



Knowledge-Driven Vision-Language Encoding



Manling Li
UIUC
(Incoming AP at Northwestern)



Xudong Lin
Columbia



Jie Lei
Meta AI



Mohit Bansal
UNC



Carl Vondrick
Columbia



Shih-Fu Chang
Columbia



Heng Ji
UIUC

How humans learn about knowledge?



Human

Interaction



External World

Knowledge



How humans learn about knowledge?



Human

Interaction



Knowledge

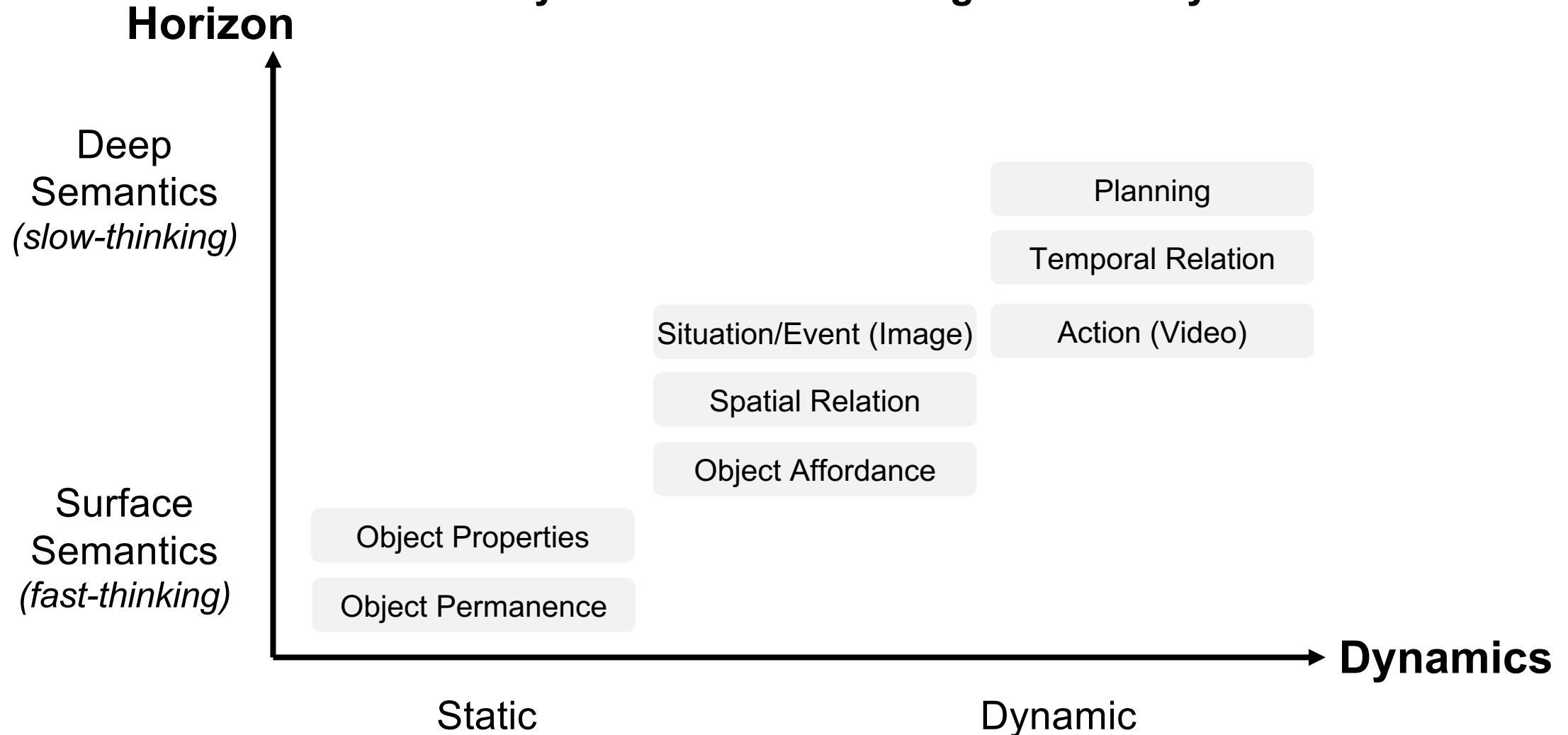


External World

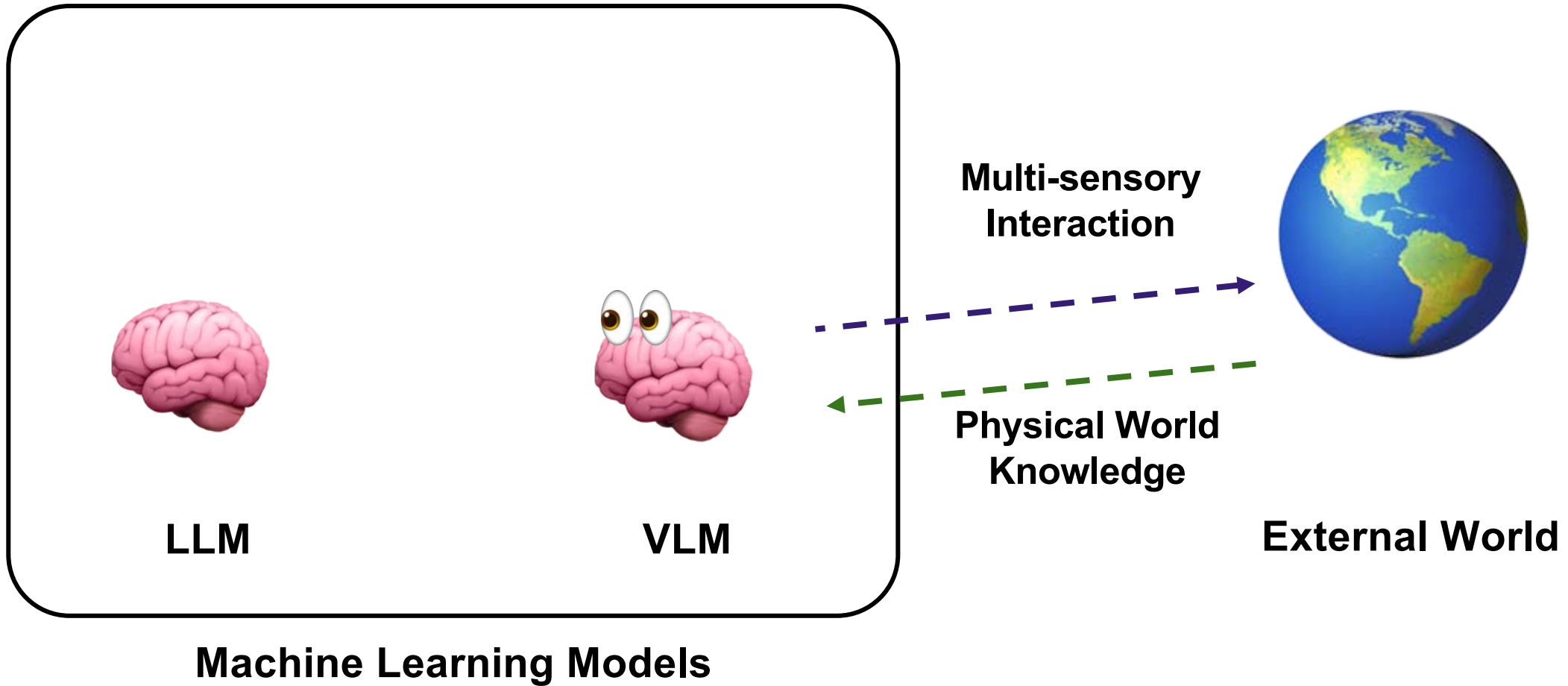
What is the knowledge that humans aim to learn?



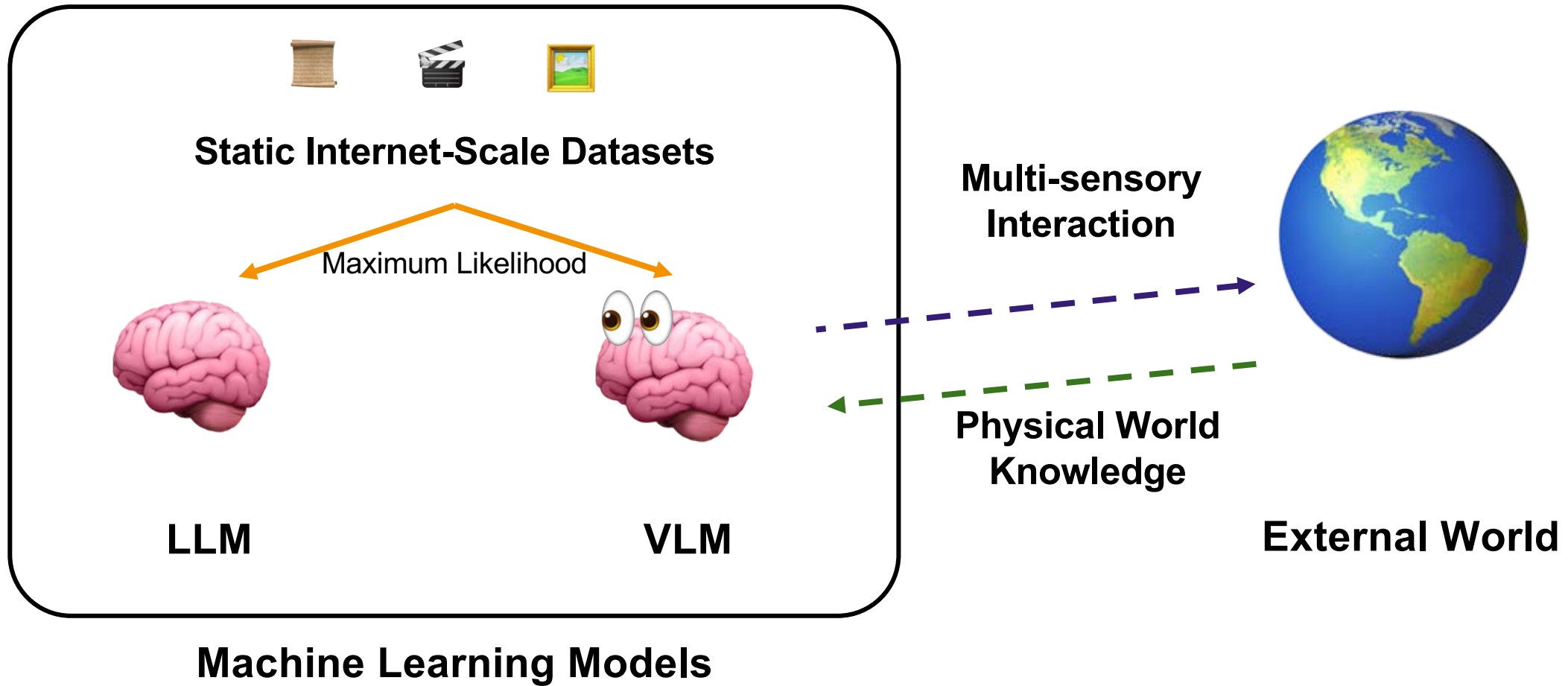
Physical World Knowledge Taxonomy



Machines learn knowledge through multi-sensory interactions



Machines learn knowledge through multi-sensory interactions



Video: A “Visual Recording” of World State Changes



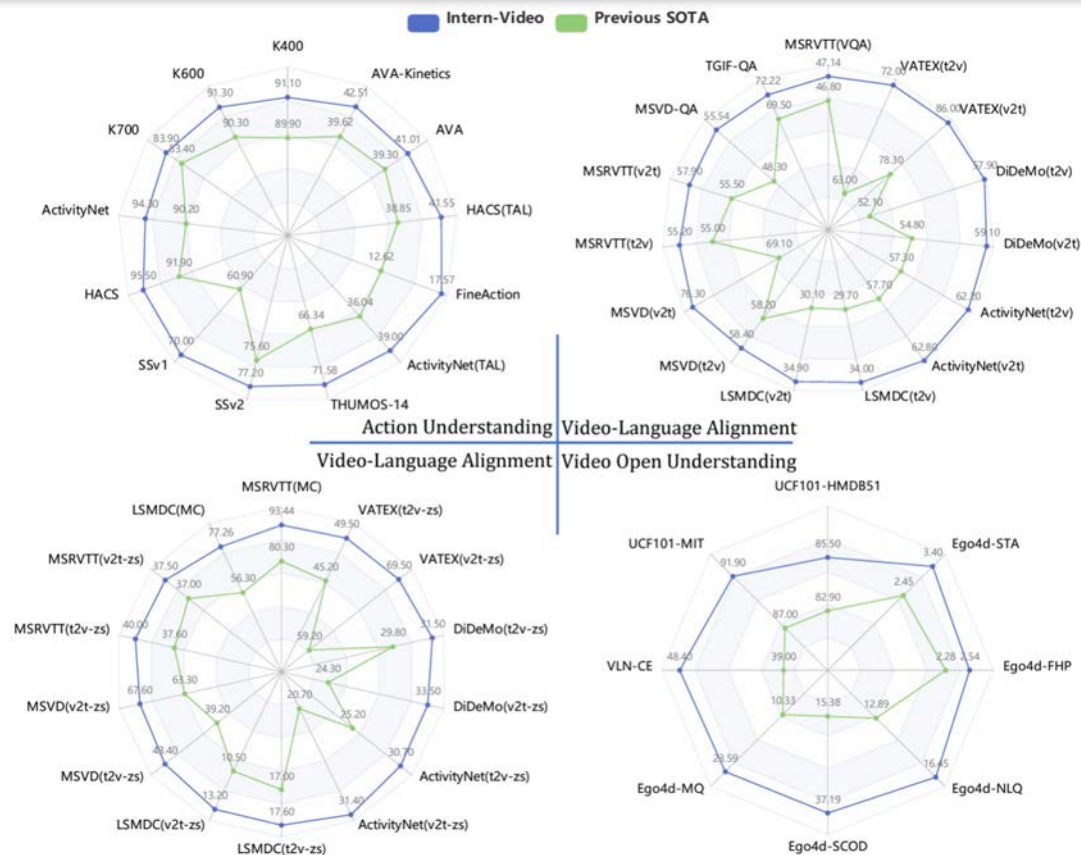
Where can we learn such physical world knowledge without interactive data \$\$\$?

proxy



“Book falling like a rock”

Video-Language Datasets



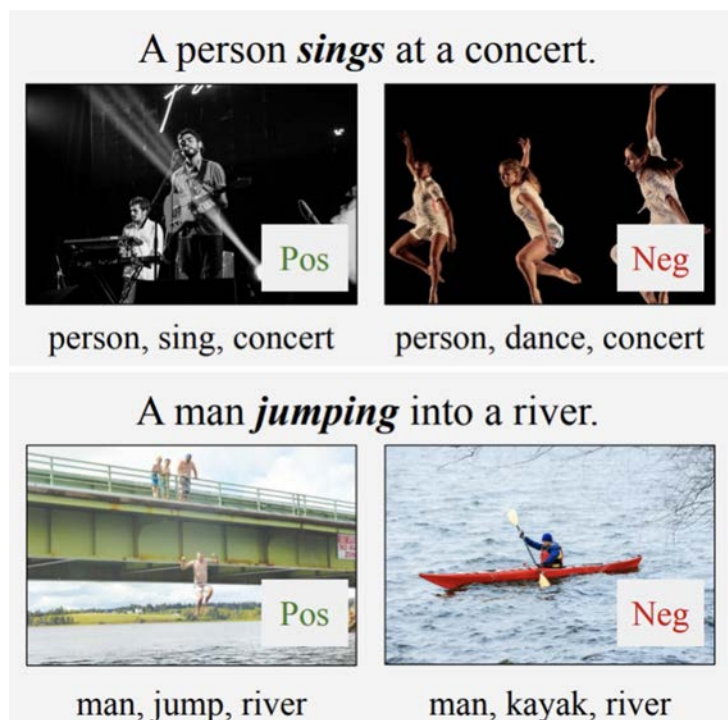
[InternVideo](#) (2023)

Video-Language Foundation Models

Existing V+L Models



Current models rely on object-centric abilities as a **shortcut** for V+L understanding.



Model	Verb Accuracy
MMT	60.8
Merged–MMT	60.7
Lang–MMT	64.5
Image–MMT	59.7

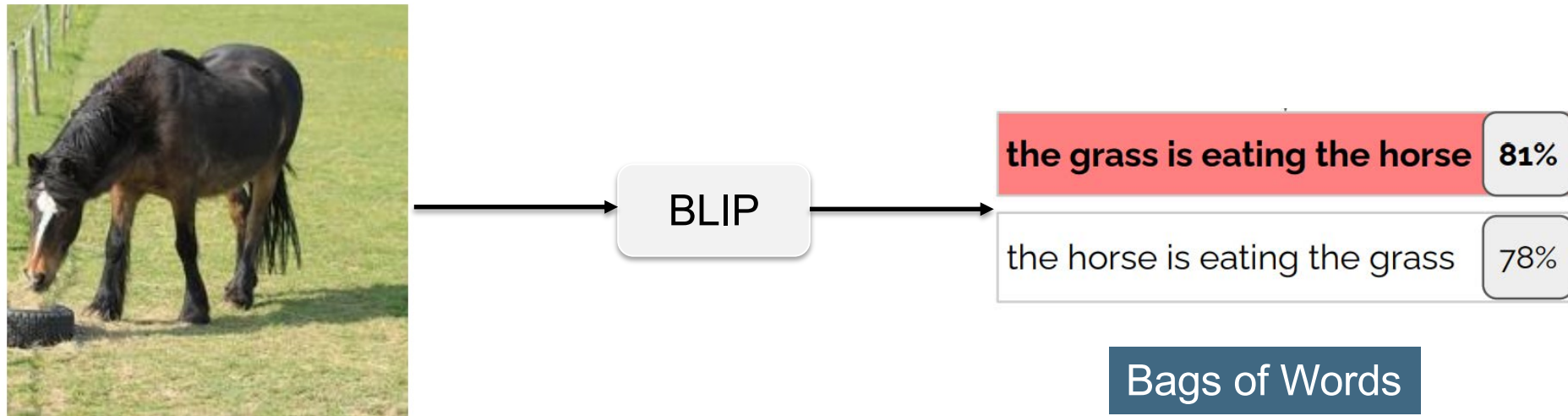
Low Verb Performance

“[Probing Image–Language Transformers for Verb Understanding](#)” Lisa Hendricks, et al. (arXiv 2021)

Existing V+L Models



Current models rely on object-centric abilities as a **shortcut** for V+L understanding.



[“When and why vision-language models behave like bags-of-words, and what to do about it”](#) Mert Yuksekgonul, et al. (ICLR 2023)

Existing V+L Models



Current models rely on object-centric abilities as a **shortcut** for V+L understanding.

Method	#PT	#Train Frame	MSRVTT			DiDeMo			ActivityNet Cap		
			R1	R5	R10	R1	R5	R10	R1	R5	R10
HERO [37]	136M	310	20.5	47.6	60.9	-	-	-	-	-	-
ClipBERT [31]	0.2M	16/16/8	22.0	46.8	59.9	20.4	48.0	60.8	21.3	49.0	63.5
VideoCLIP [61]	136M	960	30.9	55.4	66.8	-	-	-	-	-	-
Frozen [4]	5M	4	31.0	59.5	70.5	31.0	59.8	72.4	-	-	-
AlignPrompt [34]	5M	8	33.9	60.7	73.2	35.9	67.5	78.8	-	-	-
All-in-one [58]	138M	9	34.4	65.4	75.8	32.7	61.4	73.5	22.4	53.7	67.7
CLIP4Clip [47]	400M	12/64/64	42.0	68.6	78.7	42.8	68.5	79.2	40.5	72.4	98.2
SINGULARITY	5M	1	36.8	65.9	75.5	47.4	75.2	84.0	43.0	70.6	81.3
SINGULARITY	17M	1	41.5	68.7	77.0	53.9	79.4	86.9	47.1	75.5	85.5

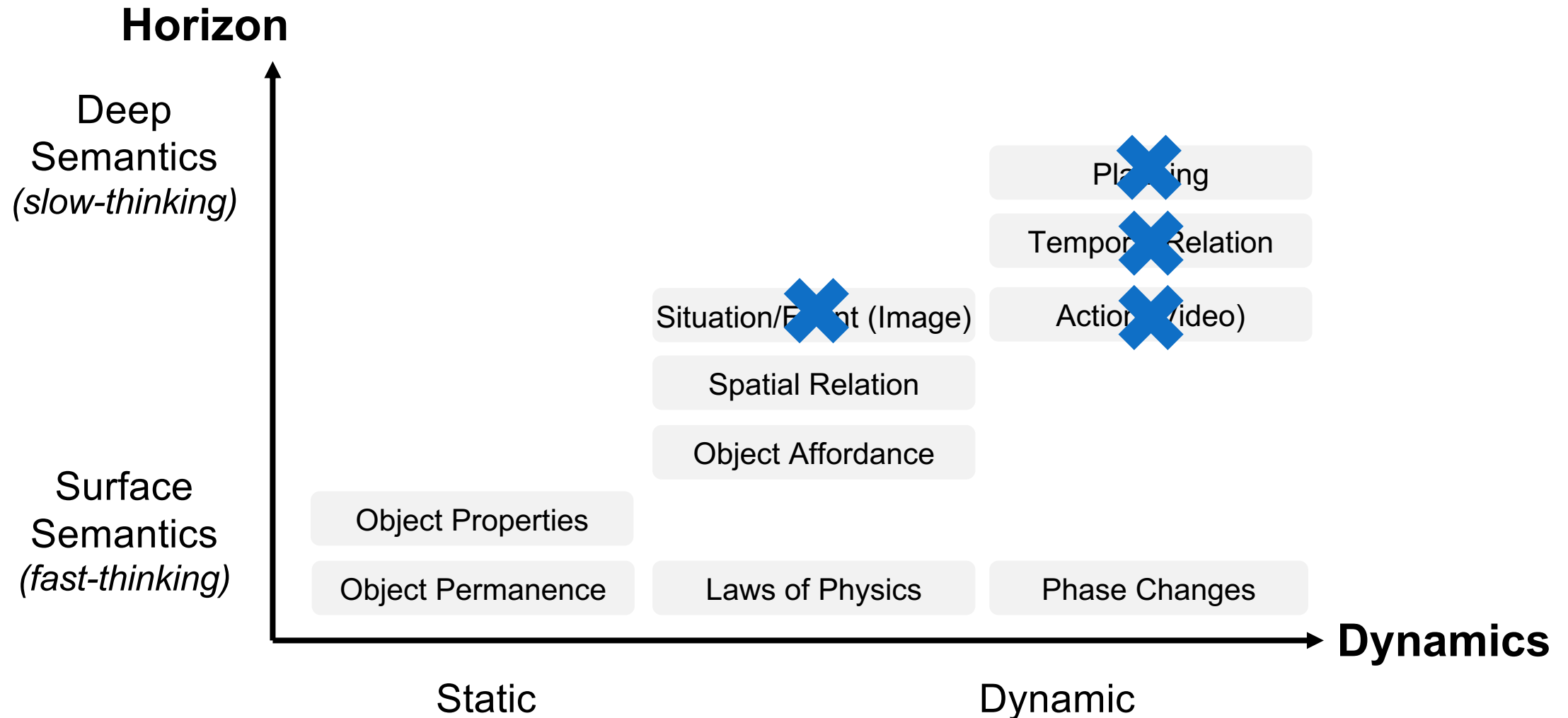
Single-Frame

[“Revealing Single Frame Bias for Video-and-Language Learning”](#) Jie Lei, et al. (ACL23)

Existing V+L Models



Physical Knowledge Taxonomy (🚧)



Surface

Deep



Surface

Deep



Object-Centric
Local
Static

Event-Centric
Situational
Dynamic

Knowledge can also help with V+L Pretraining



Compared to raw data, knowledge is **important and useful information.**

We learn three types of knowledge



Factual Knowledge are information about **instance-level facts** extracted from raw data.

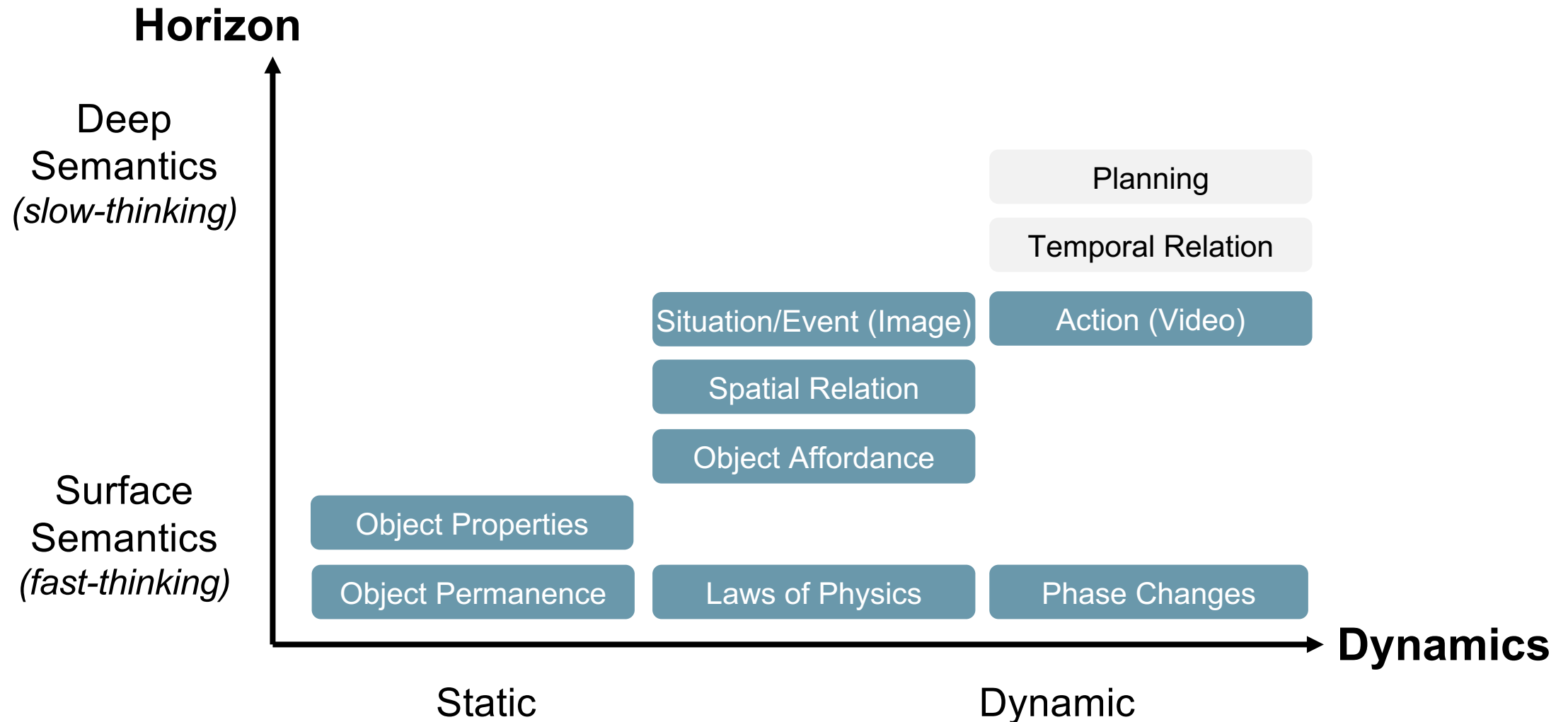
Text	Vision
Entity	Object
Relation	Scene Graph
Event	Activity/Situation
Affordance	Embodied AI

Factual Knowledge

Goal: Surface → Deep Semantic Knowledge



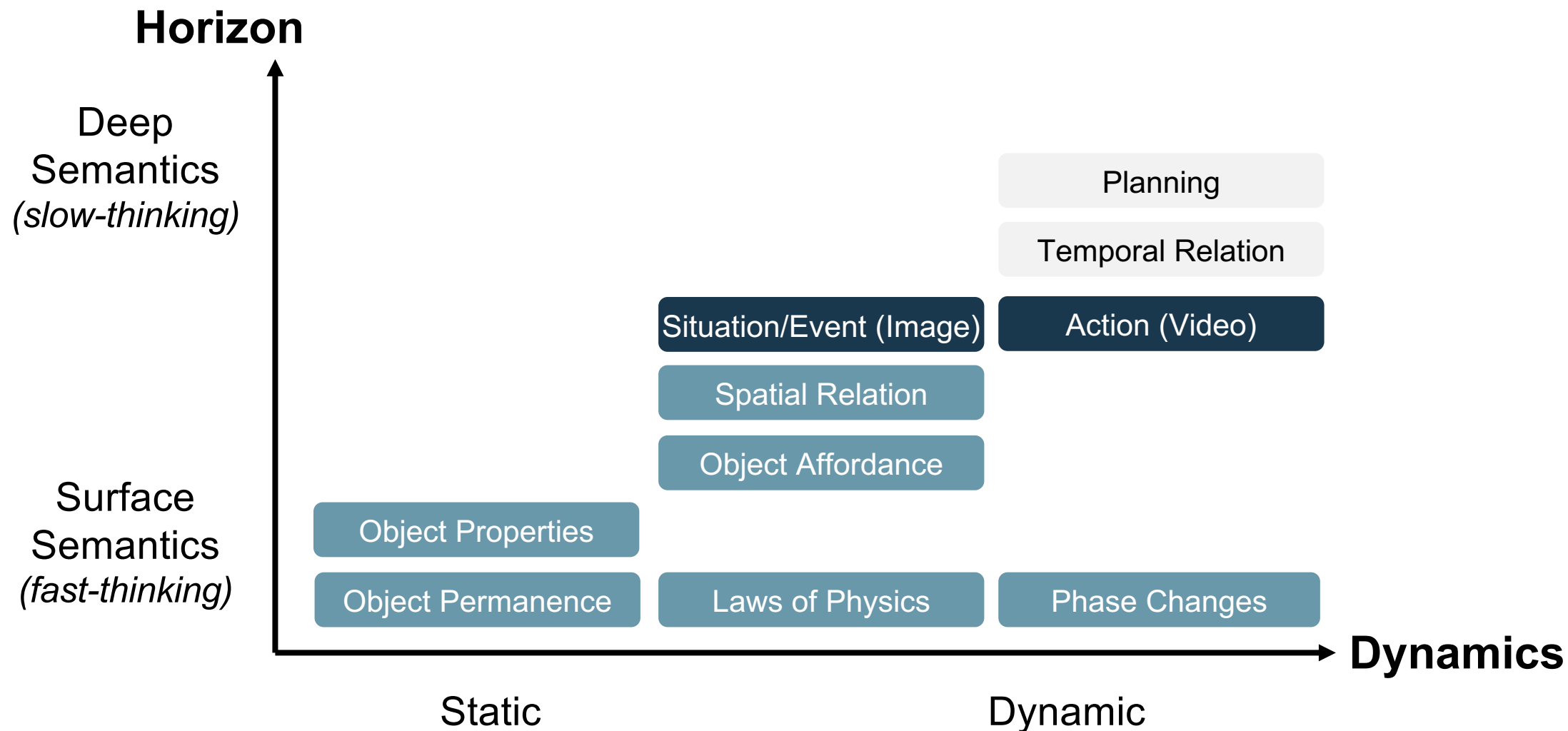
Physical Knowledge Taxonomy



Goal: Surface → Deep Semantic Knowledge



Physical Knowledge Taxonomy



Challenge 1: Complex Situation



Complex Situation

Event / Action

Semantics



Surface

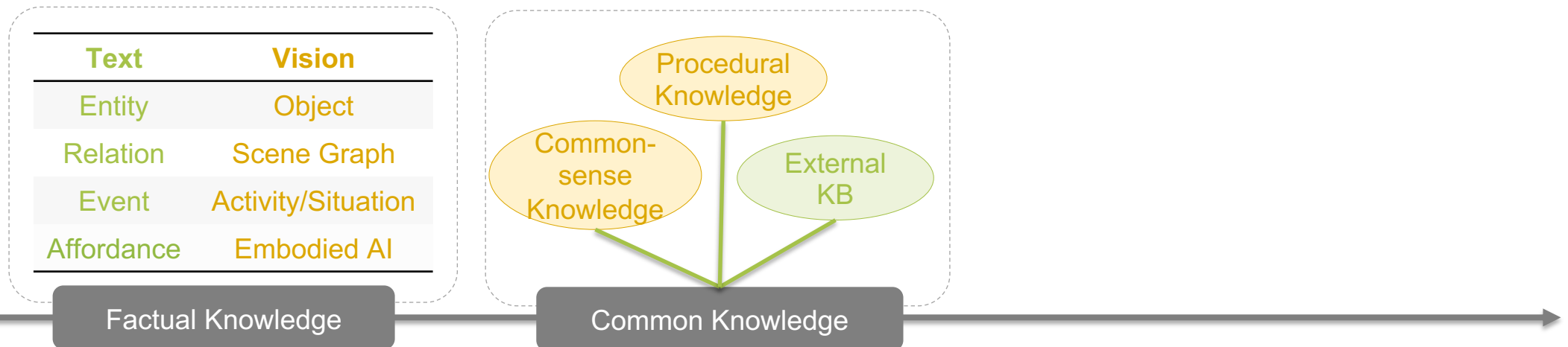


Deep

We learn three types of knowledge



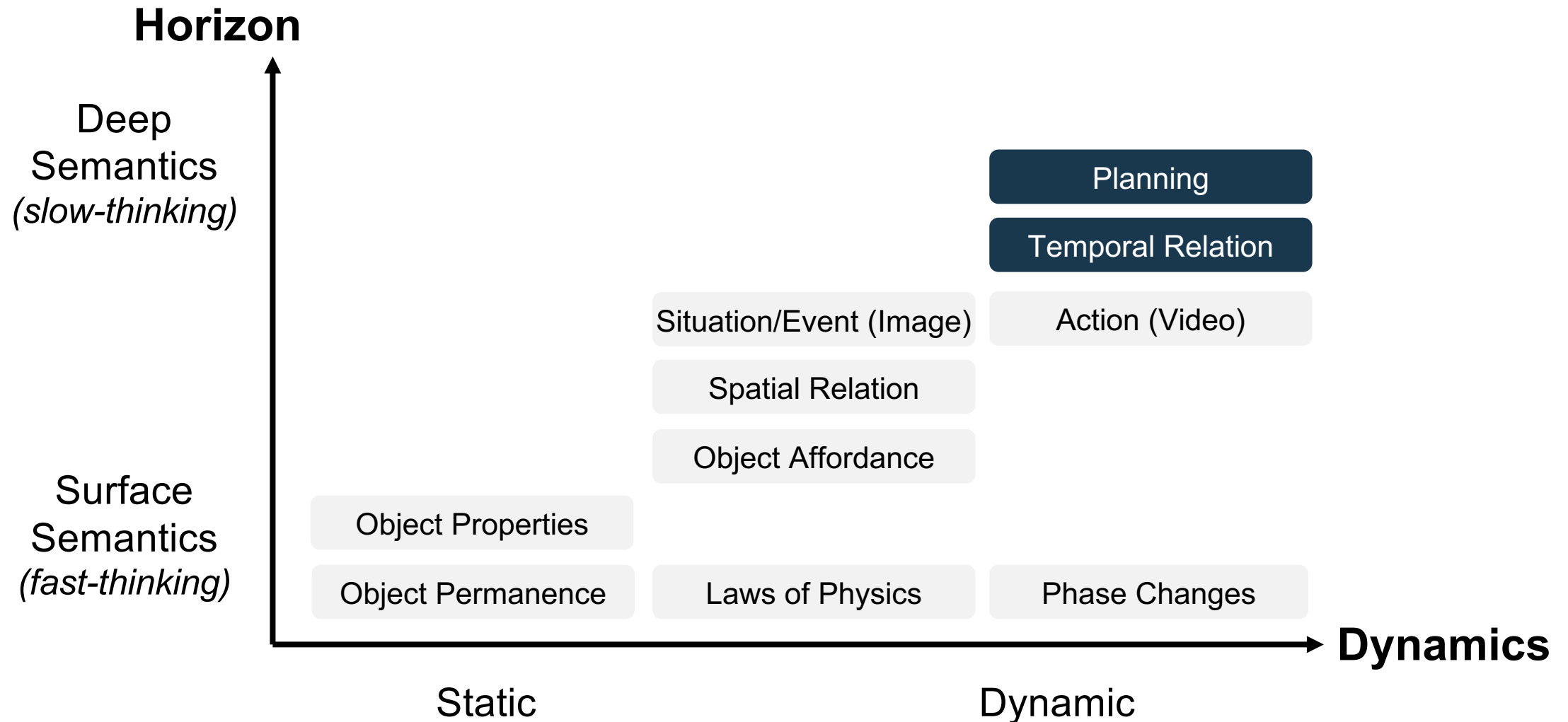
Common knowledge refers to knowledge of **common patterns** that is acquired or summarized from historical interaction with the world.



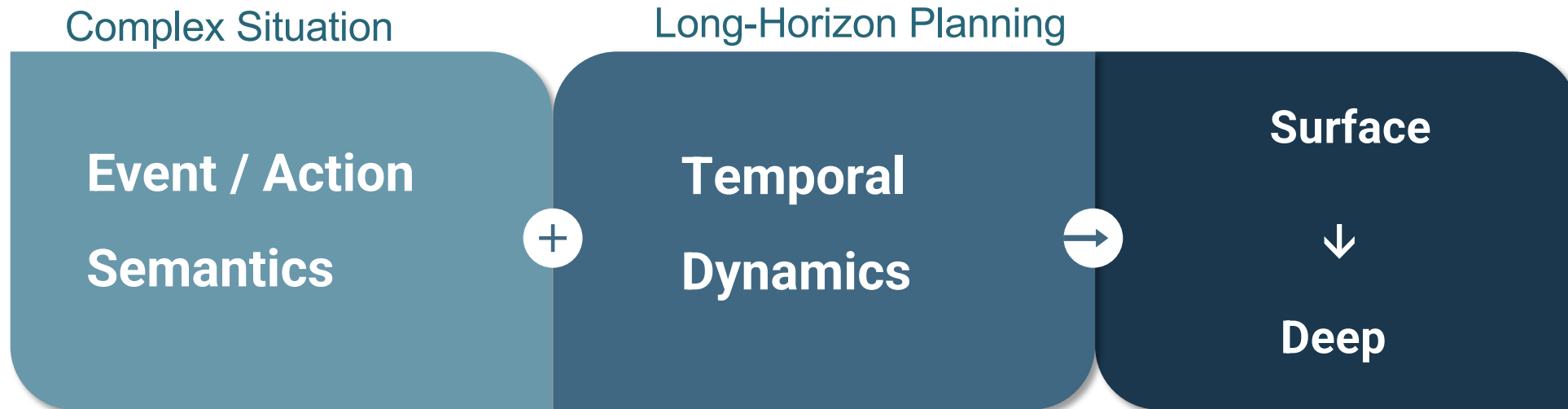
Goal: Surface → Deep Semantic Knowledge



Physical Knowledge Taxonomy



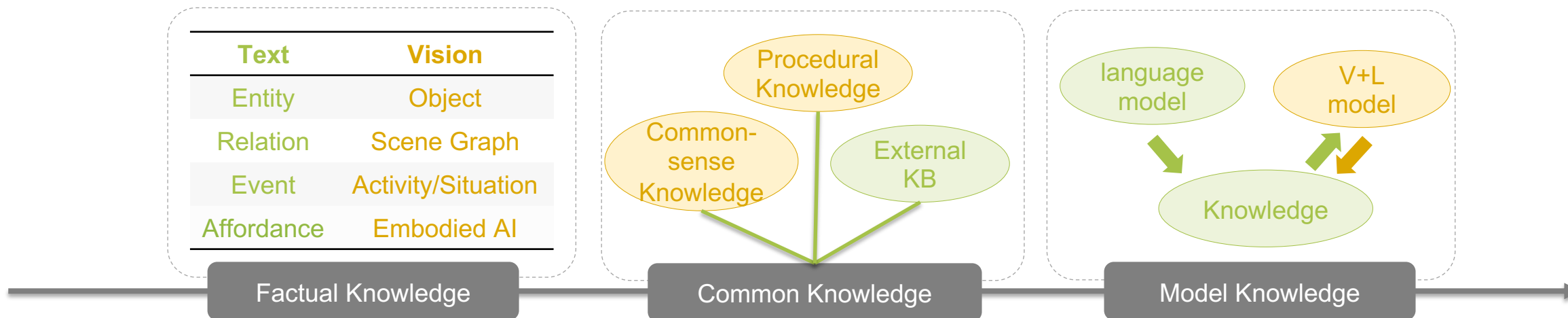
Challenge 2: Long-Horizon Planning



We learn three types of knowledge



Model Knowledge (parametric knowledge) is the knowledge embedded and encoded in models.



Can we borrow the ability from LLM?

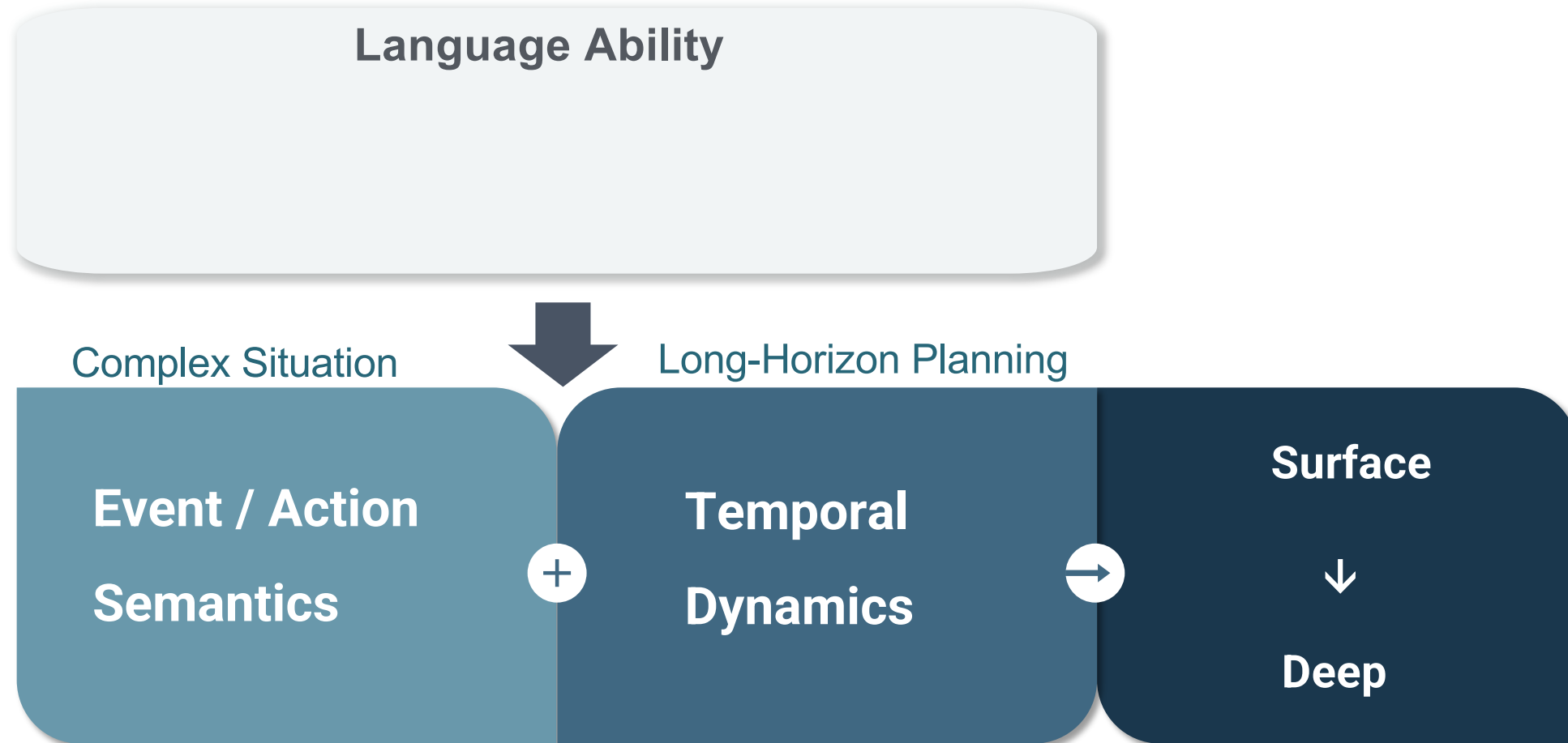


Large Language Models (LLMs) are very powerful.

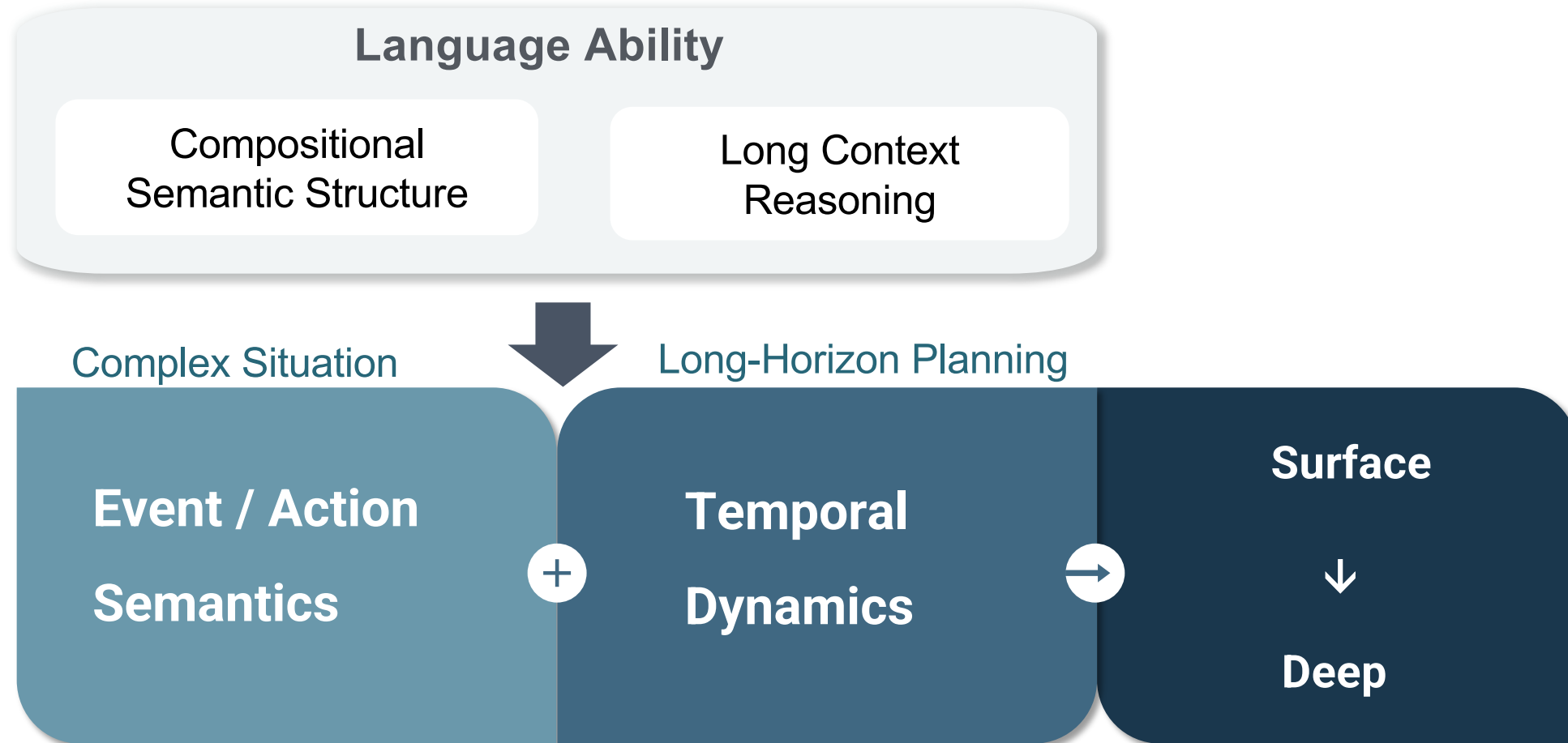


https://twitter.com/xiye_nlp

One Solution: Language As Supervision



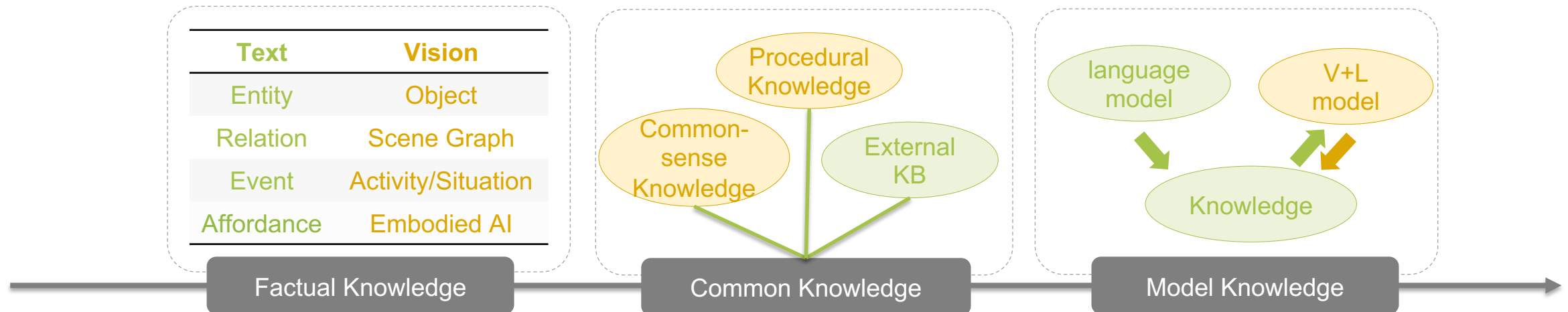
One Solution: Language As Supervision



We inject knowledge to V+L foundation models



We patch three types of knowledge into V+L foundation models.



Outline: Factual Knowledge



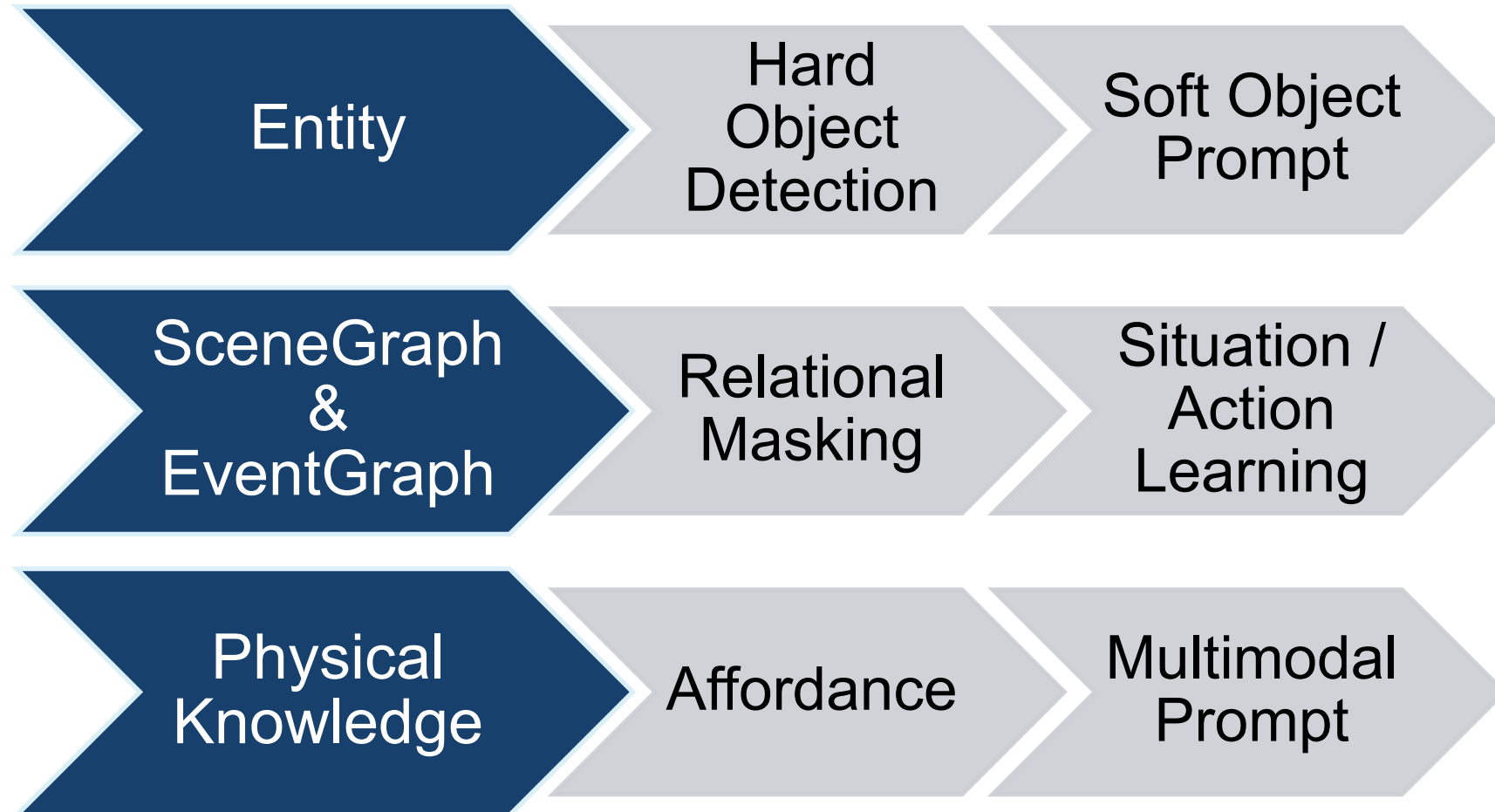
Factual Knowledge are information about instances extracted from raw data.

Text	Vision
Entity	Object
Relation	Scene Graph
Event	Activity/Situation
Affordance	Embodied AI

Factual Knowledge

Implicit Knowledge

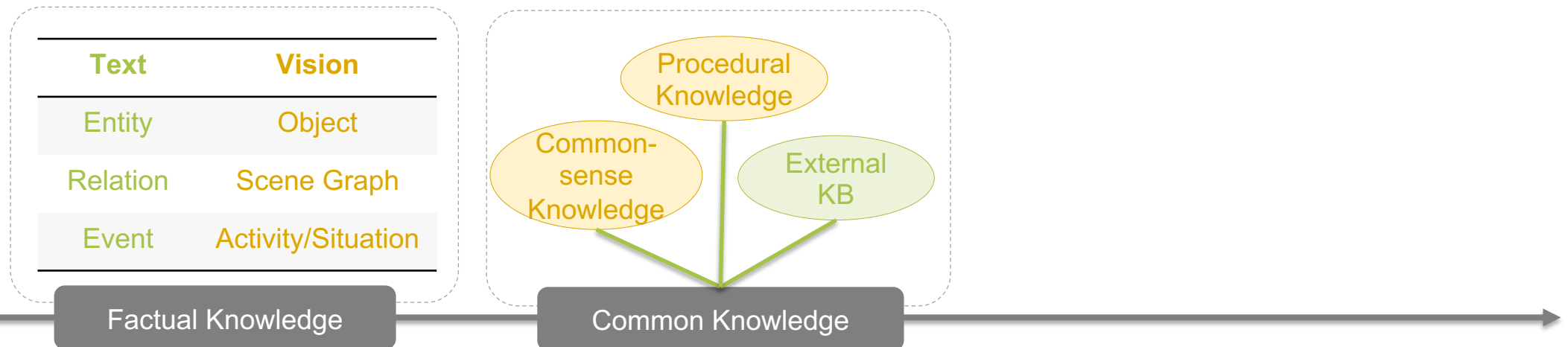
Outline: Factual Knowledge



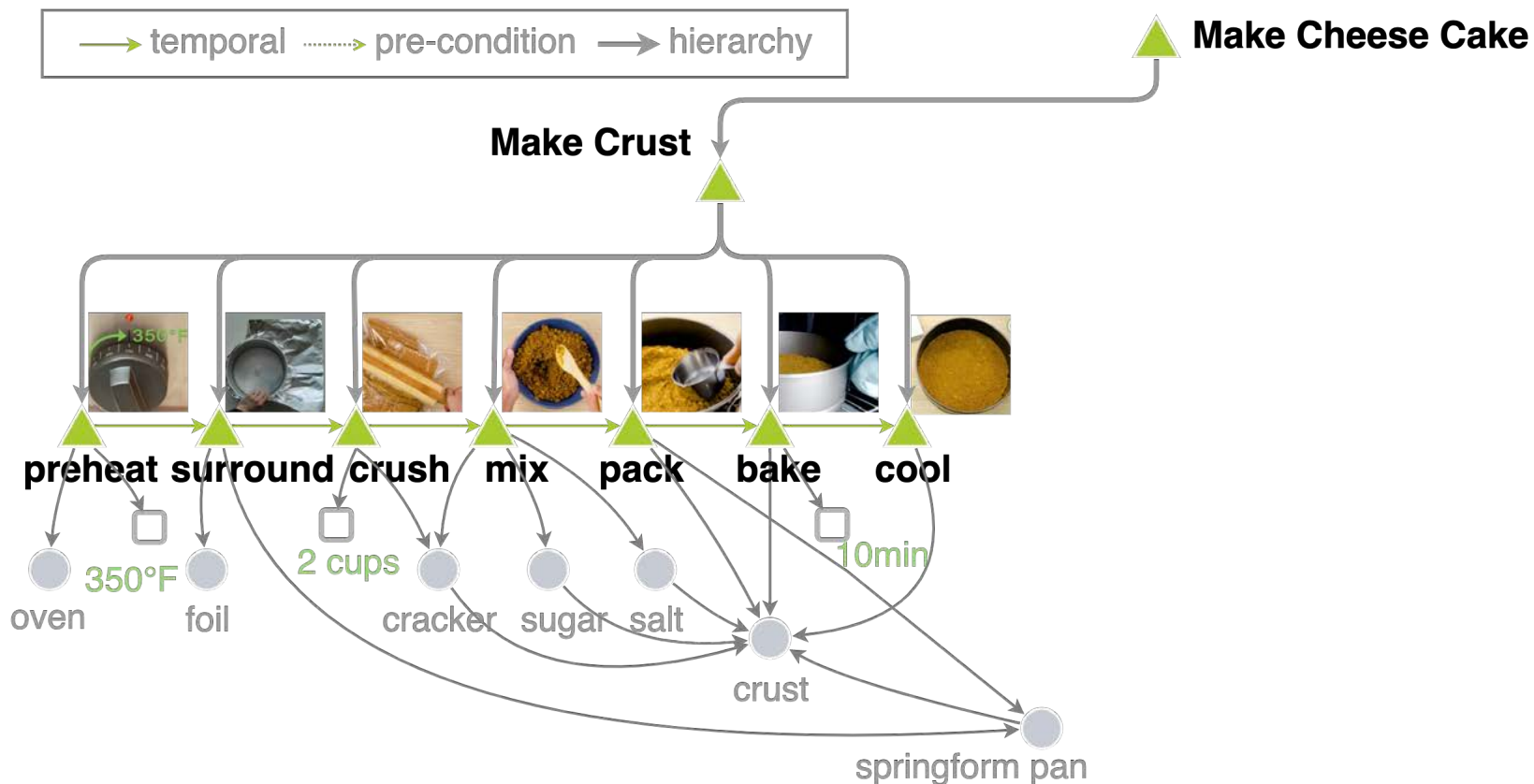
Outline: Common Knowledge



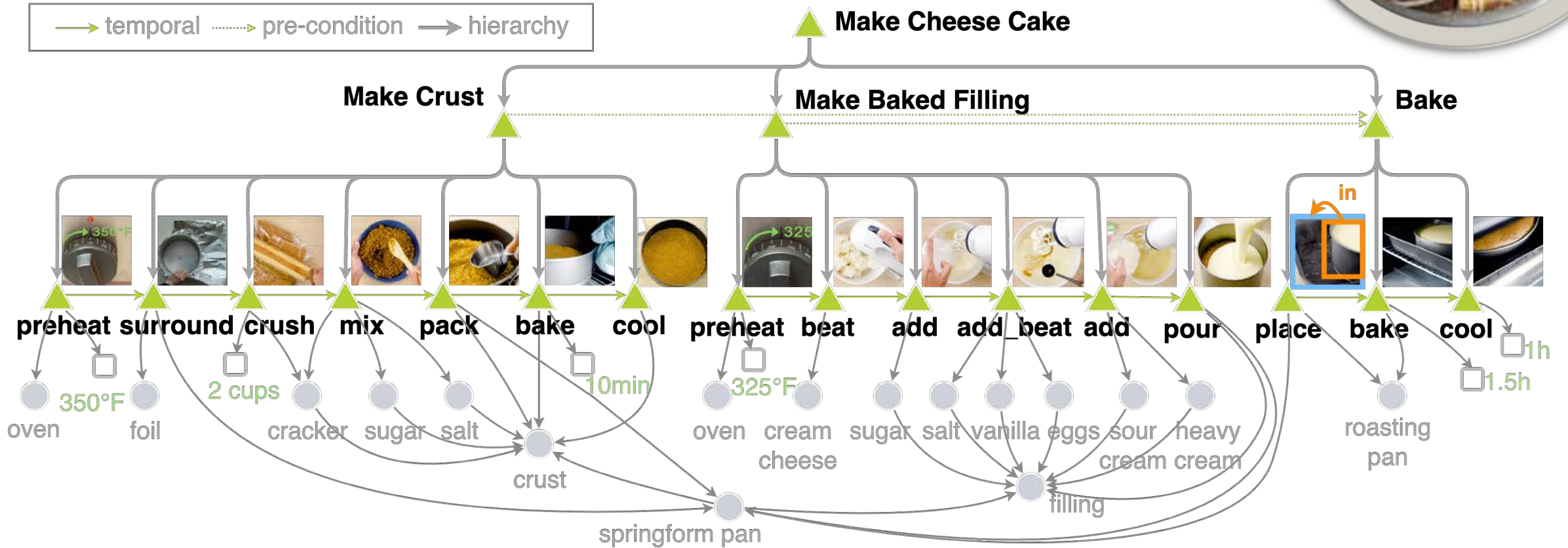
Common knowledge refers to knowledge of common patterns that is acquired or summarized from historical interaction with the world.



Common Knowledge: History repeats itself



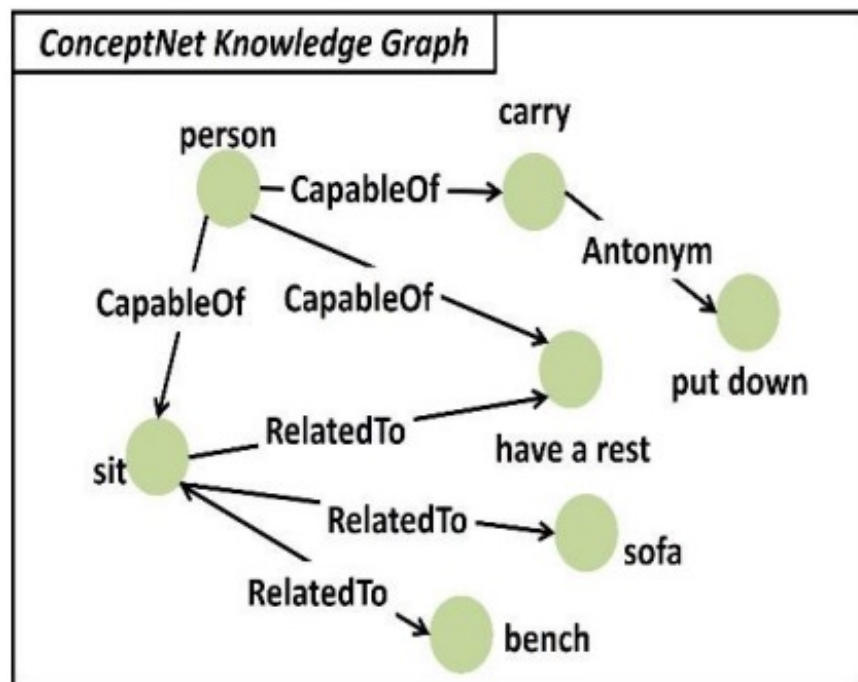
Common Knowledge: History repeats itself



Commonsense Knowledge



Commonsense knowledge includes facts about events occurring in time, about the effects of actions.

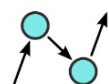


Why are [person1] and [person3] shaking hands?

- (a) [person1] and [person3] are presenting a trophy to someone.
- (b) [person1] and [person3] just made a deal.
- (c) [person1] and [person3] are old friends seeing each other for the first time in a long time.
- (d) They have just met and are greeting each other.**

I think so because ...

- (a) People like to greet each other when they meet by shaking hands.**
- (b) They look like they are shaking hands to introduce themselves.
- (c) They are meeting each other for the first time.
- (d) Some people shake hands to greet one another by grasping each others' arms.



ConceptNet
An open, multilingual knowledge graph



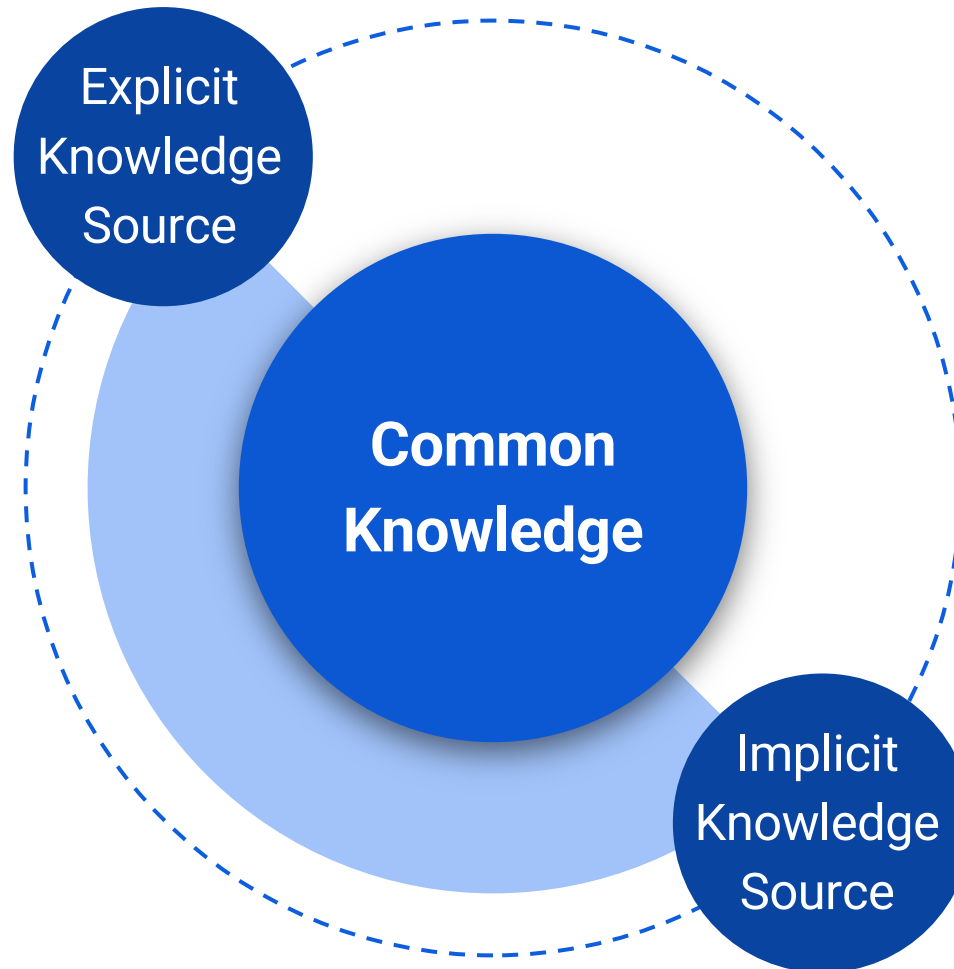
How to add common knowledge?



- Two ways to learn procedural knowledge

Use knowledge:

- As Data
- As Supervision
- In Model

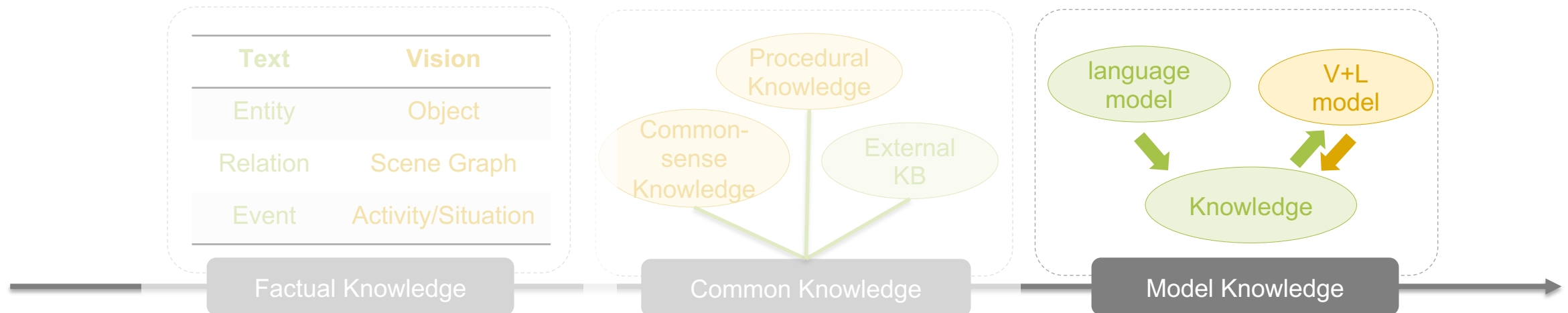


Learning from massive data

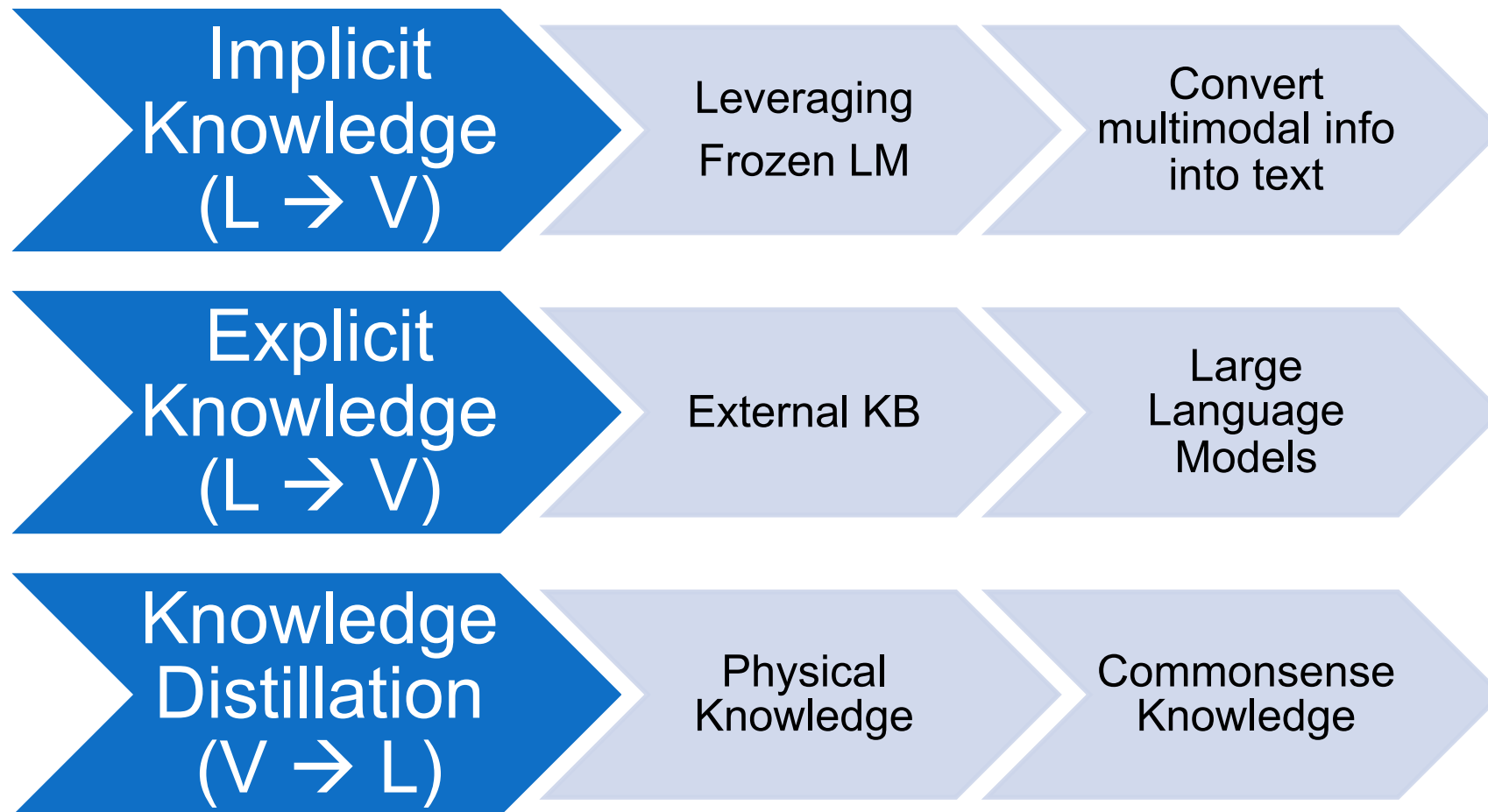
Model Knowledge



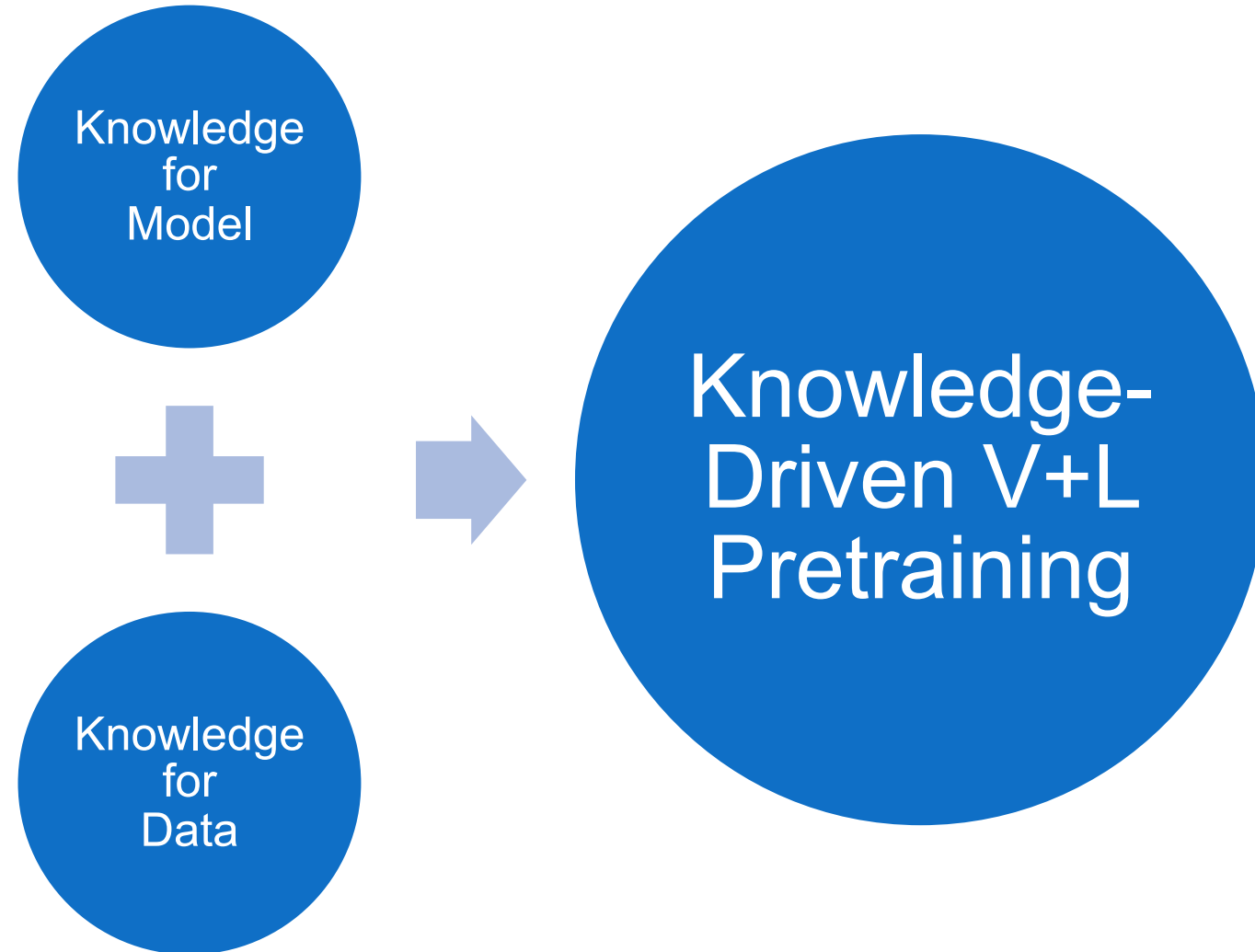
Model Knowledge is the knowledge embedded and encoded in models.



Model Knowledge



Summary: How to learn multimedia embedding?



- On the model side, adding knowledge can guide the model where to focus.
 - [Compositional Multi-Granularity Semantic Knowledge](#) (such as verb, adjectives, etc)
 - [Long Horizon Reasoning](#) (such as temporal dynamics, etc)
 - [Parametric Knowledge Controlling](#) (such as parameter editing, etc)
- On the data side, knowledge is useful in the following ways:
 - In-context prompt
 - Data augmentation
 - Data selection
 - Effective Feedback

Jun 2023

CVPR Tutorials

Knowledge-Driven Vision-Language Encoding

CVPR

Factual Knowledge in V+L Pretraining: Information about Instances

Knowledge-Driven Vision-Language Pretraining (Part II)

Manling Li
UIUC

manling2@illinois.edu



Northwestern
University



COLUMBIA
UNIVERSITY

 Meta AI

Timetable

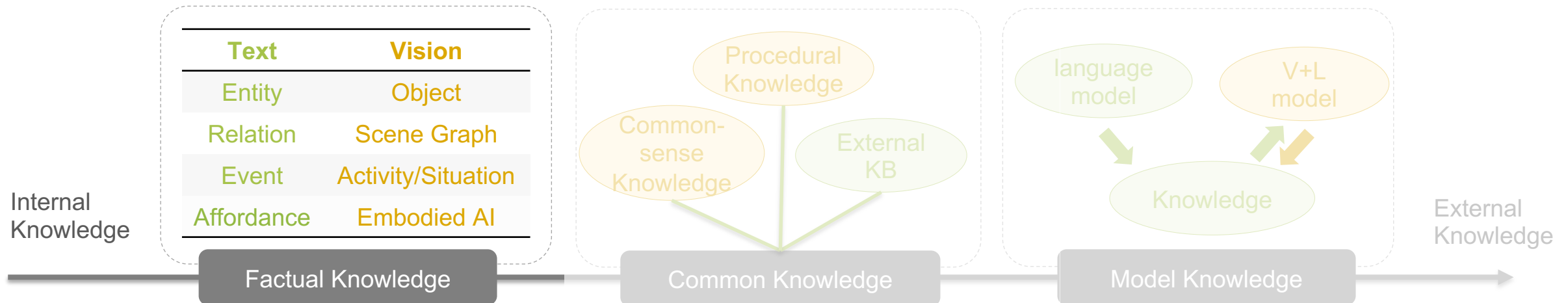


Content	Time	Presenter
Motivation and Overview	15min	Manling Li
Factual Knowledge	30min	Manling Li
Commonsense Knowledge	15min	Manling Li
Procedural Knowledge	30min	Xudong Lin
Model Knowledge	30min	Jie Lei
Panel: Knowledge vs Large Models	15min	Mohit Bansal, Carl Vondrick, Xudong Lin
Panel: LLMs for multimodal	15min	Mohit Bansal, Carl Vondrick, Jie Lei
Panel: Image vs Video vs Audio vs Others	15min	Mohit Bansal, Carl Vondrick, Xudong Lin
Panel: Open Challenges	15min	Mohit Bansal, Carl Vondrick, Jie Lei
QA	30min	All

Factual Knowledge



Compared to raw data, knowledge is **important and useful information**.



What is “Event” Knowledge?



What happened?



What is “Event” Knowledge?



What happened?

Yes! A protest.

What are they protesting for?



What is “Event” Knowledge?



What happened?

Yes! A protest.

What are they protesting for?



vaccine



What is "Event" Knowledge?



What happened?

Yes! A protest.

What are they protesting for?



vaccine



negation

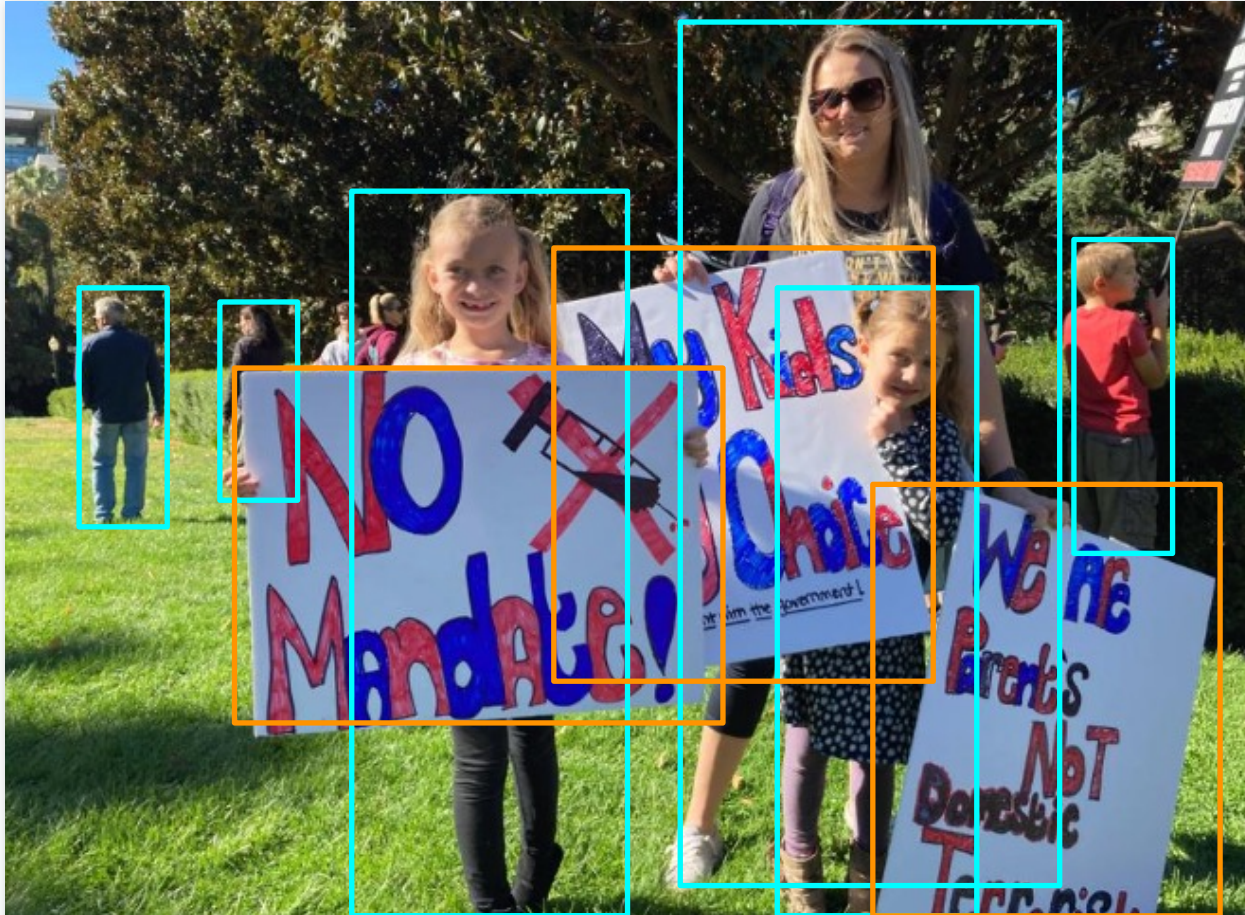


Existing object-centric info miss situational understanding



Vision

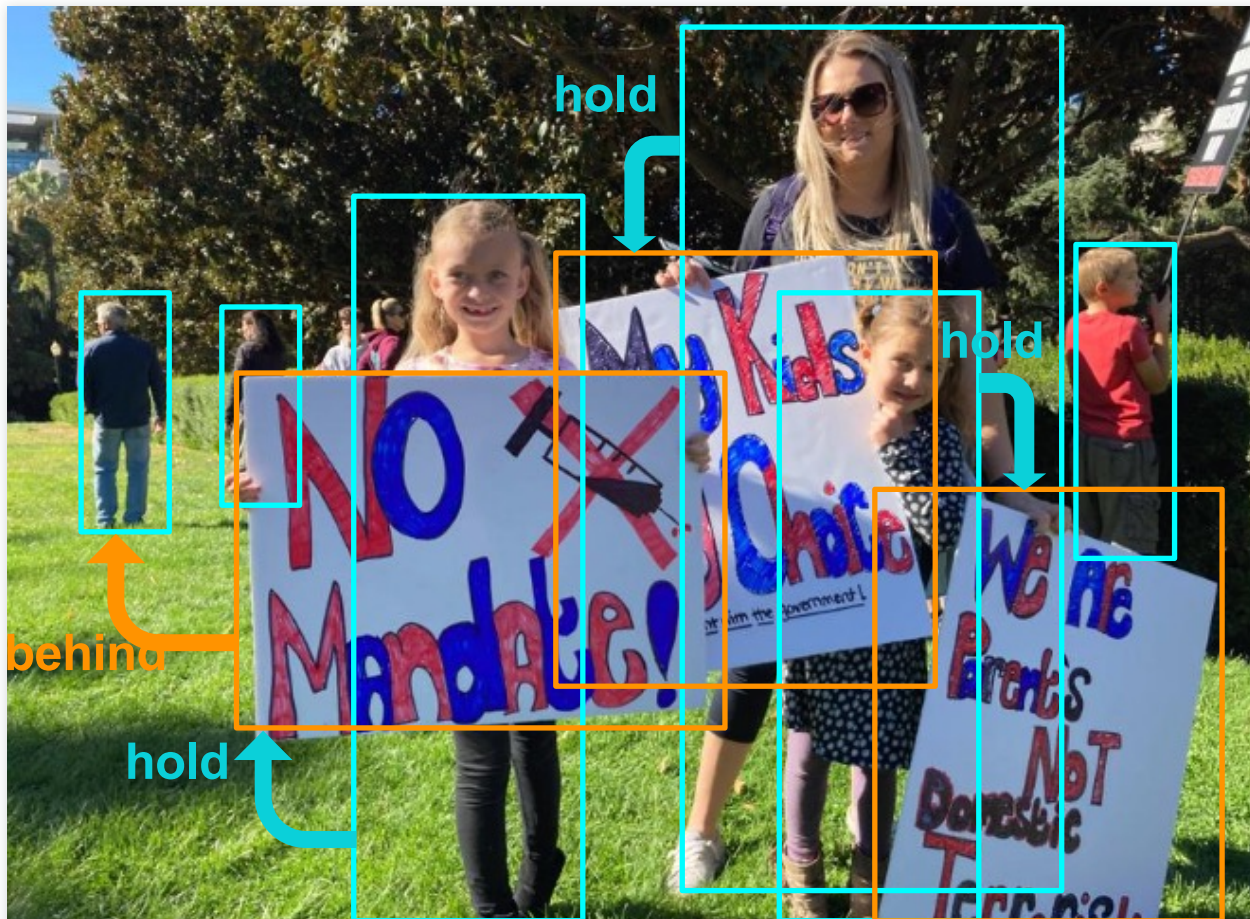
Existing object-centric info miss situational understanding



Vision

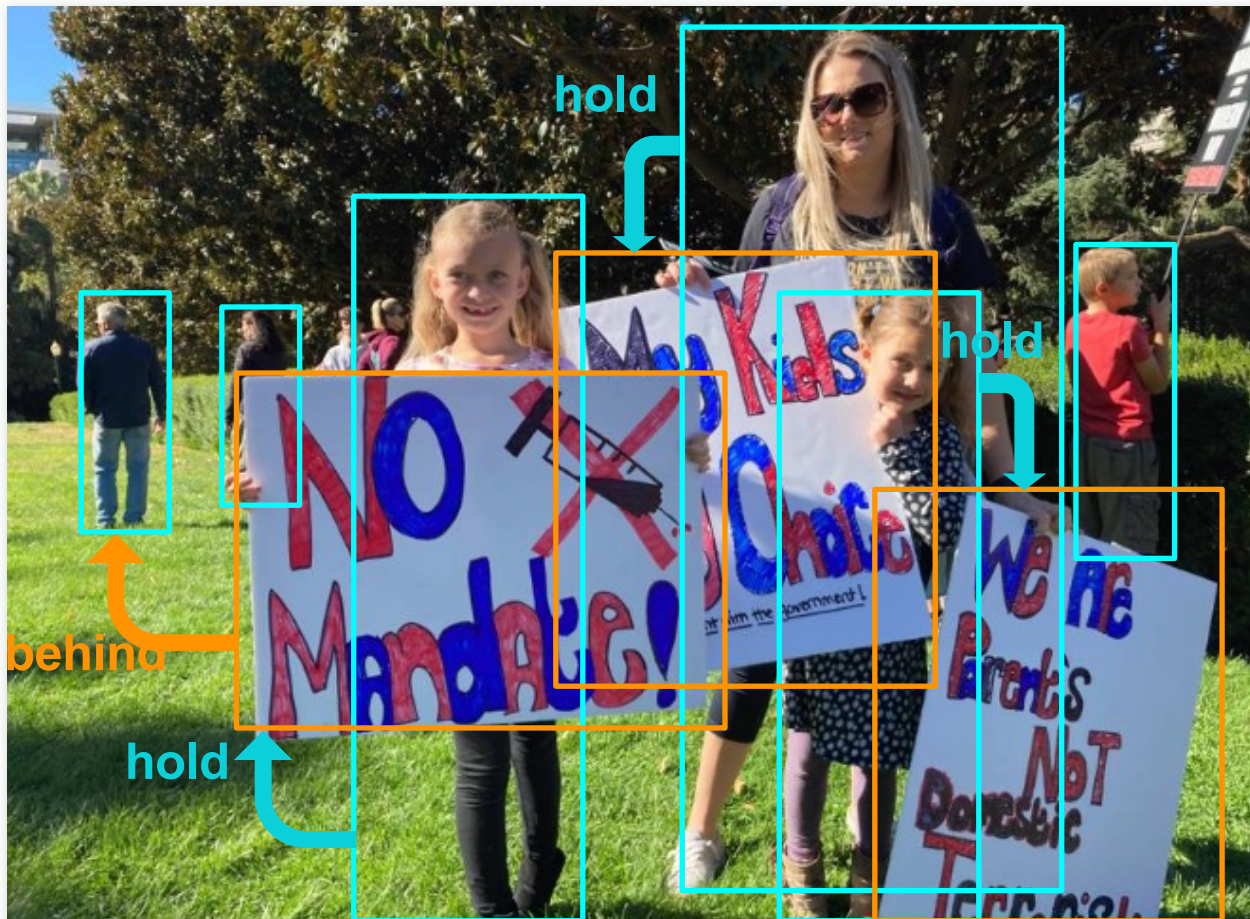
Object

Existing object-centric info miss situational understanding



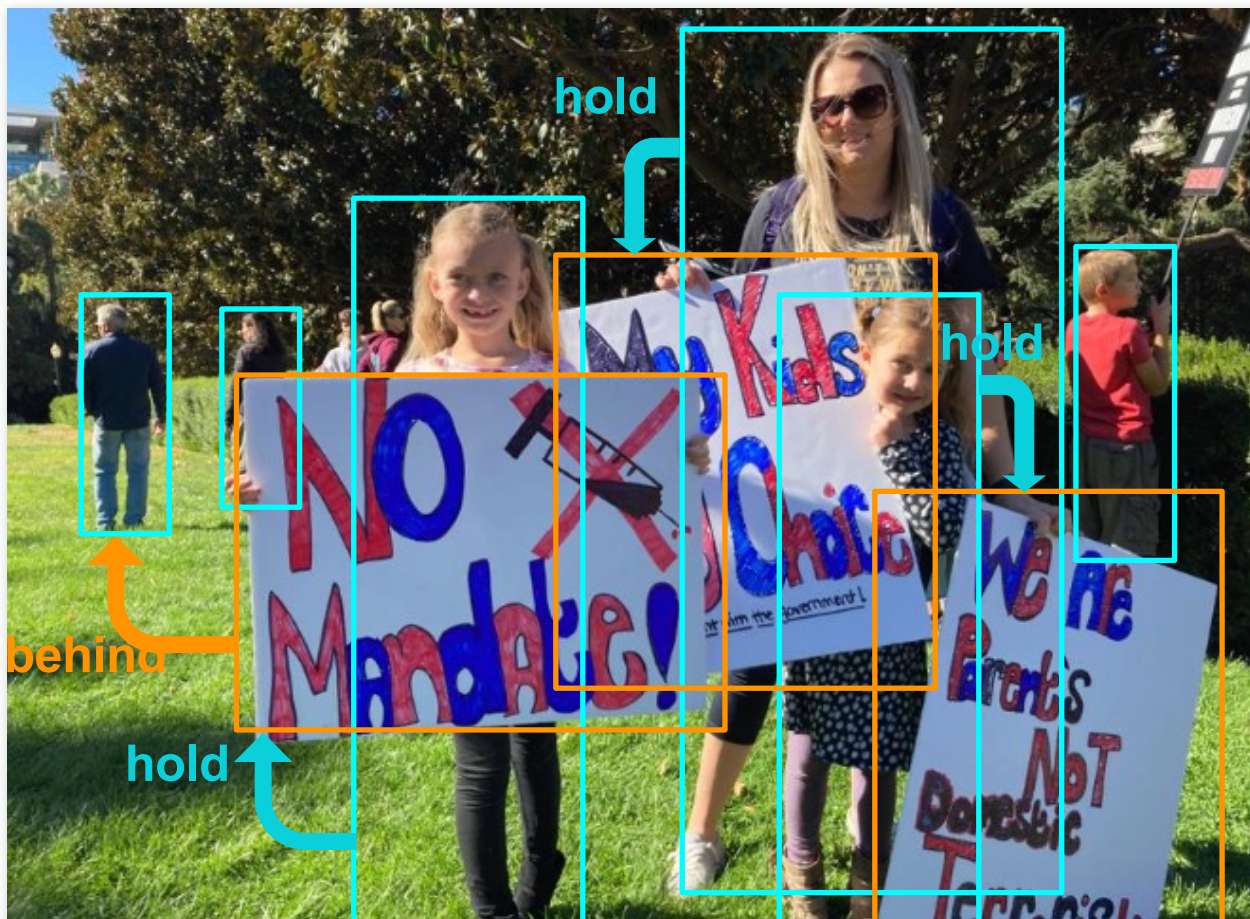
Vision
Object
Relation
Scene Graph

Existing object-centric info miss situational understanding



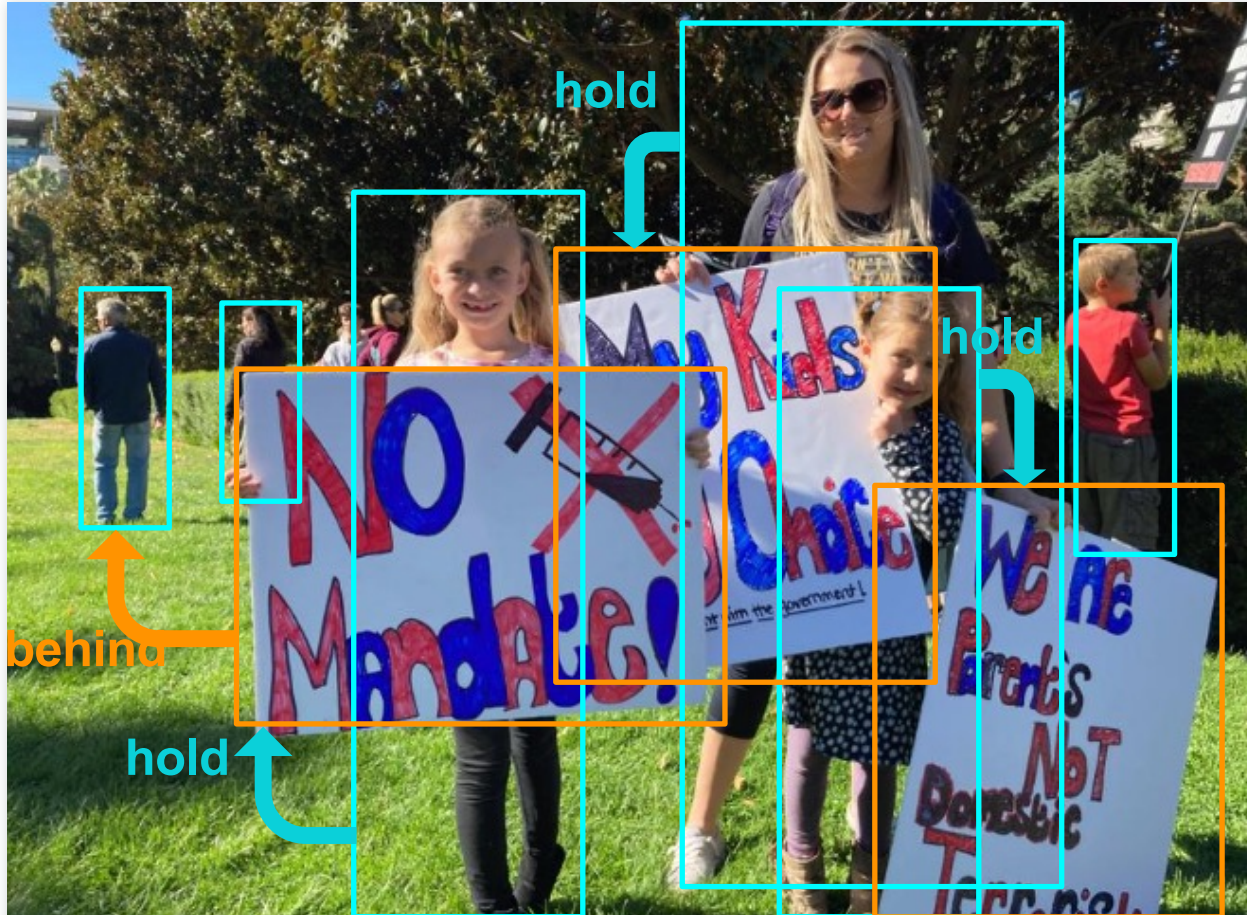
Vision	Text
Object	
Relation	
Scene Graph	

Existing object-centric info miss situational understanding



Vision	Text
Object	Entity
Relation	
Scene Graph	

Existing object-centric info miss situational understanding



Vision	Text
Object	Entity
Relation	Relation
Scene Graph	Entity-Relation Graph

Existing object-centric info miss situational understanding



	Vision	Text
Entity-centric	Object	Entity
	Relation	Relation
	Scene Graph	Entity-Relation Graph

Existing object-centric info miss situational understanding



	Vision	Text
Entity-centric	Object	Entity
	Relation	Relation
	Scene Graph	Entity-Relation Graph

State-of-the-art Captioner (Kamath et al., 2022)

Answer ▾

a woman holding a sign in front of a group of people.

a woman holding a sign while standing in a park.

a woman holding a sign in front of a crowd.

Definition of "Event"



Event

Protest

What happened?

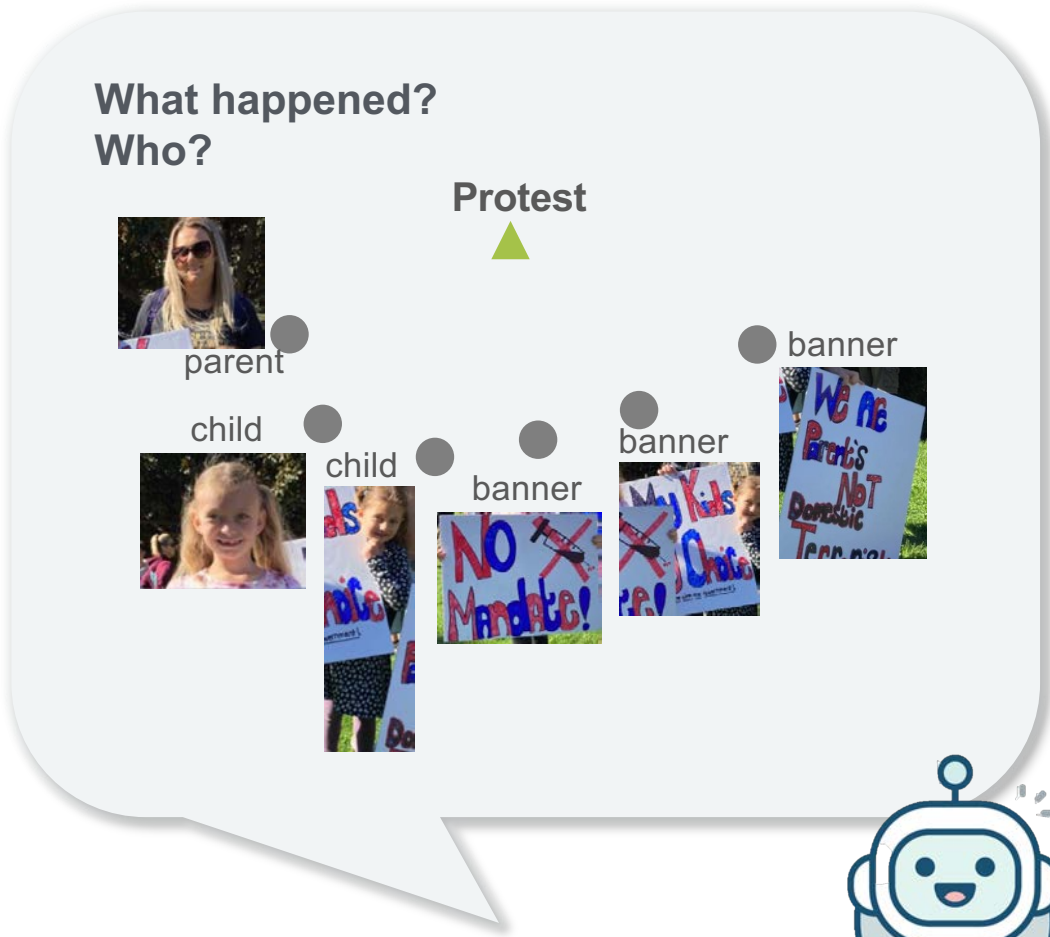
Protest



Definition of “Event”



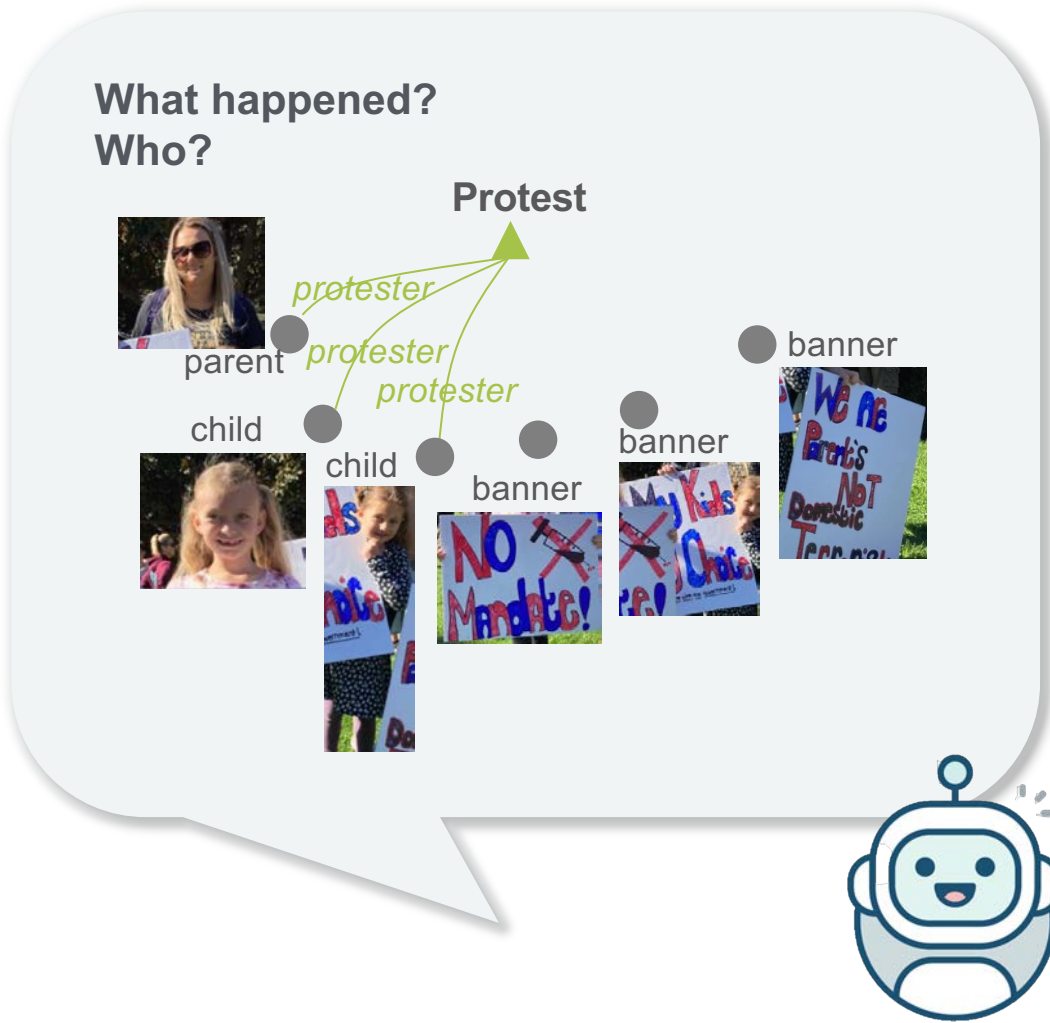
Event	Protest
	parent
	child
	banner
	No vaccine mandate for kids



Definition of “Event”



Event	Protest
Protester	parent
Protester	child
	banner
	No vaccine mandate for kids

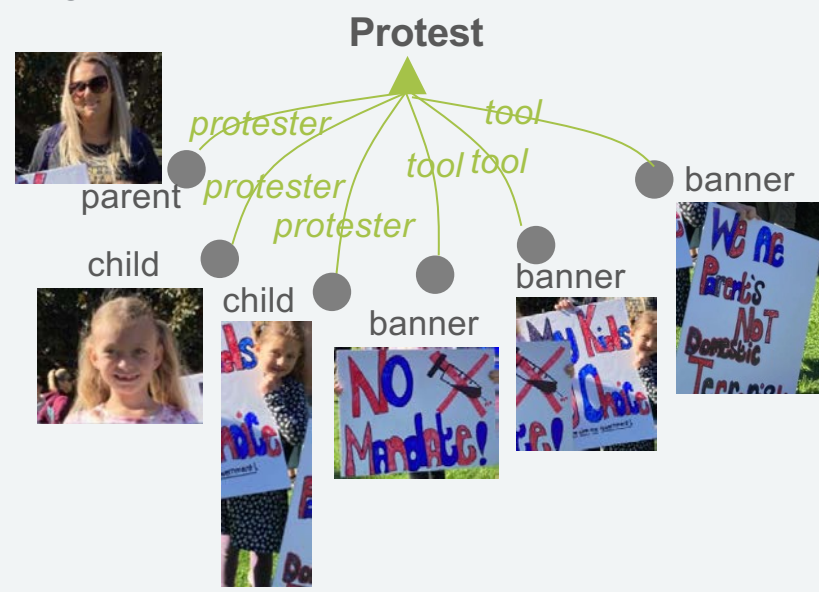


Definition of “Event”



Event	Protest
Protester	parent
Protester	child
Tool	banner
Topic	No vaccine mandate for kids

What happened?
Who?



A New Task of Multimodal Event Extraction [ACL'20]

Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. **Event Type** (*e.g., protest*)

2. **Participants** (*e.g., child*) & **Semantic Roles** (*e.g., protester*)

A New Task of Multimodal Event Extraction [ACL'20]

Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. **Event Type** (*e.g., protest*)

2. **Participants** (*e.g., child*) & **Semantic Roles** (*e.g., protester*)

What is Multimodal Event Extraction? [Li et al, ACL'20]

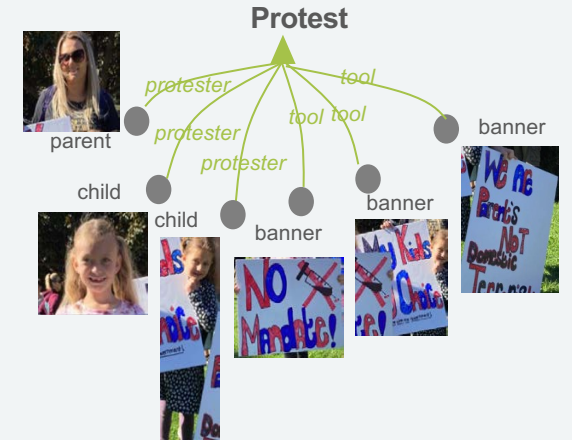
Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. Event Type (e.g., *protest*)

2. Participants (e.g., *child*) & Semantic Roles (e.g., *protester*)



What is Multimodal Event Extraction? [Li et al, ACL'20]

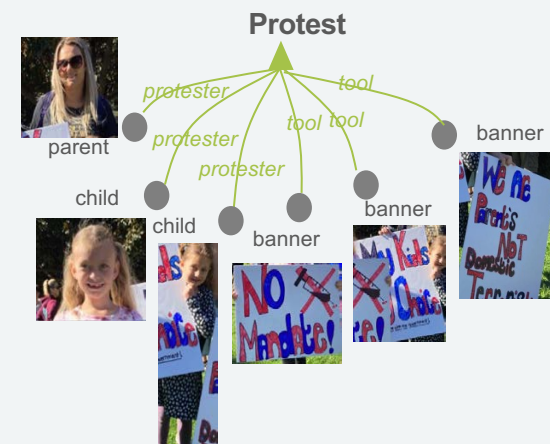
Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. **Event Type** (e.g., *protest*)

2. **Participants** (e.g., *child*) & **Semantic Roles** (e.g., *protester*)



What is Multimodal Event Extraction? [Li et al, ACL'20]

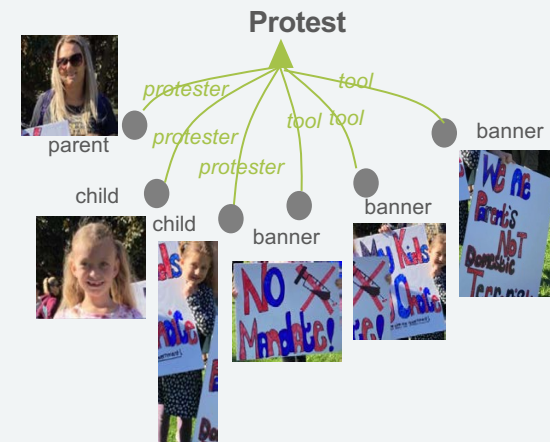
Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. **Event Type** (e.g., *protest*)

2. **Participants** (e.g., *child*) & **Semantic Roles** (e.g., *protester*)



What is Multimodal Event Extraction? [Li et al, ACL'20]

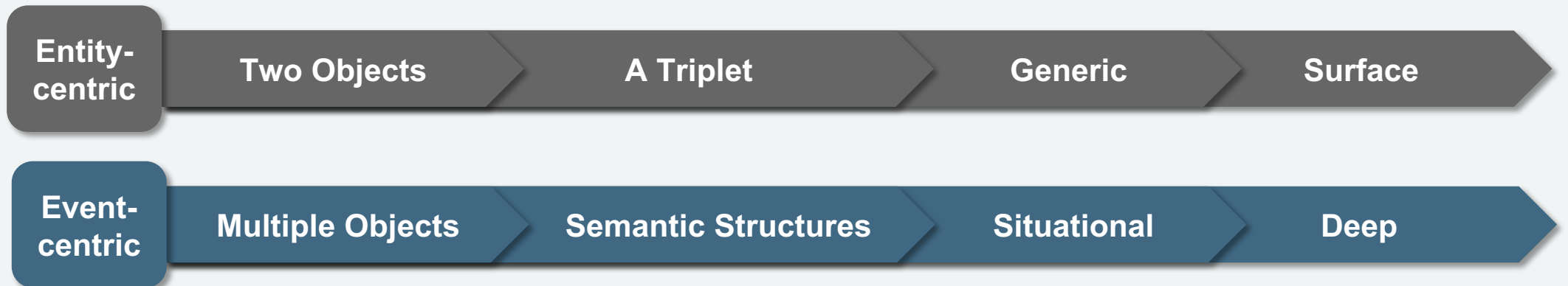
Event Extraction

Input: text, image, video, speech, ...

Output: structured knowledge

1. **Event Type** (e.g., *protest*)

2. **Participants** (e.g., *child*) & **Semantic Roles** (e.g., *protester*)



Goal: Entity-Centric → Event-Centric



	Text	Vision
	Entity	Object
Entity-centric	Relation	Relation
	Entity-Relation Graph	Scene Graph

Goal: Entity-Centric → Event-Centric



Text	Vision
Entity	Object
Relation	Relation
Entity-Relation Graph	Scene Graph
Verb	Activity
Event Structure	Image Event Graph

Event-centric

Goal: Entity-Centric → Event-Centric



	Text	Vision
	Entity	Object
Entity-centric	Relation	Relation
	Entity-Relation Graph	Scene Graph
	Verb	Activity
Event-centric	Event Structure	Image Event Graph

Goal: Entity-Centric → Event-Centric



State-of-the-art Captioner (Kamath et al., 2022)

Answer ▾

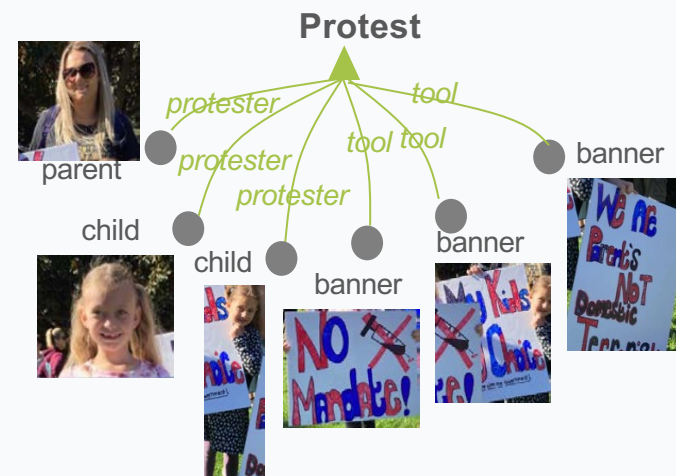
a woman holding a sign in front of a group of people.

a woman holding a sign while standing in a park.

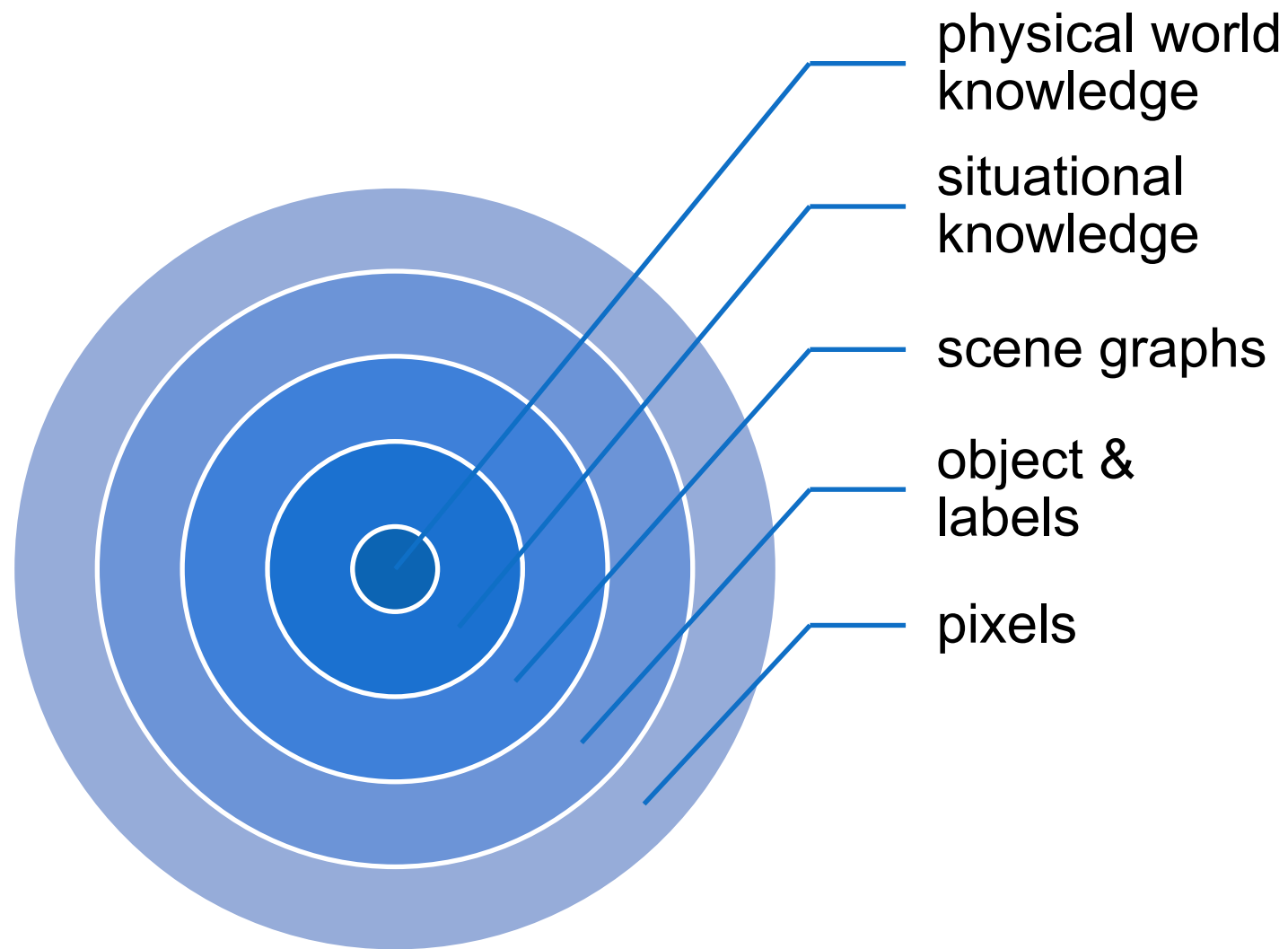
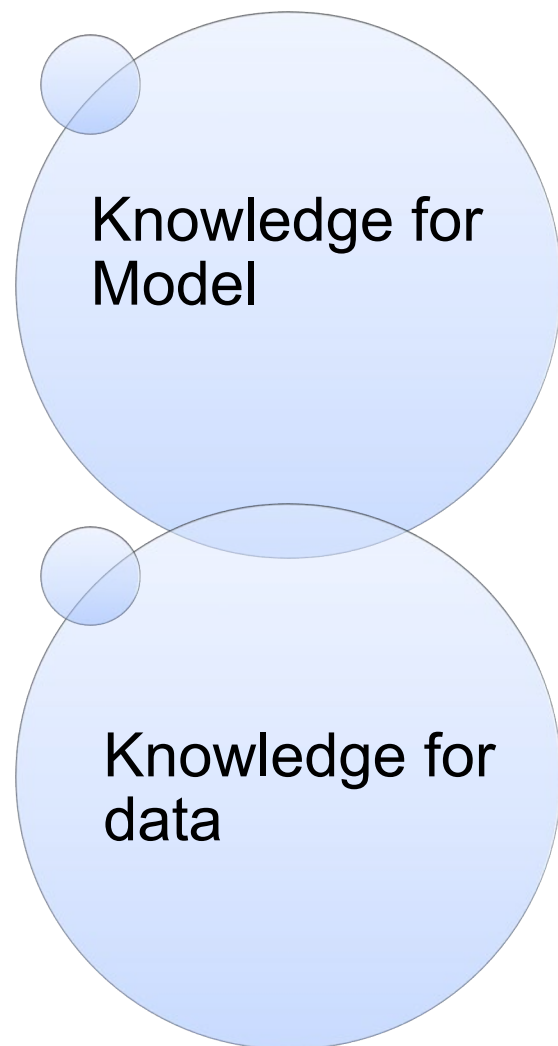
a woman holding a sign in front of a crowd.

Entity-centric

Event-centric



Adding knowledge to pretraining models



What is factual knowledge?

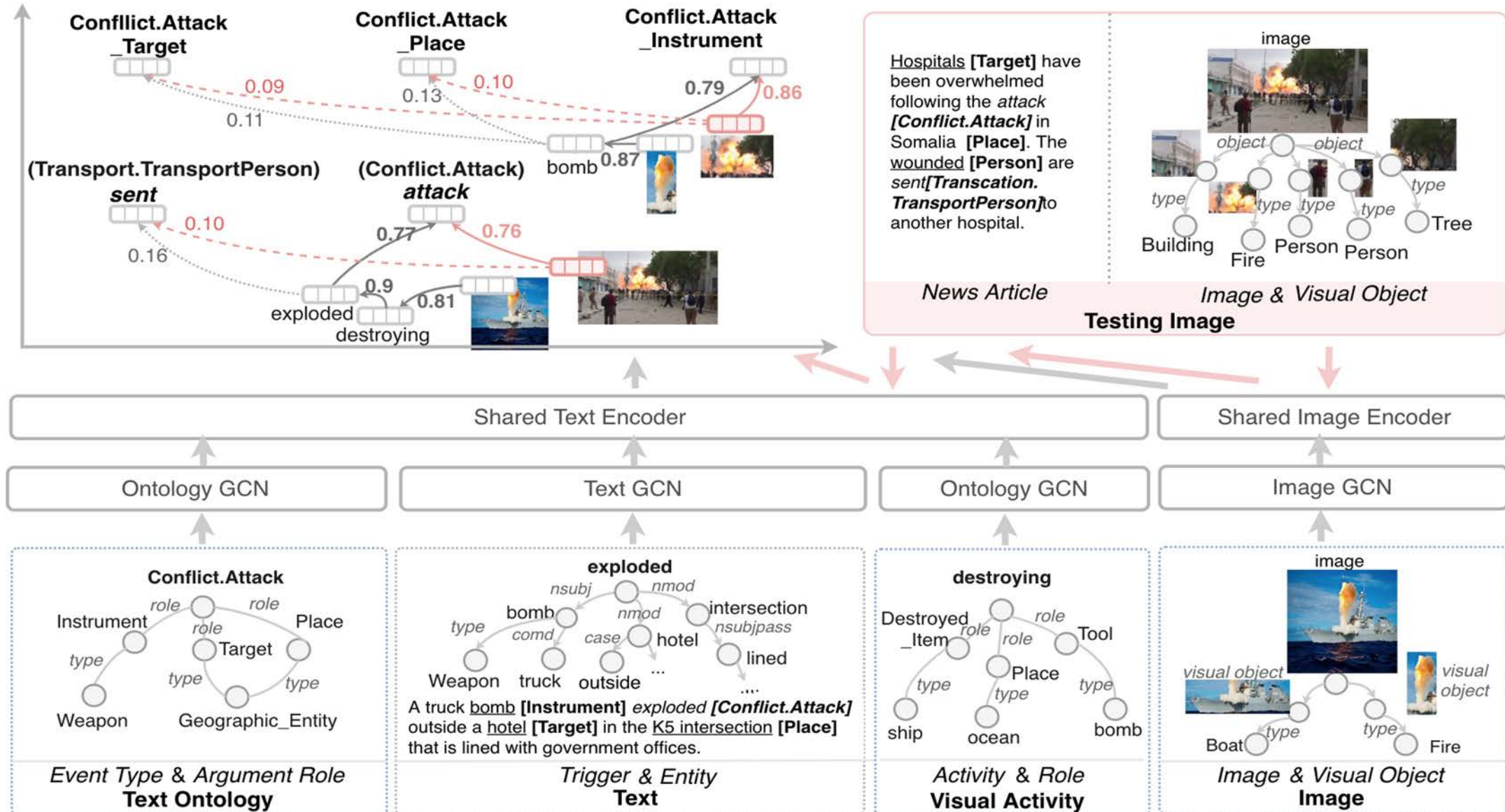
- Multimedia Knowledge Base with entities, relations and events.



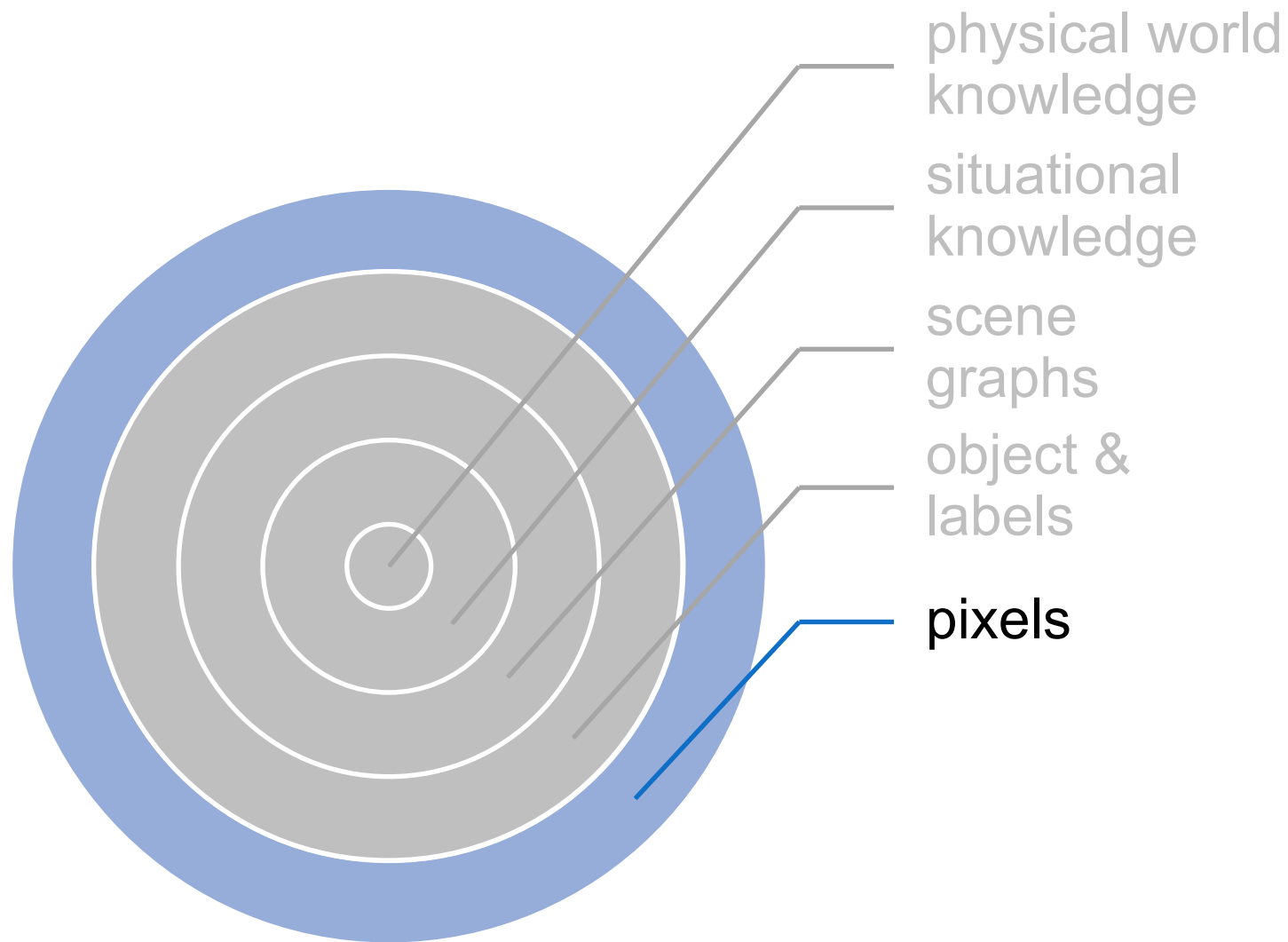
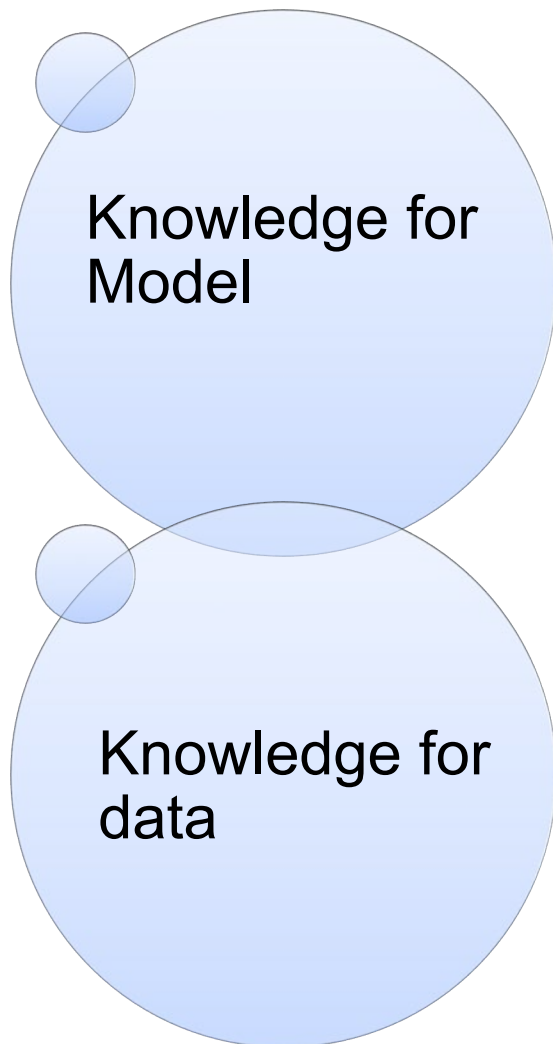
The first-ever official **visit** by a British royal to **Israel** is underway. **Prince William** the 36 year-old Duke of **Cambridge** and second in line to the throne will **meet** with both **Israeli** and **Palestinian** leaders over the next three days.

Contact.Meet_Participant

Goal: A joint representation of text and vision knowledge

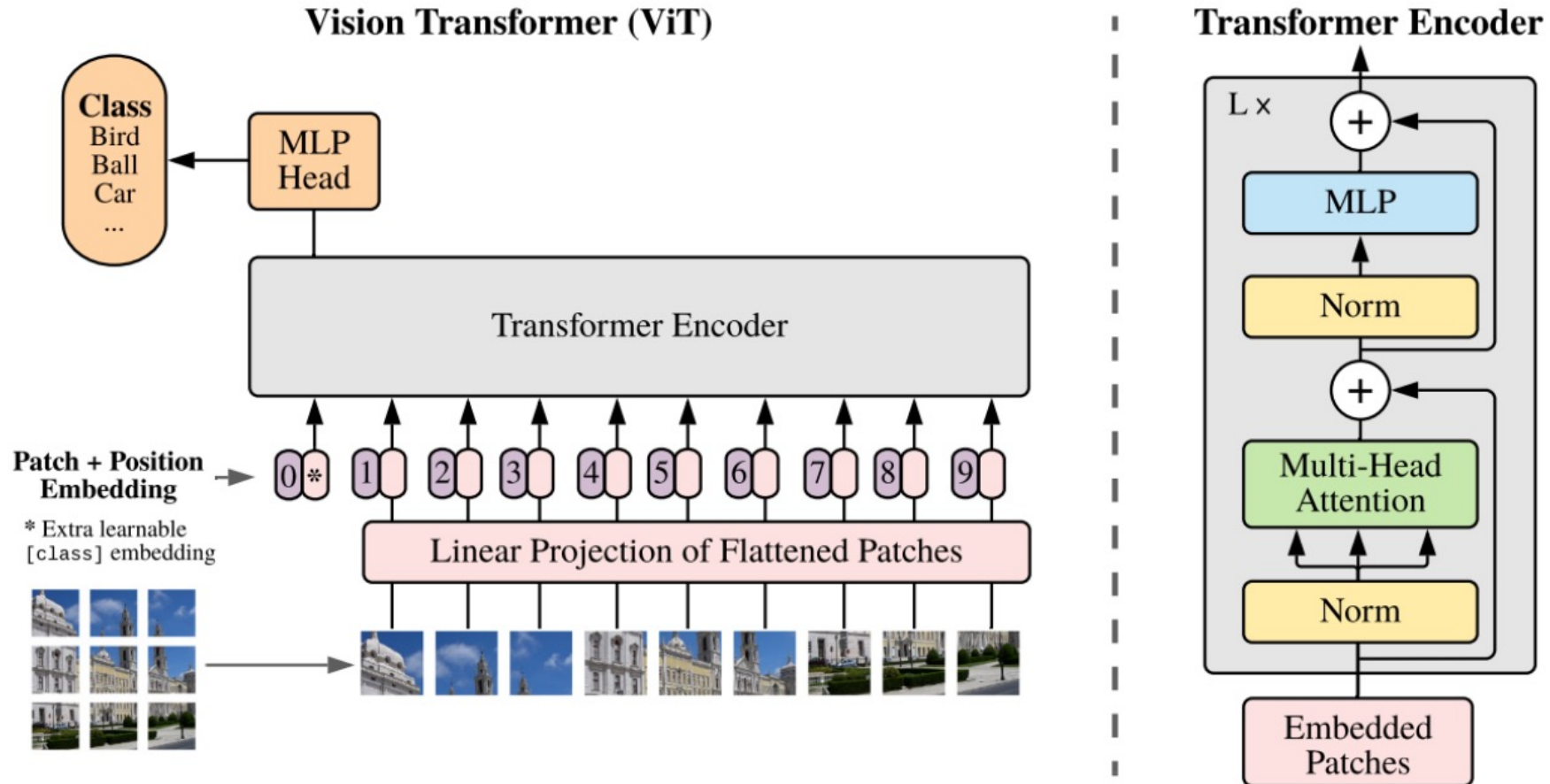


Adding knowledge to pretraining models



Pixels - An Image is Worth 16x16 Words

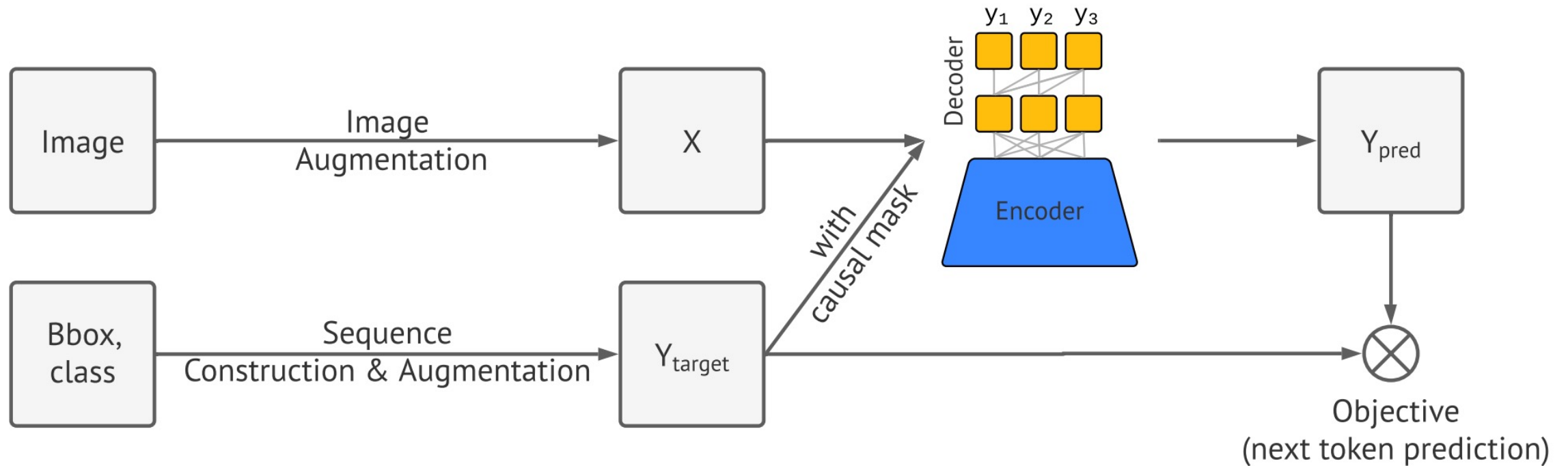
The simplest way is to split an image into patches



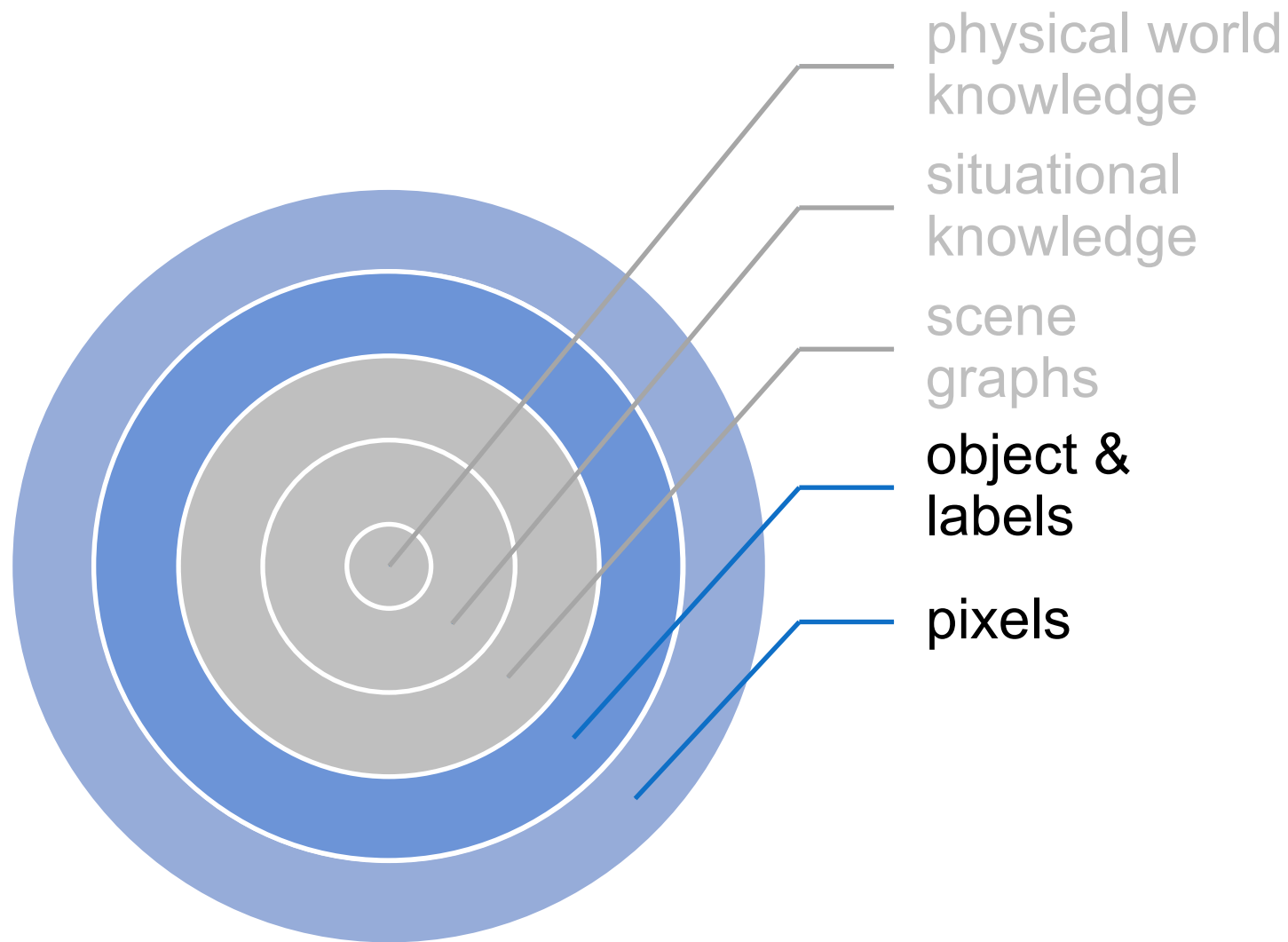
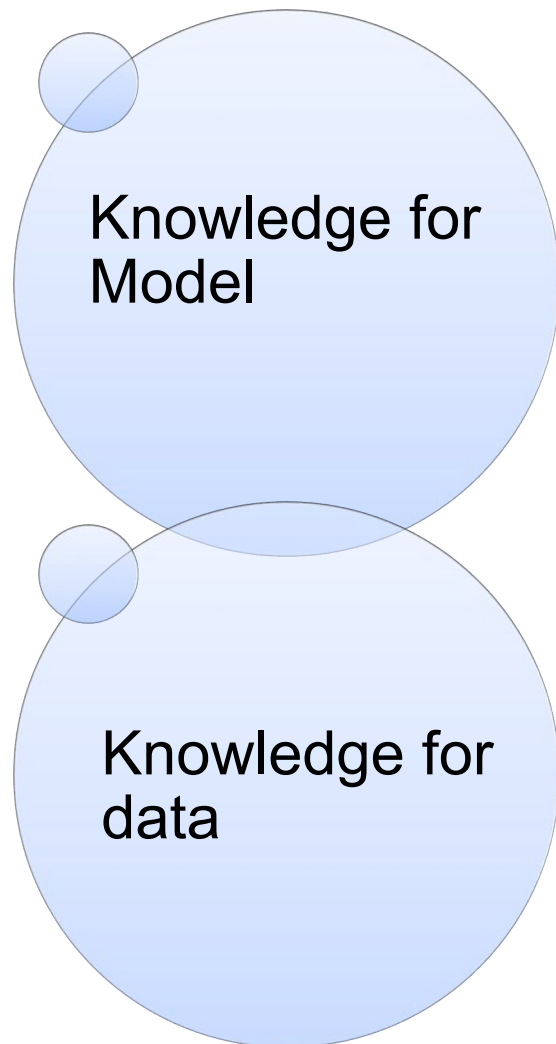
Pixels - Unified Model: Pix2Seq



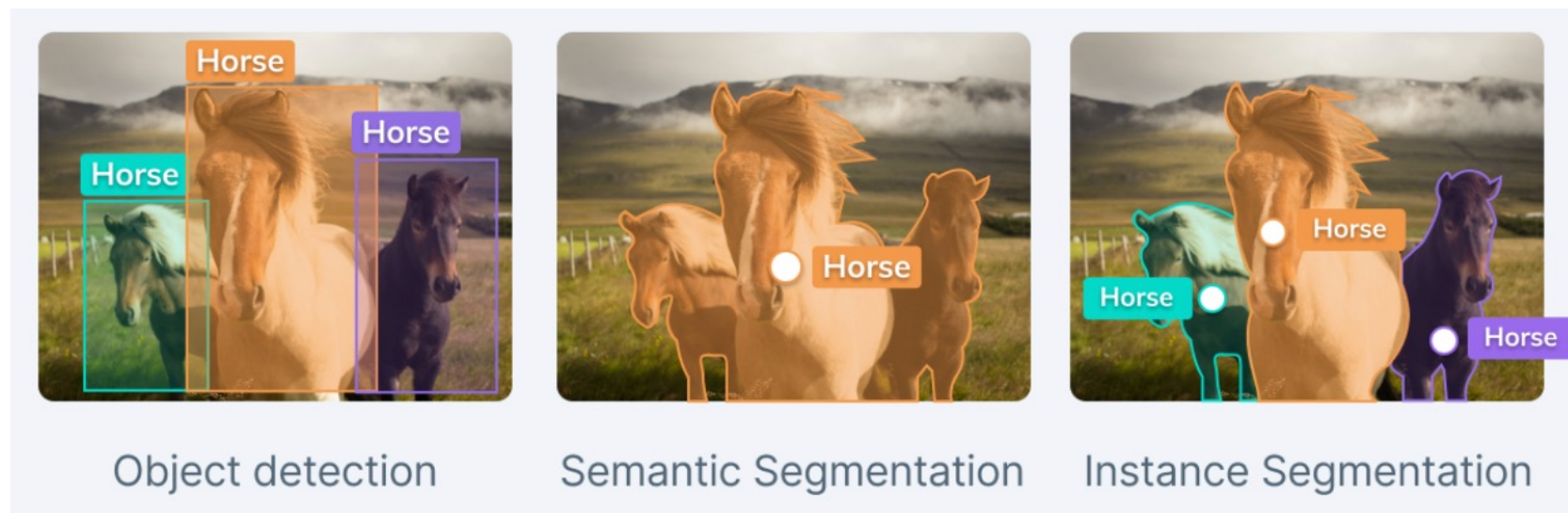
Another way is to treat pixels as tokens.



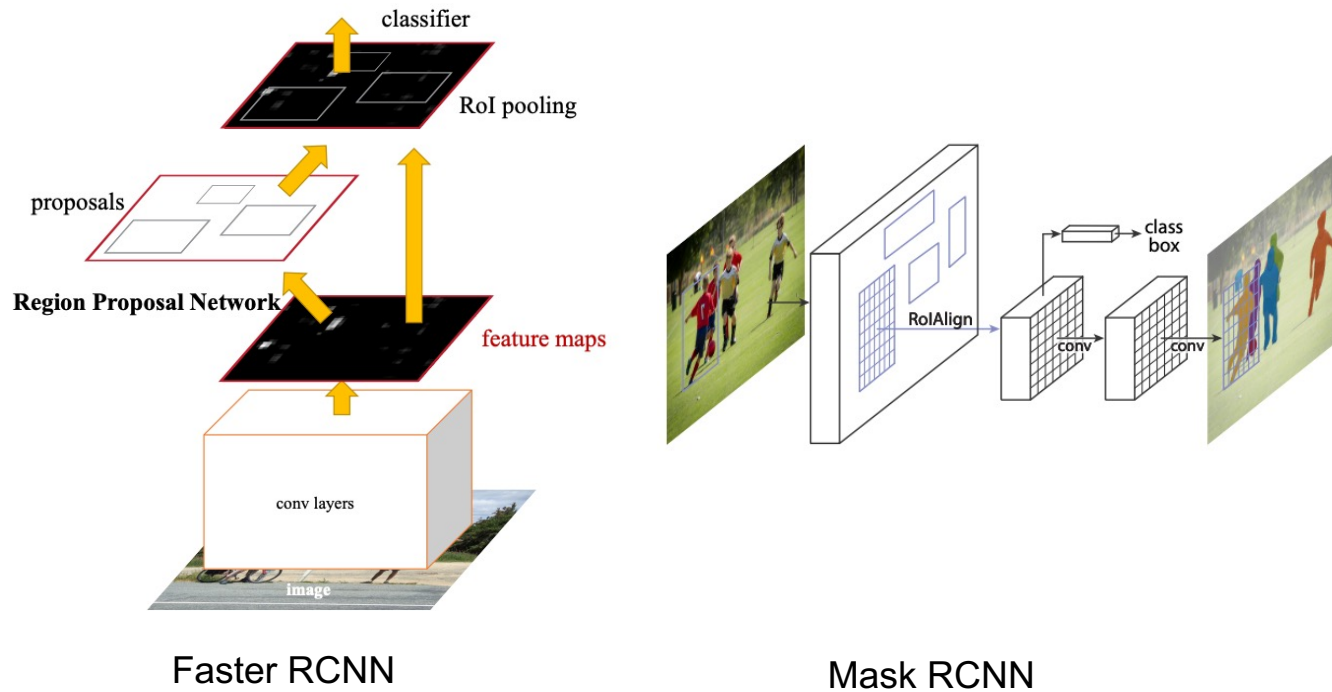
Adding knowledge to pretraining models



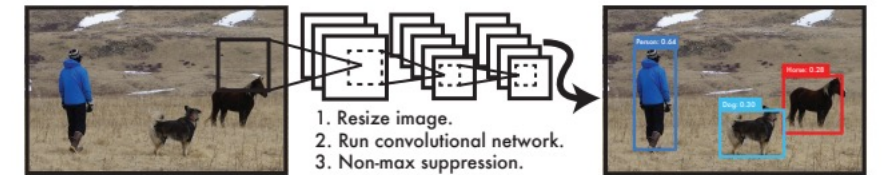
- Object Detection: Object instances at the bounding box level
- Semantic Segmentation: Object class at the pixel level
- Instance Segmentation: Object instances at the pixel level



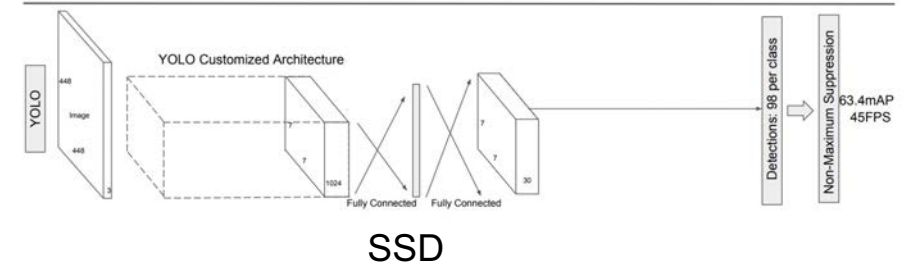
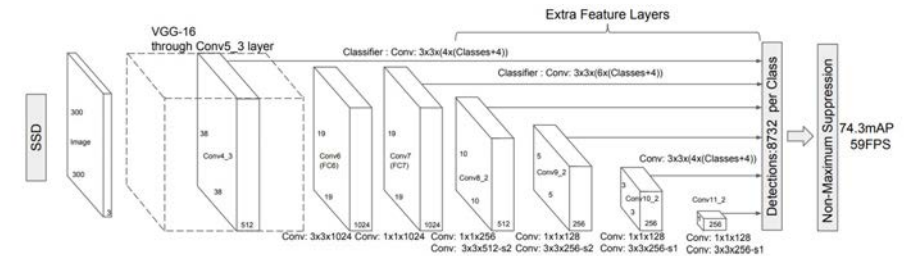
Two-Stage (With Proposal)



One-Stage (Without Proposal)



YOLO

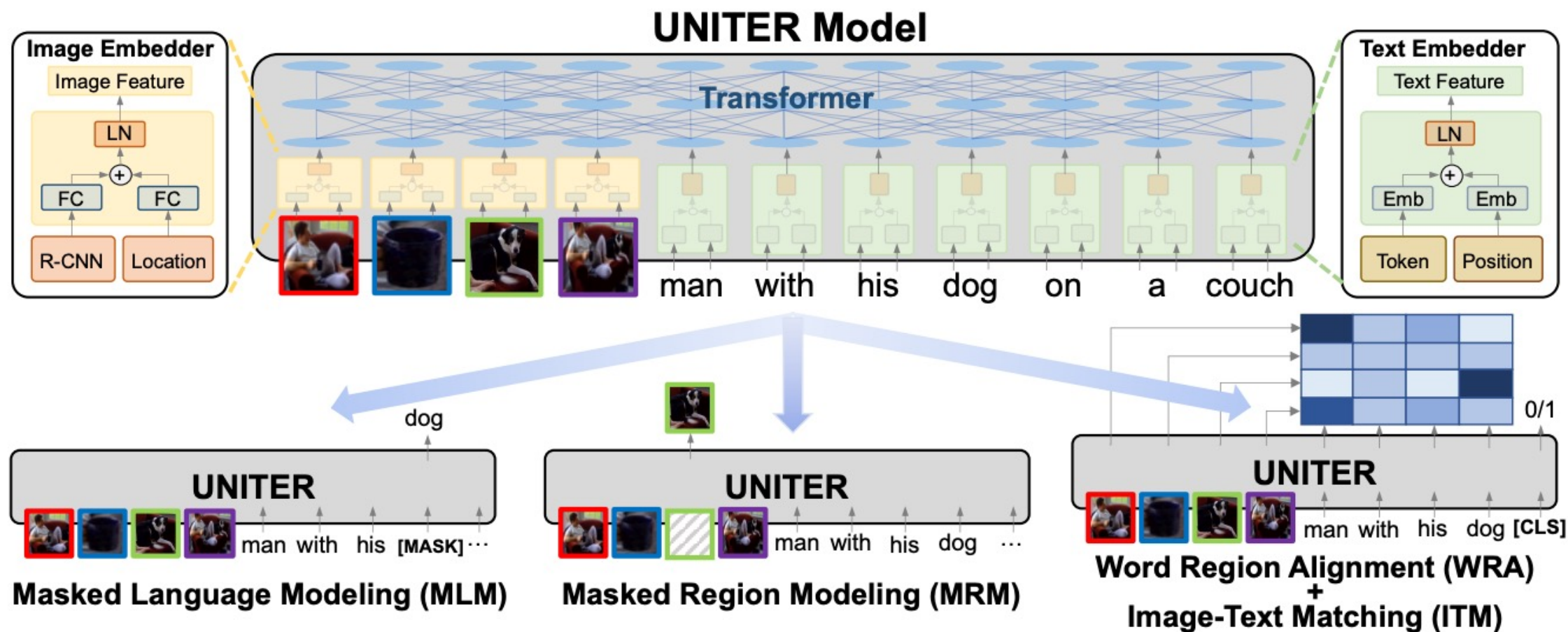


Ren, S., He, K., Girshick, R., & Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS 2015*.
 He, Kaiming, et al. "Mask r-cnn." *CVPR 2017*.
 Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." *CVPR 2016*.
 Liu, Wei, et al. "Ssd: Single shot multibox detector." *ECCV 2016*.

Adding objects to V+L Pretraining



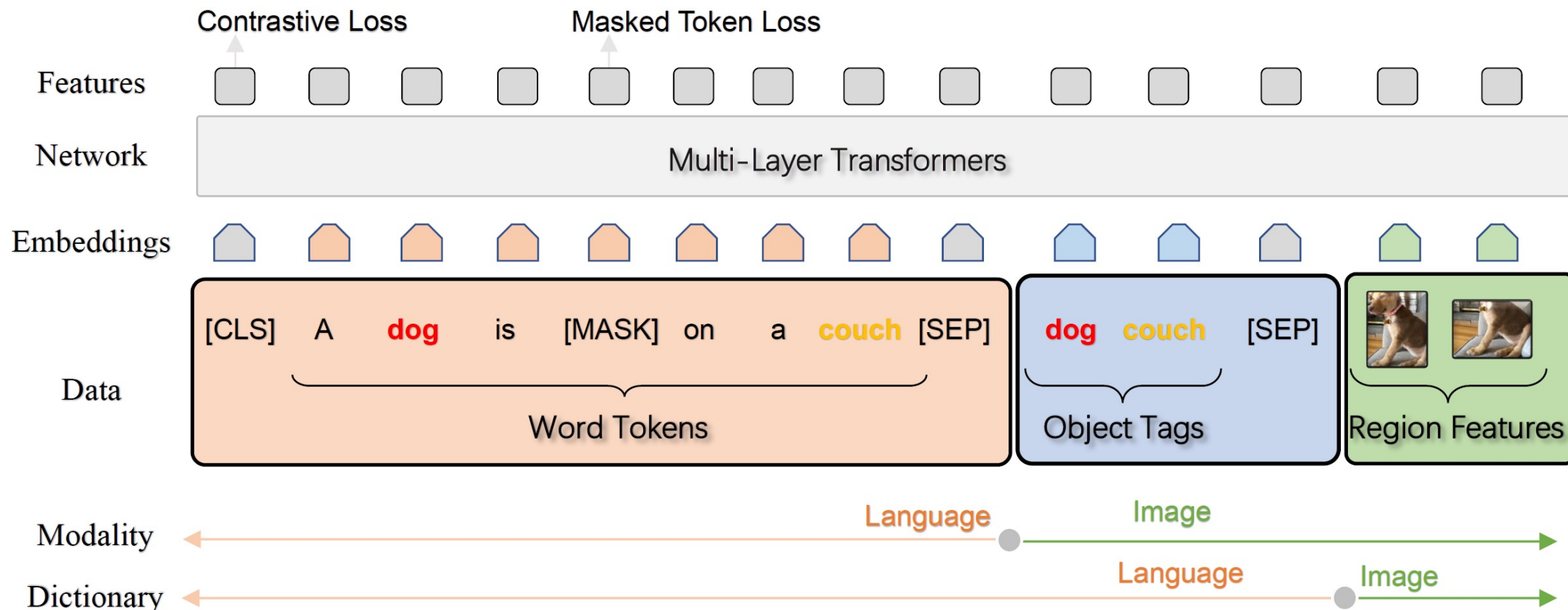
Objects are used to better mask the regions.



Oscar [ECCV 2020] and VinVL [CVPR 2021]



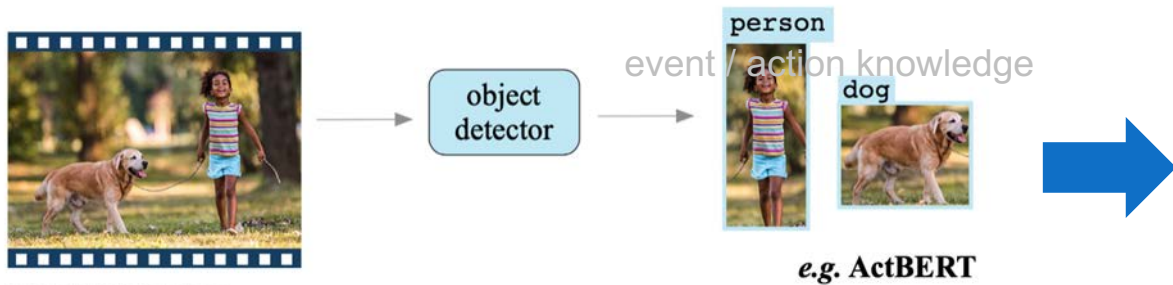
- Object knowledge is richer.
 - Add object label knowledge as anchor points



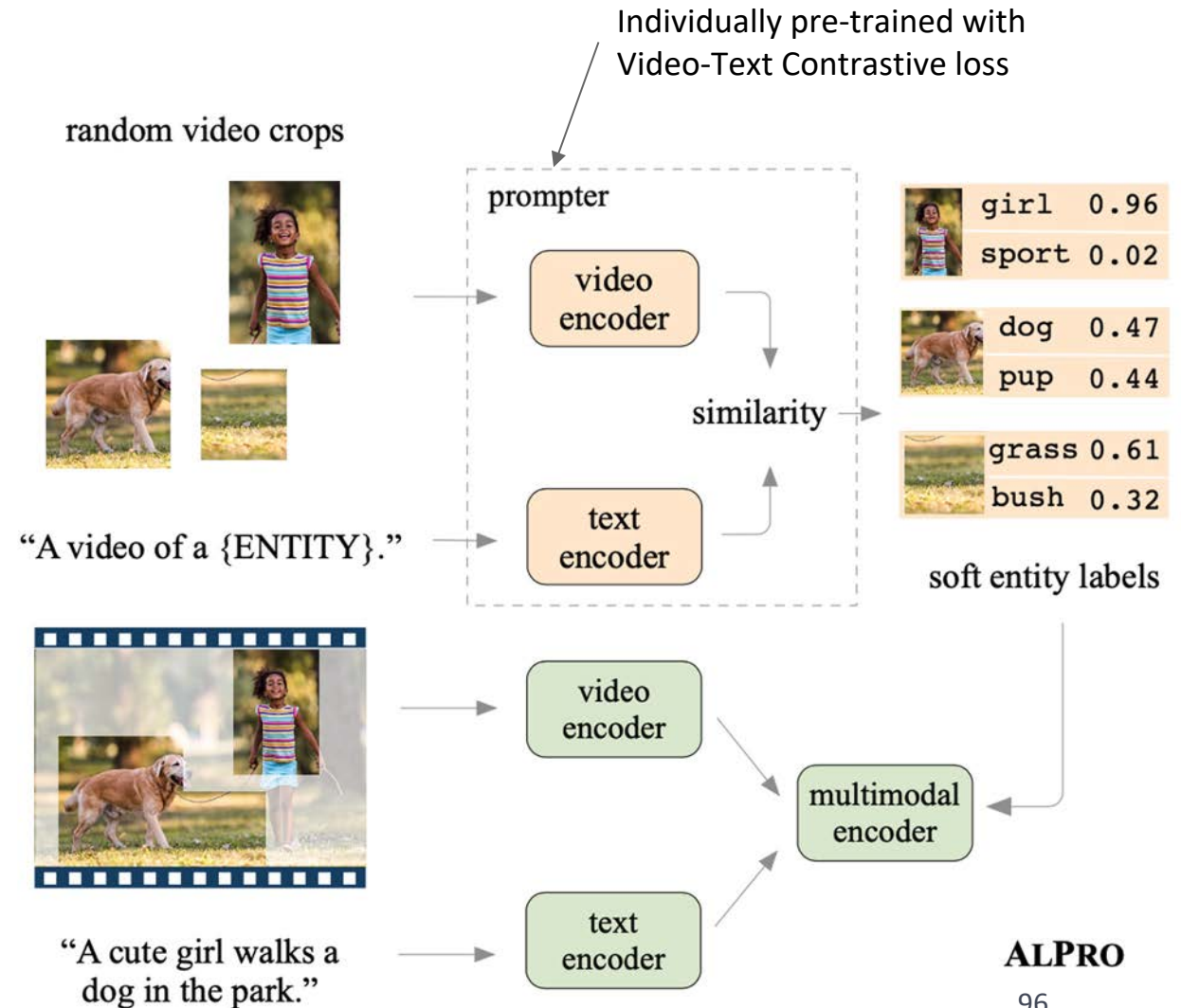
Soft Prompt Entity Knowledge [CVPR2022]



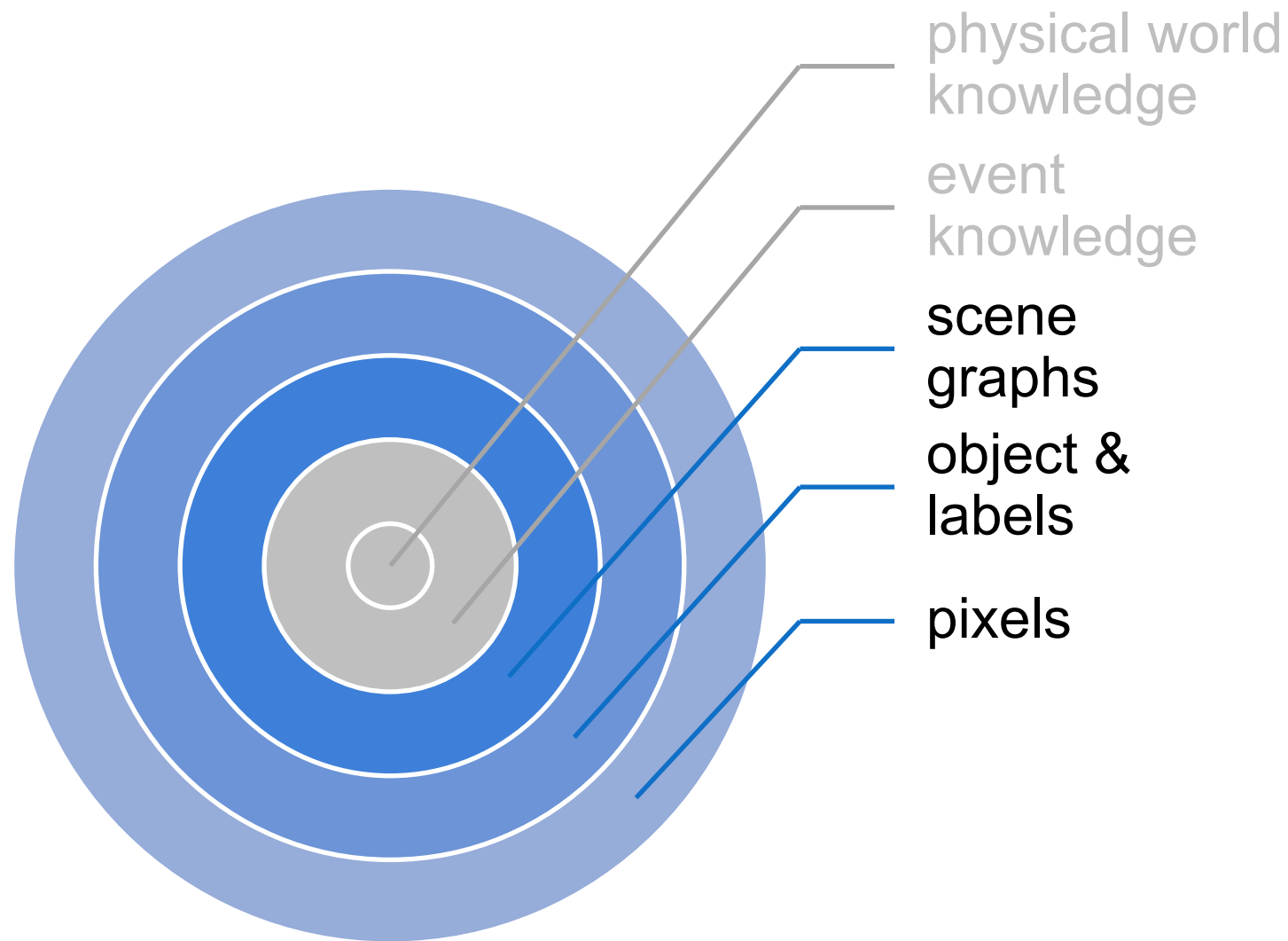
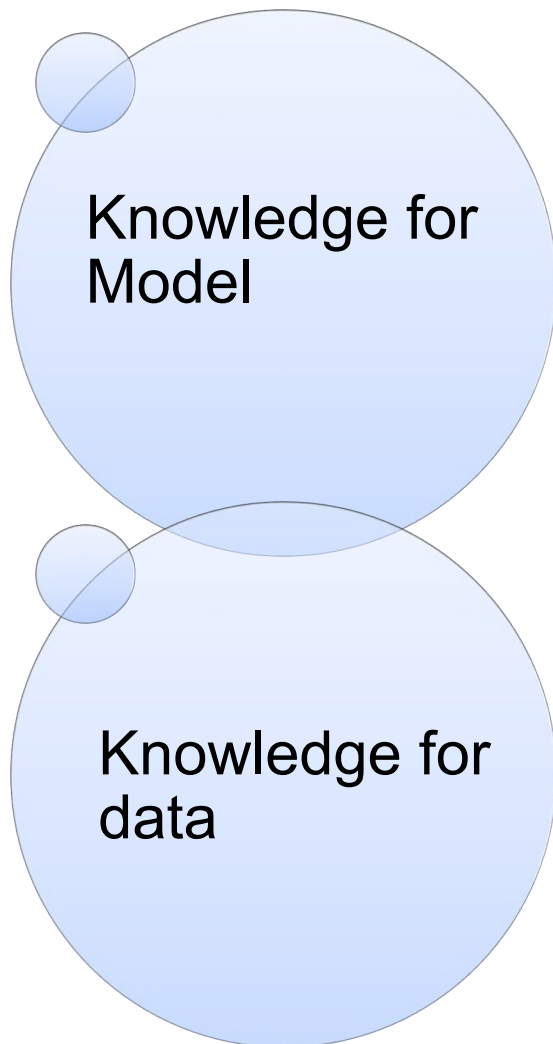
- **[Align and Prompt 2021]** Align and Prompt: Video-and-Language Pre-training with Entity Prompts
 - Adding regional entity prediction task



previous work rely on object detectors with expensive computation and limited object categories



Adding knowledge to pretraining models



- Add scene graph knowledge as downstream tasks
 - Object prediction
 - Attribute prediction
 - Relationship prediction

(a) Objects



A tan **dog** and a little girl kiss.



The little girl is kissing the brown **cat**.

(b) Attributes



A black dog playing with a **purple** toy.



A black dog playing with a **green** toy.

(c) Relationships

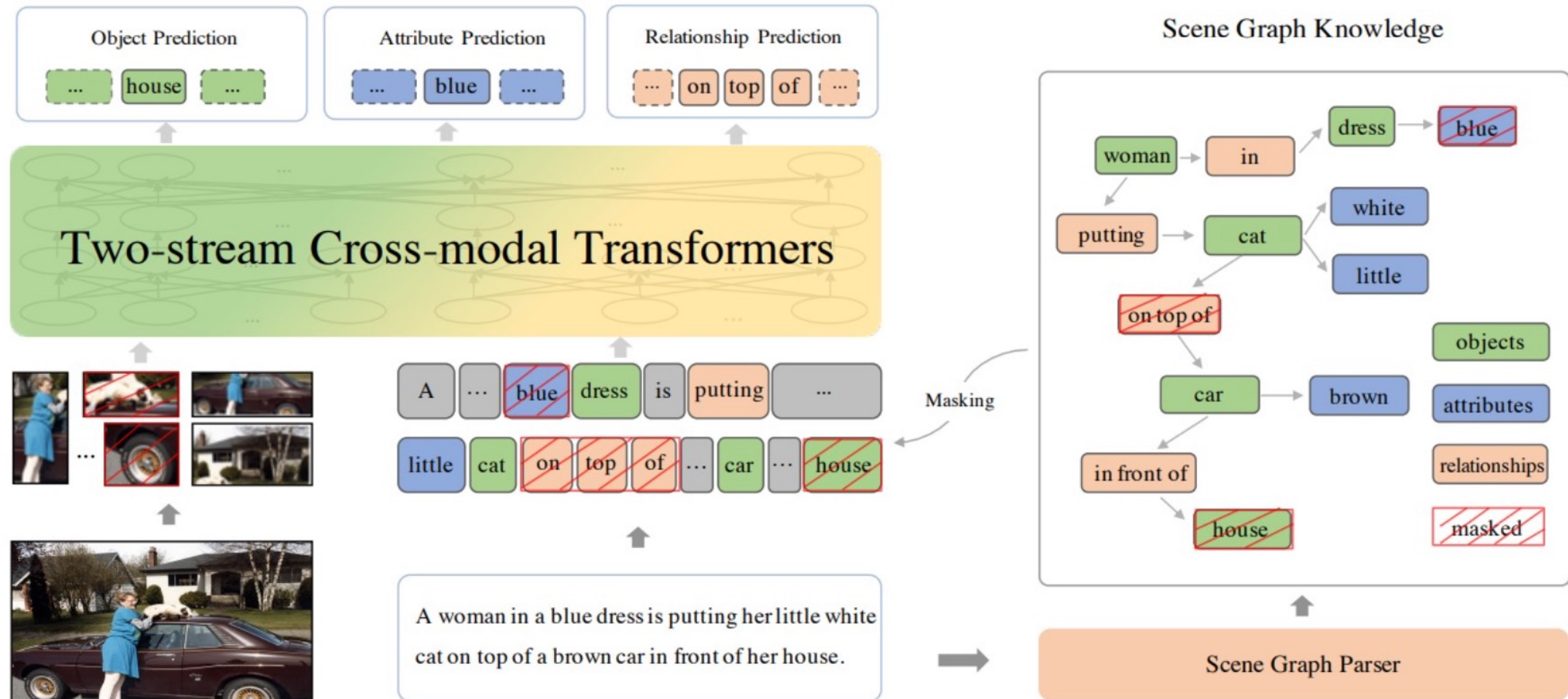


A man in red plaid **rides** his bike in a park.

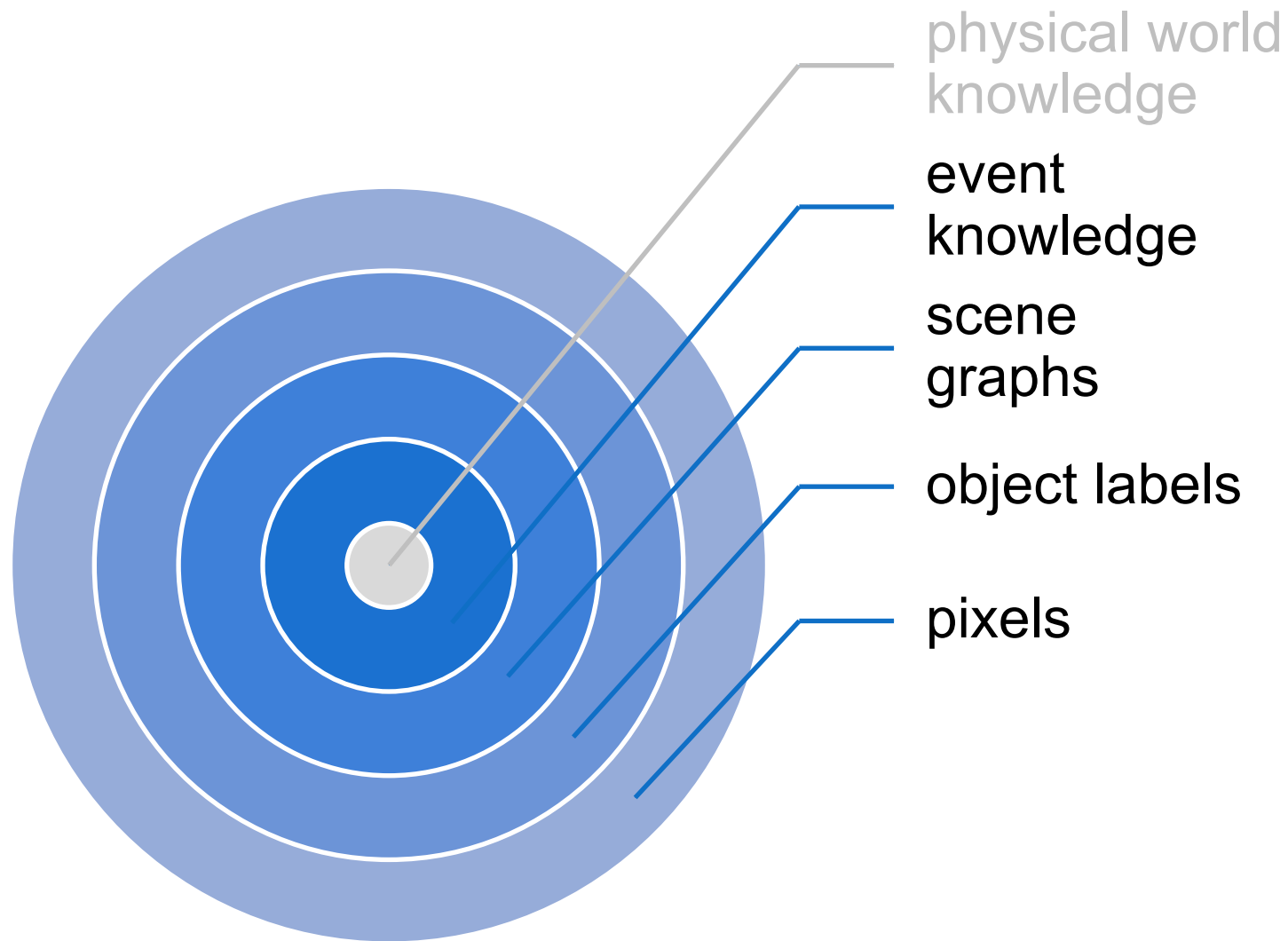
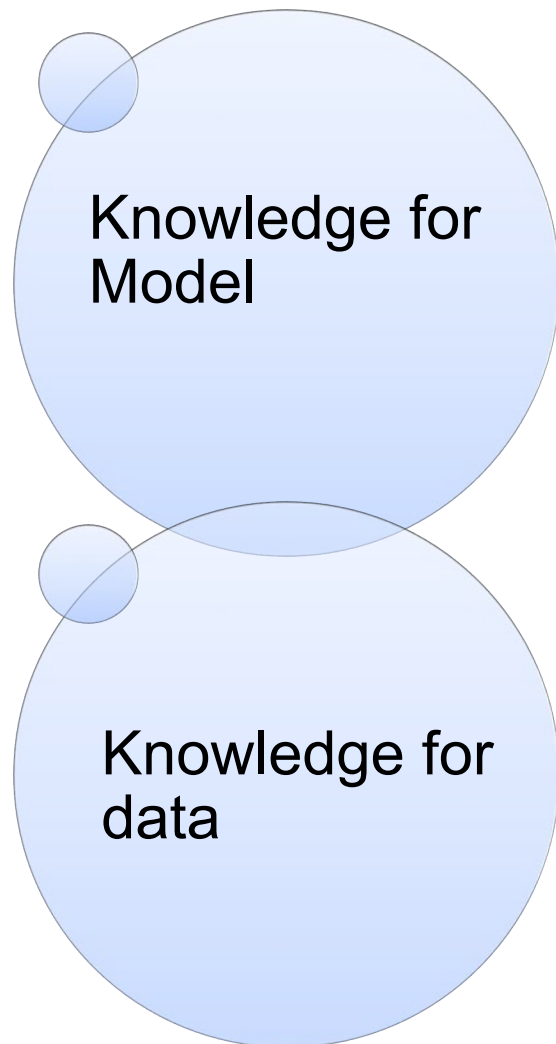


An older man **repairing** a bike tire in a park.

- Add scene graph knowledge as downstream tasks



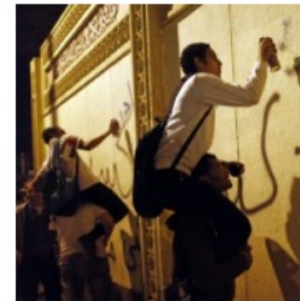
Adding knowledge to pretraining models



Vision vs. NLP for Event Extraction



- Vision does not study newsworthy, complex events
 - Focusing on daily life and sports (Perera et al., 2012; Chang et al., 2016; Zhang et al., 2007; Ma et al., 2017)
 - Without localizing a complete set of arguments for each event (Gu et al., 2018; Li et al., 2018; Duarte et al., 2018; Sigurdsson et al., 2016; Kato et al., 2018; Wu et al., 2019a)
- Most related: Situation Recognition (Yatskar et al., 2016)
 - Classify an image as one of 500+ FrameNet verbs
 - Identify 192 generic semantic roles via a 1-word description



CLIPPING	
ROLE	VALUE
AGENT	MAN
SOURCE	SHEEP
TOOL	SHEARS
ITEM	WOOL
PLACE	FIELD

CLIPPING	
ROLE	VALUE
AGENT	VET
SOURCE	DOG
TOOL	CLIPPER
ITEM	CLAW
PLACE	ROOM

JUMPING	
ROLE	VALUE
AGENT	BOY
SOURCE	CLIFF
OBSTACLE	-
DESTINATION	WATER
PLACE	LAKE

JUMPING	
ROLE	VALUE
AGENT	BEAR
SOURCE	ICEBERG
OBSTACLE	WATER
DESTINATION	ICEBERG
PLACE	OUTDOOR

SPRAYING	
ROLE	VALUE
AGENT	MAN
SOURCE	SPRAY CAN
SUBSTANCE	PAINT
DESTINATION	WALL
PLACE	ALLEYWAY

SPRAYING	
ROLE	VALUE
AGENT	FIREMAN
SOURCE	HOSE
SUBSTANCE	WATER
DESTINATION	FIRE
PLACE	OUTSIDE

Vision-only Event and Argument Extraction



- Grounded Situation Recognition adds visual argument localization [Pratt et al, 2020]

- Video Situation Recognition extends the work to videos [Sadhu et al, 2021]

Hitting	Catching
Agent: Ballplayer, Tool: Bat, Victim: Baseball, Victim Part: Ø, Place: Field	Agent: Bear, Caught Item: Fish, Tool: Mouth, Place: River
Jumping	Kneading
Agent: Female Child, Source: Sofa, Destination: Sofa, Obstacle: Ø, Place: Living Room	Agent: Person, Item: Dough, Place: Kitchen

← 2 Seconds →

Event 1 0s-2s		Verb: deflect (block, avoid) Arg0 (deflector): woman with shield Arg1 (thing deflected): boulder Scene: city park
Event 2 2s-4s		Verb: talk (speak) Arg0 (talker): woman with shield Arg2 (hearer): man with trident ArgM (manner): urgently Scene: city park
Event 3 4s-6s		Verb: leap (physically leap) Arg0 (jumper): man with trident Arg1 (obstacle): over stairs ArgM (direction): towards shirtless man ArgM (goal): to attack shirtless man Scene: city park
Event 4 6s-8s		Verb: punch (to hit) Arg0 (agent): shirtless man Arg1 (entity punched): man with trident ArgM (direction): far into distance Scene: city park
Event 5 8s-10s		Verb: punch (to hit) Arg0 (agent): shirtless man Arg1 (entity punched): woman with shield ArgM (direction): down the stairs Scene: city park

Ev3 is enabled by Ev1

Ev3 is a reaction to Ev2

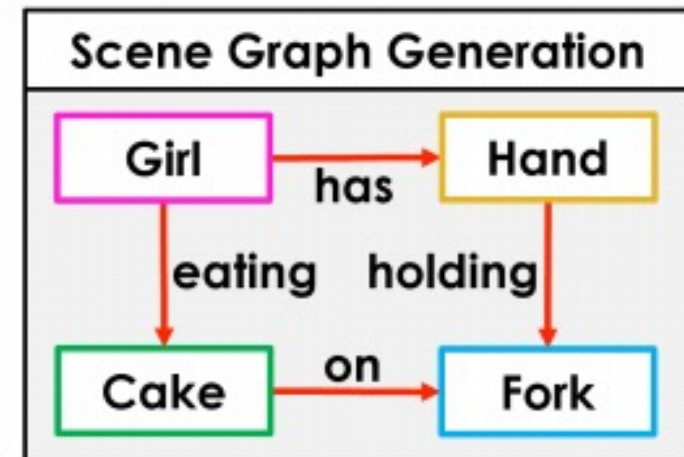
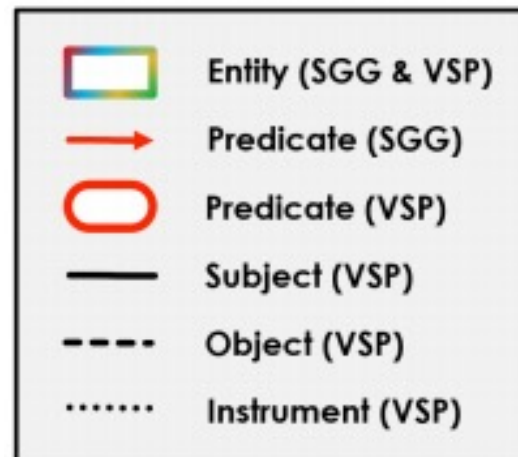
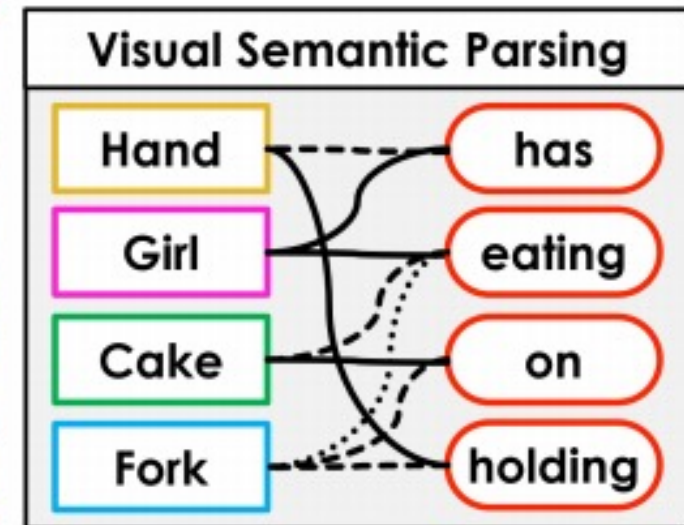
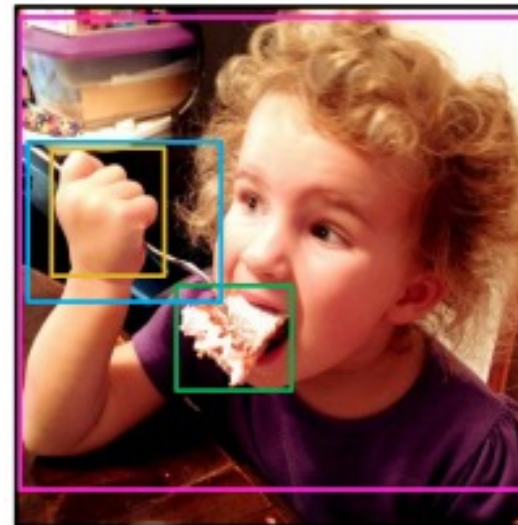
Ev4 is a reaction to Ev3

Ev5 is unrelated to Ev3

Vision-only Event and Argument Extraction



- Another line of work is based on scene graphs [Xu et al, 2017; Li et al, 2017; Yang et al, 2018; Zellers et al, 2018].
 - extracting <subject, predicate, object>
 - structure is simpler than the aforementioned multi-argument event
- Visual Semantic Parsing is using predicate as event, and subject, object, instrument as argument [Zareian et al, 2020]
 - Added bounding box grounding



Existing Work: Situation Recognition



SPRAYING

ROLE	VALUE
AGENT	MAN
SOURCE	SPRAY CAN
SUBSTANCE	PAINT
DESTINATION	WALL
PLACE	ALLEYWAY

SPRAYING

ROLE	VALUE
AGENT	FIREMAN
SOURCE	HOSE
SUBSTANCE	WATER
DESTINATION	FIRE
PLACE	OUTSIDE

CLIPPING

ROLE	VALUE
AGENT	MAN
SOURCE	SHEEP
TOOL	SHEARS
ITEM	WOOL
PLACE	FIELD

CLIPPING

ROLE	VALUE
AGENT	VET
SOURCE	DOG
TOOL	CLIPPER
ITEM	CLAW
PLACE	ROOM

(Yatskar et al., 2016, ...)

Supervised Learning

**Bottleneck:
Lack of Annotation**

Vision-Only

**Bottleneck:
Cross-modal Fusion**



Surfing

Agent	Tool	Path	Place
Man	Surfboard	Water	Ocean



Feeding

Agent	Food	Eater	Source	Place
Person	Milk	Tiger	Bottle	∅

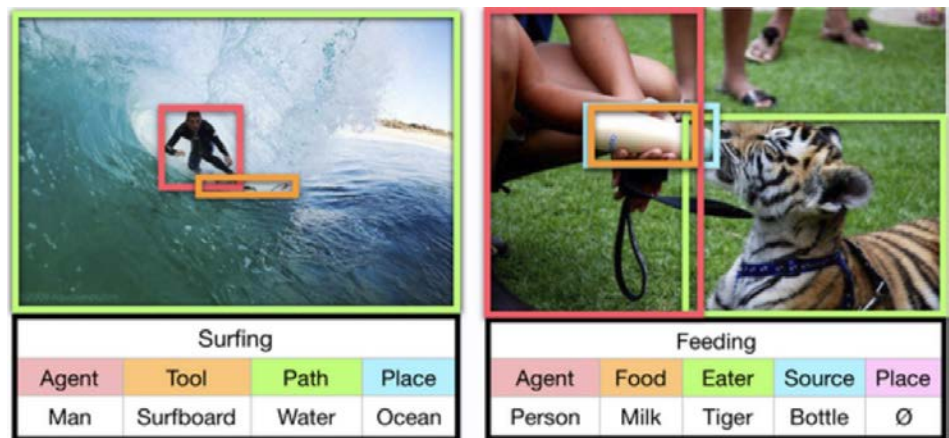
(Pratt et al., 2020, ...)

Existing Work: Situation Recognition



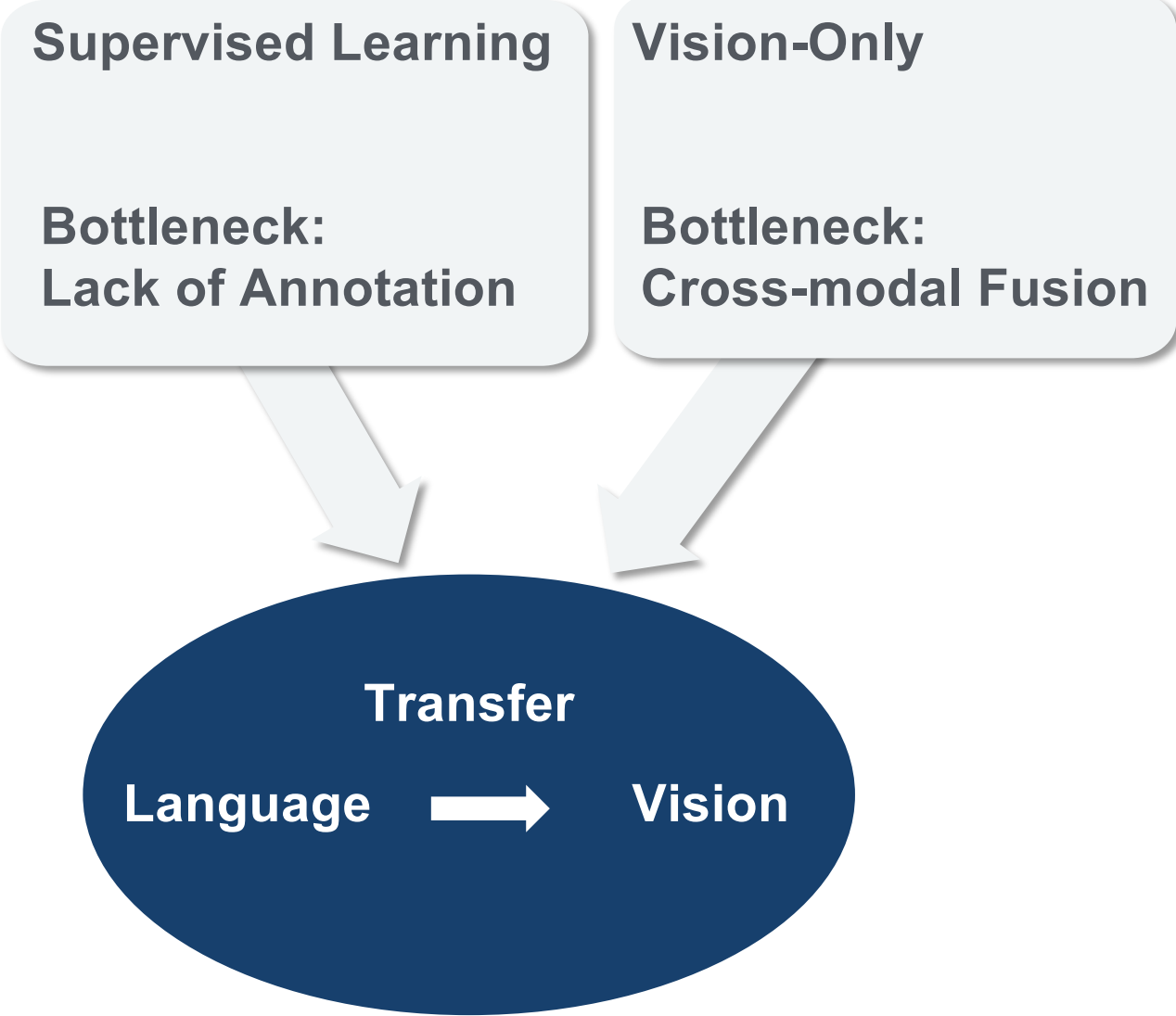
SPRAYING				CLIPPING			
ROLE	VALUE	ROLE	VALUE	ROLE	VALUE	ROLE	VALUE
AGENT	MAN	AGENT	FIREMAN	AGENT	MAN	AGENT	VET
SOURCE	SPRAY CAN	SOURCE	HOSE	SOURCE	SHEEP	SOURCE	DOG
SUBSTANCE	PAINT	SUBSTANCE	WATER	TOOL	SHEARS	TOOL	CLIPPER
DESTINATION	WALL	DESTINATION	FIRE	ITEM	WOOL	ITEM	CLAW
PLACE	ALLEYWAY	PLACE	OUTSIDE	PLACE	FIELD	PLACE	ROOM

(Yatskar et al., 2016, ...)



Surfing				Feeding				
Agent	Tool	Path	Place	Agent	Food	Eater	Source	Place
Man	Surfboard	Water	Ocean	Person	Milk	Tiger	Bottle	∅

(Pratt et al., 2020, ...)



Existing Work: Situation Recognition



SPRAYING

ROLE	VALUE
AGENT	MAN
SOURCE	SPRAY CAN
SUBSTANCE	PAINT
DESTINATION	WALL
PLACE	ALLEYWAY

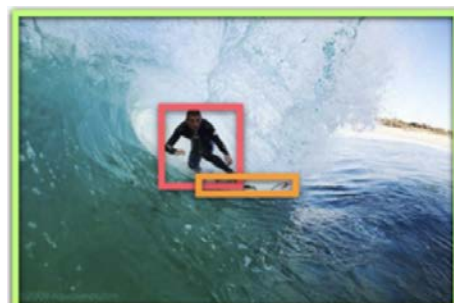
ROLE	VALUE
AGENT	FIREMAN
SOURCE	HOSE
SUBSTANCE	WATER
DESTINATION	FIRE
PLACE	OUTSIDE

CLIPPING

ROLE	VALUE
AGENT	MAN
SOURCE	SHEEP
TOOL	SHEARS
ITEM	WOOL
PLACE	FIELD

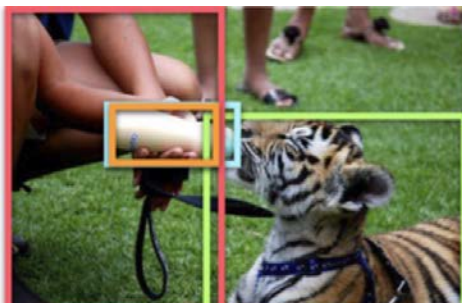
ROLE	VALUE
AGENT	VET
SOURCE	DOG
TOOL	CLIPPER
ITEM	CLAW
PLACE	ROOM

(Yatskar et al., 2016, ...)



Surfing

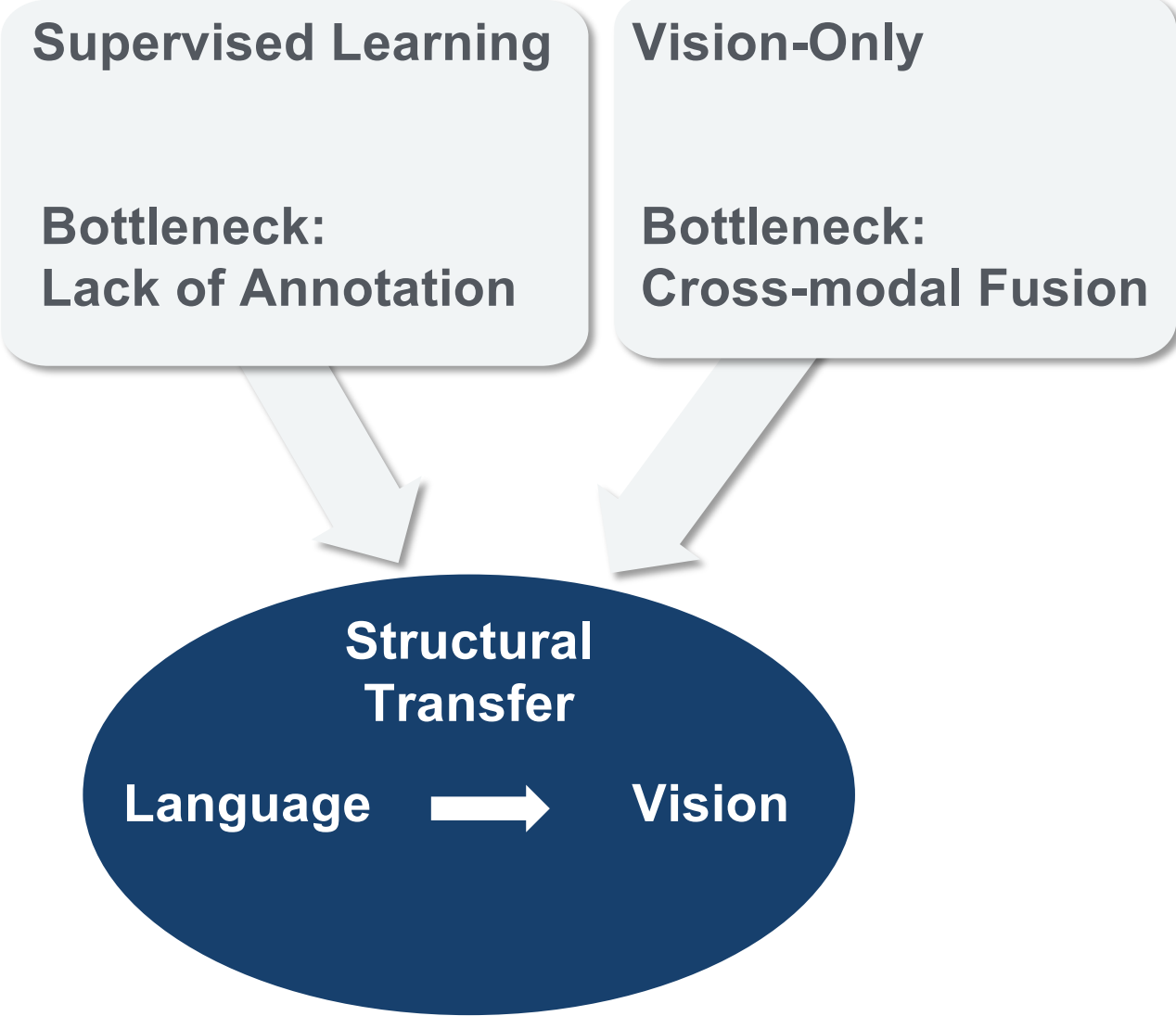
Agent	Tool	Path	Place
Man	Surfboard	Water	Ocean



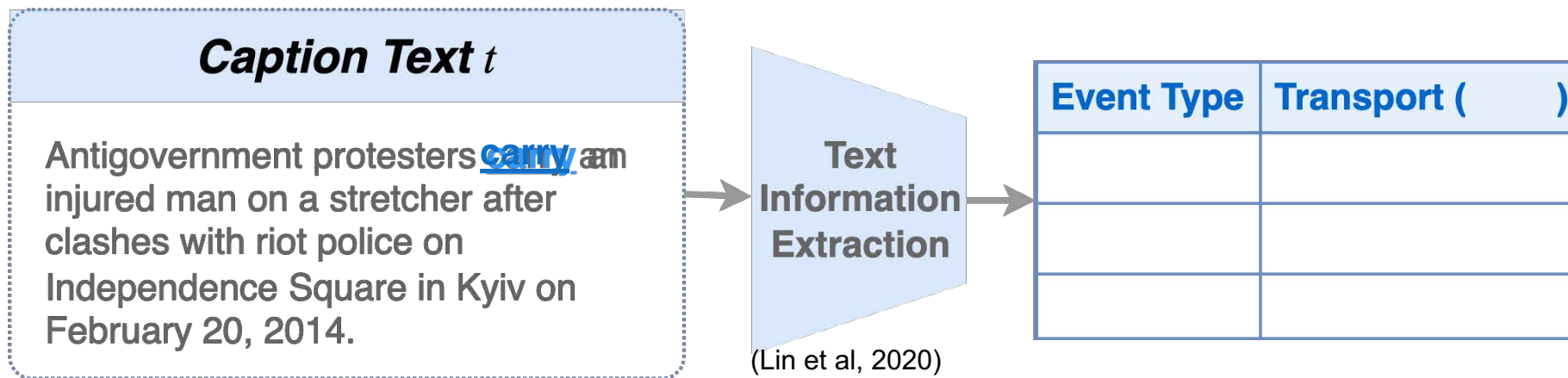
Feeding

Agent	Food	Eater	Source	Place
Person	Milk	Tiger	Bottle	∅

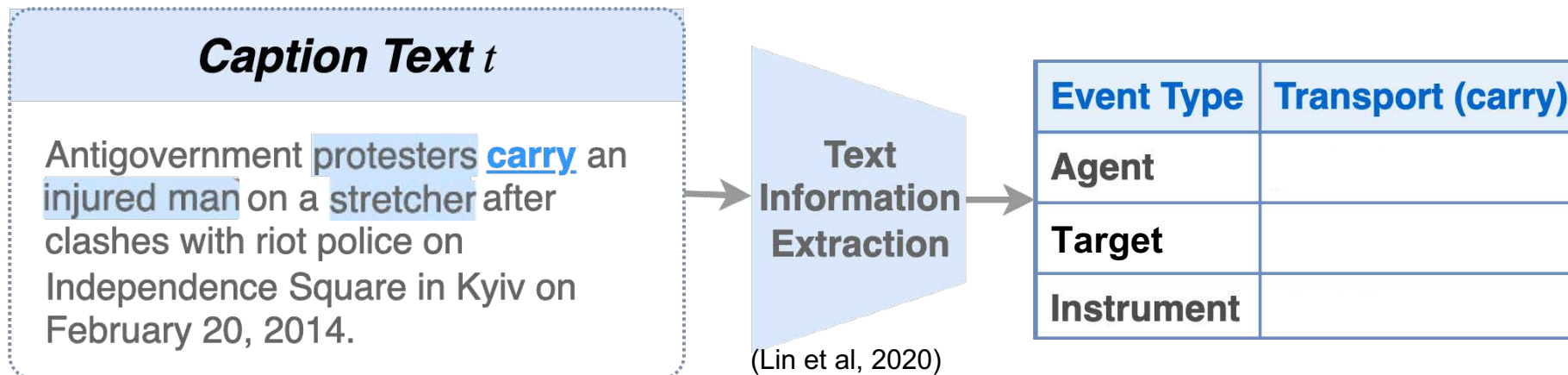
(Pratt et al., 2020, ...)



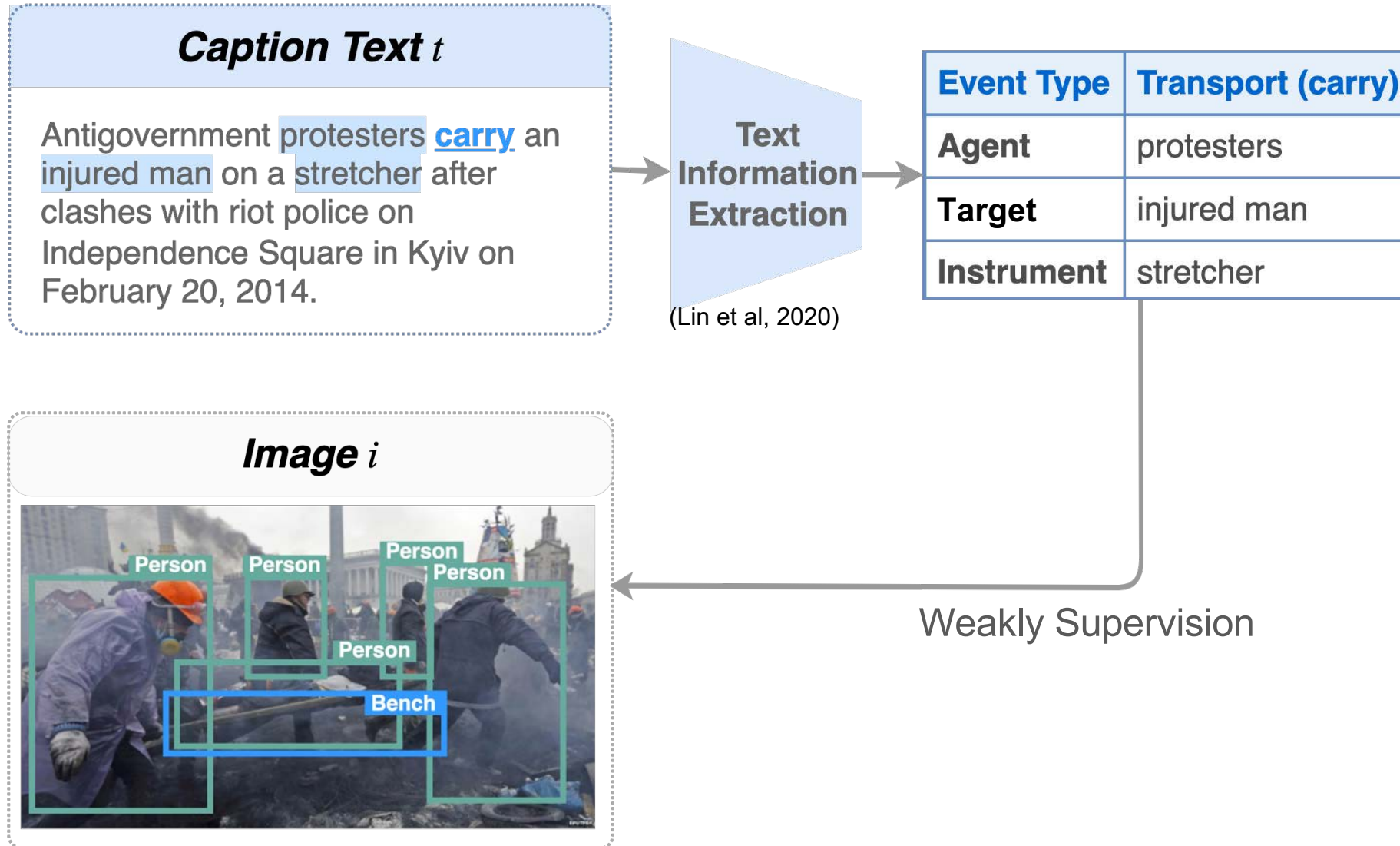
CLIP-Event: Event-Driven Vision-Language Pretraining



CLIP-Event: Event-Driven Vision-Language Pretraining



Transfer text event knowledge to images



Hard negatives via 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

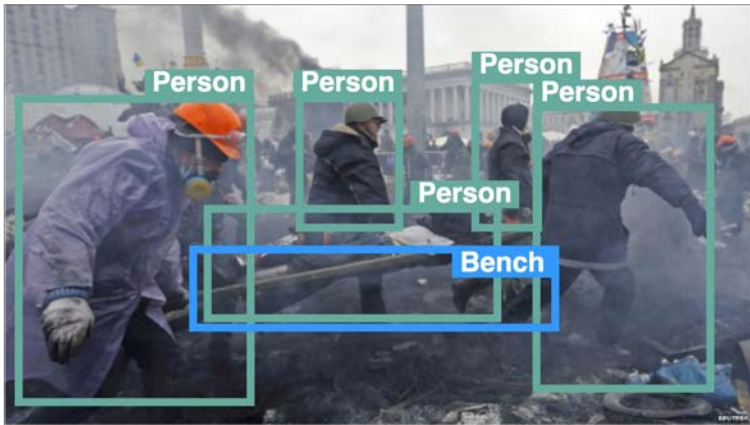
Event Type	Transport (carry)
Agent	protesters
Target	injured man
Instrument	stretcher

Negative Labels
(events)

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

Confusion Matrix of existing V+L models

Image i



Hard negatives via manipulating event structures

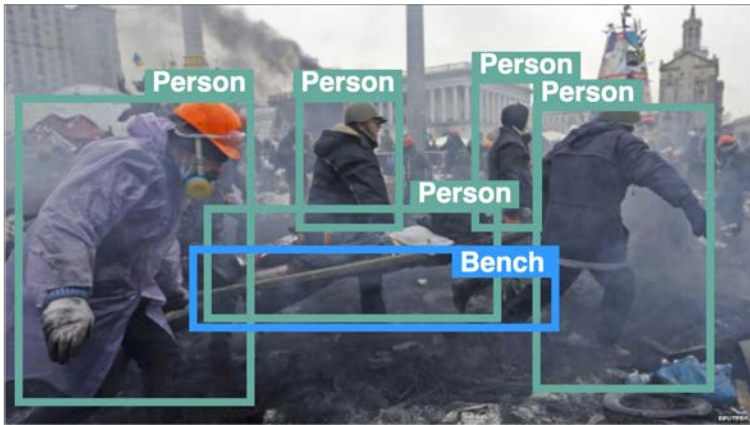
Caption Text t

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
Target	injured man
Instrument	stretcher

Image i



Negative Labels (events)

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

Negative Labels (arguments)

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

Role Switching

Hard negatives via 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.

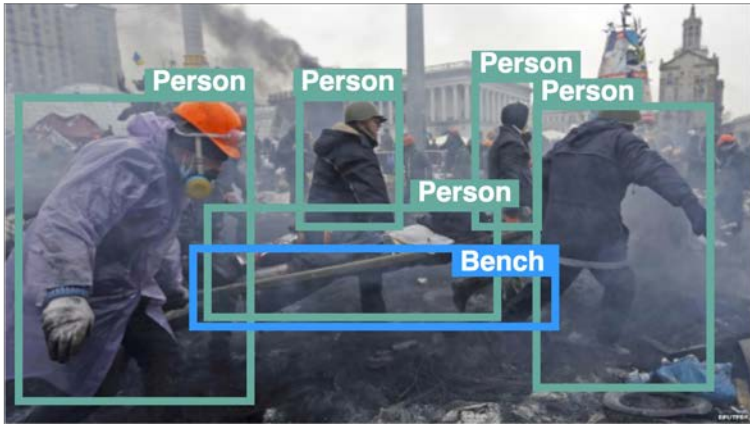
Text
Positive Labels

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

prompt

Protesters transported injured man using a stretcher.

Image i



Negative Labels (events)

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

prompt

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

Negative Labels (arguments)

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

prompt

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

Hard negatives via 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.

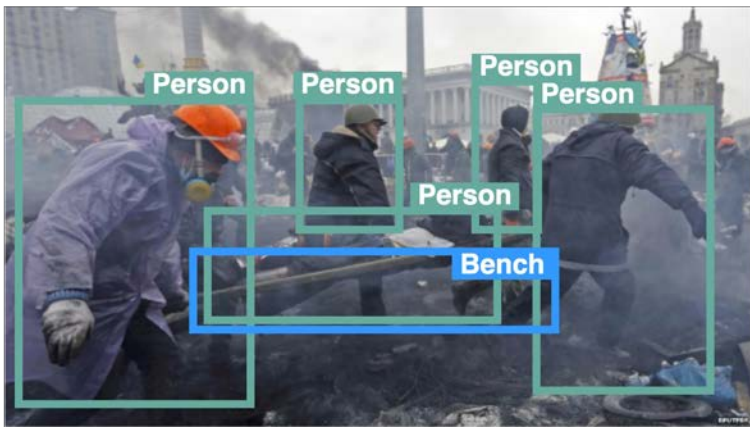
Positive Labels

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

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

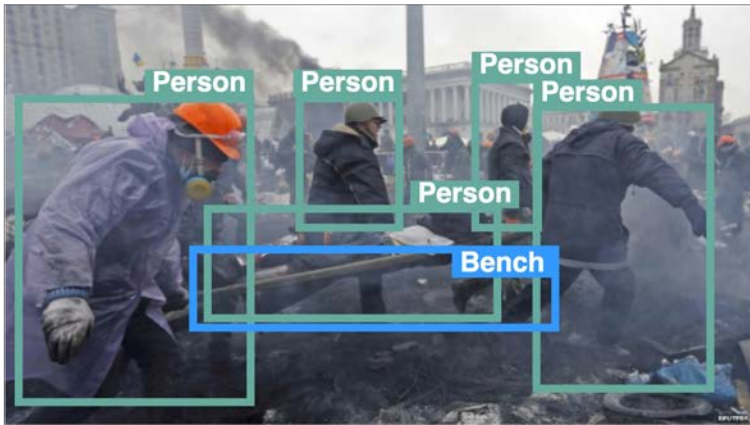
Injured man transported a **stretcher** with **protesters**.

Contrastive Learning on Event Semantics

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.

Image i



Positive Labels

Protesters transported injured man using a stretcher.

Negative Labels (events)

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

Negative Labels (arguments)

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

Text Encoder

t_0

$s(t_0, v)$

t_1

$s(t_1, v)$

t_2

$s(t_2, v)$

Image Encoder

v

Contrastive Learning

Bottlenecks of Vision Semantic Structure Learning



SPRAYING

ROLE	VALUE
AGENT	MAN
SOURCE	SPRAY CAN
SUBSTANCE	PAINT
DESTINATION	WALL
PLACE	ALLEYWAY

ROLE	VALUE
AGENT	FIREMAN
SOURCE	HOSE
SUBSTANCE	WATER
DESTINATION	FIRE
PLACE	OUTSIDE

CLIPPING

ROLE	VALUE
AGENT	MAN
SOURCE	SHEEP
TOOL	SHEARS
ITEM	WOOL
PLACE	FIELD

ROLE	VALUE
AGENT	VET
SOURCE	DOG
TOOL	CLIPPER
ITEM	CLAW
PLACE	ROOM

(Yatskar et al., 2016, ...)



Surfing

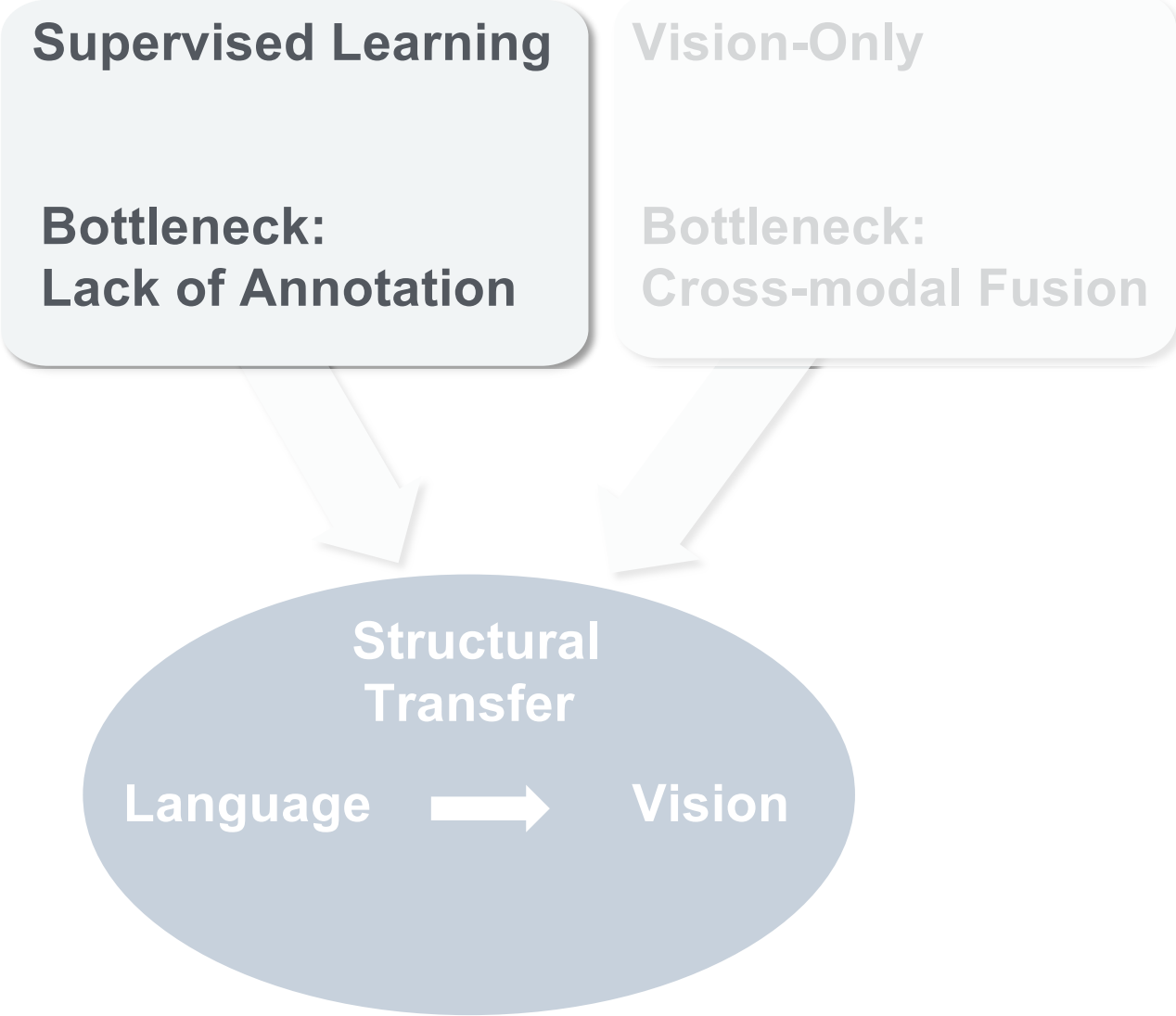
Agent	Tool	Path	Place
Man	Surfboard	Water	Ocean



Feeding

Agent	Food	Eater	Source	Place
Person	Milk	Tiger	Bottle	∅

(Pratt et al., 2020, ...)



Bottlenecks of Vision Semantic Structure Learning



SPRAYING

ROLE	VALUE
AGENT	MAN
SOURCE	SPRAY CAN
SUBSTANCE	PAINT
DESTINATION	WALL
PLACE	ALLEYWAY

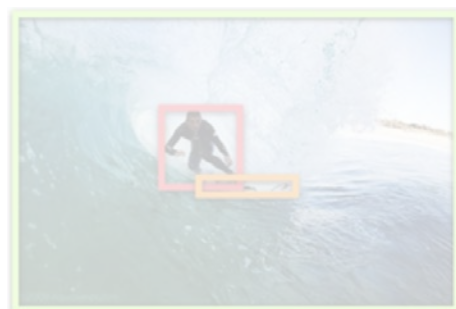
ROLE	VALUE
AGENT	FIREMAN
SOURCE	HOSE
SUBSTANCE	WATER
DESTINATION	FIRE
PLACE	OUTSIDE

CLIPPING

ROLE	VALUE
AGENT	MAN
SOURCE	SHEEP
TOOL	SHEARS
ITEM	WOOL
PLACE	FIELD

ROLE	VALUE
AGENT	VET
SOURCE	DOG
TOOL	CLIPPER
ITEM	CLAW
PLACE	ROOM

(Yatskar et al., 2016, ...)



Surfing

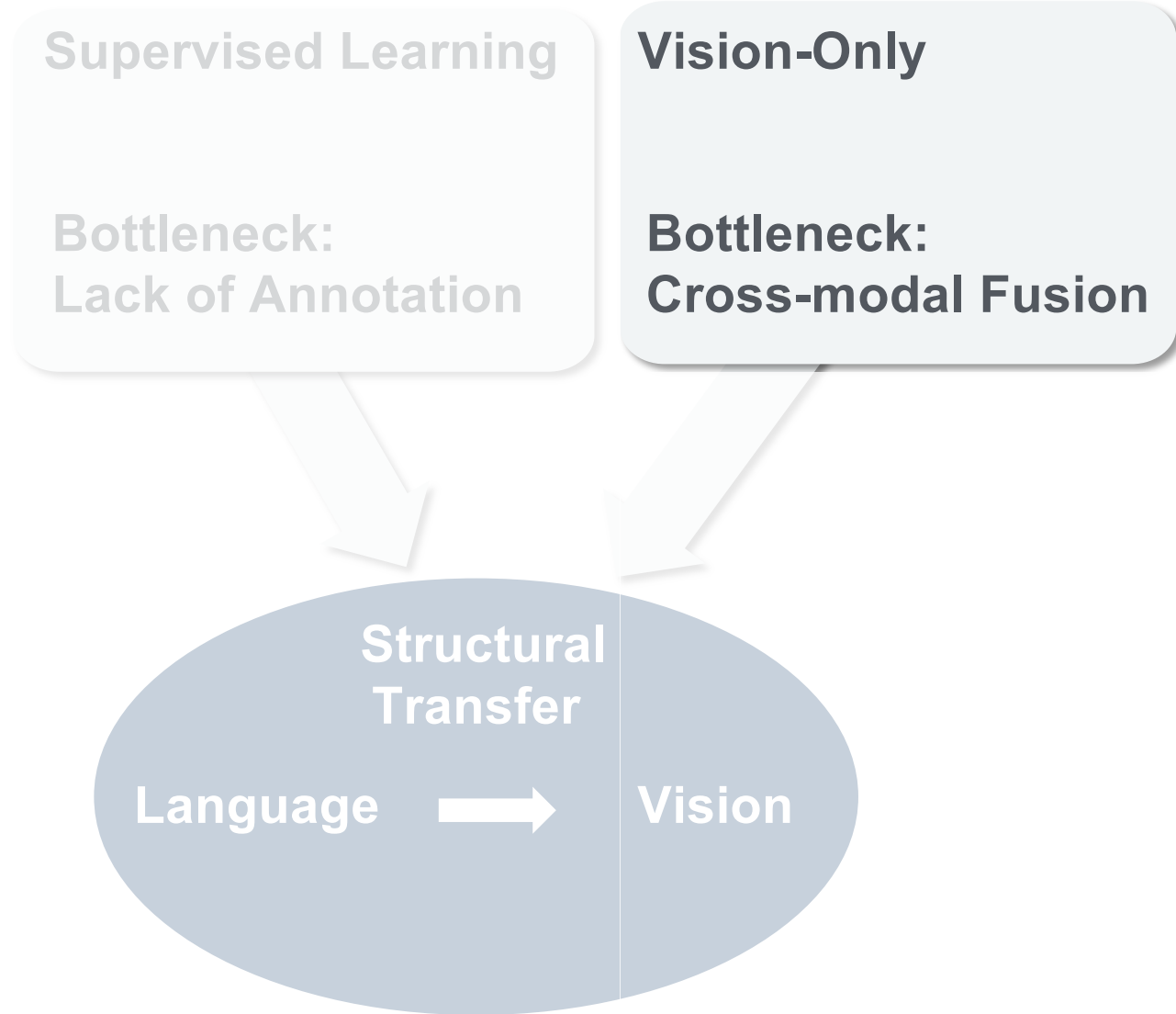
Agent	Tool	Path	Place
Man	Surfboard	Water	Ocean



Feeding

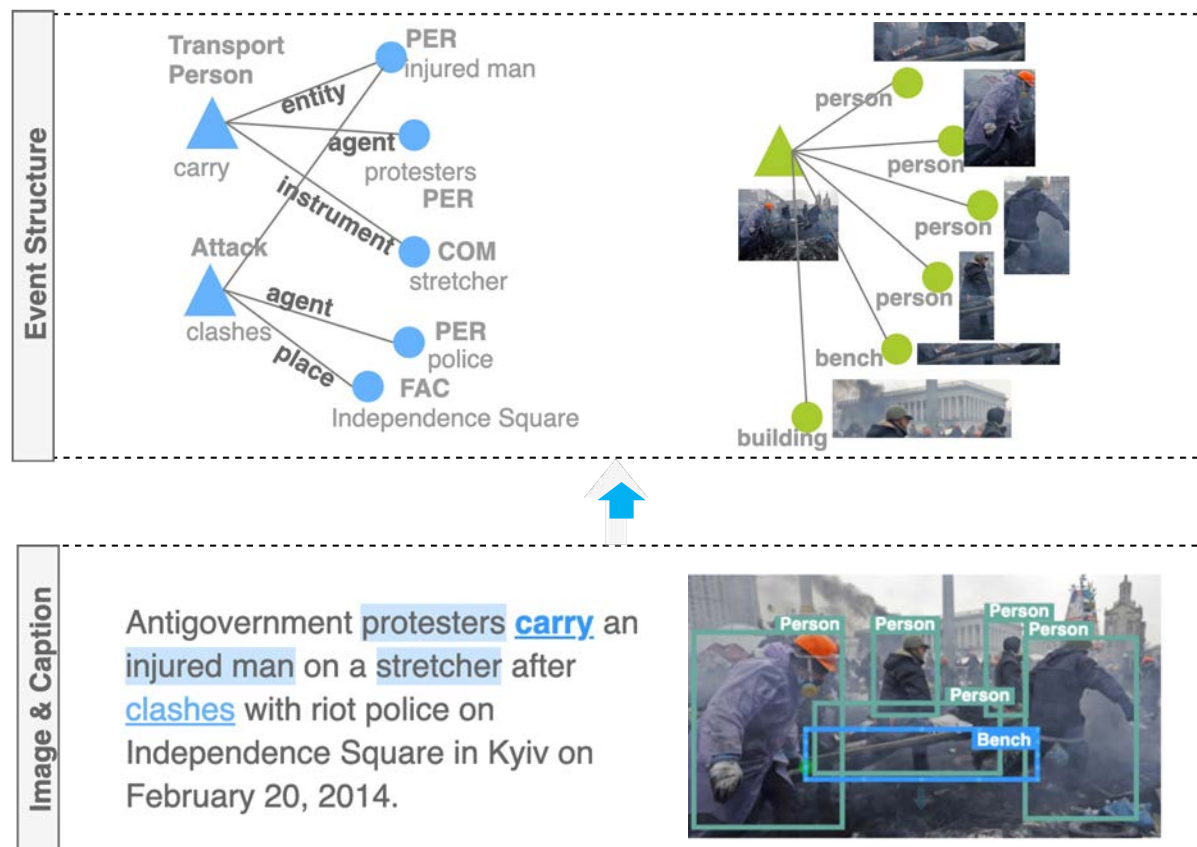
Agent	Food	Eater	Source	Place
Person	Milk	Tiger	Bottle	∅

(Pratt et al., 2020, ...)

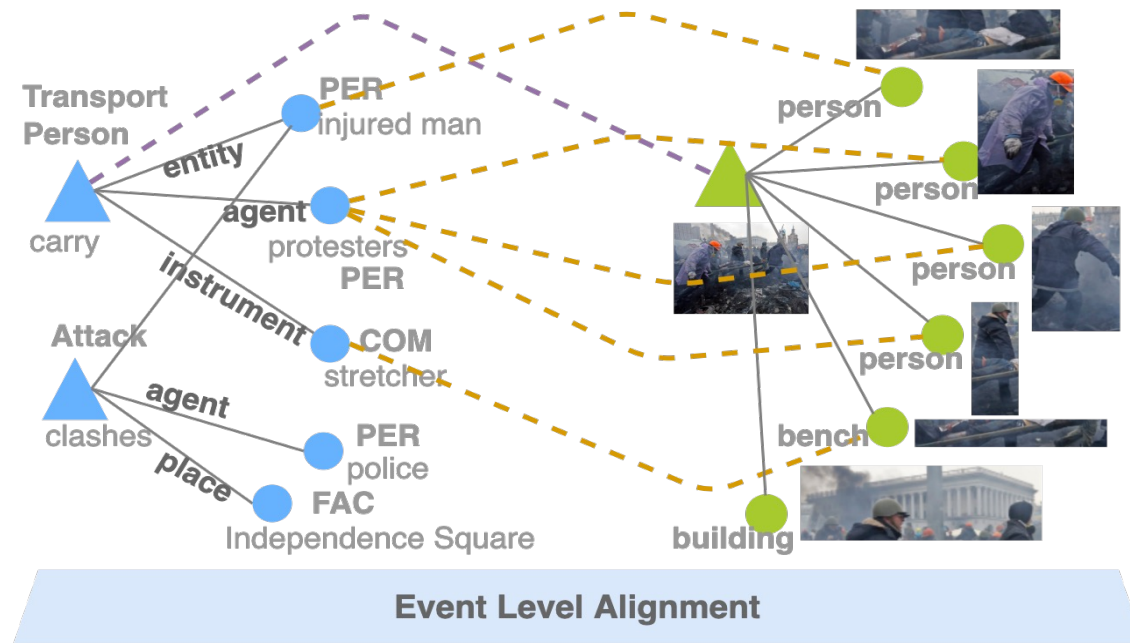


The first V+L Pretraining with Event Semantic Structures

Challenge: Structured Encoding



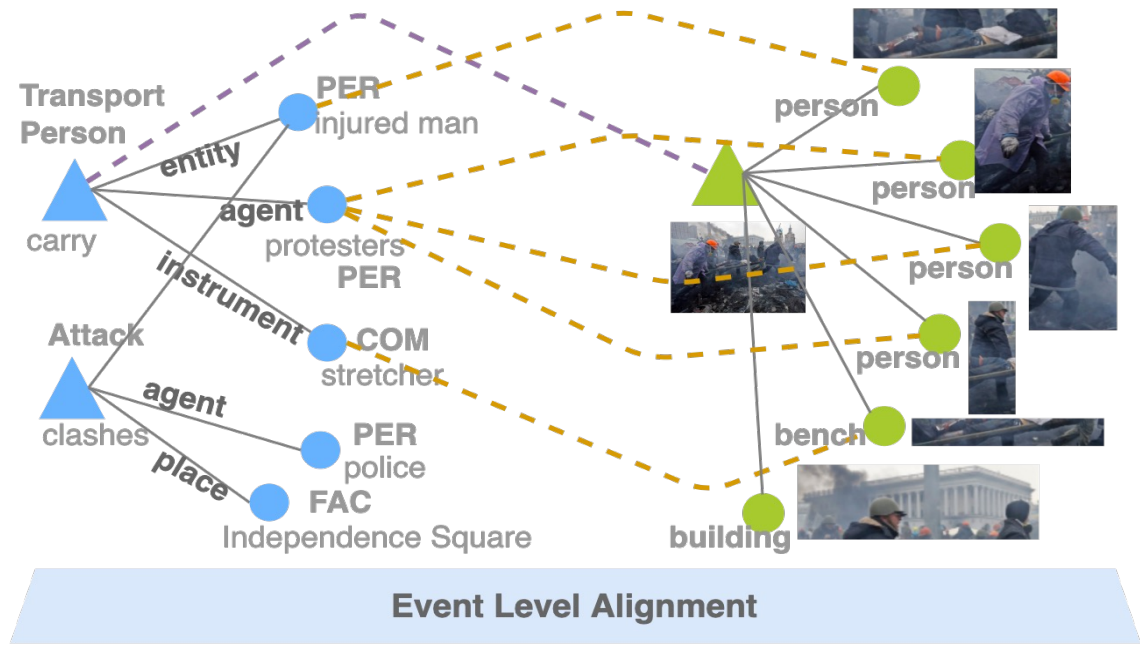
The first V+L Pretraining with Event Semantic Structures



The first V+L Pretraining with Event Semantic Structures

Structured Alignment via Optimal Transport

Text Event Graph ↔ Image Event Graph

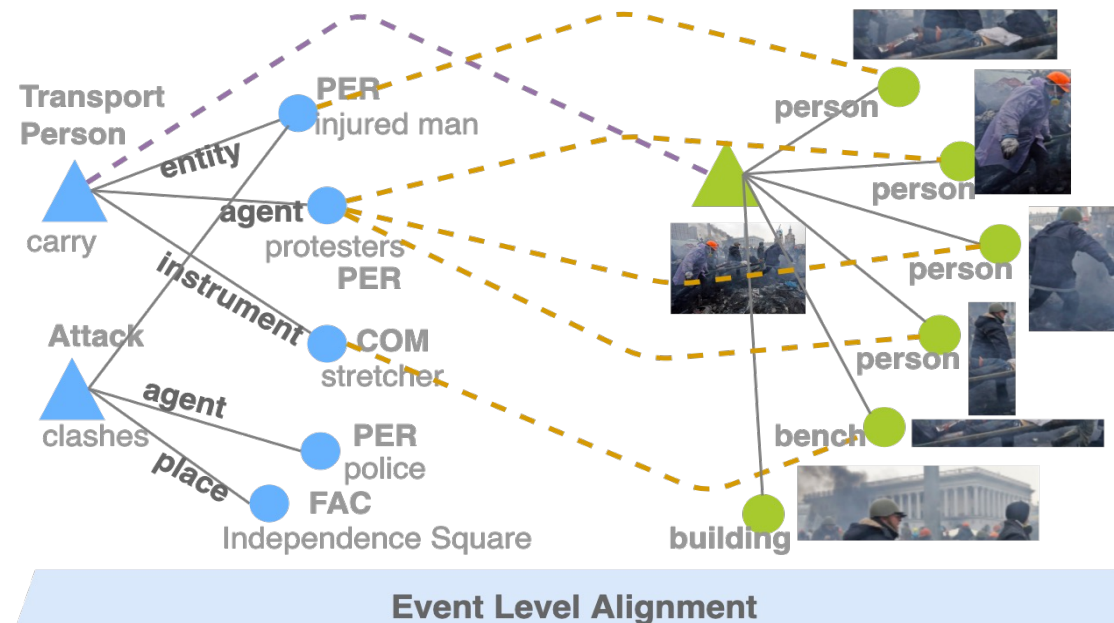


The first V+L Pretraining with Event Semantic Structures

Structured Alignment via Optimal Transport

Text Event Graph \longleftrightarrow Image Event Graph

1 Define cost matrix C (embedding similarity)



The first V+L Pretraining with Event Semantic Structures

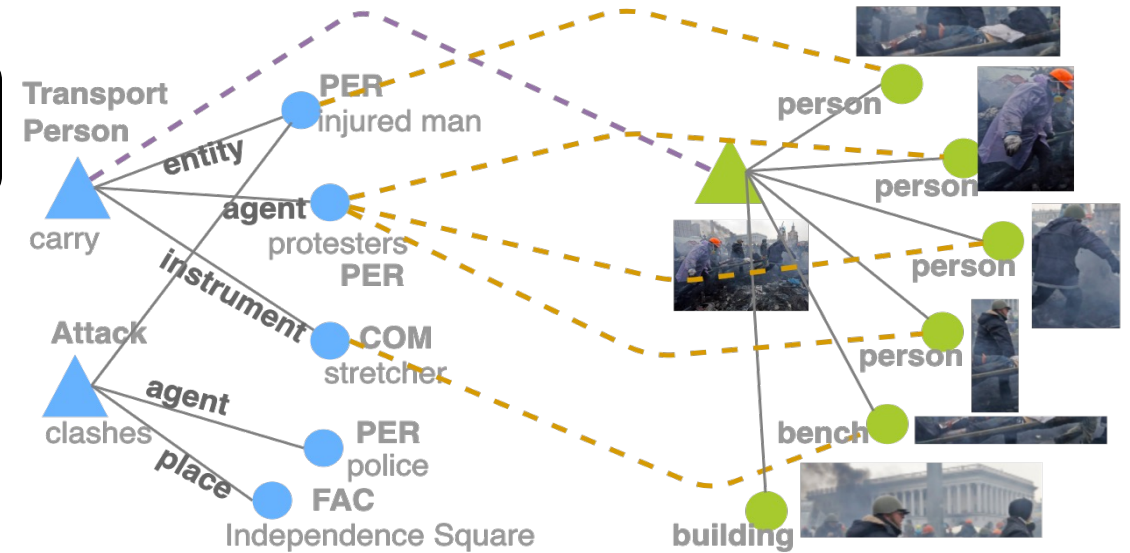
Structured Alignment via Optimal Transport

Text Event Graph \longleftrightarrow Image Event Graph

1 Define cost matrix C (embedding similarity)

2 Optimization Goal: minimize transport distance

$$D(S, T) = \min_T T \cdot C$$



Event Level Alignment

The optimal T is approximated by a differentiable Sinkhorn Knopp algorithm (Sinkhorn, 1964; Cuturi, 2013)

$$T = \text{diag}(p) \exp(-C/\gamma) \text{diag}(q)$$

for $i = 0, 1, 2, \dots$ until convergence,

$$p^{i+1} = \mathbf{1} \oslash (K q^i),$$

$$q^{i+1} = \mathbf{1} \oslash (K^\top p^{i+1}),$$

$$T^k := \text{diag}(p^k) K \text{diag}(q^k)$$

The first V+L Pretraining with Event Semantic Structures

Structured Alignment via Optimal Transport

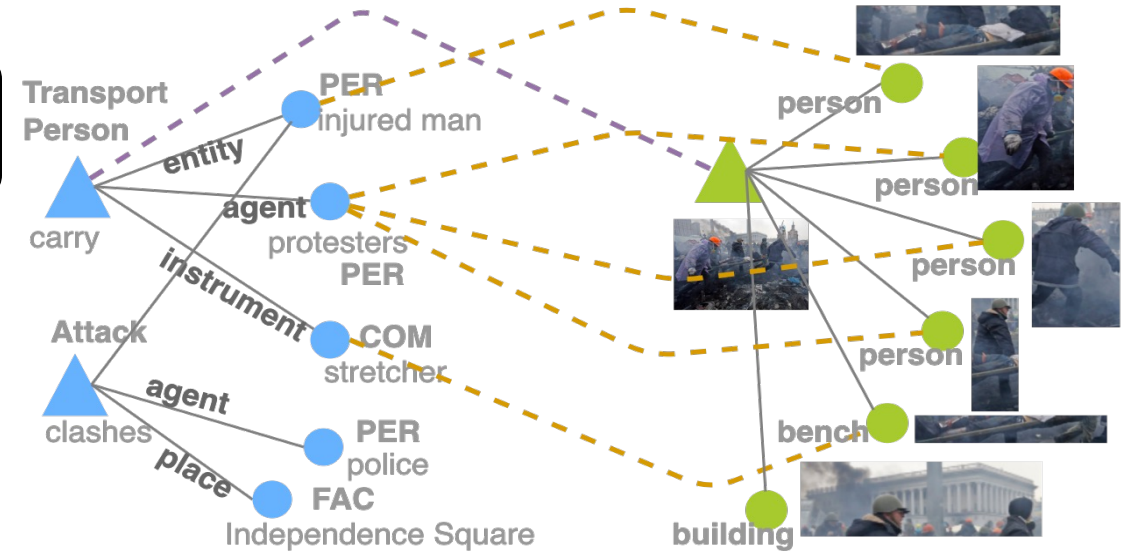
Text Event Graph \longleftrightarrow Image Event Graph

1 Define cost matrix C (embedding similarity)

2 Optimization Goal: minimize transport distance

$$D(S, T) = \min_T T \cdot C$$

3 Optimize the transport plan T within k iterations



Event Level Alignment

The optimal T is approximated by a differentiable Sinkhorn Knopp algorithm (Sinkhorn, 1964; Cuturi, 2013)

$$T = \text{diag}(p) \exp(-C/\gamma) \text{diag}(q)$$

for $i = 0, 1, 2, \dots$ until convergence,

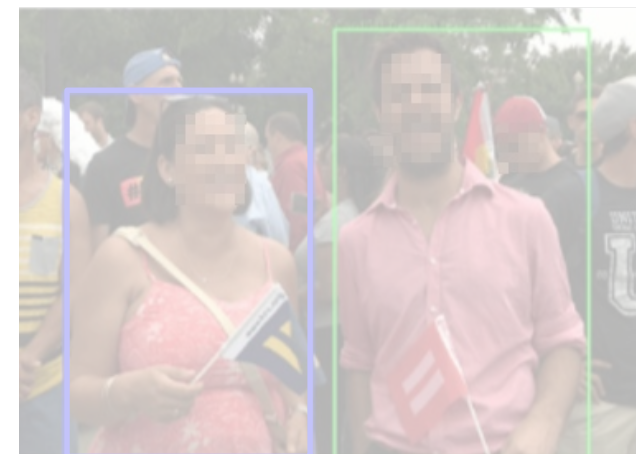
$$p^{i+1} = \mathbf{1} \oslash (Kq^i),$$

$$q^{i+1} = \mathbf{1} \oslash (K^\top p^{i+1}),$$

$$T^k := \text{diag}(p^k) K \text{diag}(q^k)$$

CLIP-Event on Visual Event Extraction

Supporting **Zero-shot** Vision Event Extraction the first time.



Vaccination



Arrest



Protest

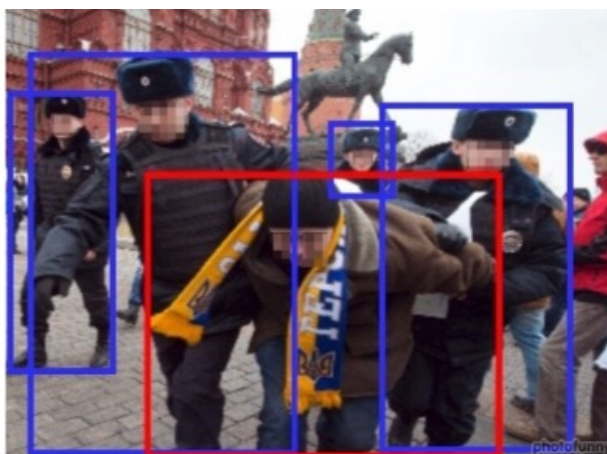


Celebration



CLIP-Event on Visual Event Extraction

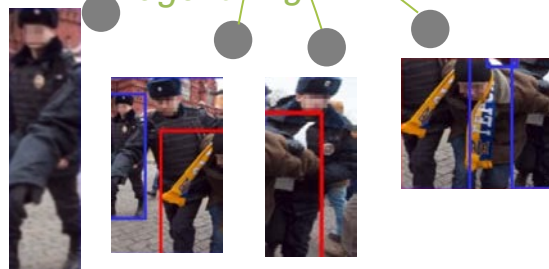
Supporting **Zero-shot** Vision Event Extraction the first time.



Vaccination



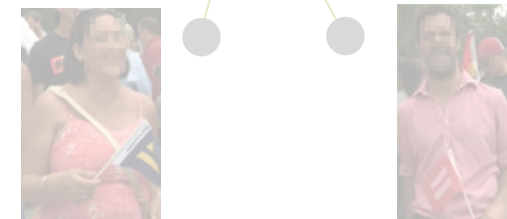
Arrest



Protest



Celebration



CLIP-Event on Visual Event Extraction

Supporting **Zero-shot** Vision Event Extraction the first time.



Vaccination



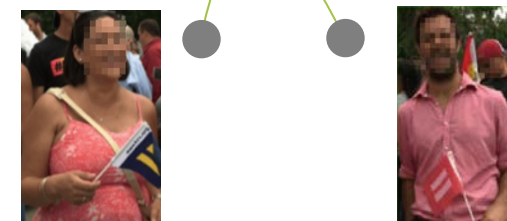
Arrest



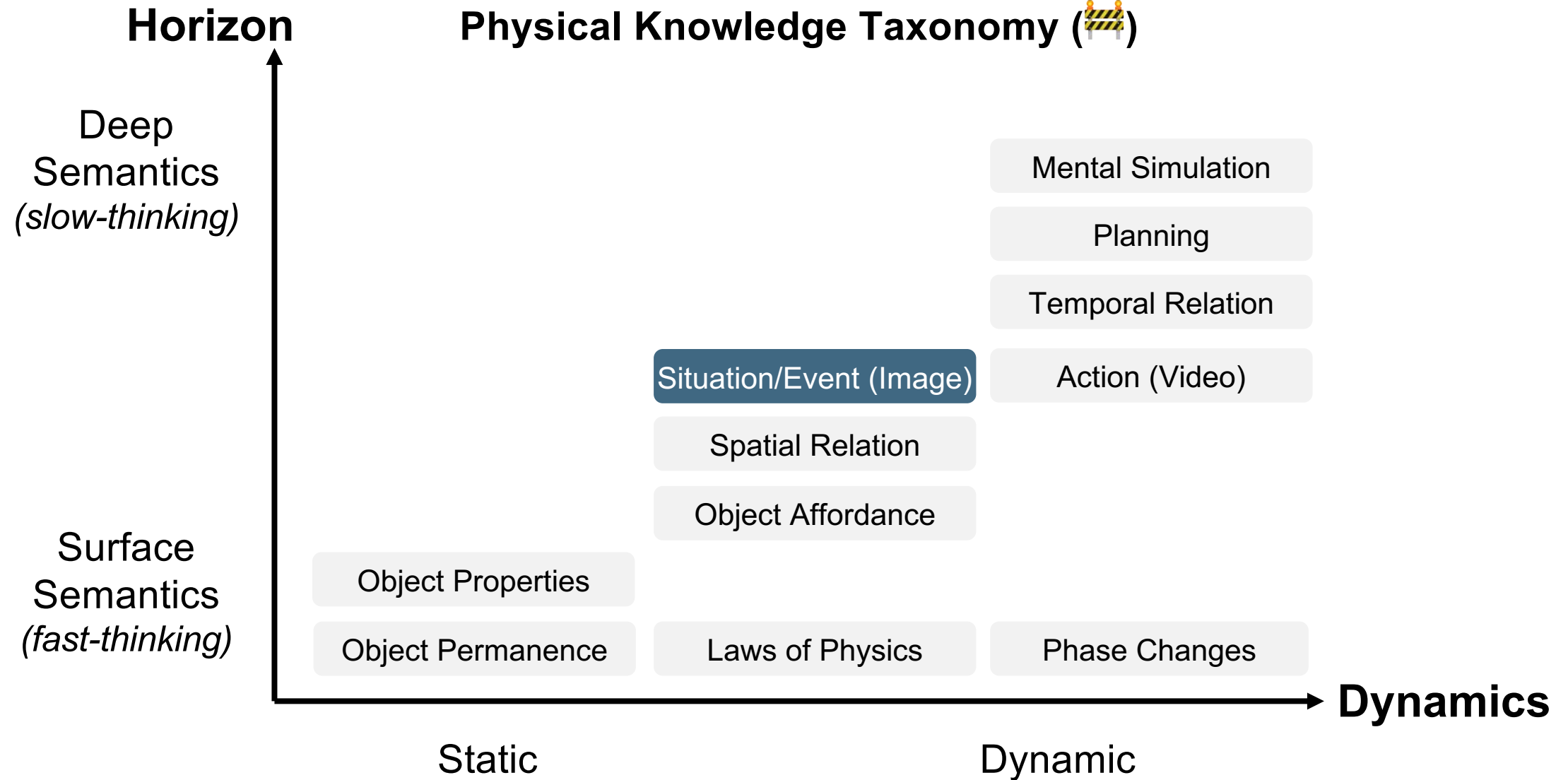
Protest



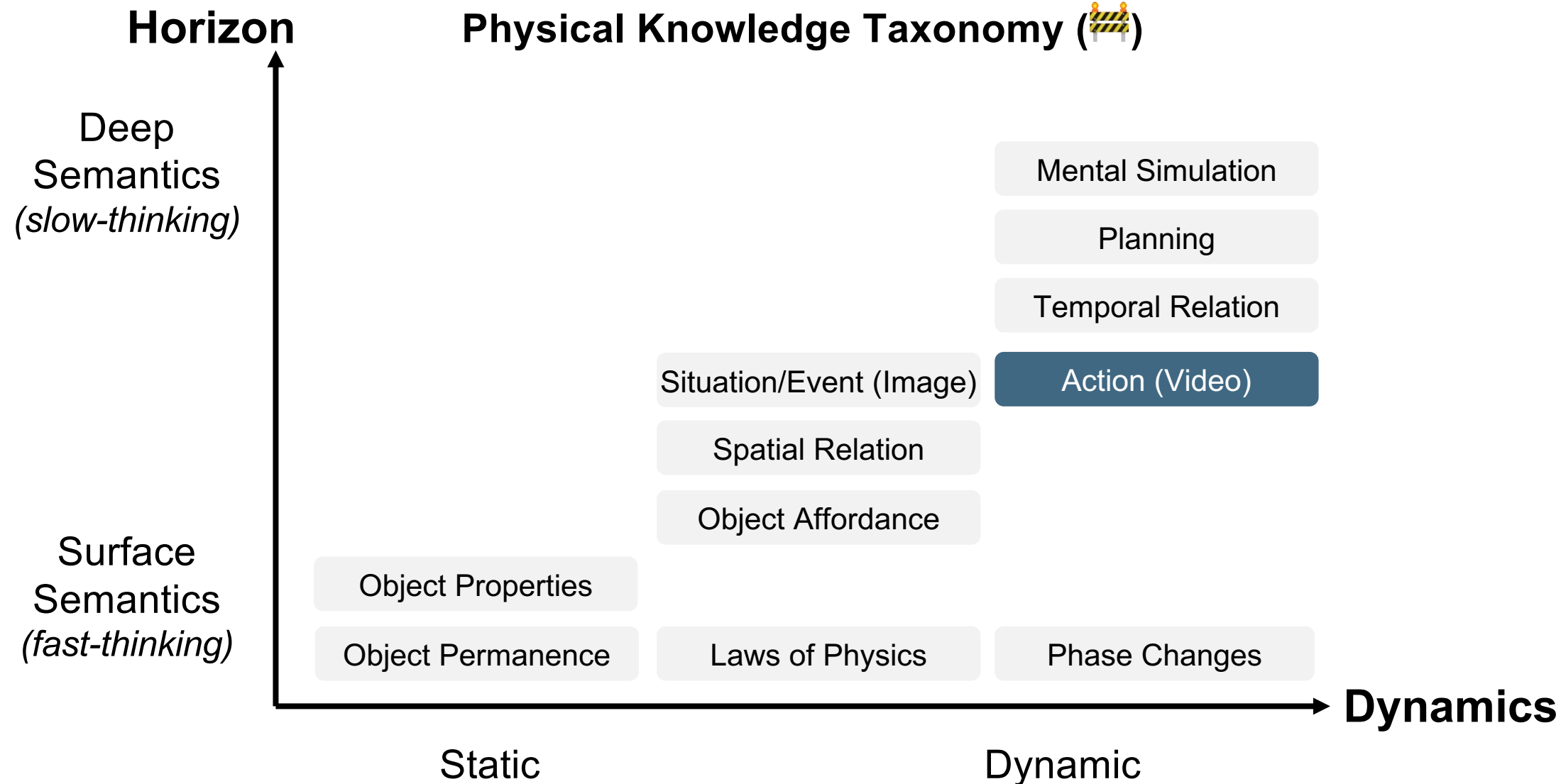
Celebration



Our Research Goal: Surface → Deep Semantics



Our Research Goal: Surface → Deep Semantics



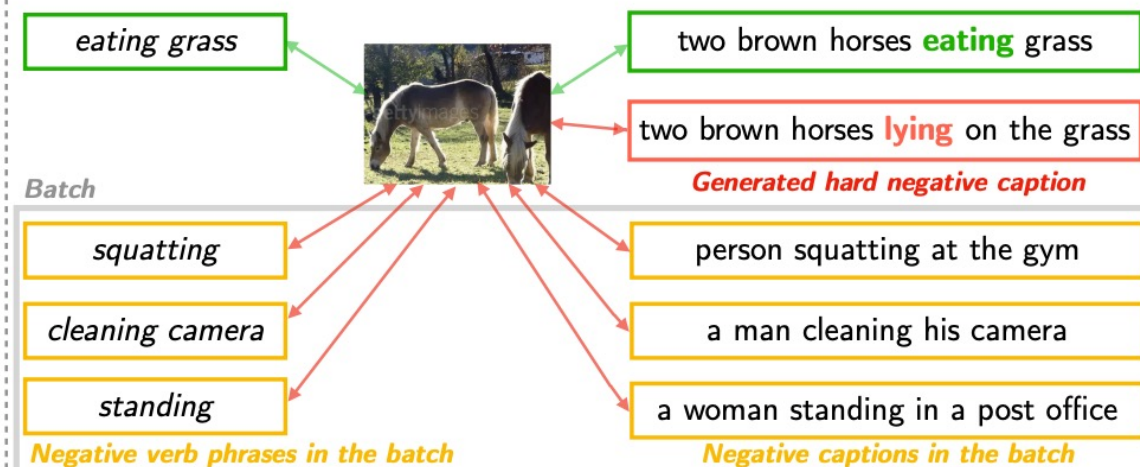
Verbs in Action: Improving verb understanding in video-language models



Video



Verb phrase Hard verb negative captions



it's a video of a bald monk **sitting** at a temple **looking** at his laptop

it's a video of a bald monk **lying** at a temple looking at his laptop

it's a video of a bald monk **standing** at a temple looking at his laptop

it's a video of a bald monk **dancing around** a temple **holding** his laptop

it's a video of a bald monk **jumping up** at a temple **closing** his laptop

it's a video of a bald monk **running in** a temple **searching for** his laptop



a person **draws** a dragon

a person **carves** a dragon

a person **paints** a dragon

a person **doodles** a dragon

a person **sculpts** a dragon

a person **destroys** a dragon



a girl **skateboarding** in a public place

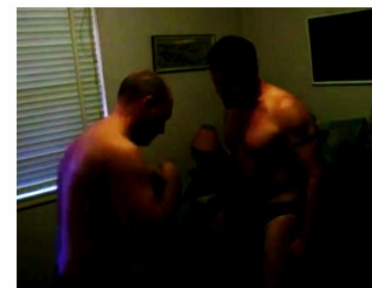
a girl **dancing** in a public place

a girl **running** in a public place

a girl **singing** in a public place

a girl **sitting on her skateboard** in a public place

a girl **falling off her skateboard** in a public place



man is **punching** another man in the dark

man is **arguing with** another man in the dark

man is **kissing** another man in the dark

man is **talking to** another man in the **daylight**

man is **kicking** another man in the **light**

man is **hugging** another man in the dark

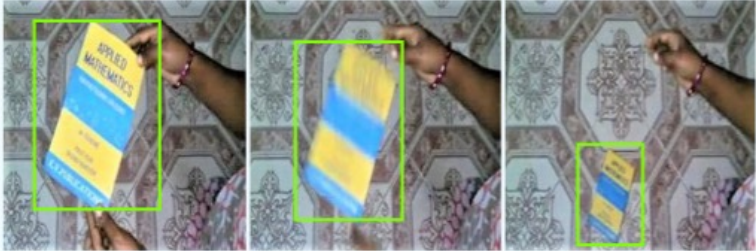
Video: A “Visual Recording” of World State Changes

Do SOTA Video-Language Models (VLM) possess fundamental **Action** Knowledge?

Video: A “Visual Recording” of World State Changes

Do SOTA Video-Language Models (VLM) possess fundamental **Action** Knowledge?

Probing Task: Action Antonym (AA)



"Book **falling** like a rock"
Original Action Text



GT VidLM Result
23.2%

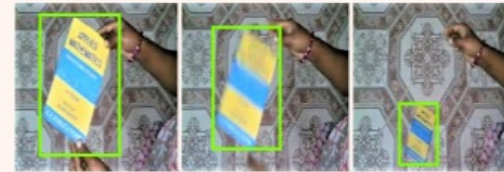
"Book **rising** like a rock"
Action Antonym Text



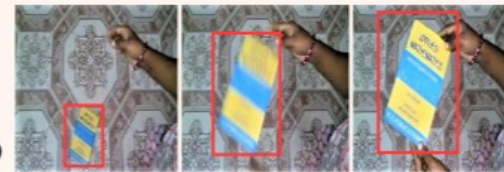
76.8%

Probing Task: Video Reversal (VR)

"Book **falling** like a rock"
Original Action Text



GT VidLM Result
49.9%



50.1%

Baseline Task: Object Replacement (OR)



"**Book** falling like a rock"
Original Action Text



GT VidLM Result
77.9%

"**Cellphone** falling like a rock"
Object Replaced Text

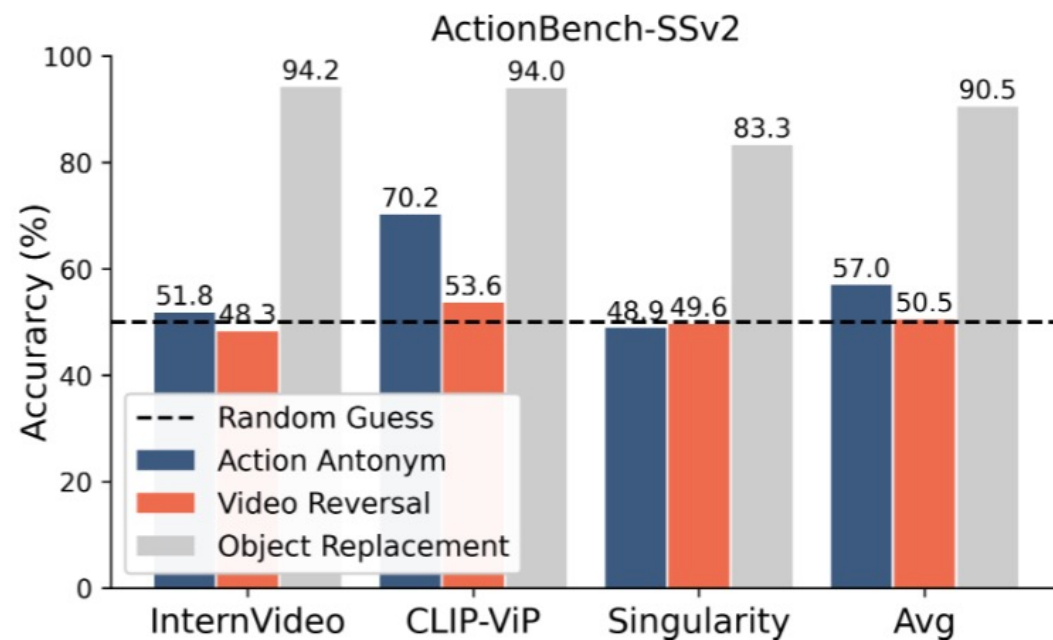
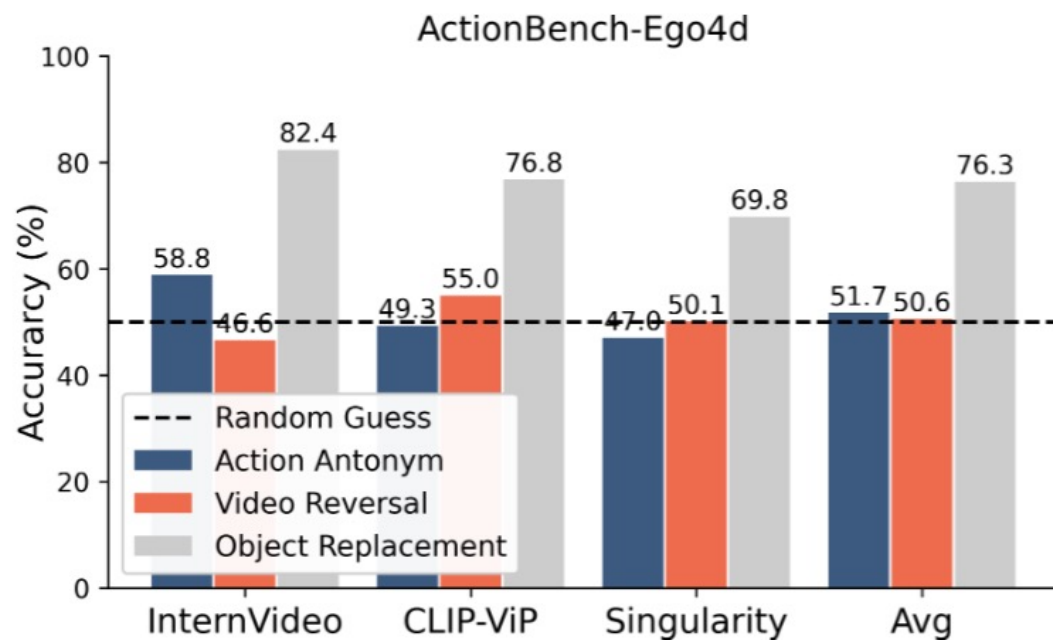


22.1%

Video: A “Visual Recording” of World State Changes

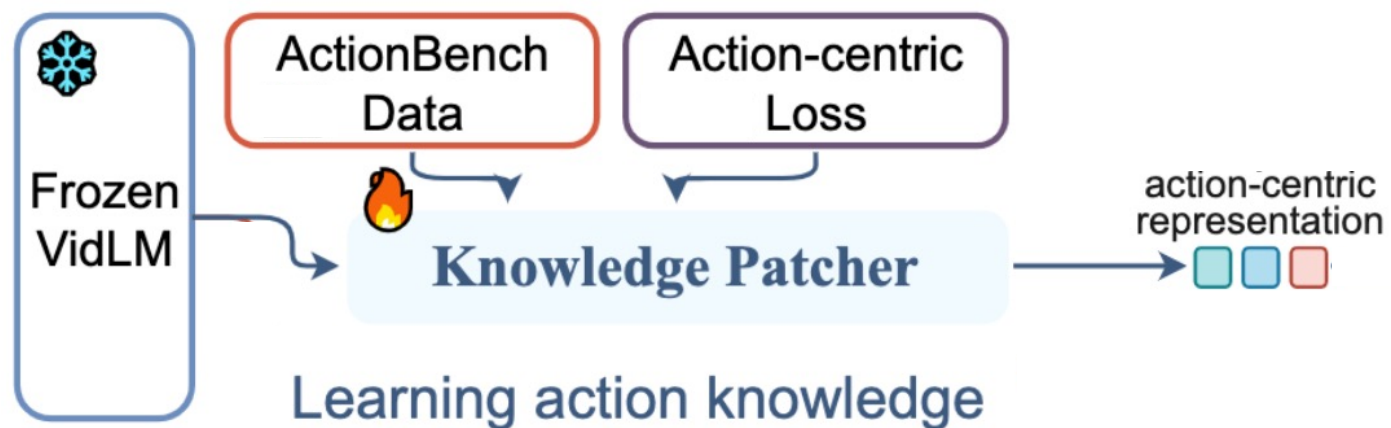
Do SOTA Video-Language Models (VLM) possess fundamental **Action** Knowledge?

- Near **random** performance on Action Antonym (AA) and Video Reversal (VR)
- Clear **biases towards objects** compared to actions



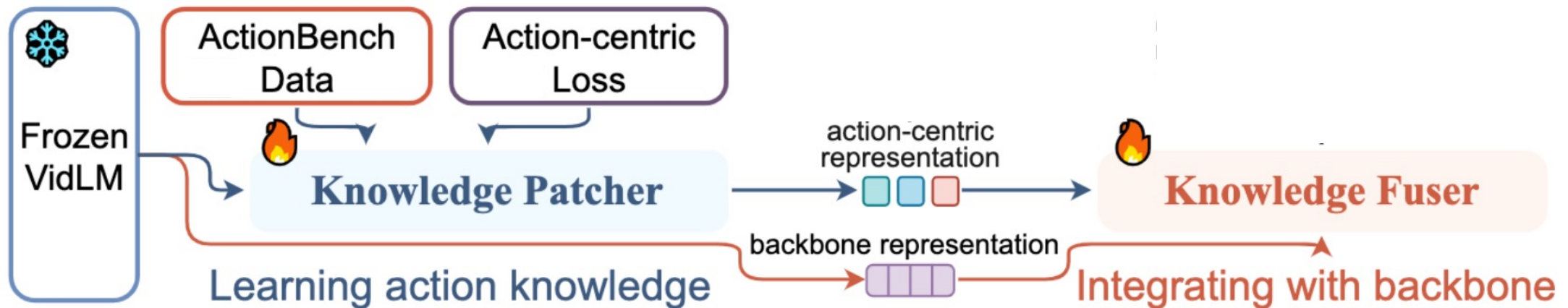
Patch & Fuse: Patching frozen VLMs with Action Knowledge

Patch frozen VidLMs **with action knowledge** without hurting their general VL capabilities.



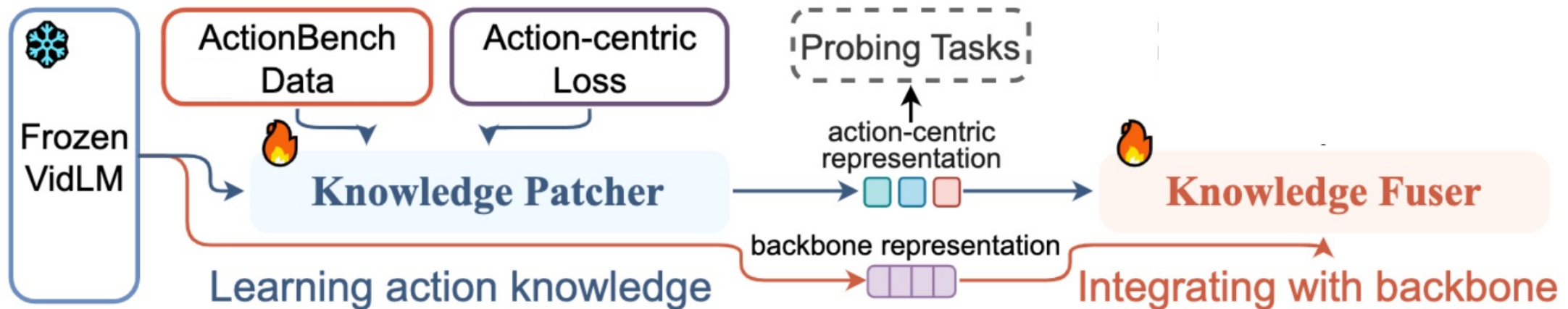
Patch & Fuse: Patching frozen VLMs with Action Knowledge

Patch frozen VidLMs with action knowledge *without hurting their general VL capabilities*.



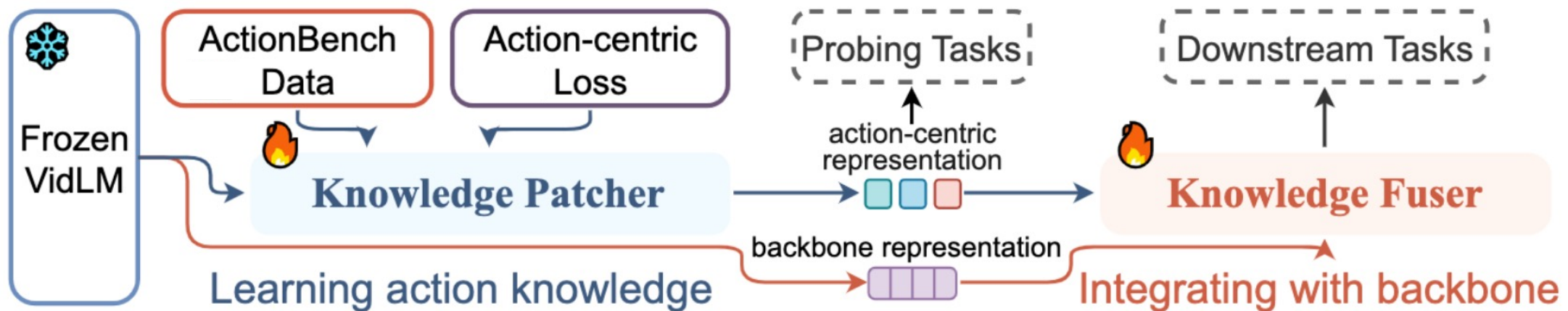
Patch & Fuse: Patching frozen VLMs with Action Knowledge

Patch frozen VidLMs with action knowledge **without hurting their general VL capabilities.**



Patch & Fuse: Patching frozen VLMs with Action Knowledge

Patch frozen VidLMs with action knowledge **without hurting their general VL capabilities.**

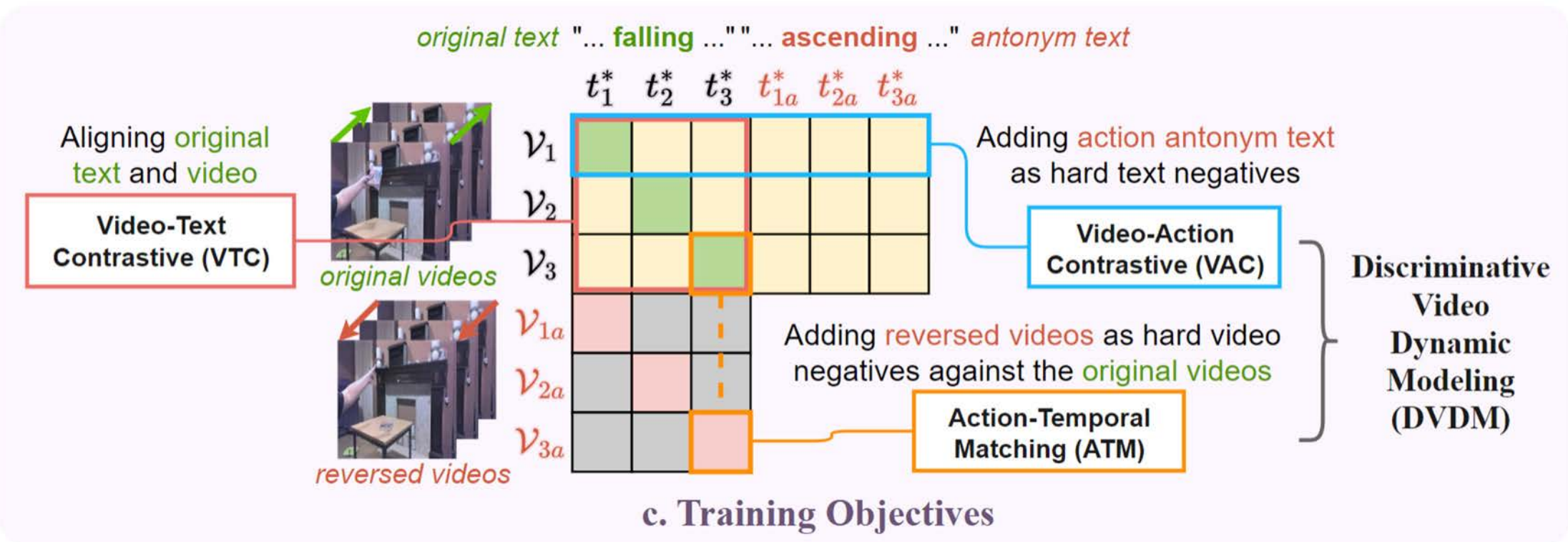


Patch & Fuse: Patching frozen VLMs with Action Knowledge



Video-Action Contrastive (VAC): encourages learning the **alignment** between the **video** and the **action verbs**

Action-Temporal Matching: encourages learning the correct **temporal ordering** implied by the **action text**

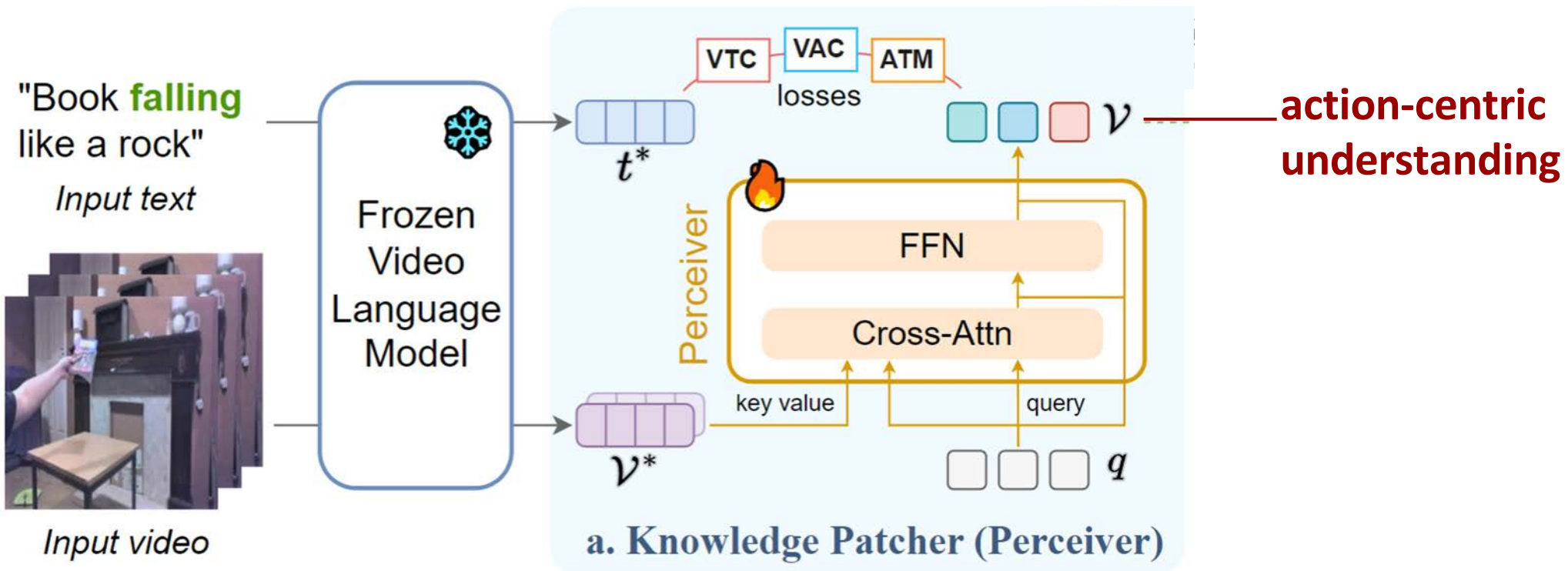


Patch & Fuse: Patching frozen VLMs with Action Knowledge



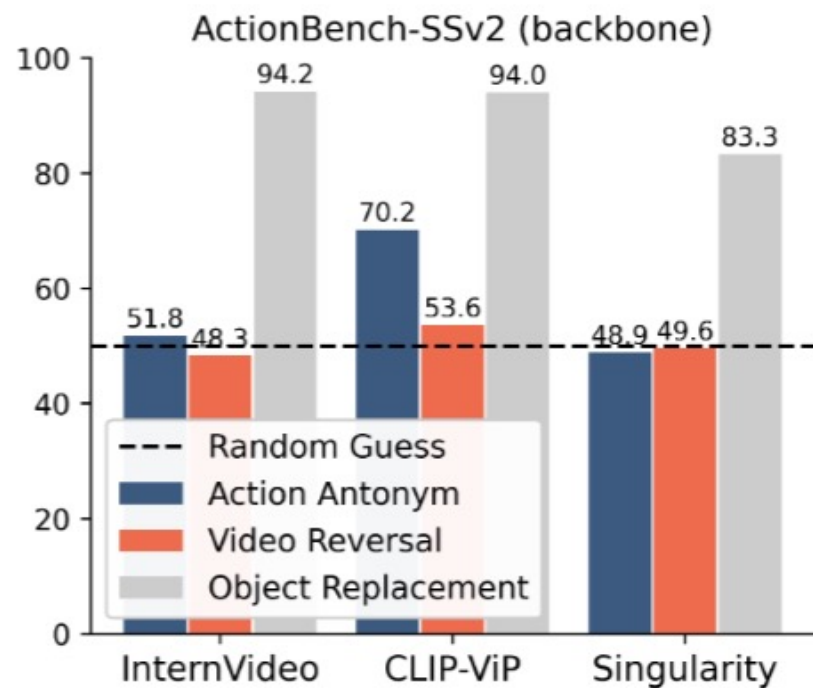
Video-Action Contrastive (VAC): encourages learning the **alignment** between the **video** and the **action verbs**

Action-Temporal Matching: encourages learning the correct **temporal ordering** implied by the **action text**



Results regarding Patch: Before vs After adding Knowledge Patcher

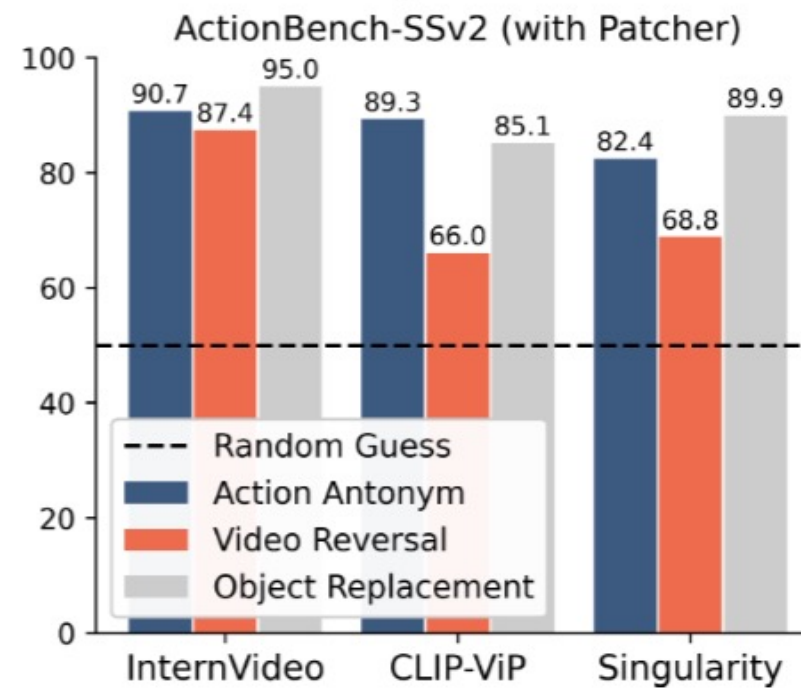
Adding the Knowledge Patcher nearly doubles the performance.



+ Knowledge Patcher



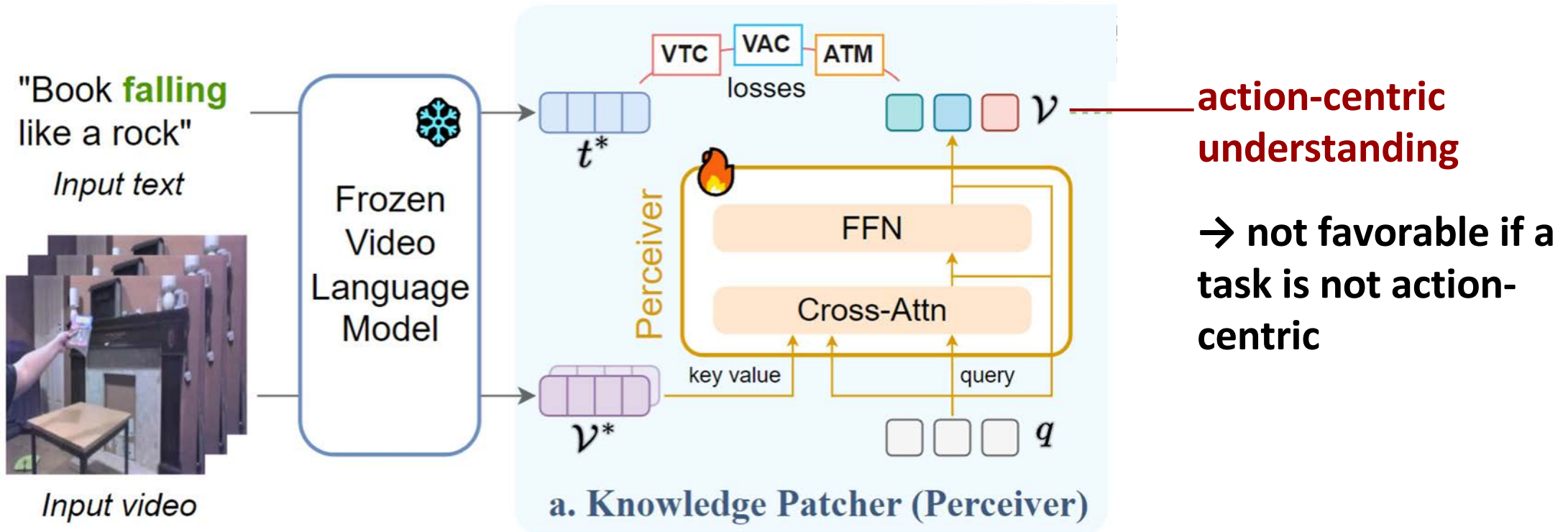
Trained with VTC+DVDM



Patch & Fuse: Retaining VL Capabilities

Video-Action Contrastive (VAC): encourages learning the **alignment** between the **video** and the **action verbs**

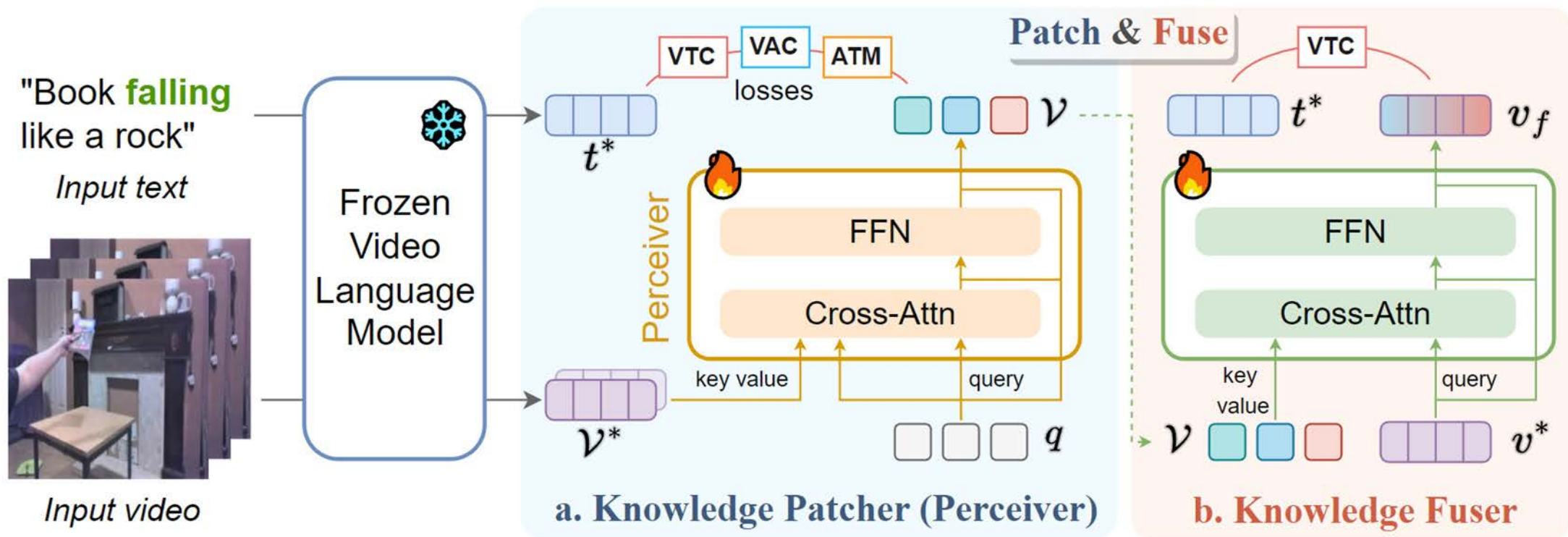
Action-Temporal Matching: encourages learning the correct **temporal ordering** implied by the **action text**



Patch & Fuse: Retaining VL Capabilities



A **unified representation** that has good understanding of both actions and objects.



Patch & Fuse: Retaining VL Capabilities



Video-Text Retrieval

SSv2-Label

Causal-Temporal VQA

NExT-QA

Video-to-Action Retrieval

SSv2-Template

Temporal-SSv2

← More object-centric

Require joint understanding of objects and actions

More action-centric →

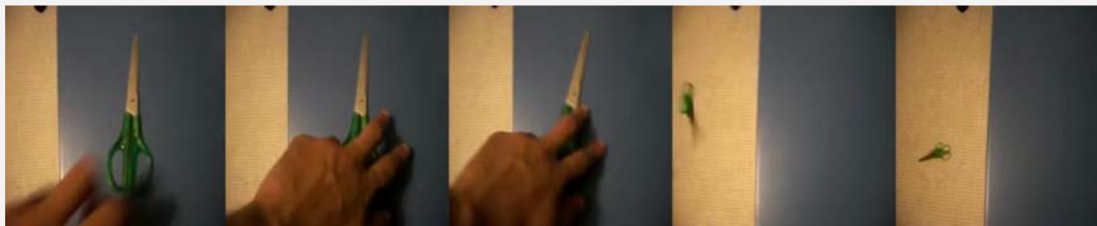
Example



Video-Text Retrieval: SSv2-label

"pushing scissors so that it falls off the table"

Example



Video-to-Action Retrieval: SSv2-template (where the main object is obfuscated)

"pushing **something** so that it falls off the table"

Results regarding Fuse: Retaining VL Capabilities



Video-Text Retrieval

SSv2-Label

Causal-Temporal VQA

NExT-QA

Video-to-Action Retrieval

SSv2-Template

Temporal-SSv2

← More object-centric

Require joint understanding of objects and actions

More action-centric →

Method [Patcher Training Loss]	Video-Text Retrieval SSv2-label			
	$R1_{v2t}$	$R5_{v2t}$	$R1_{t2v}$	$R5_{t2v}$
InternVideo Backbone	18.8	39.9	19.9	40.0
KP-Transformer FT [VTC]	24.1	50.0	21.7	46.0
KP-Perceiver FT [VTC]	27.0	57.4	27.1	56.8
Side-Tuning [83] [VTC+DVDM]	30.9	59.2	26.6	53.1
PATCH & FUSE [VTC+DVDM]	32.3	61.2	28.0	54.3

Causal-Temporal VQA NExT-QA	
Val (Acc)	Test (Acc)
43.2	44.3
48.1	49.6
48.0	49.5
56.3	56.4
56.9	56.6

Video-to-Action Retrieval			
SSv2-template		Temporal-SSv2	
$R1$	$R5$	$R1$	$R5$
5.6	15.9	11.2	35.8
21.1	55.9	41.1	88.9
24.8	59.7	42.5	91.3
22.2	55.1	50.2	90.9
26.9	61.5	51.2	91.9

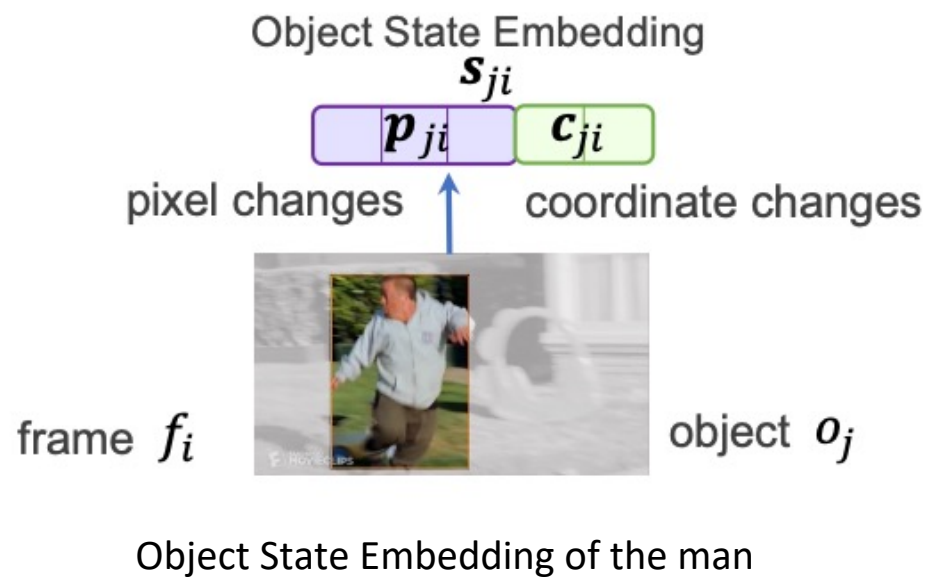
Performs competitively on **both object-centric and action-centric** tasks.

Video Events as Argument State Changes



Video Event =
Status Changes of Arguments

Status Changes of an object =
Displacement (movement of bounding box)
+
Pixel Changes (intra-boundingbox changing)

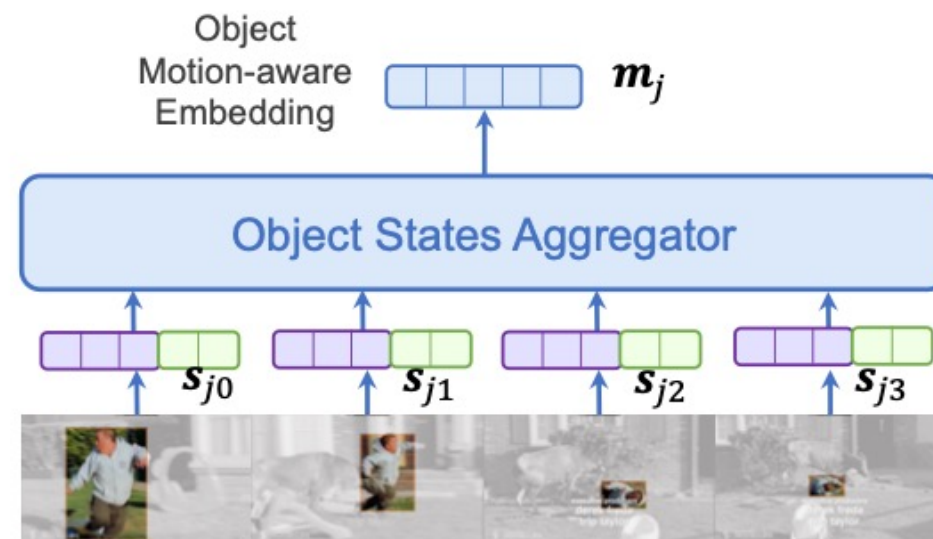


Video Events as Argument State Changes



**Video Event =
Status Changes of Arguments**

**Status Changes of an object =
Displacement (movement of bounding box)
+
Pixel Changes (intra-boundingbox changing)**

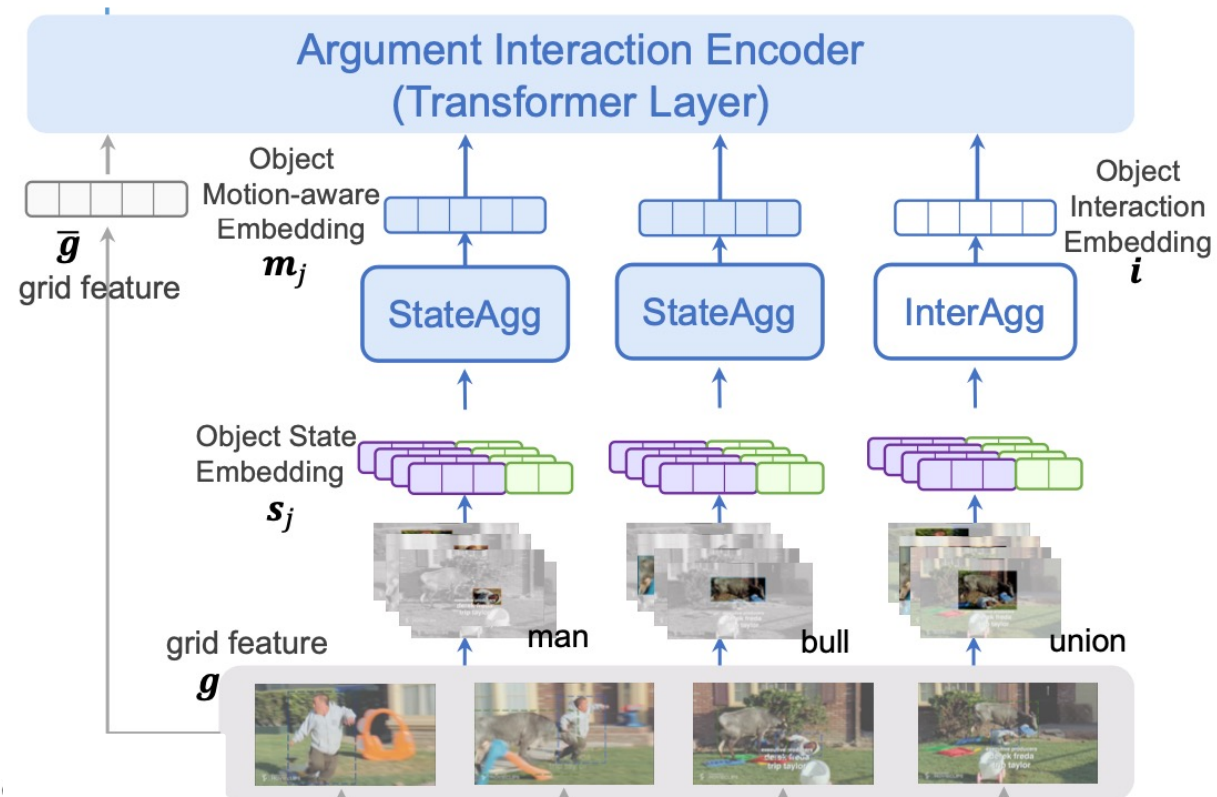


Video Events as Argument State Changes



Video Event =
Status Changes of Arguments

Status Changes of an object =
Displacement (movement of bounding box)
+
Pixel Changes (intra-boundingbox changing)

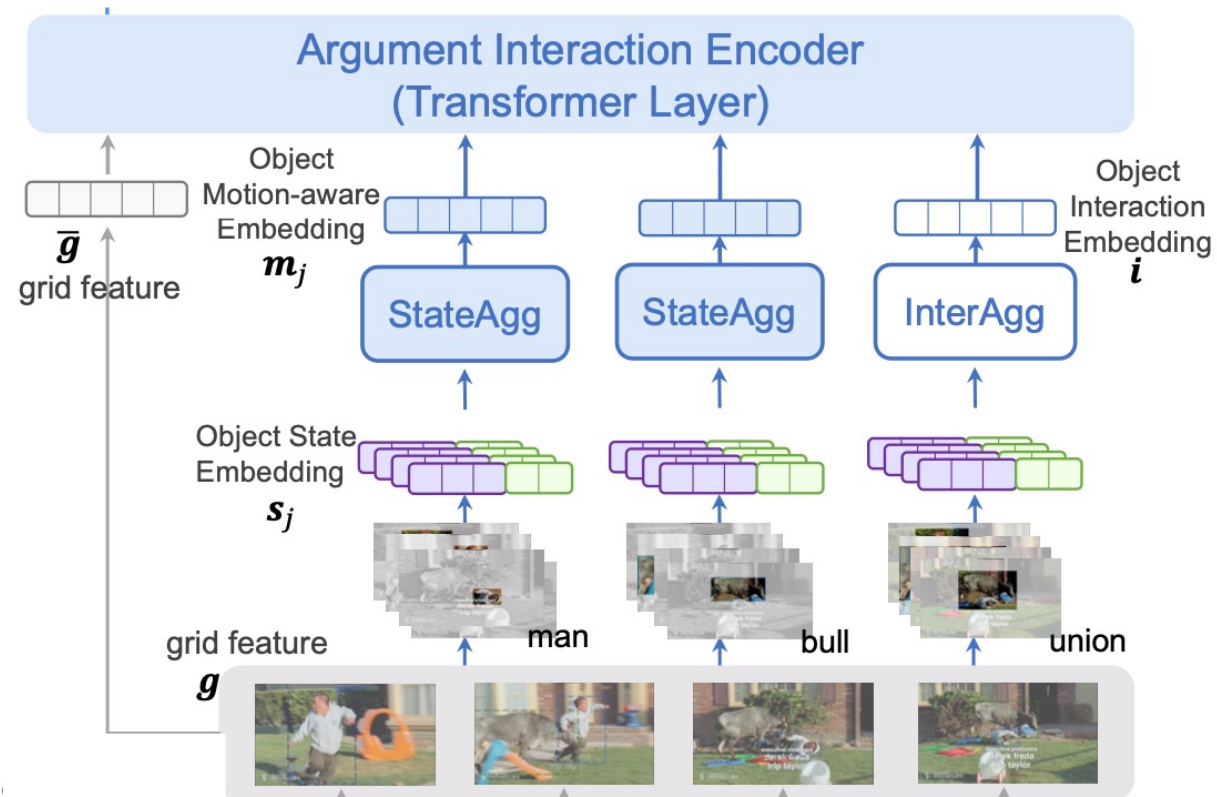


Video Events as Argument State Changes

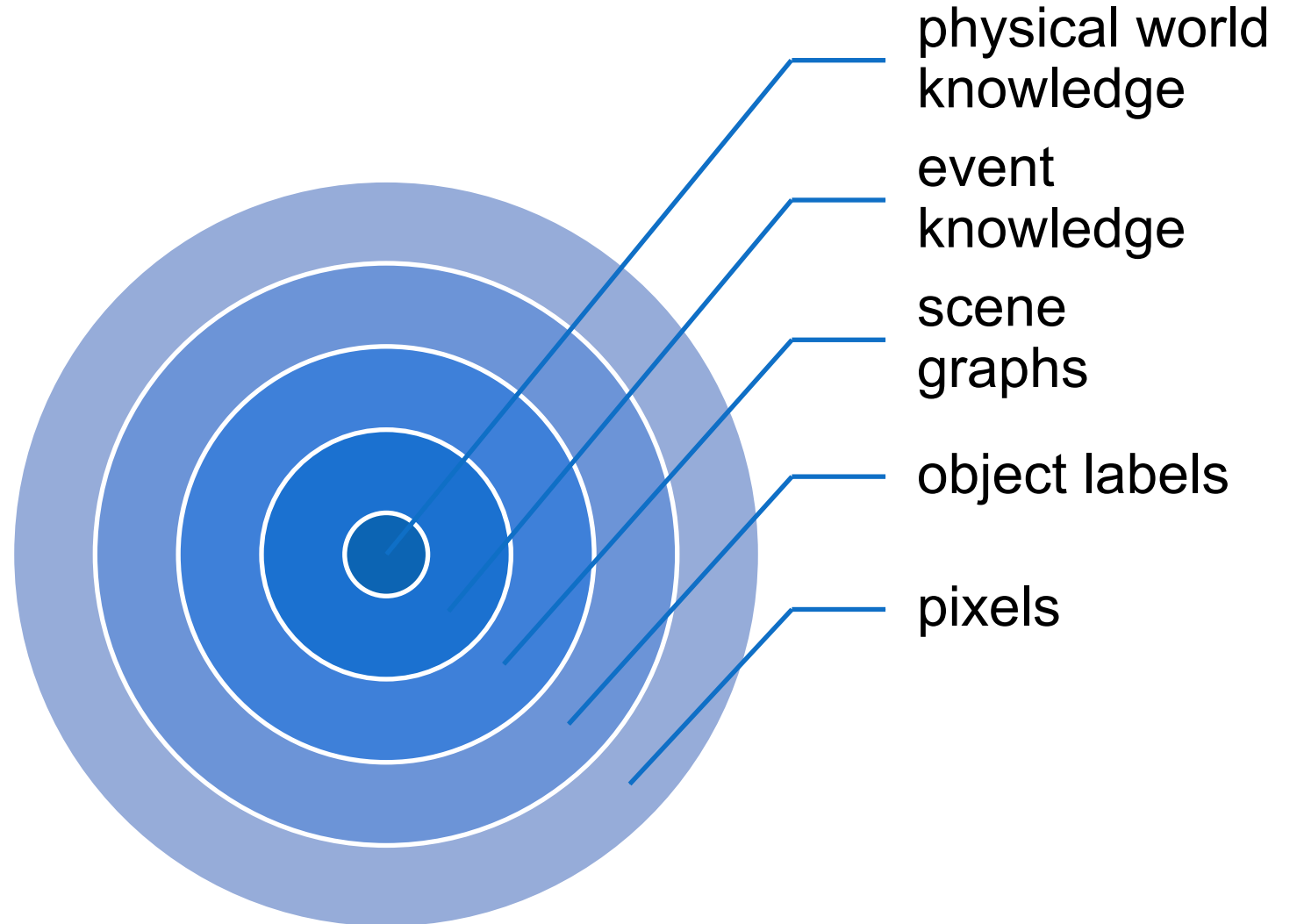
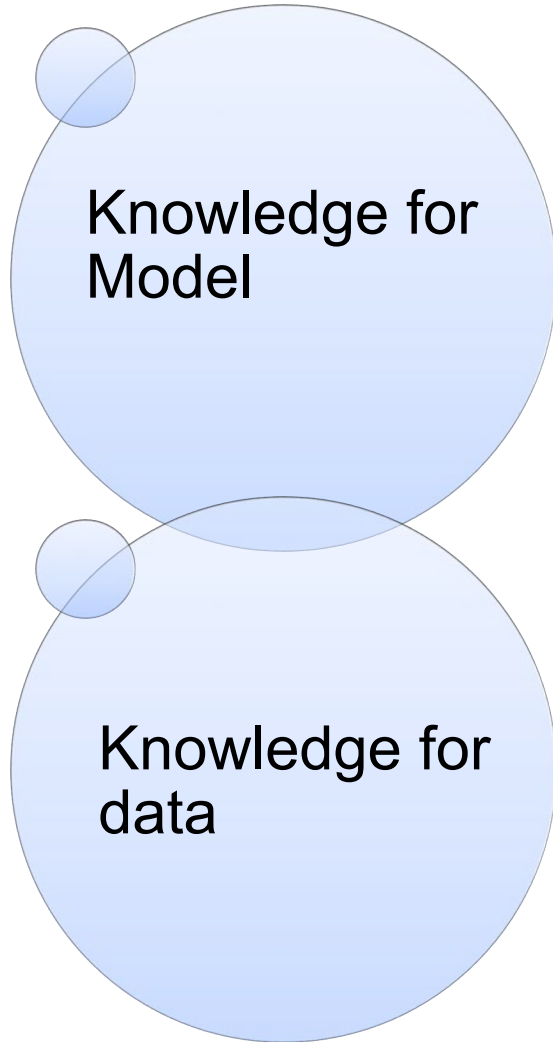


Video Event =
Status Changes of Arguments

Status Changes of an object =
Displacement (movement of bounding box)
+
Pixel Changes (intra-boundingbox changing)



Adding knowledge to pretraining models



What is embodied AI?



Aiming at the creation of an **embodied agent** (e.g., a robot) which learns, through **interaction and exploration**, to creatively solve challenging tasks within its environment.



- See: perceive their environment through vision or other senses.
- Talk: hold a natural language dialog grounded in their environment.
- Listen: understand and react to audio input anywhere in a scene.
- Act: navigate and interact with their environment to accomplish goals.
- Reason: consider and plan for the long-term consequences of their actions.

SayCan: Grounding Language in Affordances



SayCan: Grounding Language in Affordances



We need to ground the language model in tasks that are feasible within a specific real-world context.

I spilled my drink, can you help?

GPT3

You could try using a vacuum cleaner.

LaMDA

Do you want me to find a cleaner?

FLAN

I'm sorry, I didn't mean to spill it.

I spilled my drink, can you help?

LLM

"find a cleaner"

"find a sponge"

"go to the trash can"

"pick up the sponge"

"try using the vacuum"

Value Functions

"find a cleaner"

"find a sponge"

"go to the trash can"

"pick up the sponge"

"try using the vacuum"



SayCan

"find a cleaner"

"find a sponge"

"go to the trash can"

"pick up the sponge"

"try using the vacuum"



I would:

1. find a sponge
2. pick up the sponge
3. come to you
4. put down the sponge
5. done

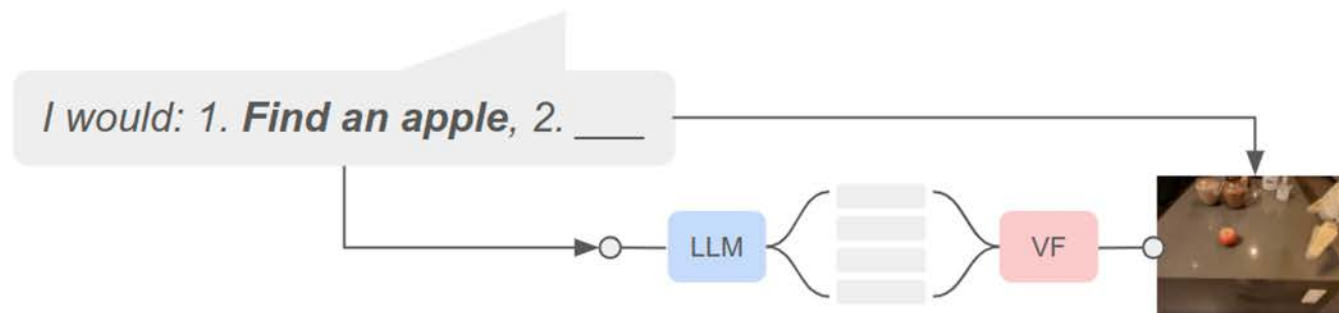
SayCan: Grounding Language in Affordances



Instruction Relevance with LLMs Combined Task Affordances with Value Functions



Note that the LM scores options instead of doing generation



Chain-of-thought prompting can improve examples with **negation**.

CHAIN-OF-THOUGHT PLANNING ROLLOUT WITH PaLM-SAYCAN.

Human: Can you bring a fruit-flavored drink without caffeine?

Explanation: The user has asked for a drink that is fruit-flavored and does not have caffeine, I will bring the lime soda.

Robot: 1. find a lime soda, 2. pick up the lime soda, 3. bring it to you, 4. put down the lime soda, 5. done

Human: Bring me the spicy chips.

Explanation: The user has asked for chips that are spicy, I will bring the Jalapeno chips.

Robot: 1. find the jalapeno chips, 2. pick up the jalapeno chips, 3. bring it to you, 4. put down the jalapeno chips, 5. done

Human: The snack is too small, bring me something more filling.

Explanation: The user has asked for a snack that is more filling, I will bring the multigrain chips.

Robot: 1. find the multigrain chips, 2. pick up the multigrain chips, 3. bring it to you, 4. put down the multigrain chips, 5. done

Table 4: Chain-of-thought planning rollout with PaLM-SayCan. The highlighted part is the chain of thought generated by PaLM-SayCan.

PaLM-E: An Embodied Multimodal Language Model



Encoding embodied observations as language tokens.

Mobile Manipulation



Human: Bring me the rice chips from the drawer. Robot: 1. Go to the drawers, 2. Open top drawer. I see ****. 3. Pick the green rice chip bag from the drawer and place it on the counter.

Visual Q&A, Captioning ...



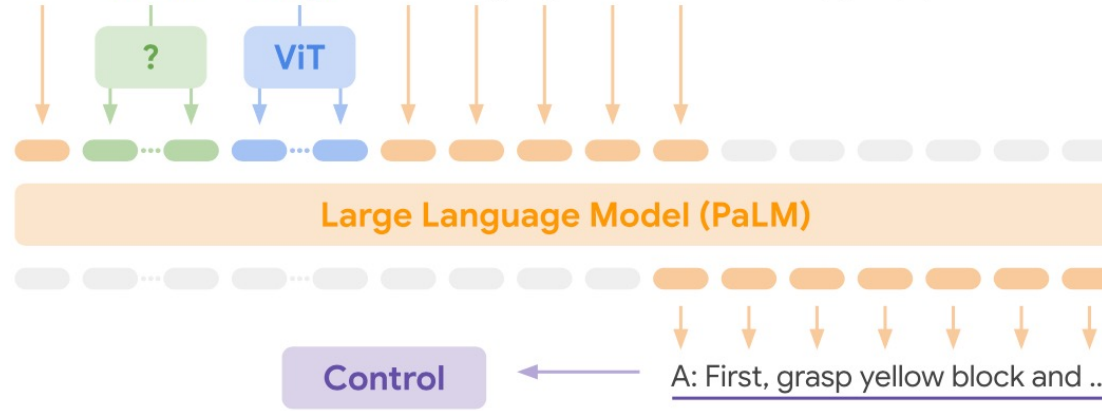
Given ****. Q: What's in the image? Answer in emojis. A: 🍏 🍌 🍇 🍏 🍏 🍏 🍒.



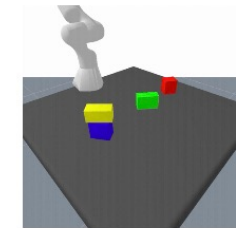
Describe the following ****:
A dog jumping over a hurdle at a dog show.

PaLM-E: An Embodied Multimodal Language Model

Given **<emb>** ... **** Q: How to grasp blue block? A: First, grasp yellow block

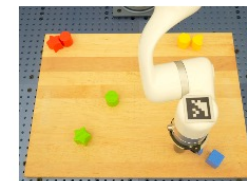


Task and Motion Planning



Given **<emb>** Q: How to grasp blue block? A: First grasp yellow block and place it on the table, then grasp the blue block.

Tabletop Manipulation

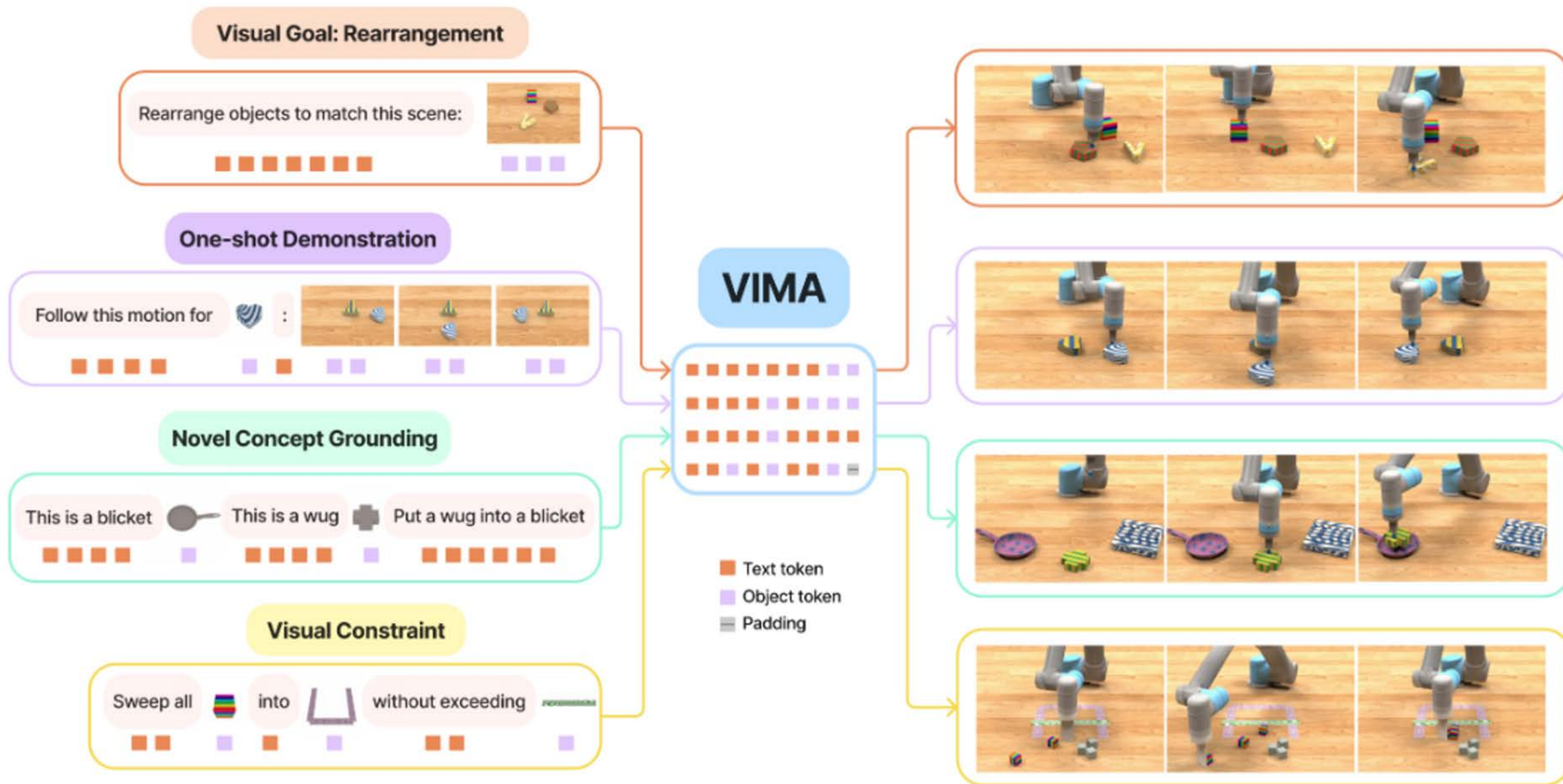


Given **** Task: Sort colors into corners. Step 1. Push the green star to the bottom left. Step 2. Push the green circle to the green star.

Language Only Tasks

Q: Miami Beach borders which ocean? A: Atlantic. Q: What is 372 x 18? A: 6696. Q: Write a Haiku about embodied LLMs. A: Embodied language. Models learn to understand. The world around them.

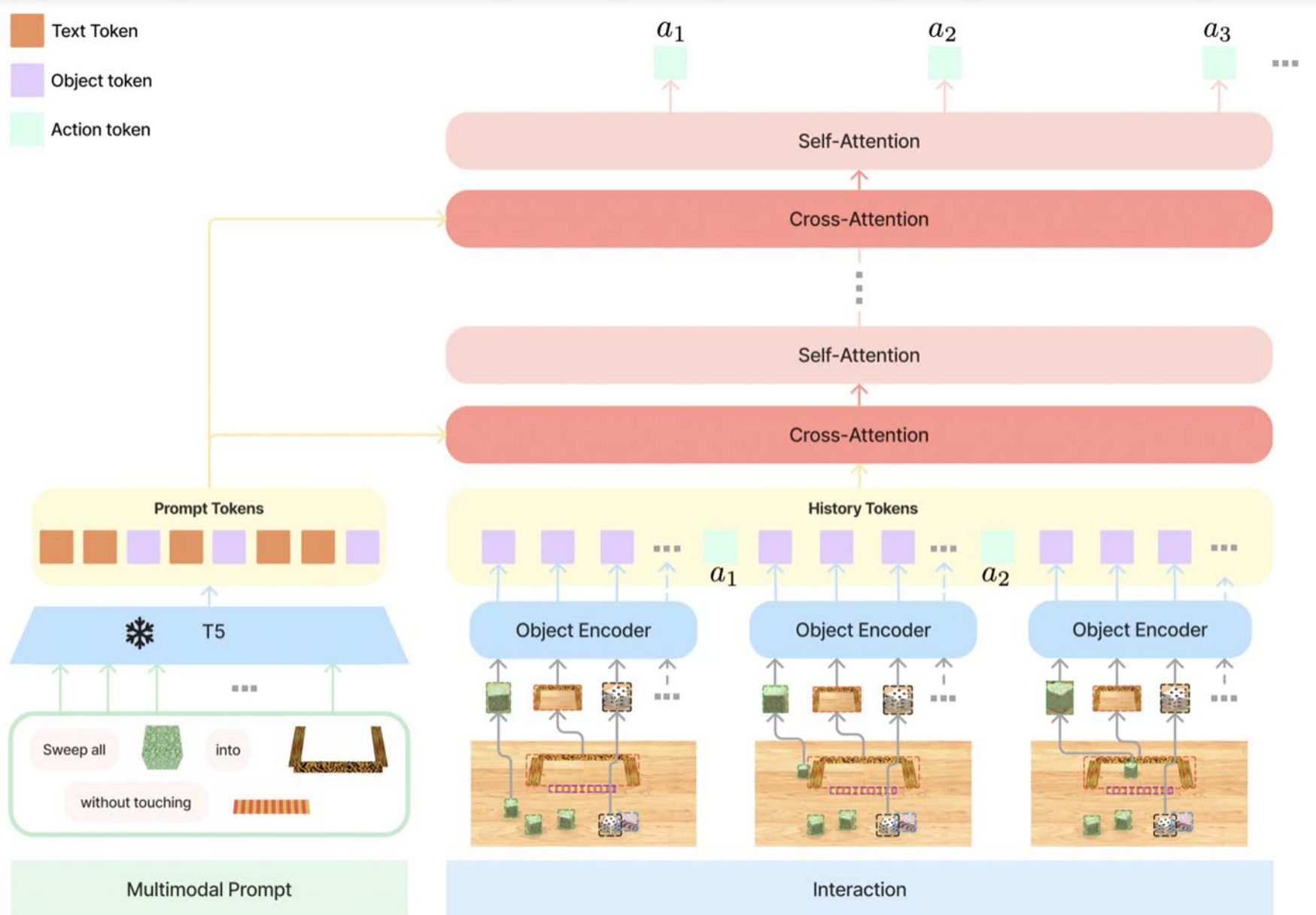
VIMA: Robot Manipulation with Multimodal Prompts



VIMA: Robot Manipulation with Multimodal Prompts



- Text Token
- Object token
- Action token

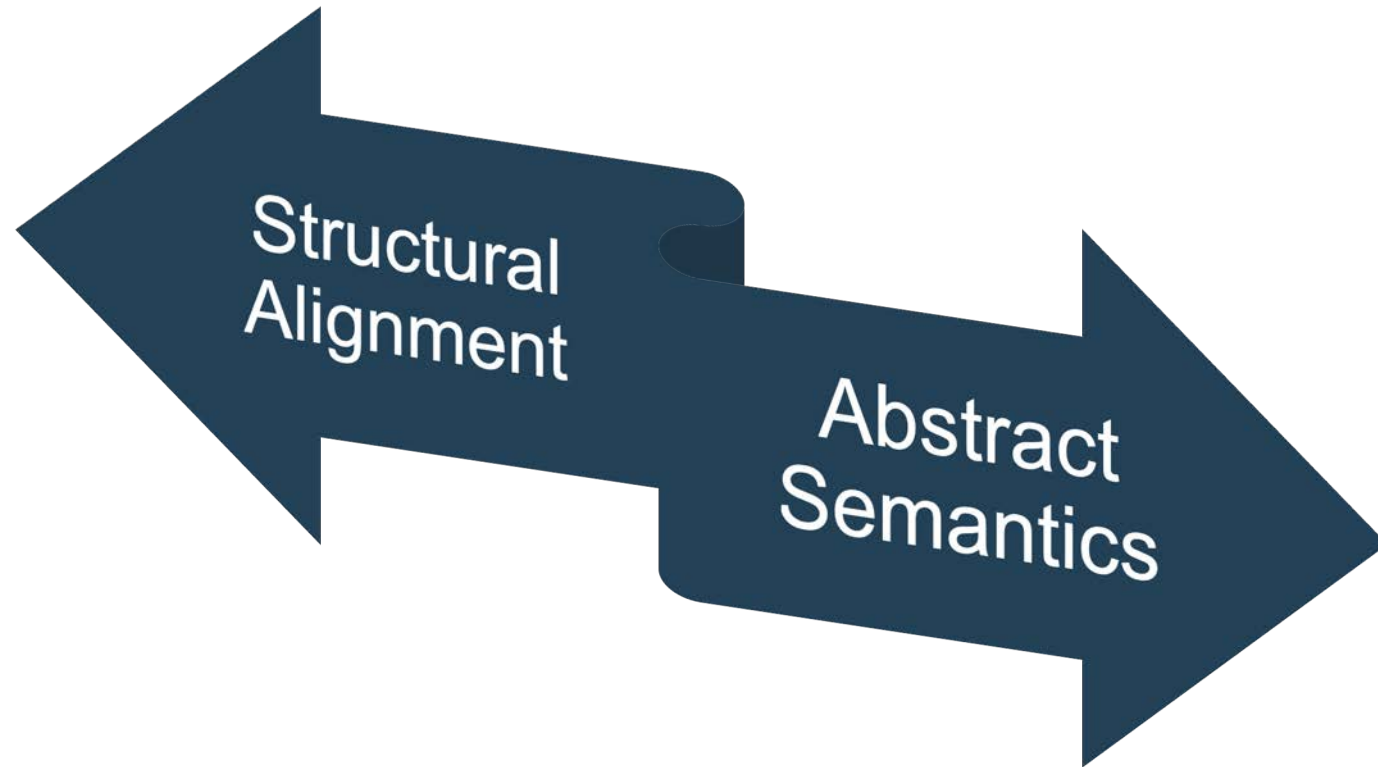


- *We note that this can only be achieved with both **cross-attention** and **object token sequence representation** — altering any component will degrade the performance significantly, especially in the low model capacity regime.*
- *The **data efficiency** can be attributed to VIMA's **object-centric representation**, which is less prone to overfitting than learning directly from pixels in the low-data regime.*

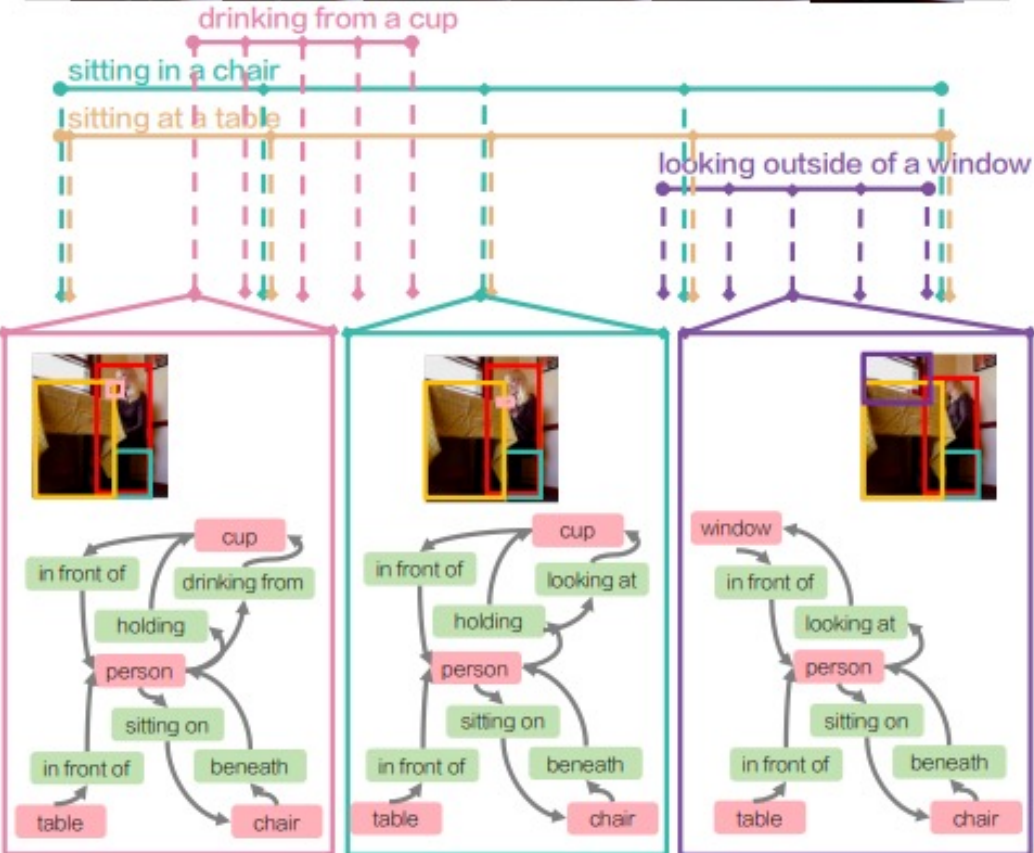
Future Challenges



- Structured: Capturing semantic structure
- Abstract: Understanding abstract and complicated concepts

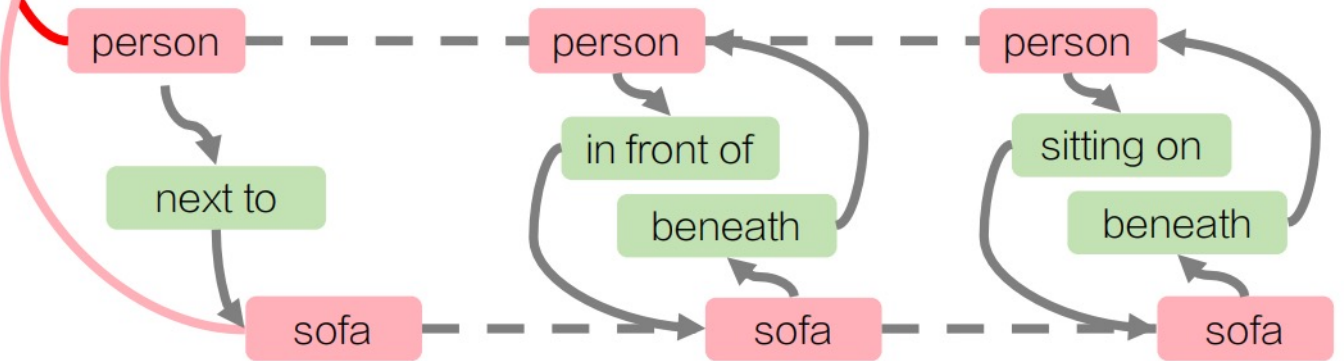
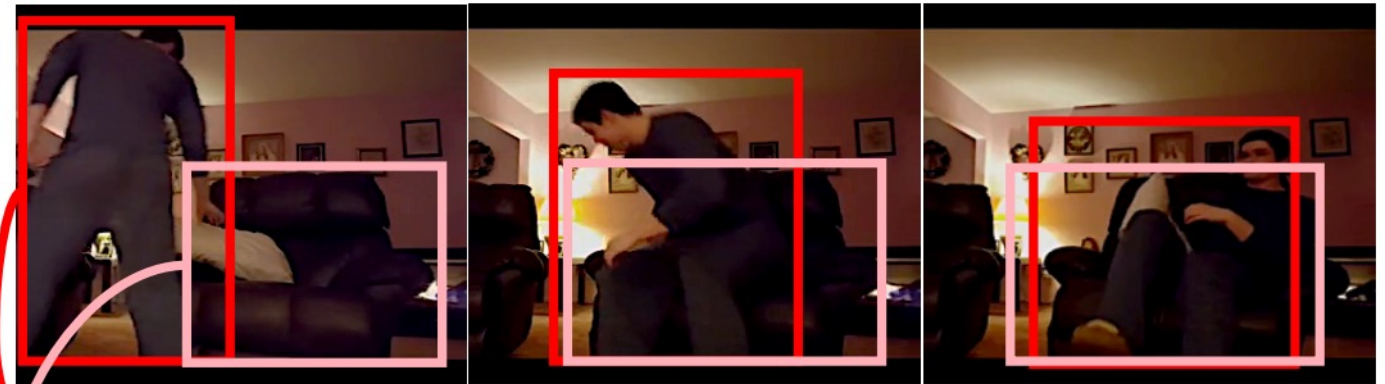


Future Direction 1: Structure-Aware Encoding



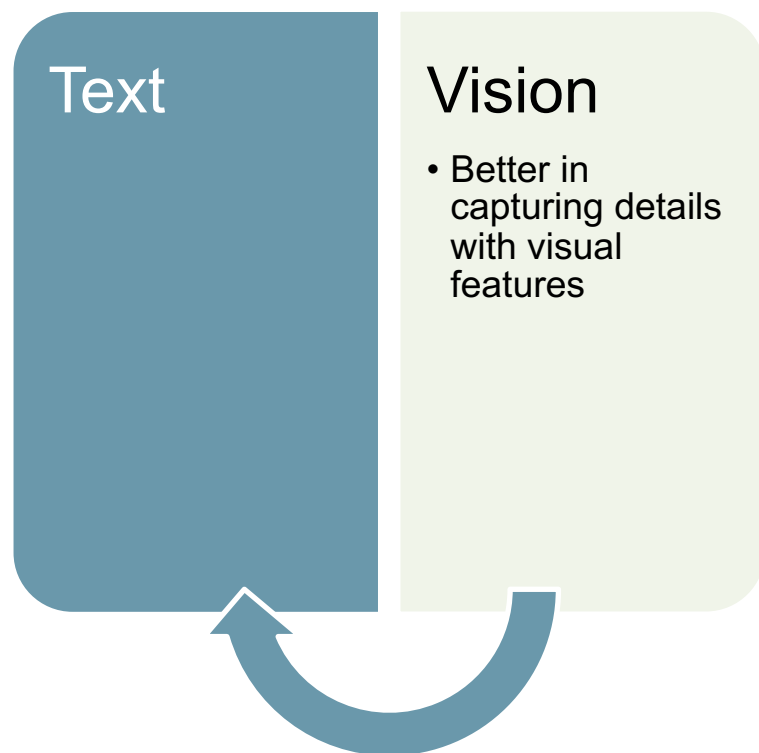
Action: "Sitting on a sofa"

time →



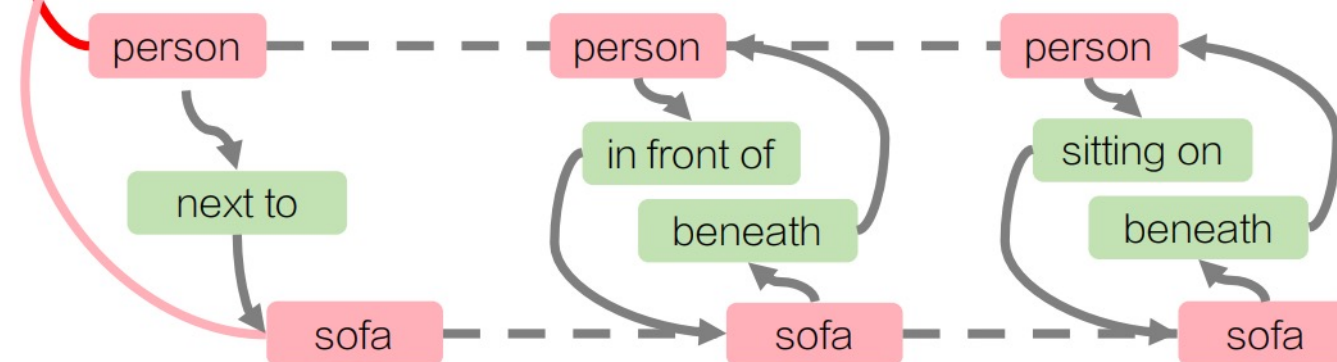
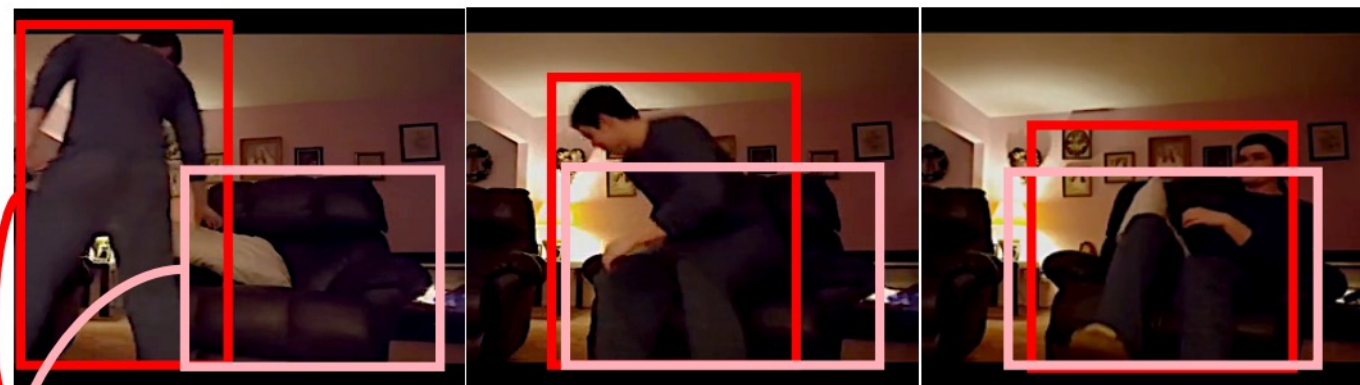
Spatio-temporal scene graphs

Future Direction 1: Structure-Aware Encoding



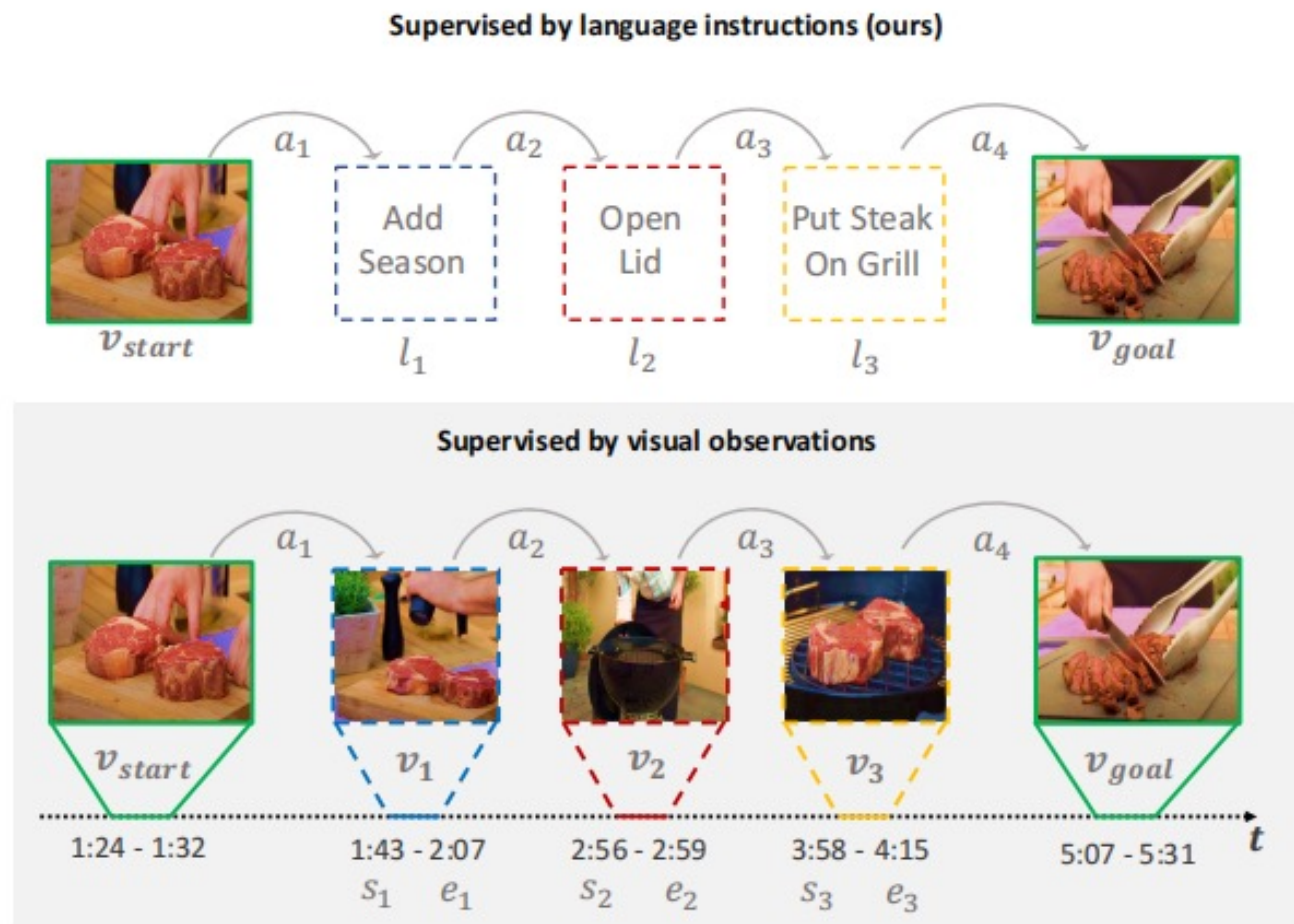
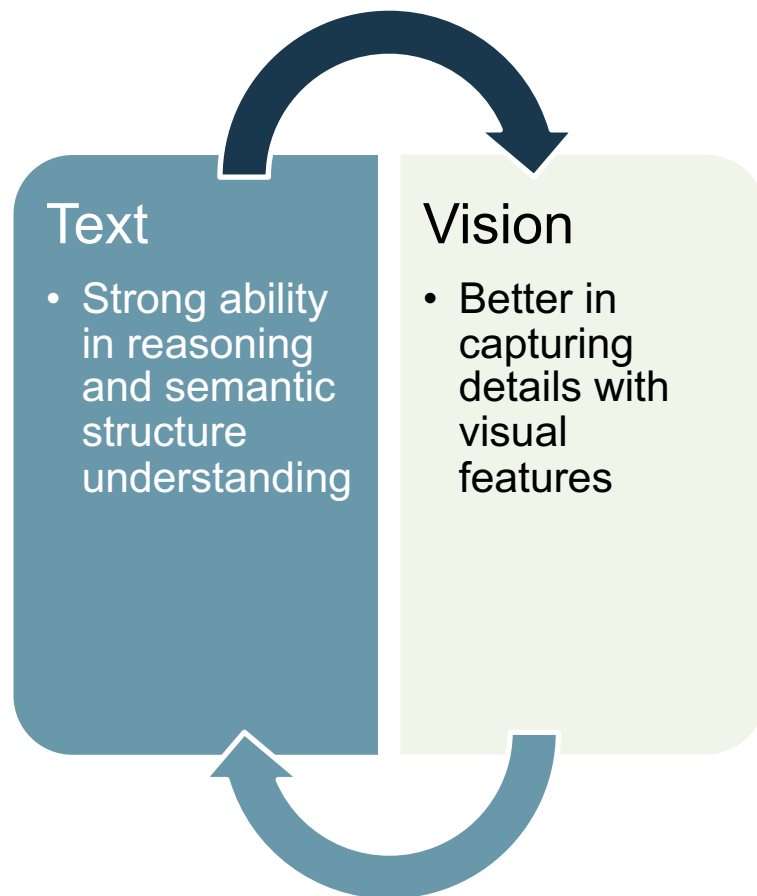
Action: "Sitting on a sofa"

time



Spatio-temporal scene graphs

Future Direction 1: Structure-Aware Encoding



Deep Semantic Understanding:

Discover knowledge (important information) that humans are actively seeking or communicating.

Future Direction 2: Abstract Semantics



Text generation paradigm (e.g., GPT-3) is taking over the NLP world.
But it is flat and surface-to-surface.

Bounded Knowledge

Short Context

Surface-to-Surface

Future Direction 2: Abstract Semantics



Text generation paradigm (e.g., GPT-3) is taking over the NLP world.
But it is flat and surface-to-surface.

Bounded Knowledge

Short Context

Surface-to-Surface

Surface → Deep

Concrete → Abstract

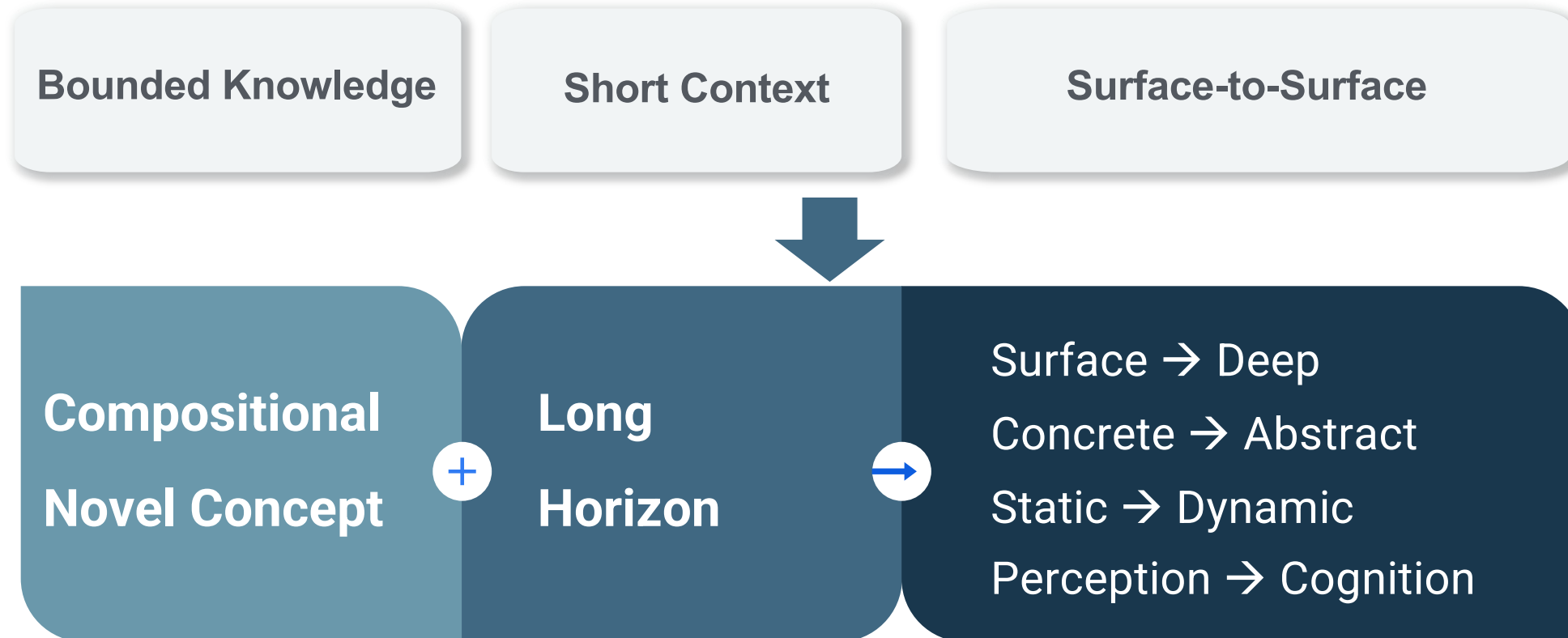
Static → Dynamic

Perception → Cognition

Future Direction 2: Abstract Semantics



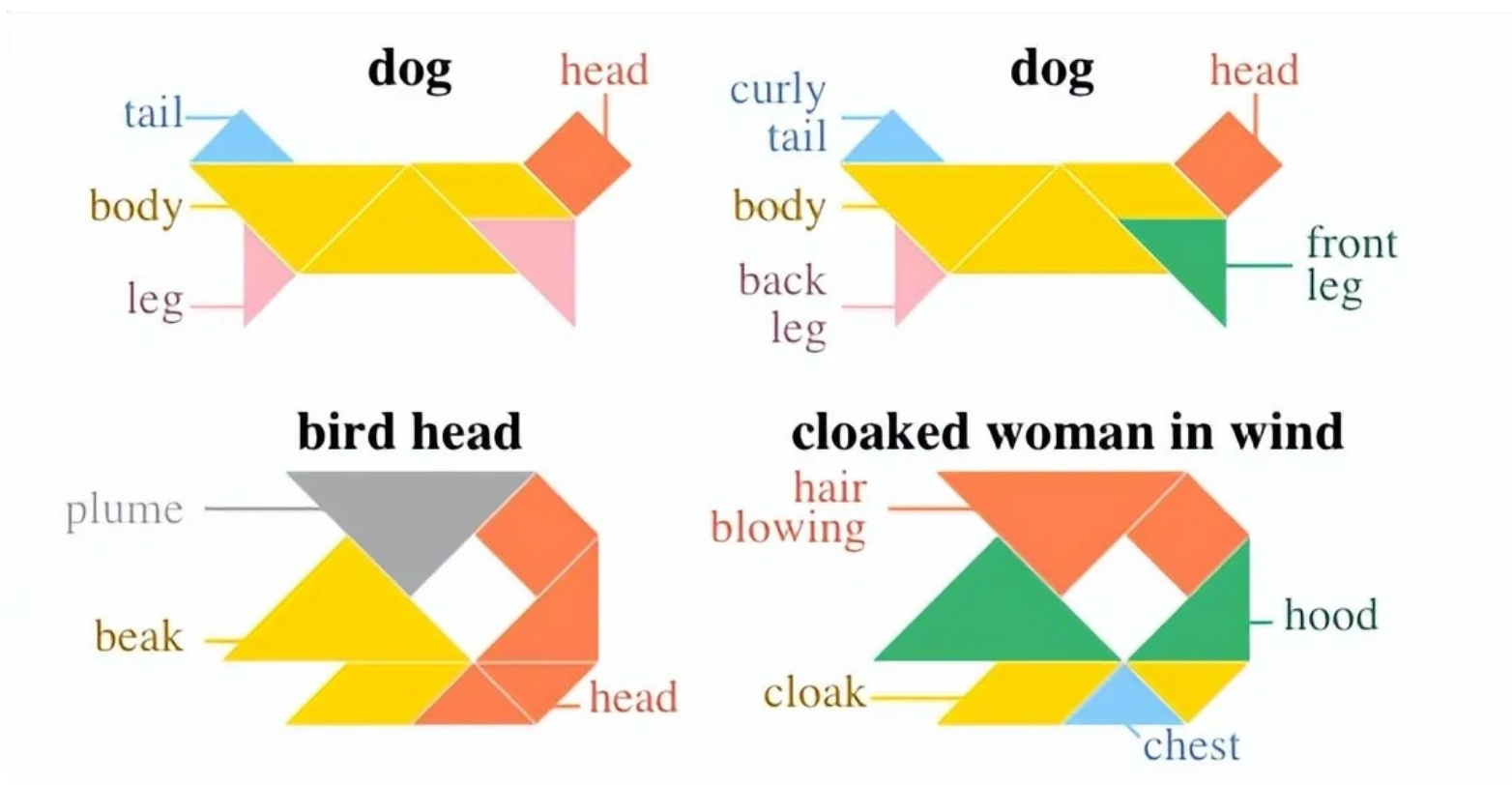
Text generation paradigm (e.g., GPT-3) is taking over the NLP world.
But it is flat and surface-to-surface.



Future Direction 2: Abstract Semantics



Abstract



Future Direction 2: Abstract Semantics



Abstract



Love



Happiness

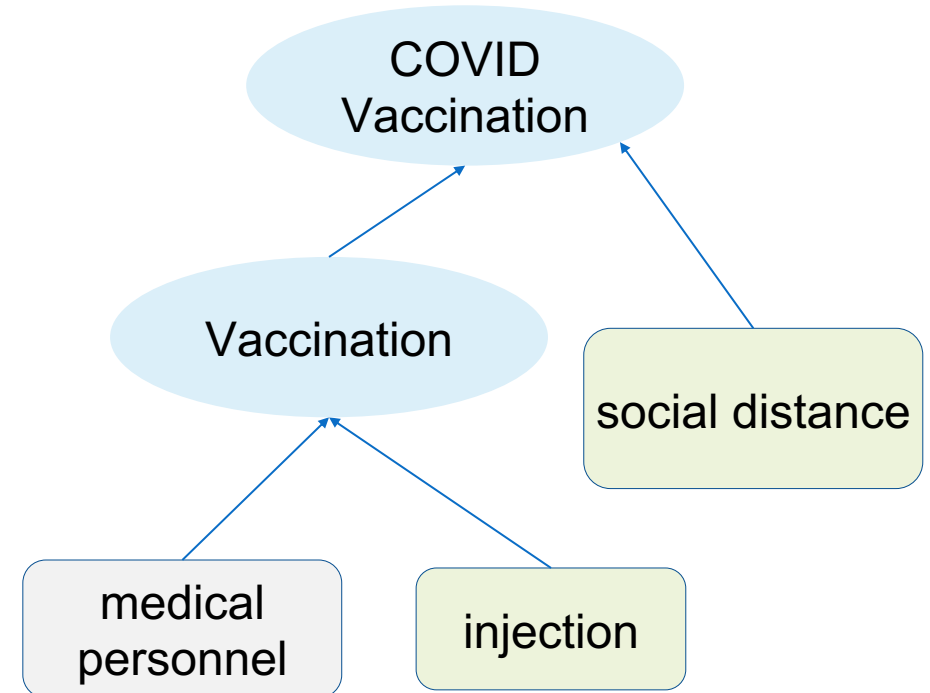


Emotion \leftrightarrow Music

Future Direction 2: Compositional Semantics



Compositional



Future Direction 2: Abstract Semantics

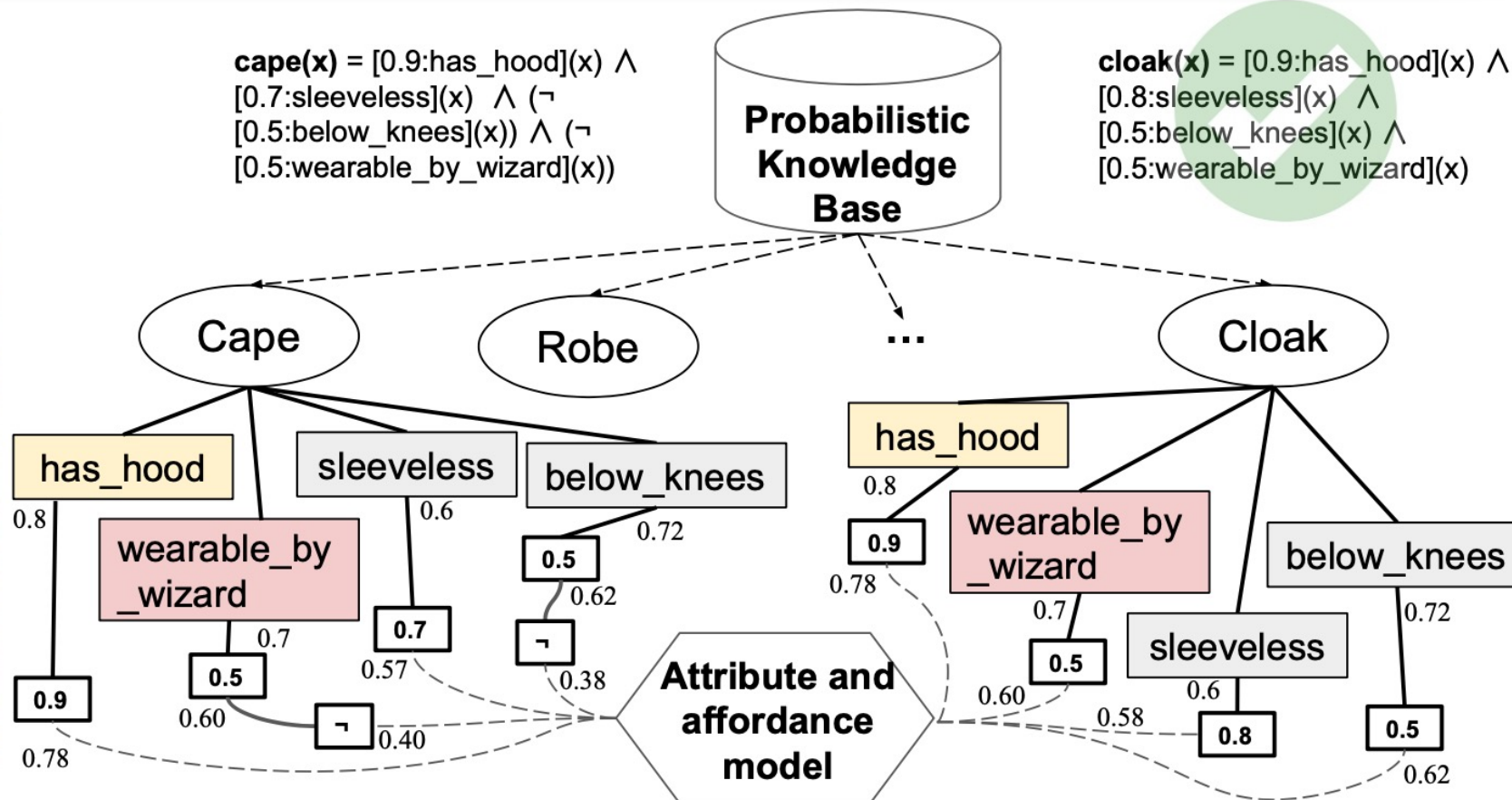
Reasoning & Explainability

Image and object bounding box



$\text{cape}(x) = [0.9:\text{has_hood}](x) \wedge$
 $[0.7:\text{sleeveless}](x) \wedge (\neg$
 $[0.5:\text{below_knees}](x)) \wedge (\neg$
 $[0.5:\text{wearable_by_wizard}](x))$

$\text{cloak}(x) = [0.9:\text{has_hood}](x) \wedge$
 $[0.8:\text{sleeveless}](x) \wedge$
 $[0.5:\text{below_knees}](x) \wedge$
 $[0.5:\text{wearable_by_wizard}](x)$



Alternative Interpretation & Propaganda



 New York Post Lean Right

Zelensky describes Ukraine war as 'horror film' at Venice Film Festival opening

 Scroll India Lean Left

Ukrainian President addresses Venice Film Festival, calls war 'drama based on real-life'

Information Surgery

Video Framing

Propagation Path

Propaganda Detection

Misinformation Detection

Reframing

LLMs have **good compositional ability** from low-level observations to high level concepts.



Z

We observe the following features in an image: "Dog standing, looking back," "Man in camouflage uniform holding baby in warm clothes," "Rubble and overturned vehicle in background," "Outdoor scene with water, mud, and man-made objects".

what is the activity happening in the image?



Based on the provided features, it is difficult to determine the exact activity happening in the image. However, we can make some observations and assumptions:

- There is a dog standing and looking back, suggesting that the dog may be accompanying the man and baby.
- The man is wearing a camouflage uniform and holding a baby in warm clothes, which suggests that he may be a military or security personnel.
- There is rubble and an overturned vehicle in the background, suggesting that there may have been some kind of conflict or disaster in the area.
- The outdoor scene includes water, mud, and man-made objects, which could suggest that the area is a construction site, a disaster zone, or a military training ground.

Based on these observations, it is possible that the man is a military or security personnel who is holding the baby while on duty in a conflict or disaster zone. The dog may be accompanying them on their mission. However, without more context, it is impossible to determine the exact activity happening in the image.

- Evaluation:
 - Is the performance boost because of adding knowledge?
- Explanability and Reasoning
 - How can large models leverage knowledge?
- Continuous learning for large models
 - How can large models continue learning more complicated concepts?
- Bias in learning knowledge
 - Debiasing with knowledge guidance

Jun 2023

CVPR Tutorials

Knowledge-Driven Vision-Language Encoding

CVPR

Commonsense Knowledge in V+L Pretraining

Knowledge-Driven Vision-Language Pretraining (Part III)

Manling Li
UIUC

manling2@illinois.edu



Northwestern
University



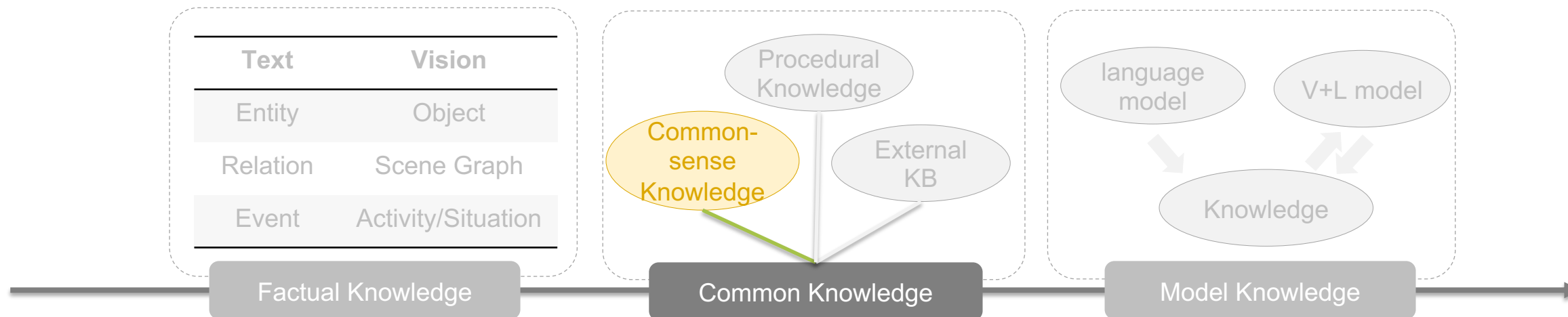
COLUMBIA
UNIVERSITY

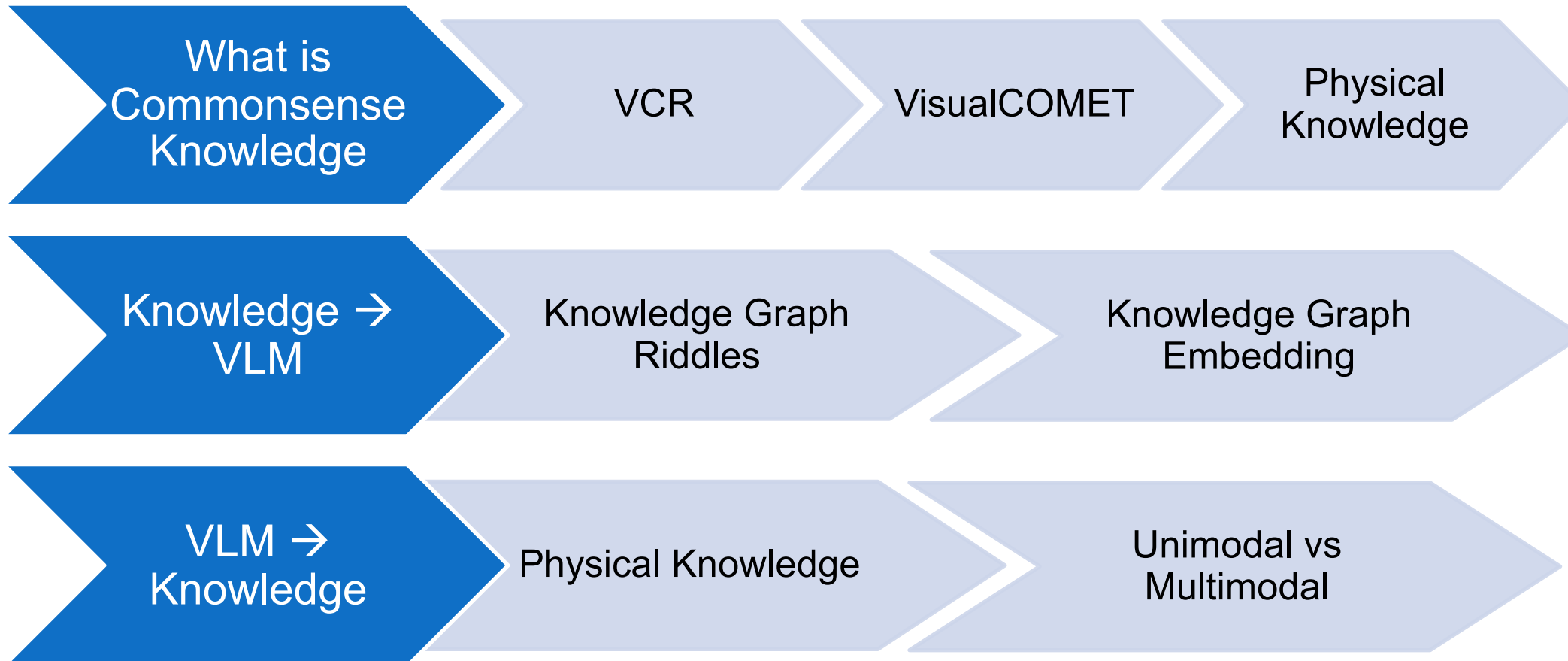


Commonsense Knowledge



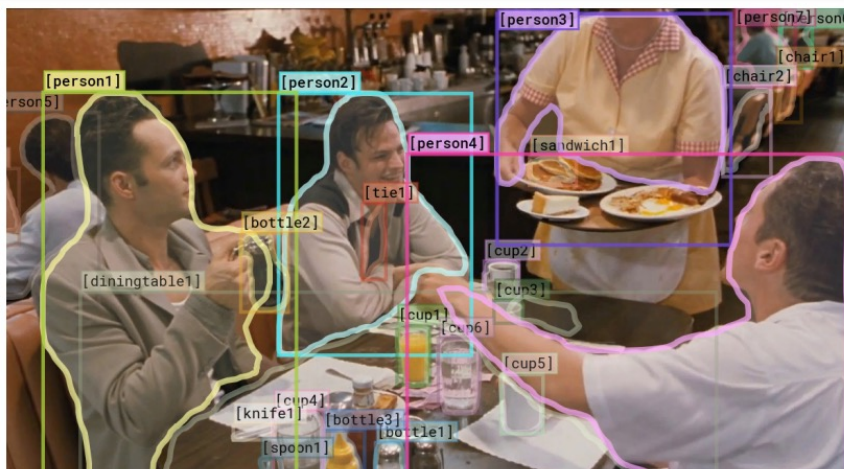
Commonsense Knowledge is the basic facts and behaviors of the everyday world.





Part 1: What is Visual Commonsense Knowledge?

Visual Commonsense Reasoning (VCR): From Recognition to Cognition



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.

I chose a) 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 and both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.



How did [person2] get the money that's in front of her?

- a) [person2] is selling things on the street.
- b) [person2] earned this money playing music.
- c) She may work jobs for the mafia.
- d) She won money playing poker.

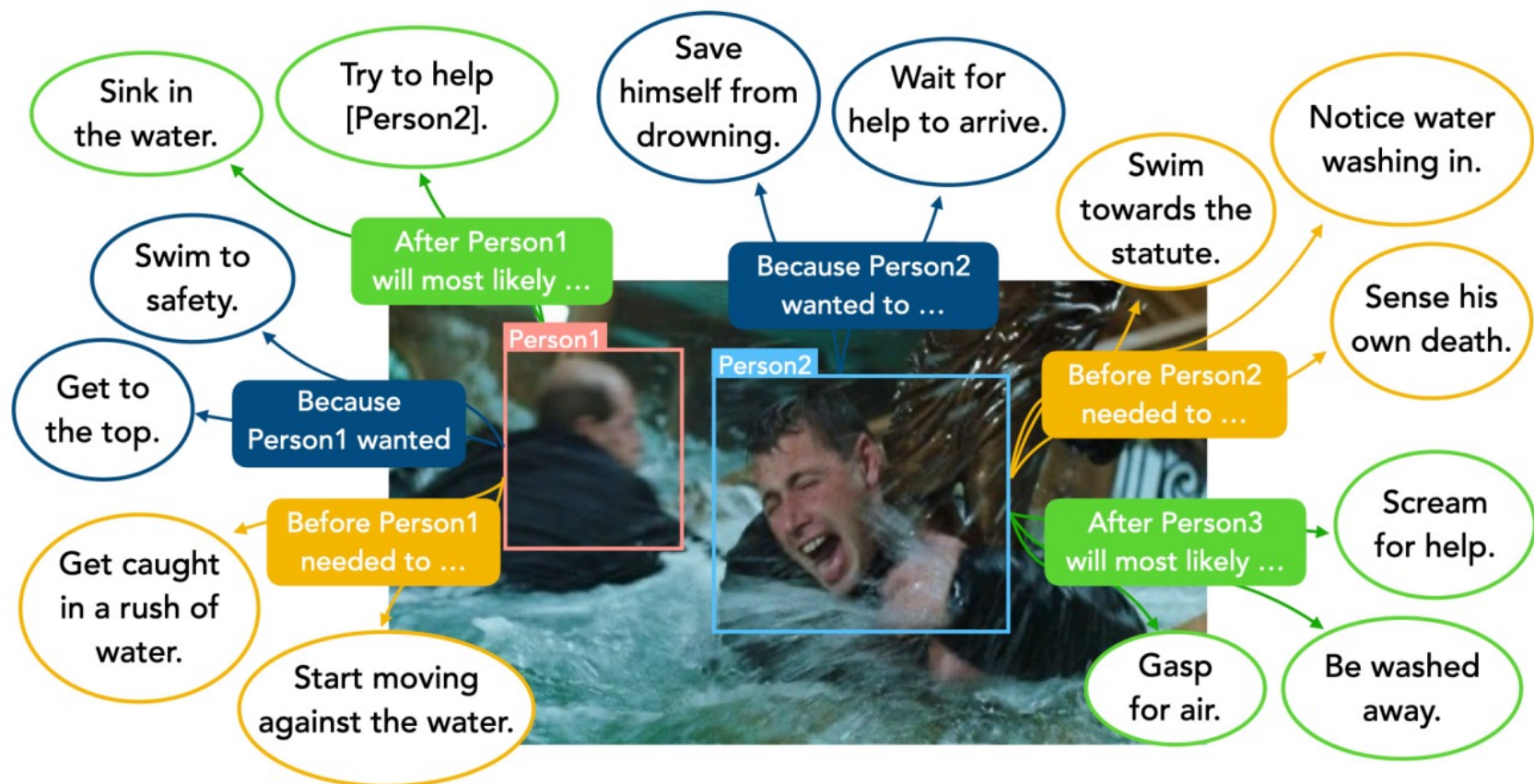
I chose b) because...

- a) She is playing guitar for money.
- b) [person2] is a professional musician in an orchestra.
- c) [person2] and [person1] are both holding instruments, and were probably busking for that money.
- d) [person1] is putting money in [person2]'s tip jar, while she plays music.

Visual Commonsense Knowledge



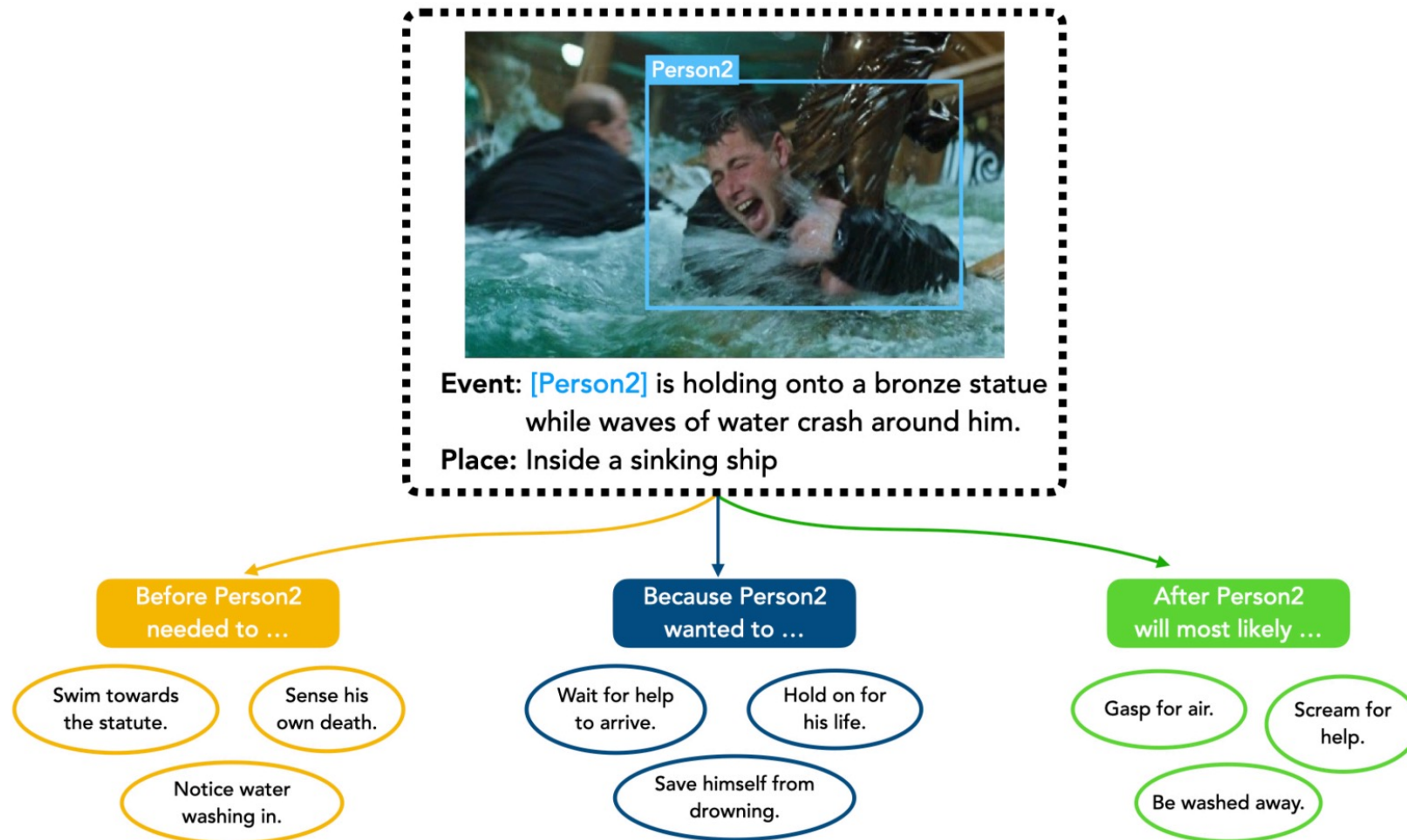
VisualCOMET: Cognitive Image Understanding via Visual Commonsense Graphs



Visual Commonsense Knowledge



VisualCOMET Task Formulation: Generate the entire visual commonsense graph



Visual Commonsense Knowledge



Large Dataset Collection: There are in total 139,377 distinct Visual Commonsense Graphs over 59,356 images involving 1,465,704 commonsense inferences.

	Train	Dev	Test	Total
# Images/Places	47,595	5,973	5,968	59,356
# Events at Present	111,796	13,768	13,813	139,377
# Inferences on Events Before	467,025	58,773	58,413	584,211
# Inferences on Events After	469,430	58,665	58,323	586,418
# Inferences on Intents at Present	237,608	28,904	28,568	295,080
# Total Inferences	1,174,063	146,332	145,309	1,465,704

Physical Commonsense Knowledge can be learned via natural language.

a. Shape, Material, and Purpose

- [Goal] Make an outdoor pillow
- [Sol1] Blow into a **tin can** and tie with rubber band ✗
- [Sol2] Blow into a **trash bag** and tie with rubber band ✓
- [Goal] To make a hard shelled taco,
- [Sol1] put seasoned beef, cheese, and lettuce **onto** the hard shell ✗
- [Sol2] put seasoned beef, cheese, and lettuce **into** the hard shell ✓
- [Goal] How do I find something I lost on the carpet?
- [Sol1] Put a **solid seal** on the end of your vacuum and turn it on ✗
- [Sol2] Put a **hair net** on the end of your vacuum and turn it on. ✓

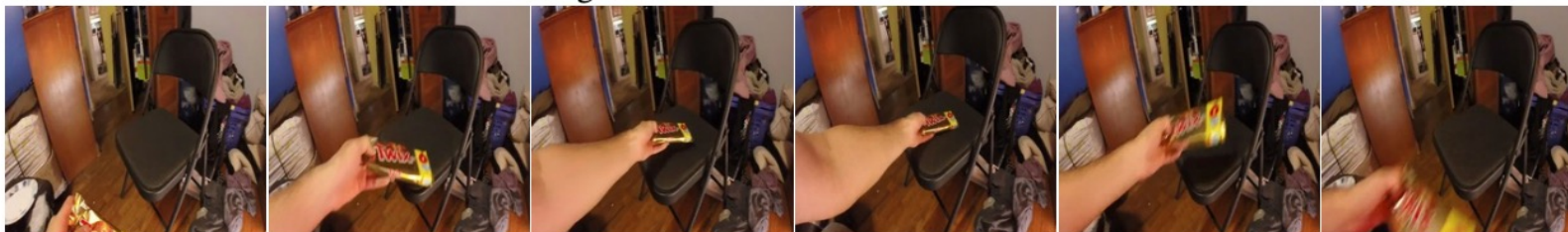
b. Commonsense Convenience

- [Goal] How to make sure all the clocks in the house are set accurately?
- [Sol1] Get a solar clock for a reference and place it just outside a window that gets lots of sun. Use a system of call and response once a month, having one person stationed at the solar clock who yells out the correct time and have another person move to each of the indoor clocks to check if they are showing the right time. Adjust as necessary. ✗
- [Sol2] Replace all wind-ups with digital clocks. That way, you set them once, and that's it. Check the batteries once a year or if you notice anything looks a little off. ✓

The “Something Something” Dataset



Putting a white remote into a cardboard box



Pretending to put candy onto chair



Pushing a green chilli so that it falls off the table



Moving puncher closer to scissor

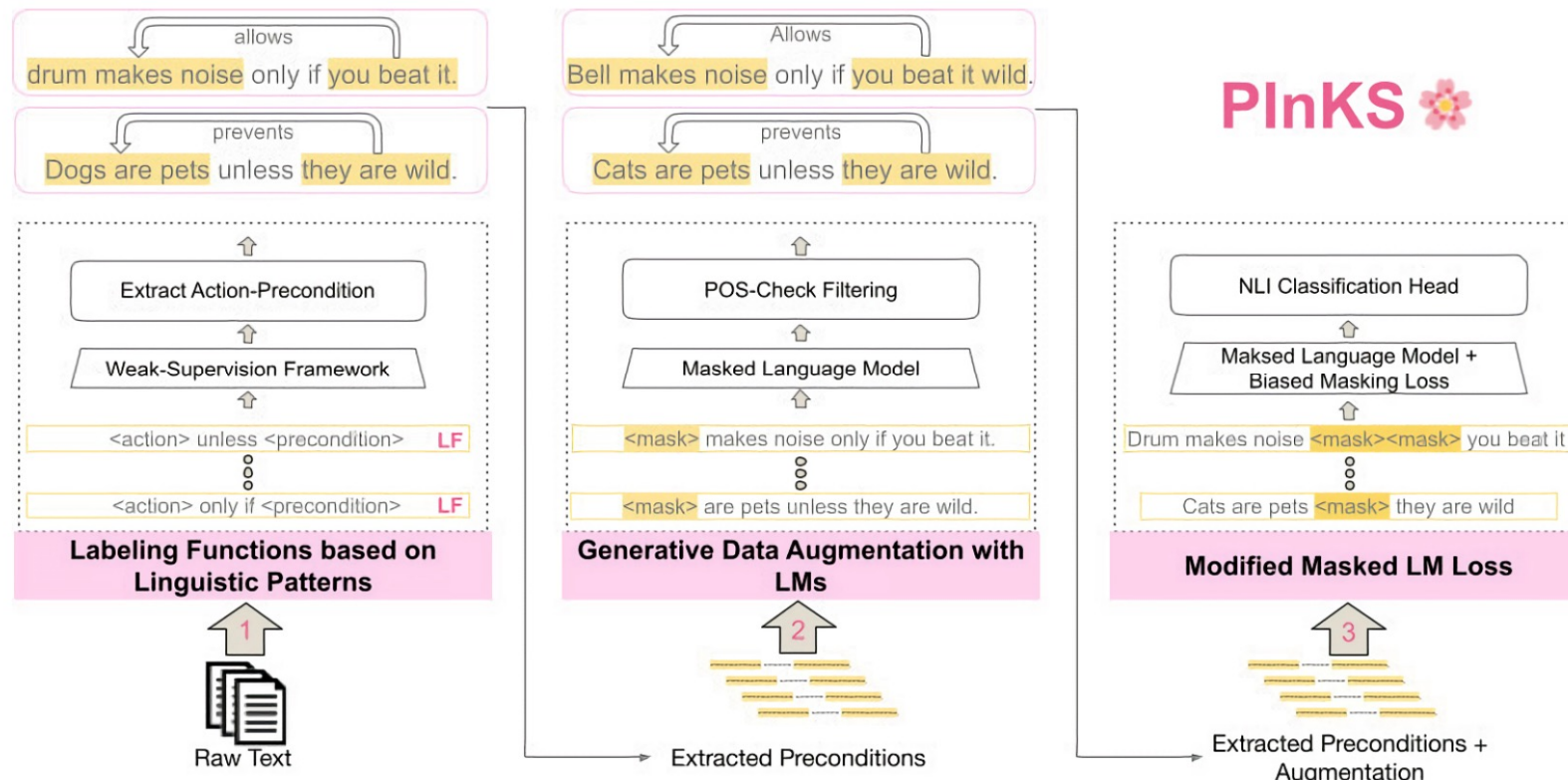
10 selected classes

- Dropping [something]
- Moving [something] from right to left
- Moving [something] from left to right
- Picking [something] up
- Putting [something]
- Poking [something]
- Tearing [something]
- Pouring [something]
- Holding [something]
- Showing [something] (almost no hand)

PIInKS: Preconditioned Commonsense Inference

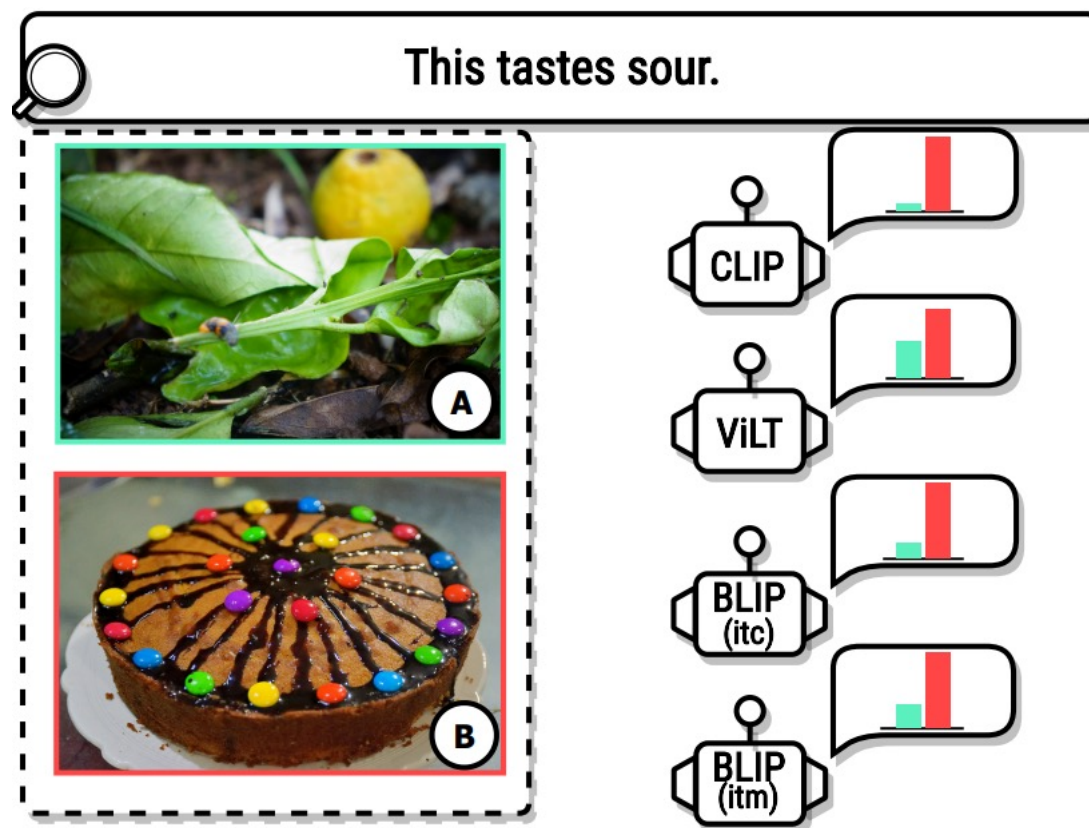


Text	Label	Action	Precondition
A drum makes noise only if you beat it.	Allow	A drum makes noise	you beat it.
Your feet might come into contact with something if it is on the floor.	Allow	Your feet might come into contact with something	it is on the floor.
Pears will rot if not refrigerated	Prevent	Pears will rot	refrigerated
Swimming pools have cold water in the winter unless they are heated.	Prevent	Swimming pools have cold water in the winter	they are heated.



Part 2: How can commonsense knowledge be learned via V+L pretraining?

Current V+L models lack abilities to capture commonsense knowledge:



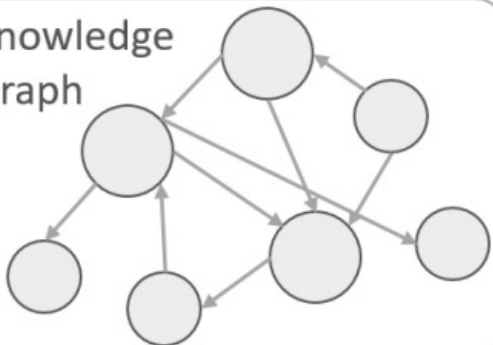
DANCE: Data Augmentation with kNowledge graph linearization for CommonsenseE capability

Original image-text pair



*A **cat** with a **box** in an **office**.*

Knowledge
Graph



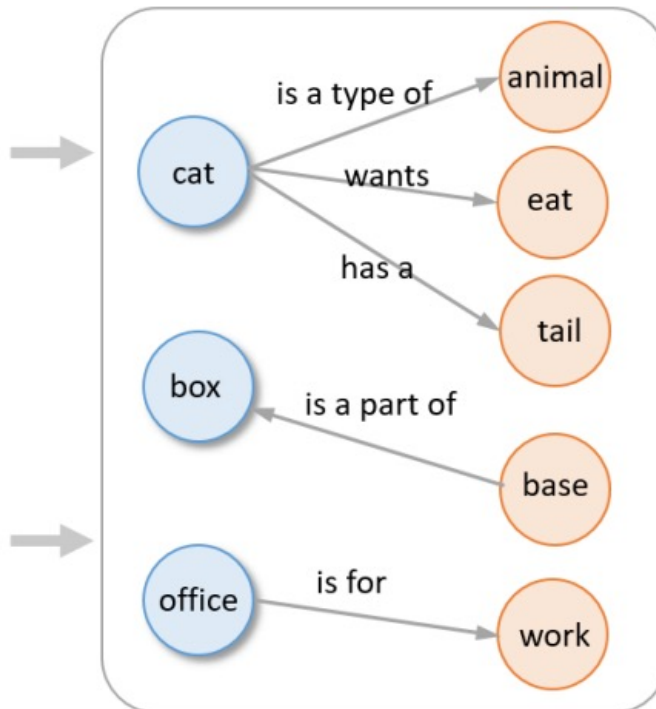
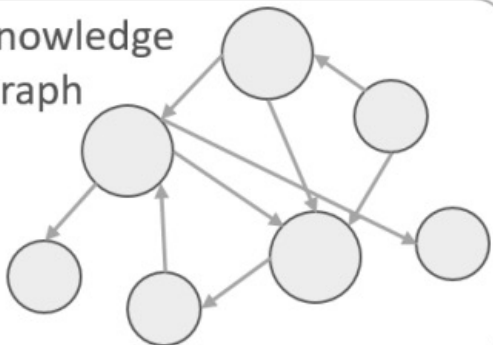
DANCE: Data Augmentation with kNowledge graph linearization for CommonsenseE capability

Original image-text pair



*A **cat** with a **box** in an **office**.*

Knowledge Graph



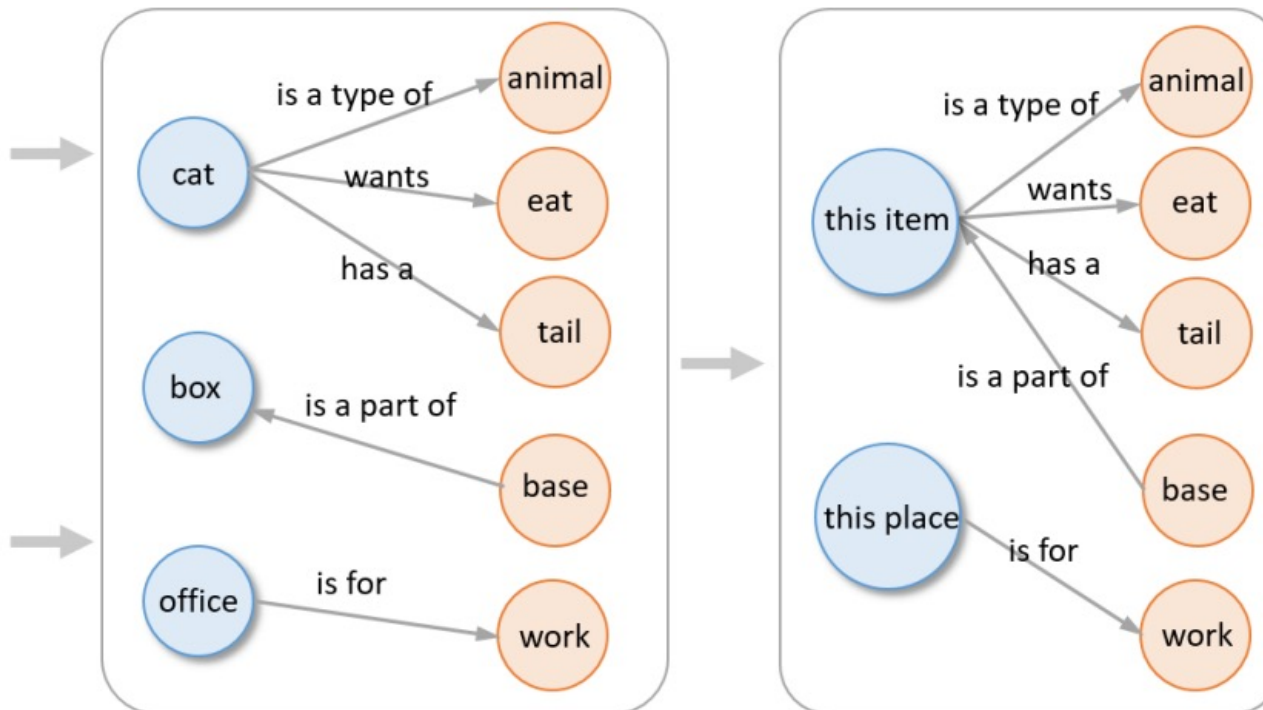
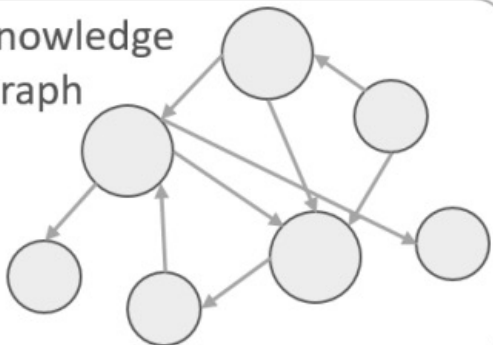
DANCE: Data Augmentation with kNowledge graph linearization for CommonsenseE capability

Original image-text pair




*A **cat** with a **box** in an **office**.*

Knowledge Graph

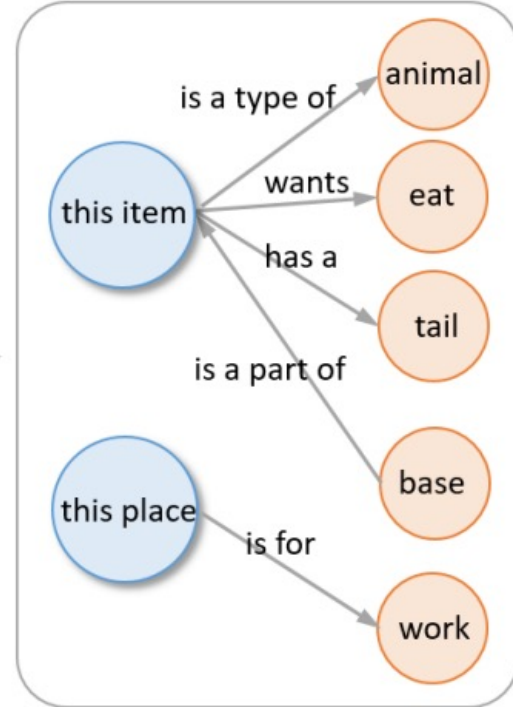
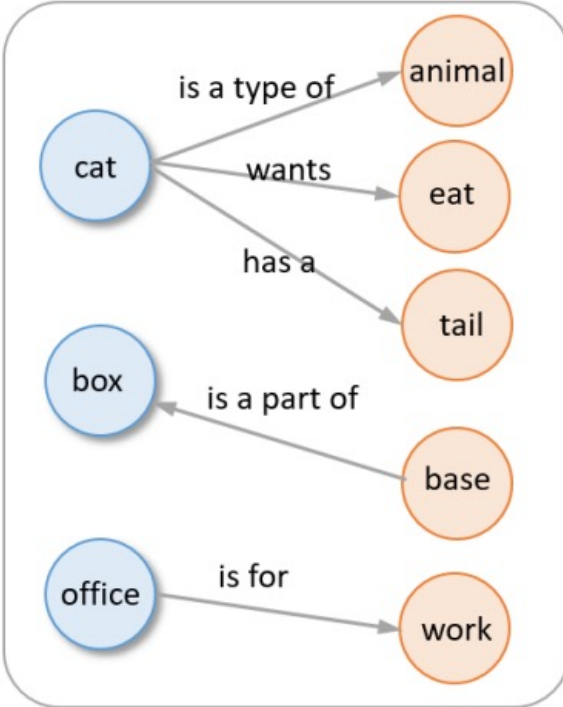
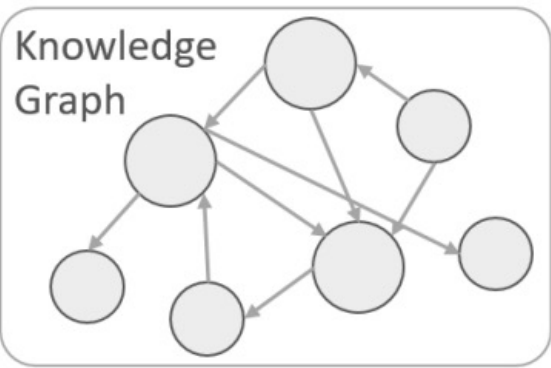


DANCE: Data Augmentation with kNowledge graph linearization for CommonsenseE capability

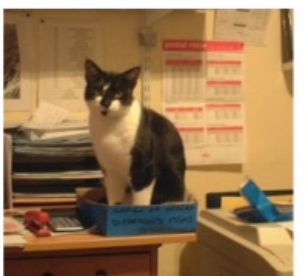
Original image-text pair



*A **cat** with a **box** in an **office**.*



Augmented image-riddle pairs



***This item** is a type of animal.*

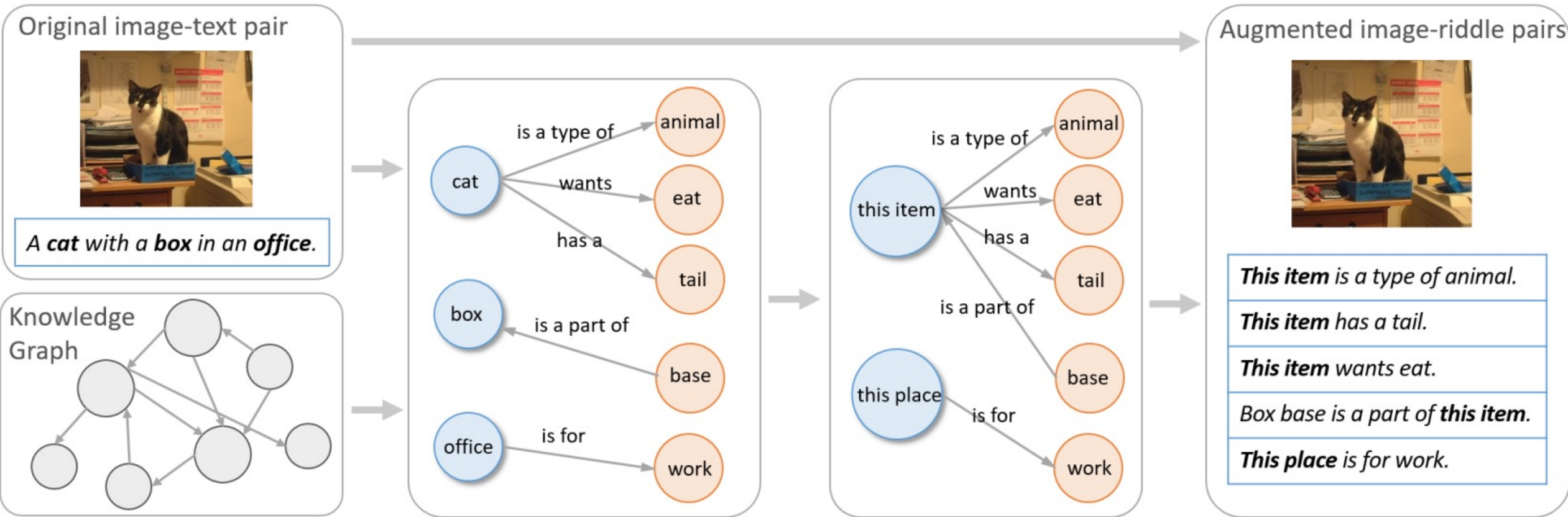
***This item** has a tail.*

***This item** wants eat.*

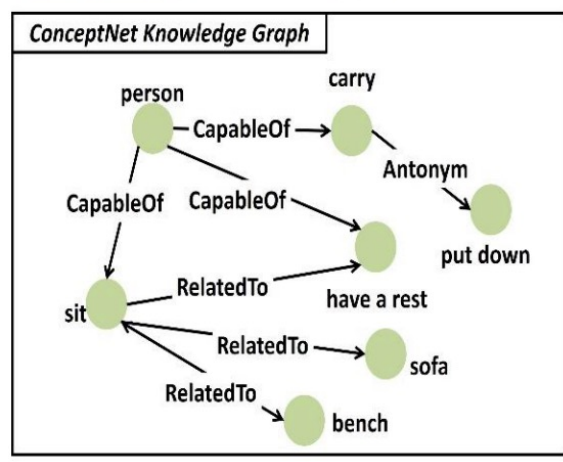
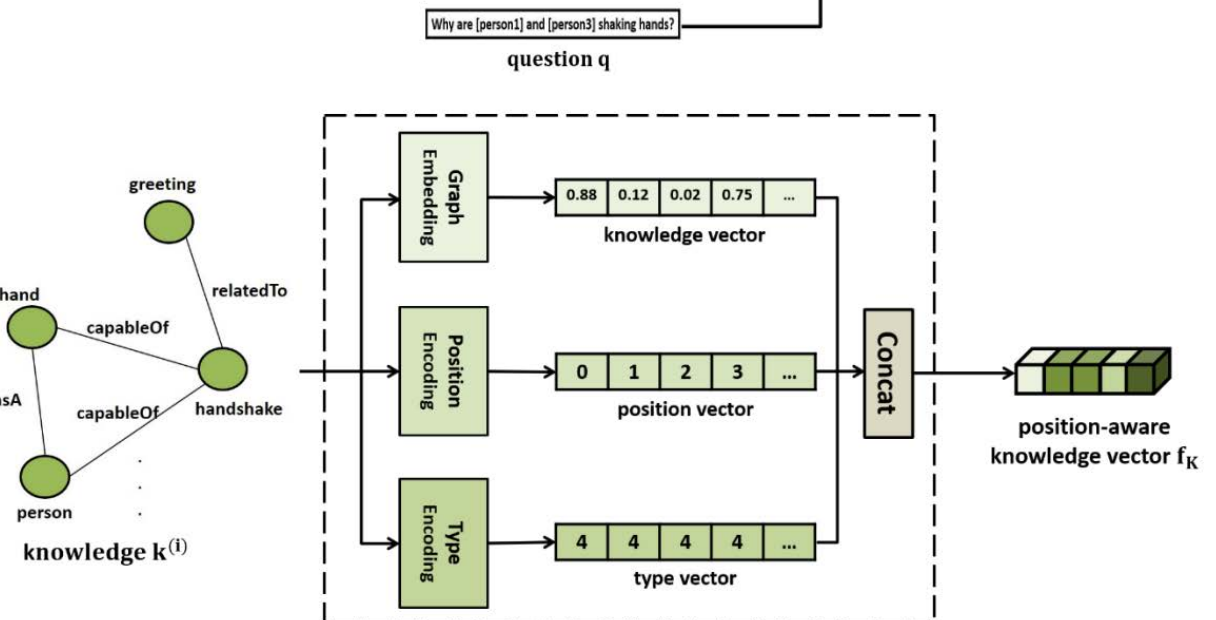
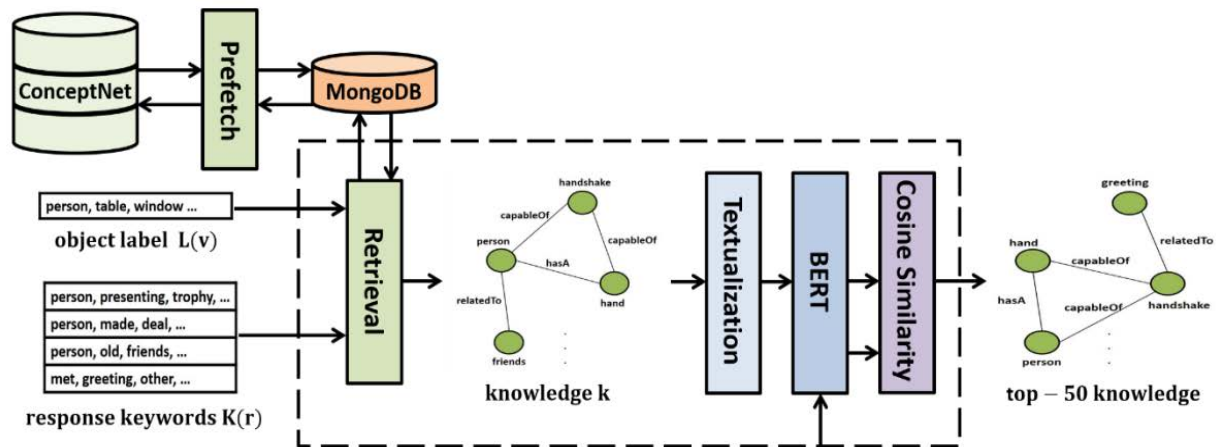
*Box base is a part of **this item**.*

***This place** is for work.*

DANCE: Data Augmentation with kNowledge graph linearization for CommonsenseE capability



Vision–Language Knowledge Co-Embedding



What is [person1] doing now?








- (a) He is riding on a horse. 0.1%
- (b) [person1] is being handed a fork to eat his food. 0.2%
- (c) He is sitting and resting after carrying a heavy item into the room. 98%**
- (d) He is in the process of turning to look at [person1]. 1.7%

Because...









- (a) He is sitting next to a cart with a large object on it. He is tired from pulling the kart because it is heavy. 0.2%
- (b) He is using a brace, used to carry things. 14%
- (c) He has brought many items into the fort which make him feel at ease. 0.8%
- (d) He had a bench and has put it down. 85%**

Borrowing Knowledge from Language



Input Prompt				
	<p>Question: Explain why this photo is funny? Answer:</p>	<p>Question: Why did the little boy cry? Answer:</p>	<p>Question: What is the hairstyle of the blond called? Answer:</p>	<p>Question: When will the movie be released? Answer:</p>
Completion	<p>The cat is wearing a mask that gives the cat a smile.</p>	<p>Because his scooter broke.</p>	<p>pony tail</p>	<p>On June 27</p>
	(1)	(2)	(3)	(4)
Input Prompt	<p>TorchScale - A Library for Transformers at (Any) Scale</p> <p>TorchScale is a PyTorch library that allows researchers and developers to scale-up Transformers efficiently and effectively. It has the implementation of fundamental research to improve modeling generality and capability as well as training stability and efficiency of scaling Transformers.</p> <ul style="list-style-type: none"> • Stability - <i>Deepen</i> scaling Transformers to 1000 layers and beyond • Generality - <i>Foundation Transformers (Magical)</i> towards true general purpose modeling across tasks and modalities (including language, vision, speech, and multimodal) • Capability - <i>A Length-Extrapolatable Transformer</i> • Efficiency - <i>8.8Mx</i> smaller & faster than open SOTA (GPT-2) 			
	<p>What is TorchScale?</p>	<p>Question: The result is? Answer:</p>	<p>Question: What is the heart rate in this picture? Answer:</p>	<p>The time now is</p>
Completion	<p>A library that allows transformers to work efficiently and effectively at any scale.</p>	<p>5 + 4 = 9</p>	<p>57 bpm</p>	<p>10:10 on a large clock</p>
	(5)	(6)	(7)	(8)

Borrowing Knowledge from Language

Input Prompt	 <p>Question: Explain why this photo is funny? Answer:</p>	 <p>Question: Why did the little boy cry? Answer:</p>	 <p>Question: What is the hairstyle of the blond called? Answer:</p>	 <p>Question: When will the movie be released? Answer:</p>
Completion	<p>The cat is wearing a mask that gives the cat a smile.</p> <p>(1)</p>	<p>Because his scooter broke.</p> <p>(2)</p>	<p>pony tail</p> <p>(3)</p>	<p>On June 27</p> <p>(4)</p>
Input Prompt	 <p>What is TorchScale?</p>	 <p>Question: The result is? Answer:</p>	 <p>Question: What is the heart rate in this picture? Answer:</p>	 <p>The time now is</p>
Completion	<p>A library that allows transformers to work efficiently and effectively at any scale.</p> <p>(5)</p>	<p>$5 + 4 = 9$</p> <p>(6)</p>	<p>57 bpm</p> <p>(7)</p>	<p>10:10 on a large clock</p> <p>(8)</p>

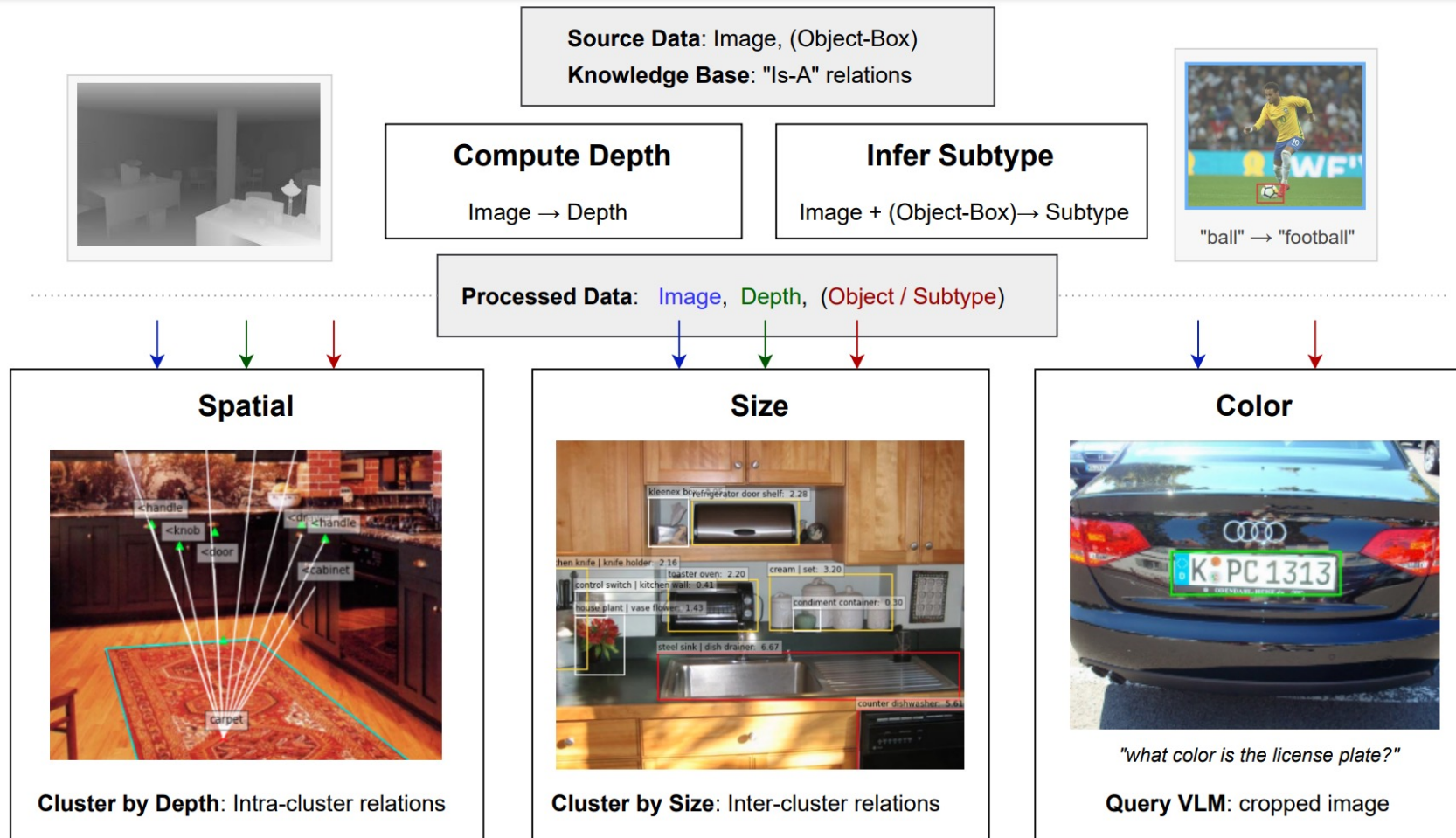
- Language tasks
 - Language understanding
 - Language generation
 - OCR-free text classification
- Cross-modal transfer
 - Commonsense reasoning
- Nonverbal reasoning
 - IQ Test (Raven’s Progressive Matrices)
- Perception-language tasks
 - Image captioning
 - Visual question answering
 - Web page question answering
- Vision tasks
 - Zero-shot image classification
 - Zero-shot image classification with descriptions

Part 3: Are VLMs commonsense KBs?

Probing "Visible" Physical Commonsense Knowledge



Visually accessible knowledge representing color, size and space



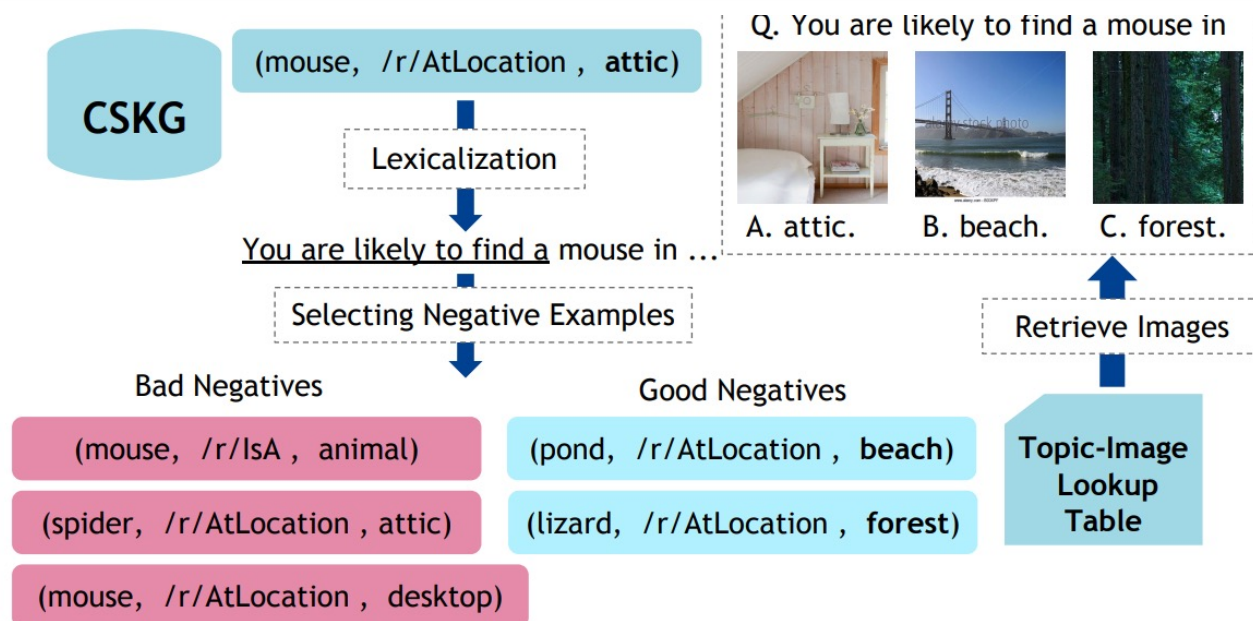
Visually accessible knowledge representing color, size and space

Task	Setting	Prompt
Color	ZS	O is of [MASK] color
	FT	[CLS] color of O
	QA	What is the color of O ? (a) .. (b) ..
Size	ZS	O_1 is [MASK] than O_2 in size
	FT	[CLS] size of O_1 in comparison to O_2
	QA	what is the size of O_1 in comparison to O_2 ? (a) .. (b) ..
Spatial	ZS	in a S , the O_1 is located [MASK] the O_2
	FT	[CLS] in a S , the O_1 is located in comparison to O_2
	QA	in a S , where is O_1 is located in comparison to O_2 ? (a) .. (b) ..

Are Visual-Linguistic Models Commonsense KBs?



CS dimension	Starting prompt	Answer candidates	# Instances
part-whole	<u>Furry animals</u> have	A ₁ : effect of <u>chilling innovation</u> . A ₂ : millions of <u>hair</u> . A ₃ : <u>hole in</u> .	1,165
taxonomic	<u>Recruit</u> is a way to	A ₁ : <u>rate</u> . A ₂ : enlist . A ₃ : <u>slope</u> .	1,323
distinctness	<u>Shade</u> is not	A ₁ : <u>flat</u> . A ₂ : <u>postal worker</u> . A ₃ : sunny .	828
similarity	<u>Throw up</u> is a synonym of	A ₁ : <u>rutinic acid</u> . A ₂ : <u>random</u> . A ₃ : vomit .	644
quality	A <u>wet floor</u> is	A ₁ : slippery . A ₂ : <u>light brown</u> . A ₃ : <u>abbreviated to unido</u> .	1,840
utility	A <u>fork</u> is used for	A ₁ : <u>speed of transit</u> . A ₂ : <u>confuse voters</u> . A ₃ : picking up food .	2,090
creation	<u>Music</u> is created by	A ₁ : <u>olive oil mill</u> . A ₂ : <u>mapping process</u> . A ₃ : instruments .	100
temporal	Going for a <u>haircut</u> requires	A ₁ : finding barber . A ₂ : <u>hard examinations</u> . A ₃ : <u>write persuasively</u> .	1,889
spatial	You are likely to find a <u>document folder</u> in	A ₁ : file drawer . A ₂ : <u>madagascar jungle</u> . A ₃ : <u>minerals</u> .	1,599
desire	You would thank someone because you want to	A ₁ : <u>accomplish mutual goal</u> . A ₂ : feel good . A ₃ : <u>cool off</u> .	1,781



Are Visual-Linguistic Models Commonsense KBs?



CS dimension	Starting prompt	Answer candidates	# Instances
part-whole	<u>Furry animals</u> have	A ₁ : effect of <u>chilling innovation</u> . A ₂ : millions of <u>hair</u> . A ₃ : <u>hole in</u> .	1,165
taxonomic	<u>Recruit</u> is a way to	A ₁ : <u>rate</u> . A ₂ : enlist . A ₃ : <u>slope</u> .	1,323
distinctness	<u>Shade</u> is not	A ₁ : <u>flat</u> . A ₂ : <u>postal worker</u> . A ₃ : sunny .	828
similarity	<u>Throw up</u> is a synonym of	A ₁ : <u>rutinic acid</u> . A ₂ : <u>random</u> . A ₃ : vomit .	644
quality	A <u>wet floor</u> is	A ₁ : slippery . A ₂ : <u>light brown</u> . A ₃ : <u>abbreviated to unido</u> .	1,840
utility	A <u>fork</u> is used for	A ₁ : <u>speed of transit</u> . A ₂ : <u>confuse voters</u> . A ₃ : picking up food .	2,090
creation	<u>Music</u> is created by	A ₁ : <u>olive oil mill</u> . A ₂ : <u>mapping process</u> . A ₃ : instruments .	100
temporal	Going for a <u>haircut</u> requires	A ₁ : finding barber . A ₂ : <u>hard examinations</u> . A ₃ : <u>write persuasively</u> .	1,889
spatial	You are likely to find a <u>document folder</u> in	A ₁ : file drawer . A ₂ : <u>madagascar jungle</u> . A ₃ : <u>minerals</u> .	1,599
desire	You would thank someone because you want to	A ₁ : <u>accomplish mutual goal</u> . A ₂ : feel good . A ₃ : <u>cool off</u> .	1,781



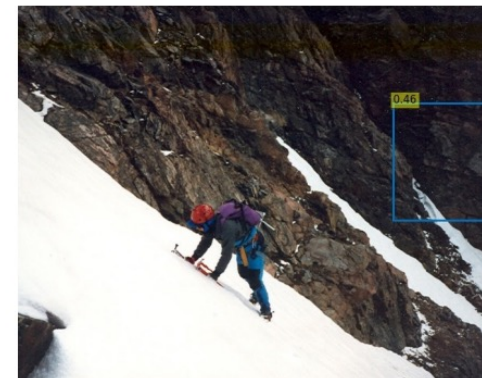
dim.: spatial

You are likely to find vegetables in:
A. workplace.
B. stationary shop.
C. **garden**.



dim.: part-whole

A boat has:
A. reached legal age.
B. **sails**
C. different rules.



dim.: quality

A hill can be:
A. **steep**.
B. about to change.
C. important for normal living.

Are Visual-Linguistic Models Commonsense KBs?



Visual Commonsense Knowledge is more difficult than textual knowledge.

row	Images	part-whole 1, 165	taxonomic 1, 323	distinctness 828	similarity 644	quality 1, 840	utility 2, 090	creation 100	temporal 1, 189	spatial 1, 599	desire 1, 781	All 13, 259	
1	RoBERTa	–	68.5	61.8	80.2	67.4	69.7	74.2	72.0	60.9	54.8	65.9	67.5
2	BERT	–	62.8	71.2	80.1	54.8	68.1	72.4	74.0	53.7	52.4	60.4	65.0
3	BERT _{CC}	–	68.4	62.0	66.6	51.1	66.0	65.4	62.0	53.6	63.7	58.3	61.9
4	UNITER _{BERT_T}	–	70.1	74.5	81.4	62.4	72.0	73.8	79.0	54.5	53.9	61.5	66.5
5	UNITER _T	–	70.9	59.8	71.3	51.2	69.9	71.5	71.0	52.7	61.5	62.5	64.0
6	VILBERT _T	–	63.9	60.3	64.9	46.7	66.1	71.2	58.0	52.2	61.0	62.8	60.7
7	UNITER _{TV}	retrieved	<u>63.0</u>	54.0	<u>65.9</u>	<u>46.4</u>	62.4	65.4	<u>62.0</u>	49.2	57.4	58.5	<u>58.4</u>
8	VILBERT _{TV}	retrieved	55.0	49.9	55.9	42.2	57.4	60.5	52.0	47.2	52.9	56.6	53.0
9	UNITER _{T\tilde{V}}	dummy	61.5	51.6	63.4	42.2	<u>63.6</u>	<u>66.4</u>	55.0	<u>49.4</u>	<u>58.2</u>	59.7	57.1
10	VILBERT _{T\tilde{V}}	dummy	60.4	<u>58.9</u>	64.9	43.9	63.4	65.5	55.0	48.4	56.8	<u>62.0</u>	57.9
11	UNITER _V	retrieved	36.4	<u>36.6</u>	40.1	38.5	34.2	<u>36.6</u>	32.0	<u>34.8</u>	36.2	<u>34.3</u>	36.0
12	VILBERT _V	retrieved	<u>37.8</u>	35.1	37.7	39.8	<u>36.8</u>	35.7	<u>41.0</u>	33.0	<u>37.6</u>	34.0	<u>36.8</u>
13	UNITER _{\tilde{V}}	dummy	30.8	26.3	45.7	28.6	29.2	28.7	19.0	28.7	29.6	30.7	29.7
14	VILBERT _{\tilde{V}}	dummy	34.8	35.8	<u>50.5</u>	<u>40.4</u>	30.4	31.1	30.0	29.4	33.5	30.1	34.6

Unimodal vs Multimodal models?

Unimodal and multimodal models' abilities to capture visual commonsense knowledge



Does the **model** know ...

Unimodal
BERT, ...

VS.

Multimodal
Oscar, ...

Penguins are a group of aquatic flightless birds.

The word penguin first appears in the 16th century as a name for the great auk.

...



A girl is looking at the penguin.



A plastic penguin is sitting on a chair.

...

what is the **color** of a penguin?



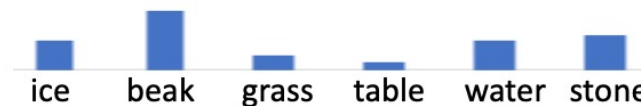
what is the **shape** of a penguin?



what is the **material** of a penguin?



what are the **co-occurring** objects of a penguin?



what is the **size** of a penguin?

It is larger than:



It is smaller than:



Unimodal vs Multimodal models?



ViComTe dataset on five relation types: color, shape, material, size, and visual co-occurrence

Relation	# Classes	# (subj, obj) Pairs	Ex Template	Ex (subj, obj) Pair
color	12	2877	[subj] <i>can be of color</i> [obj]	(<i>sky, blue</i>)
shape	12	706	[subj] <i>has shape</i> [obj] .	(<i>egg, oval</i>)
material	18	1423	[subj] <i>is made of</i> [obj] .	(<i>sofa, cloth</i>)
size (smaller)	107	2000	[subj] <i>is smaller than</i> [obj] .	(<i>book, elephant</i>)
size (larger)	107	2000	[subj] <i>is larger than</i> [obj] .	(<i>face, spoon</i>)
co-occurrence	5939	2108	[subj] <i>co-occurs with</i> [obj] .	(<i>fence, horse</i>)

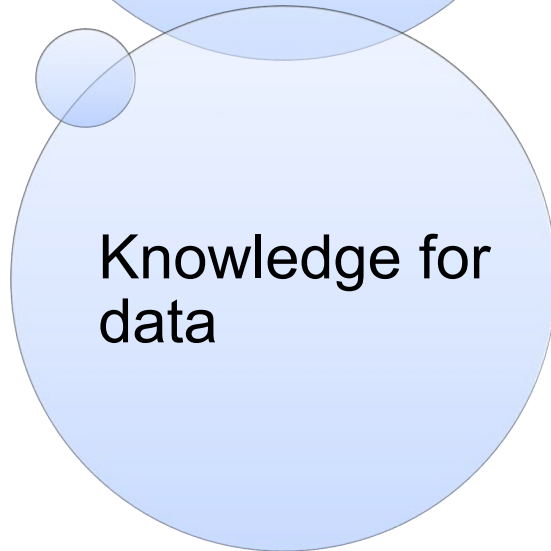
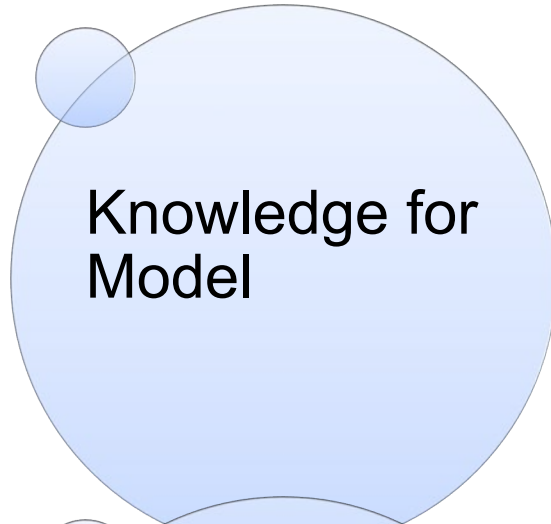
Unimodal vs Multimodal models?



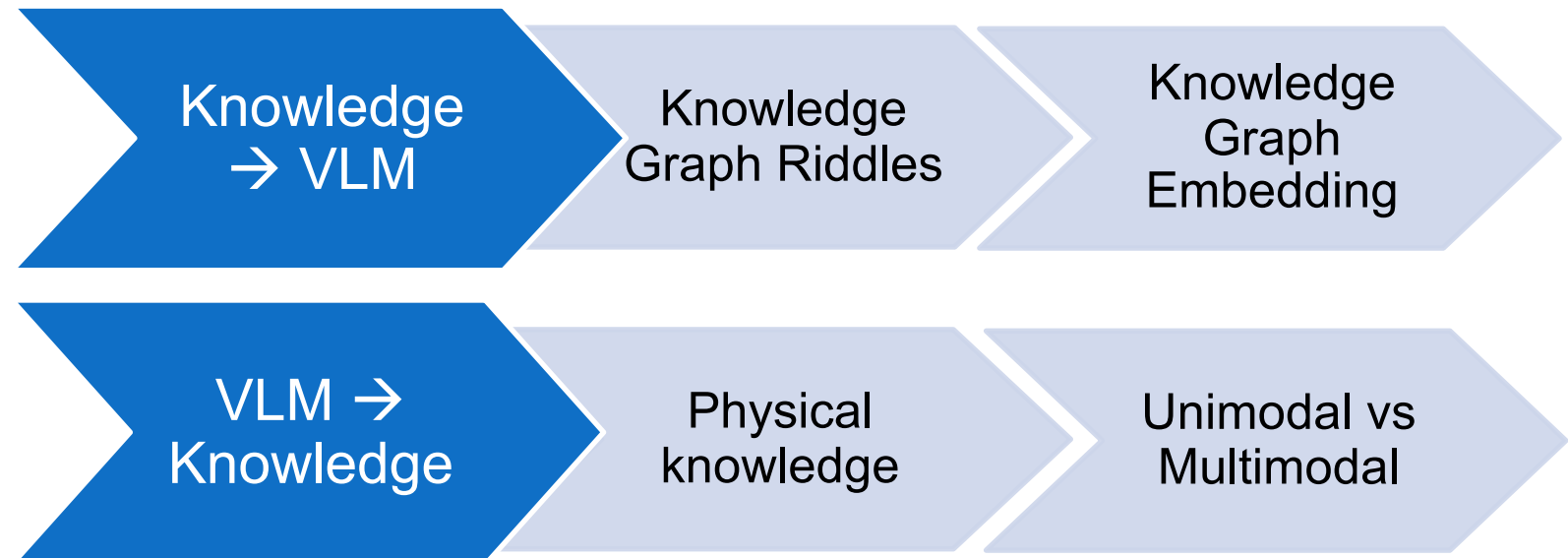
Unimodal and multimodal models' abilities to capture visual commonsense knowledge

Source	Group	Spearman ρ	# Subjs	Avg # Occ	Top5 # Occ	Btm5 # Occ	Acc@1
VG	All	64.3 \pm 23.9	355	1252.6	64.6	308.6	
	SINGLE	62.2 \pm 24.0	131	494.9	64.6	1181.6	80.2
	MULTI	69.3 \pm 20.7	136	1156.1	2062.2	347.0	
	ANY	58.4 \pm 27.1	88	2529.6	8452.4	1213.4	
Wikipedia	All	33.4 \pm 30.6	302	543.6	1758.0	49.8	
	SINGLE	29.6 \pm 29.9	110	352.2	345.8	35.0	35.5
	MULTI	33.9 \pm 30.9	119	500.8	1242.0	27.6	
	ANY	38.2 \pm 30.4	73	902.0	3000.2	161.2	

Future Direction: Adding commonsense knowledge to pretraining



- In-context prompt
- data augmentation
- data selection



Future Direction: Physical Knowledge Enhanced LM/VLM



Humans learn a huge amount of knowledge about the external world via **multisensory experience and interactions**, however, current **LLM/VLM** are trained with **static datasets**, thus **lacks understanding of the physical world**.


Spatial
Relation

Z Put object A to the left of object B. Then, put object B in front of object A. Then, put object C to the left of object A. Which object is directly behind object B?

 Object C is directly behind object B.  

Knowledge
requiring
embodiment

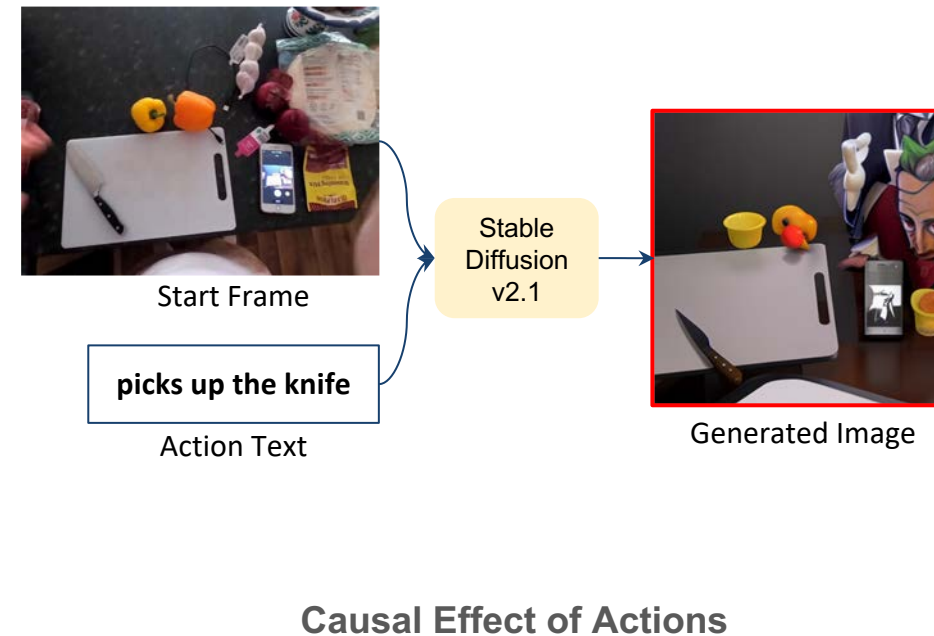
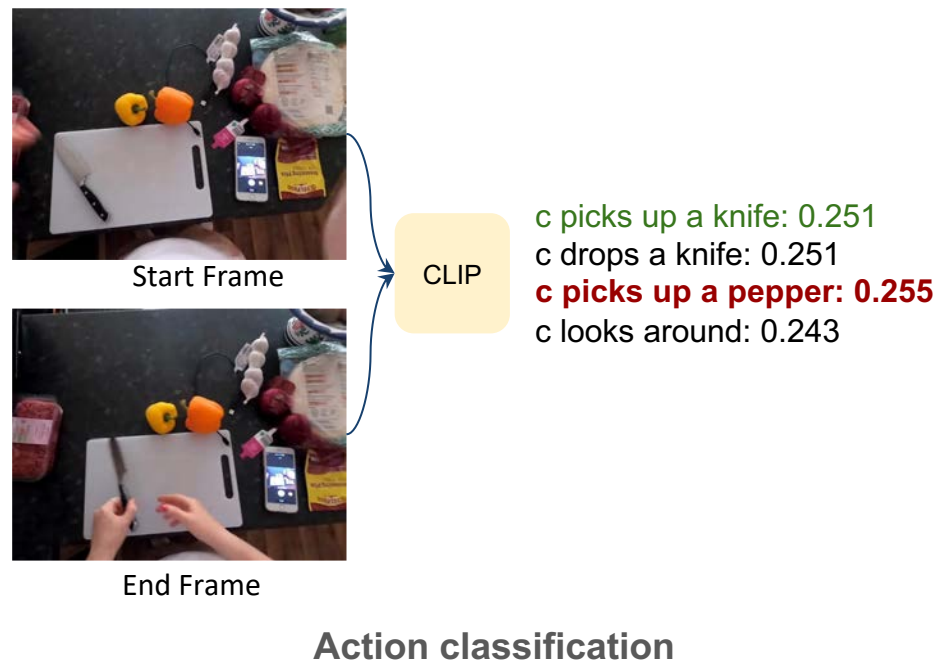
Z Imagine you are a human being. Put your left hand on the back of your head. Can you still see your left hand?

 Yes, I can still see my left hand as it is positioned on the back of my head.

Future Direction: Physical Knowledge Enhanced LM/VLM



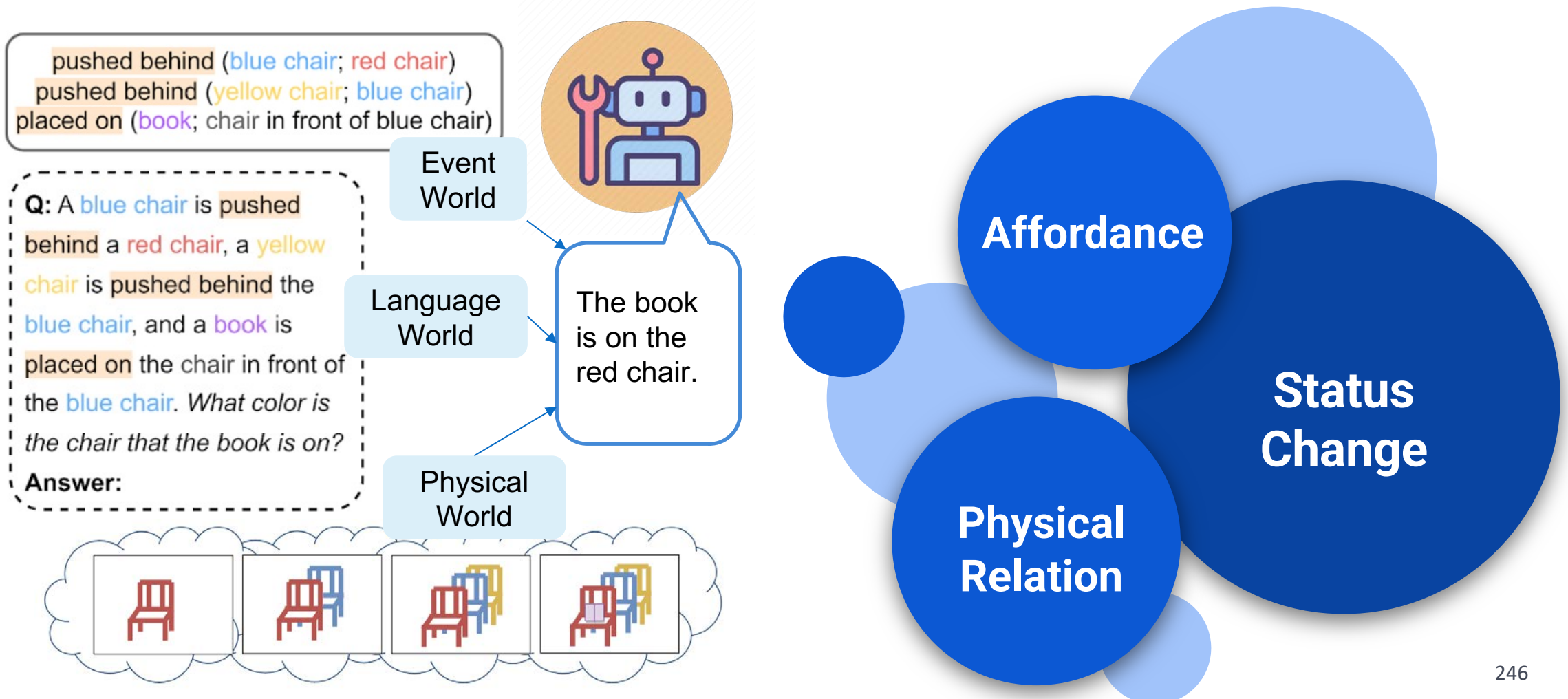
Humans learn a huge amount of knowledge about the external world via **multisensory experience and interactions**, however, current **LLM/VLM** are trained with **static datasets**, thus **lacks understanding of the physical world**.



Future Direction: Physical Knowledge Enhanced LM/VLM



From Reading/Seeing to Doing: From passive perception to interaction with the world.





Disentangling Perception & Reasoning

Initial Exploration: ViperGPT



VIPER-style reasoning shows the potential of **treating perception models as *tools*** and **LLMs as *reasoner*** to solve difficult problems.

Query: How many muffins can each kid have for it to be fair?



Generated Code

```
def execute_command(image):  
    image_patch = ImagePatch(image)  
    muffin_patches = image_patch.find("muffin")  
    kid_patches = image_patch.find("kid")  
    return str(len(muffin_patches) // len(kid_patches))
```

Execution



```
kid_patches =  
image_patch.find("kid")  
► len(muffin_patches)=8  
► len(kid_patches)=2
```

► $8 / 2 = 4$
Result: 4

Query: Drink with zero alcohol



```
def execute_command(image):  
    image_patch = ImagePatch(image)  
    drink_patches = image_patch.find("drink")  
    for drink_patch in drink_patches:  
        drink_name = drink_patch.simple_query("What is this?")  
        alcoholic = llm_query(f"Does the {drink_name} have alcohol?")  
        if alcoholic == "no":  
            return drink_patch  
    return None
```

► drink_patches=



```
► drink_name = 'tullamore dew'  
► alcoholic = 'yes'  
  
► drink_name = 'bacardi'  
► alcoholic = 'yes'  
  
► drink_name = 'gin'  
► alcoholic = 'yes'  
  
► drink_name = 'dr pepper'  
► alcoholic = 'no'
```

Result:

Initial Exploration: ViperGPT



Bottleneck:

- Reasoning is **limited** to the *commonsense* and *code generation* capabilities of LLMs;
- Can only handle a single image context.

Query: How many muffins can each kid have for it to be fair?



Generated Code

```
def execute_command(image):  
    image_patch = ImagePatch(image)  
    muffin_patches = image_patch.find("muffin")  
    kid_patches = image_patch.find("kid")  
    return str(len(muffin_patches) // len(kid_patches))
```

Execution



```
kid_patches =  
image_patch.find("kid")
```



```
► len(muffin_patches)=8  
► len(kid_patches)=2
```

```
► 8//2 = 4
```

Result: 4

Query: Drink with zero alcohol



```
def execute_command(image):  
    image_patch = ImagePatch(image)  
    drink_patches = image_patch.find("drink")  
    for drink_patch in drink_patches:  
        drink_name = drink_patch.simple_query("What is this?")  
        alcoholic = llm_query(f"Does the {drink_name} have alcohol?")  
        if alcoholic == "no":  
            return drink_patch  
    return None
```

► drink_patches=



```
► drink_name = 'tullamore dew'  
► alcoholic = 'yes'
```

```
► drink_name = 'bacardi'  
► alcoholic = 'yes'
```

```
► drink_name = 'gin'  
► alcoholic = 'yes'
```

```
► drink_name = 'dr pepper'  
► alcoholic = 'no'
```

Result:



Example: Answering Questions using Tools



Q: At which festival can you see a castle in the background: Oktoberfest in Domplatz Austria or Tanabata festival in Hiratsuka, Japan?



In 1938, after Hitler had annexed Austria and won the Sudetenland via the Munich Agreement, Oktoberfest was renamed to Großdeutsches Volksfest (Greater German folk festival), and as a showing of strength, the Nazi regime transported people from Sudetenland to the Wiesn by the score.

Large-scale Tanabata festivals are held in many places in Japan, mainly along shopping malls and streets, which are decorated with large, colorful streamers. The most famous Tanabata festival is held in Sendai from 6 to 8 August.



Calella - Catalonia, Spain - 11 Aug. 2009

For the Oktoberfest Löwenbräu brews a special Märzen beer called Oktoberfestbier or Wiesnbier ("meadow beer," referring to the Bavarian name of the festival site, the "Wiesn").

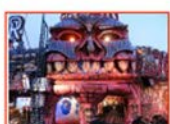
In the summer, the Sendai Tanabata Festival, the largest Tanabata festival in Japan, is held. In winter, the trees are decorated with thousands of lights for the Pageant of Starlight, lasting through most of December.



Masskruege Four mugs of beer at Oktoberfest 2008.



Fussa Tanabata Festival-Tokyo



Ghost train on the Munich Oktoberfest.

A: You can see a castle in the background at Oktoberfest in Domplatz, Austria

WebQA

```
question = "At which festival can you see a castle in the background: Oktoberfest in Domplatz Austria or Tanabata festival in Hiratsuka, Japan?"
input_data = [(img1, "J24 029 Dom, Oktoberfest"), ...]
```

```
# In[1]:
# Filter out irrelevant information for the question
input_data = [
    data_instance
    for data_instance in input_data
    if solve(f"is {data_instance[1]} relevant to the question: {question}")
    # recursive call
    # we can even offload these to a better model (e.g., GPT-3.5)
    # if ask_gpt_yes_or_no(f"is {data_instance[1]} relevant to the question:
{question}")
]
# Out[1]:
input_data == [(img1, "J24 029 Dom, Oktoberfest"), (img6, "Tanabata festival in Hiratsuka")]
```

```
# In[2]:
# solve("Which data instance with a image has a castle on the background?")
from multimodal_models import CLIP
img_features = [CLIP.image_encoder(img) for img, text in input_data]
text = "There is a castle in the background"
text_feature = CLIP.text_encoder(text)
has_castle = [cosine_similarity(text_feature, feat) for feat in img_features]
idx_more_likely_to_have_castle = argmax(has_castle)
# Out[2]:
idx_more_likely_to_have_castle == 0
input_data[idx_more_likely_to_have_castle] == (img1, "J24 029 Dom, Oktoberfest")
```

```
# In[3]:
# Synthesize solution from caption
answer = "You can see a castle in the background at Oktoberfest in Domplatz, Austria"
```



- **Perception:**
 - Vision-only model (e.g., object detection)
 - Vision-language model (e.g., captioning, QA)
- **Reasoner:** Language-only model – We need **divide-and-conquer!**
 - (1) decompose a problem (e.g., a sub-function call)
 - (2) use tool to solve a problem (e.g., access external database, fetch relevant information)
 - (3) update the conclusion (e.g., store something back into the database)

Jun 2023

CVPR Tutorials

Knowledge-Driven Vision-Language Encoding

CVPR

Procedural Knowledge

Knowledge-Driven Vision-Language Encoding (Part IV)

Xudong Lin

Columbia University

xudong.lin@columbia.edu



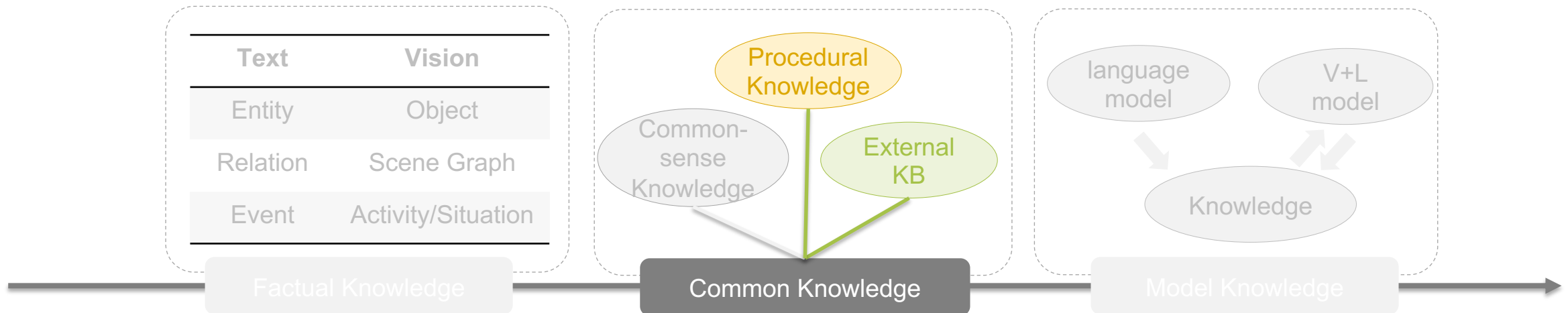
Northwestern
University



COLUMBIA
UNIVERSITY

 Meta AI

Learning patterns of procedure with human-curated patterns and data.



Agenda

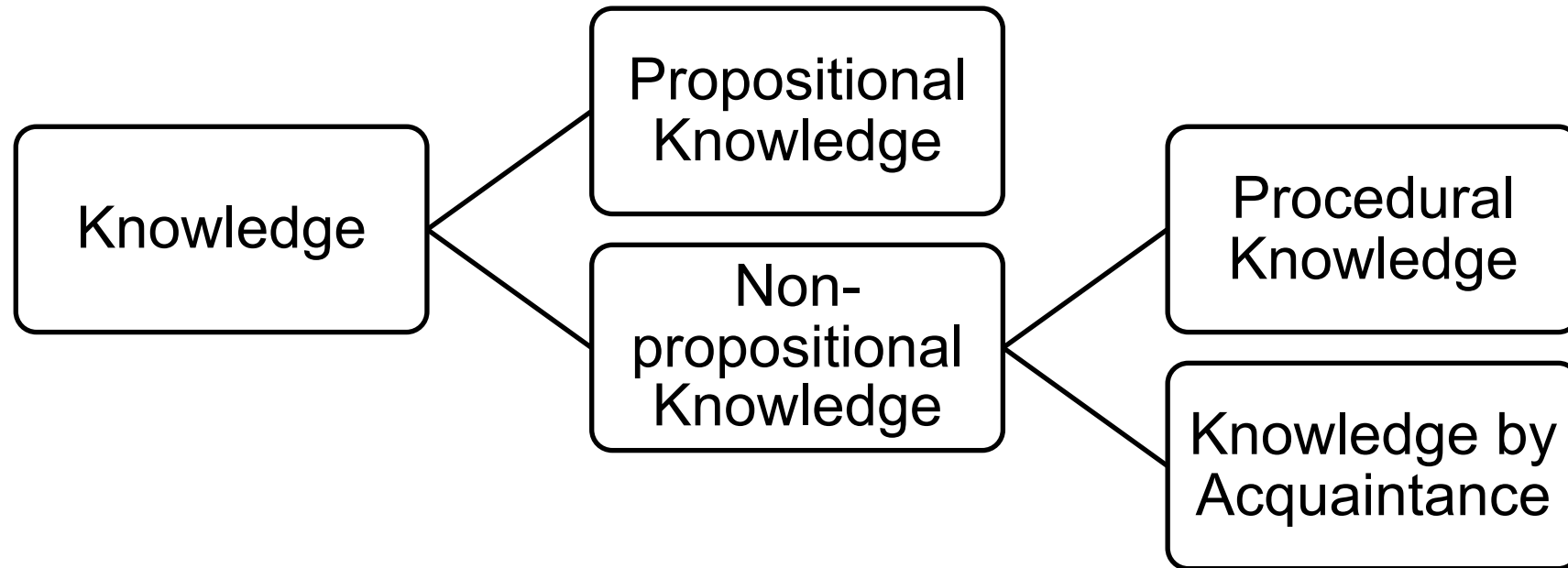


- What is Procedural Knowledge?
- Tasks requiring Procedural knowledge.

What is Procedural Knowledge?



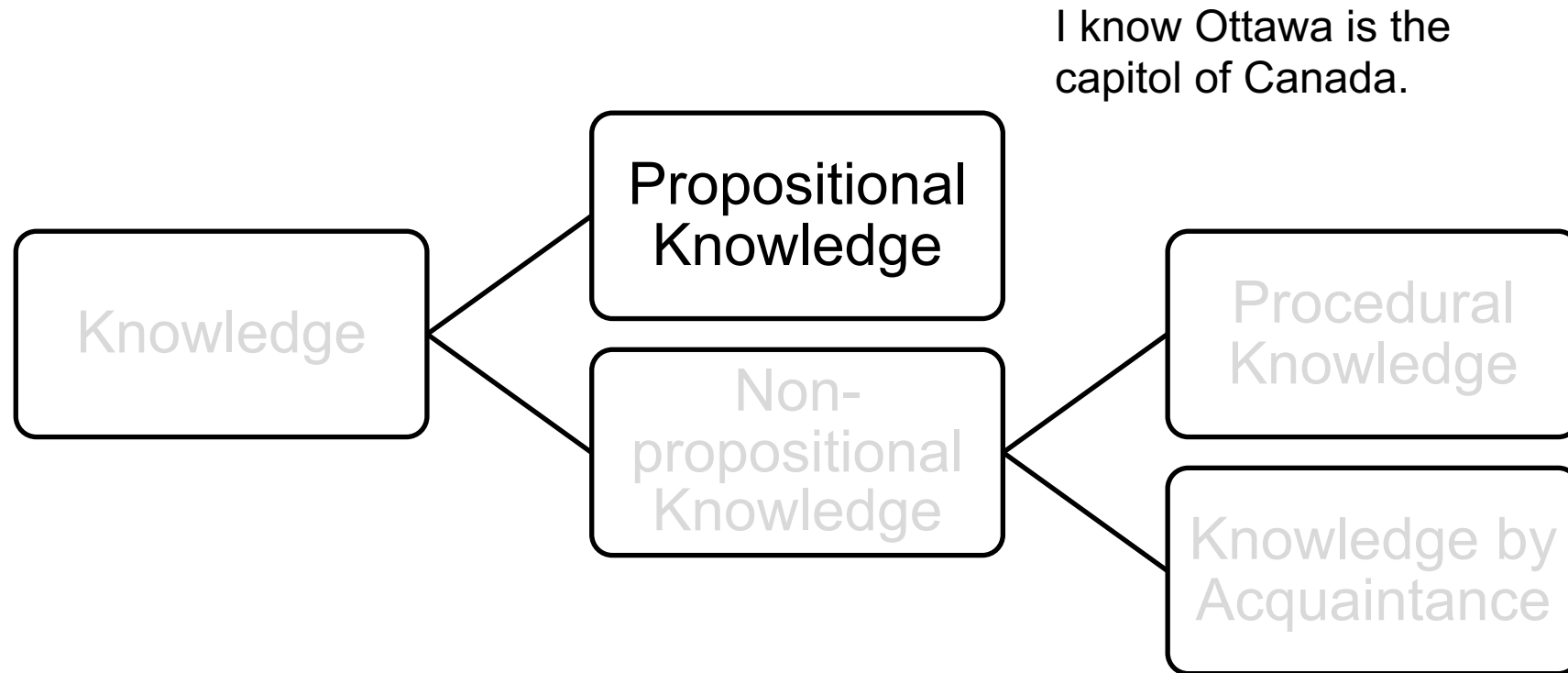
- Psychology View



What is Procedural Knowledge?



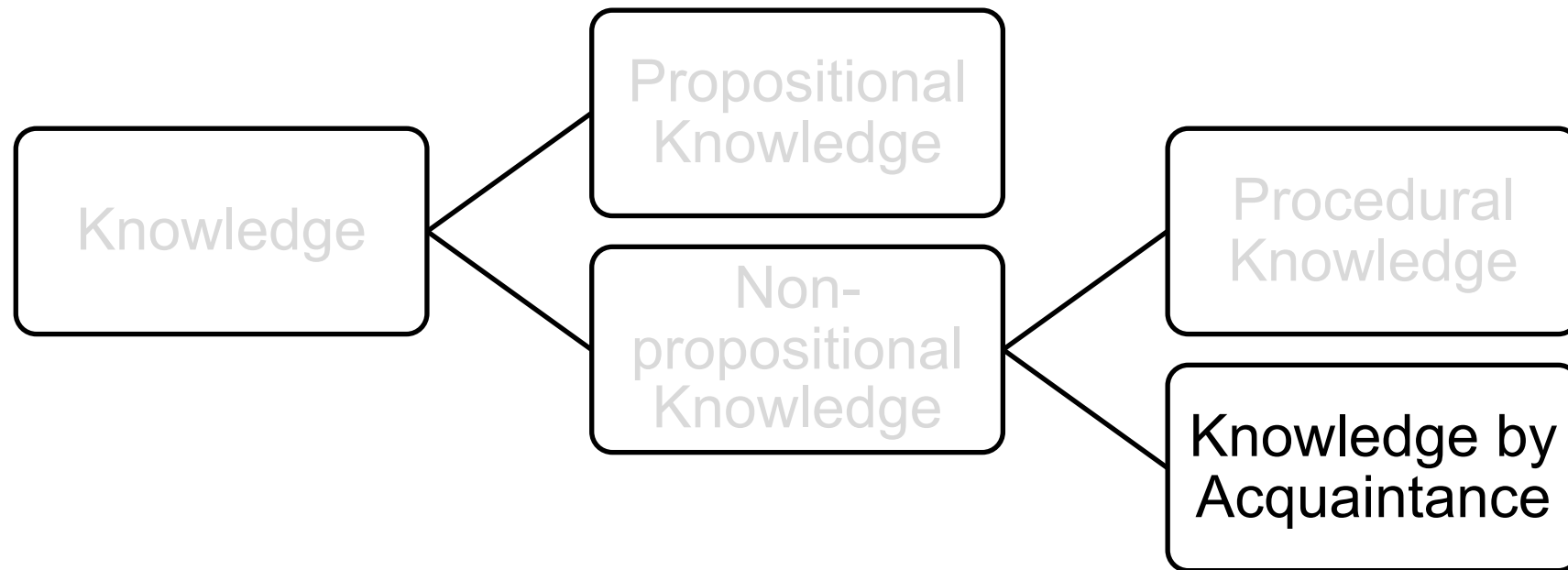
- Psychology View



What is Procedural Knowledge?



- Psychology View

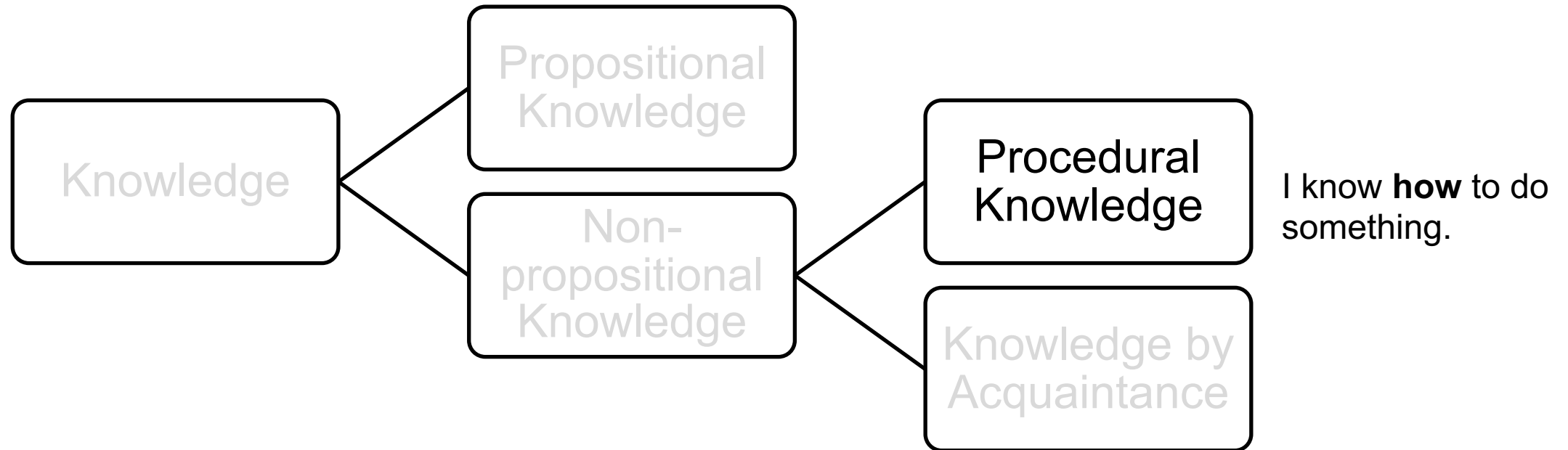


I know someone.

What is Procedural Knowledge?



- Psychology View



Tasks Requiring Procedural Knowledge



- Procedural planning

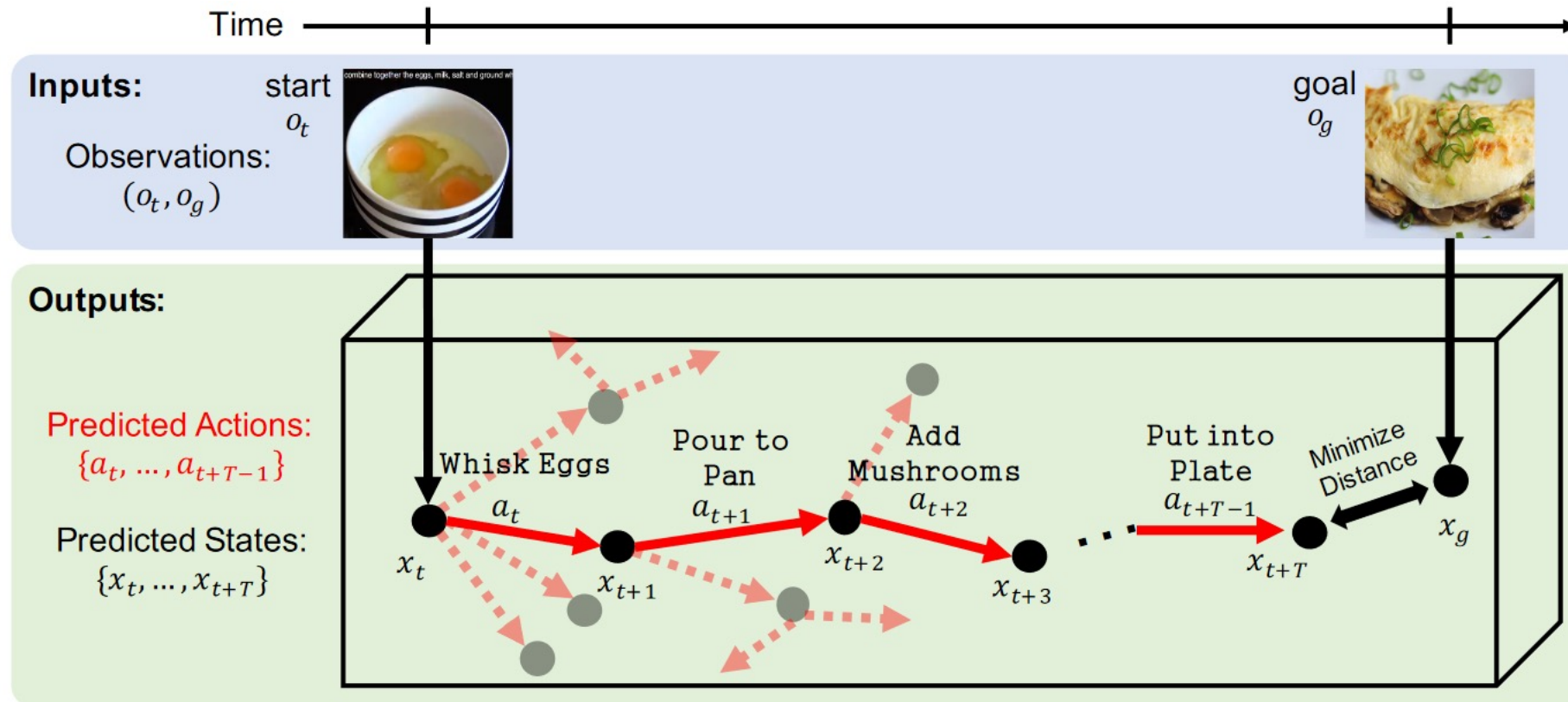


Given a start image and an end image, generate **a sequence of actions**.

Tasks Requiring Procedural Knowledge



- Procedural planning



Given a start image and an end image, generate **a sequence of actions**.

Tasks Requiring Procedural Knowledge



- Step forecasting



What is the next step?

Given the historical video, predict **the next step**.

Frames are from Gordon Ramsay's **Fillet of Beef Wellington**

Sener, Fadime, and Angela Yao. "Zero-shot anticipation for instructional activities." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019.

Lin, Xudong, et al. "Learning to recognize procedural activities with distant supervision." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

Tasks Requiring Procedural Knowledge



- Step forecasting

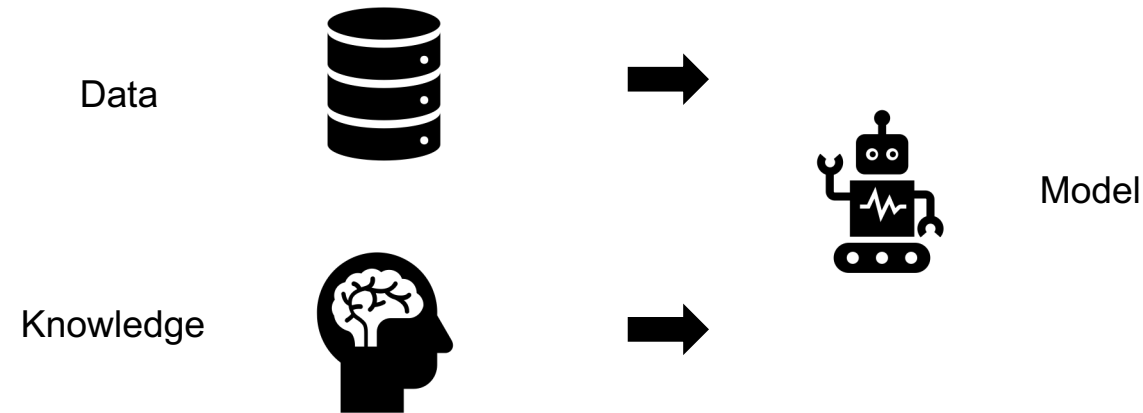


What is the next step?

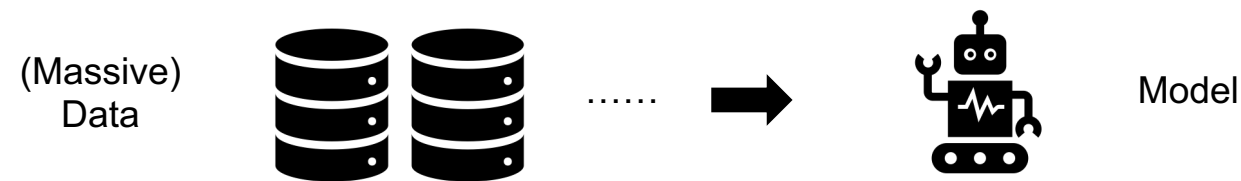
Assembling: Shingle the prosciutto on the plastic wrap; Spread mushroom over prosciutto; ...

Given the historical video, predict **the next step**.

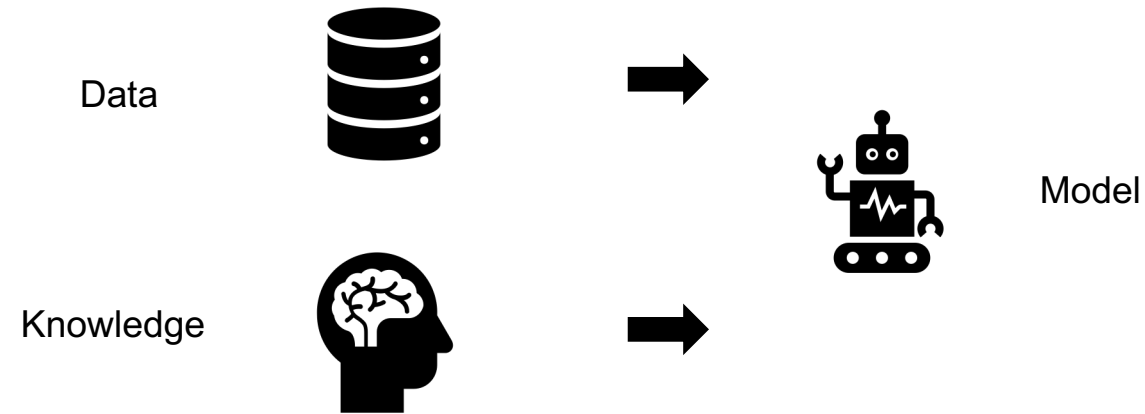
- Explicit Knowledge Source: Learning with the help of external knowledge



- Implicit Knowledge Source: Learning procedural knowledge from data



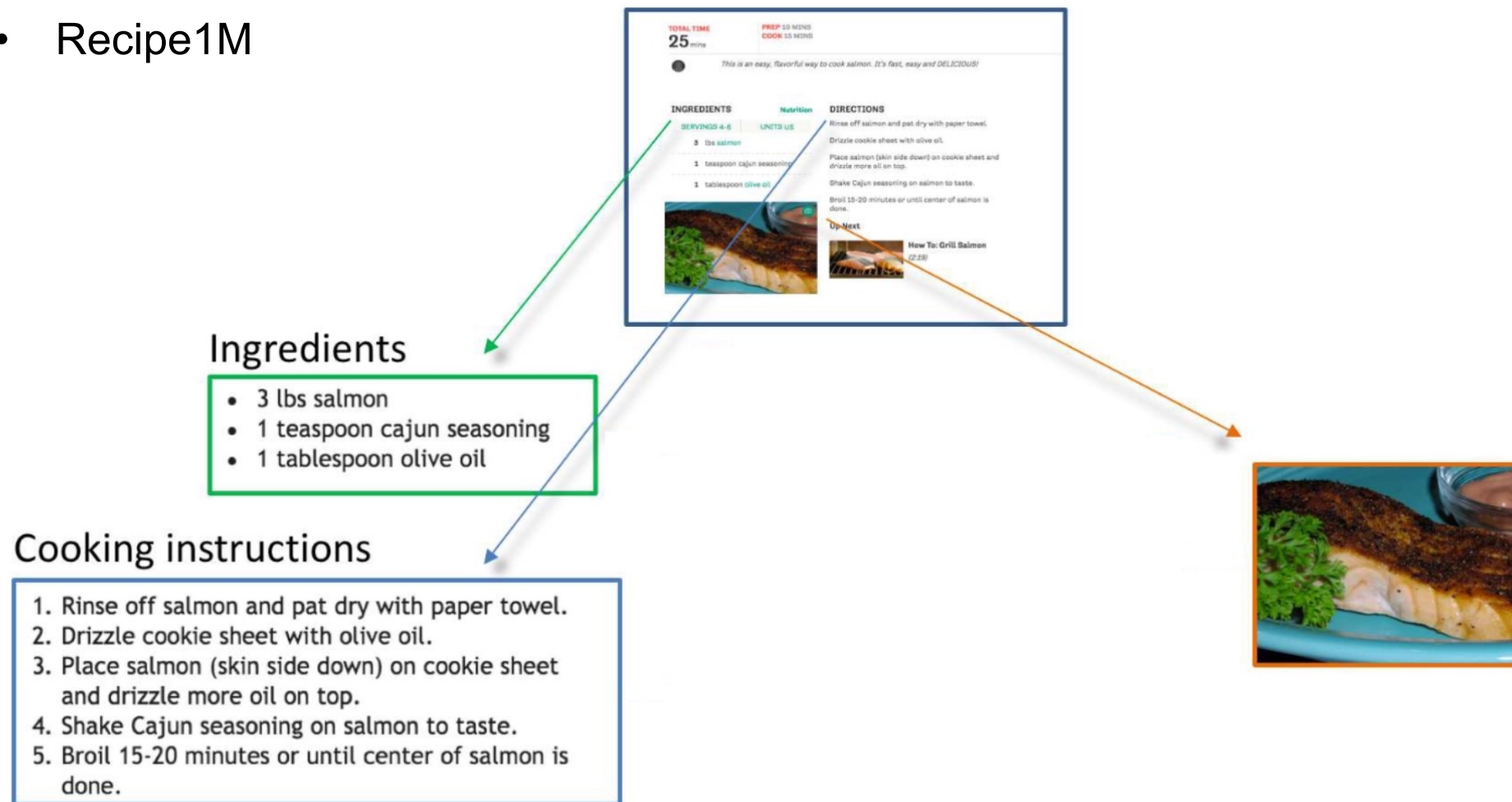
- Explicit Knowledge Source: Learning with the help of external knowledge



Explicit Knowledge Source



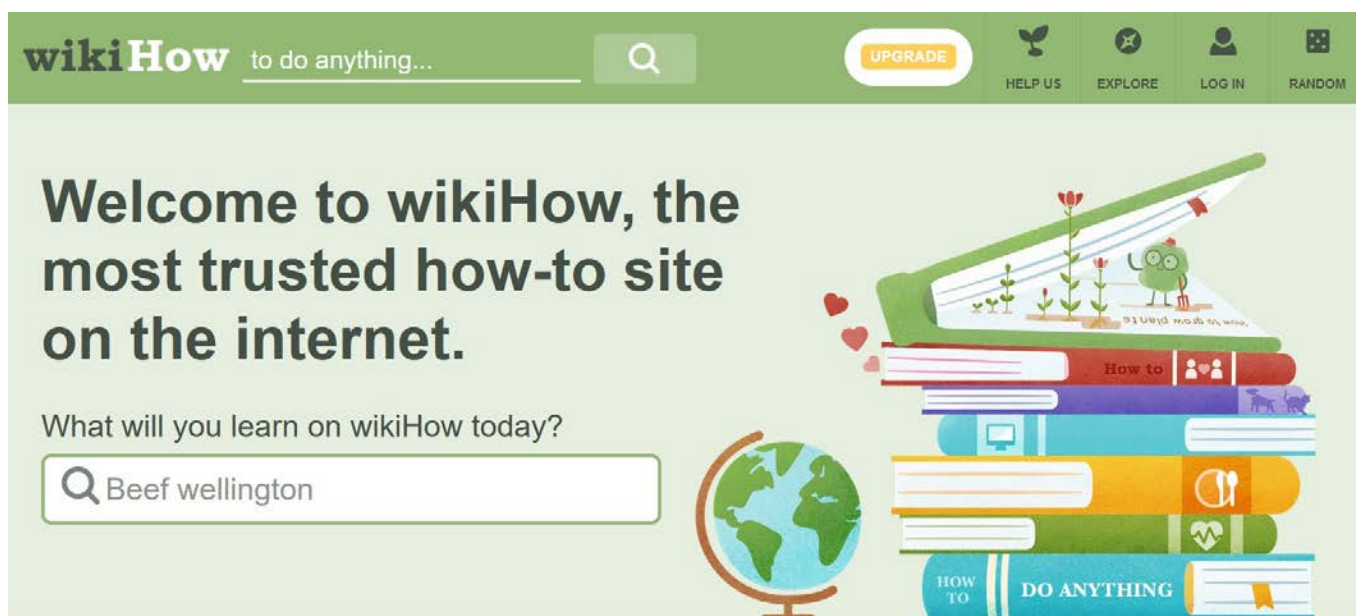
- Procedural knowledge can be easily curated from the Internet
 - Recipe1M



Explicit Knowledge Source



- Procedural knowledge can be easily curated from the Internet
 - Recipe1M
 - wikiHow



Step 1. Sear the fillet mignon to brown.

Over high heat, coat bottom of a heavy skillet with olive oil. Once pan is nearly smoking, sear tenderloin until well-browned on all sides.

Step 2. Fry the mushroom until they are dried.

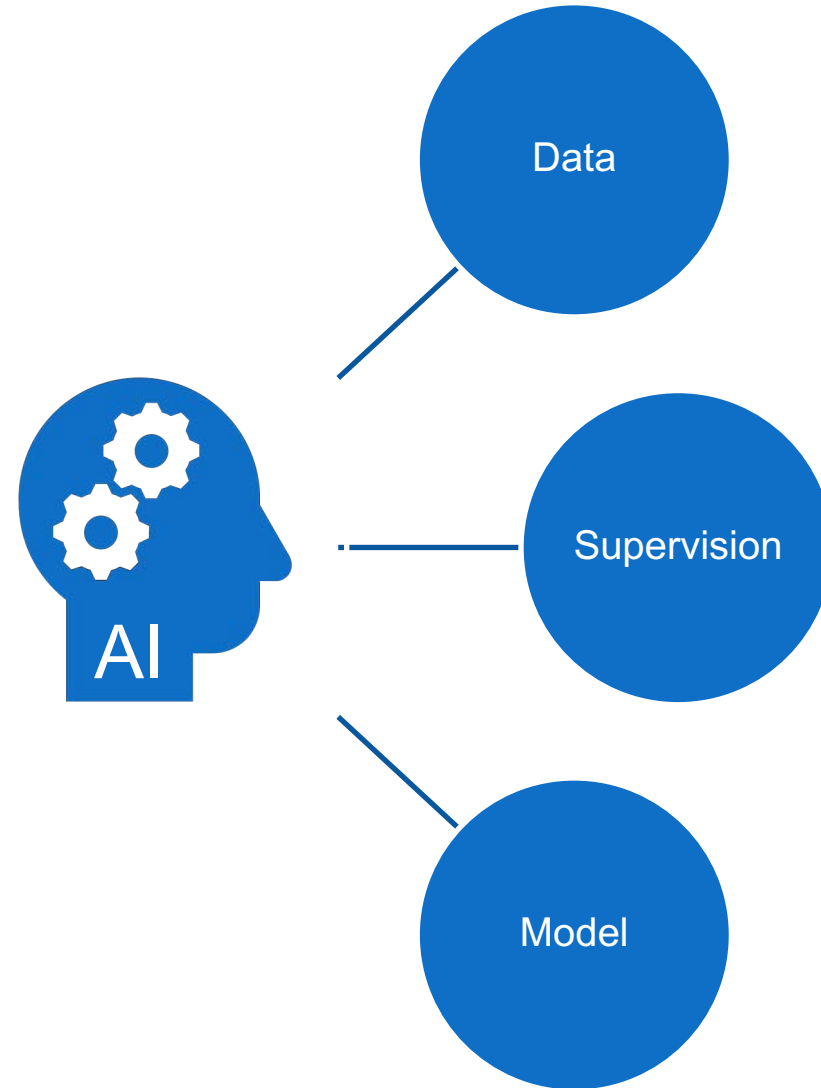
To skillet, add butter and melt over medium heat. Add mushroom mixture and cook until liquid has evaporated.

Step 3. Assembling.

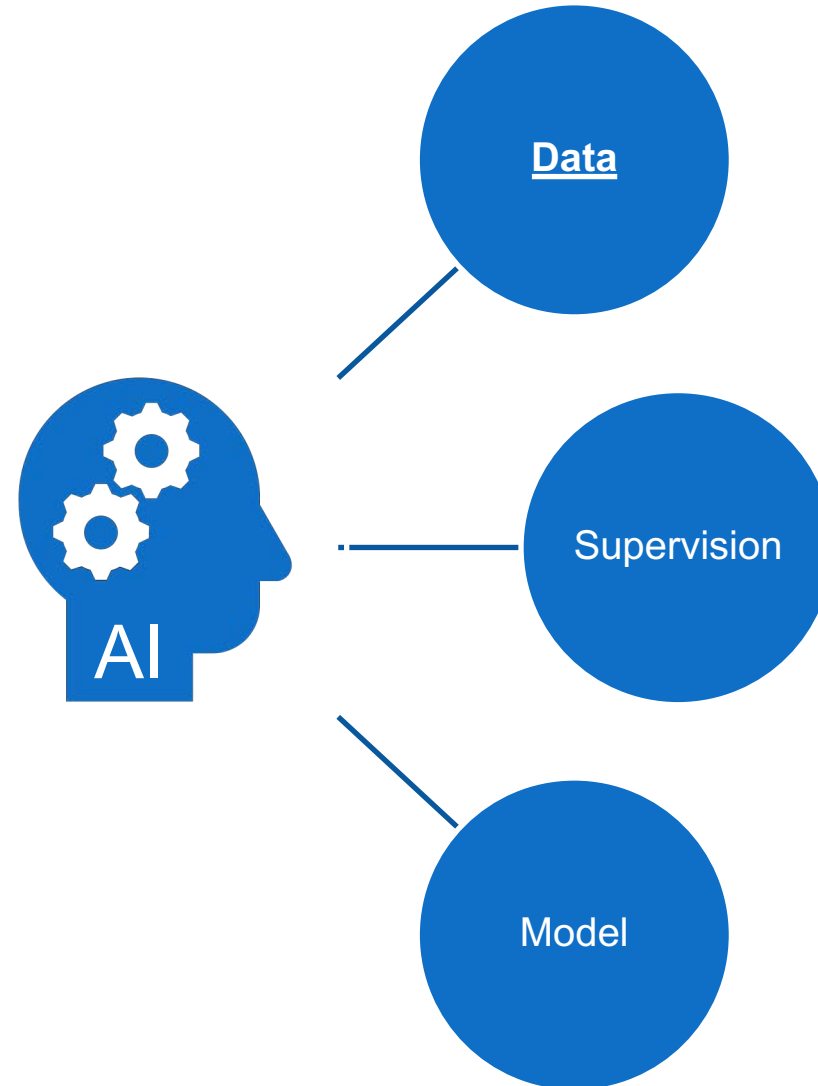
Shingle the prosciutto on the plastic wrap into a rectangle that's big enough to cover the whole tenderloin. Spread the duxelles evenly and thinly over the prosciutto.

.....

How to Utilize the Knowledge Source?



How to Utilize the Knowledge Source?

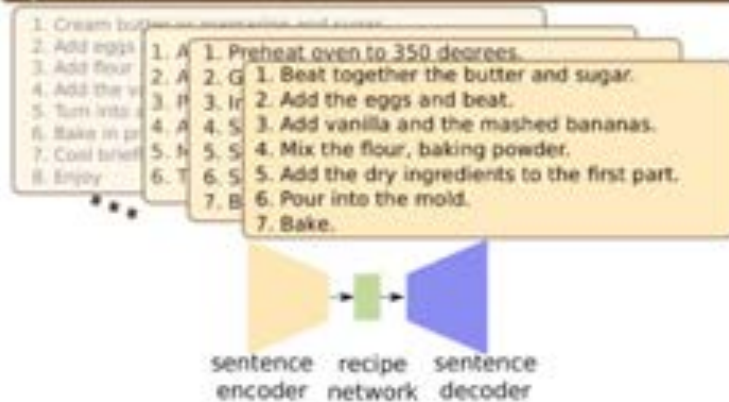


Zero-Shot Anticipation for Instructional Activities



- **Key Idea: Obtain training data from knowledge base.**

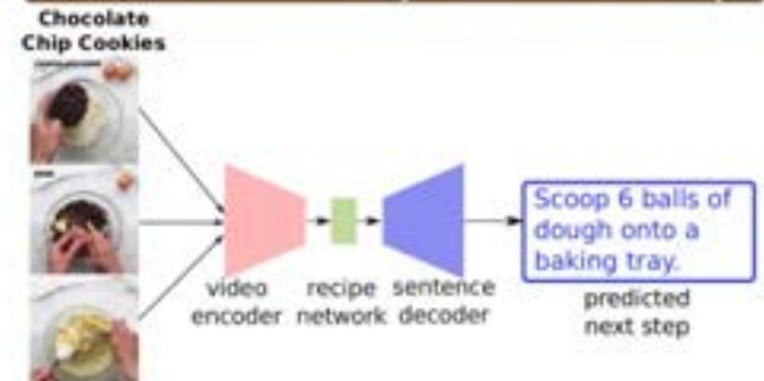
1) Learn instructional tasks from text



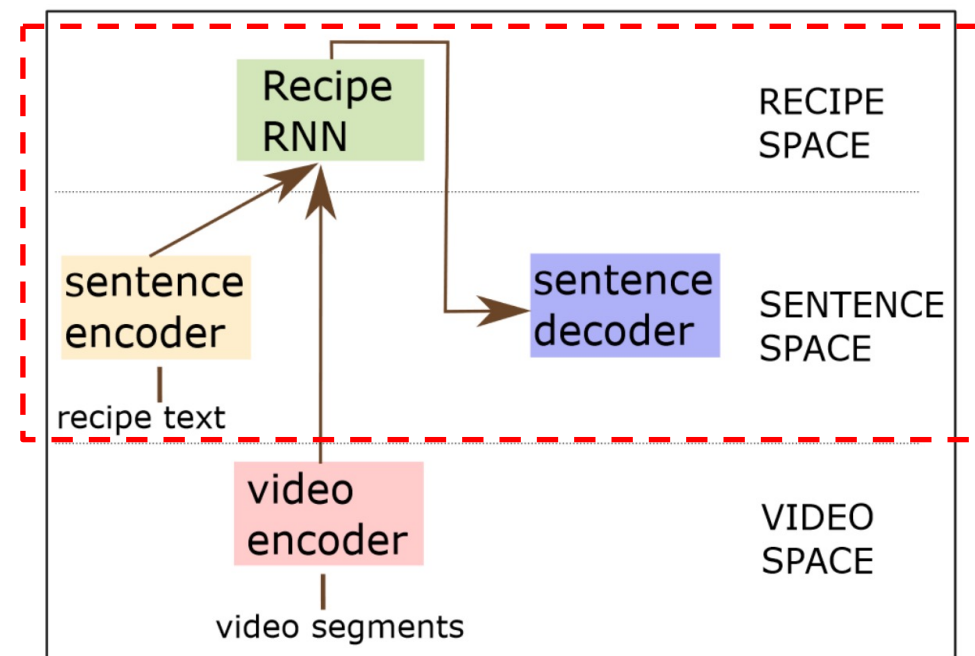
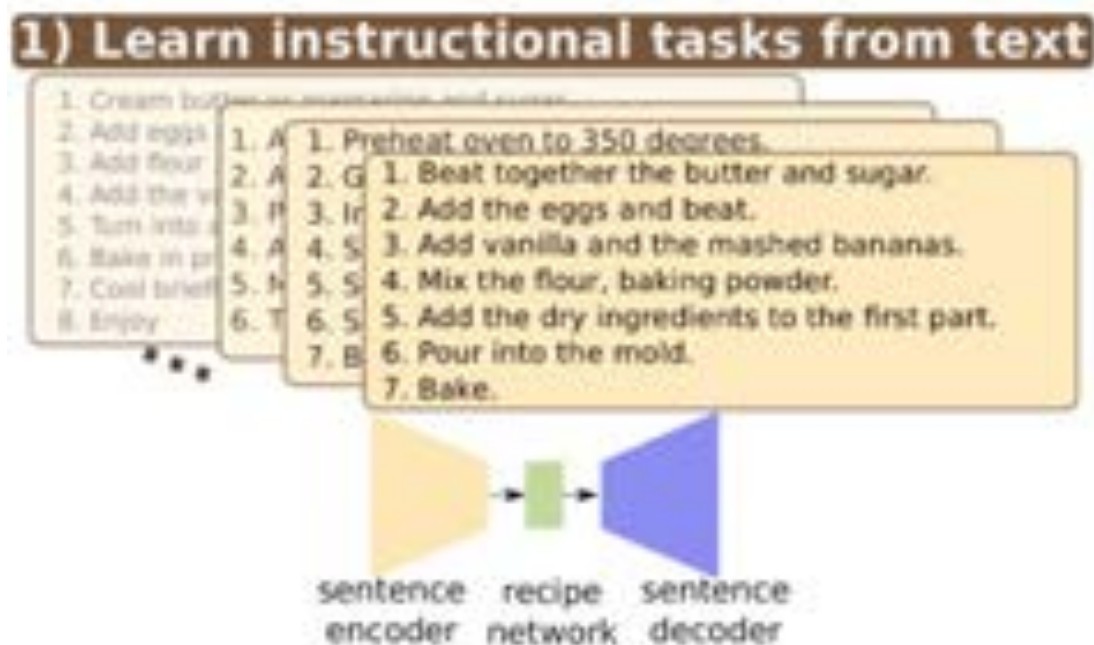
2) Transfer knowledge to video



3) Zero-shot Task: predict next steps



- Sentence encoder encodes a step sentence into a step vector.
- Recipe network is a RNN modeling procedures.
- Sentence decoder decodes step sentences.

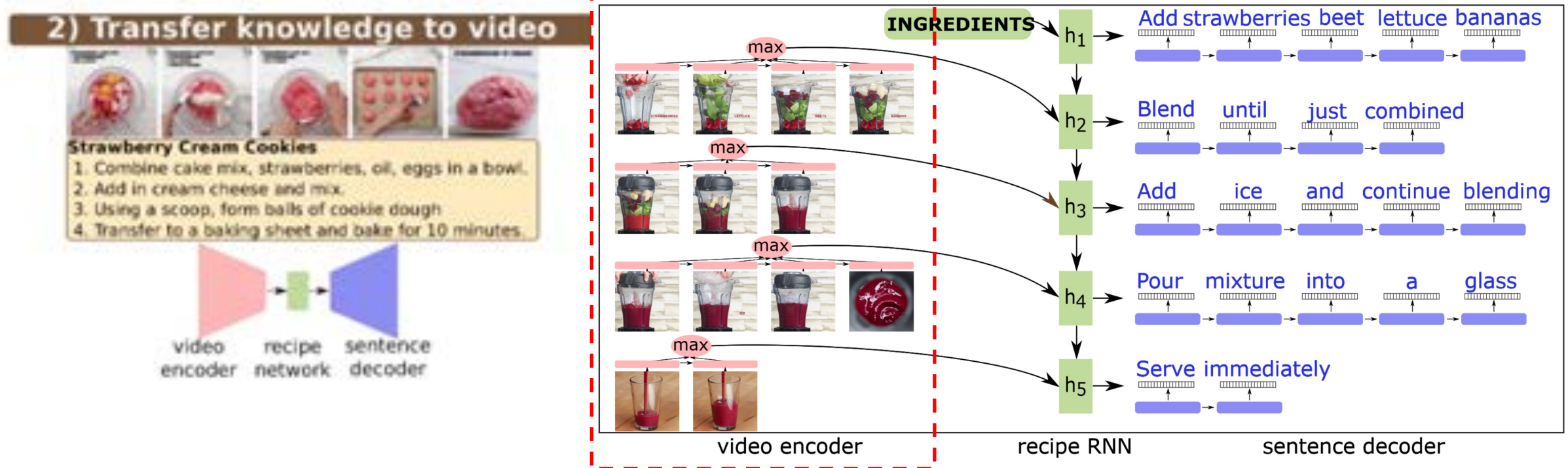


Model Overview

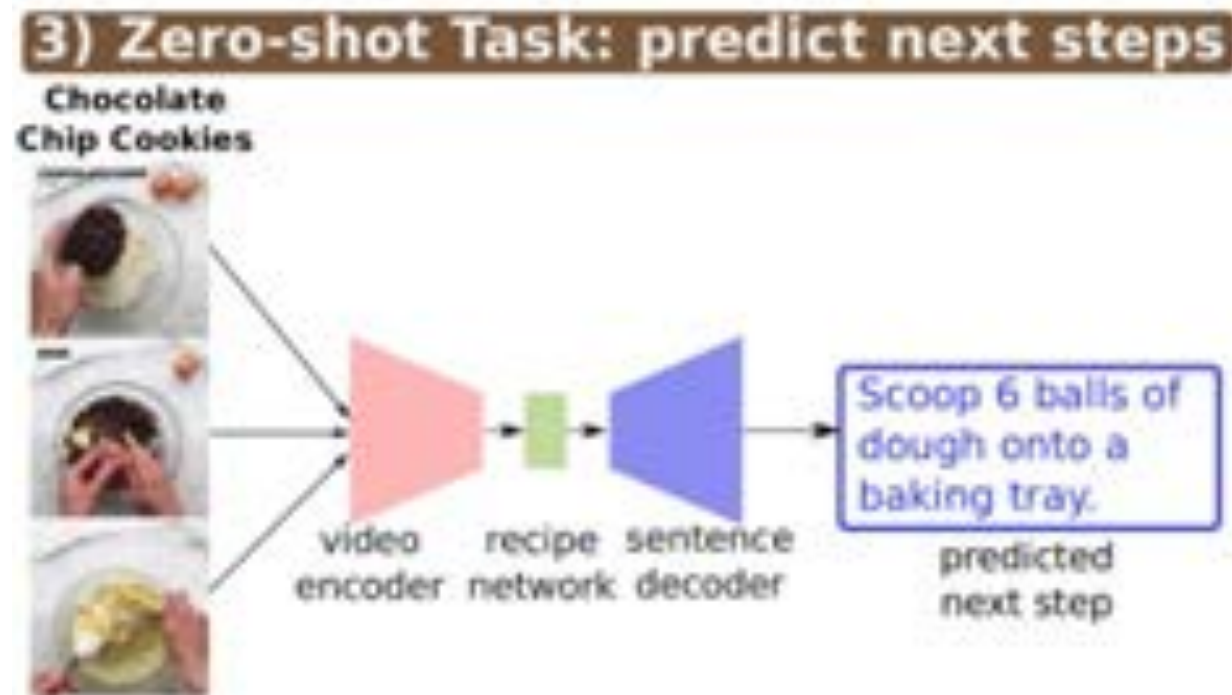
Zero-Shot Anticipation for Instructional Activities



- Only train the video encoder to project video into step vectors with annotated data.



- Generalize on new tasks.



- Strong zero-shot performance on the proposed Tasty video dataset

The larger knowledge base used, the better!

Method	ING	VERBS	BLEU1	BLEU4	METEOR
S2VT [53] (GT)	7.59	19.18	18.03	1.10	9.12
S2VT [53], next (GT)	1.54	10.66	9.14	0.26	5.59
End-to-end [60]	-	-	-	0.54	5.48
Ours Visual (GT)	20.40	19.18	19.05	1.48	11.78
Ours Visual	16.66	17.08	17.59	1.23	11.00
Ours Text (100%)	26.09	27.19	26.78	3.30	17.97
Ours Text (50%)	23.01	24.90	25.05	2.42	16.98
Ours Text (25%)	19.43	23.83	23.54	2.03	16.05
Ours Text (0%)	5.80	9.42	10.58	0.24	6.80
Ours Text noING	9.04	22.00	20.11	0.92	13.07
Ours joint video-text	22.27	23.35	21.75	2.33	14.09

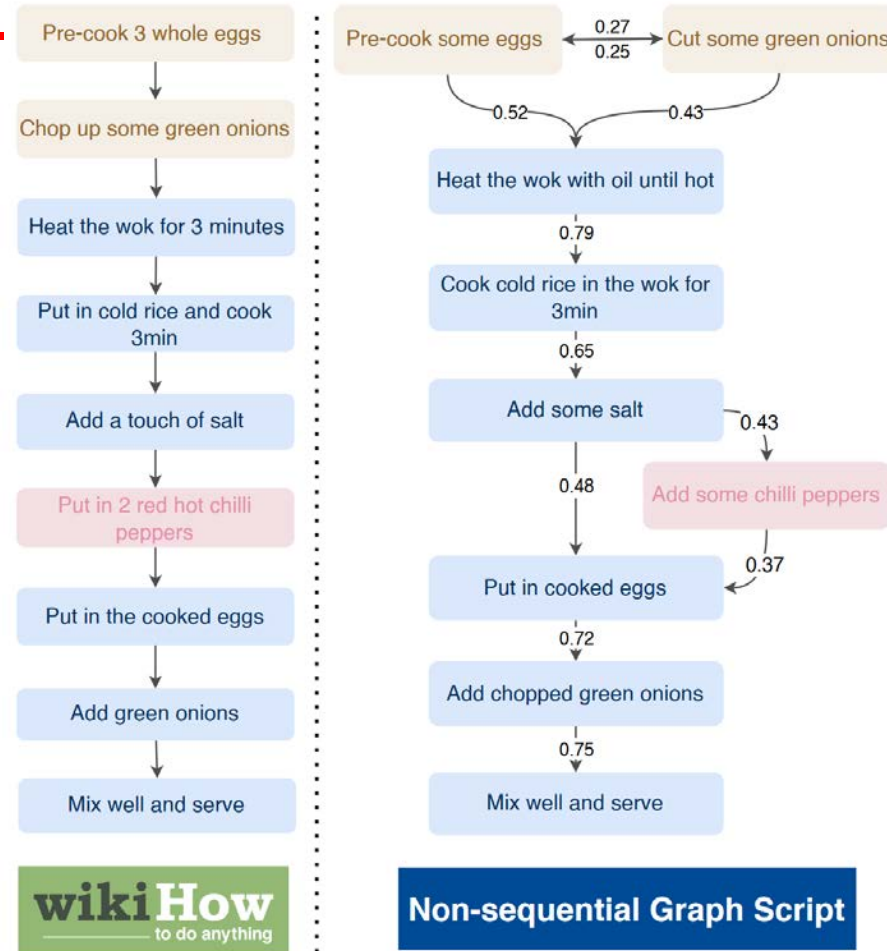
- Limitation

- Domain is limited to cooking.
- Rely on annotated data samples for training video encoder.

Non-Sequential Graph Script Induction via Multimedia Grounding



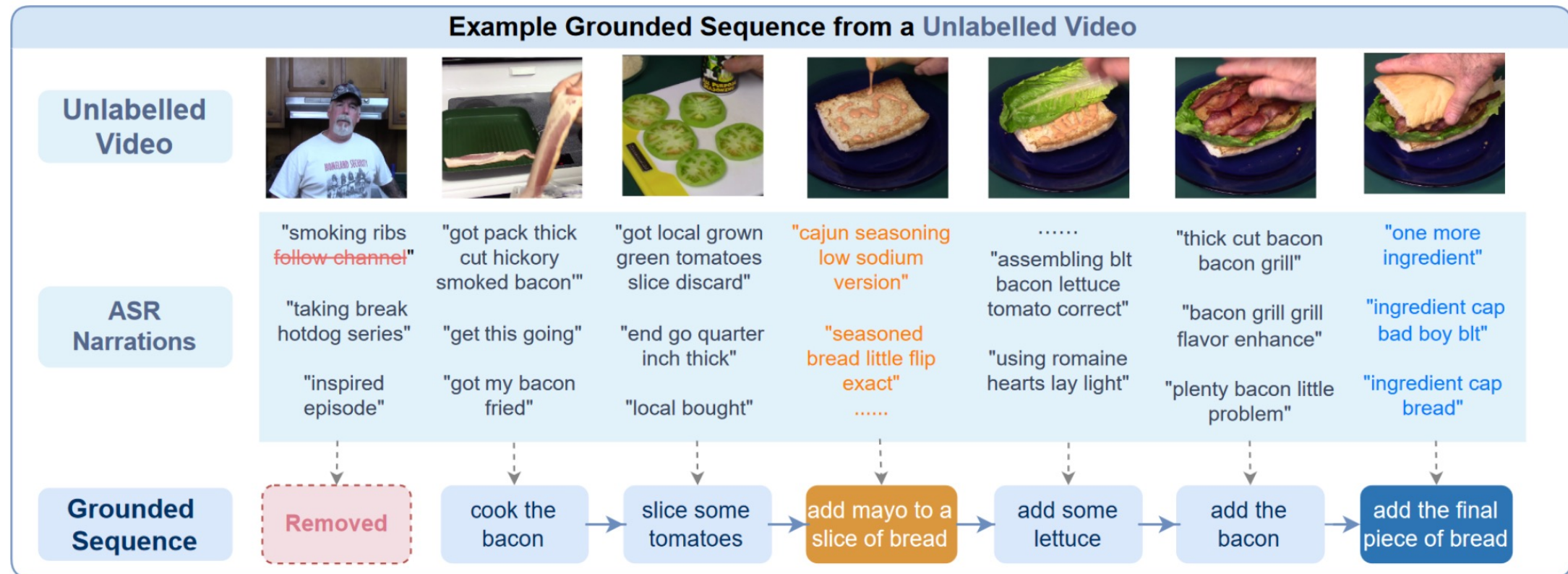
- **Key Idea: Obtain non-sequential script by grounding wikiHow steps to video observations.**



Non-Sequential Graph Script Induction via Multimedia Grounding



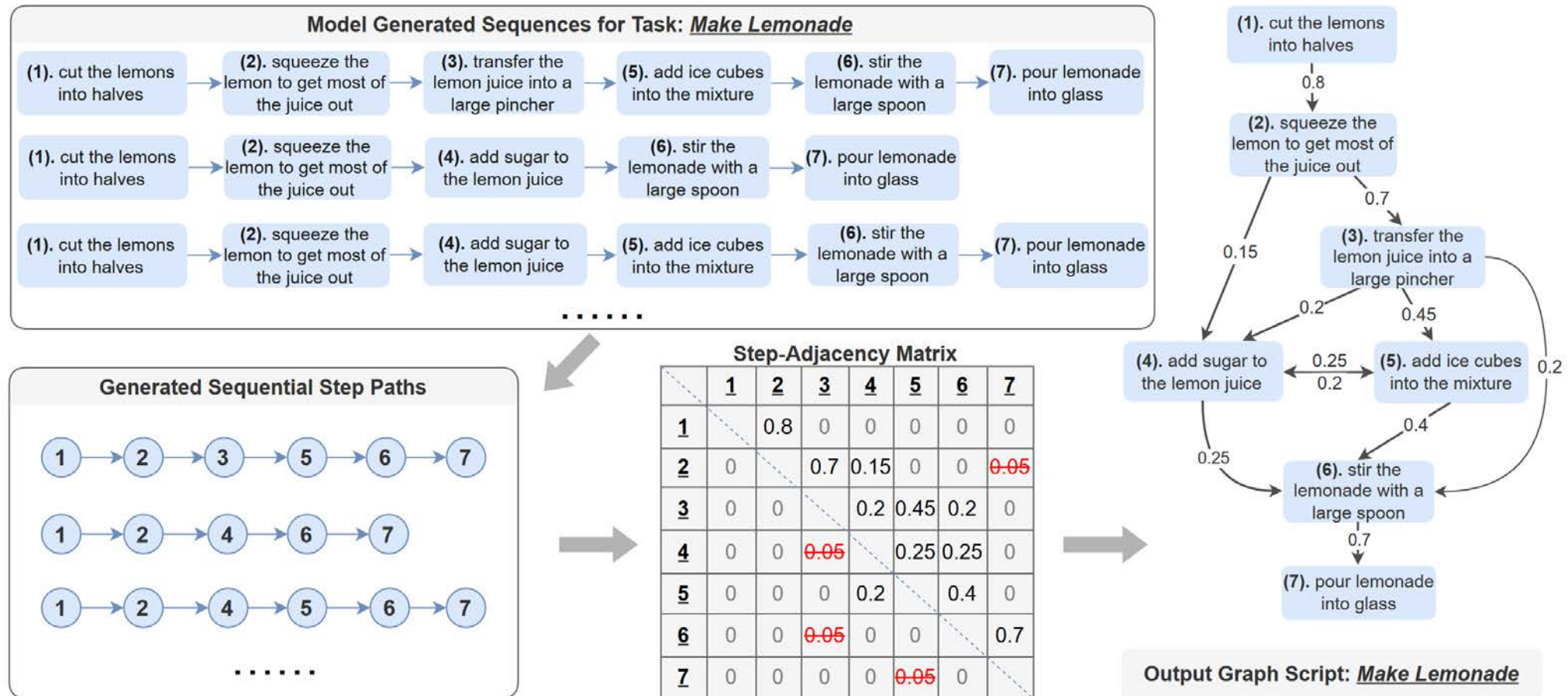
- Video observations contain real-world variance in procedure, which makes wikiHow scripts non-sequential.



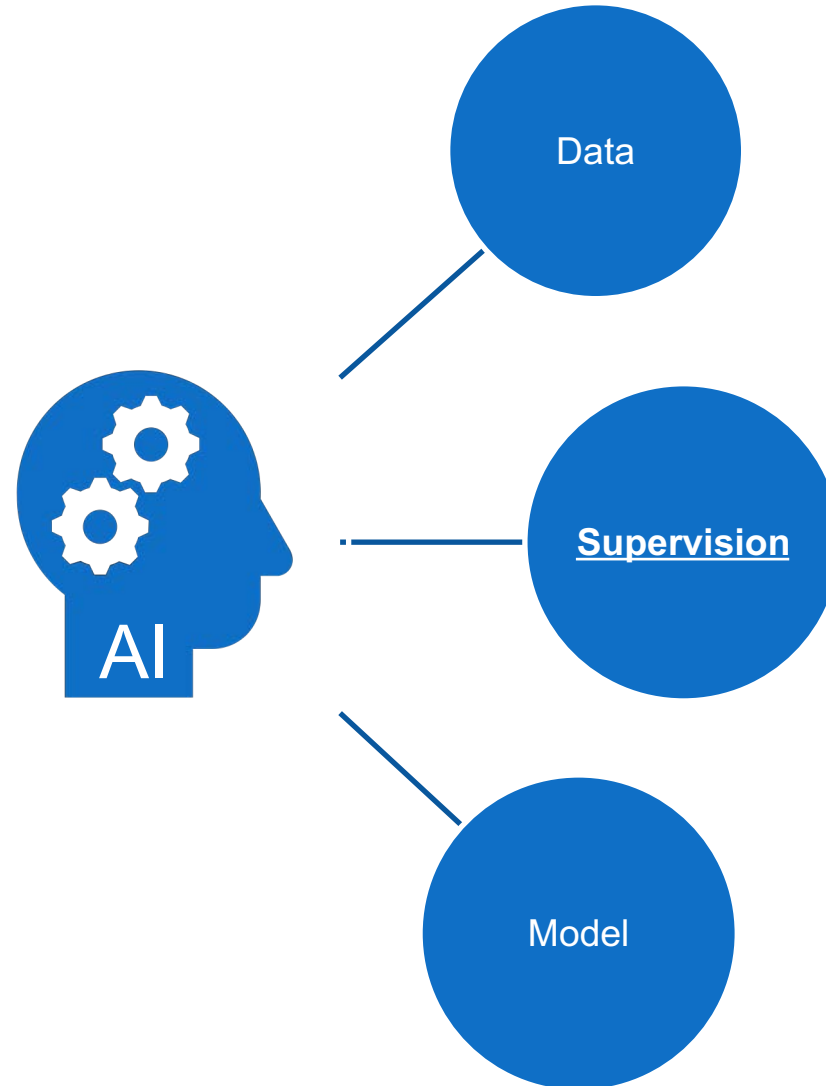
Non-Sequential Graph Script Induction via Multimedia Grounding



- Graphs can be constructed by merging multiple decoded sequence.
- Limitation: closed-vocabulary; text-only graph.



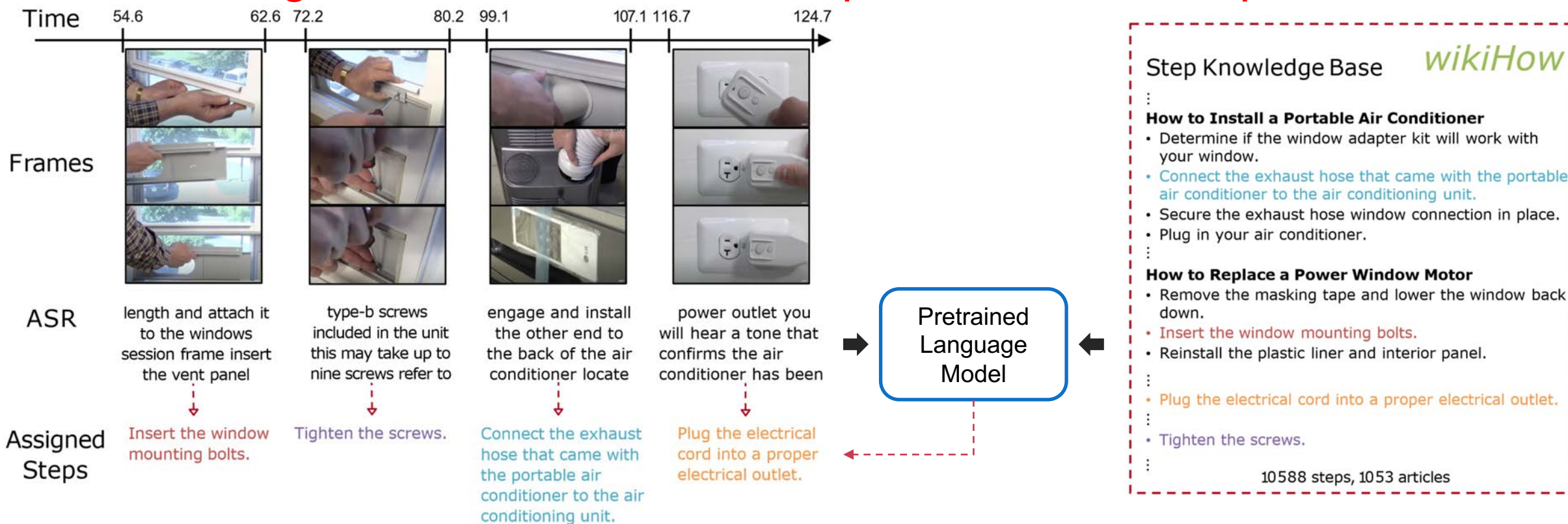
How to Utilize the Knowledge Source?



Learning To Recognize Procedural Activities with Distant Supervision



- Key Idea: Leverage pretrained language model to align knowledge base and videos with speech to obtain supervision.**



- Step Knowledge Base Construction

- Use 1053 tasks, each of which has at least 100 examples in the HowTo100M dataset
- Find the corresponding articles on WikiHow
- Collect sentences for each step in each of the tasks

Step Knowledge Base *wikiHow*

⋮

How to Install a Portable Air Conditioner

- Determine if the window adapter kit will work with your window.
- Connect the exhaust hose that came with the portable air conditioner to the air conditioning unit.
- Secure the exhaust hose window connection in place.
- Plug in your air conditioner.

⋮

How to Replace a Power Window Motor

- Remove the masking tape and lower the window back down.
- Insert the window mounting bolts.
- Reinstall the plastic liner and interior panel.

⋮

- Plug the electrical cord into a proper electrical outlet.

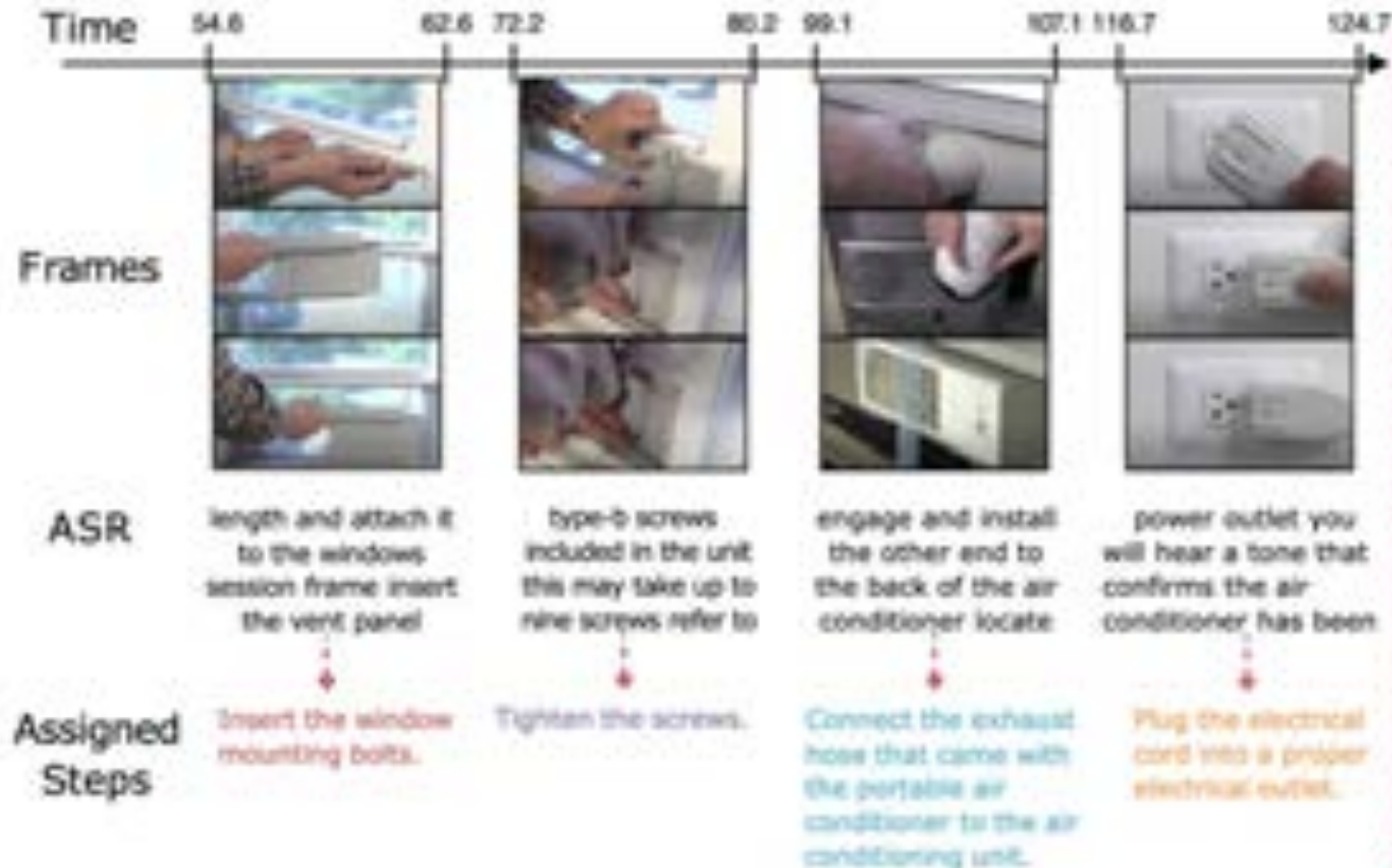
⋮

- Tighten the screws.

⋮

10588 steps, 1053 articles

Learning To Recognize Procedural Activities with Distant Supervision

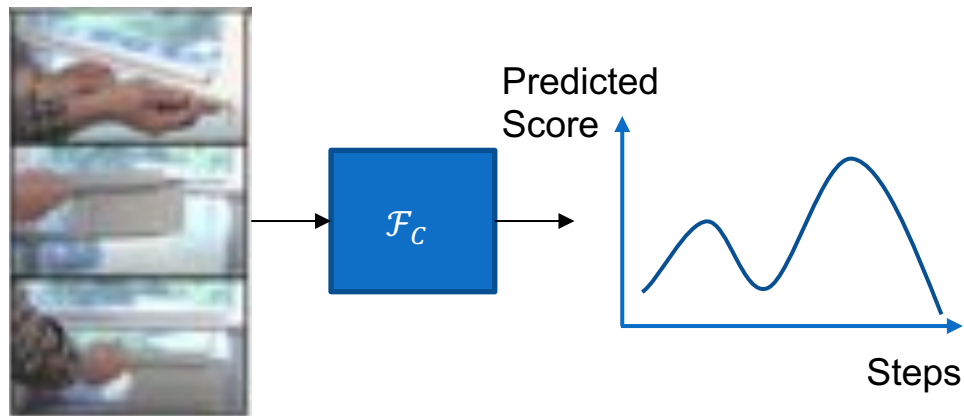


- Distant supervision creation
 - Leverage a pretrained language model to produce embeddings for both **steps** and **ASR sentences** from the video.
 - Then calculate similarity between each ASR sentence and all the steps.

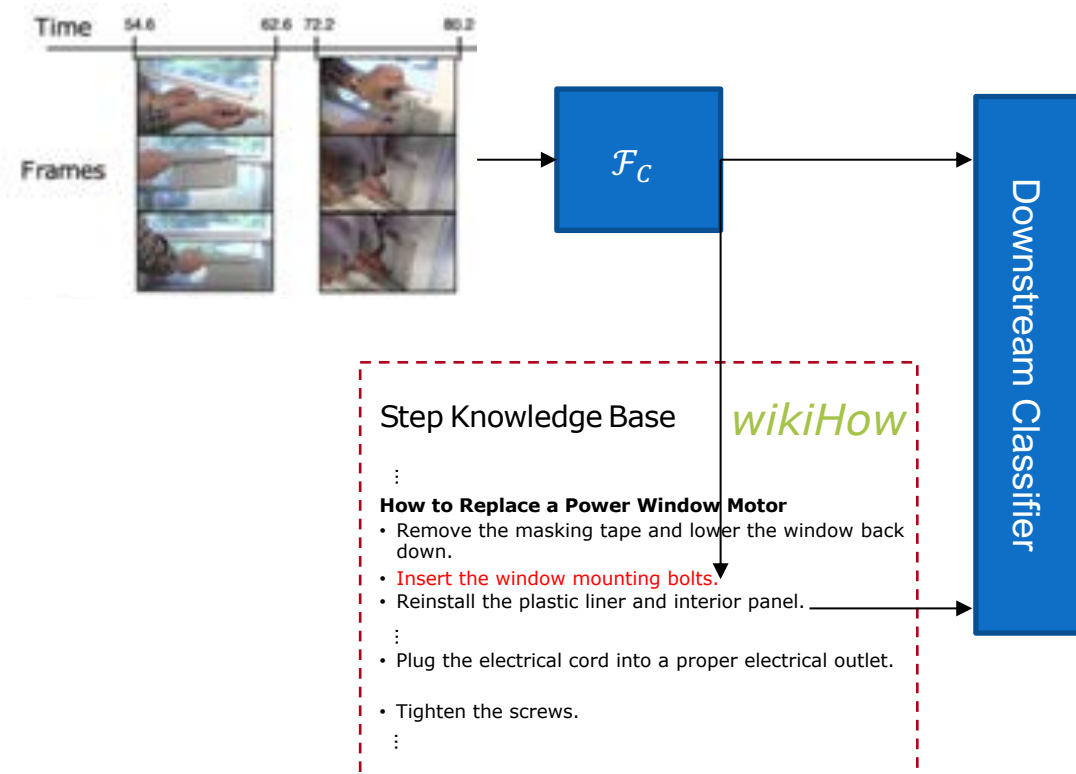
Learning To Recognize Procedural Activities with Distant Supervision



Pretraining: Learning to align videos and the step knowledge base



Finetuning: Training a classifier with both step-level video representation and ordering information from the knowledge base



Learning To Recognize Procedural Activities with Distant Supervision



- Step Forecasting on COIN
 - **Wikihow Knowledge provides high-quality distant supervision!**
 - **Ordering information in the knowledge base further helps!**

Long-term Model	Segment Model	Pretraining Supervision	Pretraining Dataset	Acc (%)
Basic Transformer	S3D [39]	Unsupervised: MIL-NCE on ASR	HT100M	28.1
Basic Transformer	SlowFast [17]	Supervised: action labels	Kinetics	25.6
Basic Transformer	TimeSformer [8]	Supervised: action labels	Kinetics	34.7
Basic Transformer	TimeSformer [8]	Unsupervised: k -means on ASR	HT100M	34.0
Basic Transformer	TimeSformer	Unsupervised: distant supervision (ours)	HT100M	38.2
Transformer w/ KB Transfer	TimeSformer	Unsupervised: distant supervision (ours)	HT100M	39.4

- The supervision from the wikihow knowledge base also helps

Recognition of procedural activities on COIN

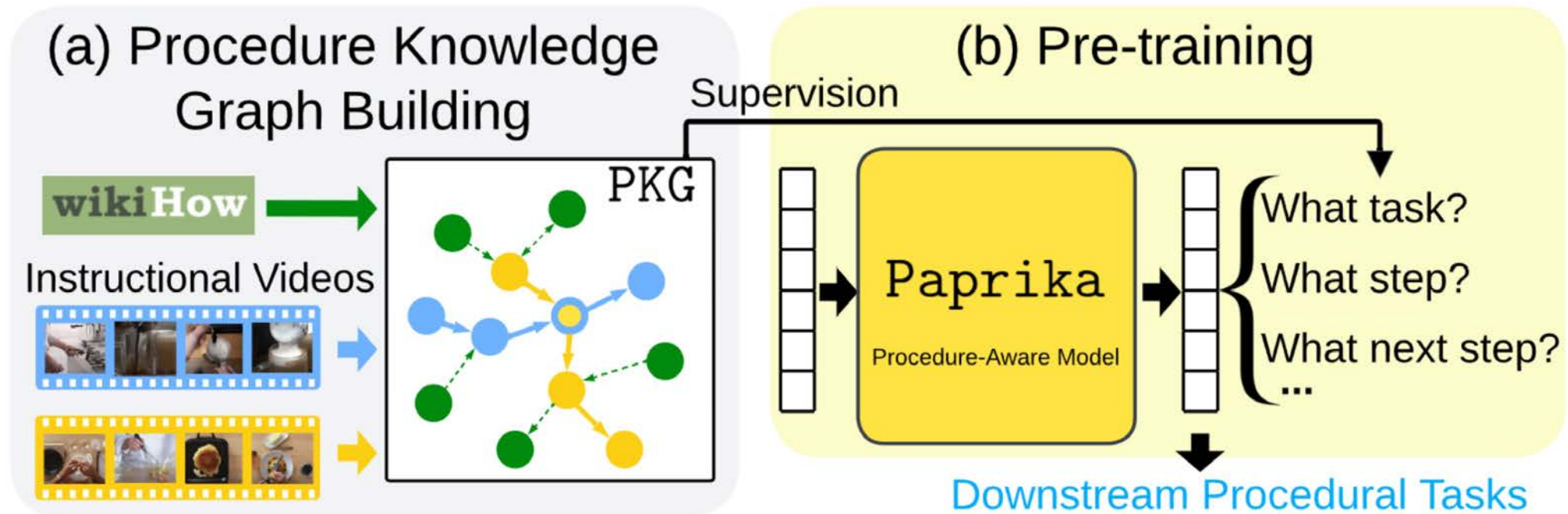
Long-term Model	Segment Model	Pretraining Supervision	Pretraining Dataset	Acc (%)
TSN (RGB+Flow) [57]	Inception [54]	Supervised: action labels	Kinetics	73.4*
Basic Transformer	S3D [39]	Unsupervised: MIL-NCE on ASR	HT100M	70.2*
Basic Transformer	TimeSformer	Unsupervised: distant supervision (ours)	HT100M	88.9
Transformer w/ KB Transfer	TimeSformer	Unsupervised: distant supervision (ours)	HT100M	90.0

Egocentric video classification

Segment Model	Pretraining Supervision	Pretraining Dataset	Action (%)	Verb (%)	Noun (%)
ViViT-L [6]	Supervised: action labels	Kinetics	44.0	66.4	56.8
TimeSformer [8]	Supervised: action labels	Kinetics	42.3	66.6	54.4
TimeSformer	Unsupervised: distant supervision (ours)	HT100M	44.4	67.1	58.1

- Limitation: Didn't employ ordering information in the pretraining model.

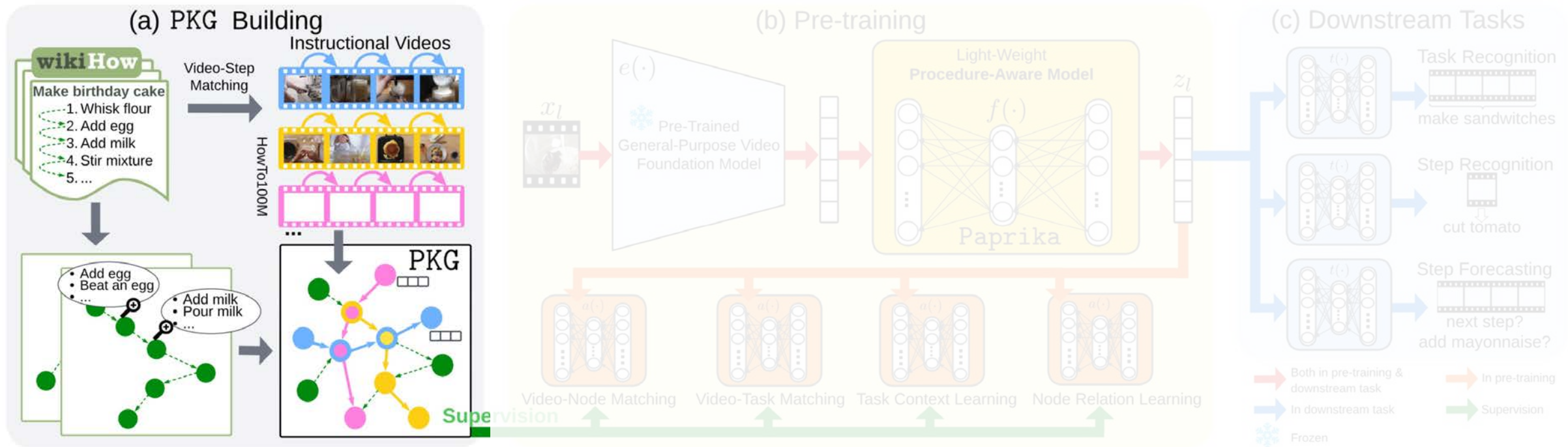
- **Key Idea: Construct procedural knowledge graph and then use it to obtain supervision.**



Procedure-Aware Pretraining for Instructional Video Understanding



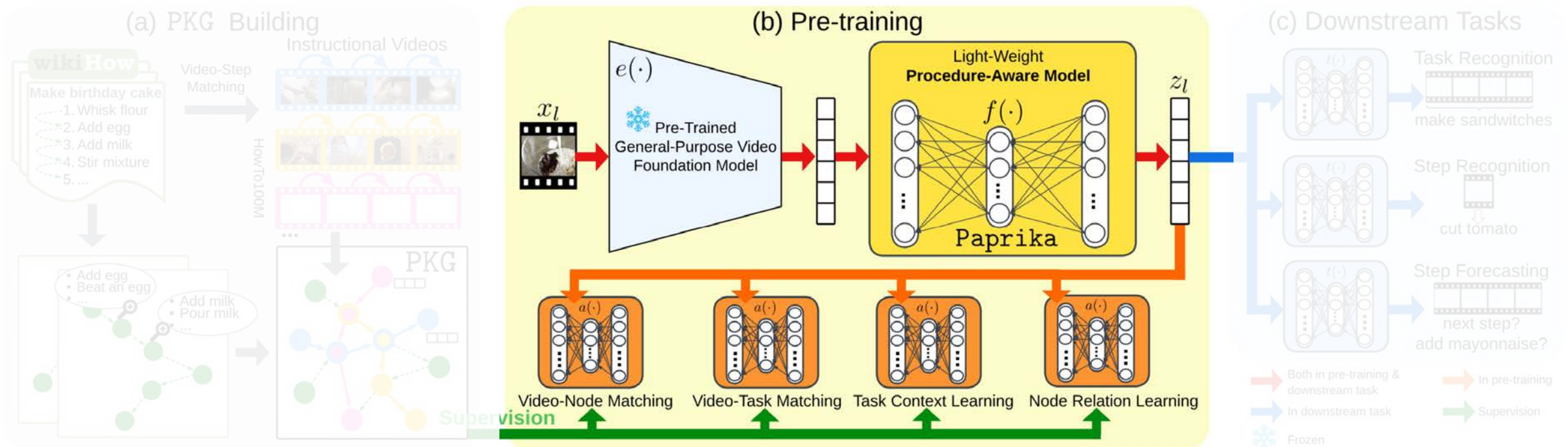
- Construct procedural knowledge graph by grounding wikiHow steps to instructional videos;



Procedure-Aware Pretraining for Instructional Video Understanding



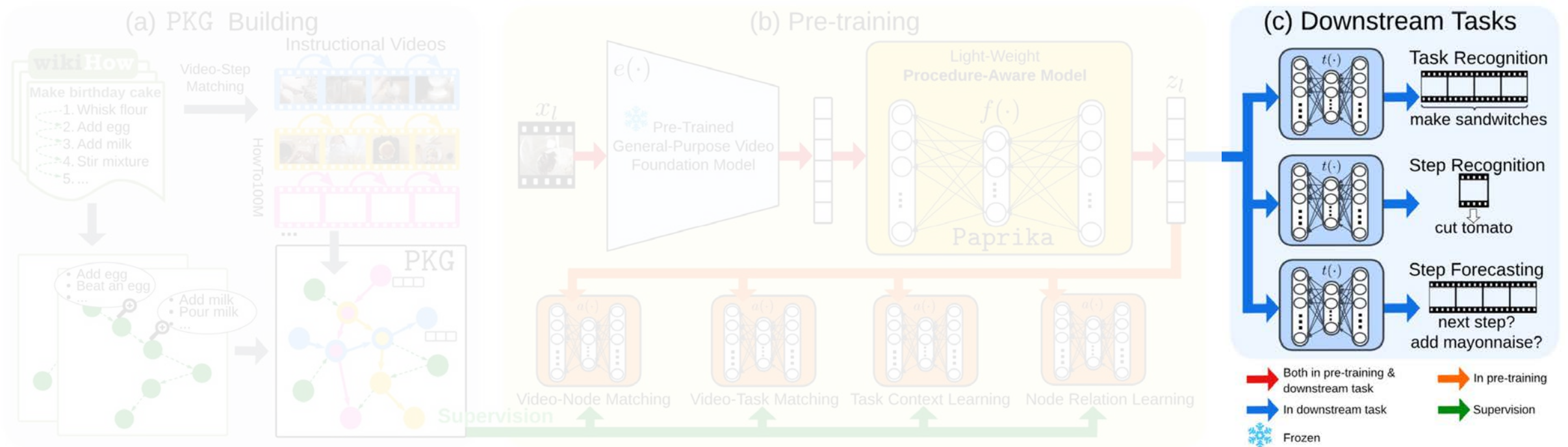
- Use procedural knowledge graph to supervise a procedure-aware model;



Procedure-Aware Pretraining for Instructional Video Understanding



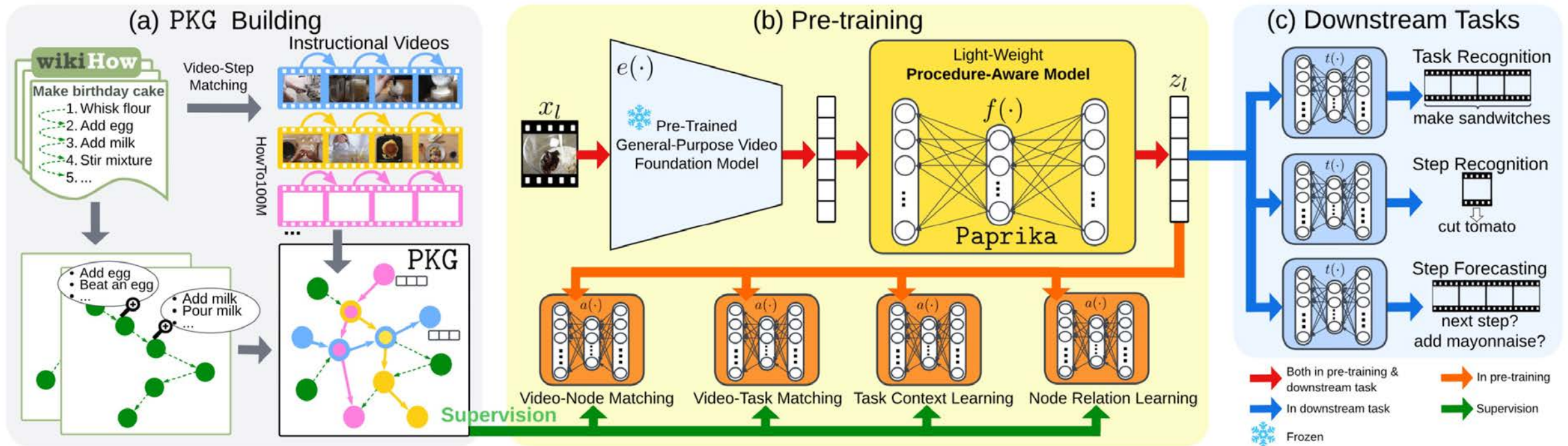
- The representation produced by the procedure-aware model can be directly used for downstream tasks.



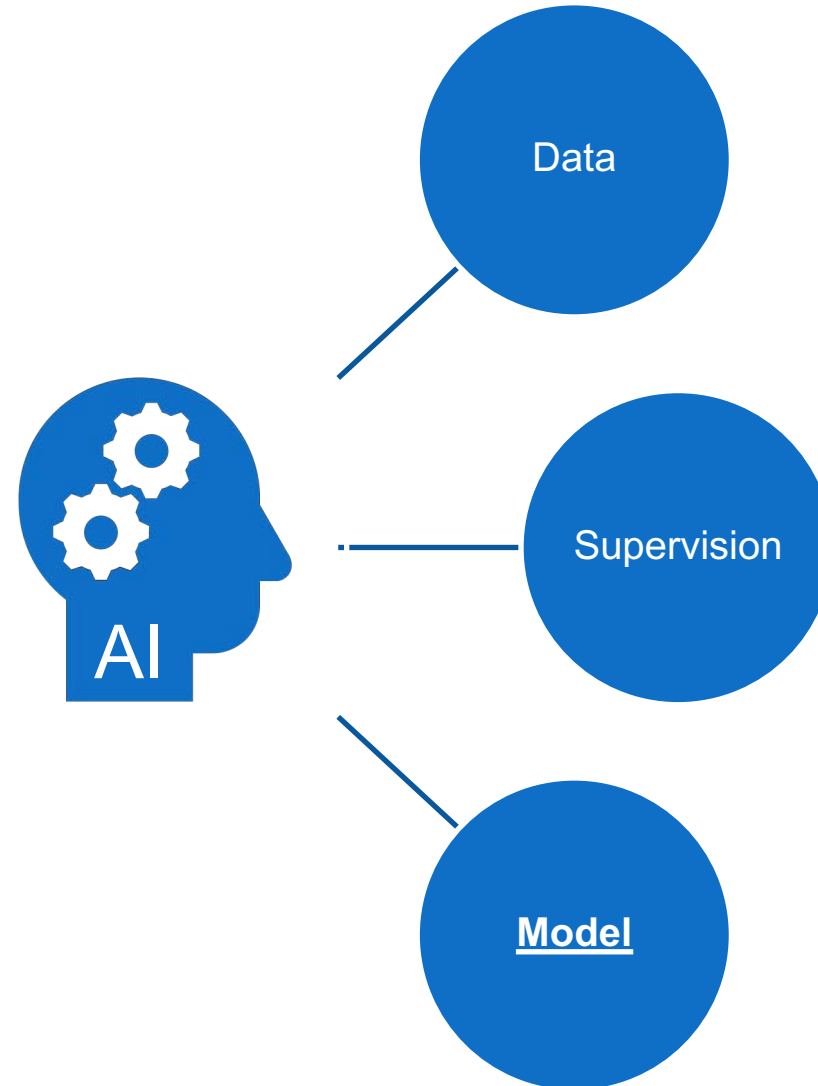
Procedure-Aware Pretraining for Instructional Video Understanding



- Limitation: closed-vocabulary; text-only graph.



How to Utilize the Knowledge Source?



Induce, Edit, Retrieve: Language Grounded Multimodal Schema for Instructional Video Retrieval



- **Key Idea: Learning multimodal schema to represent procedural knowledge.**

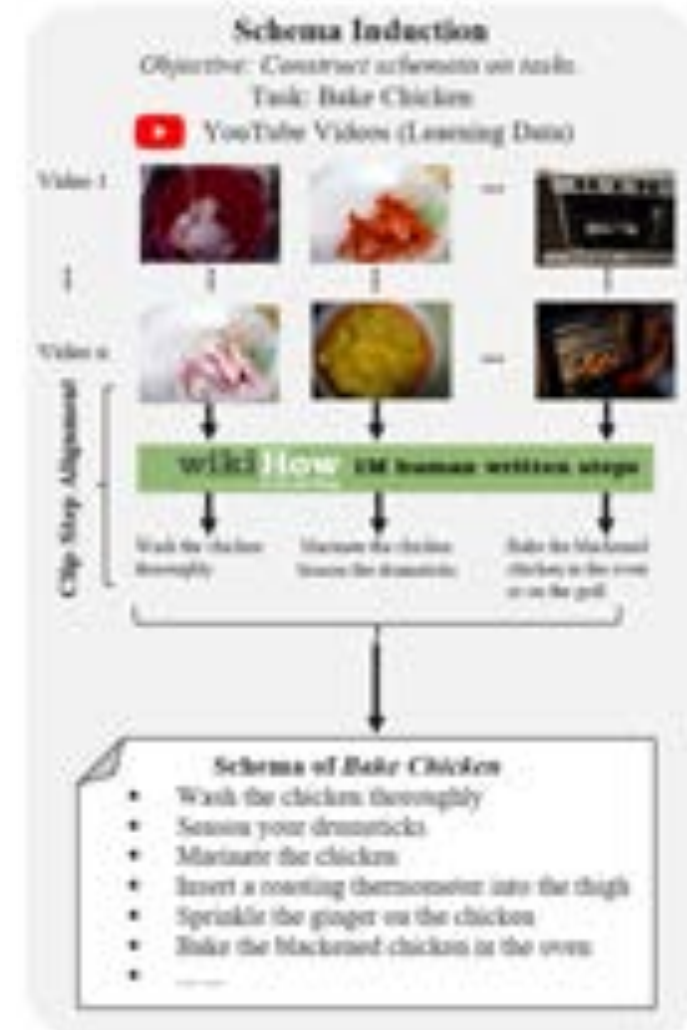


Induce, Edit, Retrieve: Language Grounded Multimodal Schema for Instructional Video Retrieval



• Schema Induction

- For each task, find corresponding steps from wikiHow and videos from YouTube.
- For each segment in each video, retrieve most relevant steps with existing video-text matching models.



Induce, Edit, Retrieve: Language Grounded Multimodal Schema for Instructional Video Retrieval



- Schema Editing

- For an unseen task, find the most similar seen task based on both textual and visual similarity.



Induce, Edit, Retrieve: Language Grounded Multimodal Schema for Instructional Video Retrieval



- Schema Editing

- For an unseen task, find the most similar seen task based on both textual and visual similarity
- Replace object towards the unseen task.



- Schema Editing

- For an unseen task, find the most similar seen task based on both textual and visual similarity
- Replace object towards the unseen task.
- Delete steps that are not relevant in the new task with a pretrained language model.



- **Schema Editing**

- For an unseen task, find the most similar seen task based on both textual and visual similarity
- Replace object towards the unseen task.
- Delete steps that are not relevant in the new task with a pretrained language model.
- Replace tokens least likely associated with the task in each step by prompting a pretrained language model.



Induce, Edit, Retrieve: Language Grounded Multimodal Schema for Instructional Video Retrieval



- The learned schema provides step-level information to better retrieve videos.

Method	Howto-GEN					COIN					Youcook2				
	P@1↑	R@5↑	R@10↑	Med r↓	MRR↑	P@1↑	R@5↑	R@10↑	Med r↓	MRR↑	P@1↑	R@5↑	R@10↑	Med r↓	MRR↑
MIL-NCE [31]	45.2	31.0	43.1	15.0	.198	48.3	37.1	52.8	9.5	.227	27.0	18.2	26.5	32.0	.126
T5 [30]	44.0	29.9	41.0	19.0	.190	46.1	35.3	50.7	10.0	.219	21.3	16.0	24.7	61.5	.108
GPT-2 [39]	46.0	31.5	43.3	16.0	.200	48.9	39.2	53.4	8.0	.233	31.5	19.0	27.3	44.5	.130
GPT-3 [2]	49.3	33.3	45.7	13.0	.211	53.3	42.1	59.0	8.0	.252	37.1	22.4	34.6	27.0	.160
GOSC [30]	54.7	37.0	49.8	11.0	.231	53.9	41.6	55.1	8.0	.248	30.3	20.7	34.8	28.0	.146
wikiHow	51.9	35.4	47.8	11.0	.222	53.9	40.8	56.1	7.0	.246	31.5	21.0	34.2	24.5	.149
IER (Ours)	54.4	37.3	50.1	10.0	.231	57.2	42.2	57.8	7.0	.256	41.6	25.8	38.8	20.0	.175
IER ³ (Ours)	55.0	37.4	50.6	10.0	.234	56.1	42.3	59.1	8.0	.258	40.4	25.1	38.8	20.0	.172
Oracle	56.5	38.0	50.8	10.0	.237	60.0	43.4	59.3	7.0	.262	52.8	33.5	47.1	14.0	.215

Even comparable with oracle (using manual step annotation for each query)



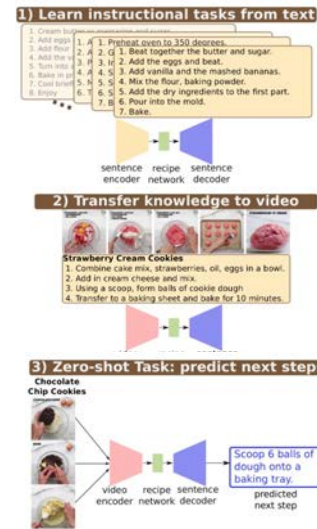
Limitation

- Schema is restricted to step sequence without considering graph structures, e.g., optional/exchangeable steps.
- Only evaluated on text-video retrieval.

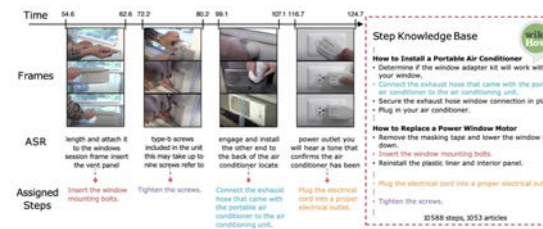
Summary of Methods Using Explicit Knowledge



Sener & Yao ICCV 2019



Lin et al. CVPR 2022



Yang et al. 2021



Knowledge as data	✓		
Knowledge as supervision		✓	
Knowledge for model		✓	✓
Sequential knowledge	✓	✓	✓
Multimodal knowledge			✓

Summary of Methods Using Explicit Knowledge



	Sener & Yao ICCV 2019	Lin et al. CVPR 2022	Yang et al. 2021	Zhou et al. ACL 2023	Zhou et al. CVPR 2023
Knowledge as data	✓			✓	
Knowledge as supervision		✓			✓
Knowledge for model		✓	✓		
Sequential knowledge	✓	✓	✓		
Multimodal knowledge			✓		

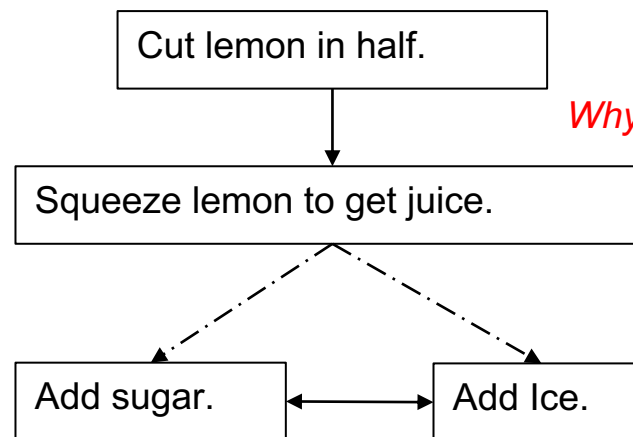
- What is next?

Future Challenge: Interpret but Not Memorize



- Do models understand **why** the steps are ordered as in the knowledge base?

How to make lemonade?



Why these two steps cannot be exchanged?

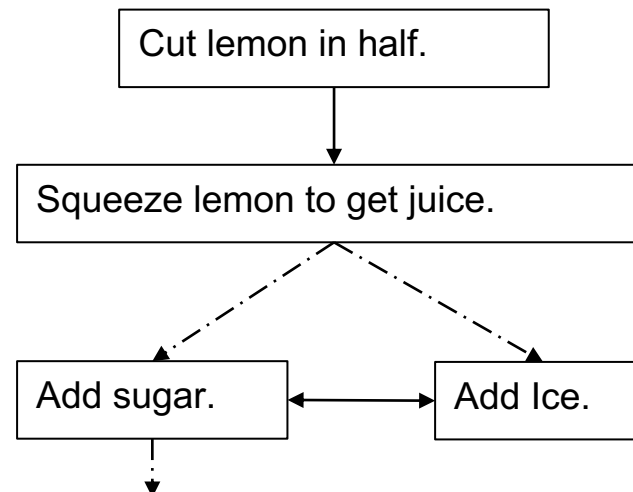
What is the intent of this step?

Open Problem: Open-vocabulary



- Can the knowledge be automatically extended to open-vocabulary?
 - **Generalize to new tasks;**
 - **Discover new steps and add them in the knowledge base...**

How to make lemonade?



Generalize to new tasks:
Make Berrynade

New step
discovered from
video: Carbonate



<https://www.youtube.com/watch?v=GaC14YpDJHw>

Agenda



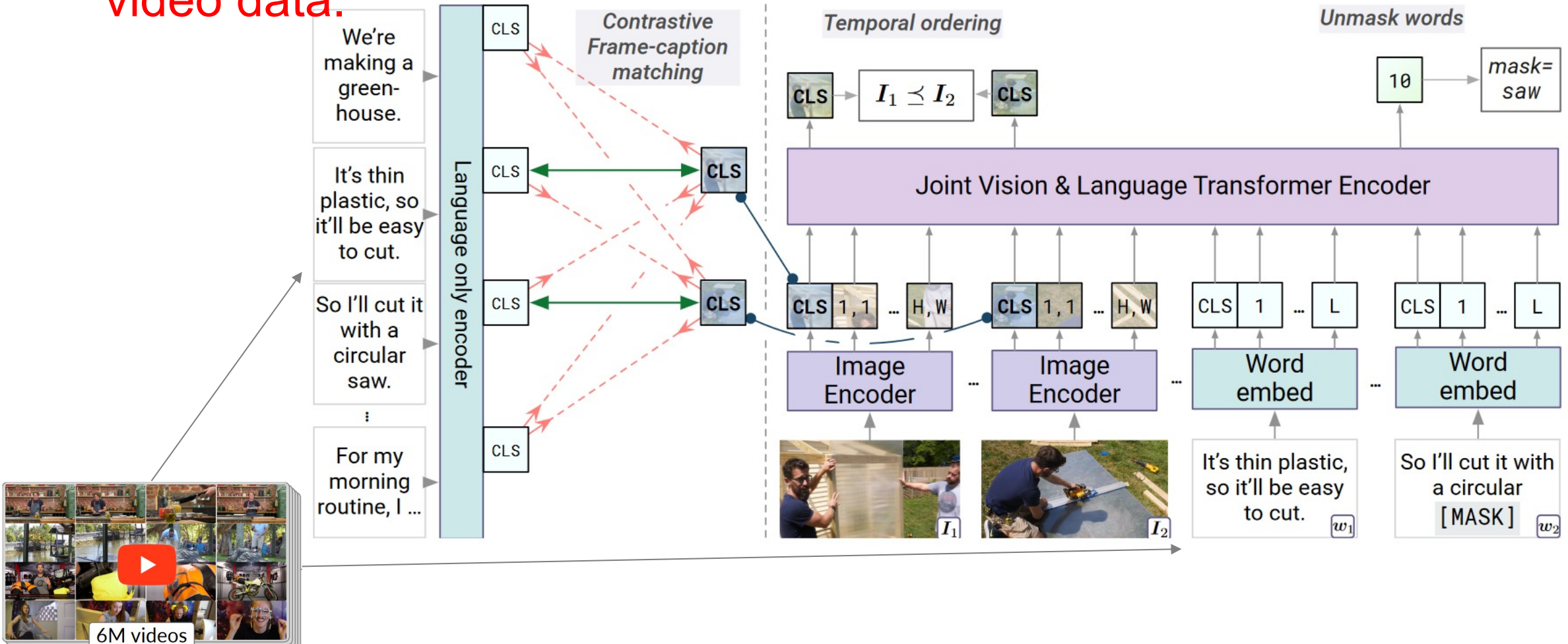
- Explicit Knowledge Source: Learning with the help of external knowledge
- Implicit Knowledge Source: Learning procedural knowledge from data



MERLOT: Multimodal Neural Script Knowledge Models



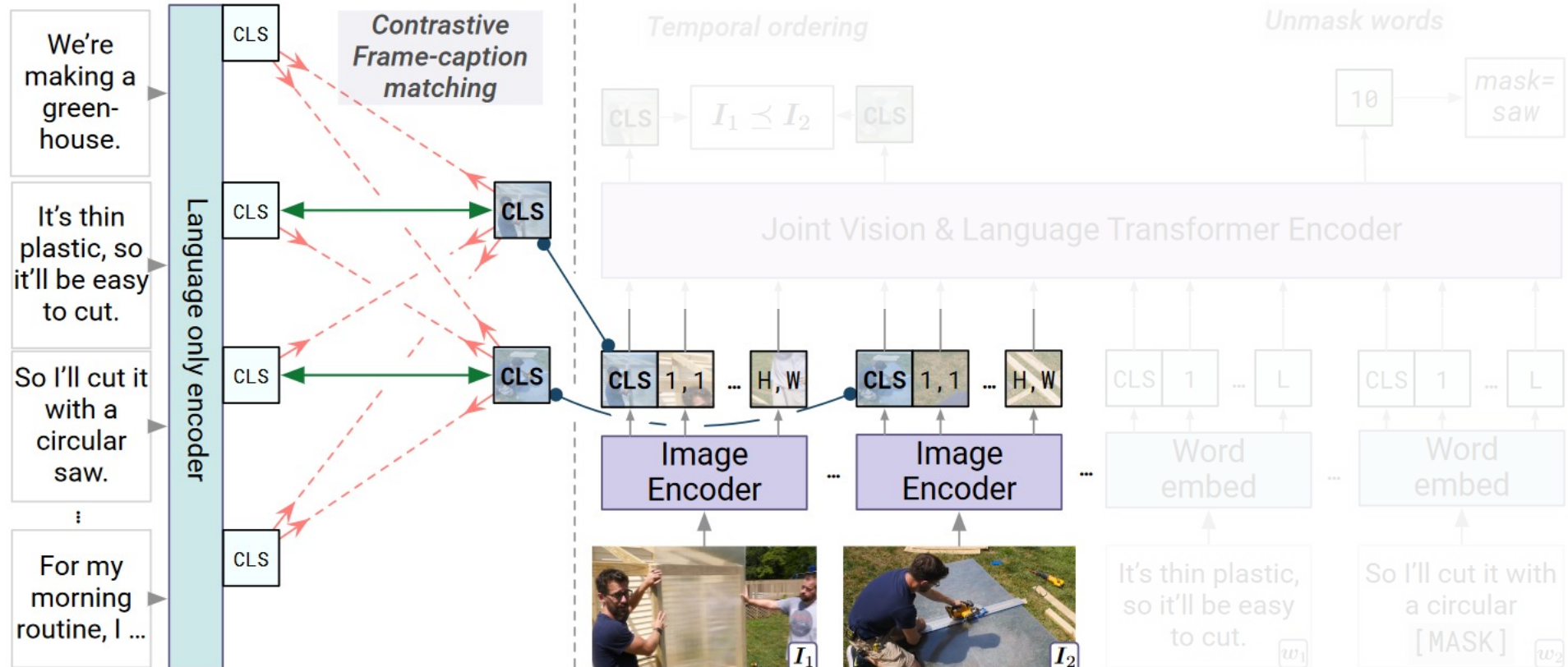
- **Key Idea: Learning temporal reasoning ability through massive video data.**



MERLOT: Multimodal Neural Script Knowledge Models



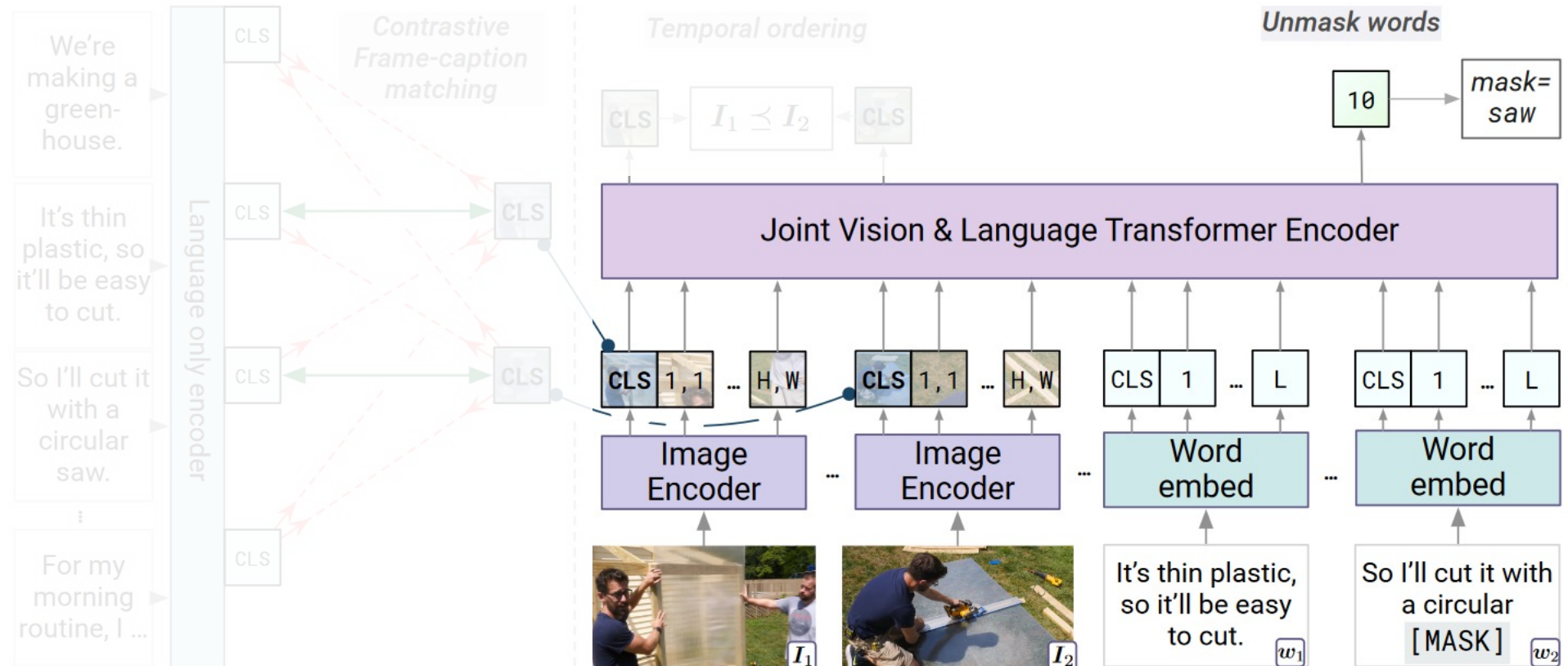
- Objective 1: Alignment between frame representations and text representations



MERLOT: Multimodal Neural Script Knowledge Models



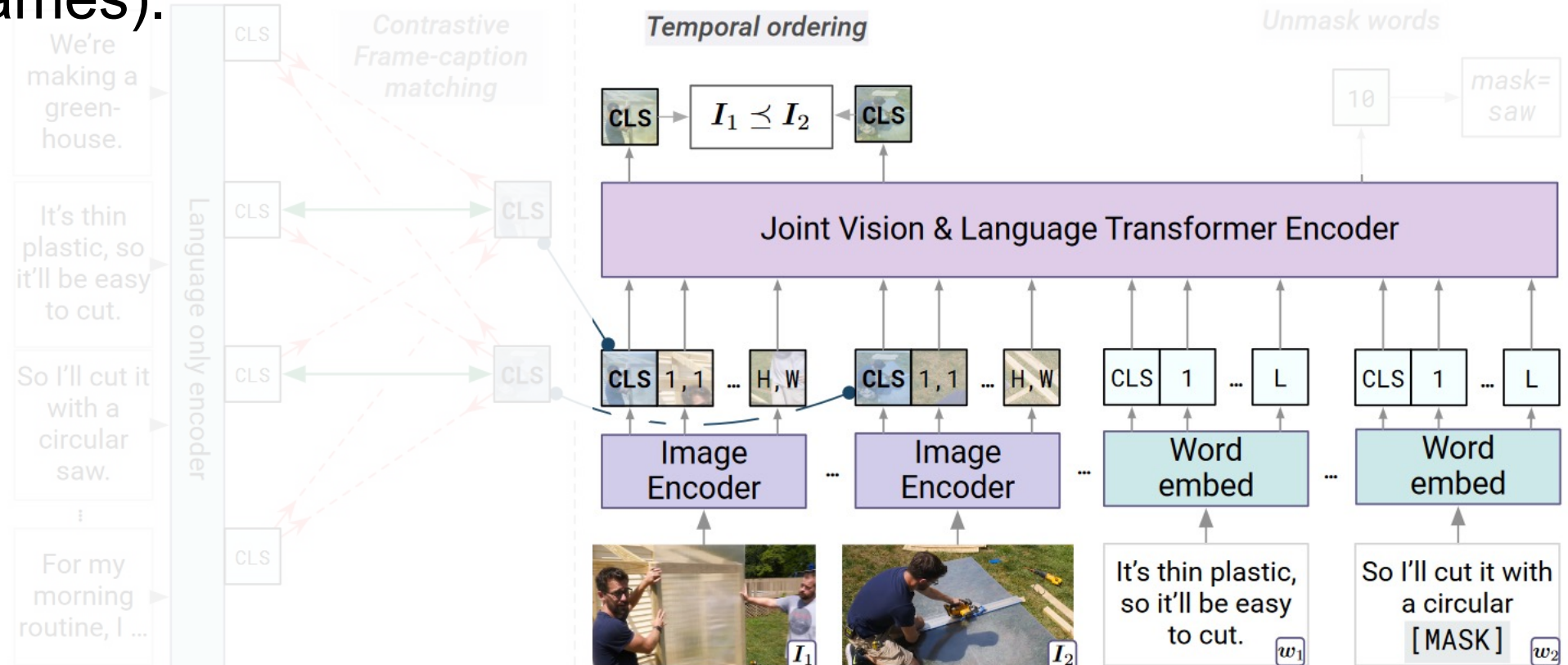
- Objective 2: Masked Token Modeling.



MERLOT: Multimodal Neural Script Knowledge Models



- Objective 3: Temporal Ordering (Binary classification between each pair of frames).



MERLOT: Multimodal Neural Script Knowledge Models



- The model learns strong temporal reasoning ability and joint video-language reasoning ability.

Ordering Images from Visual Stories

	Spearman (↑)	Pairwise acc (↑)	Distance (↓)
CLIP [89]	.609	78.7	.638
UNITER [22]	.545	75.2	.745
MERLOT	.733	84.5	.498

State-of-the-art over various video-language tasks

Tasks	Split	Vid. Length	ActBERT [127]	ClipBERT _{8x2} [67]	SOTA	MERLOT
MSRVTT-QA	Test	Short	-	37.4	41.5 [118]	43.1
MSR-VTT-MC	Test	Short	88.2	-	88.2 [127]	90.9
TGIF-Action	Test	Short	-	82.8	82.8 [67]	94.0
TGIF-Transition	Test	Short	-	87.8	87.8 [67]	96.2
TGIF-Frame QA	Test	Short	-	60.3	60.3 [67]	69.5
LSMDC-FiB QA	Test	Short	48.6	-	48.6 [127]	52.9
LSMDC-MC	Test	Short	-	-	73.5 [121]	81.7
ActivityNetQA	Test	Long	-	-	38.9 [118]	41.4
Drama-QA	Val	Long	-	-	81.0 [56]	81.4
TVQA	Test	Long	-	-	76.2 [56]	78.7
TVQA+	Test	Long	-	-	76.2 [56]	80.9
VLEP	Test	Long	-	-	67.5 [66]	68.4



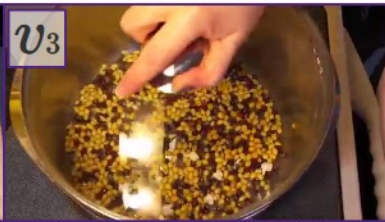

↓
Predict future event given historical videos

- Limitation: short temporal span; importance of the temporal ordering loss is unclear.


MERLOT Reserve: Neural Script Knowledge through Vision and Language and Sound



- **Key Idea: Jointly learn script knowledge with video, language and audio.**

			
w_1 Add a third of a cup of popcorn	w_2 Now turn the heat on high	w_3 Add a lid, and then	[MASKed span]
a_1 *pouring sound*	a_2 *sizzling*	a_3 *lid clinking*	

...



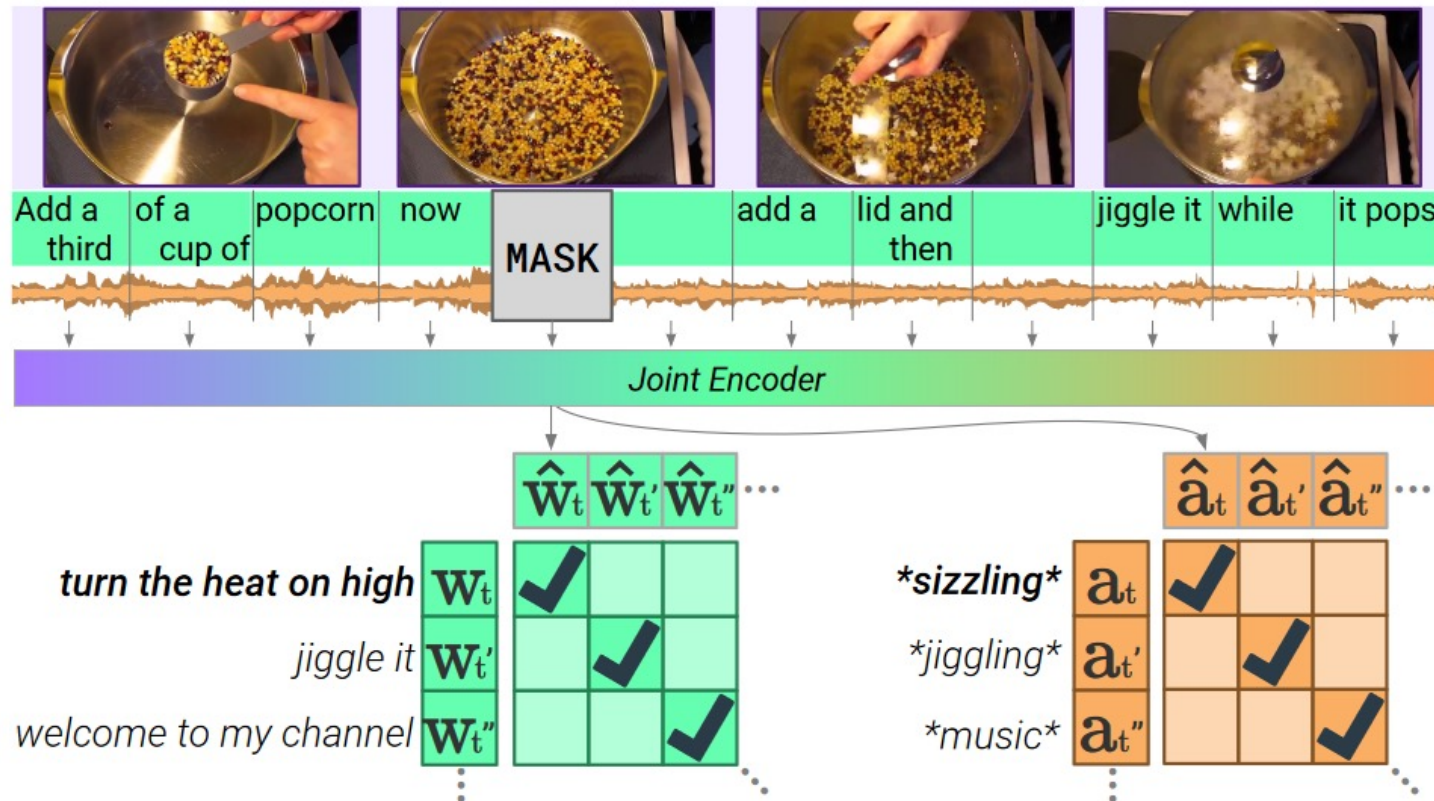
w_4 jiggle it while it pops

a_4 *jiggling, popcorn popping*

MERLOT Reserve: Neural Script Knowledge through Vision and Language and Sound



- Key objective design: contrastive loss between predicted and actual representation of the masked audio/text



MERLOT Reserve: Neural Script Knowledge through Vision and Language and Sound



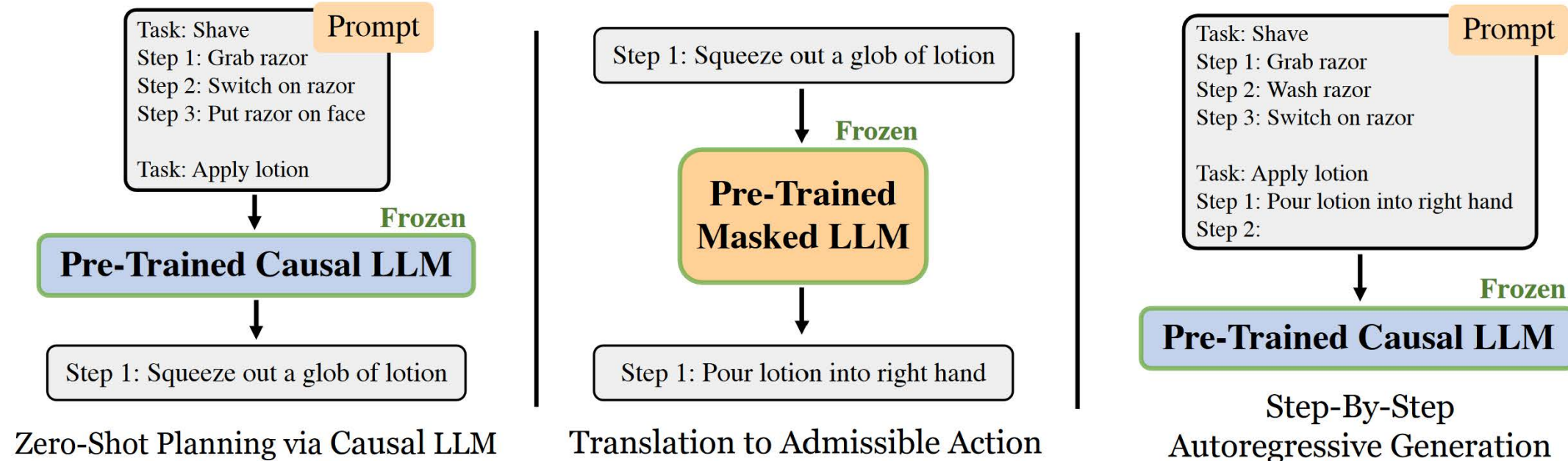
- Audio brings extra supervision and information towards stronger video understanding and video-language performance.
- Limitation: improvement on learned procedural knowledge may be less significant.

Requiring understand procedures of actions/objects

Model	Kinetics-600 (%)	
	Top-1	Top-5
VATT-Base[2]	80.5	95.5
VATT-Large [2]	83.6	96.6
TimeSFormer-L [9]	82.2	95.6
Florence [125]	87.8	97.8
MTV-Base [122]	83.6	96.1
MTV-Large [122]	85.4	96.7
MTV-Huge [122]	89.6	98.3
RESERVE-B	88.1	95.8
RESERVE-L	89.4	96.3
+Audio		
RESERVE-B	89.7	96.6
RESERVE-L	91.1	97.1

Model	(test acc; %)				Overall
	Interaction	Sequence	Prediction	Feasibility	
Supervised SoTA	39.8	43.6	32.3	31.4	36.7
Random	25.0	25.0	25.0	25.0	25.0
CLIP (ViT-B/16) [92]	39.8	40.5	35.5	36.0	38.0
CLIP (RN50x16) [92]	39.9	41.7	36.5	37.0	38.7
Just Ask (ZS)[123]					
zero-shot					
RESERVE-B	44.4	40.1	38.1	35.0	39.4
RESERVE-L	42.6	41.1	37.4	32.2	38.3
RESERVE-B (+audio)	44.8	42.4	38.8	36.2	40.5
RESERVE-L (+audio)	43.9	42.6	37.6	33.6	39.4


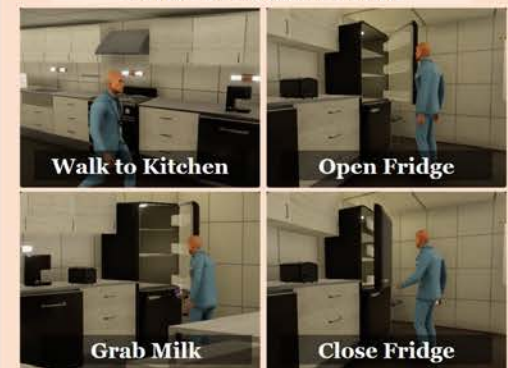
- **Key Idea:** Large language models learn rich procedural knowledge and such knowledge could be extracted.



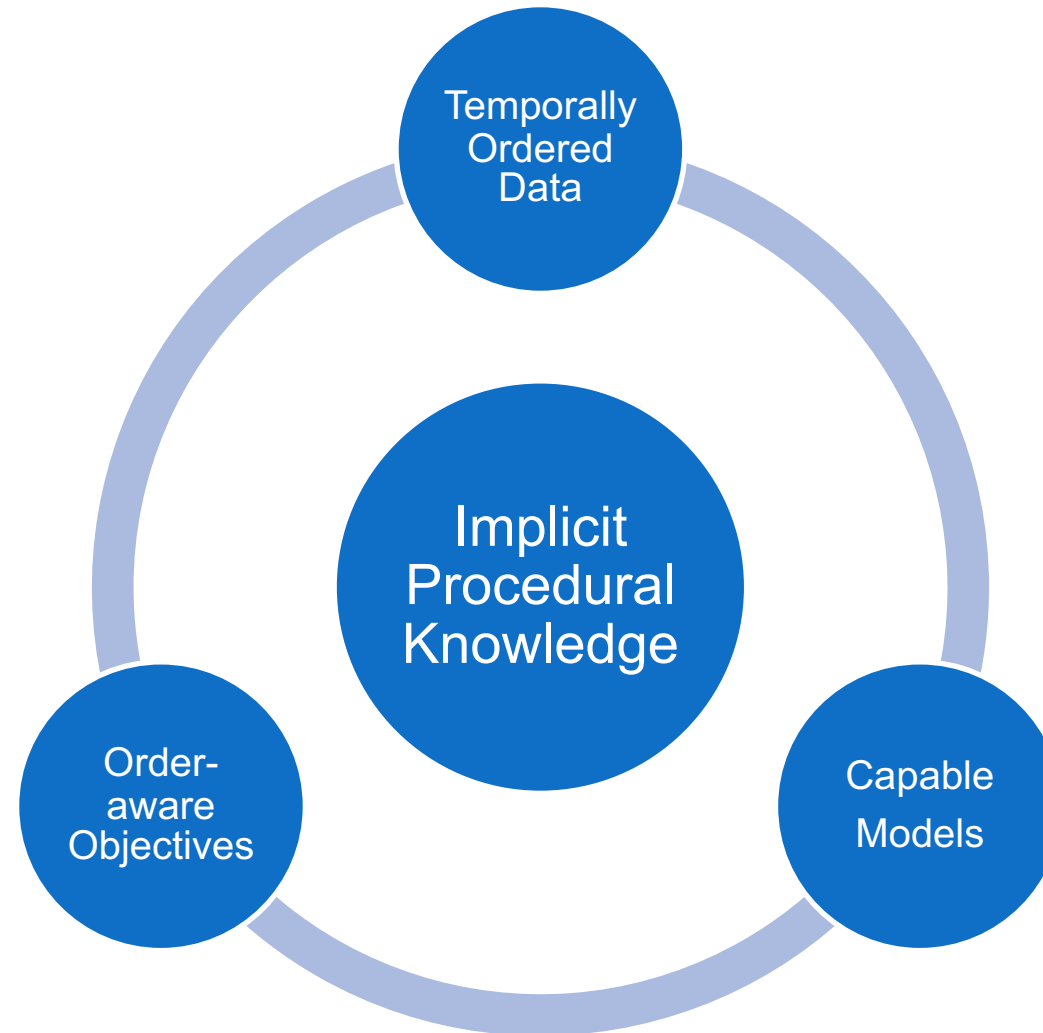
Language models as zero-shot planners: Extracting actionable knowledge for embodied agents



- Impressive results.
- Challenge: verification; groundability to real-world videos.

Human	GPT-3 175B	Translated GPT-3 175B	Task: Complete Amazon Turk Surveys
Task: Throw away paper Step 1: Walk to home office Step 2: Walk to desk Step 3: Find desk Step 4: Turn to desk Step 5: Find chair Step 6: Sit on chair Step 7: Find check Step 8: Grab check Step 9: Squeeze check Step 10: Stand up Step 11: Walk to trashcan Step 12: Put check on trashcan	Task: Throw away paper Step 1: Walk to home office Step 2: Walk to wastebasket Step 3: Find wastebasket Step 4: Grab wastebasket Step 5: Walk to desk Step 6: Drop paper in wastebasket	Task: Throw away paper Step 1: Walk to home office Step 2: Walk to table Step 3: Find table Step 4: Turn to table Step 5: Find paper Step 6: Grab paper Step 7: Walk to trashcan Step 8: Open trashcan Step 9: Put paper on trashcan Step 10: Close trashcan	
GPT-2 1.5B	Codex 12B	Translated Codex 12B	Task: Get Glass of Milk
Task: Brush teeth Step 1: Go to bathroom	Task: Brush teeth Step 1: Walk to bathroom Step 2: Walk to sink Step 3: Find toothbrush Step 4: Pick up toothbrush Step 5: Put toothbrush in mouth Step 6: Move brush around mouth for two minutes Step 7: Spit out toothpaste and brush into sink Step 8: Turn on water in sink and rinse brush for one minute Step 9: Turn off water in sink and return brush to cupboard	Task: Brush teeth Step 1: Walk to bathroom Step 2: Open door Step 3: Walk to sink Step 4: Put pot on sink Step 5: Put brush on toothbrush Step 6: Turn to toothpaste Step 7: Put toothpaste on toothbrush Step 8: Put teeth on toothbrush	

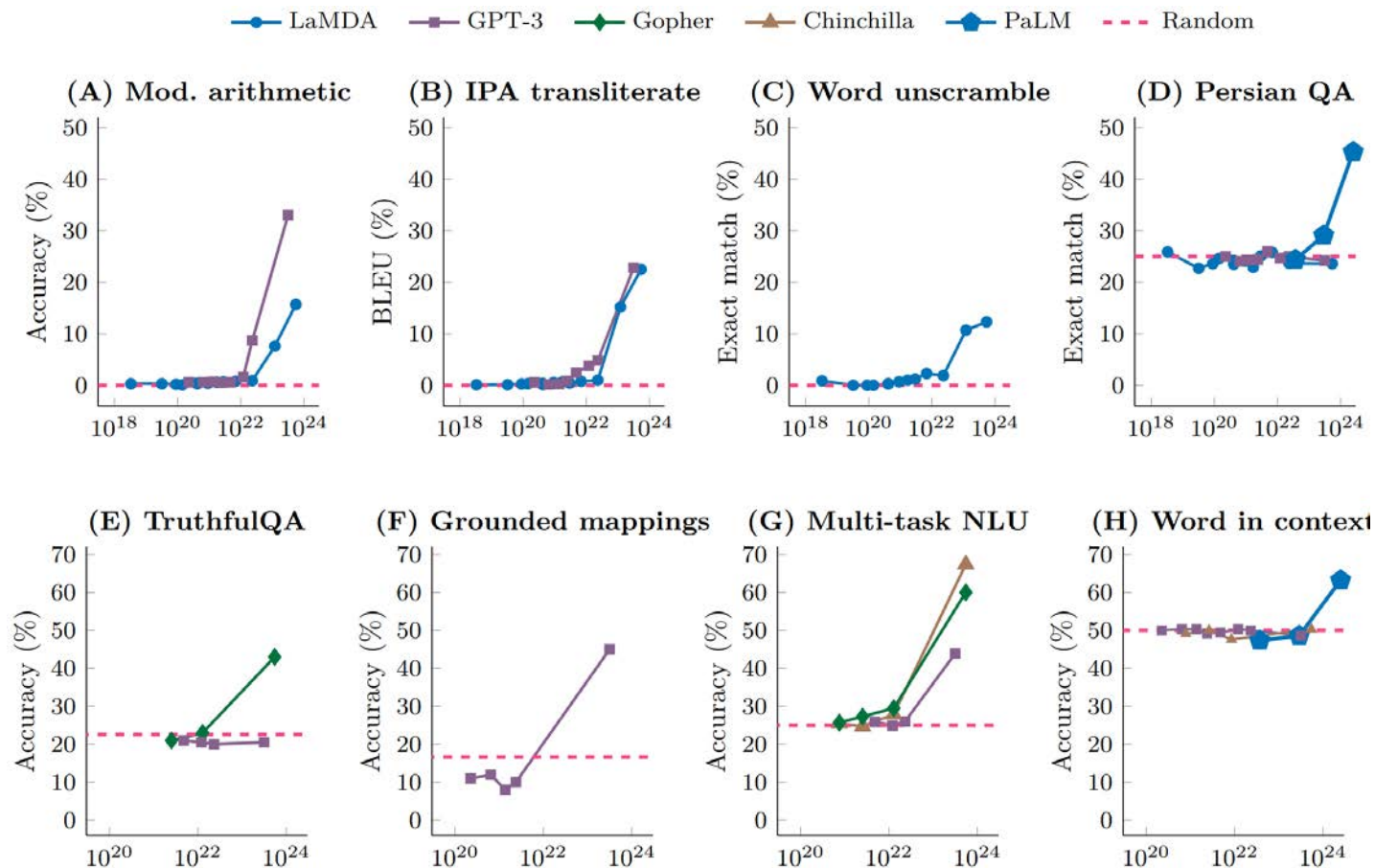
Summary of Methods Learning Implicit Knowledge



Future Challenge: Is there a critical point on scale?



- Can models learn procedural knowledge with a limited scale?

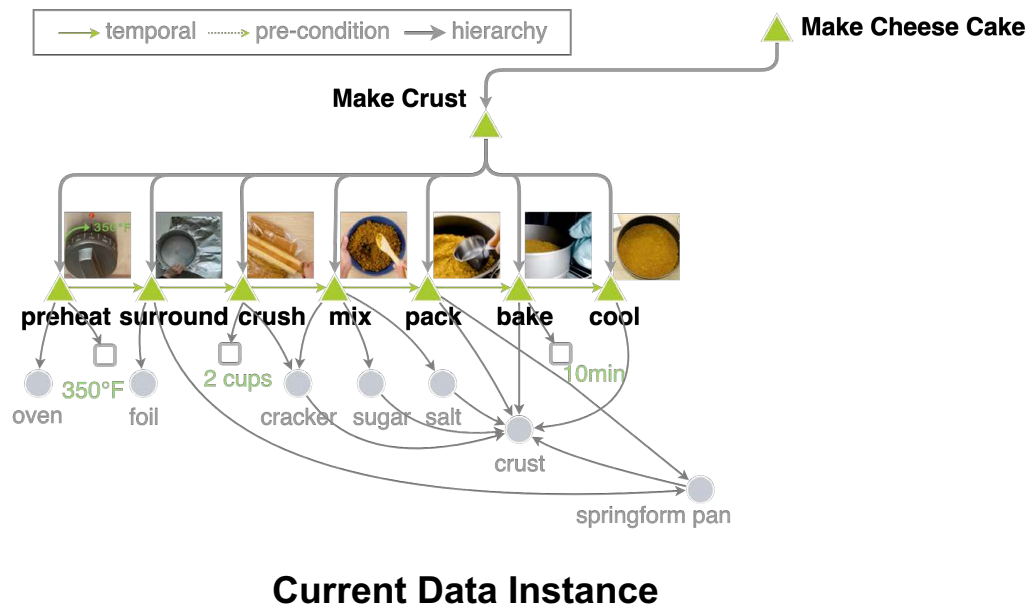


Many reasoning ability of large language models emerge when the model scale is larger than a critical point.

Future Challenge: From an instance to a set



- Can models learn from temporally ordered sets of instances?




Here are the NYC zip codes targeted for new COVID-19 lockdown

By *Kate Sheehy* October 4, 2020 | 4:16pm | Updated


| Sign up for our *special edition newsletter* to get a daily update on the coronavirus pandemic.

There are nine COVID-riddled neighborhoods that New York City *wants to return to lockdown* Wednesday.




Amid Ongoing COVID-19 Pandemic, Governor Cuomo Issues Executive Order Requiring All People in New York to Wear Masks or Face Coverings in Public

APRIL 15, 2020 | Albany, NY



FDA Takes Key Action in Fight Against COVID-19 By Issuing Emergency Use Authorization for First COVID-19 Vaccine

Action Follows Thorough Evaluation of Available Safety, Effectiveness, and Manufacturing Quality Information by FDA Career Scientists, Input from Independent Experts



...

Real-world complex task

https://en.wikipedia.org/wiki/COVID-19_pandemic_in_New_York_City#cite_note-48

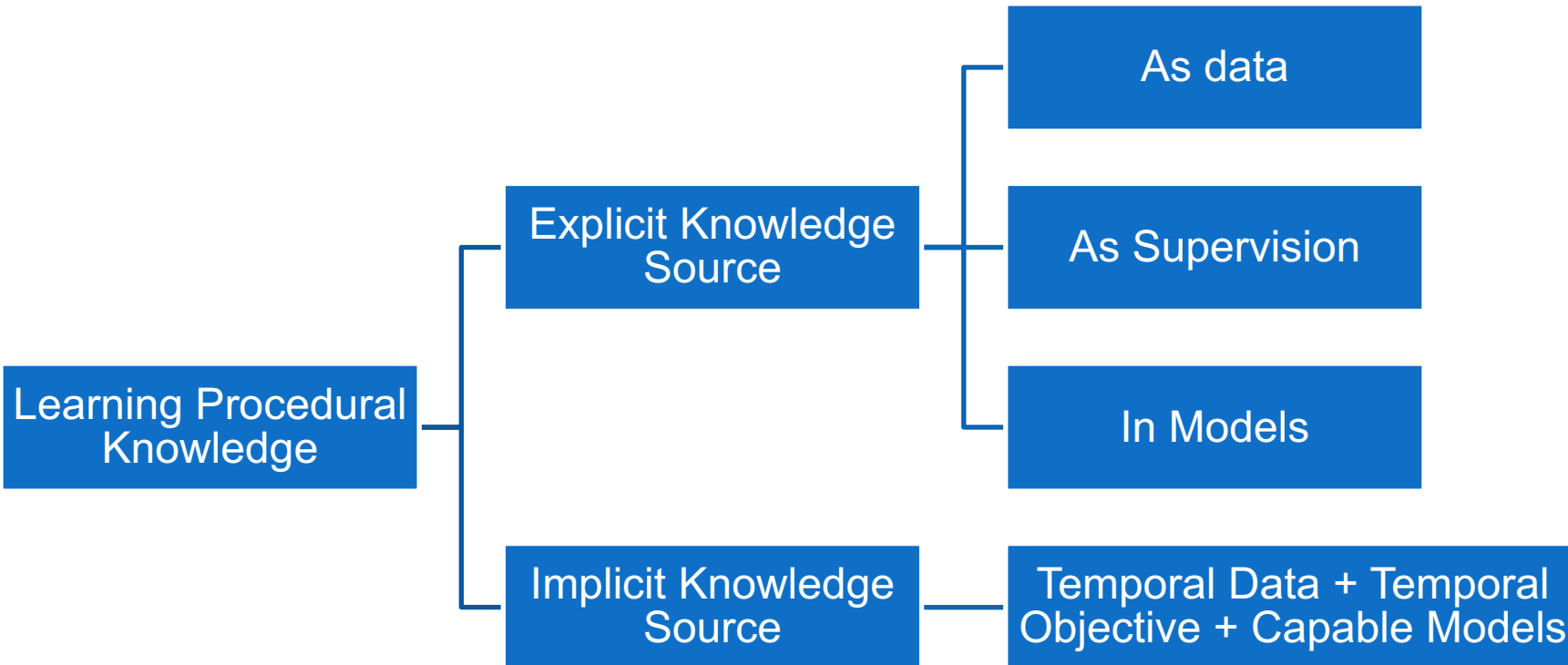
<https://nypost.com/2020/10/04/here-are-the-nyc-zip-codes-targeted-for-new-covid-19-lockdown/>

<https://www.governor.ny.gov/news/amid-ongoing-covid-19-pandemic-governor-cuomo-issues-executive-order-requiring-all-people-new>

<https://www.fda.gov/news-events/press-announcements/fda-takes-key-action-fight-against-covid-19-issuing-emergency-use-authorization-first-covid-19>

<https://www.nature.com/articles/d41586-020-02684-9>

Take-away Messages



Future Challenges:

- Interpret but not memorize;
- Towards open-vocabulary;
- Learning with limited scale;
- From an instance to a set;
- ...

Jun 2023

CVPR Tutorials

Knowledge-Driven Vision-Language Encoding

CVPR

Cross-Modal Knowledge Transfer

Knowledge-Driven Vision-Language Encoding (Part V)

Jie Lei

Meta AI

jielei@meta.com



Northwestern
University



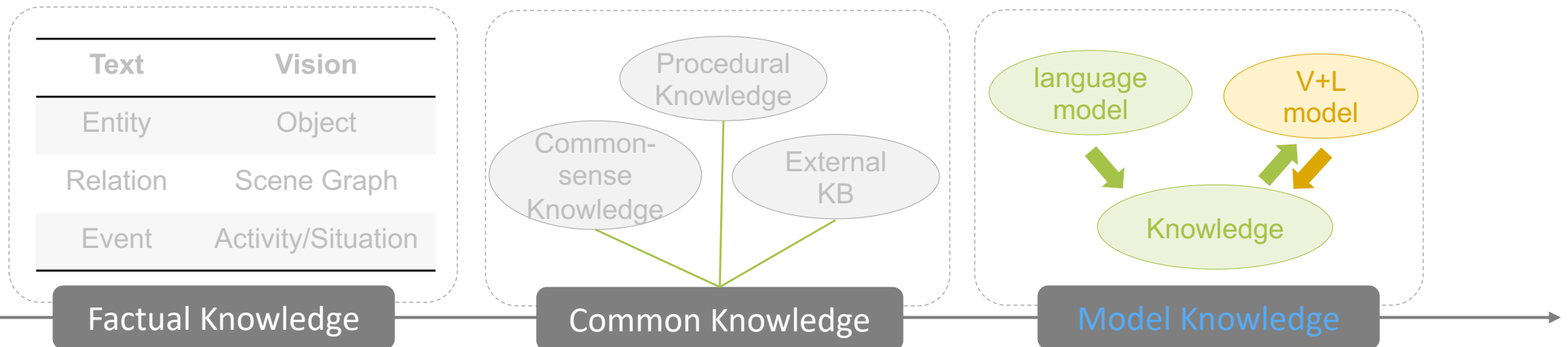
COLUMBIA
UNIVERSITY



Overview



Compared to raw data, knowledge is **important and useful information**.



Part 1. Language knowledge helps learn better vision models

- Pure vision tasks: object detection, image classification, etc.
- Multimodal tasks with vision signals: VQA, video captioning, etc.



Part 2. Vision knowledge helps learn better language models

- Human learn language by connecting the words to their visual appearance in the surrounding world.

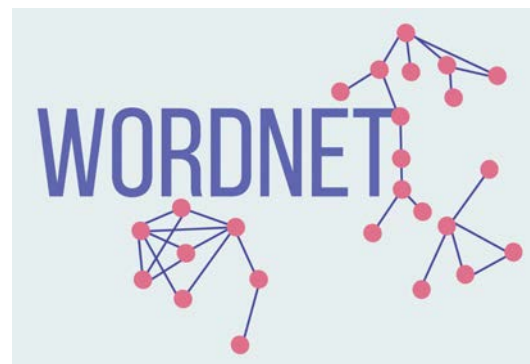
Part 1. Language → Vision





- **Implicit knowledge** from pre-trained Language Models (LM)



- **Explicit knowledge** from human curated sources (e.g., wiki) or model generated knowledge (e.g., GPT-3 generated category definitions)

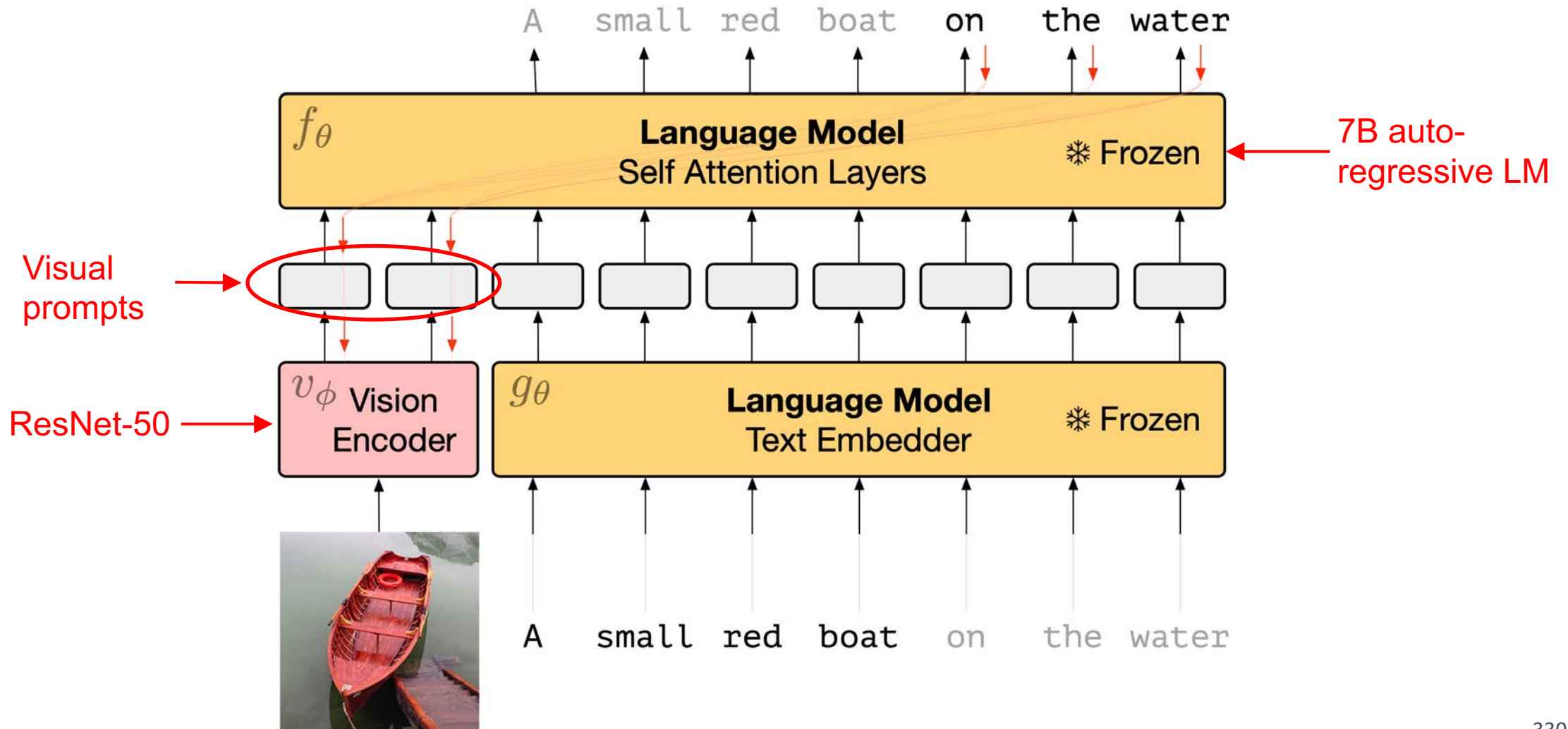


Concept name: snowberg
 **Def_wik:** None
 **GPT3 Query:**
Please explain the concept according to the context.
===
Q: ship
A: A water-borne vessel generally larger than a boat.










Part 1.1 Implicit Knowledge from Language Models

Frozen

- Preserve LM ability by **freezing** it during cross-modal model training.
- Gradient: frozen LM \rightarrow vision encoder



- Few-shot multimodal in-context learning after trained on 3M image-text pairs.

	This person is like 😊.		This person is like 😞.		This person is like	Model Completion	🤖. <EOS>
	This was invented by Zacharias Janssen.		This was invented by Thomas Edison.		This was invented by	Model Completion	the Wright brothers. <EOS>
	With one of these I can drive around a track, overtaking other cars and taking corners at speed		With one of these I can take off from a city and fly across the sky to somewhere on the other side of the world		With one of these I can	Model Completion	break into a secure building, unlock the door and walk right in <EOS>

Wiki knowledge

- Reasonably good zero/few-shot performance, but still underperform SOTA: limited multimodal data? (3M); LM is relatively small? (7B)

n-shot Acc.	n=0	n=1	n=4	τ
Frozen	29.5	35.7	38.2	✗
<i>Frozen</i> scratch	0.0	0.0	0.0	✗
<i>Frozen</i> finetuned	24.0	28.2	29.2	✗
<i>Frozen</i> train-blind	26.2	33.5	33.3	✗
<i>Frozen</i> VQA	48.4	–	–	✓
<i>Frozen</i> VQA-blind	39.1	–	–	✓
Oscar [23]	73.8	–	–	✓

n-shot Acc.	n=0	n=1	n=4	τ
Frozen	5.9	9.7	12.6	✗
<i>Frozen</i> 400mLM	4.0	5.9	6.6	✗
<i>Frozen</i> finetuned	4.2	4.1	4.6	✗
<i>Frozen</i> train-blind	3.3	7.2	0.0	✗
<i>Frozen</i> VQA	19.6	–	–	✗
<i>Frozen</i> VQA-blind	12.5	–	–	✗
MAVEx [42]	39.4	–	–	✓

Large gap w/ SOTA

Flamingo

- A frozen 70B pre-trained LM + a frozen pre-trained ResNet.
- Trained w/ image/video-text pairs, along with interleaved image-text data (M3W), which is important for in-context learning.

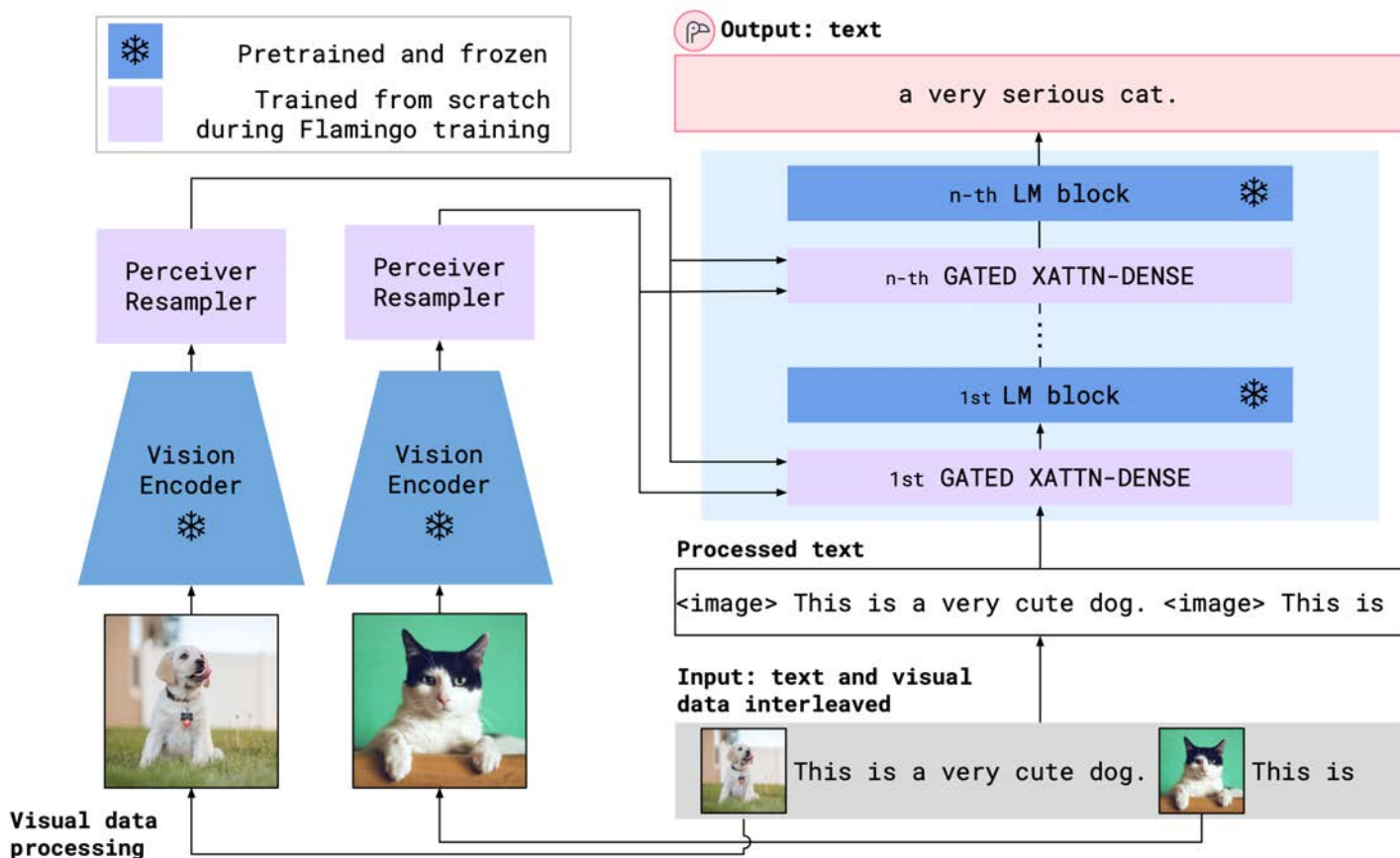


Image-Text Pairs dataset
 $[N=1, T=1, H, W, C]$
 ALIGN: 1.8B +
 LTIP: 312M images



Video-Text Pairs dataset
 $[N=1, T>1, H, W, C]$
 VTP: 27M videos



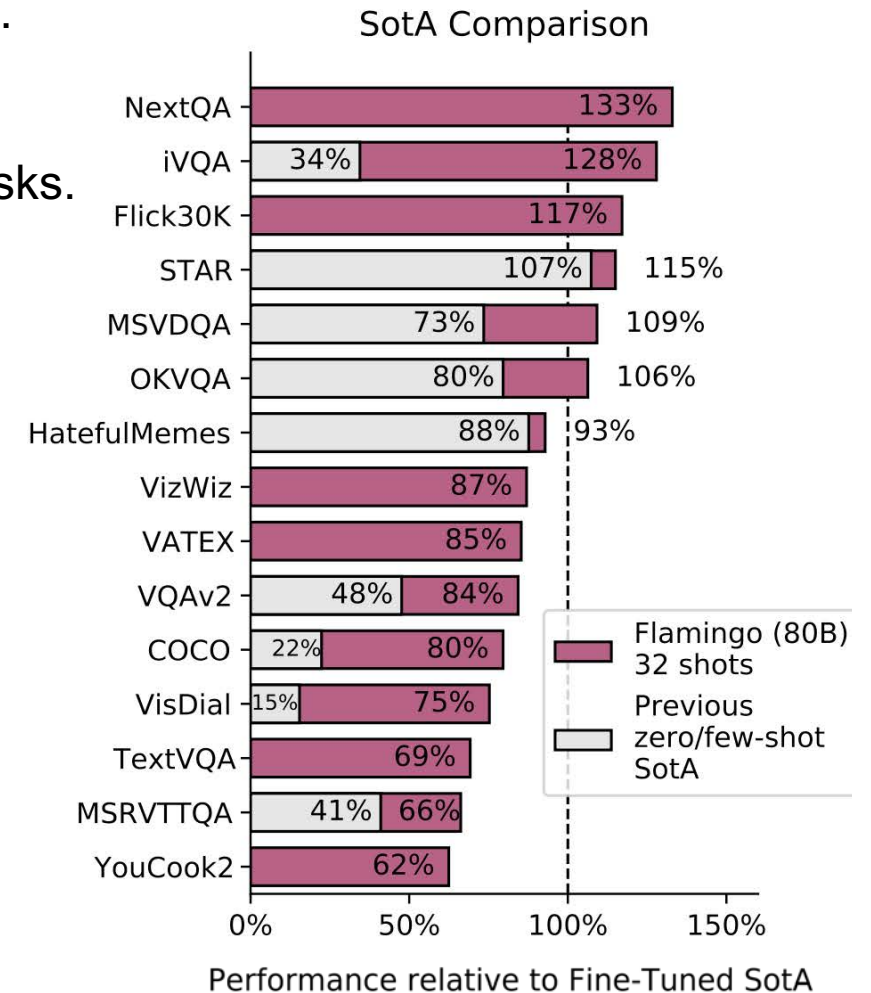
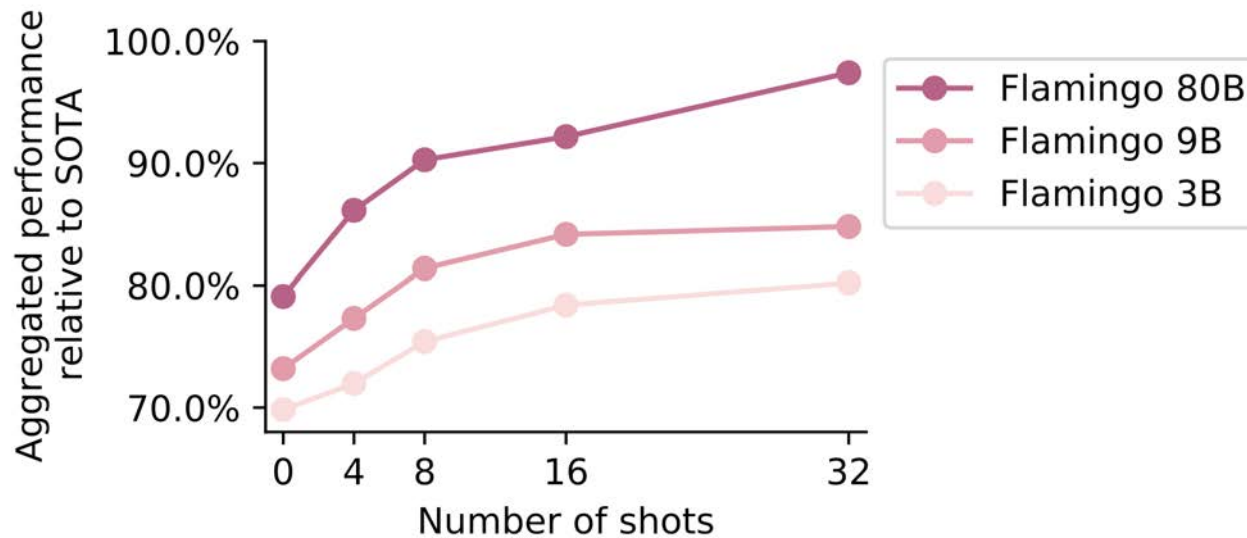
Multi-Modal Massive Web (M3W) dataset
 $[N>1, T=1, H, W, C]$

M3W: 43M webpages (185M images)

Flamingo



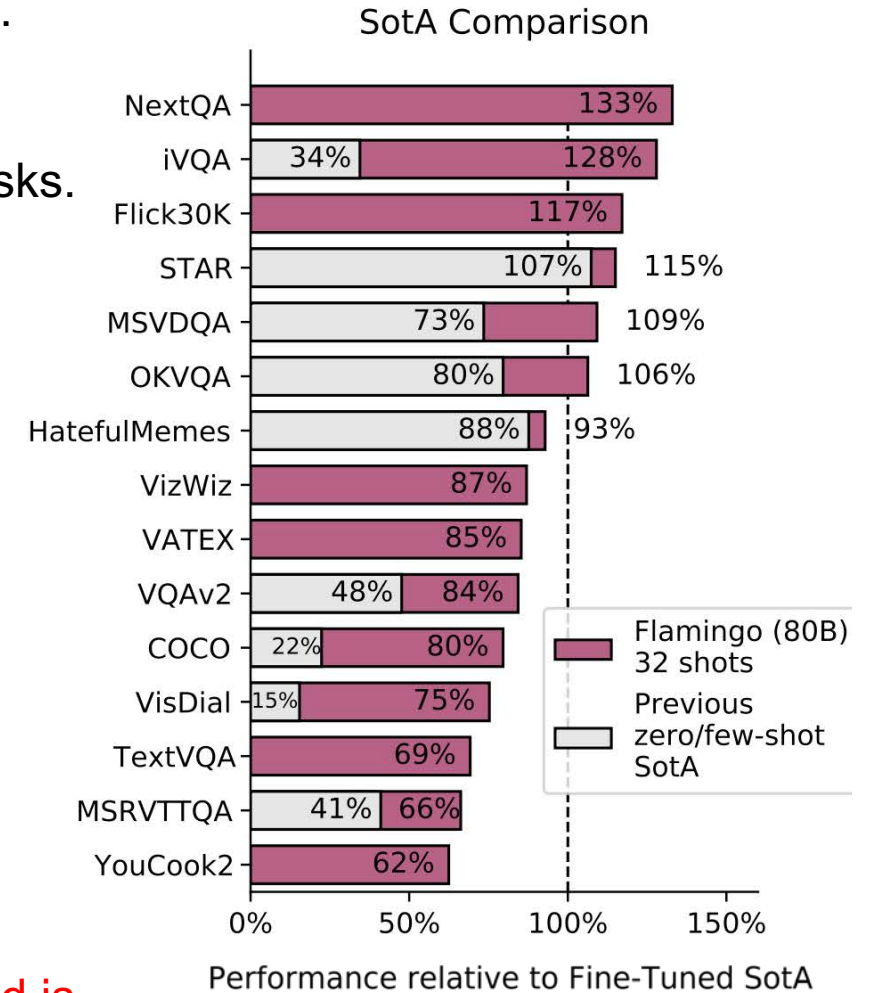
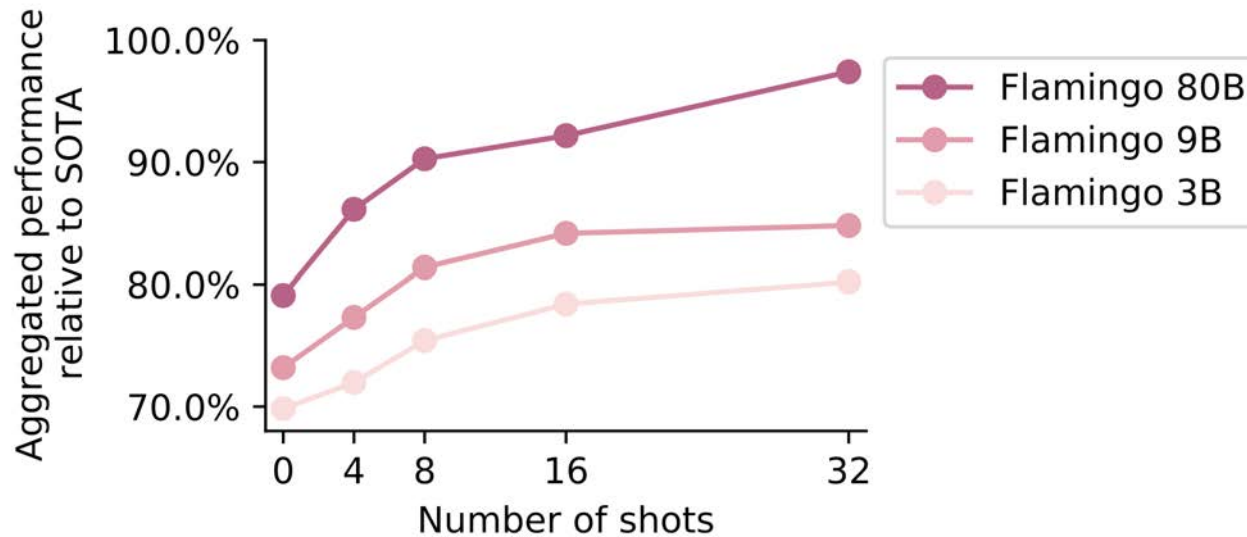
- **Left:** larger model works better; more in-context examples helps.
- **Right:** thanks to larger model and more training data, he model achieves comparable or better results than SOTA on multiple tasks.



Flamingo



- **Left:** larger model works better; more in-context examples helps.
- **Right:** thanks to larger model and more training data, he model achieves comparable or better results than SOTA on multiple tasks.

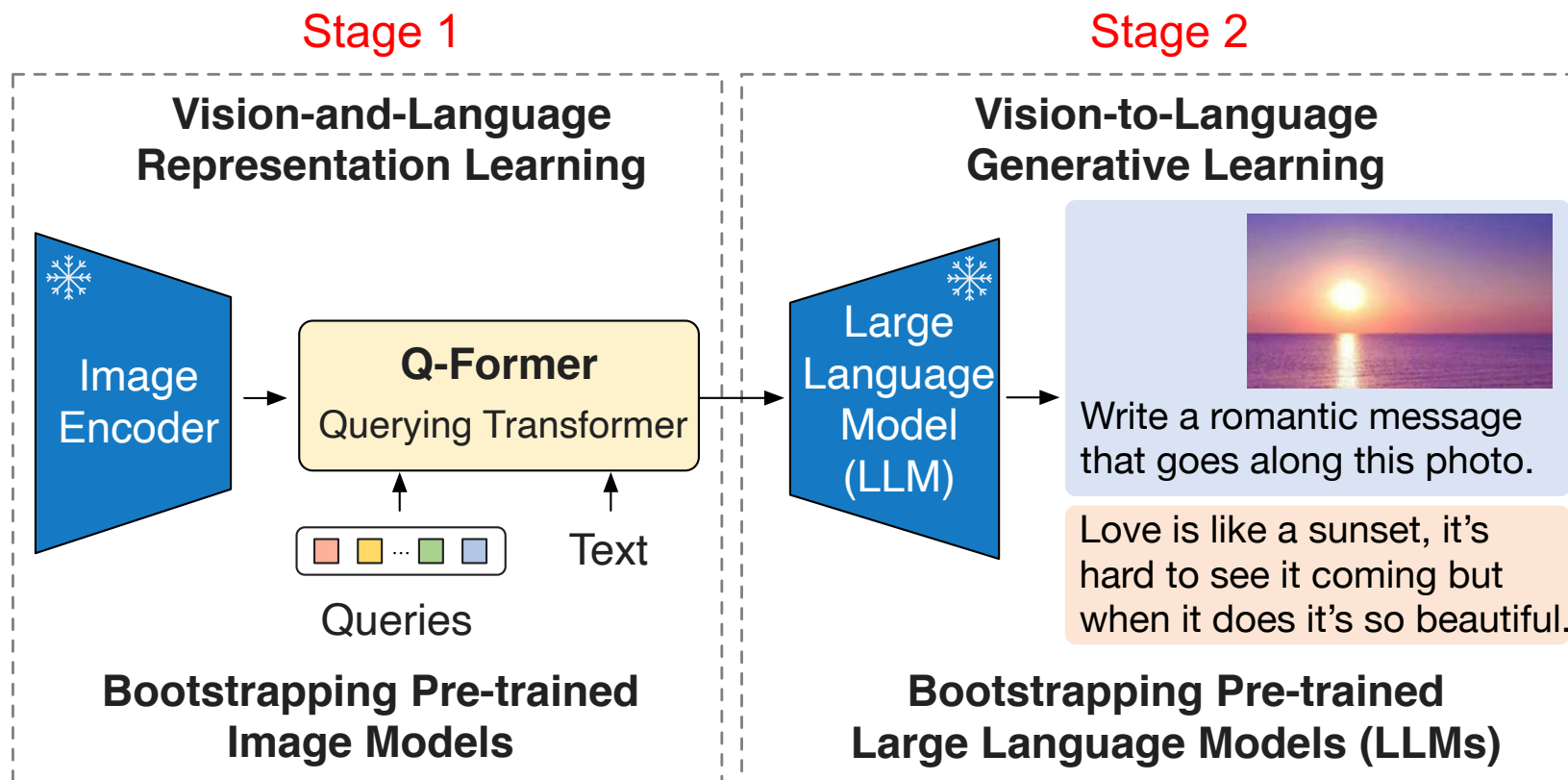


Expensive! The 80 model has 10B trainable parameters and is trained with 1536 TPUv4 chips for 15 days.

BLIP-2, Learning w/ Frozen LLM



- **Architecture:** frozen image encoder + a light-weight Q-Former + frozen LLM
- **Two stage training:**
 - **Stage 1** vision-language representation learning: image-text contrastive & matching, image captioning
 - **Stage 2** vision-language generative learning: generate text conditioned on image
- **Q-Former:** BERT-base, using learned query vectors with cross-attention to extract visual info.



BLIP-2, Learning w/ Frozen LLM



- **Architecture:** frozen image encoder + a light-weight Q-Former + frozen LLM
- **Two stage training:**
 - **Stage 1** vision-language representation learning: image-text contrastive & matching, image captioning
 - **Stage 2** vision-language generative learning: generate text conditioned on image
- **Q-Former:** BERT-base, using learned query vectors with cross-attention to extract visual info.
- Trained on 129M image-text pairs

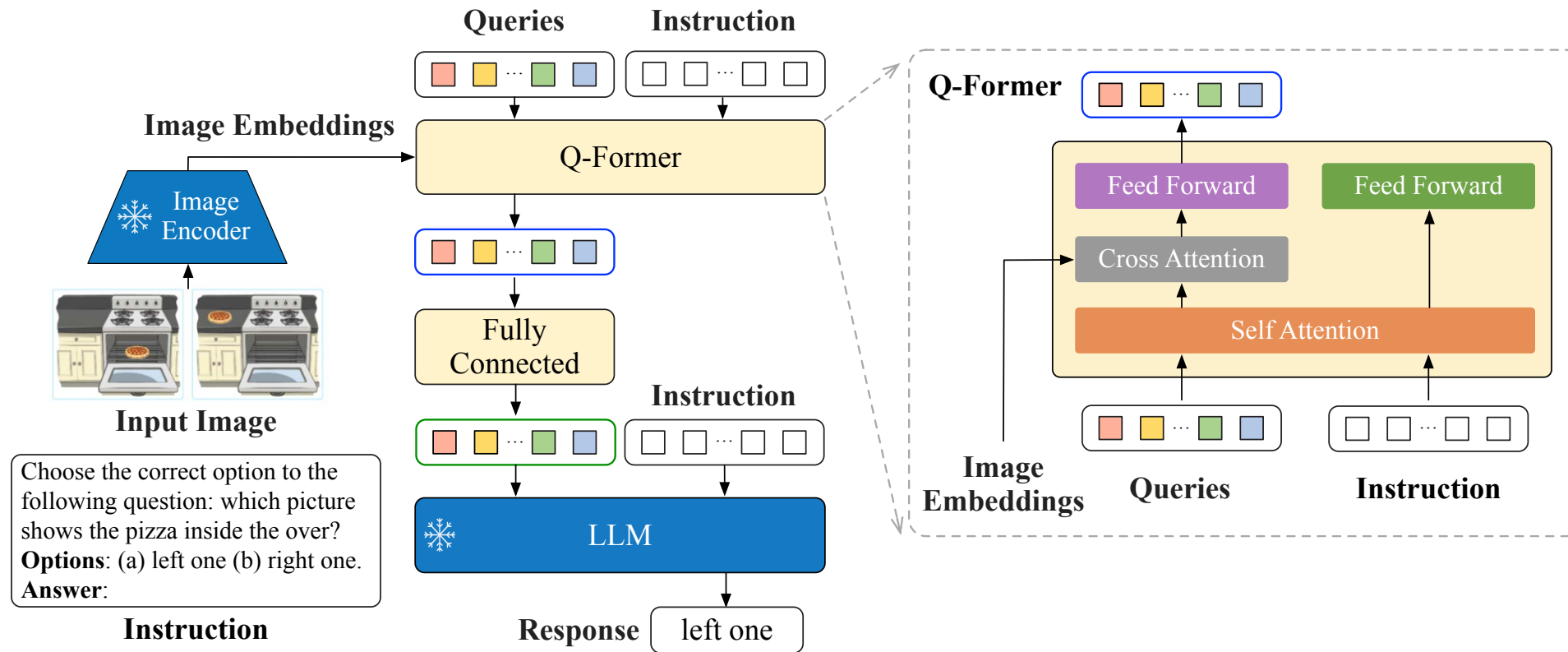
Models	#Trainable Params	Open-sourced?	Visual Question Answering	Image Captioning		Image-Text Retrieval	
			VQAv2 (test-dev) VQA acc.	NoCaps (val) CIDEr	SPICE	Flickr (test) TR@1	IR@1
BLIP (Li et al., 2022)	583M	✓	-	113.2	14.8	96.7	86.7
SimVLM (Wang et al., 2021b)	1.4B	✗	-	112.2	-	-	-
BEIT-3 (Wang et al., 2022b)	1.9B	✗	-	-	-	94.9	81.5
Flamingo (Alayrac et al., 2022)	10.2B	✗	56.3	-	-	-	-
BLIP-2	188M	✓	65.0	121.6	15.8	97.6	89.7

Results on zero-shot vision-language tasks.

InstructBLIP



- **Architecture:** same as BLIP-2, except instruction text is added to Q-Former for instruction-aware visual feature extraction
- **Training:** BLIP-2 pre-training + Instruction Finetuning on 13 held-in datasets
- **Evaluation:** on both held-in and held-out datasets



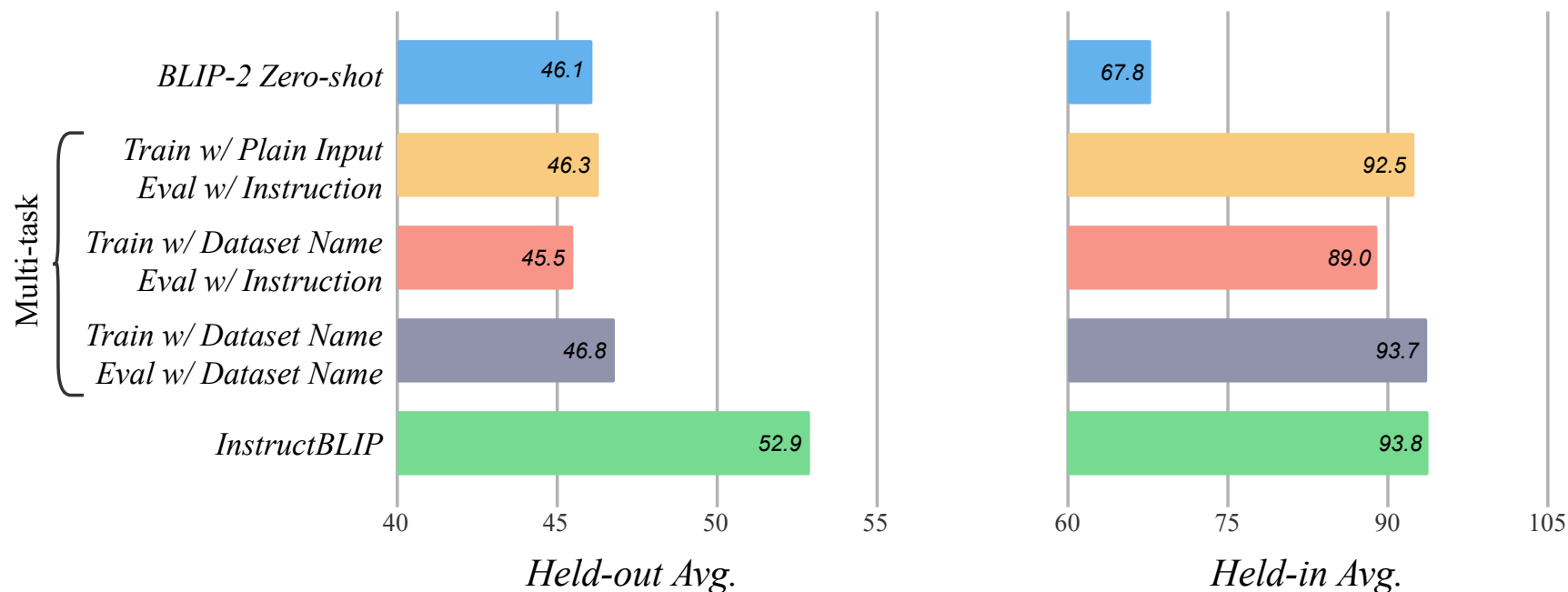
InstructBLIP Model Overview

InstructBLIP: Instruction Tuning vs. Multi-task



- Instruction tuned model excels in unseen datasets and tasks

Strategy	Template (use VQAv2 dataset as an example)
Instruction Tuning	<ul style="list-style-type: none"> • <image> Question: {question} Short answer: • <image> What is the answer to the following question? {question} • <image> Based on the image, respond to this question with a short answer: {Question}. Answer: • ...
Multi-task	<ul style="list-style-type: none"> • Plain text: {image} {question} → {answer} • Dataset Name: {image} [Visual question answering:VQAv2] {question} → {answer}

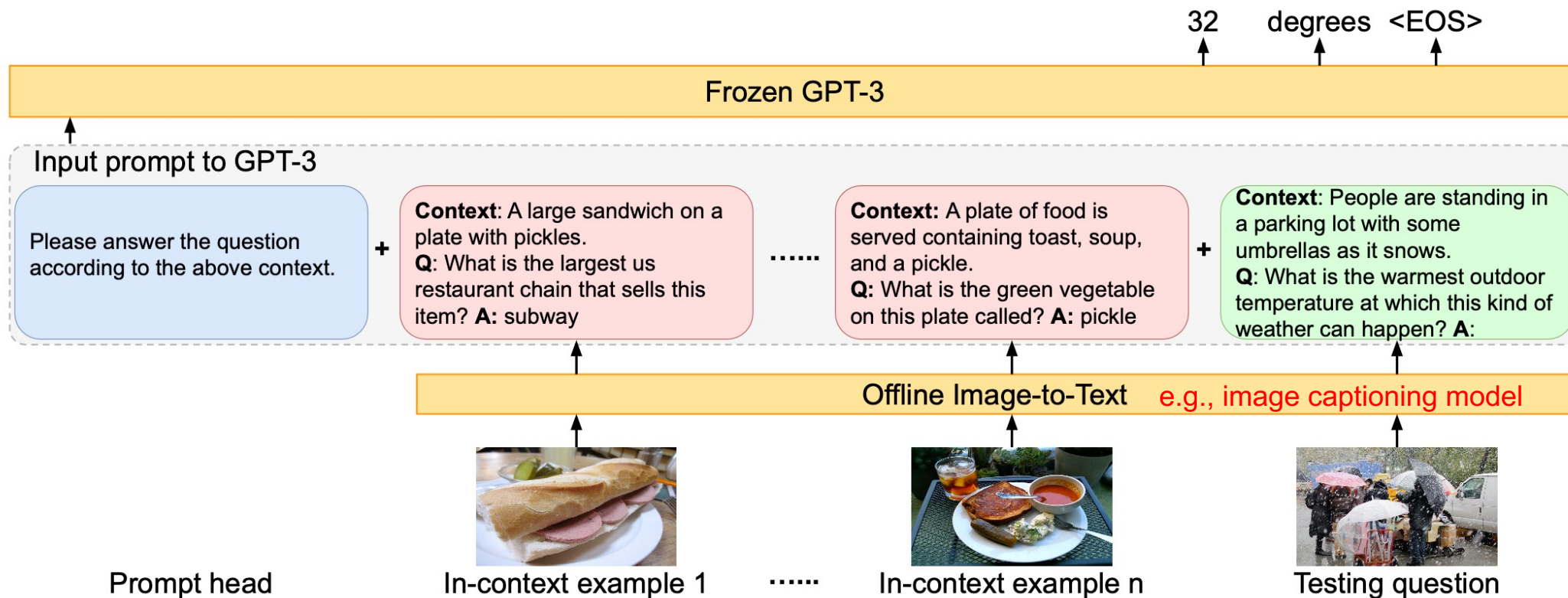


The method discussed above all require additional multi-modal pre-training, however, it is very expensive for LLMs. Is there an alternative way to utilize knowledge in LLMs?

PICa for few-shot knowledge-based VQA



- Summarize image info in text form with an image-to-text model, and prompt GPT-3 to get an answer.
 - Image QA problem is converted into a text QA problem.
 - Implicit GPT-3 knowledge \leftrightarrow previous approaches explicitly query external knowledge
 - Few-shot w/o parameter update.




PICa for few-shot knowledge-based VQA



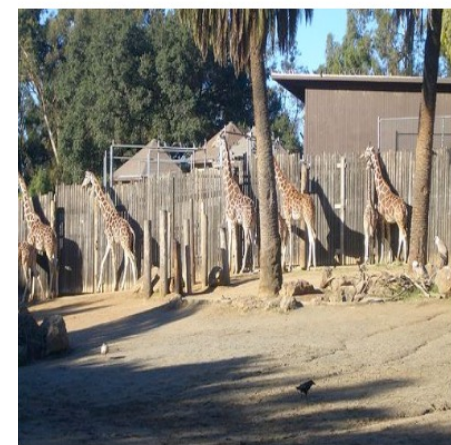
- Works better than fine-tuned models that use explicit wiki knowledge.

Method	Image Repr.	Knowledge Resources	Few-shot	Accuracy
MUTAN+AN (Ben-Younes et al. 2017)	Feature Emb.	Wikipedia	X	27.8
Mucko (Zhu et al. 2020)	Feature Emb.	Dense Captions	X	29.2
ConceptBert (Garderes et al. 2020)	Feature Emb.	ConceptNet	X	33.7
ViLBERT (Lu et al. 2019)	Feature Emb.	None	X	35.2
OKVQA KRISP (Marino et al. 2021)	Feature Emb.	Wikipedia + ConceptNet	X	38.9
MAVEx (Wu et al. 2021)	Feature Emb.	Wikipedia + ConceptNet + Google Images	X	39.4
Frozen (Tsimpoukelli et al. 2021)	Feature Emb.	Language Model (7B)	✓	12.6
PICa-Base	Caption	GPT-3 (175B)	✓	42.0
PICa-Base	Caption+Tags	GPT-3 (175B)	✓	43.3
PICa-Full	Caption	GPT-3 (175B)	✓	46.9
PICa-Full	Caption+Tags	GPT-3 (175B)	✓	48.0

- A core issue: image-to-text models are not perfect, it will cause **information loss**.



(e) What color is the man's jacket?
Context: A man flying through the air while riding a snowboard.
Answer: black
GT Answer: ['red', 'red', 'red', 'orange', 'red', 'red', 'red', 'red', 'red', 'red']
Acc.: 0.0

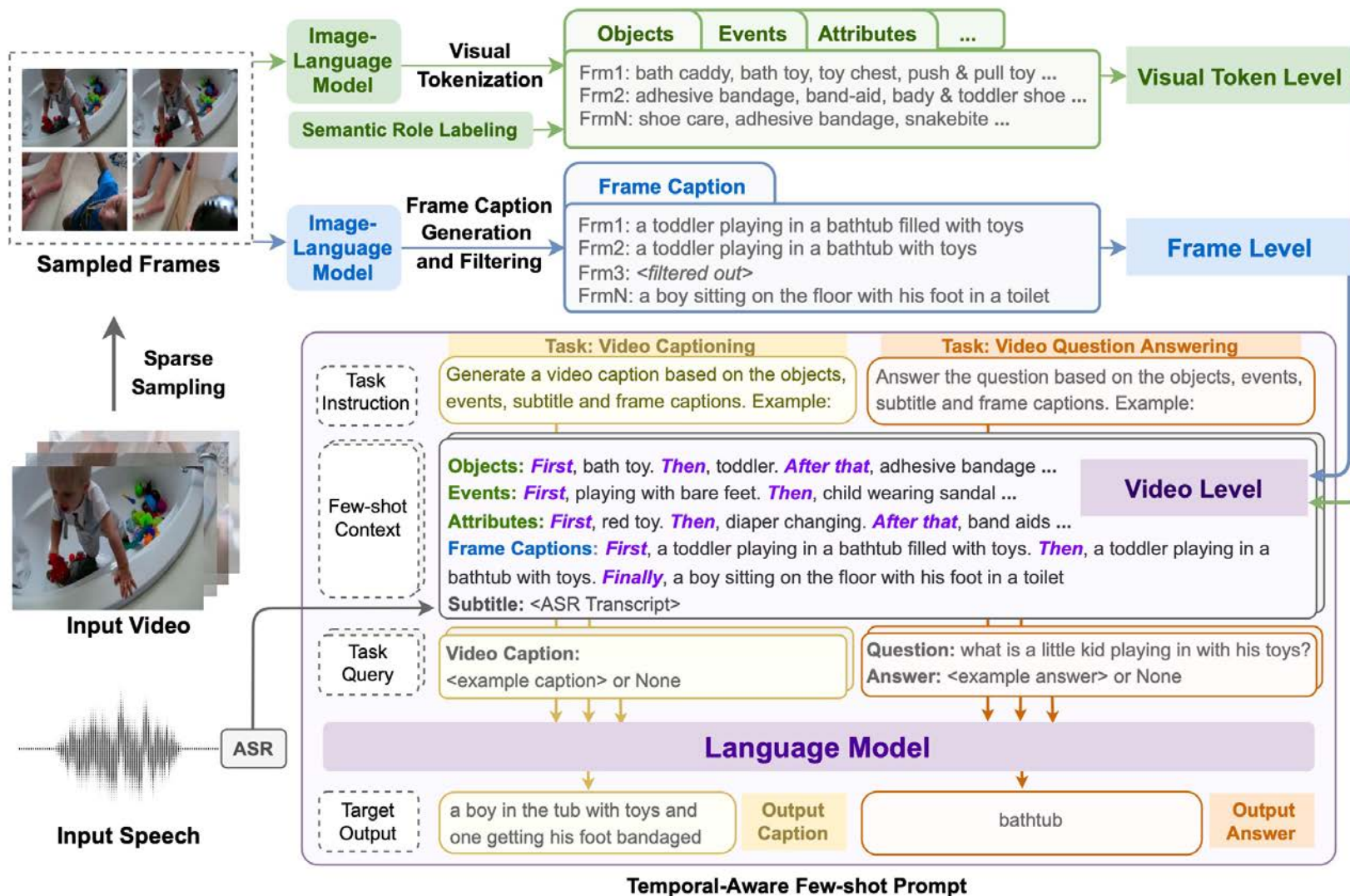


(f) How many giraffes are there?
Context: A herd of giraffe standing next to a wooden fence.
Answer: 3
GT Answer: ['6', '6', '8', '6', '8', '6', '6', '7', '8', '7']
Acc.: 0.0

VidLL: LLM video + language learning



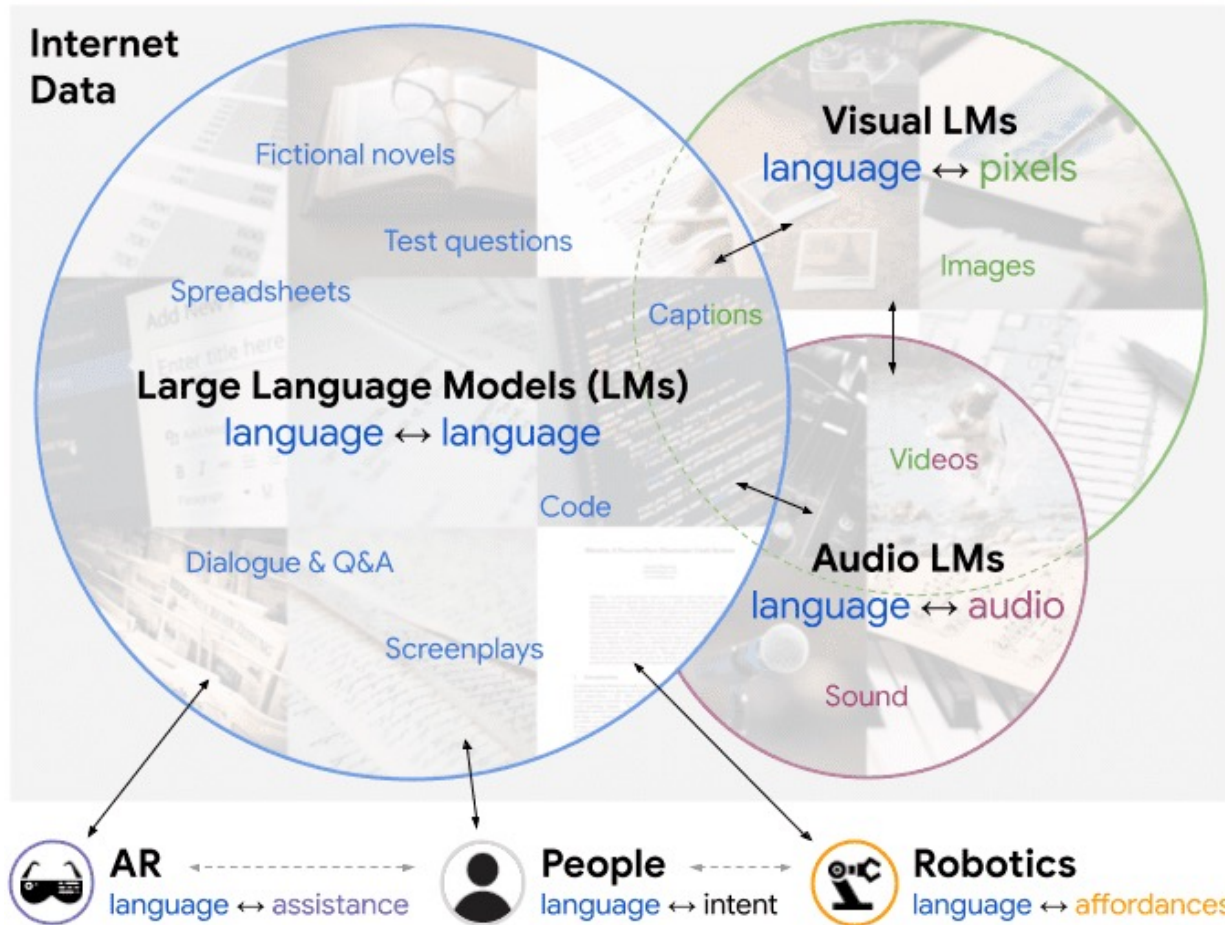
- Generate frame-level info at **various granularity**, and put them in a **temporal aware** prompt for LLM.



Socratic: Composing Multi-modality w/ LLM



- A modular framework in which multiple pretrained models may be composed zero-shot through language without training.



Visual LM

LM + Audio LM

LM + Visual LM

LM

Summarize ego-centric videos.

Socratic: Composing Multi-modality w/ LLM



- The model works well on **vision-language tasks** such as image captioning, it can also **parse & generate robot instructions** from free form human language.

Method	BLEU-4	METEOR	CIDEr	SPICE	ROUGE-L
*ClipCap [45]	40.7	30.4	152.4	25.2	60.9
†MAGIC [61]	11.4	16.4	56.2	11.3	39.0
ZeroCap [62]	0.0	8.8	18.0	5.6	18.3
SMs 0-shot (ours)	6.9	15.0	44.5	10.1	34.1
SMs 3-shot (ours)	18.3	18.8	76.3	14.8	43.7

COCO Captions

* finetuned on full training set with image-text pairs.

† finetuned on unpaired training set, zero-shot on image-text pairs.

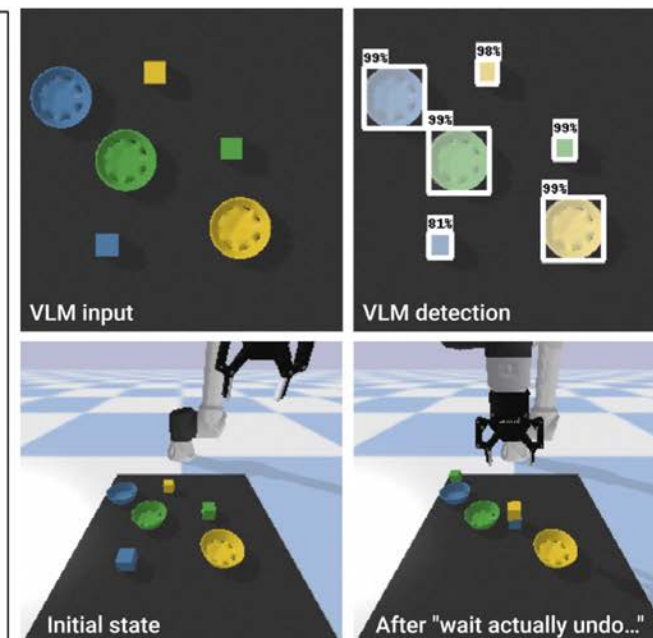
Human defines when to use which model.

Visual LM

User Instruction

LM

```
objects = ["green block", "blue block", "yellow block", "green bowl", "blue bowl", "yellow bowl"]
# move all the blocks to different corners.
Step 1. robot.pick_and_place("green block", "top left corner")
Step 2. robot.pick_and_place("blue block", "top right corner")
Step 3. robot.pick_and_place("yellow block", "bottom left corner")
# now move the blue block to the middle.
Step 1. robot.pick_and_place("blue block", "middle")
# stack the blocks on top of each other.
Step 1. robot.pick_and_place("yellow block", "blue block")
Step 2. robot.pick_and_place("green block", "yellow block")
# wait actually undo that last step.
Step 1. robot.pick_and_place("green block", "top left corner")
# put the yellow block in the bowl you think it best fits.
Step 1. robot.pick_and_place("yellow block", "yellow bowl")
# ok now sort the remaining blocks in the same way.
Step 1. robot.pick_and_place("blue block", "blue bowl")
Step 2. robot.pick_and_place("green block", "green bowl")
```



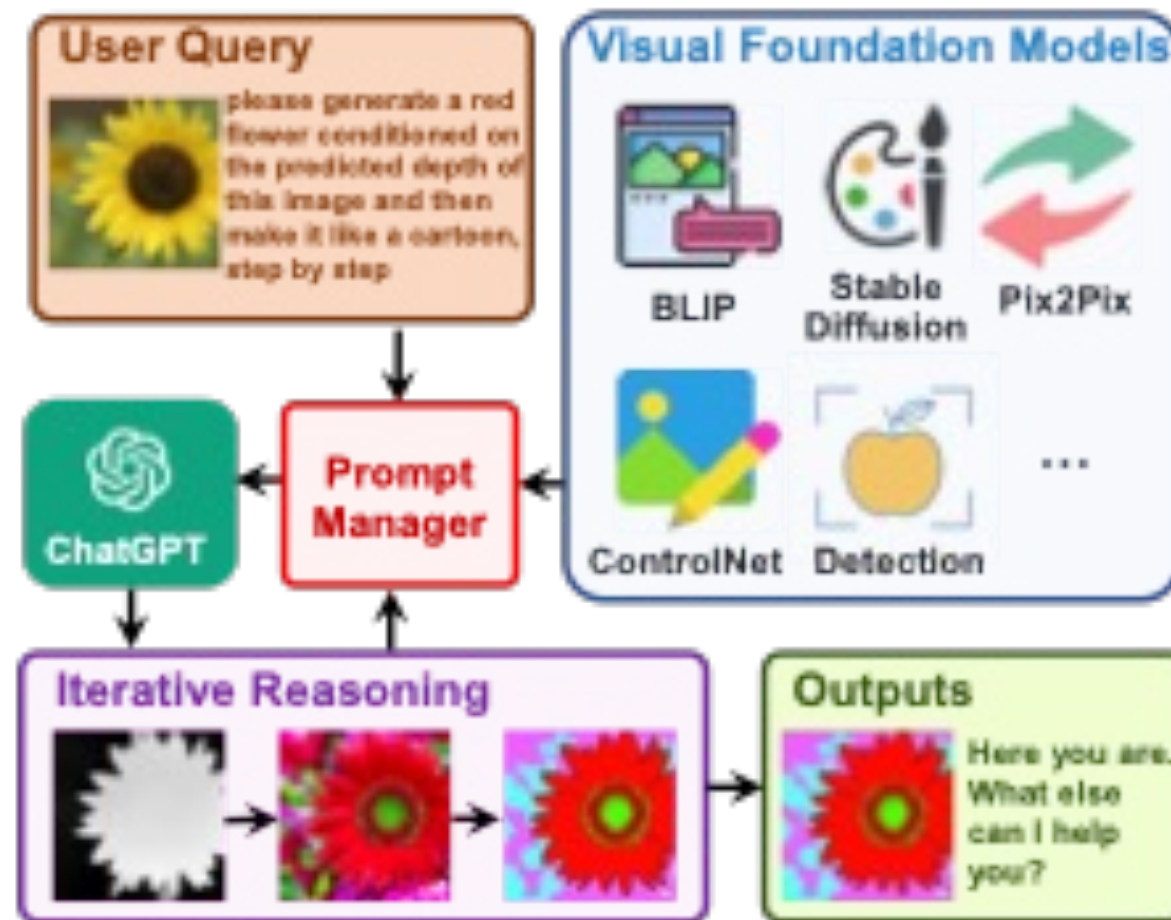
Visual ChatGPT, ViperGPT, ...



- Visual ChatGPT defines a few **system principles**, and give ChatGPT the **autonomy** execute actions:
 - System definition.
 - Define name & usage of vision models.
 - Chain-of-Thought.
 - Be strict about filename.
 - Regex to parse executable actions from language.
 - ...

Visual ChatGPT is very flexible, as the LLM controls when to use which foundation models, instead of human.

- More general **tool learning** framework
 - ViperGPT uses **generated Python code** to compose pre-defined APIs.
 - AutoGPT, New Bing, Bard, ...



Architecture of Visual ChatGPT

Pros

- It provides an **efficient** way to utilize foundation models of different modalities, no extra training required.
- The approaches are **modular**: new modules can be seamlessly plugged into the framework.

Cons

- Modality specific models are not perfect, there will be **info loss** when converted into text.
 - The lower performance vs. e2e trained Flamingo model might partly due to this info loss.



(e) What color is the man's jacket?
Context: A man flying through the air while riding a snowboard.
Answer: black
GT Answer: ['red', 'red', 'red', 'orange', 'red', 'red', 'red', 'red', 'red', 'red']
Acc.: 0.0



(f) How many giraffes are there?
Context: A herd of giraffe standing next to a wooden fence.
Answer: 3
GT Answer: ['6', '6', '8', '6', '8', '6', '6', '7', '8', '7']
Acc.: 0.0

Failure cases from the PICa model.

The use of implicit knowledge from pre-trained LMs shows strong zero-shot performance for multi-modal tasks, however, they are hard to interpret. Is there a more interpretable way of using language knowledge?

Part 1.2 Explicit Knowledge from Language

- External knowledge is useful to help the model understand **rare** concepts.



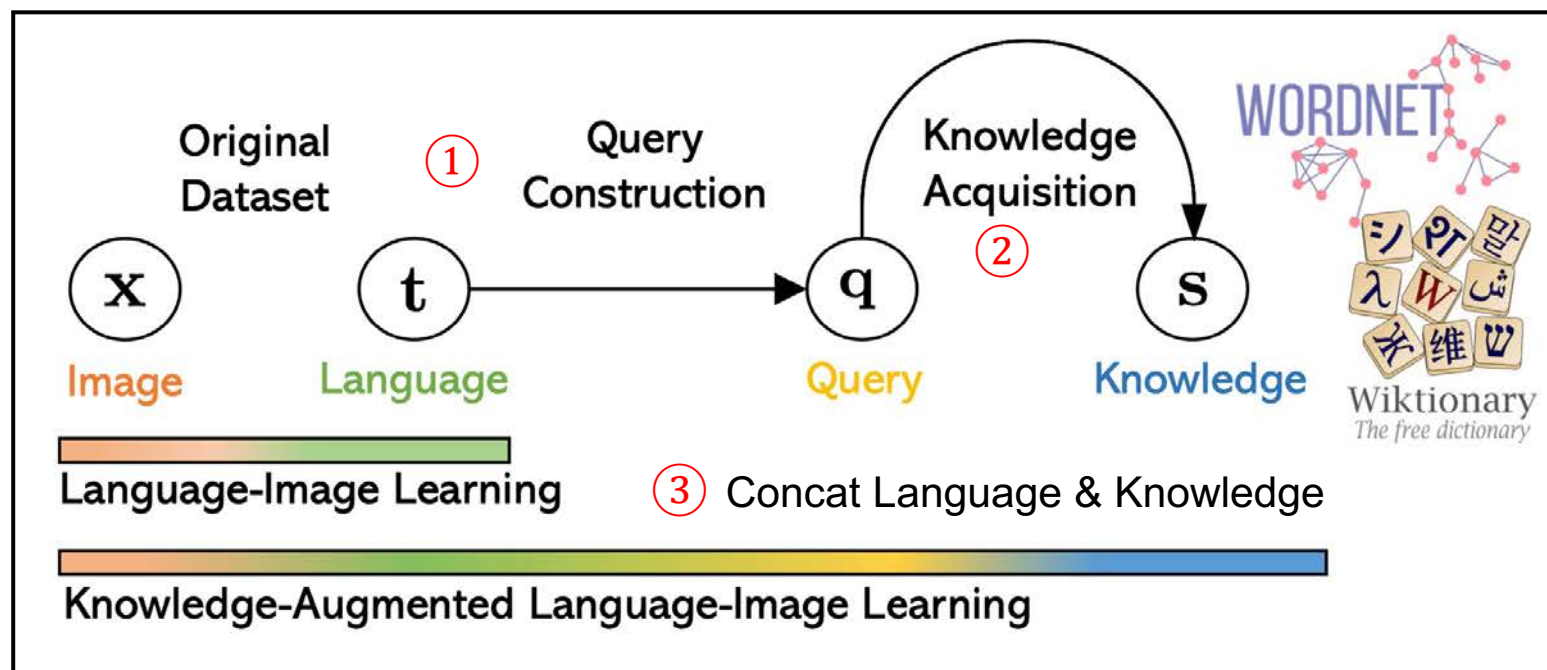
Sashimi

A dish consisting of **thin slices** or pieces of **raw fish or meat**.



Takoyaki

A **ball-shaped** Japanese **dumpling** made of batter, filled with diced octopus, **tempura scraps**, pickled ginger, and **green onion**.



- a photo of sashimi,
- a photo of takoyaki,
- ...

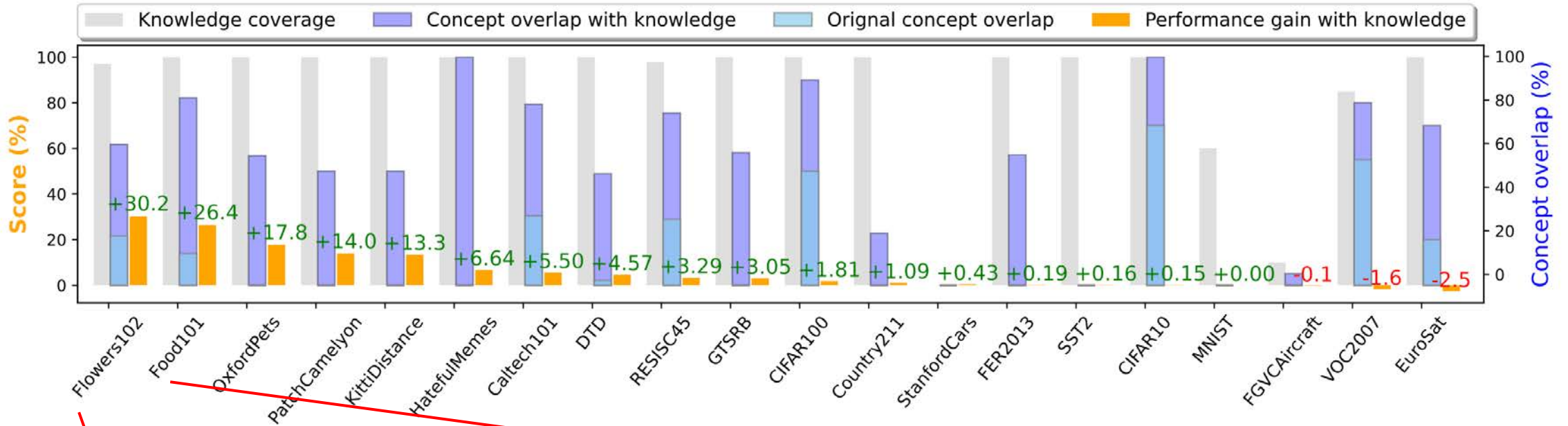


Add knowledge



- a photo of sashimi. A dish consisting of slices or
- a photo of takoyaki, a ball-shaped Japanese dumpling...
- ...

- Orange: knowledge improves zero-shot performance on 16/20 image classification datasets.



✓ **English marigold:** Any of the Old World plants, of the genus *Calendula*, with orange, yellow or reddish flowers.

✗ **Wallflower:** Any of several short-lived herbs or shrubs of the *Erysimum* genus with bright yellow to red flowers.



✓ **Lobster bisque:** A thick creamy soup made from fish, shellfish, meat or vegetables.

✗ **Hot and sour soup:** Any one of several soups, served in various Asian cuisines, which are both spicy and sour

ELEVATER



- Same K-LITE model, but with **GPT-3 knowledge**
- GPT-3 knowledge improves ZS image classification and object detection. More is better.
- GPT-3 + wiki is often better for image classification, but not for object detection.

□ **Concept name:** snowberg

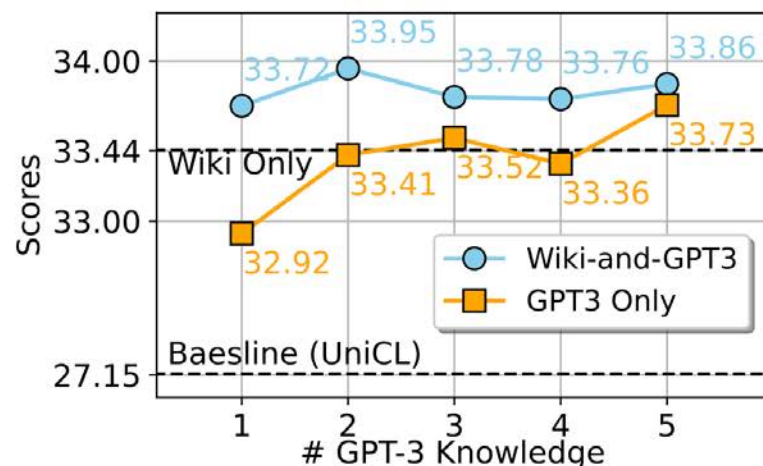
Def_wik: None

GPT3 Query:
Please explain the concept according to the context.
===

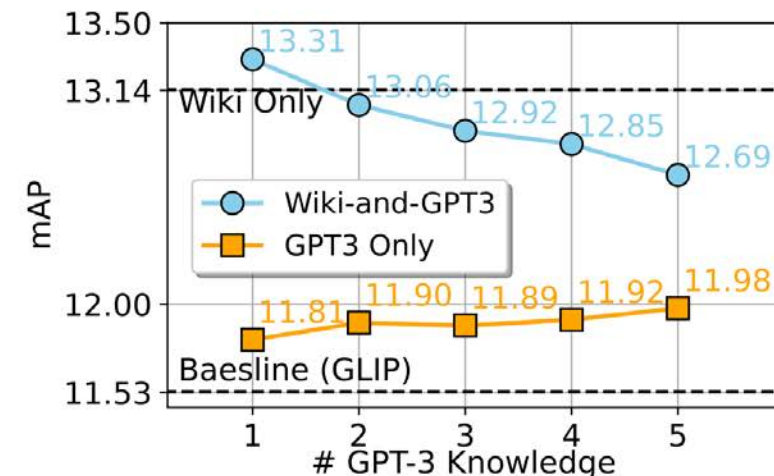
Q: ship
A: A water-borne vessel generally larger than a boat.
===

Q: storage tank
A: A closed container for liquids or gases.
===

Q: snowberg
A:
 GPT3 Answer: A large mass of ice floating in the sea.



(a) Image classification



(b) Object detection

Zero-shot performance

Could vision knowledge help learn language?

Could vision knowledge help learn language?



- Visual pointing is an essential step for most children to learn meanings of words [Bloom 2002].



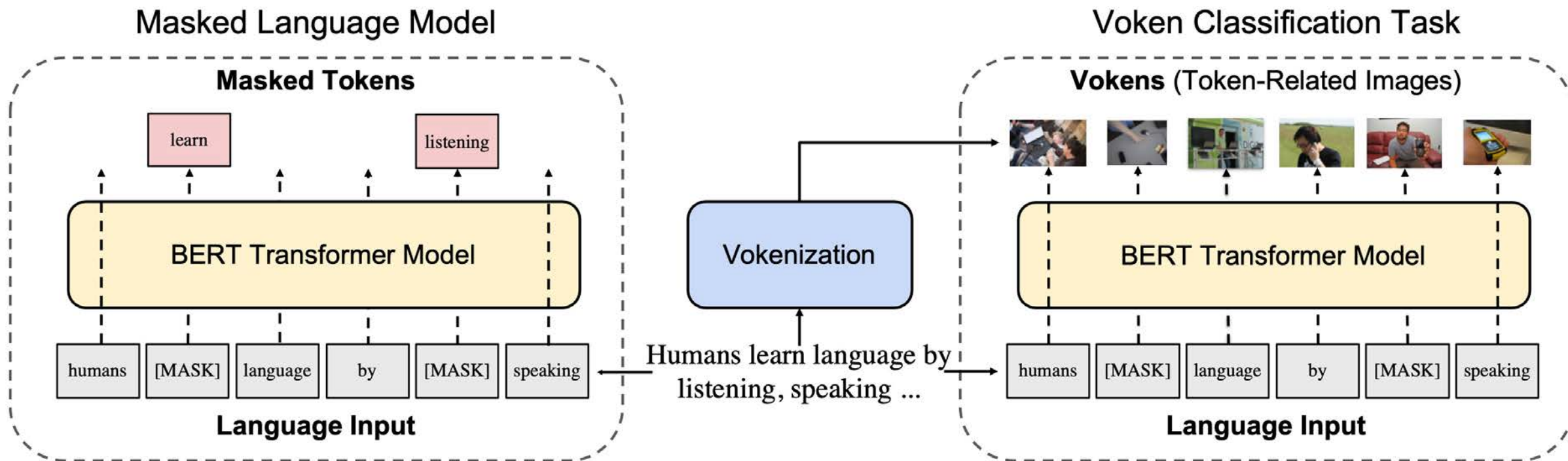
Look! This is a “cat”!



Vokenization: LM w/ Vision Supervision



- Besides standard Masked Language Modeling (MLM), the LM is also trained w/ a **voken classification** task, by assigning each text token into one of the **images (vokens)** in the pool.
- Vokens are **pre-defined**, and are obtained by using a pre-trained **image-text retrieval** model



Vokenization: LM w/ Vision Supervision



- Voken classification task **improves LM performance** on a wide range of **pure-language tasks**.
- This conclusion holds for both BERT and RoBERTa.

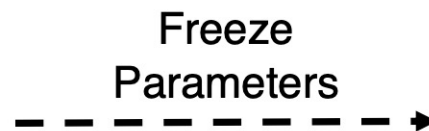
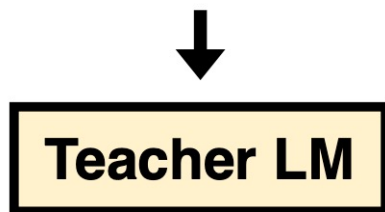
Method	SST-2	QNLI	QQP	MNLI	SQuAD v1.1	SQuAD v2.0	SWAG	Avg.
BERT _{6L/512H}	88.0	85.2	87.1	77.9	71.3/80.2	57.2/60.8	56.2	75.6
BERT _{6L/512H} + Voken-cl	89.7	85.0	87.3	78.6	71.5/80.2	61.3/64.6	58.2	76.8
BERT _{12L/768H}	89.3	87.9	83.2	79.4	77.0/85.3	67.7/71.1	65.7	79.4
BERT _{12L/768H} + Voken-cl	92.2	88.6	88.6	82.6	78.8/86.7	68.1/71.2	70.6	82.1
RoBERTa _{6L/512H}	87.8	82.4	85.2	73.1	50.9/61.9	49.6/52.7	55.1	70.2
RoBERTa _{6L/512H} + Voken-cl	87.8	85.1	85.3	76.5	55.0/66.4	50.9/54.1	60.0	72.6
RoBERTa _{12L/768H}	89.2	87.5	86.2	79.0	70.2/79.9	59.2/63.1	65.2	77.6
RoBERTa _{12L/768H} + Voken-cl	90.5	89.2	87.8	81.0	73.0/82.5	65.9/69.3	70.4	80.6

VidLanKD: LM w/ Video-Distilled Knowledge

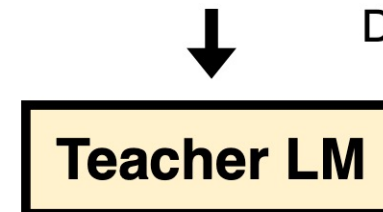
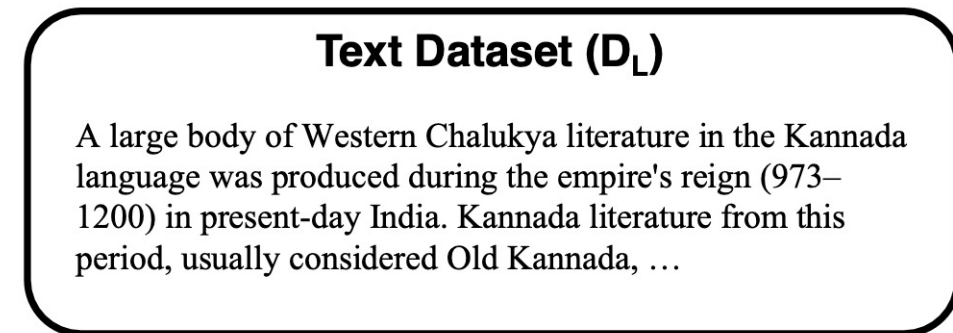


- Vokenization suffers from **approximation error** of using **finite image labels** + the **lack of vocabulary diversity** of a small image-text dataset (COCO).
- VidLanKD improves it by (1) using **knowledge distillation** instead of discrete vokenization to avoid approximation error; (2) using a **large-scale video-language dataset** HowTo100M.

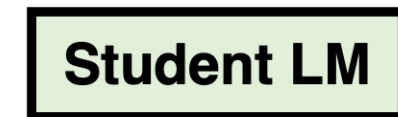
(a) Cross-modal Pretraining



(b) Knowledge Distillation



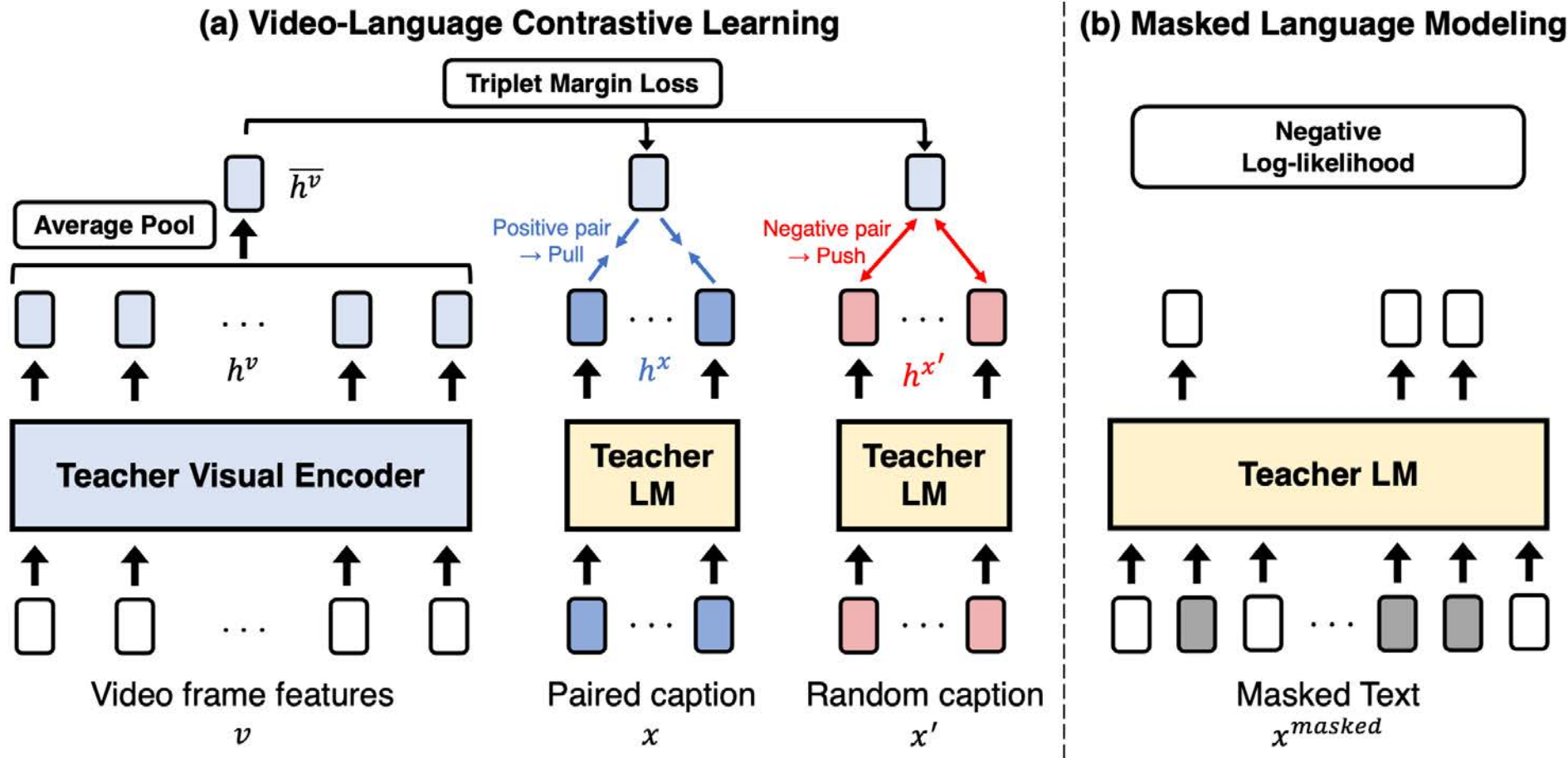
Distillation



VidLanKD: LM w/ Video-Distilled Knowledge



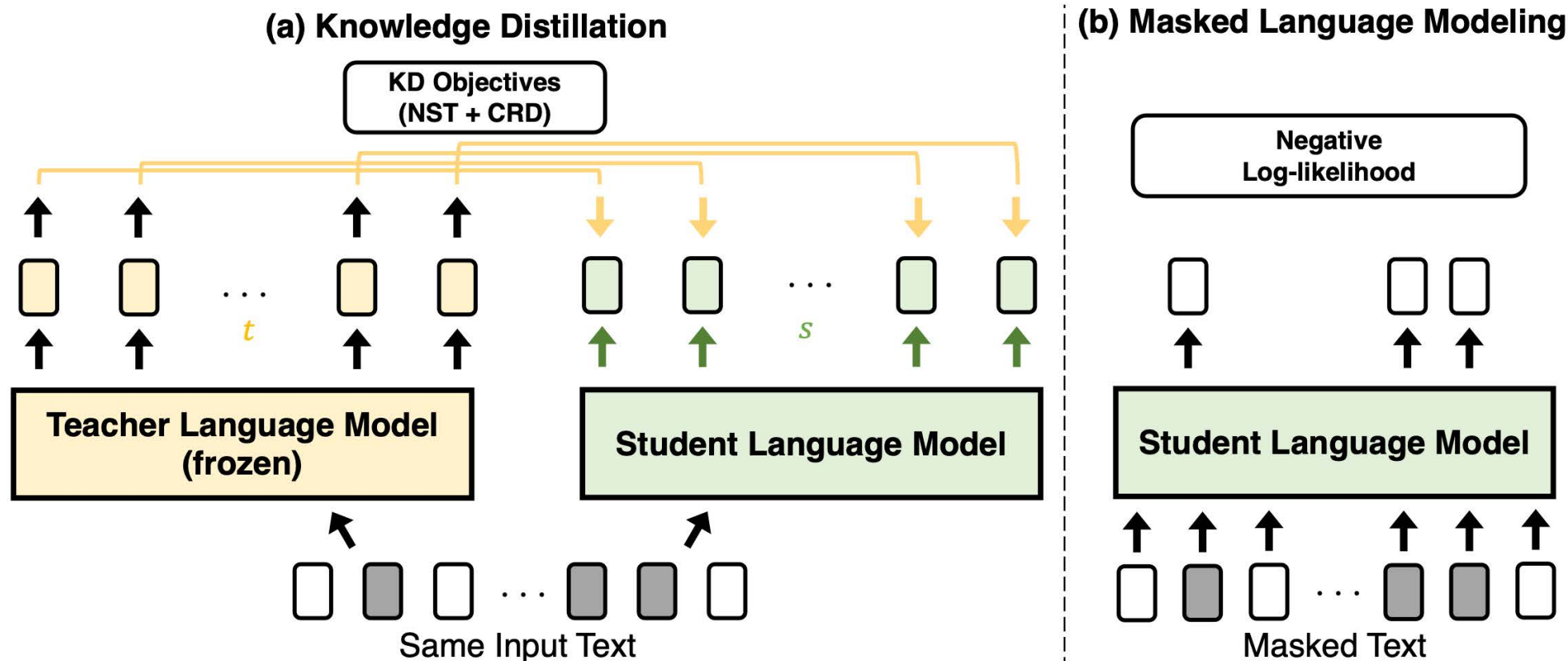
- The teacher LM is trained with (a) video-language triplet loss; + (b) masked language modeling



VidLanKD: LM w/ Video-Distilled Knowledge



- The student LM is trained with (a) knowledge distillation; + (b) masked language modeling



VidLanKD: LM w/ Video-Distilled Knowledge



- Cross-modal KD (last 2 rows) achieves better performance than image vokenization.

	SST-2 Acc	QNLI Acc	QQP Acc	MNLI Acc	SQuAD v1.1 EM [†]	SQuAD v2.0 EM	SWAG Acc	Avg.
BERT _{12L/768H} [68]	89.3	87.9	83.2	79.4	77.0	67.7	65.7	78.6
+ KD (Img-Voken) [68]	92.2	88.6	88.6	82.6	78.8	68.1	70.6	81.4
BERT _{12L/768H}	89.0	88.0	86.2	79.2	77.2	68.0	65.0	78.9
+ KD (Vid-Voken) w/ ResNet	93.4	89.2	88.7	83.0	78.9	68.7	70.0	81.7
+ KD (Vid-Voken) w/ CLIP	94.1	89.8	89.0	83.9	79.2	68.6	71.6	82.3
+ KD (NST+CRD) w/ ResNet	94.2	89.3	89.7	84.0	79.0	68.9	71.8	82.4
+ KD (NST+CRD) w/ CLIP	94.5	89.6	89.8	84.2	79.6	68.7	72.0	82.6

- Performance gain is mostly from **knowledge, physical interaction, & temporal reasoning**

	GLUE diagnostics				PIQA	TRACIE
	Lexicon	Predicate	Logic	Knowledge		
BERT _{6L/512H}	53.0	64.2	44.5	44.0	56.9	63.4
+ KD-NST	53.3 (+0.3)	63.7 (-0.5)	44.8 (+0.3)	48.6 (+4.6)	60.0 (+3.1)	66.7 (+3.3)

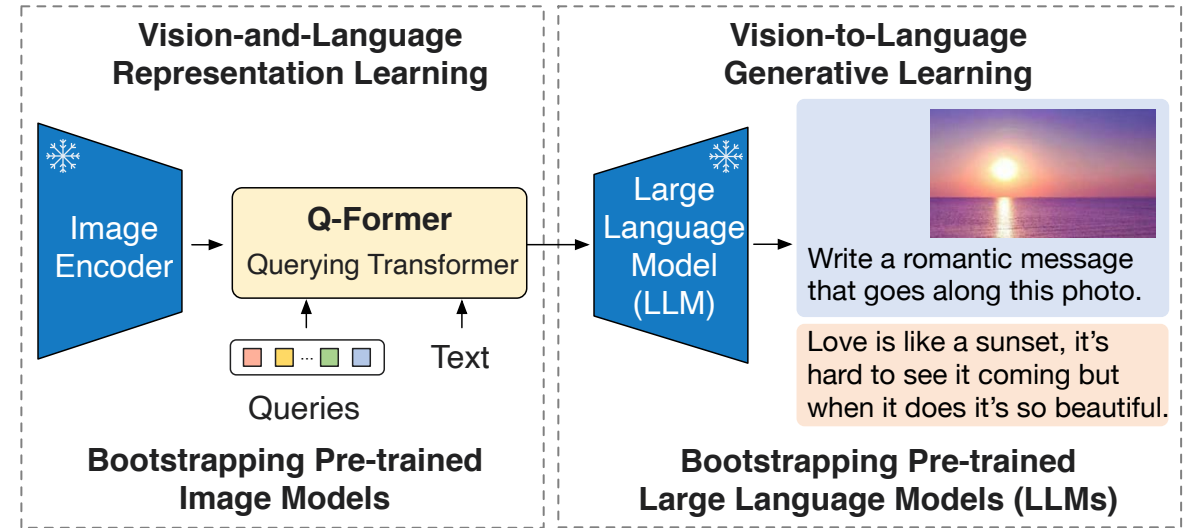
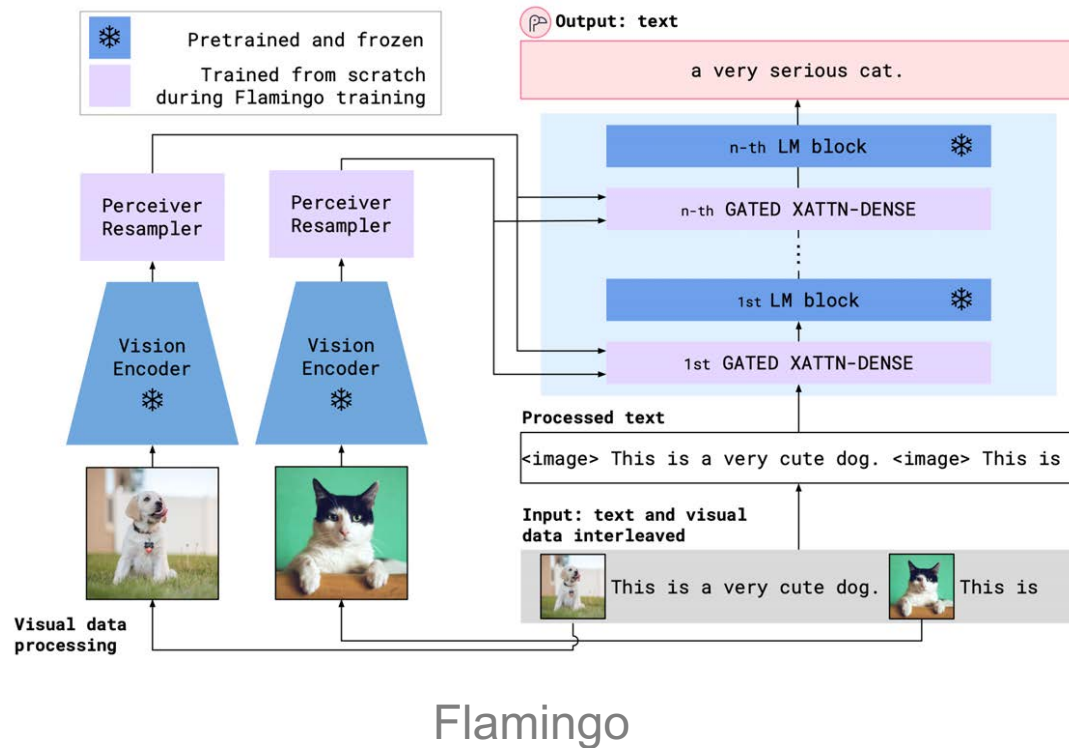
PIQA: QA w/ physical interactions + commonsense reasoning

TRACIE: a temporal reasoning benchmark

Future Work



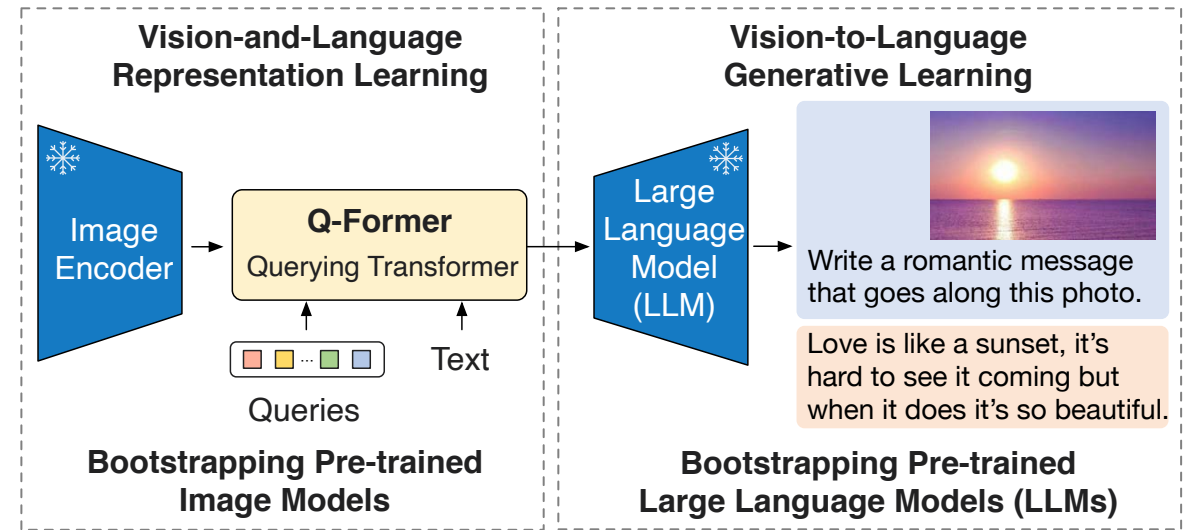
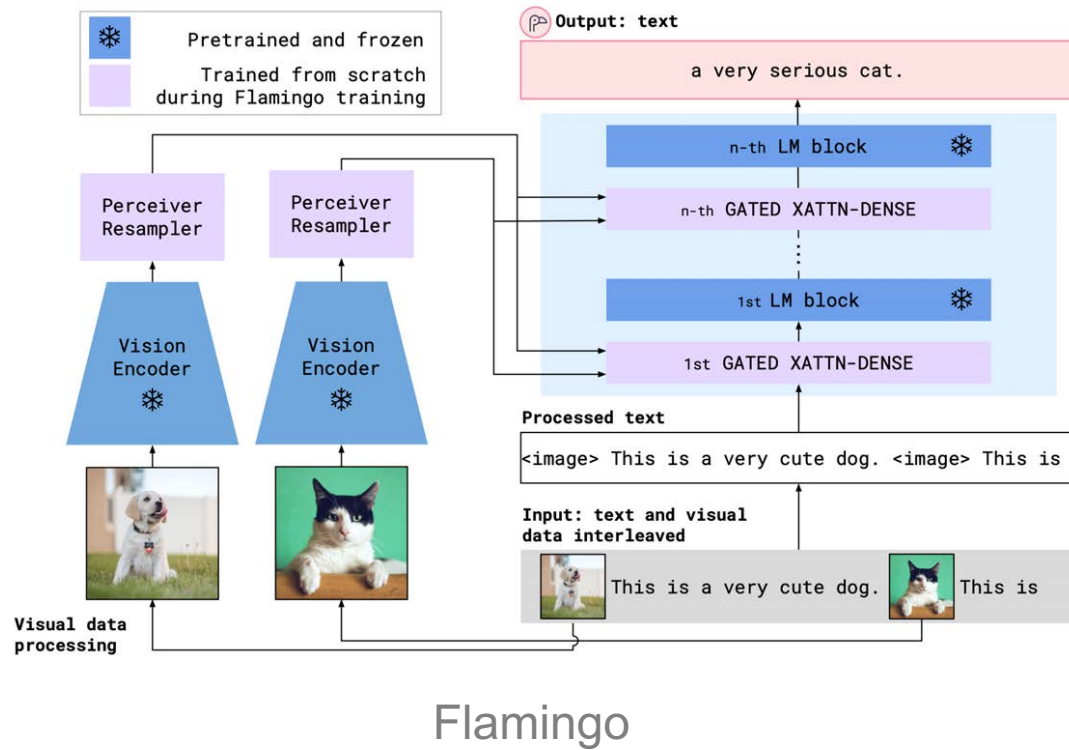
- How to better bridge LLM and other modalities?
 - Is frozen LLM the best approach?
 -



Future Work



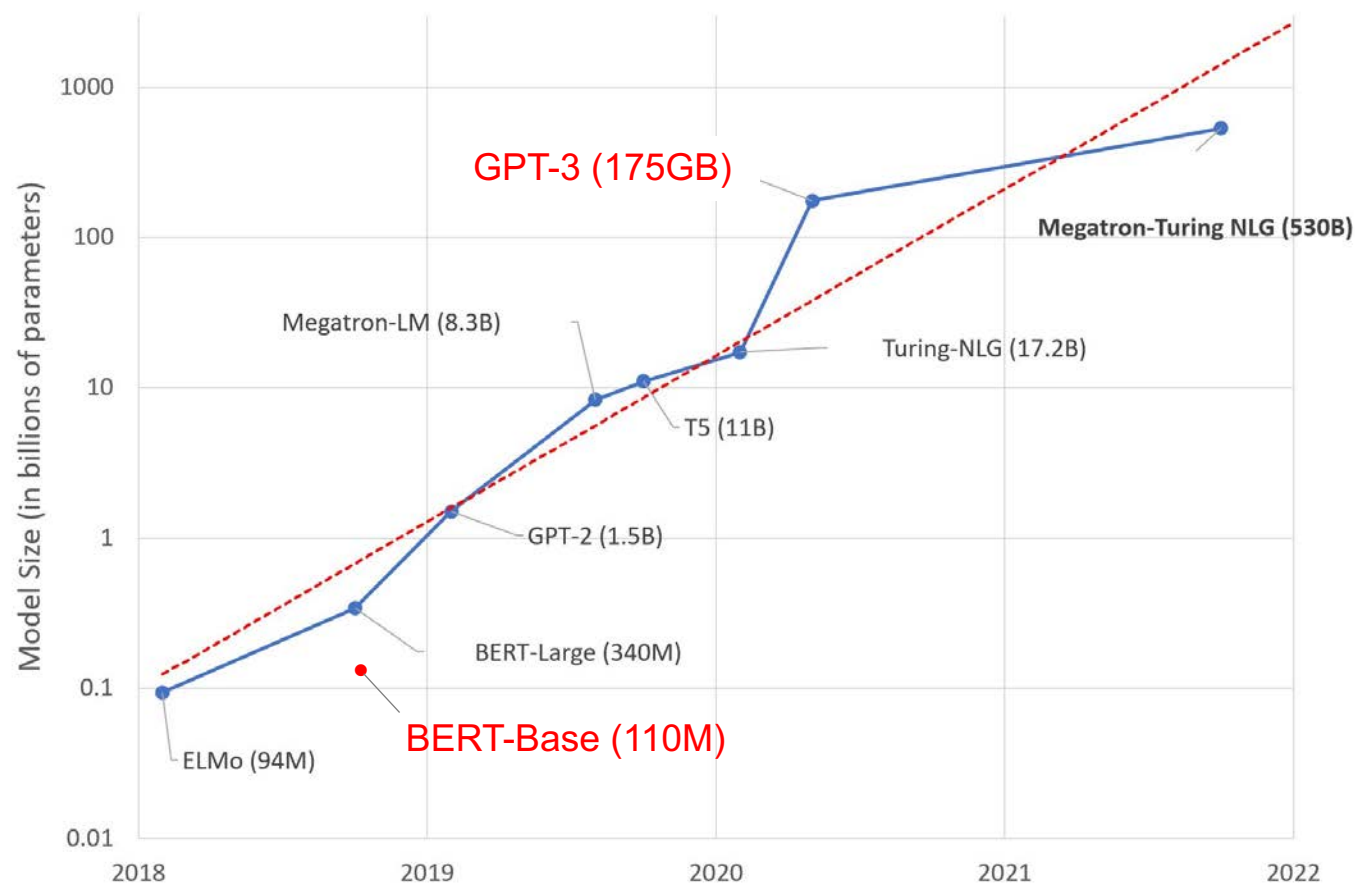
- How to better bridge LLM and other modalities?
 - Is frozen LLM the best approach?
 - If more than one modalities are needed, how to better model them together?



Future Work

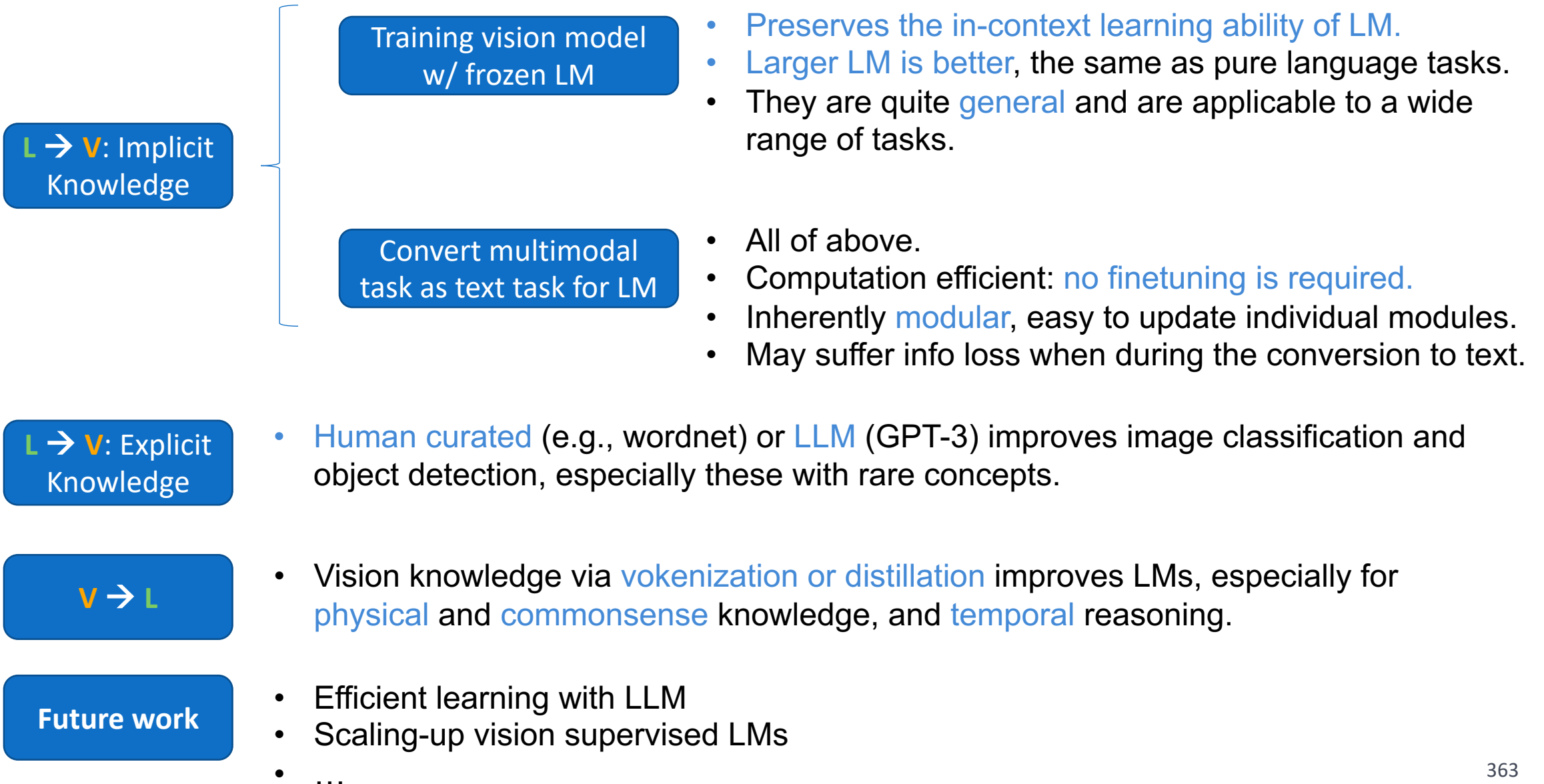


- Using vision (image or video) supervision has shown some early success.



- Bidirectional LM only, casual LM is not explored.
- Small model (up to 110M BERT-base), vs., 175B GPT-3
- How about using other modalities (audio) as supervision?

Take-way Messages



-



Thanks!



Knowledge-Driven Vision-Language Encoding

Panel 1: Explicit Knowledge vs Implicit Knowledge

What is the appropriate format of knowledge representation?

Does explicit knowledge still have value in the era of large models?



Mohit Bansal
UNC



Carl Vondrick
Columbia



Xudong Lin
Columbia



Jie Lei
Meta AI



Manling Li
UIUC



Knowledge-Driven Vision-Language Encoding

Panel 2: LLMs for Multimodality

What can we borrow from Large Language Models (LLMs)?



Mohit Bansal
UNC



Carl Vondrick
Columbia



Xudong Lin
Columbia



Jie Lei
Meta AI



Manling Li
UIUC



Knowledge-Driven Vision-Language Encoding

Panel 3: Image vs Video vs Audio vs Embodied AI

What is the bottleneck for each single modality?

What is the bottleneck to bring multiple modalities together?



Mohit Bansal
UNC



Carl Vondrick
Columbia



Xudong Lin
Columbia



Jie Lei
Meta AI



Manling Li
UIUC



Knowledge-Driven Vision-Language Encoding

Panel 4: Open Challenges

What is the recommended thesis topic for next few years?



Mohit Bansal
UNC



Carl Vondrick
Columbia



Xudong Lin
Columbia



Jie Lei
Meta AI



Manling Li
UIUC