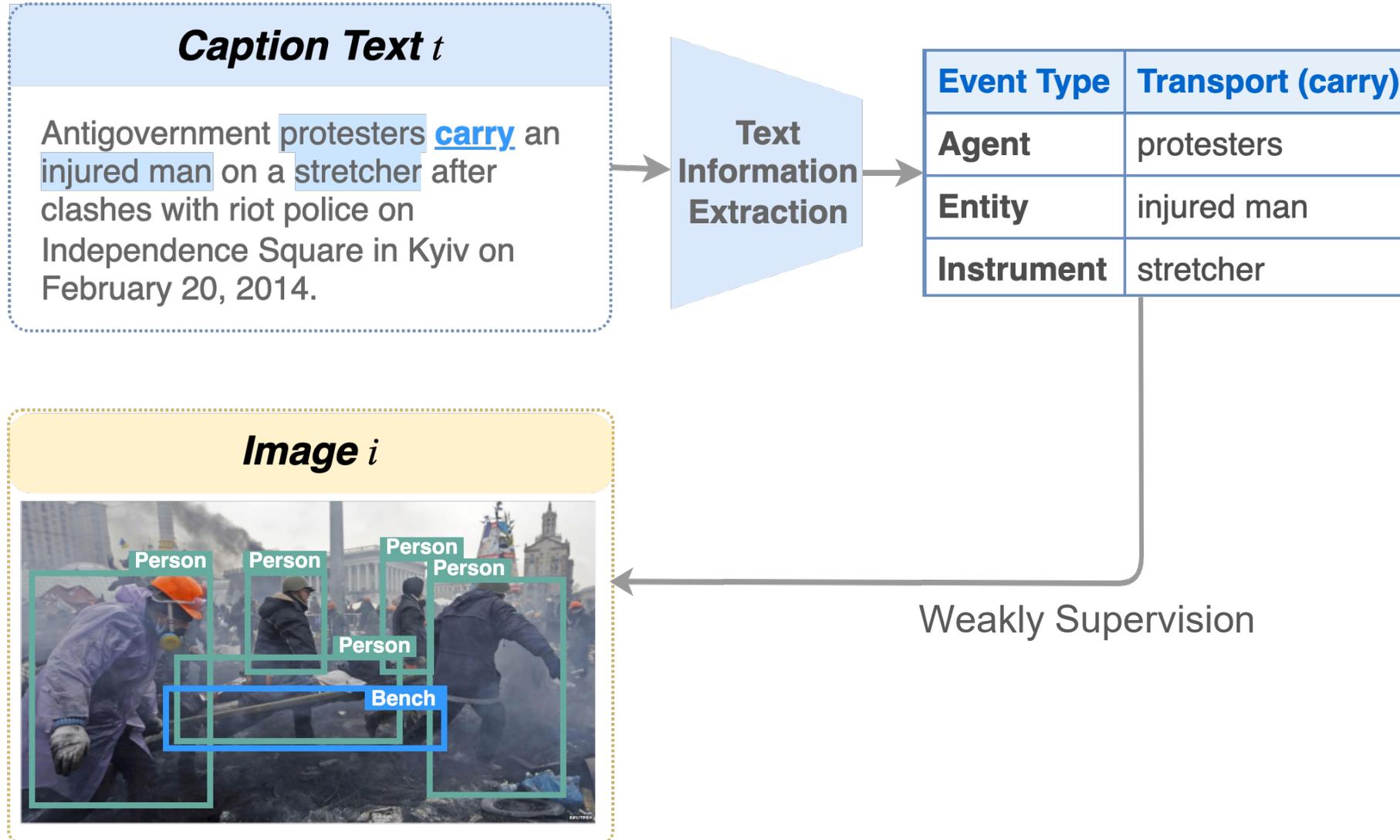
Event Extraction

Trigger Word

Argument (Participant)



Transfer Structured Knowledge Across Modalities



Construct hard negatives by manipulating event structures.

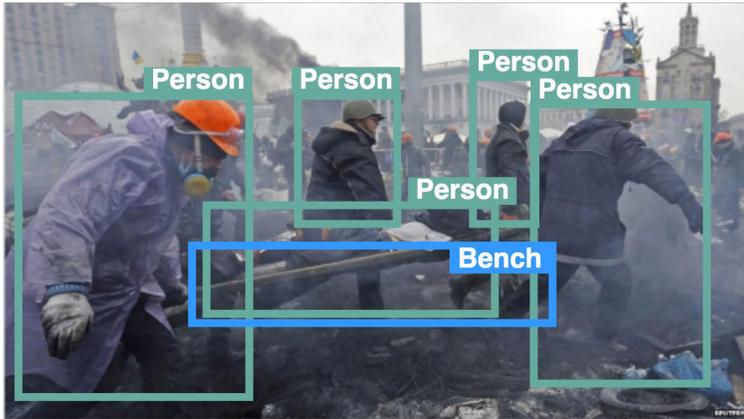
Caption Text t

Antigovernment protesters carry an injured man on a stretcher after clashes with riot police on Independence Square in Kyiv on February 20, 2014.

Positive Labels

Event Type	Transport (carry)
Agent	protesters
Entity	injured man
Instrument	stretcher

Image i



Negative Labels (events)

Event Type	Arrest (arrest)
Agent	protesters
Entity	injured man
Instrument	stretcher

Negative Labels (arguments)

Event Type	Transport (carry)
Agent	protesters
Entity	injured man
Instrument	stretcher



Construct hard negatives by manipulating event structures.

Caption Text t

Antigovernment protesters **carry** an injured man on a stretcher after clashes with riot police on Independence Square in Kyiv on February 20, 2014.

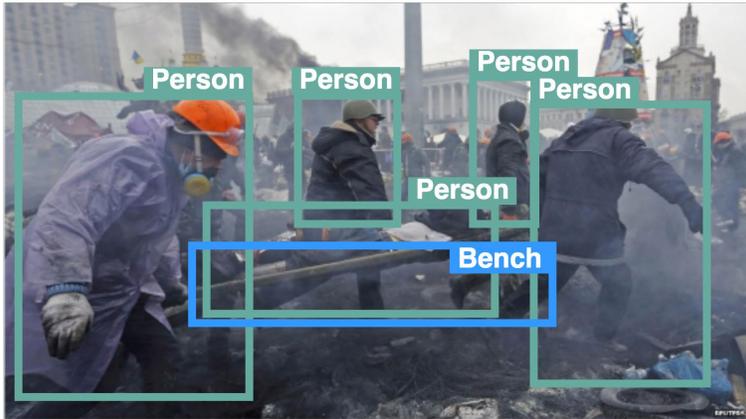
Text
Labels

Event Type	Transport (carry)
Agent	protesters
Entity	injured man
Instrument	stretcher

prompt

Protesters transported injured man using a stretcher.

Image i



Negative Labels
(events)

Event Type	Arrest (arrest)
Agent	protesters
Entity	injured man
Instrument	stretcher

prompt

Protesters **arrested** injured man using a stretcher.

Negative Labels
(arguments)

Event Type	Transport (carry)
Agent	injured man
Entity	stretcher
Instrument	protesters

prompt

Injured man transported a **stretcher** with **protesters**.



Construct hard negatives by manipulating event structures.

Caption Text t

Antigovernment protesters **carry** an injured man on a stretcher after clashes with riot police on Independence Square in Kyiv on February 20, 2014.

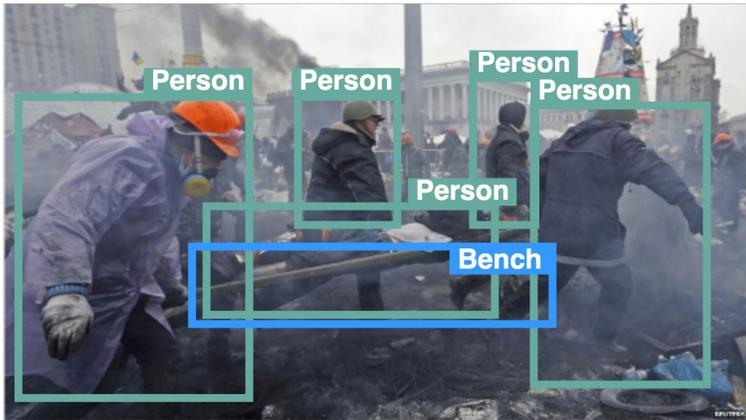
Text
Positive
Labels

Event Type	Transport (carry)
Agent	protesters
Entity	injured man
Instrument	stretcher

prompt

Protesters transported injured man using a stretcher.

Image i



Negative
Labels
(events)

Event Type	Arrest (arrest)
Agent	protesters
Entity	injured man
Instrument	stretcher

prompt

Protesters **arrested** injured man using a stretcher.

Negative
Labels
(arguments)

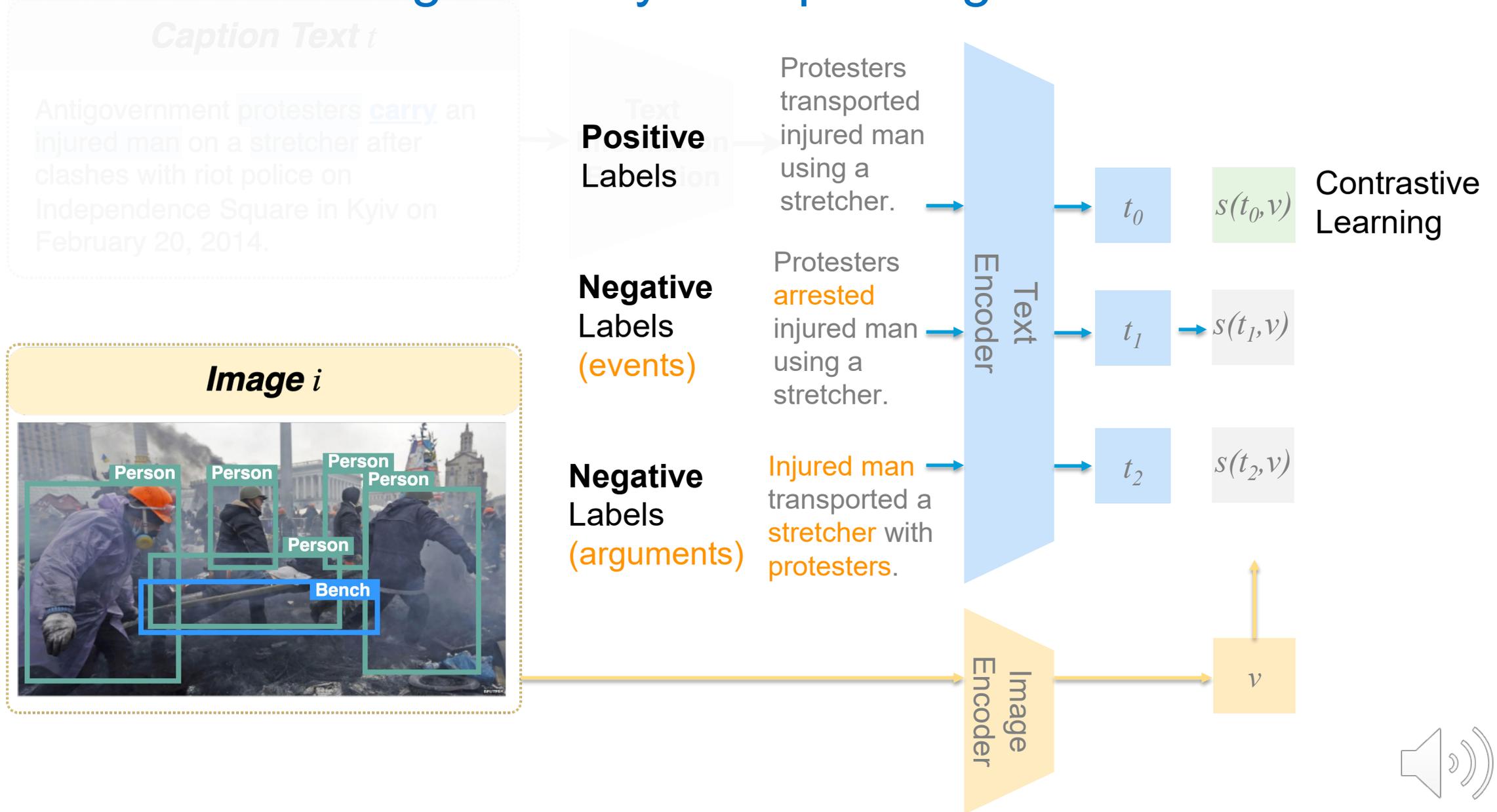
Event Type	Transport (carry)
Agent	injured man
Entity	stretcher
Instrument	protesters

prompt

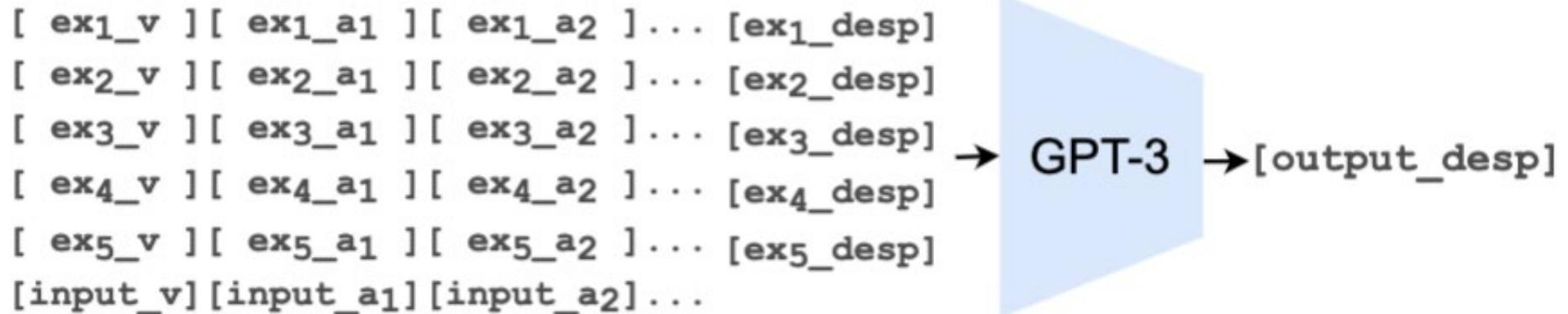
Injured man transported a **stretcher** with **protesters**.



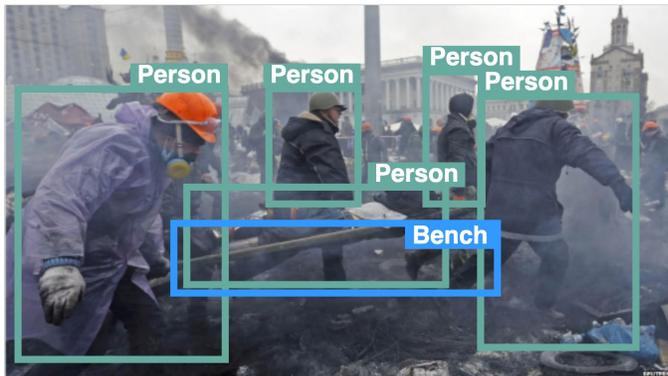
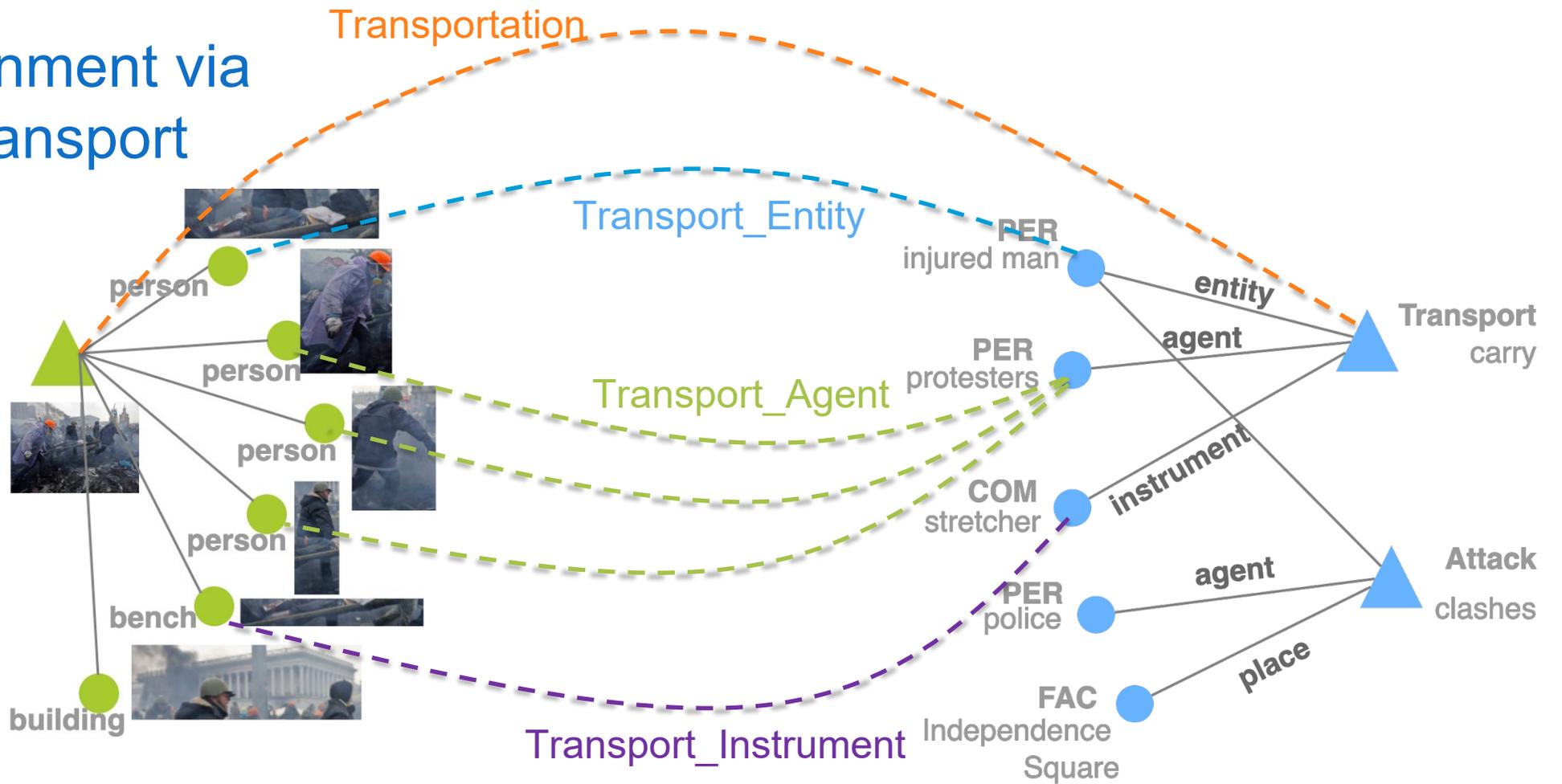
Construct hard negatives by manipulating event structures.



- From Event Structure to Natural Language Description:
 - We use five manual event description examples as few-shot prompts to control the generation.



Graph Alignment via Optimal Transport

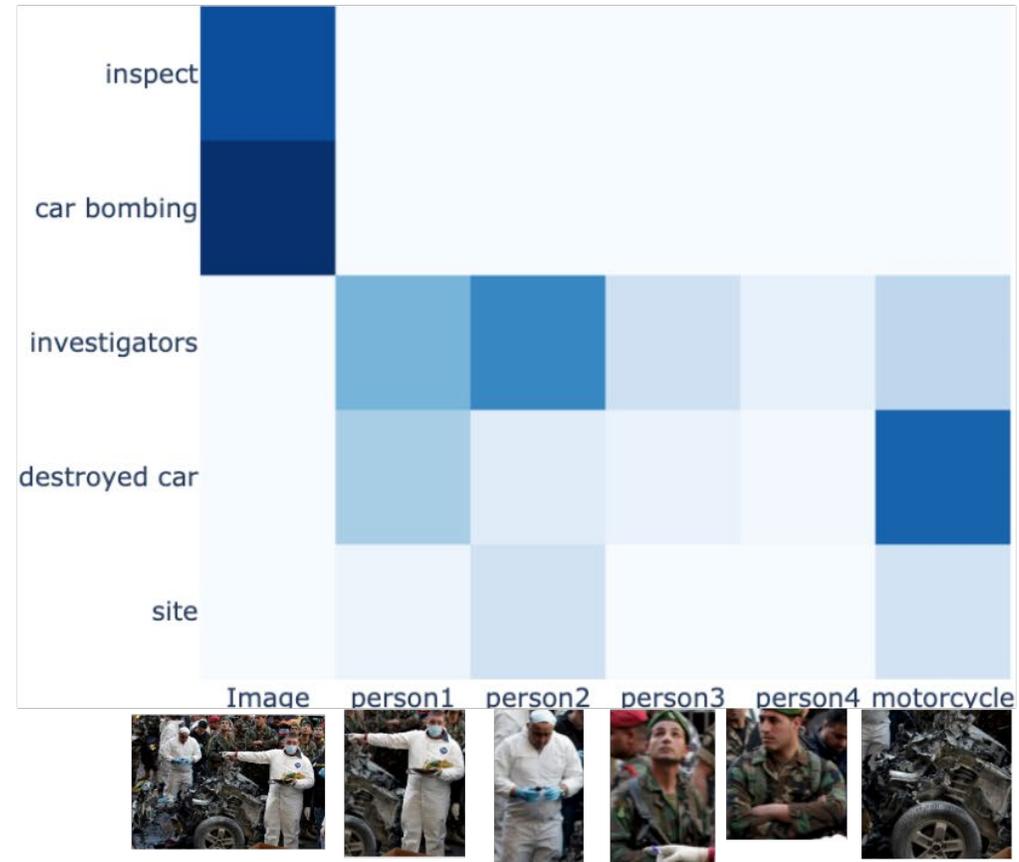


Antigovernment **protesters** **carry** an **injured man** on a **stretcher** after **clashes** with riot police on Independence Square in Kyiv on February 20, 2014.

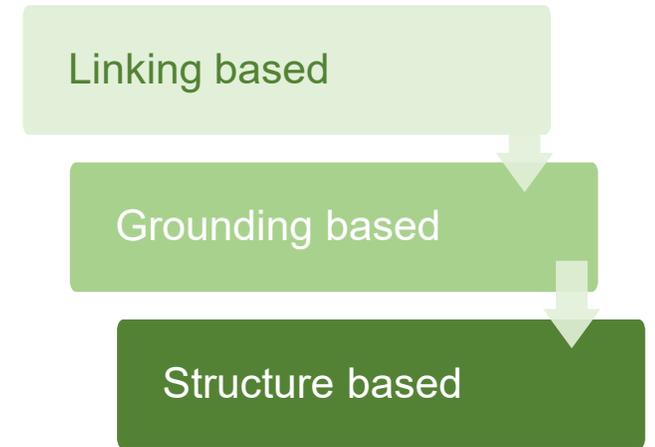
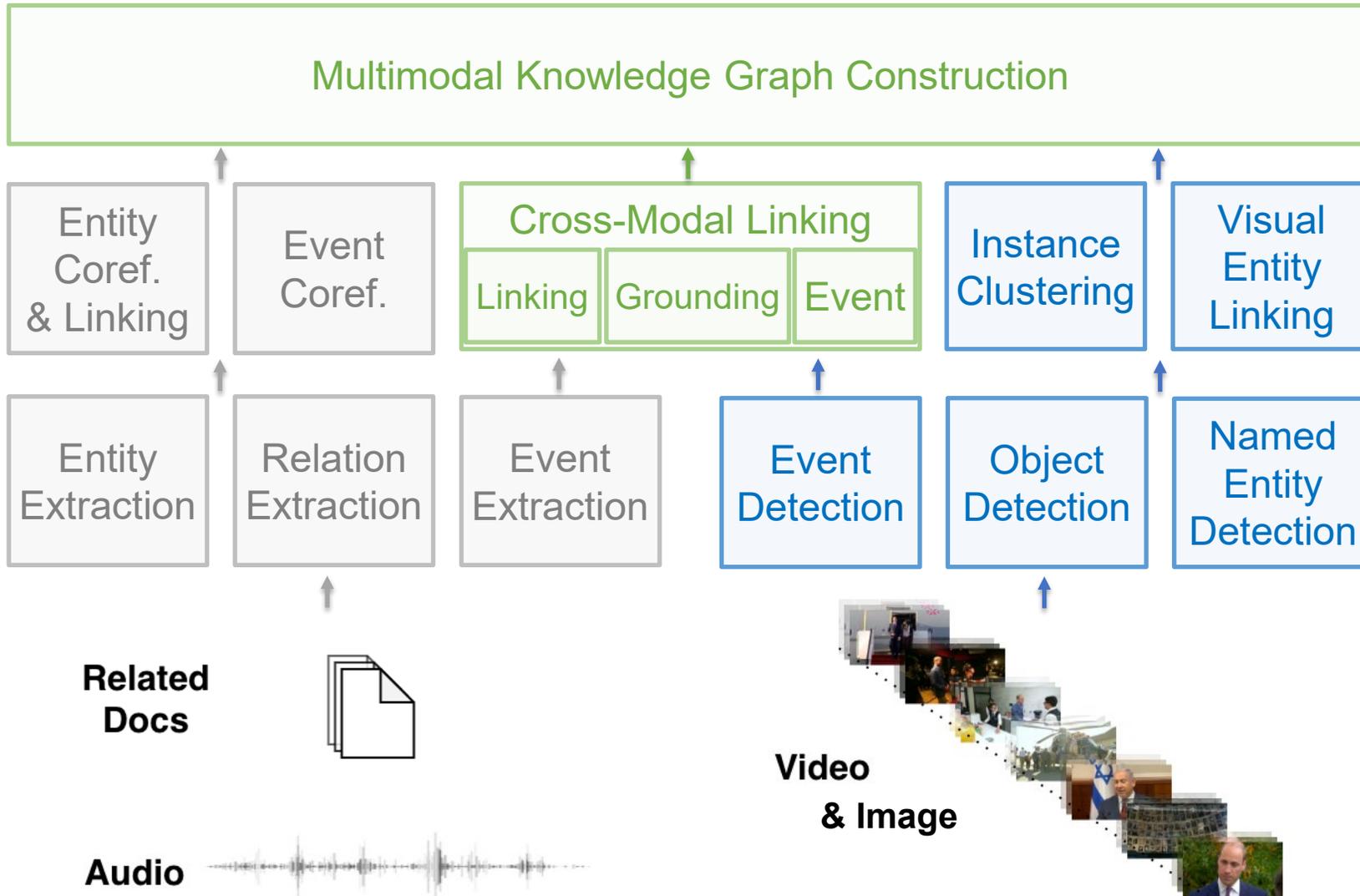


Event-level Cross-media Alignment

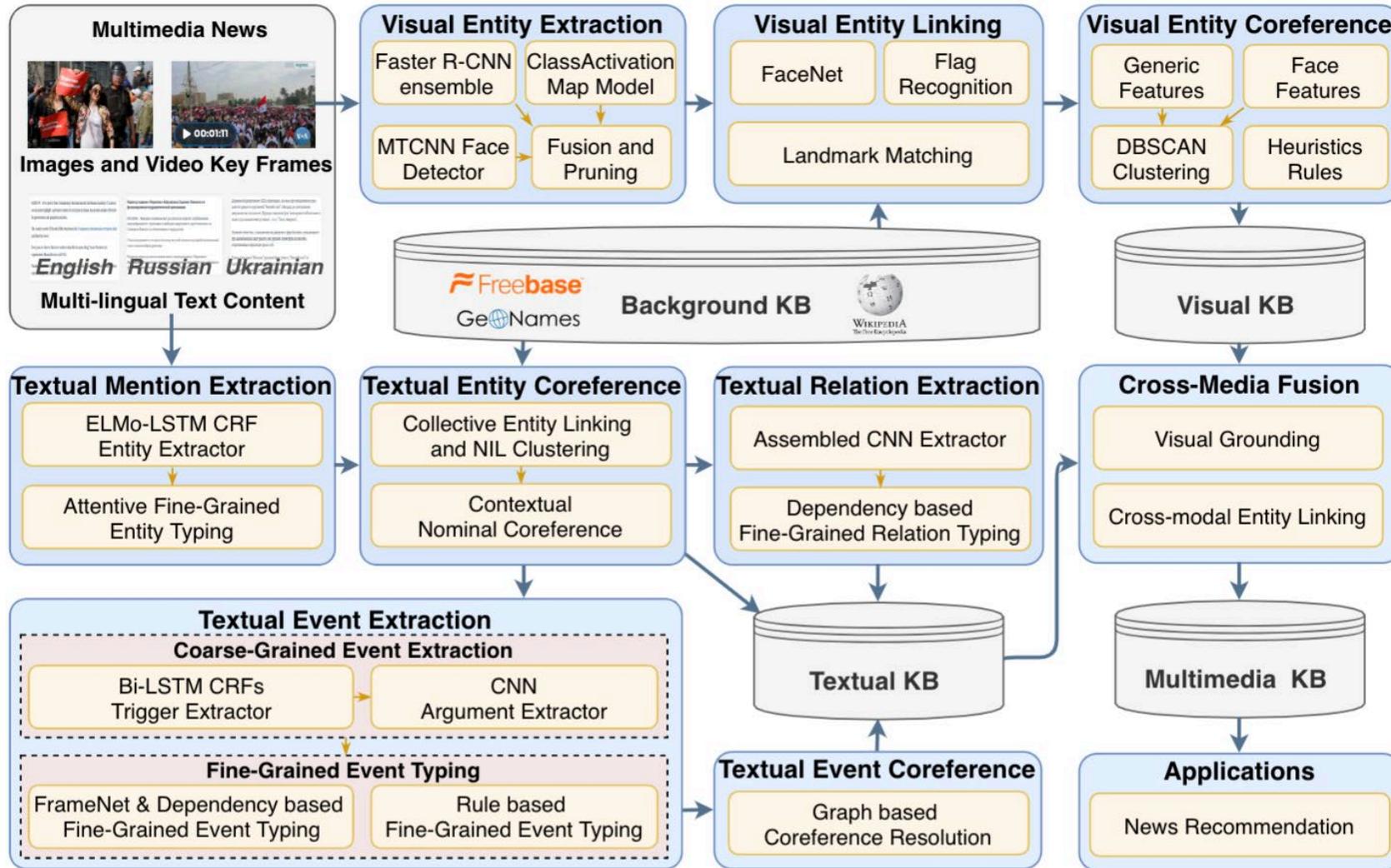
Investigators **inspect** parts of a **destroyed car** at the **site** of a **car bombing** in Beirut, Jan. 21, 2014.



Conclusions

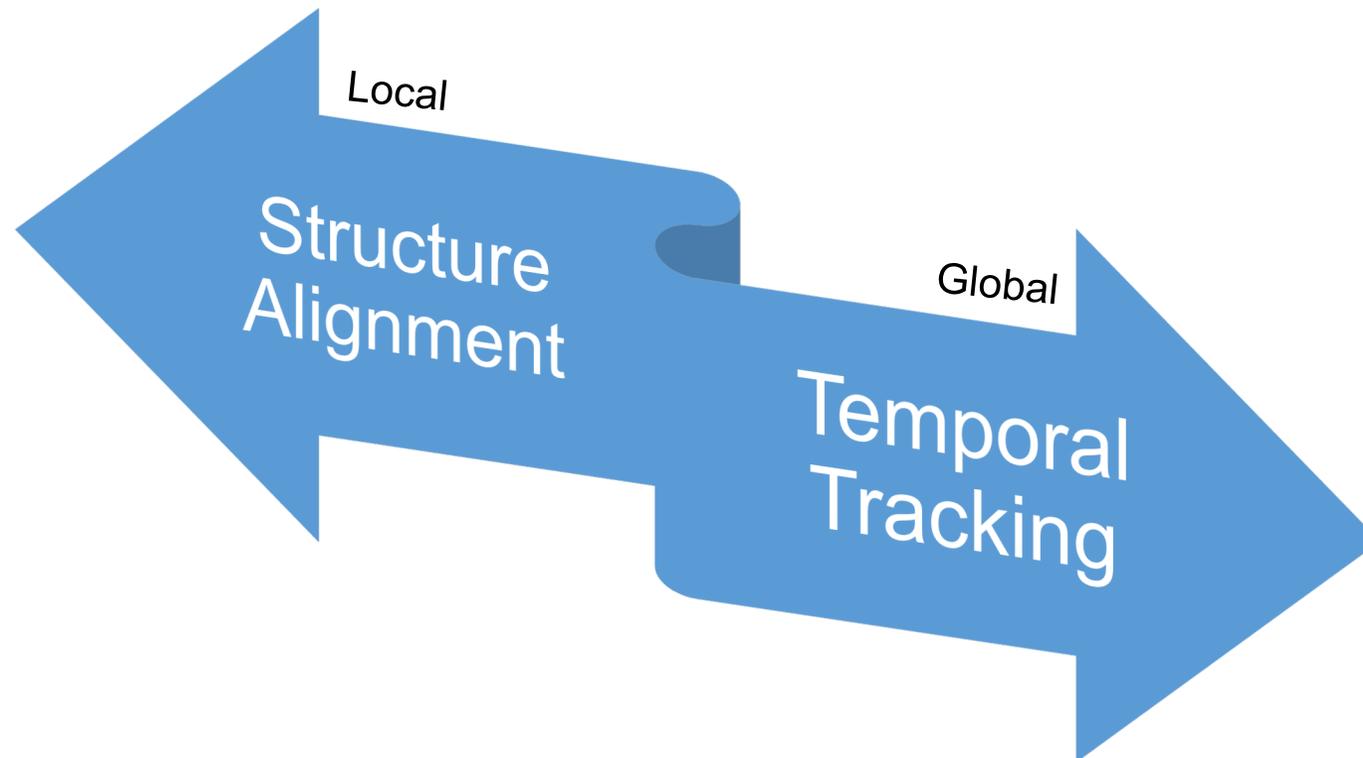


GAIA: Open-source Multimedia Knowledge Extraction System

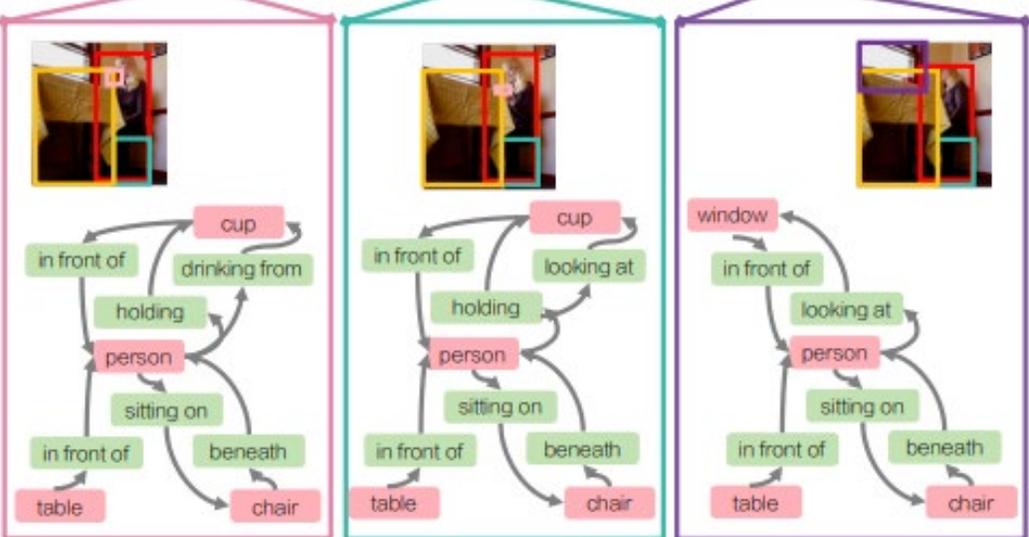
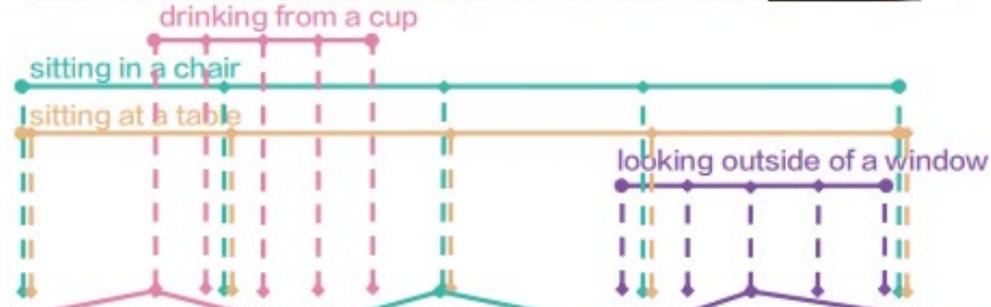


- Available at <http://blender.cs.illinois.edu/software/gaia-ie>
- Achieved best performance at TAC SM-KBP 2019 and 2020 Evaluation (10% higher than the second ranked team in TAC SM-KBP 2019)

- Local: Capturing semantic structure
- Global: Understanding a more global context of multiple entities and events

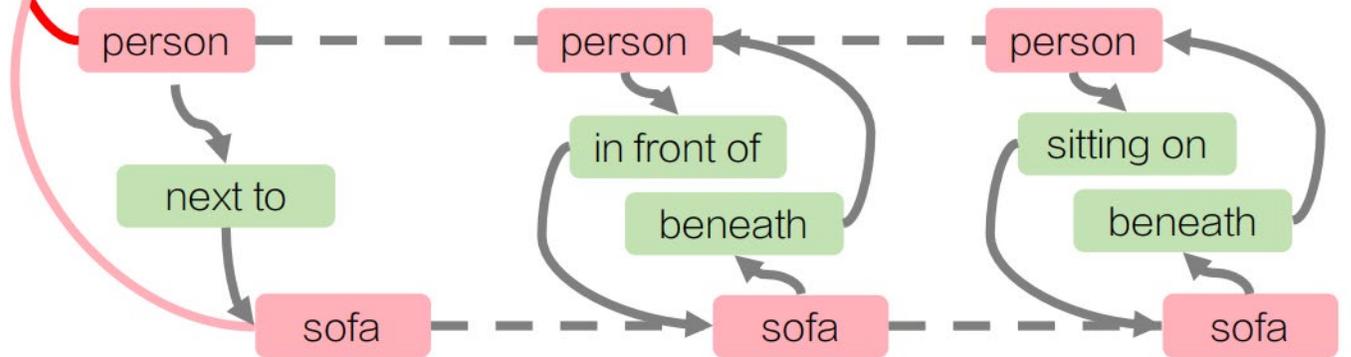


Future Direction: Event Tracking



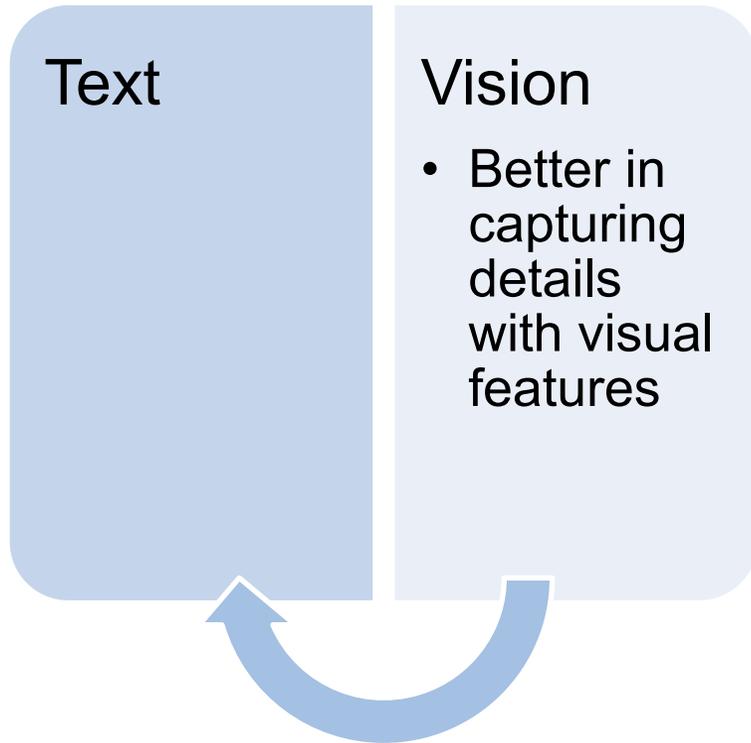
Action: "Sitting on a sofa"

time



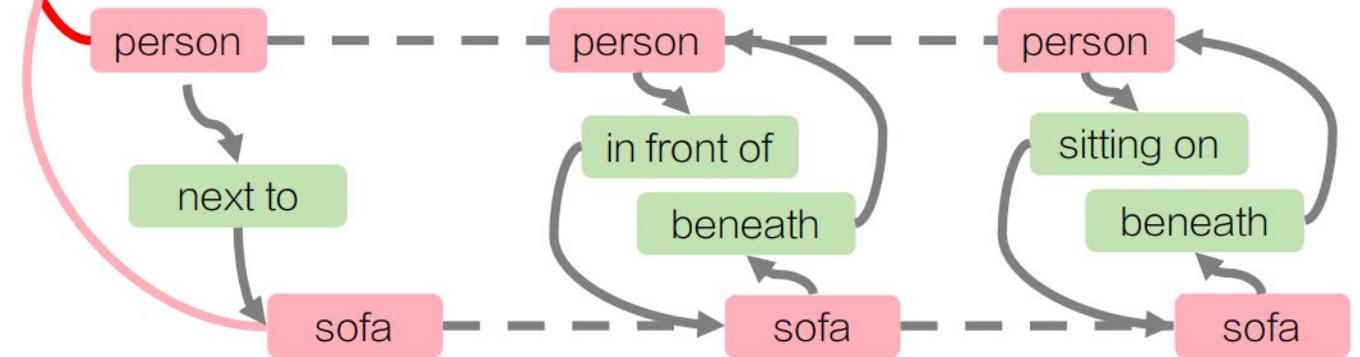
Spatio-temporal scene graphs

Future Direction: Text and Vision are Complementary



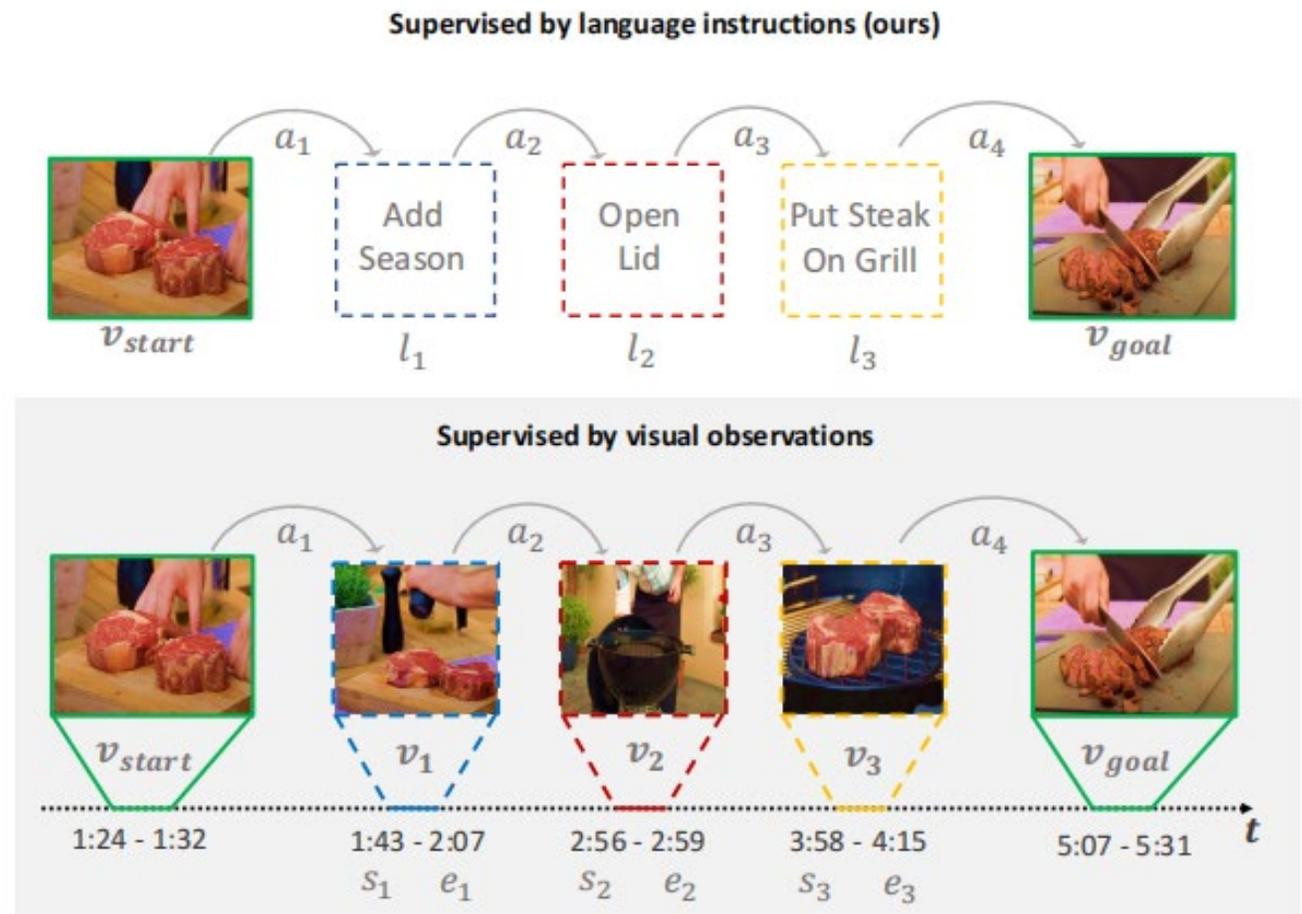
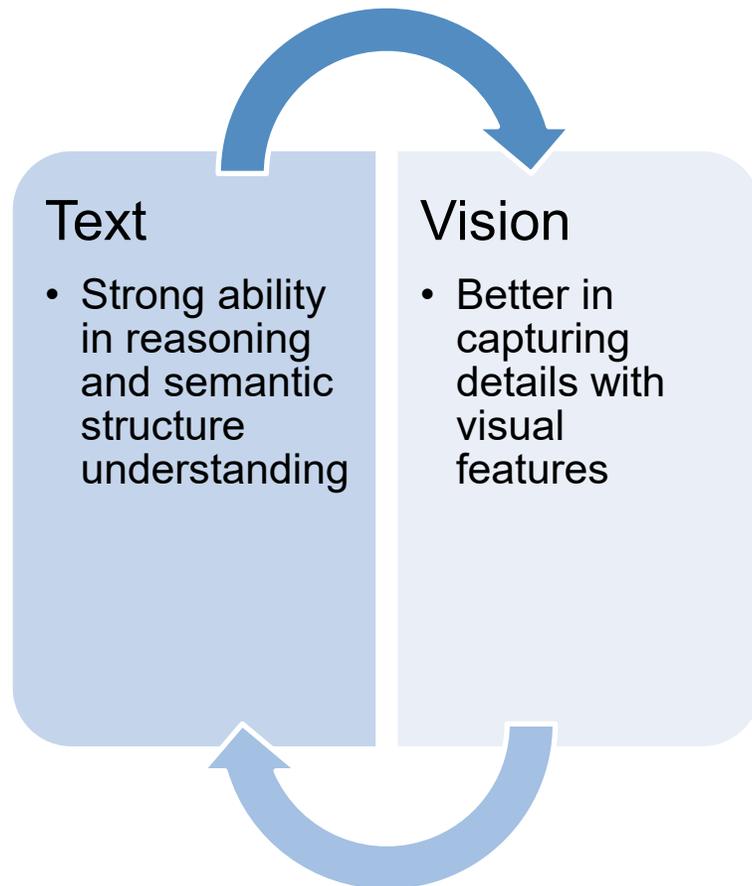
Action: "Sitting on a sofa"

time



Spatio-temporal scene graphs

Future Direction: Text and Vision are Complementary



P3IV: Probabilistic Procedure Planning from Instructional Videos with Weak Supervision

Thank You

