



Knowledge-Driven Vision-Language Pretraining



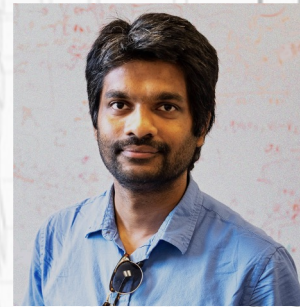
Manling Li
UIUC



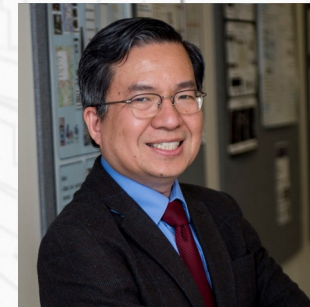
Xudong Lin
Columbia



Jie Lei
Meta AI



Mohit Bansal
UNC



Shih-Fu Chang
Columbia

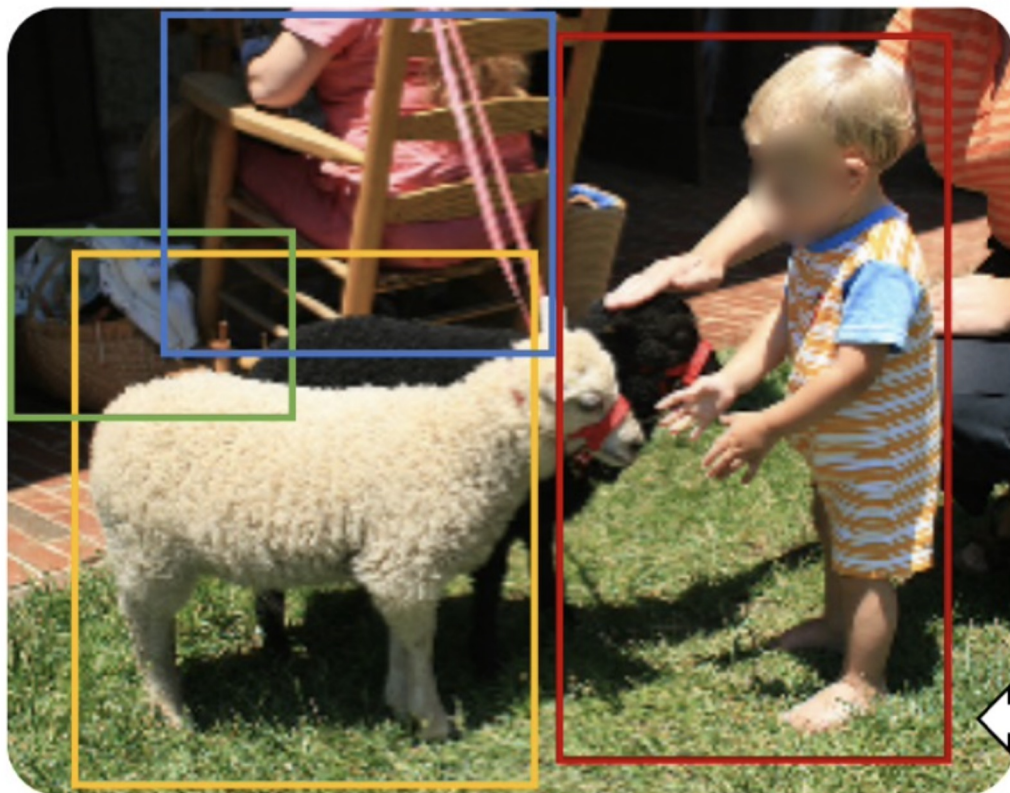


Heng Ji
UIUC

Why Vision-Language?



- Multimedia downstream tasks



Visual Question Answering

What color is the child's outfit? Orange

Referring Expressions

child sheep basket people sitting on chair

Multi-modal Verification

The child is petting a dog. **false**

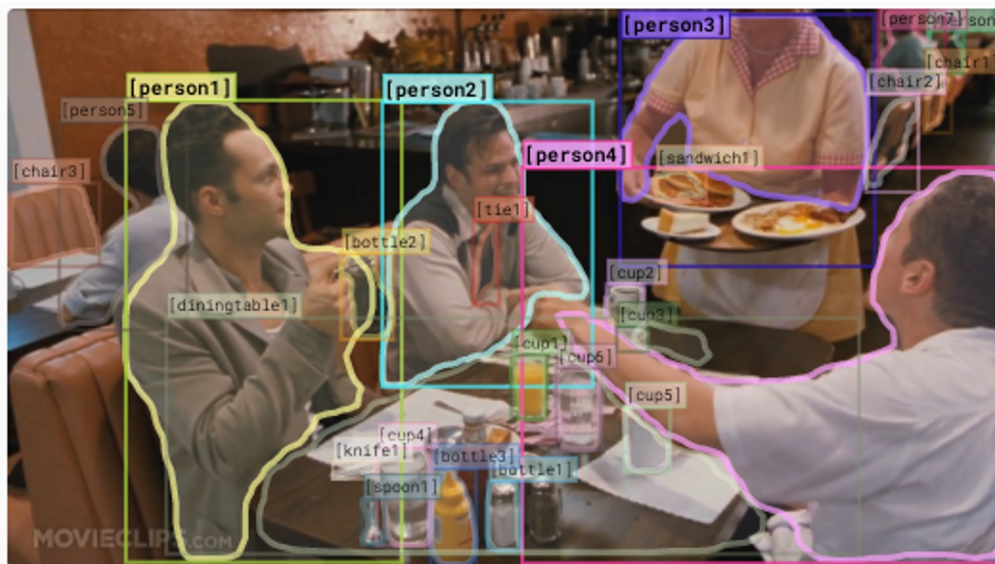
Caption-based Image Retrieval and *image captioning*

A child in orange clothes plays with sheep.

Why Vision-Language?



- VCR: Visual Commonsense Reasoning



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

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

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

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

[chair2] [chair3] [diningtable1]

Why is [person4] pointing at [person1]?

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

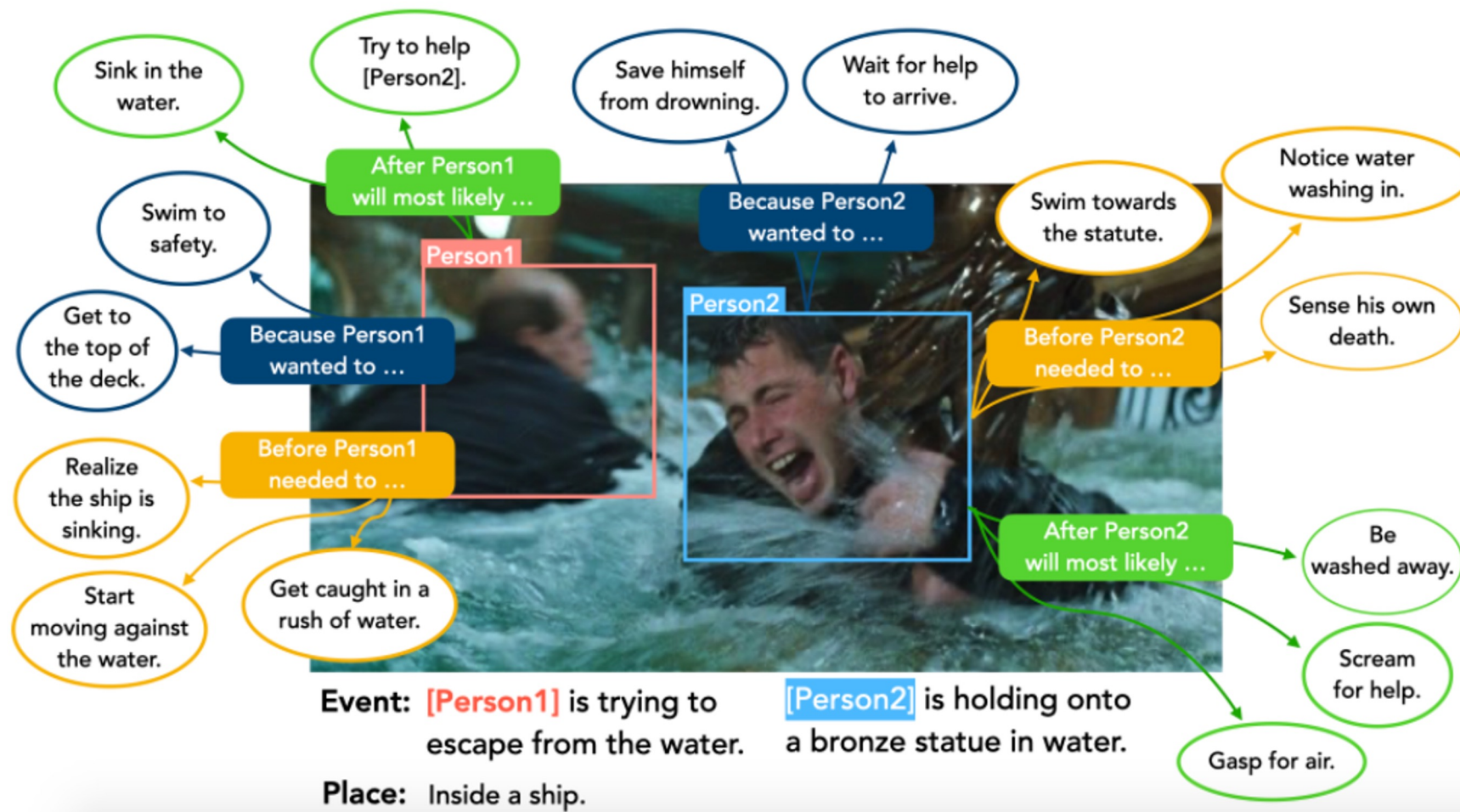
Rationale: I think so because...

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

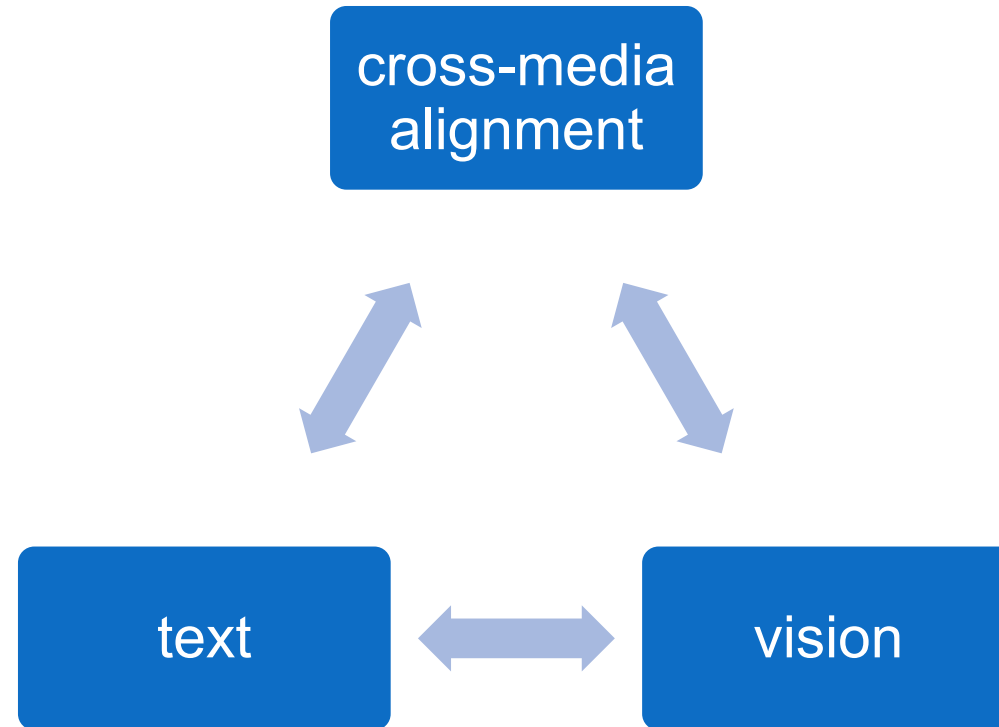
Why Vision-Language?



- VisualCOMET: Visual Commonsense Reasoning in Time

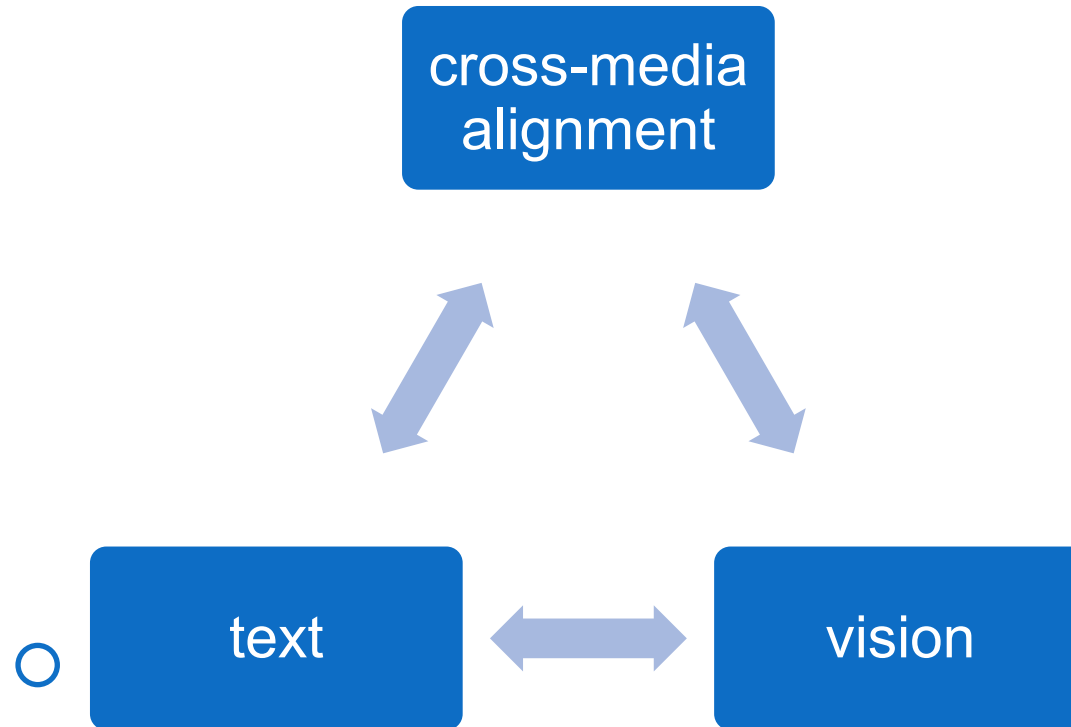


- How? Vision-language pre-training by vision-language pairs





- How? Image-text pre-training by image-text pairs

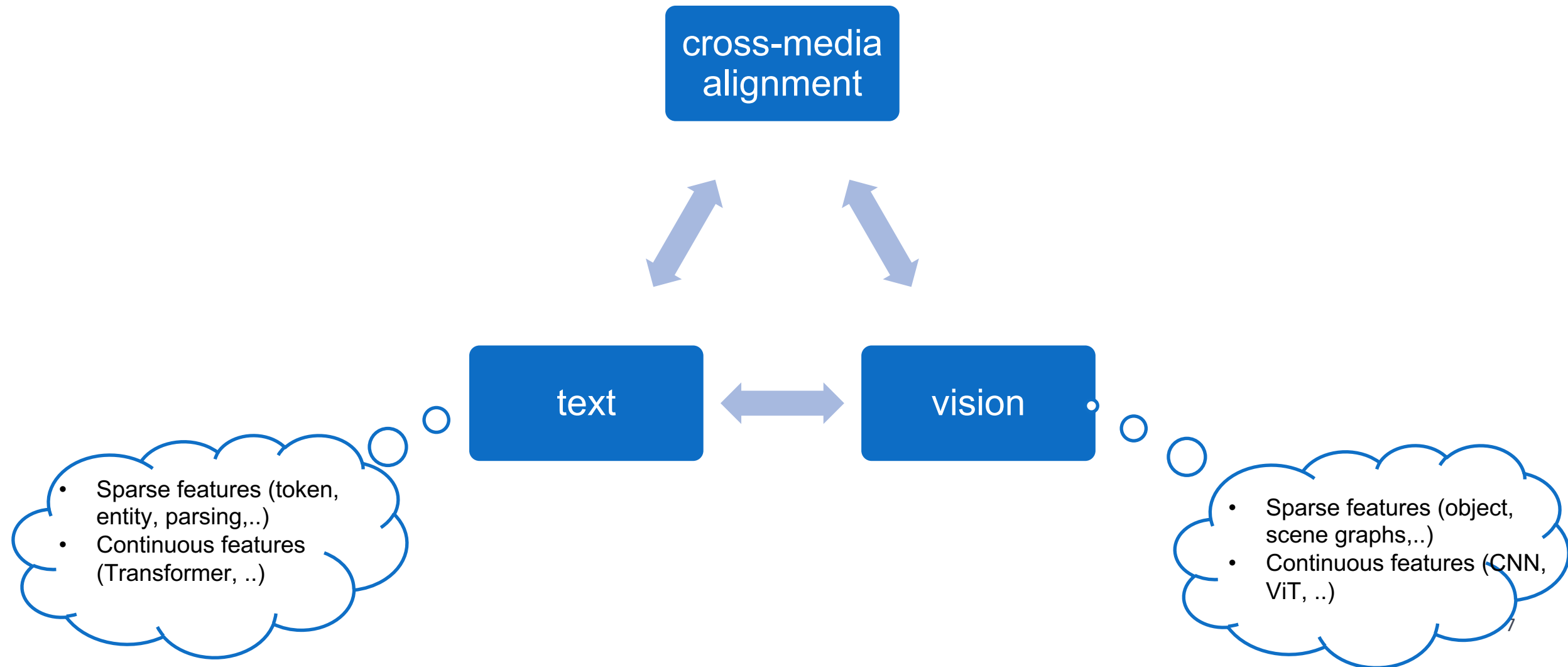


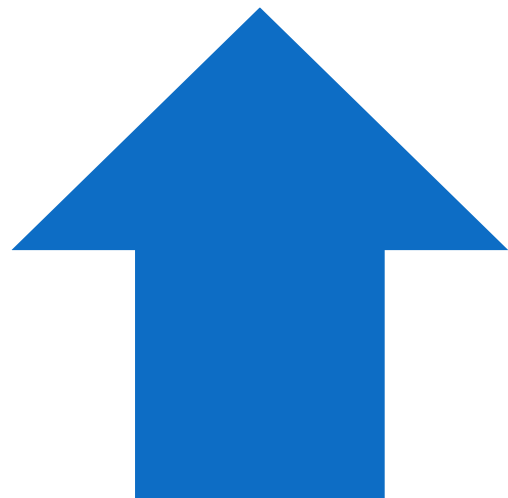
- Sparse features (token, entity, parsing,..)
- Continuous features (Transformer, ..)

Architecture of Vision-Language Pretraining



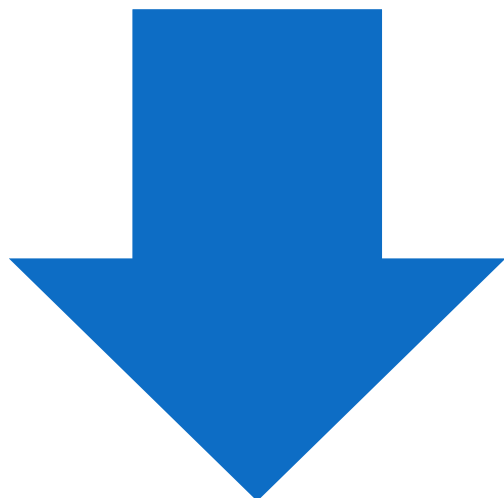
- How? Image-text pre-training by image-text pairs





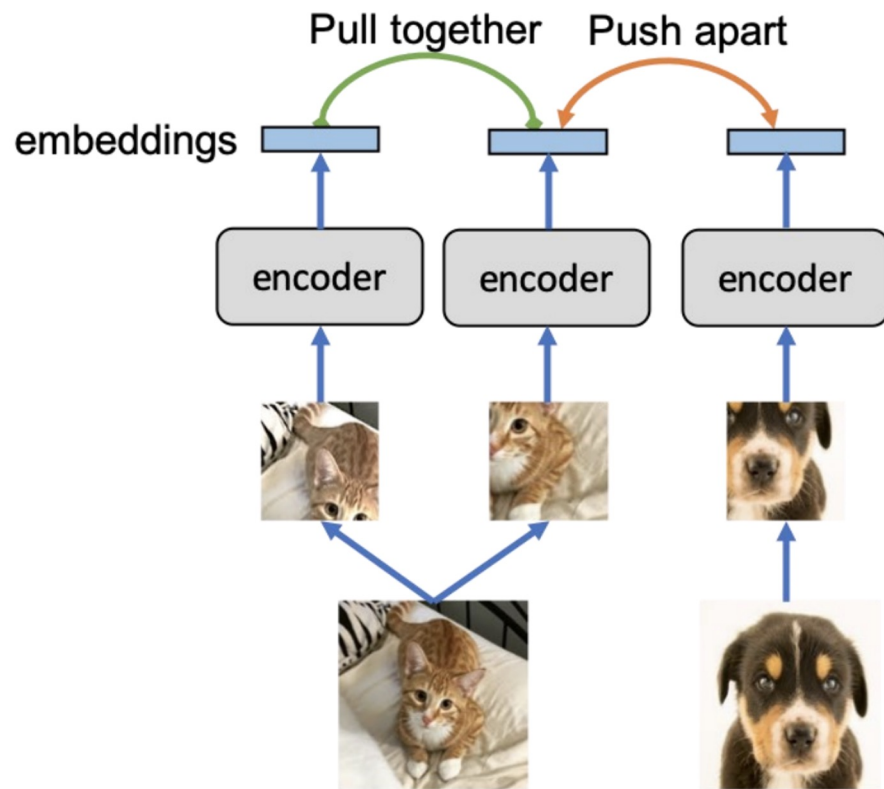
Contrastive Learning

- large batch size required
- data hungry
- hard example sensitive



Generative (Masked Prediction)

- Masked object prediction
- Masked feature regression

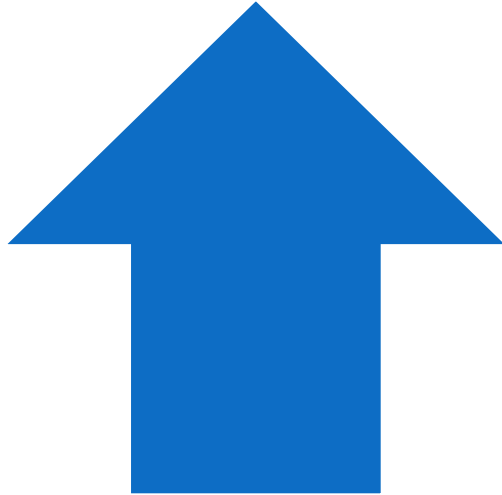


Data augmentation

SimCLR (Chen et al., 2020), MoCo (He et al., 2020), DINO (Caron et al., 2021)

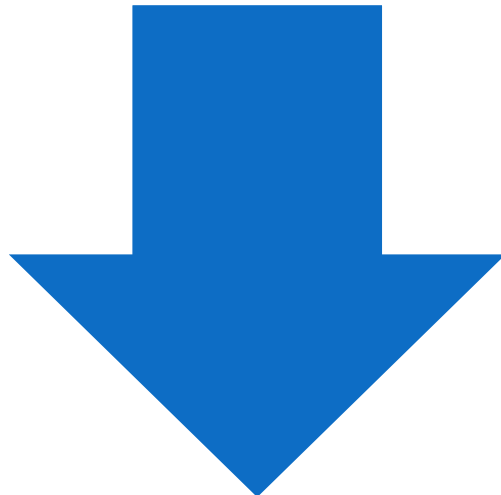
Contrastive Learning

- large batch size required
- data hungry
- hard example sensitive



Contrastive Learning

- large batch size required
- data hungry
- hard example sensitive

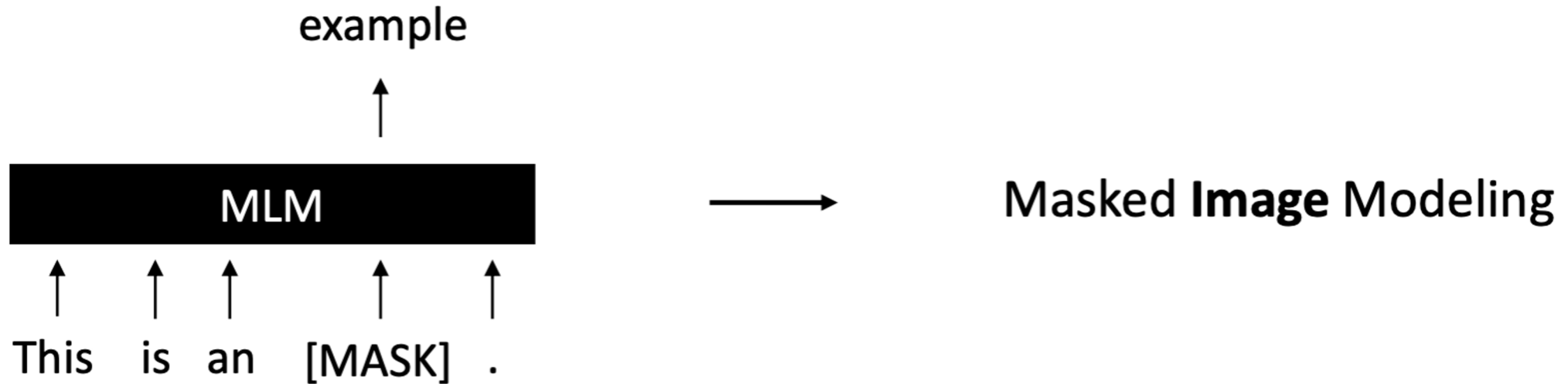


Generative (Masked Prediction)

- Masked feature regression
- Masked object prediction



- Split an image into patches

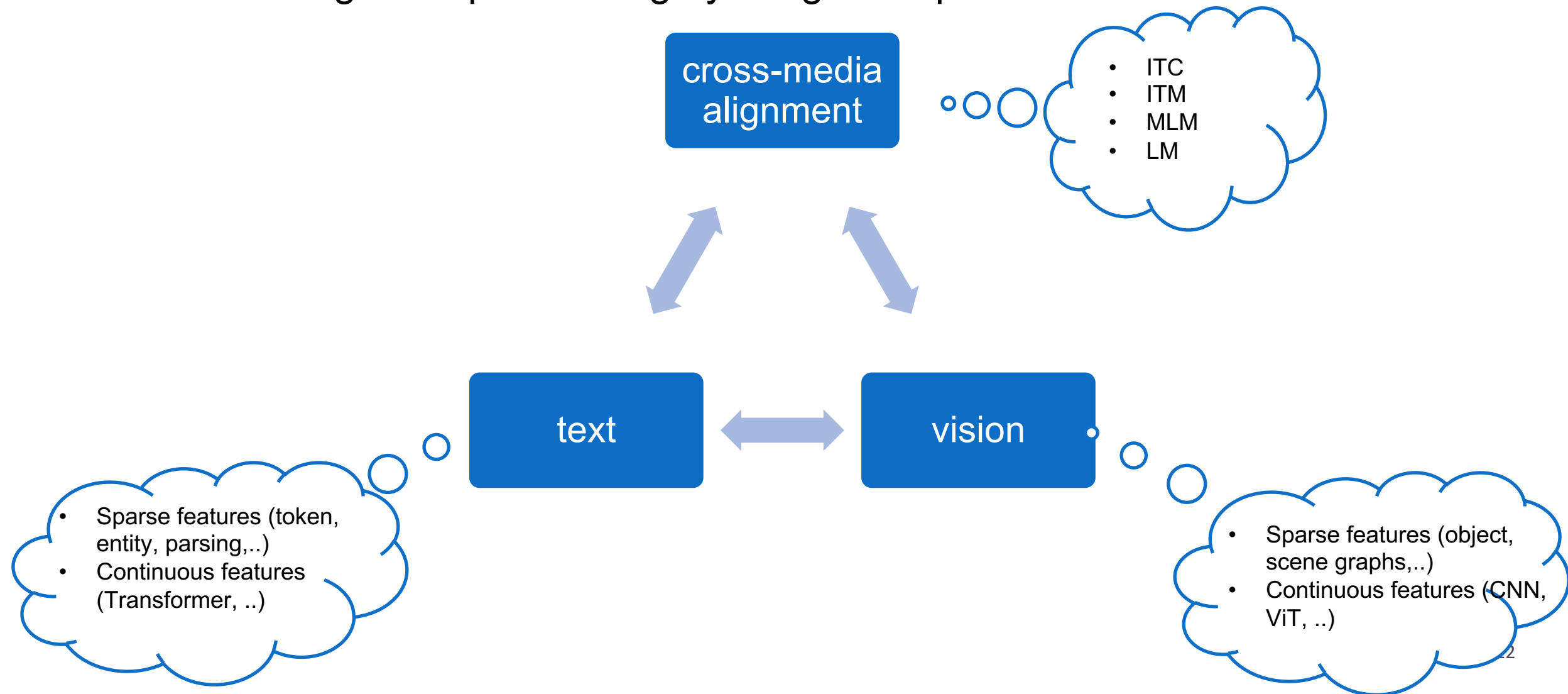


BERT Masked
Language Modeling

The Goal of Vision-Language Pretraining



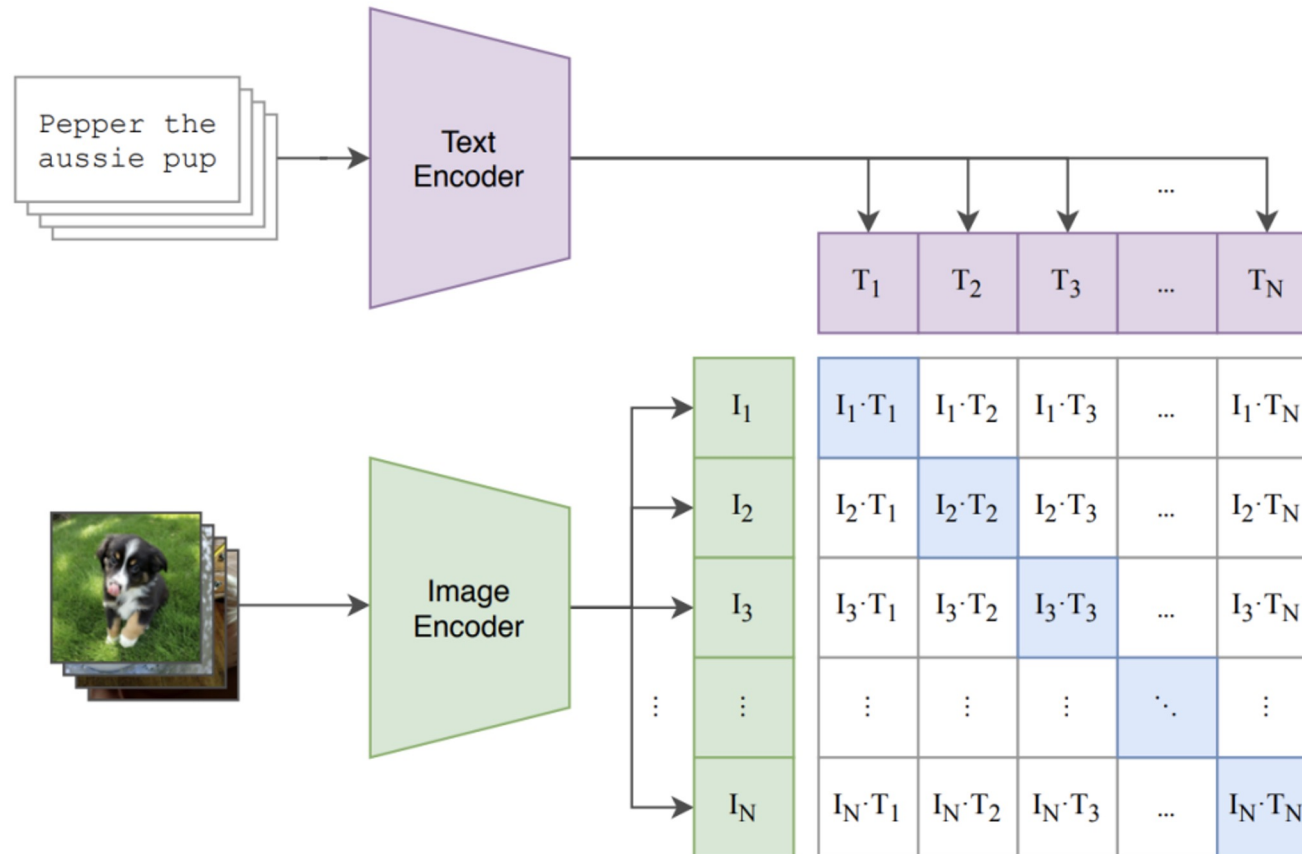
- How? Image-text pre-training by image-text pairs



Typical Loss: Image-text contrastive (ITC)



- Image-text contrastive (ITC) loss



A Simple Framework for Contrastive Learning of Visual Representations, 2020

Momentum Contrast for Unsupervised Visual Representation Learning, 2019

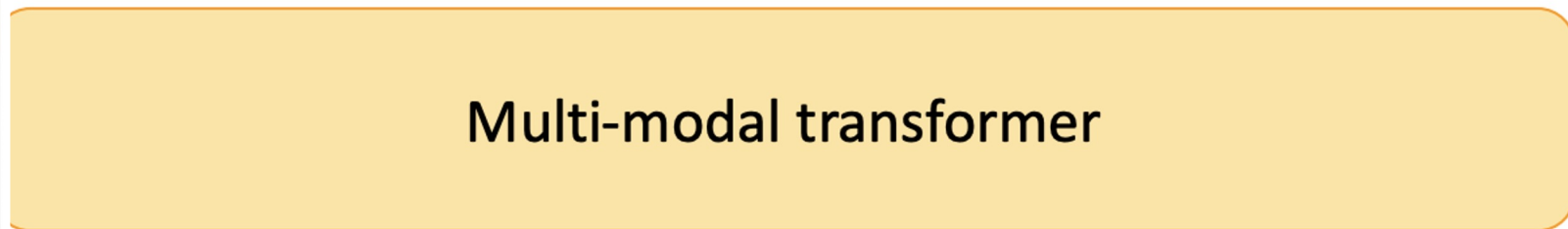
Align before Fuse: Vision and Language Representation Learning with Momentum Distillation, 2021

UFO: A UniFied TransFormer for Vision-Language Representation Learning, 2021

Typical Loss: Image-text matching (ITM) loss



Paired? Yes or no



BOS



a



dog



is



sitting



on



a

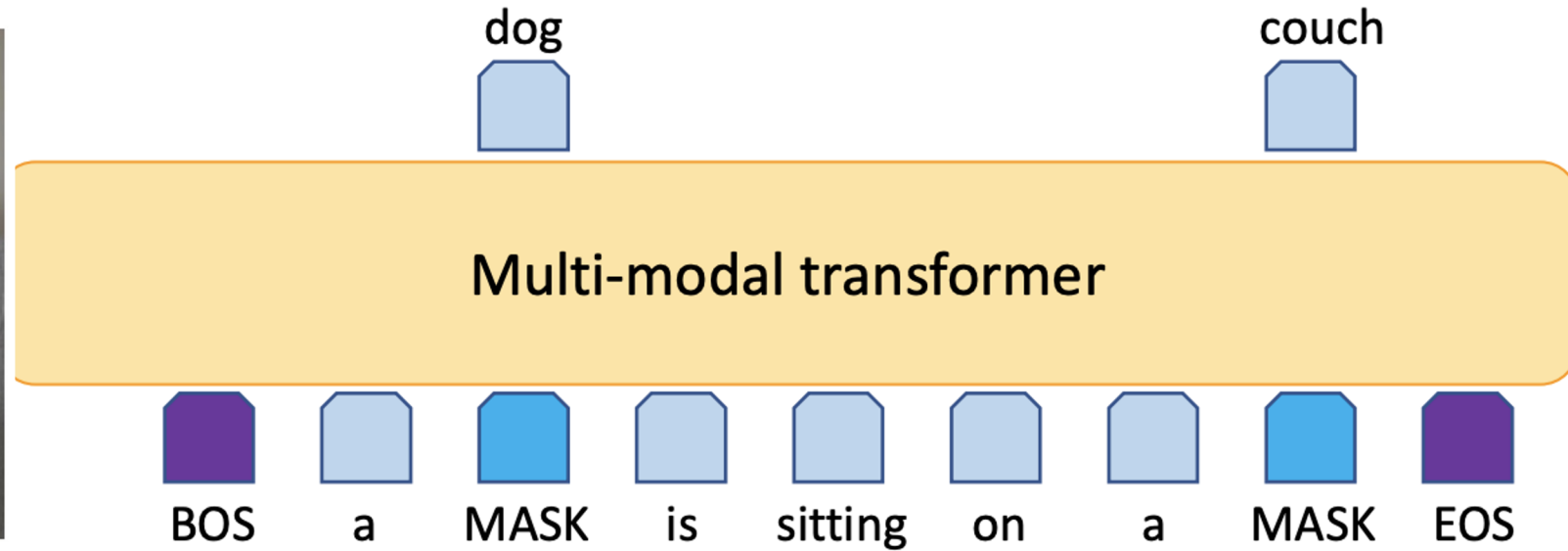


couch

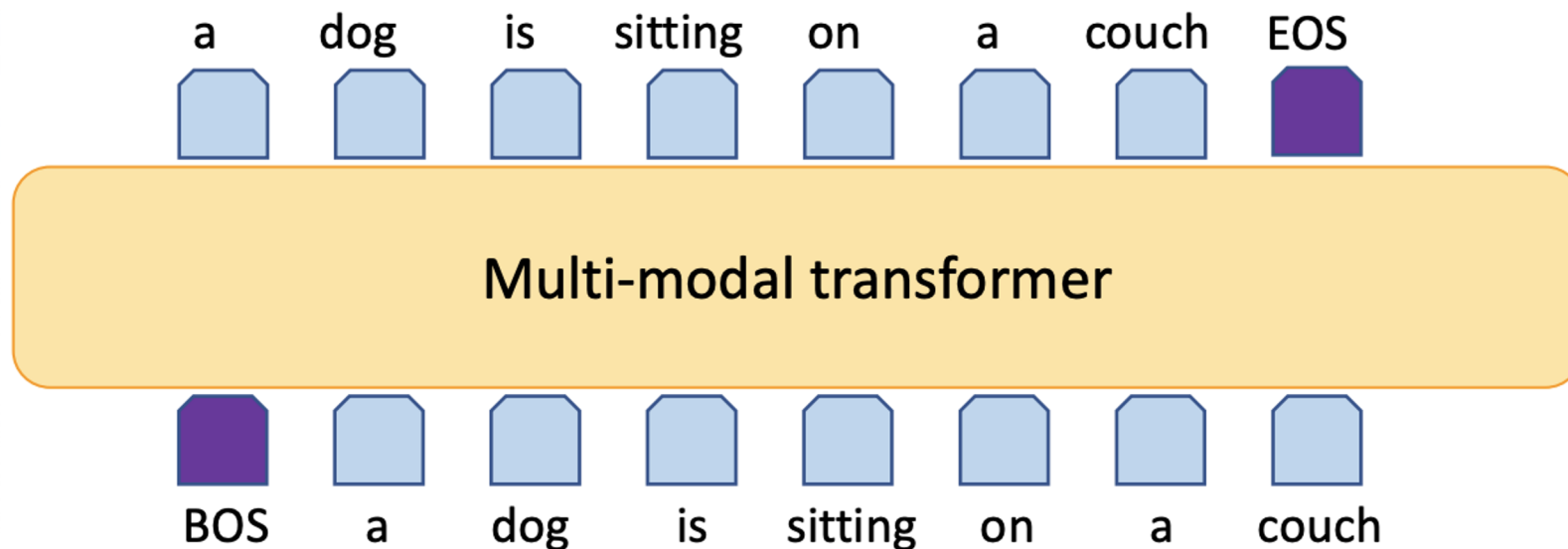


EOS

Typical Loss: Masked language modeling (MLM) loss



Typical Loss: Language modeling (LM) loss



Why Knowledge?



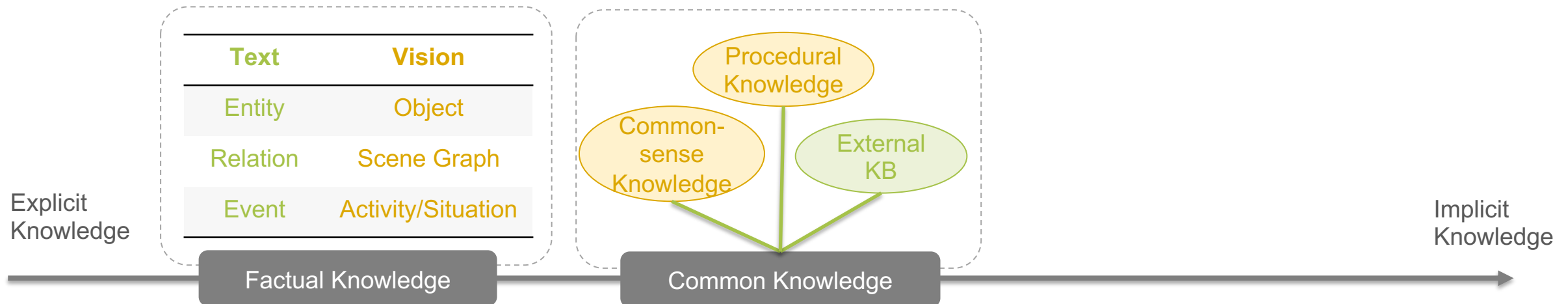
Compared to raw data, knowledge is **important and useful information**.



Why Knowledge?



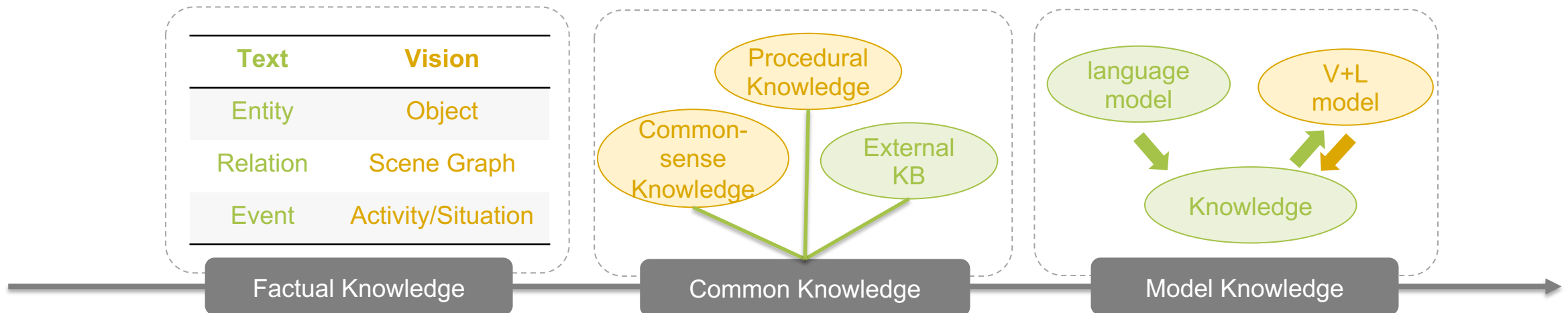
Compared to raw data, knowledge is **important and useful information**.

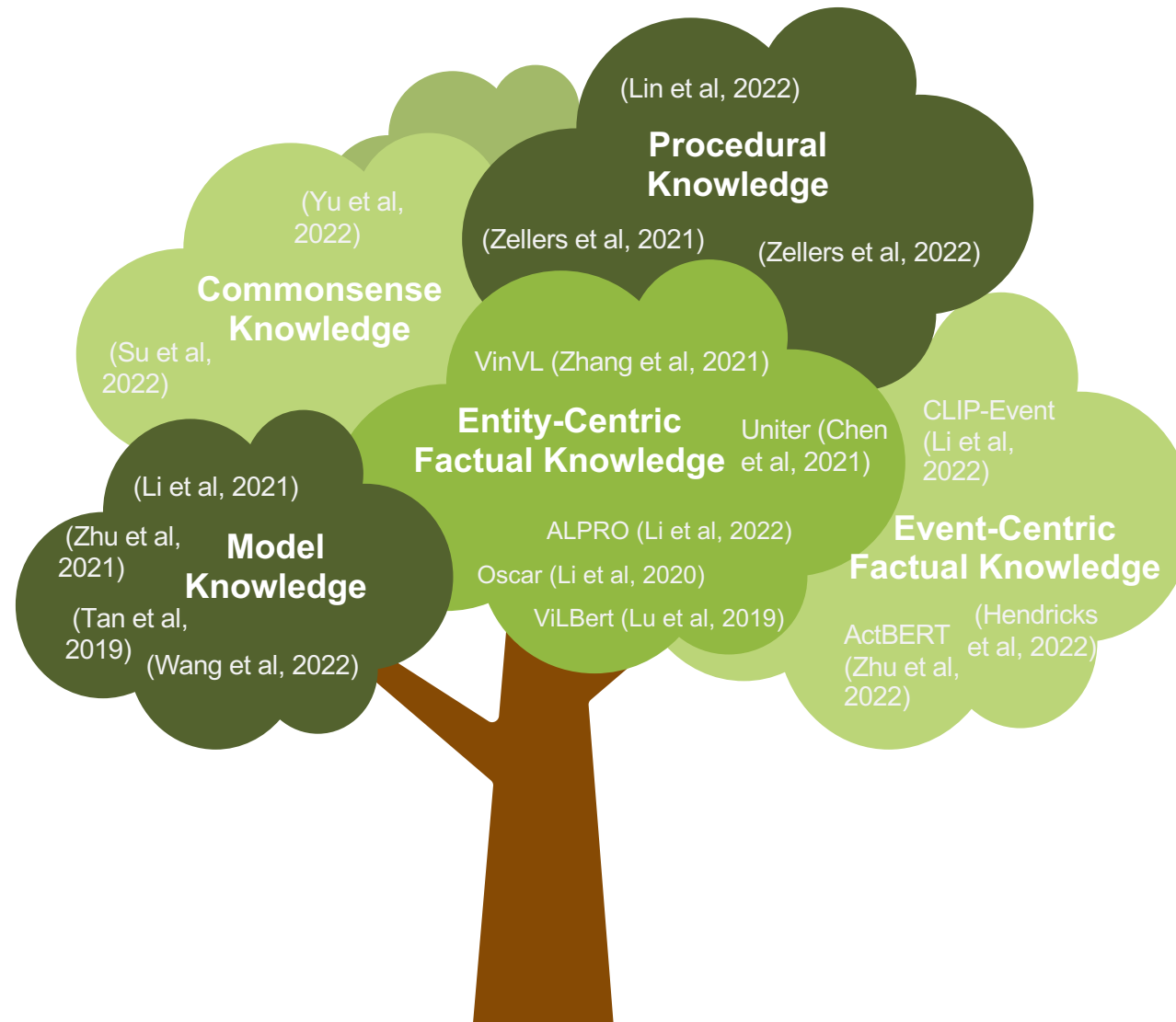


Why Knowledge?



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

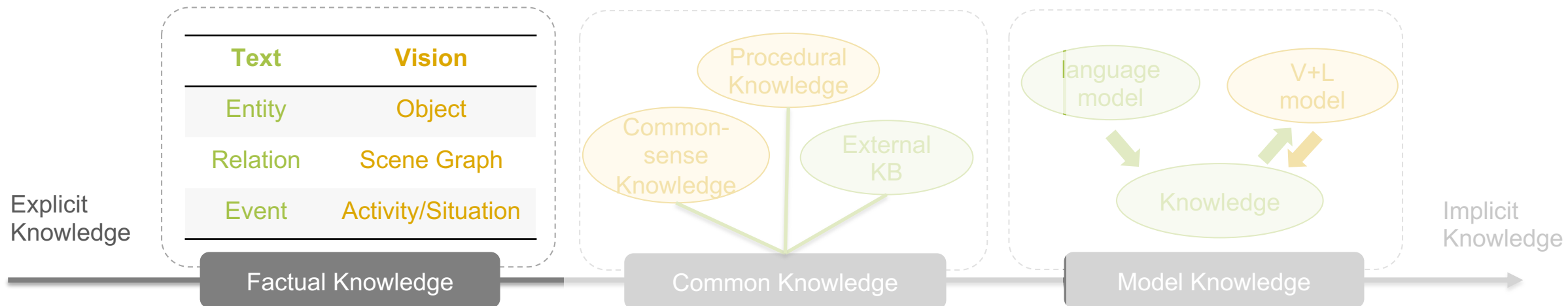


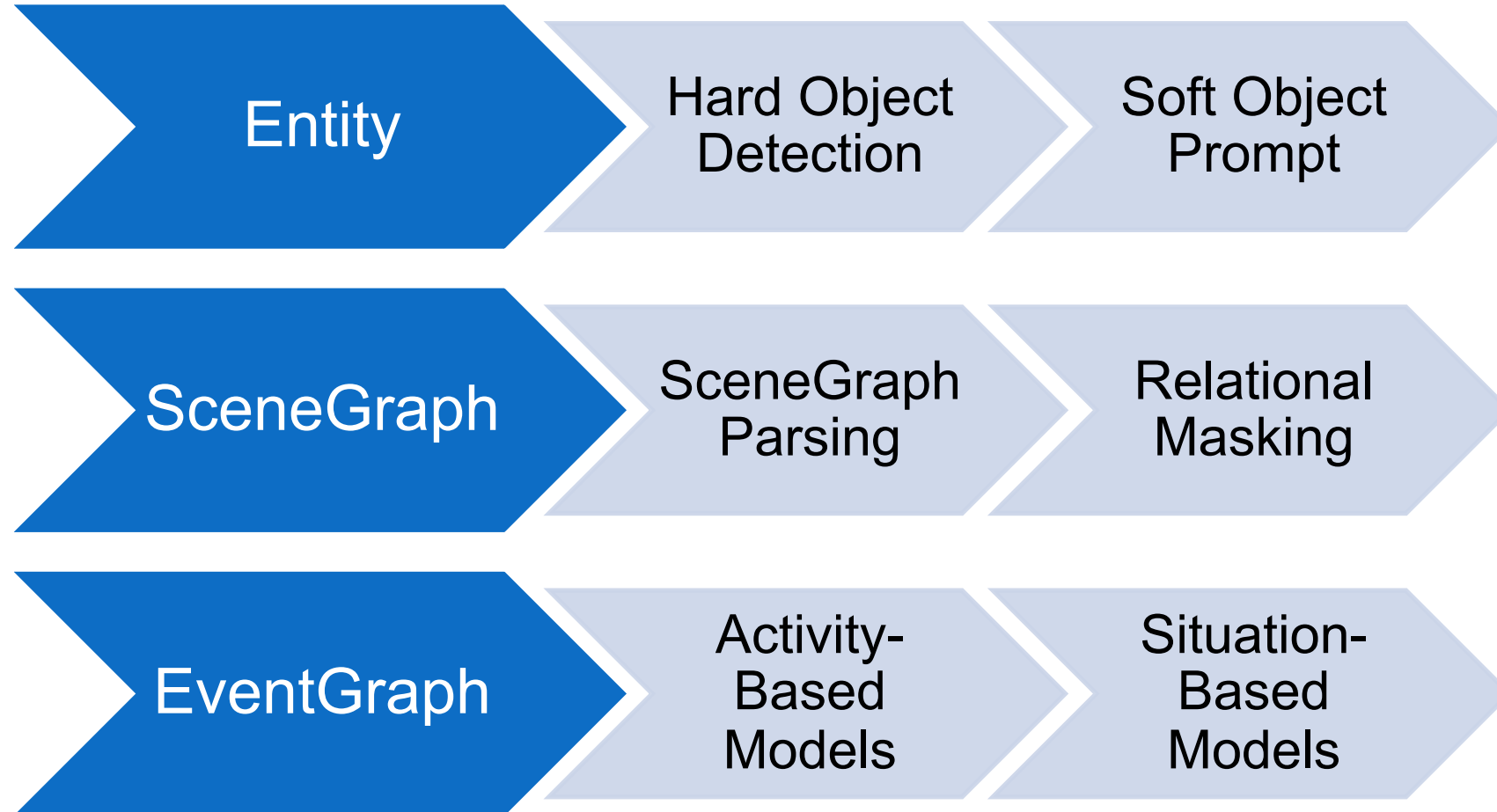


Factual Knowledge



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

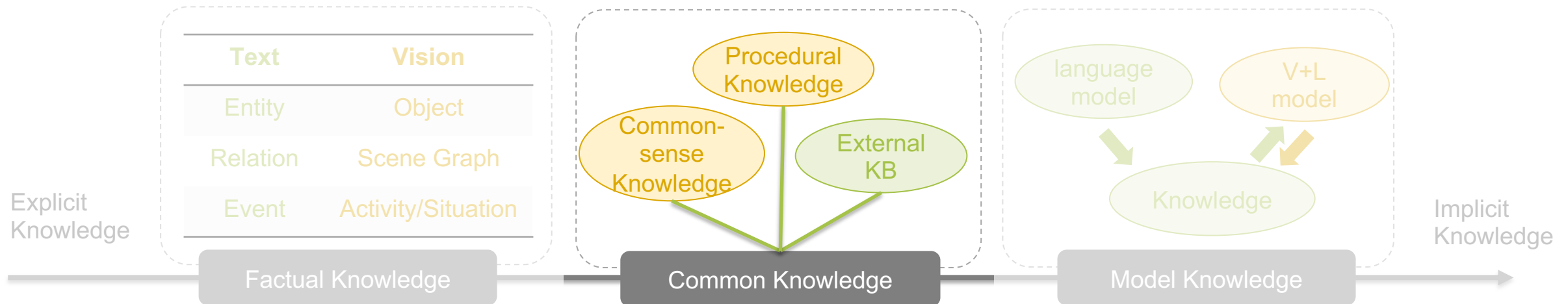




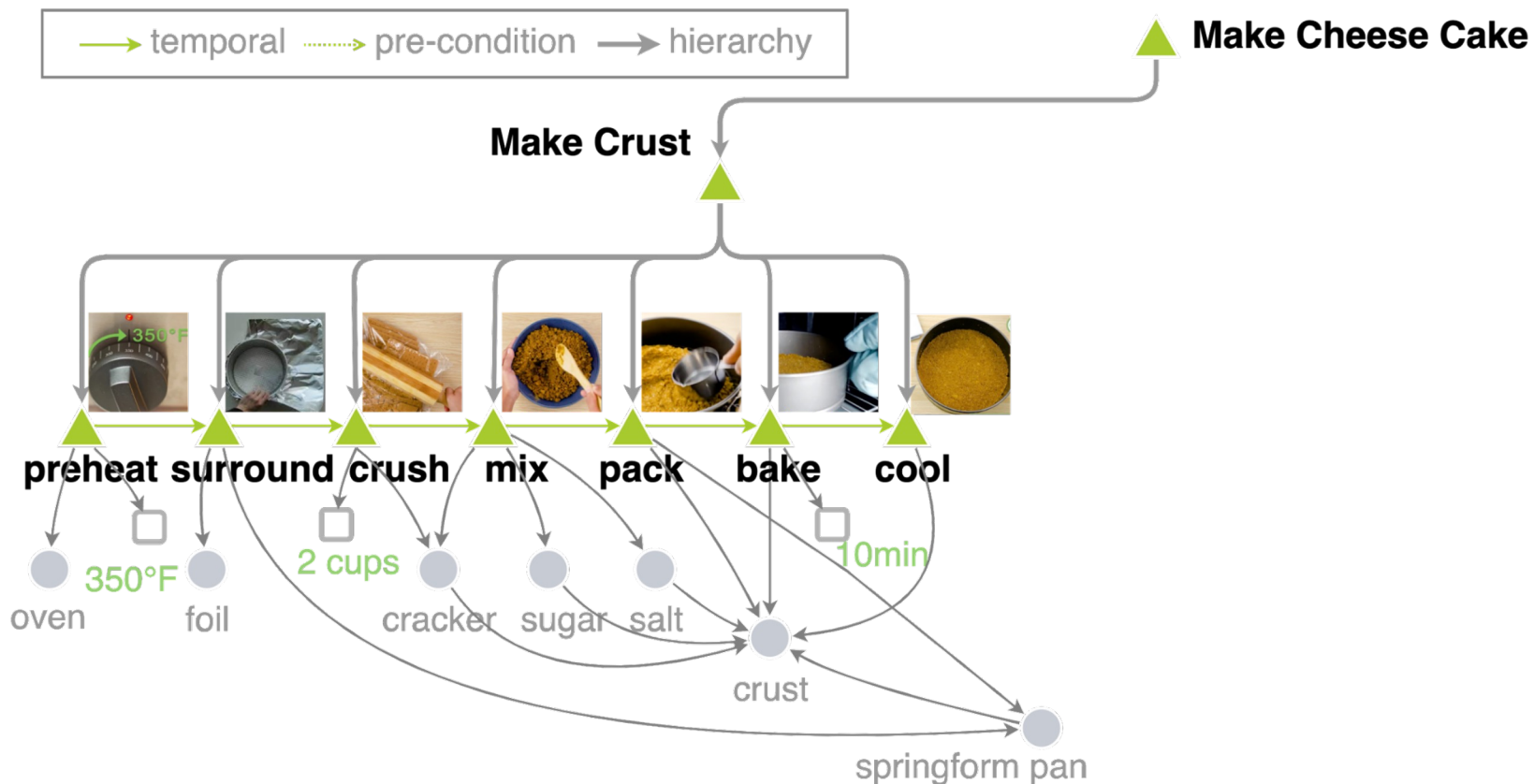
Common Knowledge



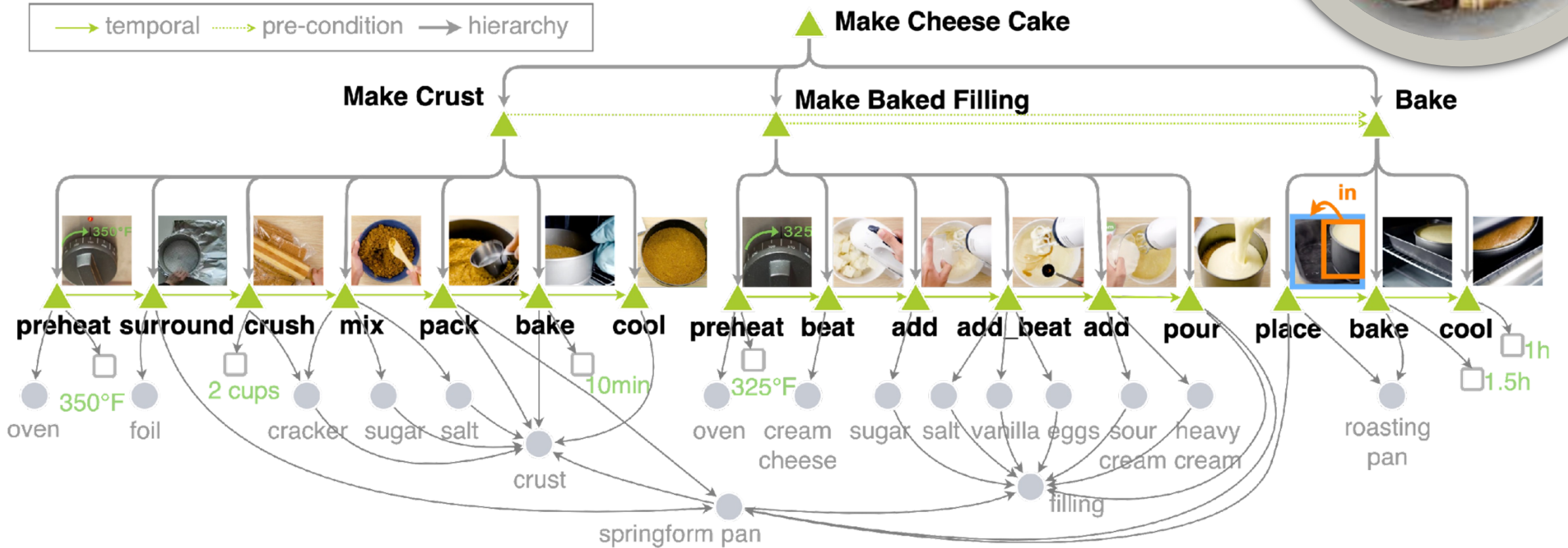
Compared to raw data, knowledge is **important and useful information**.



History repeats itself

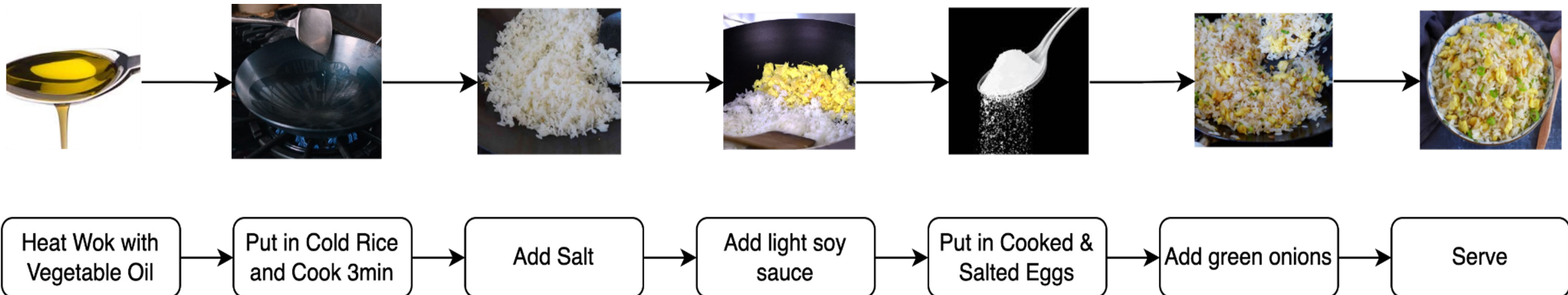


History repeats itself



Representative Resource: wikiHow

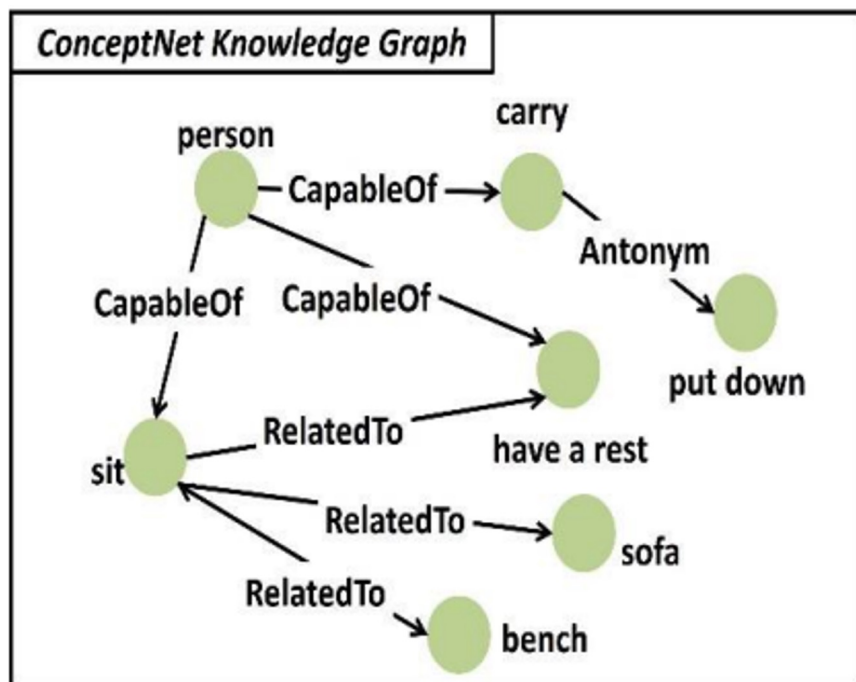
Current online instructional video/text describe a task as completing a sequential set of steps, assuming that all tasks follow a linear schema



Commonsense Knowledge



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

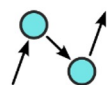


Why are [person1] and [person3] shaking hands?

- (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



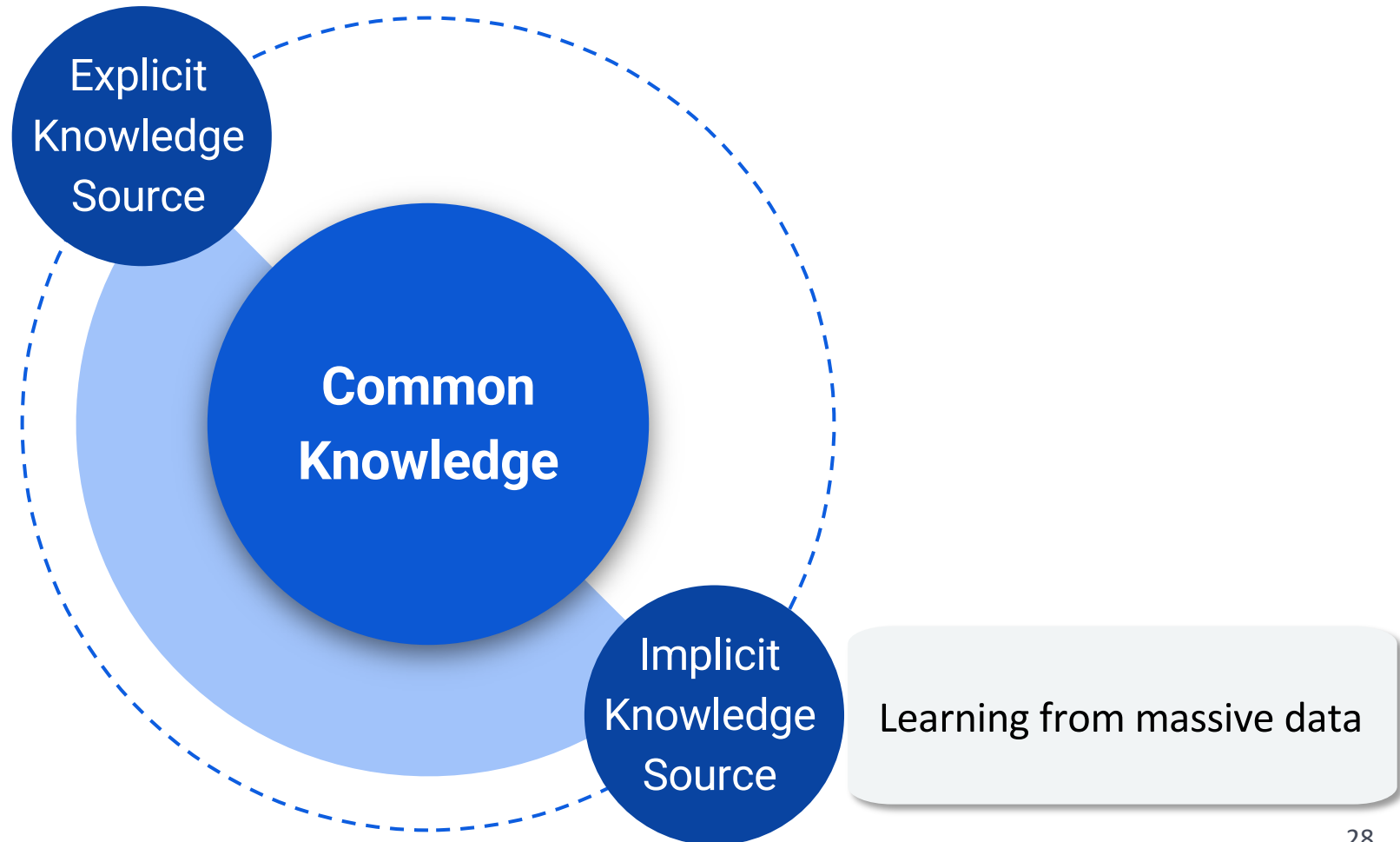
How to add common knowledge?



- Two ways to learn procedural knowledge

Use knowledge:

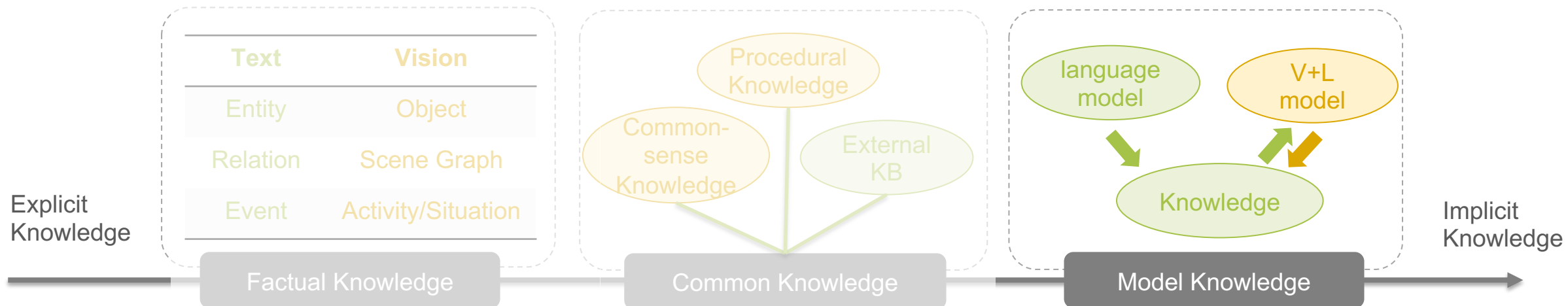
- As Data
- As Supervision
- In Model

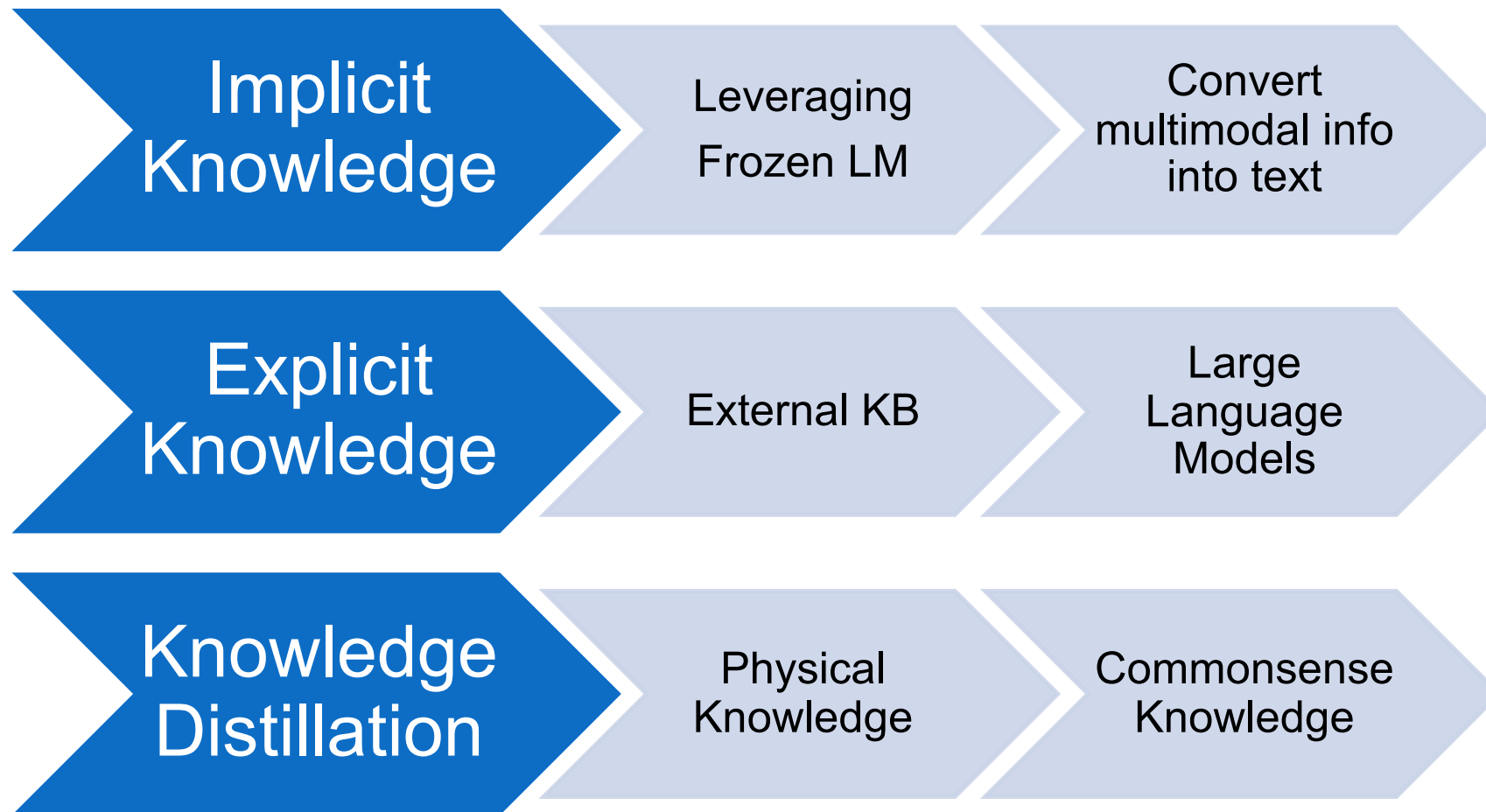


Model Knowledge

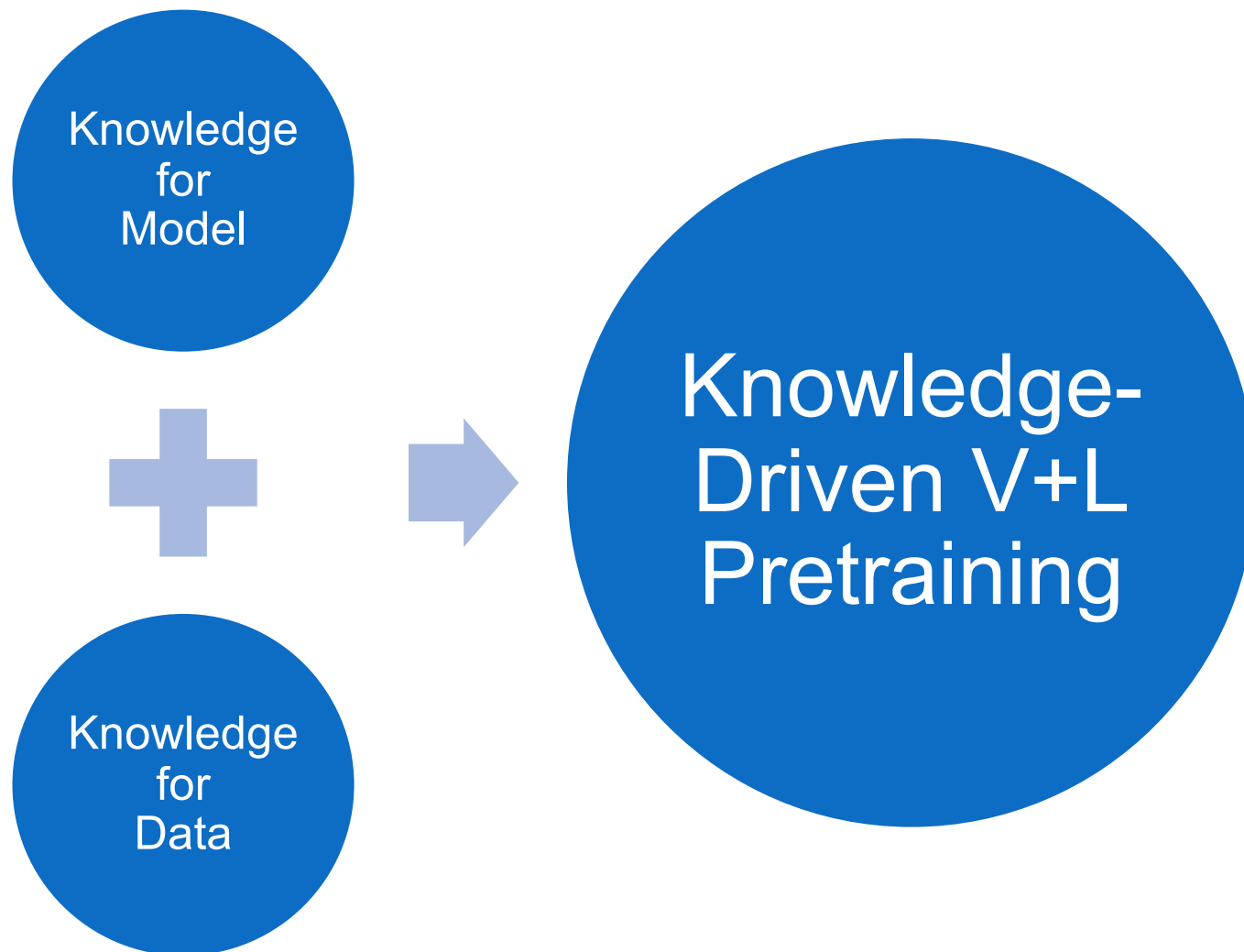


Compared to raw data, knowledge is **important and useful information**.





How to learn multimedia embedding?



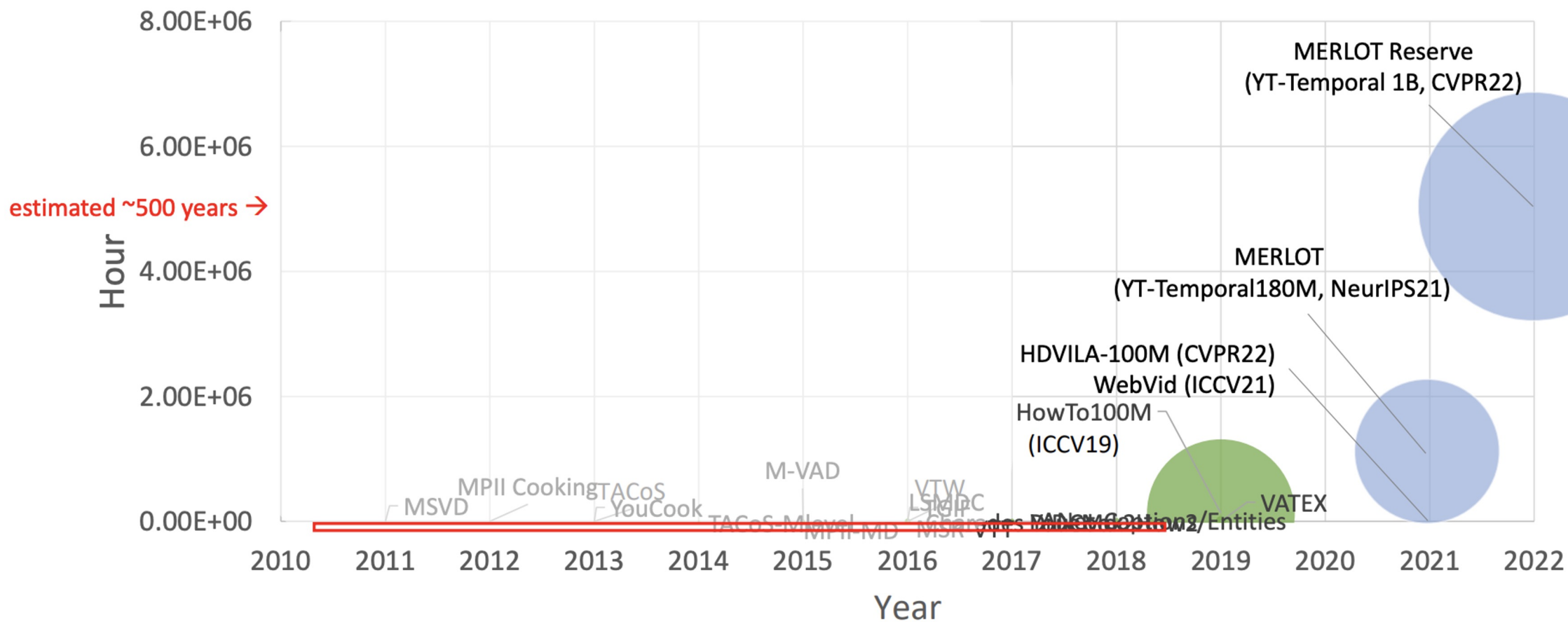


- Image-caption pairs
 - Free of cost
 - The state-of-the-art model CLIP (from OpenAI) uses ~400M pairs for model training
- Object and Scene Graph annotations in Popular datasets
 - Flickr 30K (~30K images)
 - MS COCO (~330K images)
 - Visual Genome (~108K images)
 - Conceptual Captions (~3.3M images)
 - SBU Captions (~1M images)
 - ...

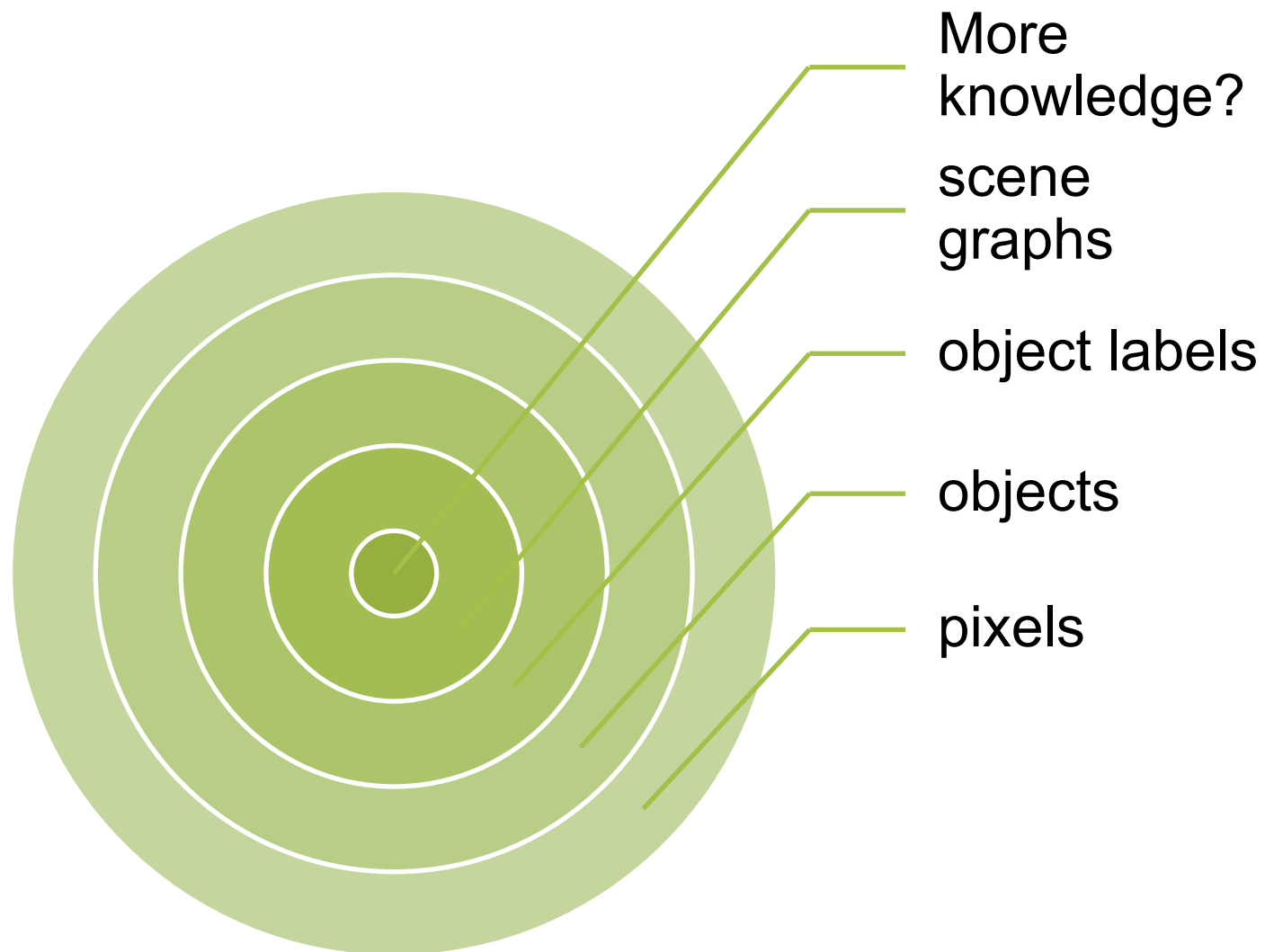
Resources: Video-and-Language Datasets



Total Video Duration



Moving towards...





- On the model side, Transformer based models using self supervision has achieved great success in multiple downstream tasks.
 - Adding knowledge can guide the model where to focus.
 - **Structured knowledge** (such as event graph structure) and **abstract word understanding** (such as verb, adjectives, etc) are still lack of exploration.
- On the data side, knowledge is useful in the following ways:
 - In-context prompt
 - data augmentation
 - data selection

Content	Time	Presenter
Motivation and Overview	15min	Manling Li
Factual Knowledge	15min	Manling Li
Procedural Knowledge	30min	Xudong Lin
Commonsense Knowledge and Model Knowledge	30min	Jie Lei
Panel: Knowledge vs Large Models	20min	Heng Ji, Mohit Bansal, Shih-Fu Chang
Panel: Text vs Image vs Video vs Others	20min	Heng Ji, Mohit Bansal, Shih-Fu Chang
Panel: Open Challenges	20min	Heng Ji, Mohit Bansal, Shih-Fu Chang
QA	30min	All