



ACL 2020



THE OHIO STATE UNIVERSITY



NATIONWIDE CHILDREN'S  
When your child needs a hospital, everything matters.

# Rationalizing Medical Relation Prediction from Corpus-level Statistics

**Zhen Wang**  
**The Ohio State University**

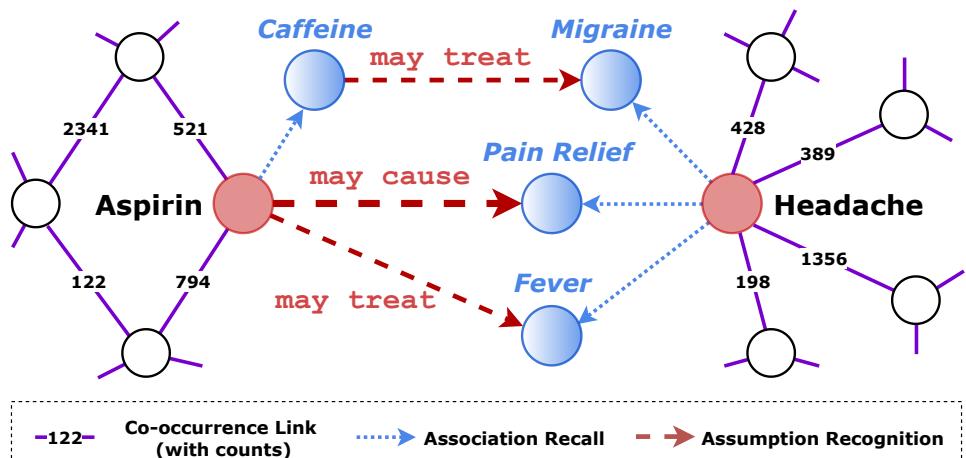
In collaboration with Jennifer Lee (NCH), Simon Lin (NCH),  
Huan Sun (OSU)

# One-Minute Summary

## What is this paper about?

**Task:** Predicting relations between two given terms from a text corpus

**Goal:** Make accurate prediction & Provide justifications for it (Rationalization)



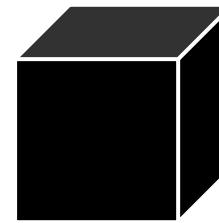
1. Draw inspirations from human memory **recall** and
2. **Recall** global **associations** (blue) from **corpus-level statistics** and **recognize** meaningful relational  (red) between them
3. Justify the prediction (*Aspirin may treat headache*) by highlighting important associations and assumptions

# Black-Box Relation Classifier

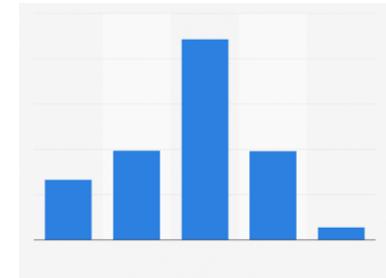
Medical Records



Input Pair



Black-Box Classifier



Class Distribution

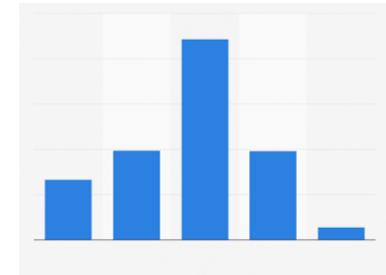
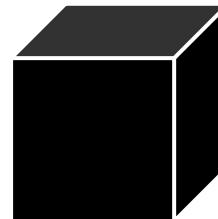
# Black-Box Relation Classifier

Medical Records



Input Pair

Aspirin  
Headache



Class Distribution

Black-Box Classifier

**High Risk!**



Medicine



Finance



Judiciary

# Can We Open the Black-Box?

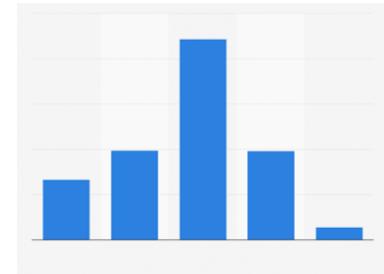
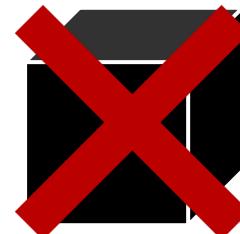
Medical Records



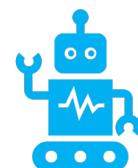
Aspirin

Headache

Input Pair



Class Distribution



# Can We Open the Black-Box?

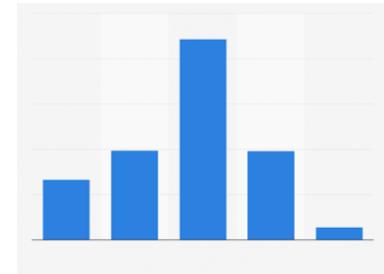
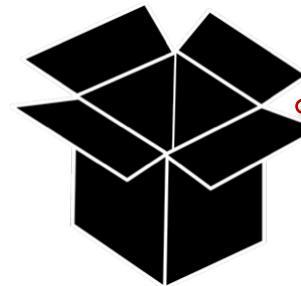
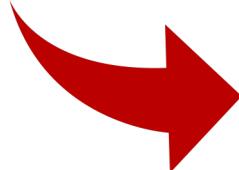
Medical Records



Aspirin

Headache

Input Pair



Class Distribution



To solve the problem, we draw inspirations from human memory theories – Recall and Recognition!

# Memory Recall and Recognition

## Recall

Retrieve association information  
from long-term memory

## Recognition

Identify previously learned  
information

# Memory Recall and Recognition

## Recall

Retrieve association information  
from long-term memory



Recall my  
friend's name

*What is his name?*

## Recognition

Identify previously learned  
information

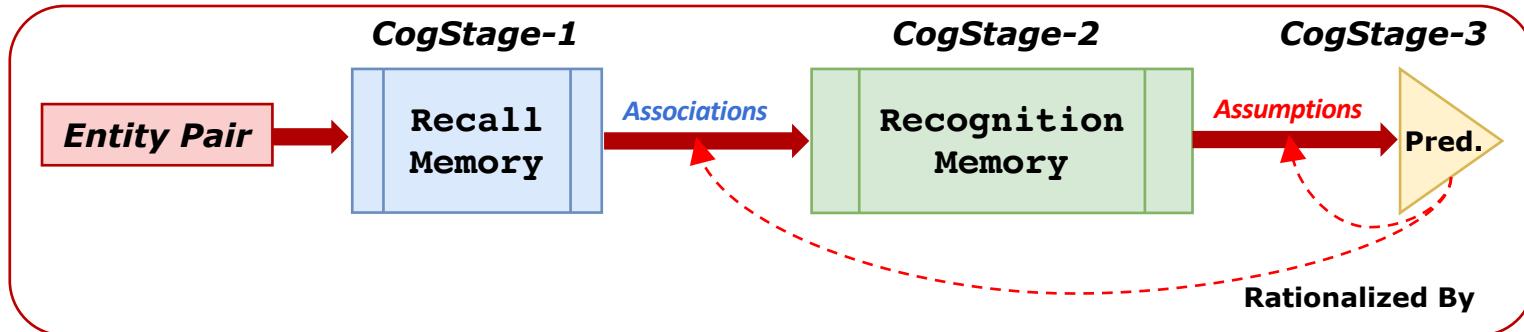


Recognize my  
friend's face

*Hello, my friend!*

Example from <http://webdesign-review.blogspot.com/2016/04/recognition-vs-recall-in-mobile-web.html>

# Rationalizing Medical Relation Prediction



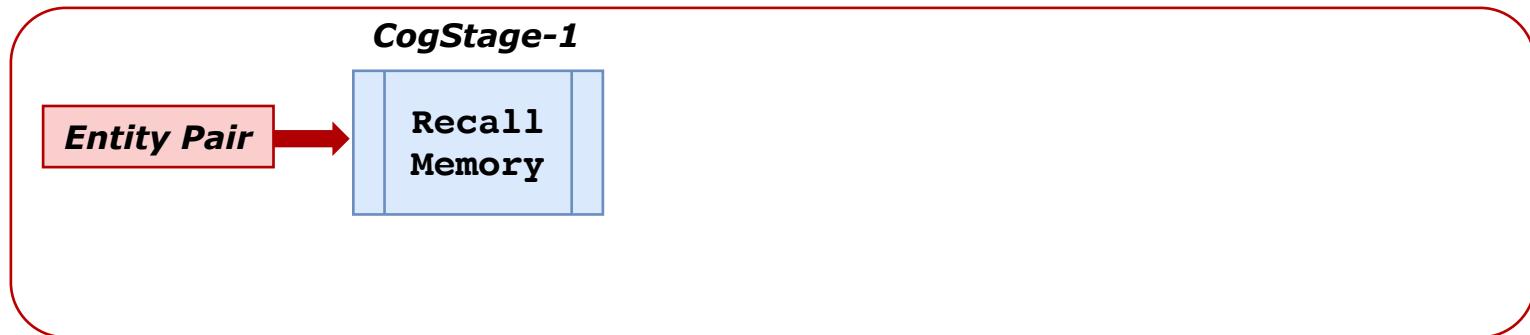
# Rationalizing Medical Relation Prediction

*Entity Pair*

Aspirin 

 Headache

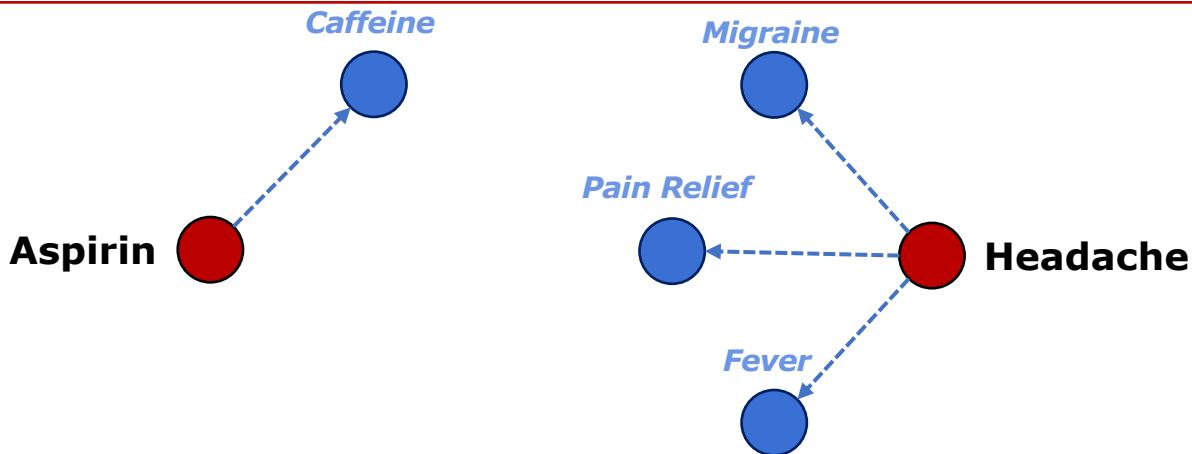
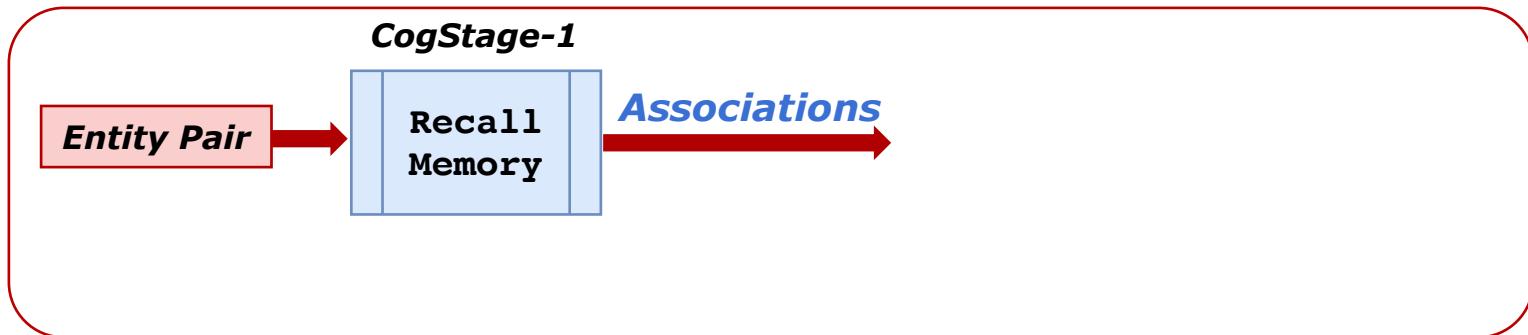
# Rationalizing Medical Relation Prediction



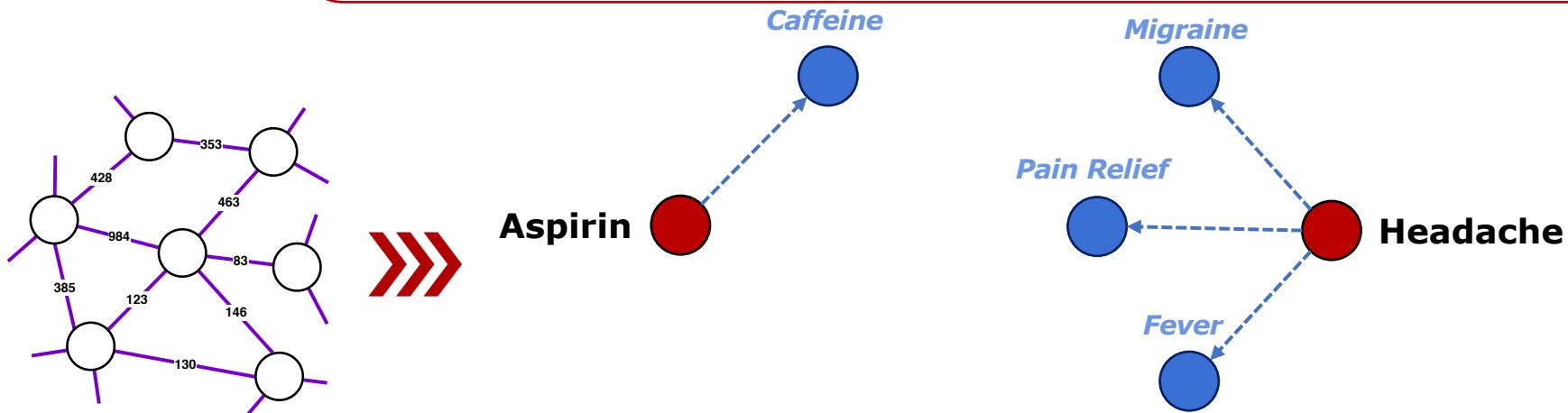
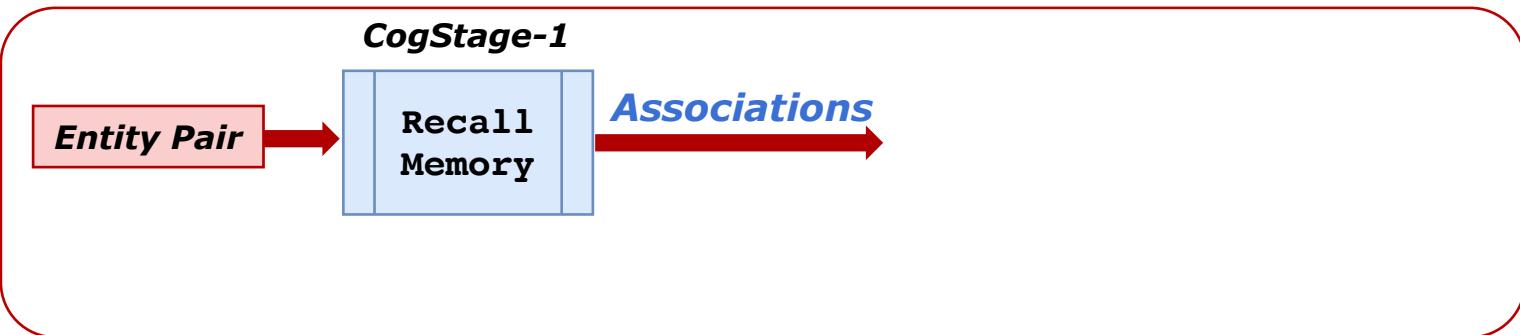
Aspirin 

 Headache

# Rationalizing Medical Relation Prediction

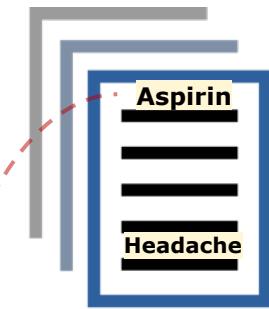


# Rationalizing Medical Relation Prediction



# Rationalizing Medical Relation Prediction

Medical Records

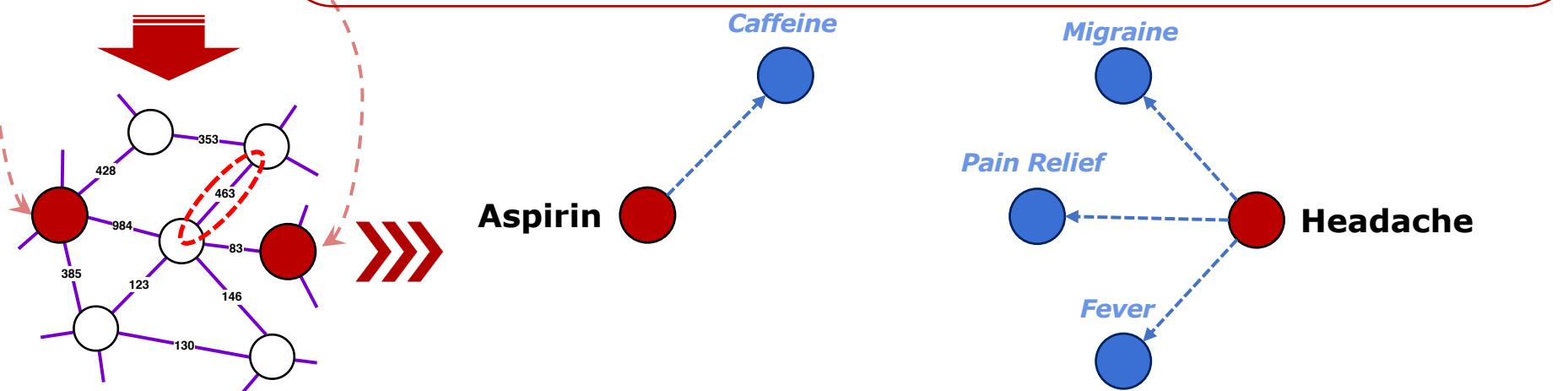


CogStage-1

*Entity Pair*

Recall  
Memory

*Associations*

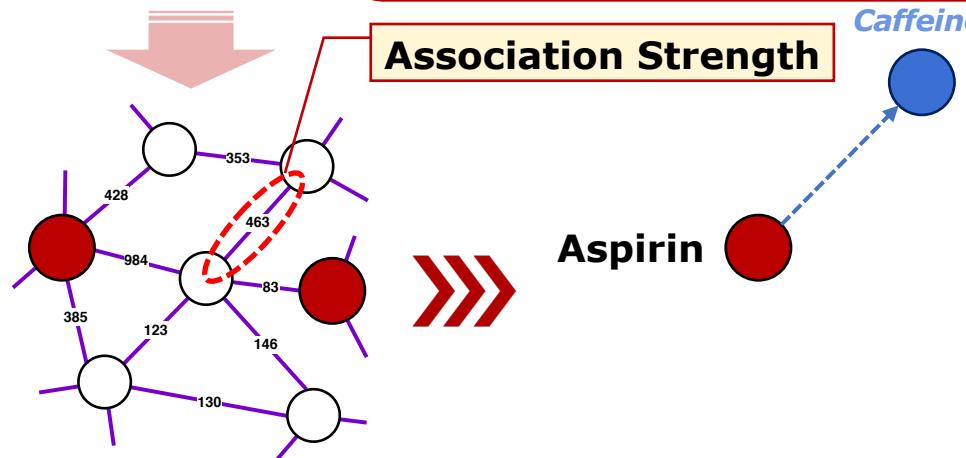


# Rationalizing Medical Relation Prediction

Medical Records



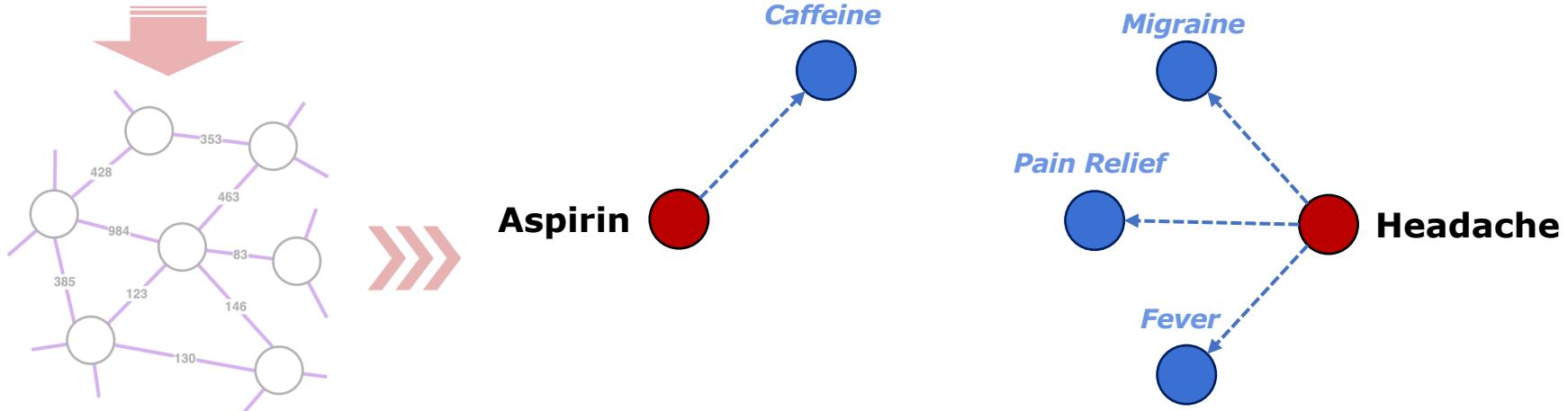
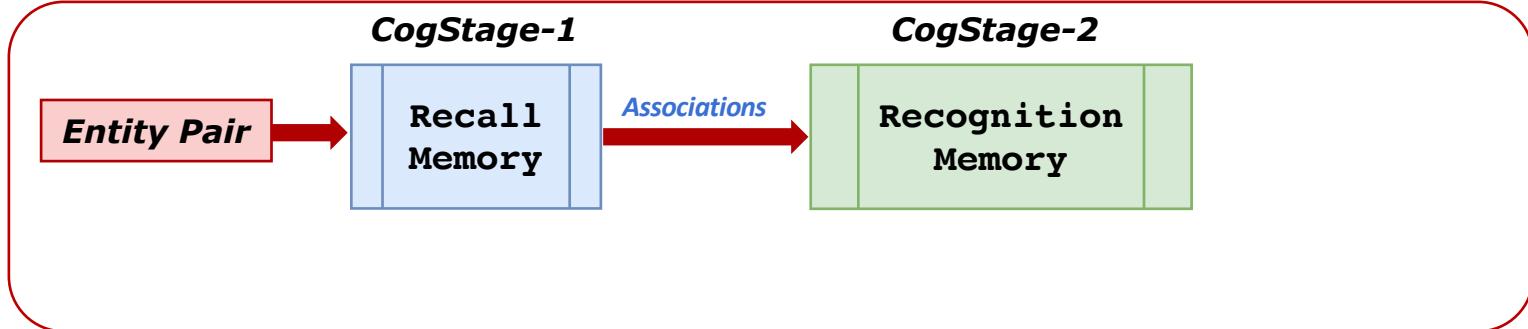
**Association Strength**



Corpus-level Statistics

# Rationalizing Medical Relation Prediction

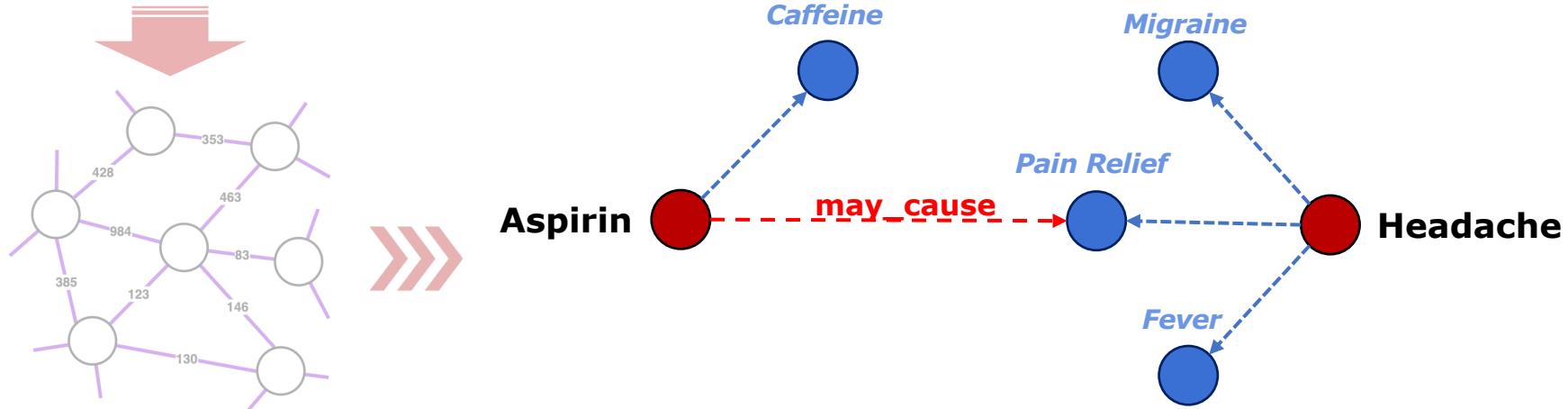
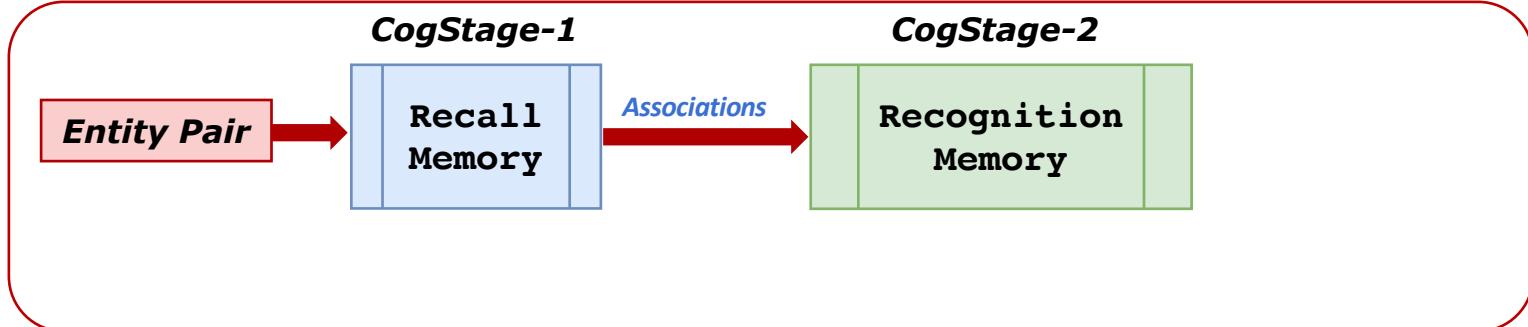
Medical Records



Corpus-level Statistics

# Rationalizing Medical Relation Prediction

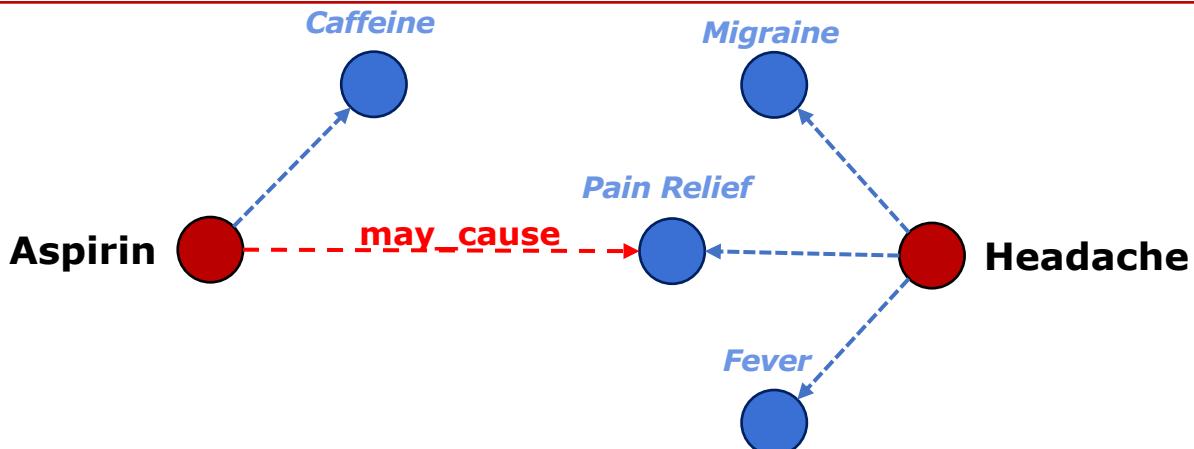
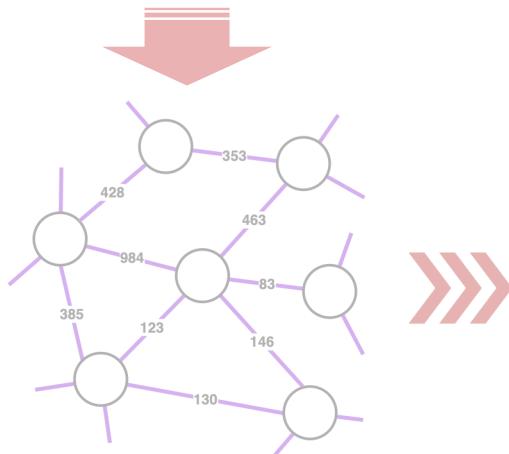
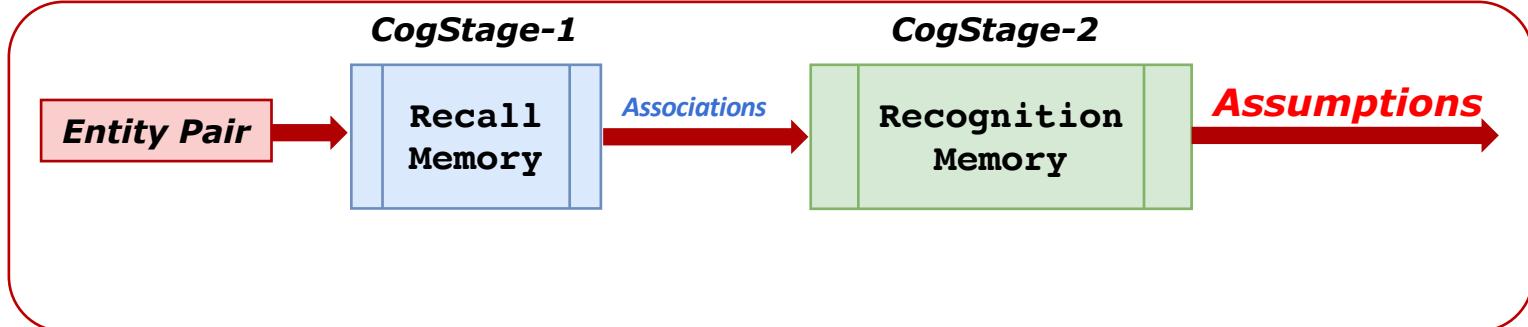
Medical Records



Corpus-level Statistics

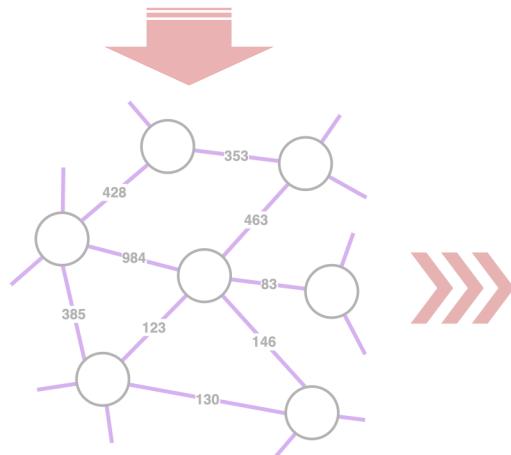
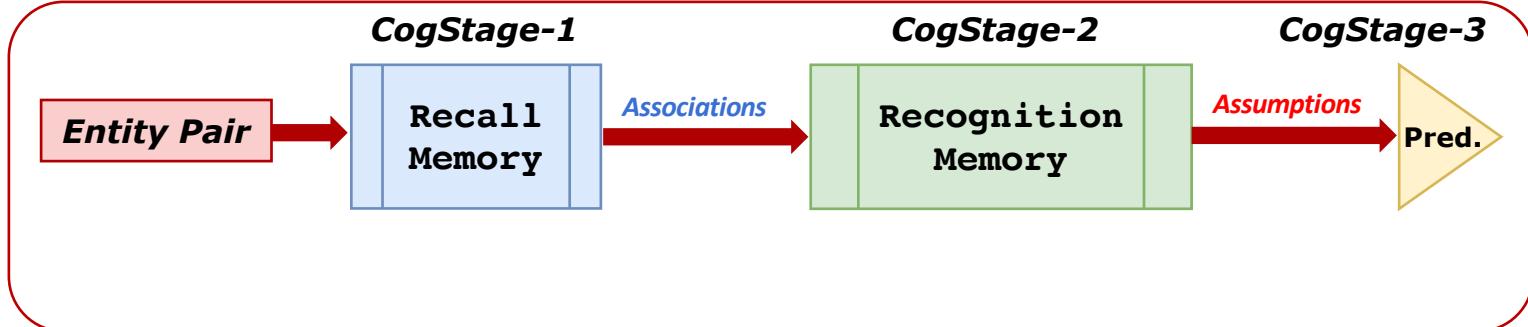
# Rationalizing Medical Relation Prediction

Medical Records



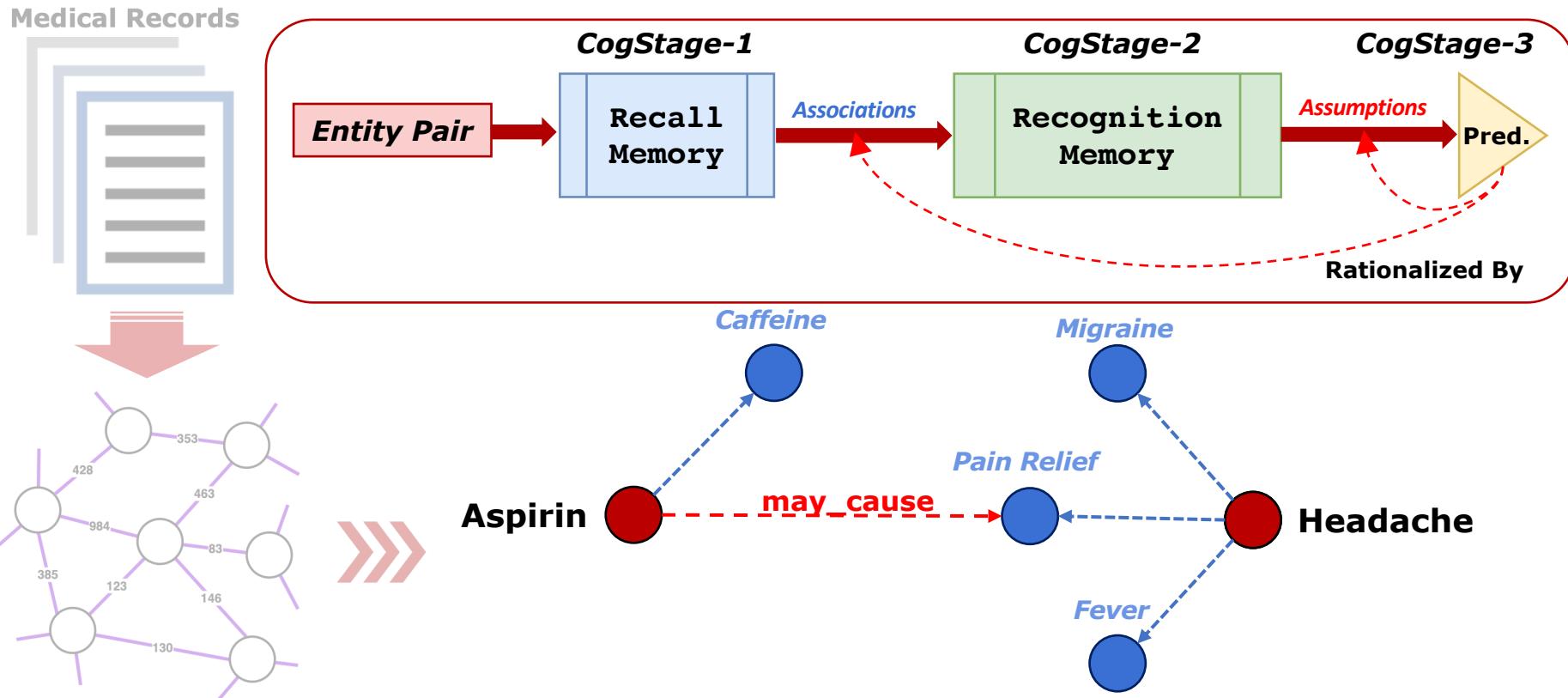
# Rationalizing Medical Relation Prediction

Medical Records



Corpus-level Statistics

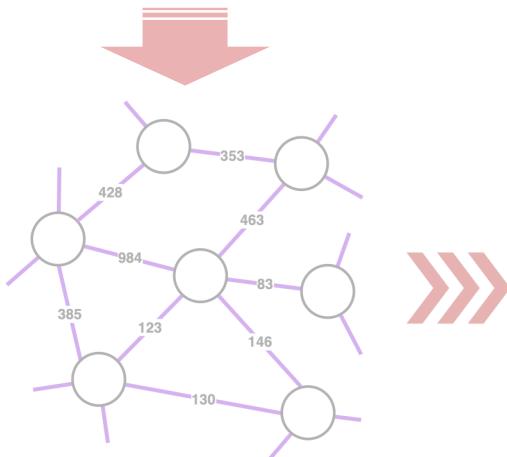
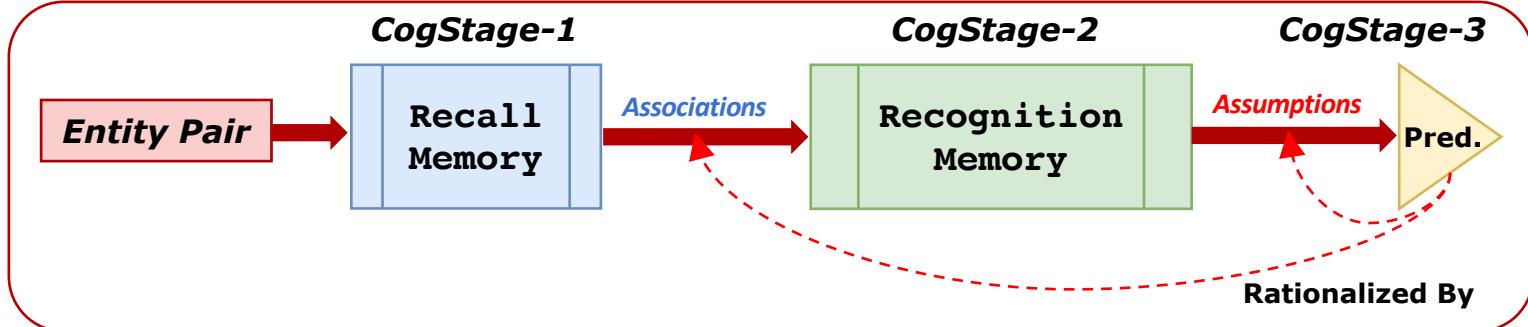
# Rationalizing Medical Relation Prediction



## Corpus-level Statistics

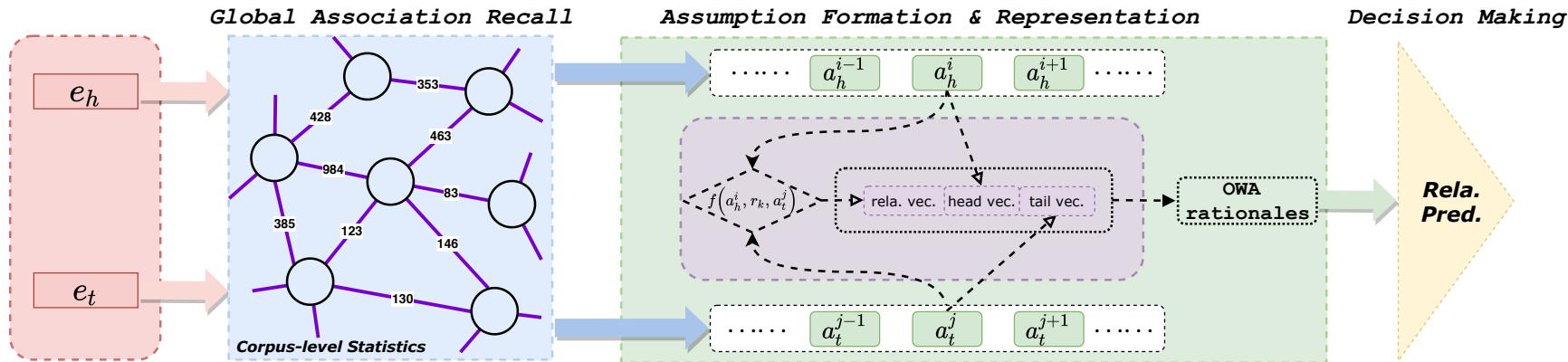
# Rationalizing Medical Relation Prediction

Medical Records

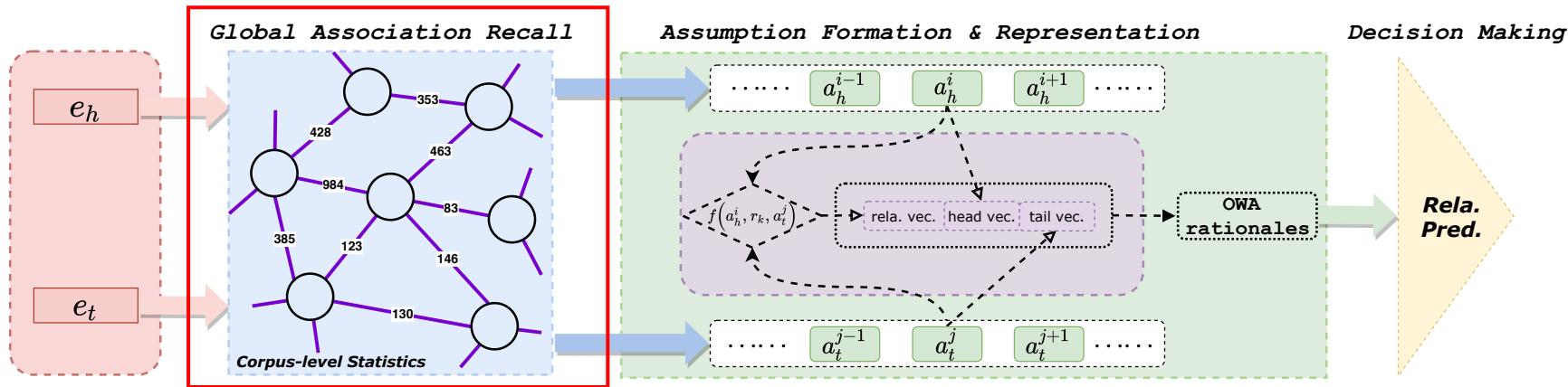


Corpus-level Statistics

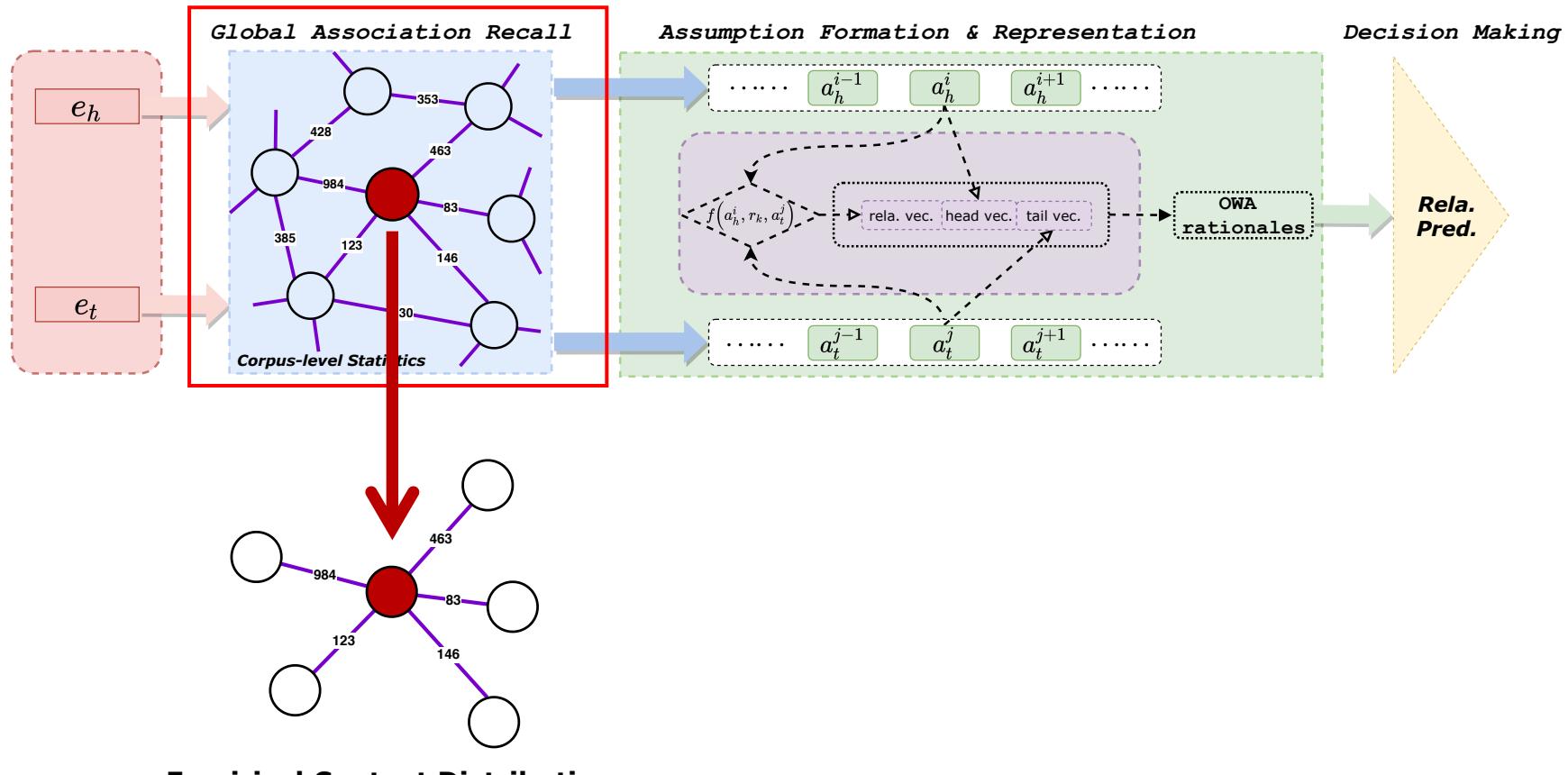
# Framework Overview



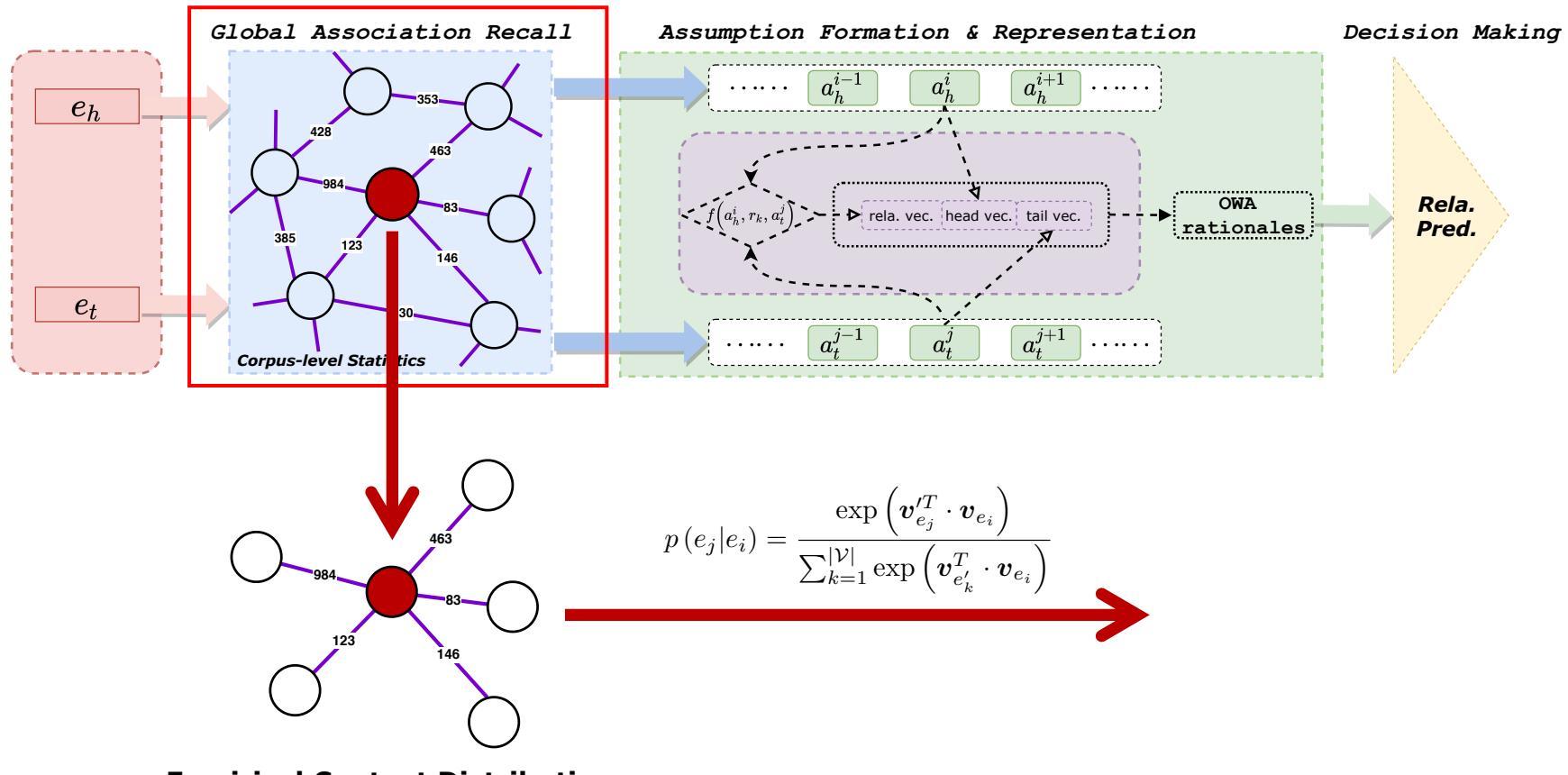
# CogStage-1: Global Association Recall



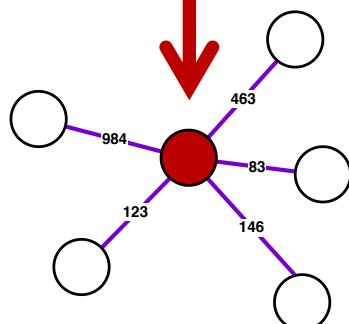
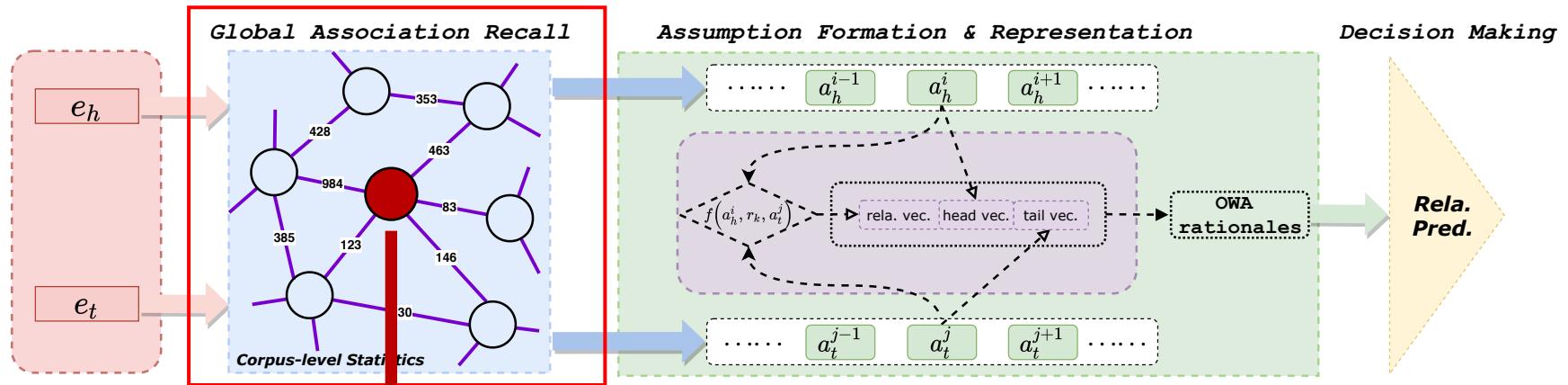
# CogStage-1: Global Association Recall



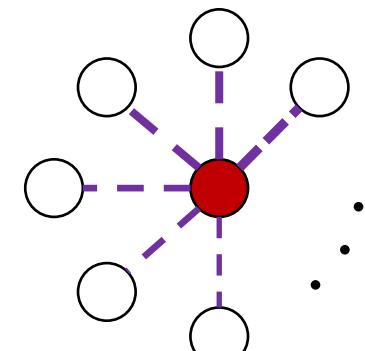
# CogStage-1: Global Association Recall



# CogStage-1: Global Association Recall



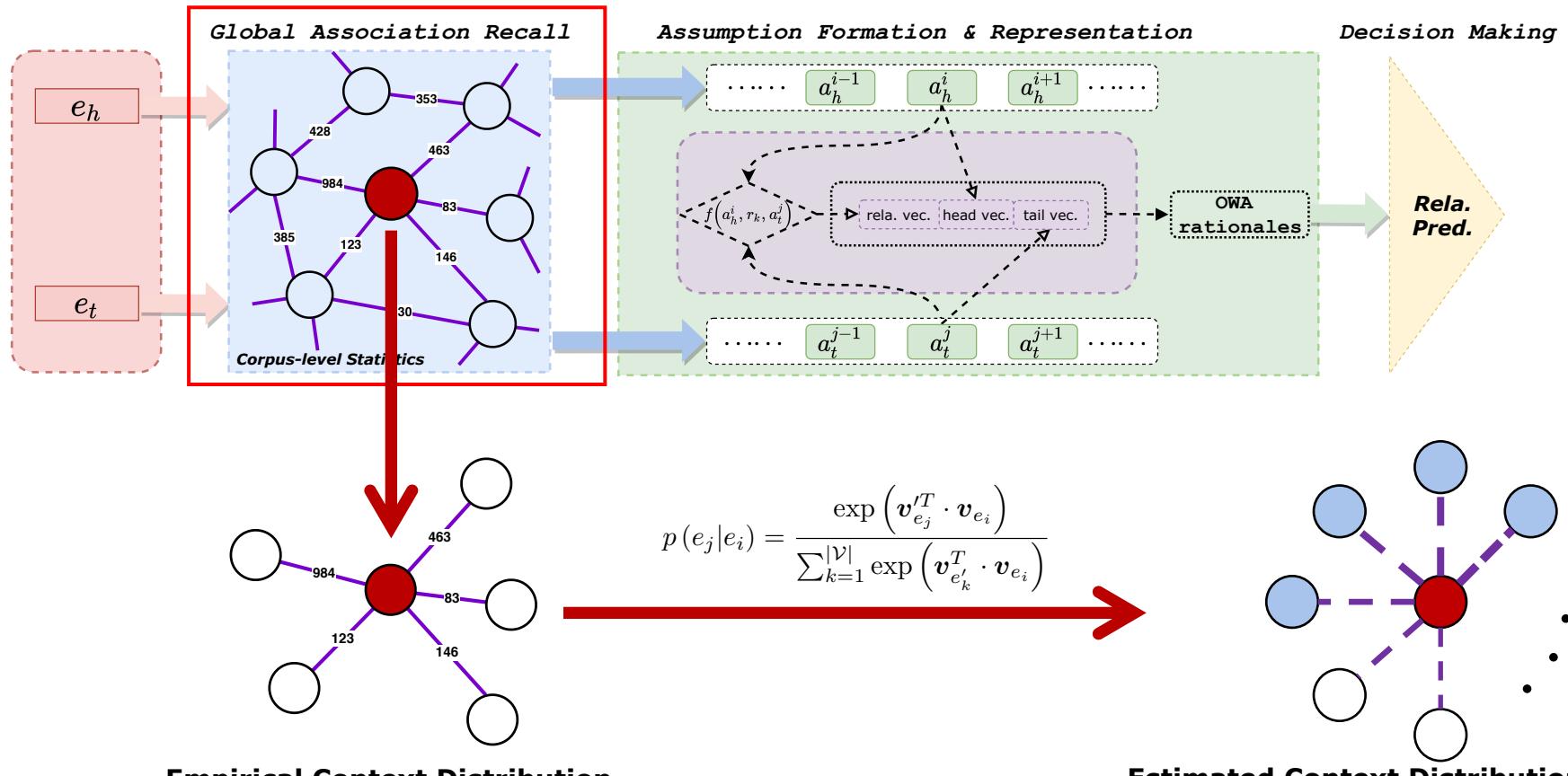
$$p(e_j|e_i) = \frac{\exp(\mathbf{v}_{e_j}^T \cdot \mathbf{v}_{e_i})}{\sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{v}_{e_k}^T \cdot \mathbf{v}_{e_i})}$$



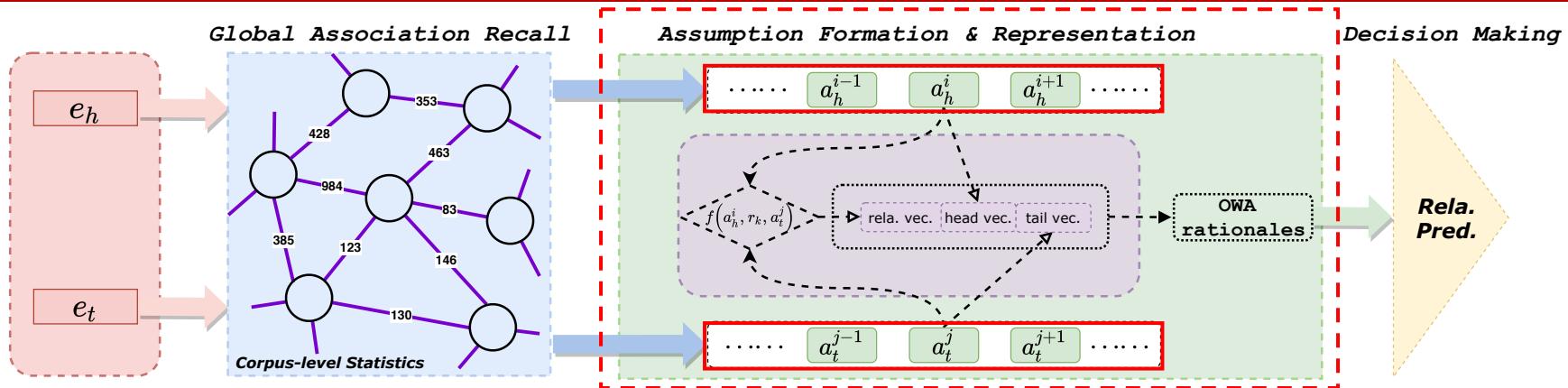
Empirical Context Distribution

Estimated Context Distribution

# CogStage-1: Global Association Recall

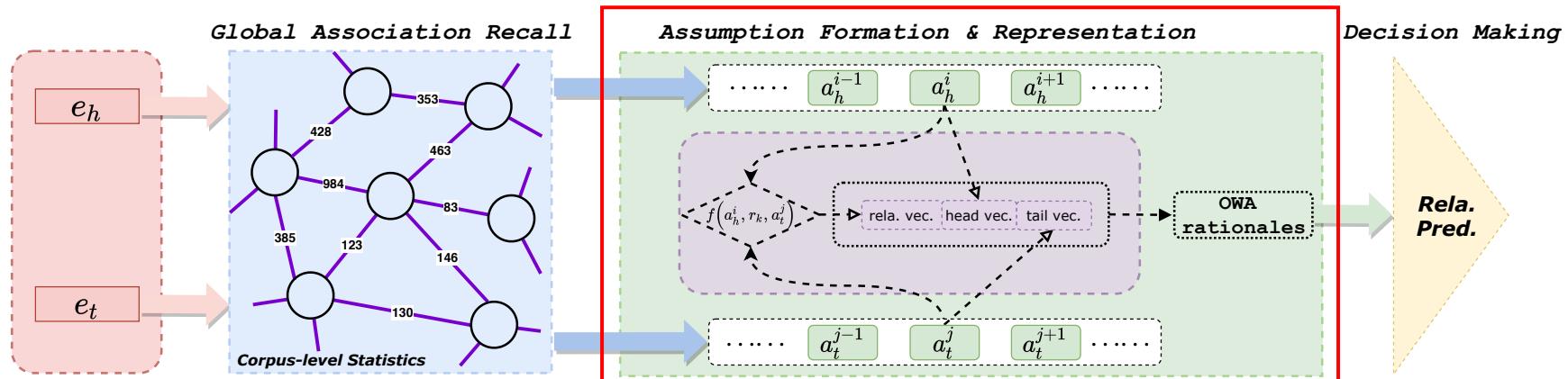


# CogStage-2: Assumption Formation and Representation



**Assumption:** Recognize whether context  $a_h^i$  and  $a_t^j$  hold a relationship  $r_k$ ?

# CogStage-2: Assumption Formation and Representation

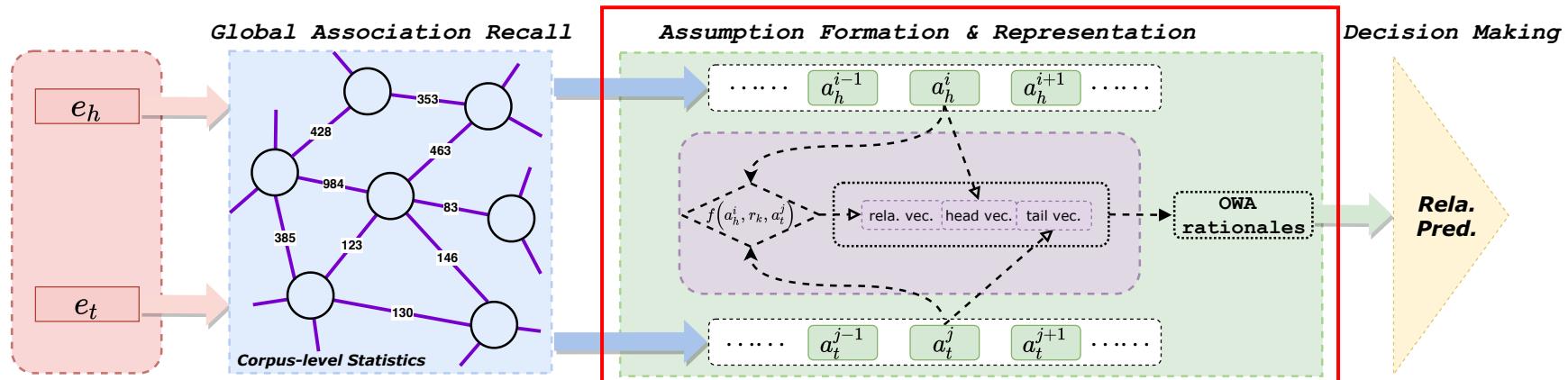


**Assumption:** Recognize whether context  $a_h^i$  and  $a_t^j$  hold a relationship  $r_k$ ?

**Closed World Assumption (CWA):**

**Open World Assumption (OWA):**

# CogStage-2: Assumption Formation and Representation



**Assumption:** Recognize whether context  $a_h^i$  and  $a_t^j$  hold a relationship  $r_k$ ?

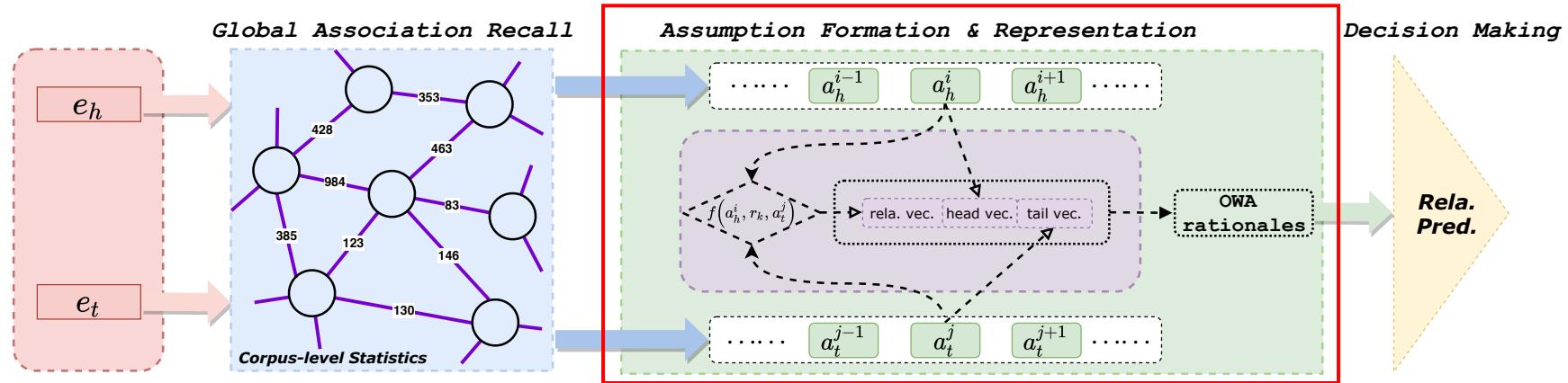
**Closed World Assumption (CWA):**

Only consider the facts that exist in KBs

**Open World Assumption (OWA):**

Consider all possible combinations and select the best ones.

# CogStage-2: Assumption Formation and Representation



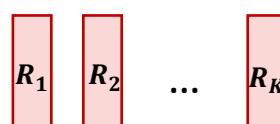
Head Vector



Tail Vector

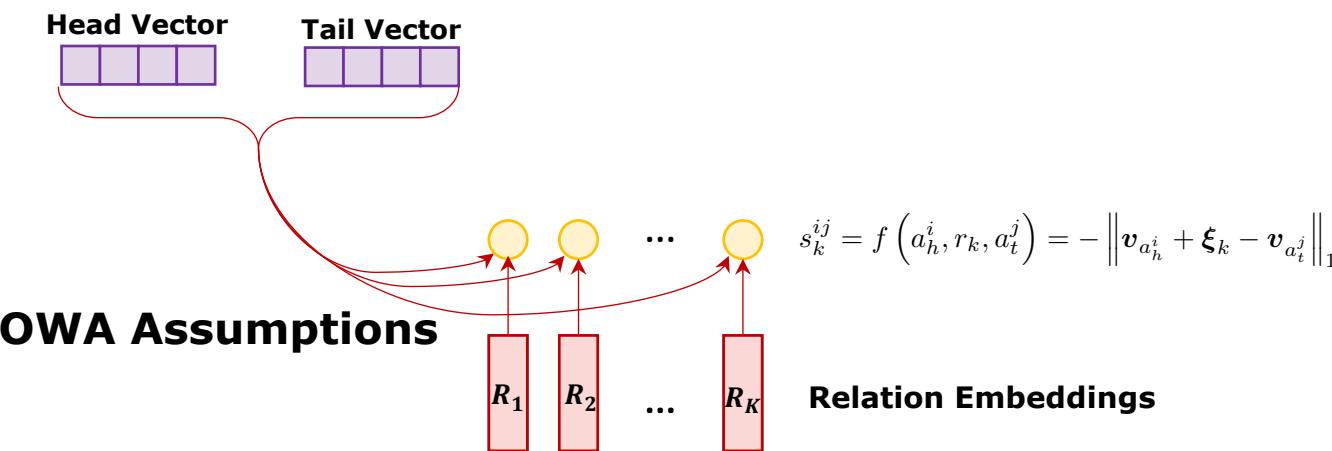
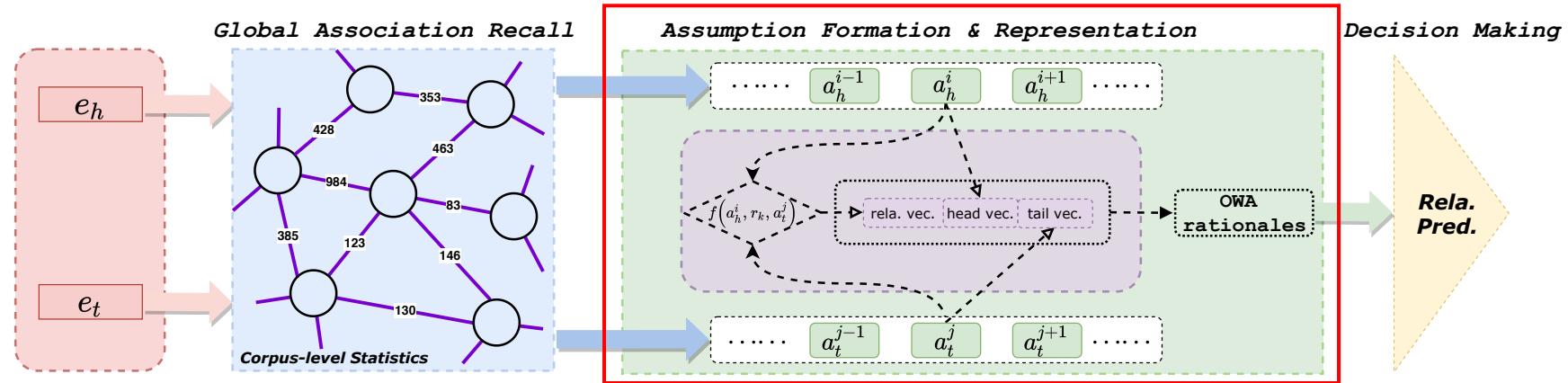


## OWA Assumptions

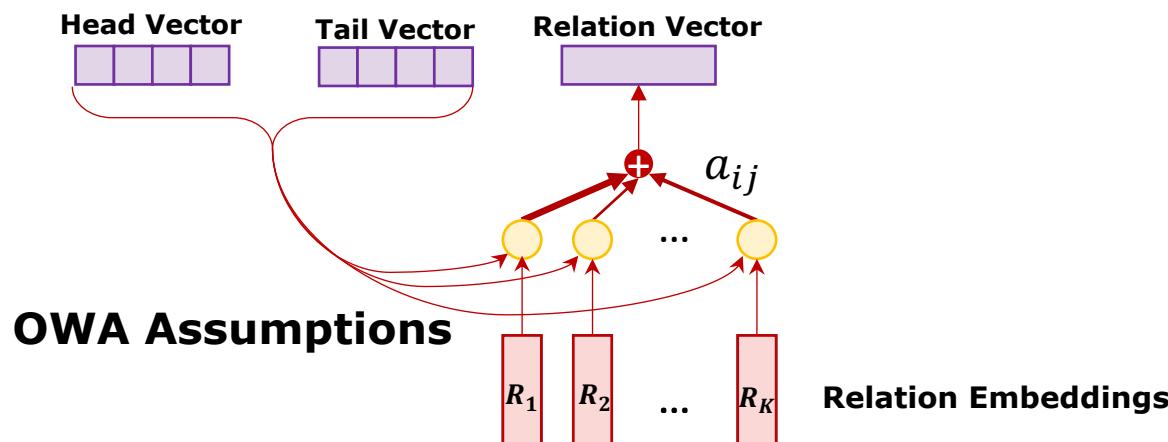
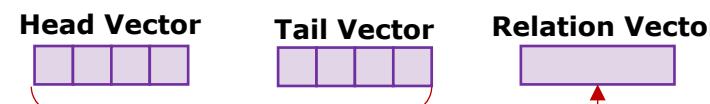
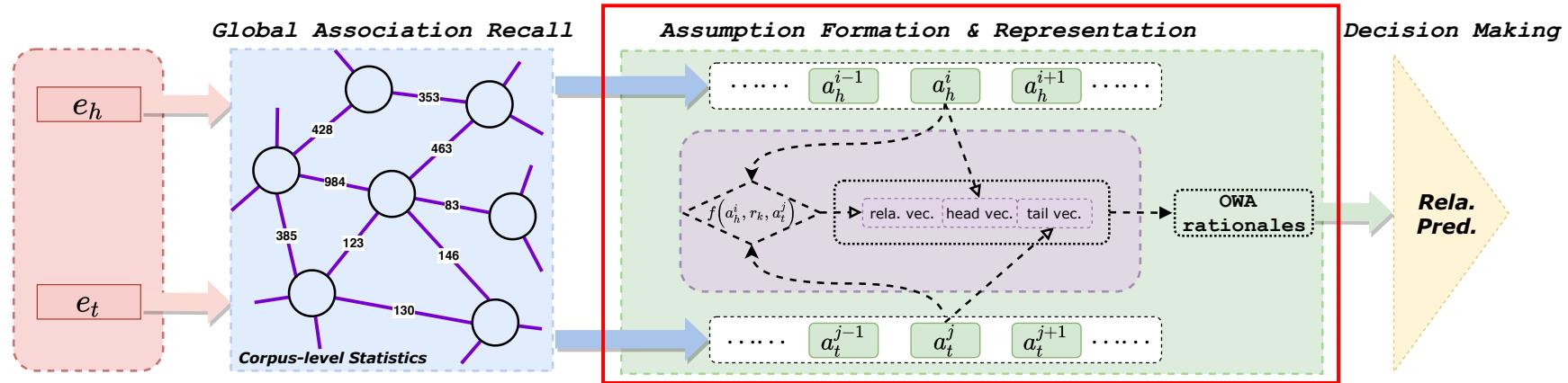


Relation Embeddings

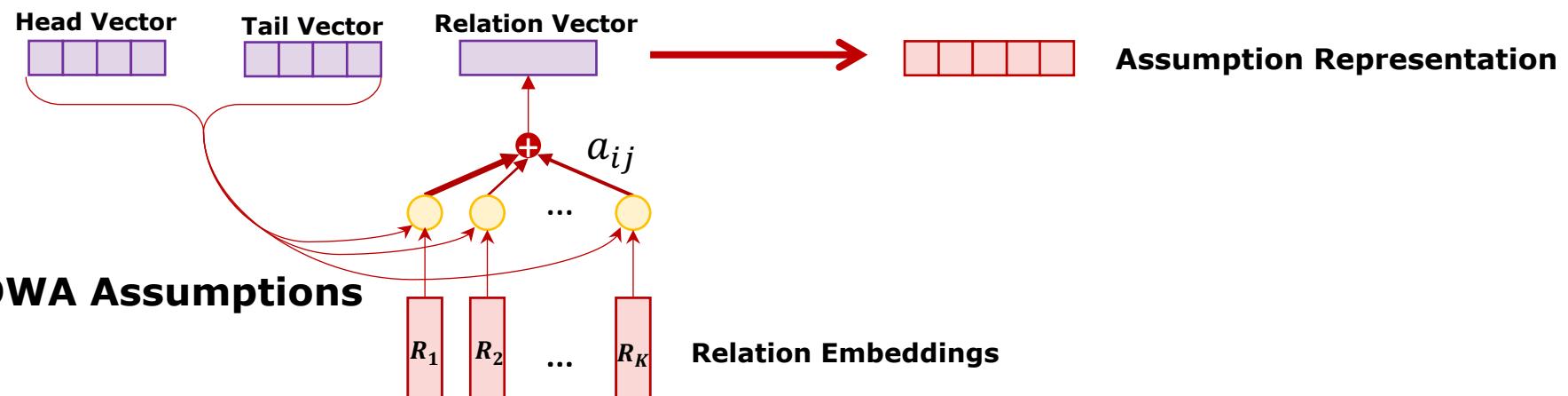
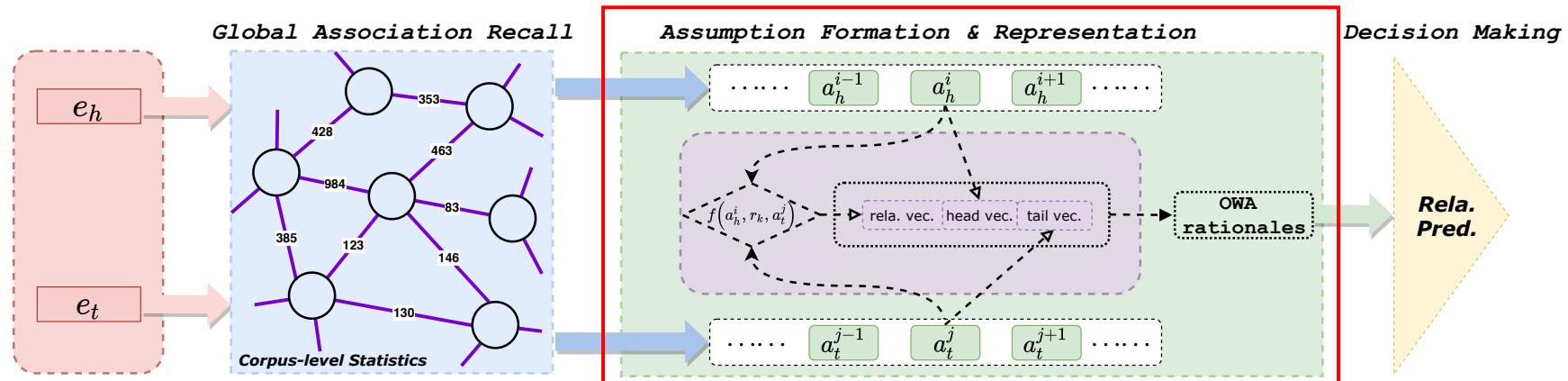
# CogStage-2: Assumption Formation and Representation



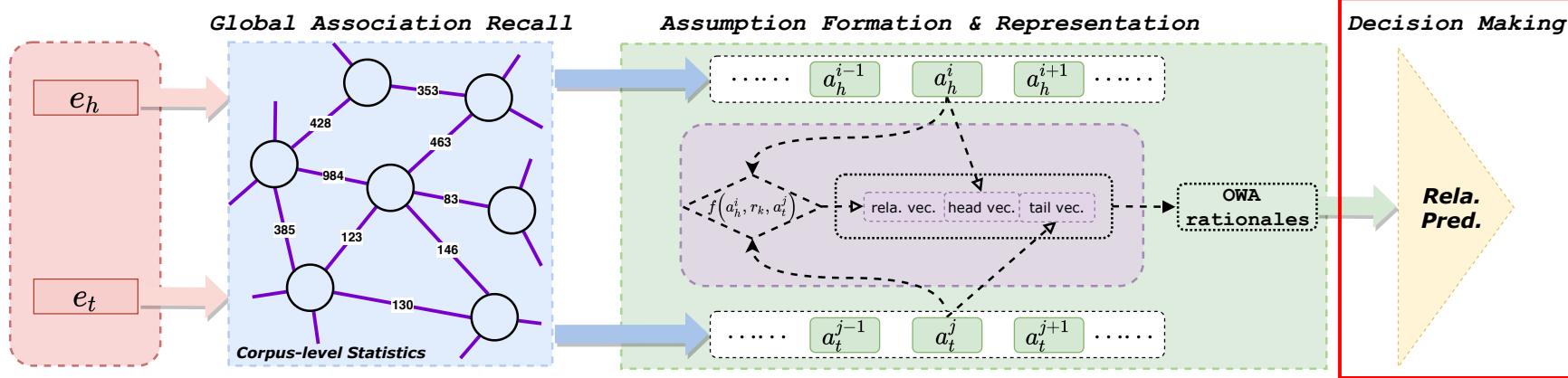
# CogStage-2: Assumption Formation and Representation



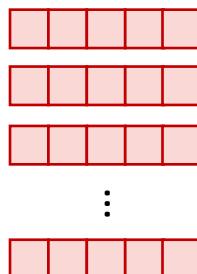
# CogStage-2: Assumption Formation and Representation



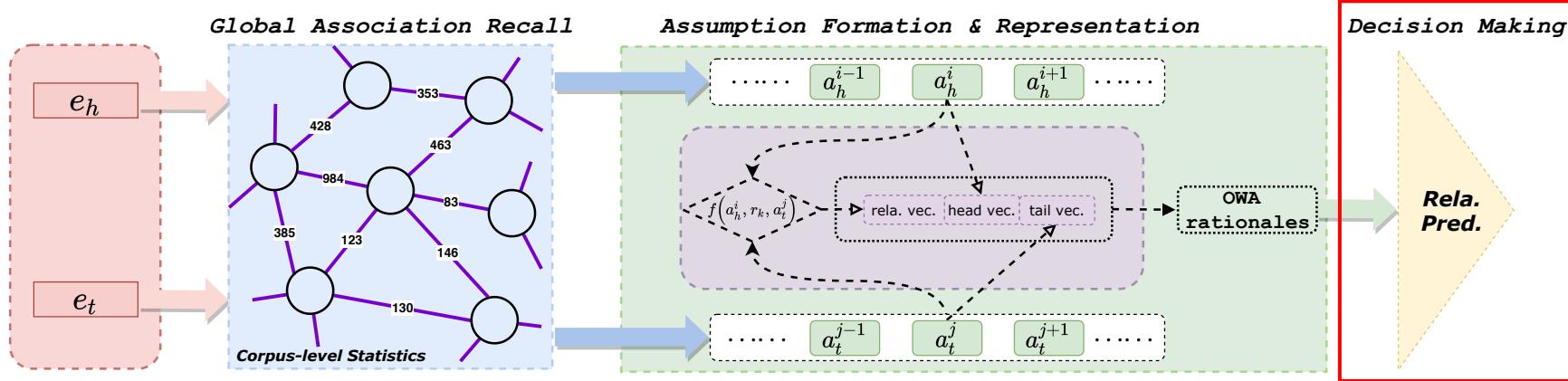
# CogStage-3: Prediction Decision Making



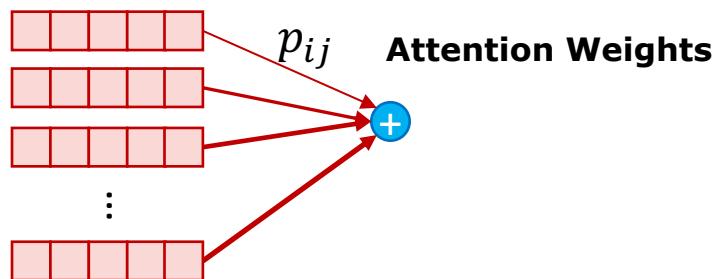
**Representations for all assumptions**



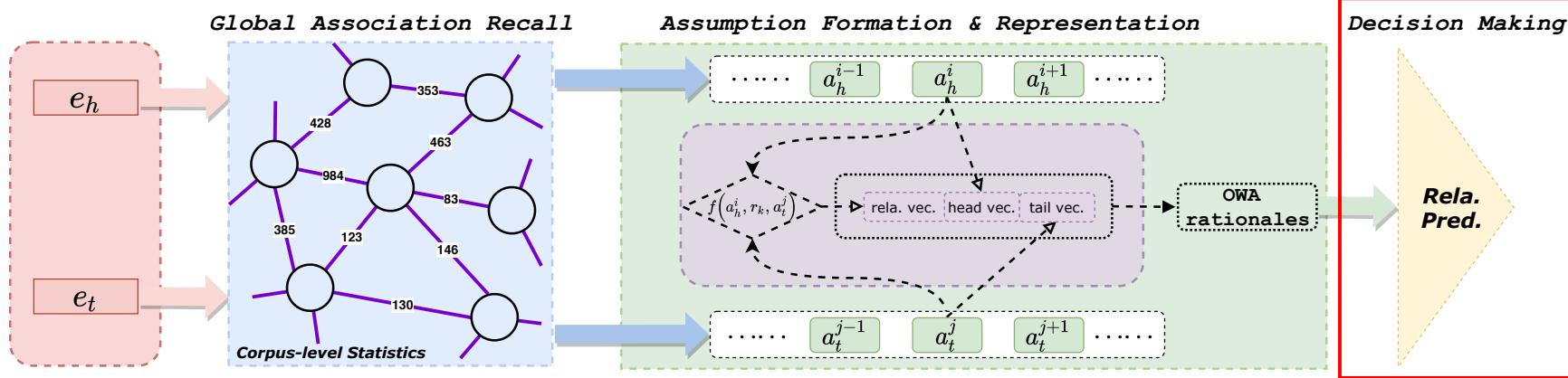
# CogStage-3: Prediction Decision Making



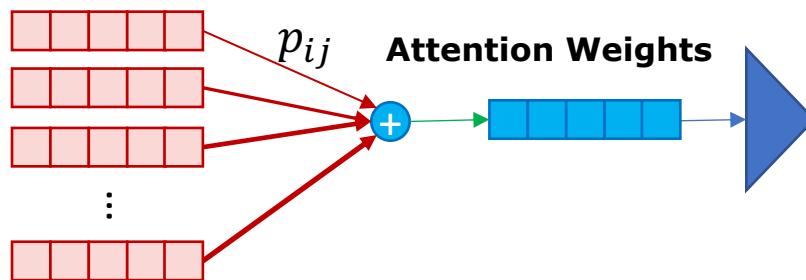
Representations for all assumptions



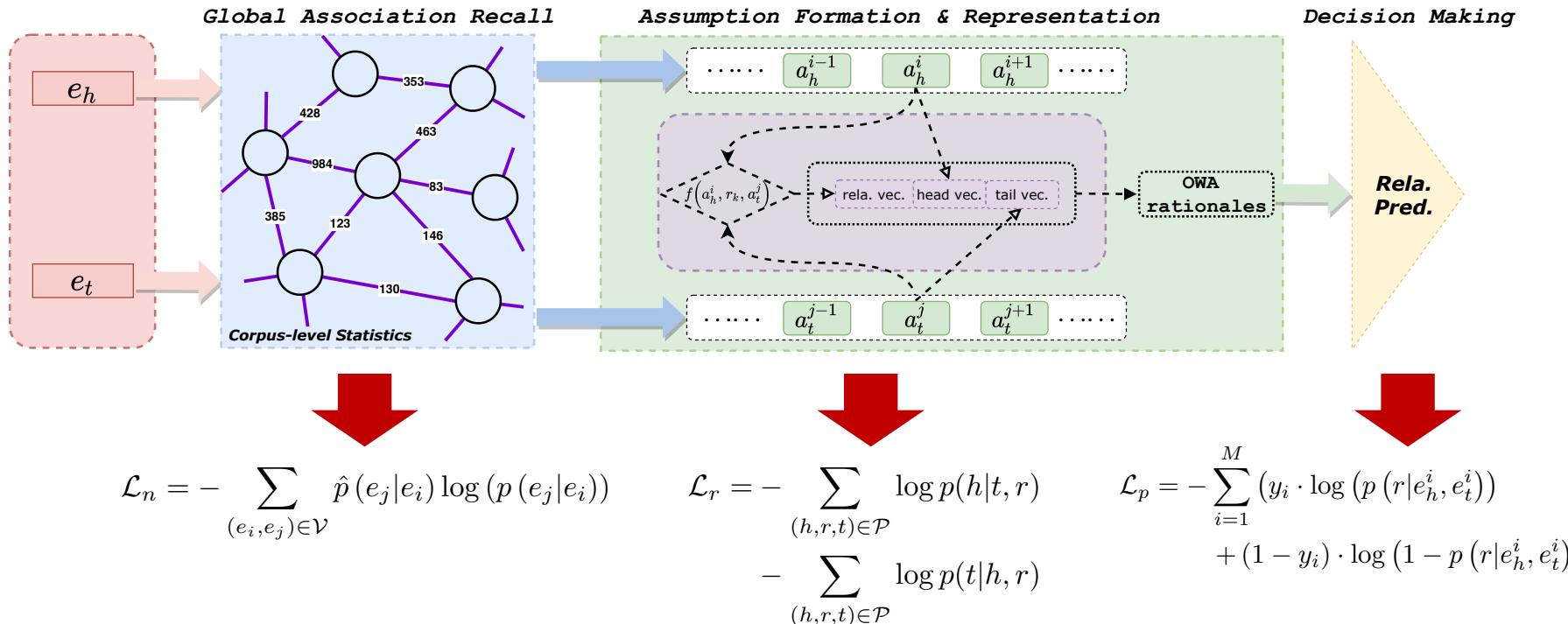
# CogStage-3: Prediction Decision Making



Representations for all assumptions



# Framework Training



# Experimental Setup

- Datasets
  - Medical Term-Term Co-occurrence Graph
  - 20 Million Clinical Notes from Stanford Hospital and Clinics since 1995
  - 52,804 Nodes, 16,197,319 Edges

# Experimental Setup

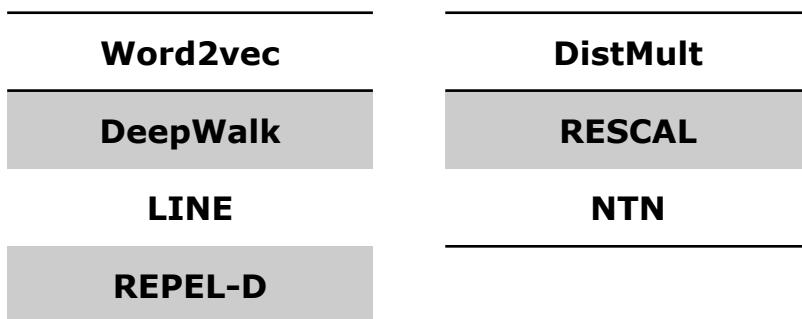
- Datasets
  - Medical Term-Term Co-occurrence Graph
  - 20 Million Clinical Notes from Stanford Hospital and Clinics since 1995
  - 52,804 Nodes, 16,197,319 Edges
  - Five Popular Medical Relations

Med Relations	Train	Dev	Test
Symptom of	14,326	3,001	3,087
May treat	12,924	2,664	2,735
Contraindicates	10,593	2,237	2,197
May prevent	2,113	440	460
Causes	1,389	305	354
Total	41.3k	8.6k	8.8k

Table 1: Dataset Statistics.

# Experimental Setup

- Datasets
  - Medical Term-Term Co-occurrence Graph
  - 20 Million Clinical Notes from Stanford Hospital and Clinics since 1995
  - 52,804 Nodes, 16,197,319 Edges
  - Five Popular Medical Relations
- Baseline Methods
  - Entity Encoder + Relation Scoring



Med Relations	Train	Dev	Test
Symptom of	14,326	3,001	3,087
May treat	12,924	2,664	2,735
Contraindicates	10,593	2,237	2,197
May prevent	2,113	440	460
Causes	1,389	305	354
Total	41.3k	8.6k	8.8k

Table 1: Dataset Statistics.

# Experiment Results: Predictive Performance

Methods	MAY_TREAT	CONTRAIN.	SYMPTOM_OF	MAY_PREVENT	CAUSES	Avg.
Word2vec + DistMult	0.767 ( $\pm 0.008$ )	0.777 ( $\pm 0.013$ )	0.815 ( $\pm 0.005$ )	0.649 ( $\pm 0.018$ )	0.671 ( $\pm 0.015$ )	0.736
Word2vec + RESCAL	0.743 ( $\pm 0.010$ )	0.767 ( $\pm 0.003$ )	0.808 ( $\pm 0.009$ )	0.658 ( $\pm 0.023$ )	0.659 ( $\pm 0.039$ )	0.727
Word2vec + NTN	0.693 ( $\pm 0.013$ )	0.758 ( $\pm 0.005$ )	0.808 ( $\pm 0.004$ )	0.605 ( $\pm 0.022$ )	0.631 ( $\pm 0.017$ )	0.699
DeepWalk + DistMult	0.740 ( $\pm 0.003$ )	0.776 ( $\pm 0.004$ )	0.805 ( $\pm 0.003$ )	0.608 ( $\pm 0.014$ )	0.650 ( $\pm 0.018$ )	0.716
DeepWalk + RESCAL	0.671 ( $\pm 0.010$ )	0.778 ( $\pm 0.003$ )	0.800 ( $\pm 0.003$ )	0.600 ( $\pm 0.023$ )	<b>0.708 (<math>\pm 0.011</math>)</b>	0.711
DeepWalk + NTN	0.696 ( $\pm 0.006$ )	0.778 ( $\pm 0.005$ )	0.787 ( $\pm 0.005$ )	0.614 ( $\pm 0.016$ )	0.674 ( $\pm 0.024$ )	0.710
LINE + DistMult	0.767 ( $\pm 0.003$ )	0.783 ( $\pm 0.002$ )	0.795 ( $\pm 0.003$ )	0.621 ( $\pm 0.015$ )	0.641 ( $\pm 0.024$ )	0.721
LINE + RESCAL	0.725 ( $\pm 0.003$ )	0.771 ( $\pm 0.002$ )	0.801 ( $\pm 0.001$ )	0.613 ( $\pm 0.013$ )	0.694 ( $\pm 0.015$ )	0.721
LINE + NTN	0.733 ( $\pm 0.002$ )	0.773 ( $\pm 0.003$ )	0.800 ( $\pm 0.001$ )	0.601 ( $\pm 0.015$ )	0.706 ( $\pm 0.013$ )	0.723
REPEL-D + DistMult	0.784 ( $\pm 0.002$ )	0.797 ( $\pm 0.002$ )	0.809 ( $\pm 0.003$ )	0.681 ( $\pm 0.010$ )	0.694 ( $\pm 0.022$ )	0.751
REPEL-D + RESCAL	0.726 ( $\pm 0.003$ )	0.780 ( $\pm 0.002$ )	0.776 ( $\pm 0.002$ )	<b>0.685 (<math>\pm 0.010</math>)</b>	<b>0.708 (<math>\pm 0.003</math>)</b>	0.737
REPEL-D + NTN	0.736 ( $\pm 0.004$ )	0.780 ( $\pm 0.002$ )	0.773 ( $\pm 0.001$ )	0.667 ( $\pm 0.015$ )	0.694 ( $\pm 0.024$ )	0.731
Ours (w/ CWA)	0.709 ( $\pm 0.005$ )	0.751 ( $\pm 0.009$ )	0.744 ( $\pm 0.007$ )	0.667 ( $\pm 0.008$ )	0.661 ( $\pm 0.032$ )	0.706
Ours	<b>0.805 (<math>\pm 0.017</math>)</b>	<b>0.811 (<math>\pm 0.006</math>)</b>	<b>0.816 (<math>\pm 0.004</math>)</b>	0.676 ( $\pm 0.020$ )	0.684 ( $\pm 0.017$ )	<b>0.758</b>

Table 2: Comparison of model predictive performance. We run all methods for five times and report the averaged F1 scores with standard deviations.

# Experiment Results: Predictive Performance

Methods	MAY_TREAT	CONTRAIN.	SYMPTOM_OF	MAY_PREVENT	CAUSES	Avg.
Word2vec + DistMult	0.767 ( $\pm 0.008$ )	0.777 ( $\pm 0.013$ )	0.815 ( $\pm 0.005$ )	0.649 ( $\pm 0.018$ )	0.671 ( $\pm 0.015$ )	0.736
Word2vec + RESCAL	0.743 ( $\pm 0.010$ )	0.767 ( $\pm 0.003$ )	0.808 ( $\pm 0.009$ )	0.658 ( $\pm 0.023$ )	0.659 ( $\pm 0.039$ )	0.727
Word2vec + NTN	0.693 ( $\pm 0.013$ )	0.758 ( $\pm 0.005$ )	0.808 ( $\pm 0.004$ )	0.605 ( $\pm 0.022$ )	0.631 ( $\pm 0.017$ )	0.699
DeepWalk + DistMult	0.740 ( $\pm 0.003$ )	0.776 ( $\pm 0.004$ )	0.805 ( $\pm 0.003$ )	0.608 ( $\pm 0.014$ )	0.650 ( $\pm 0.018$ )	0.716
DeepWalk + RESCAL	0.671 ( $\pm 0.010$ )	0.778 ( $\pm 0.002$ )	0.800 ( $\pm 0.002$ )	0.600 ( $\pm 0.022$ )	0.708 ( $\pm 0.011$ )	0.711
DeepWalk + NTN	0.696 ( $\pm 0.006$ )	Competitive predictive performance compared with a comprehensive list of baselines				
LINE + DistMult	0.767 ( $\pm 0.003$ )					
LINE + RESCAL	0.725 ( $\pm 0.003$ )	0.771 ( $\pm 0.002$ )	0.801 ( $\pm 0.001$ )	0.613 ( $\pm 0.013$ )	0.694 ( $\pm 0.015$ )	0.721
LINE + NTN	0.733 ( $\pm 0.002$ )	0.773 ( $\pm 0.003$ )	0.800 ( $\pm 0.001$ )	0.601 ( $\pm 0.015$ )	0.706 ( $\pm 0.013$ )	0.723
REPEL-D + DistMult	0.784 ( $\pm 0.002$ )	0.797 ( $\pm 0.002$ )	0.809 ( $\pm 0.003$ )	0.681 ( $\pm 0.010$ )	0.694 ( $\pm 0.022$ )	0.751
REPEL-D + RESCAL	0.726 ( $\pm 0.003$ )	0.780 ( $\pm 0.002$ )	0.776 ( $\pm 0.002$ )	<b>0.685 (<math>\pm 0.010</math>)</b>	<b>0.708 (<math>\pm 0.003</math>)</b>	0.737
REPEL-D + NTN	0.736 ( $\pm 0.004$ )	0.780 ( $\pm 0.002$ )	0.773 ( $\pm 0.001$ )	0.667 ( $\pm 0.015$ )	0.694 ( $\pm 0.024$ )	0.731
Ours (w/ CWA)	0.709 ( $\pm 0.005$ )	0.751 ( $\pm 0.009$ )	0.744 ( $\pm 0.007$ )	0.667 ( $\pm 0.008$ )	0.661 ( $\pm 0.032$ )	0.706
Ours	<b>0.805 (<math>\pm 0.017</math>)</b>	<b>0.811 (<math>\pm 0.006</math>)</b>	<b>0.816 (<math>\pm 0.004</math>)</b>	0.676 ( $\pm 0.020$ )	0.684 ( $\pm 0.017$ )	<b>0.758</b>

Table 2: Comparison of model predictive performance. We run all methods for five times and report the averaged F1 scores with standard deviations.

# Experiment Results: Predictive Performance

Methods	MAY_TREAT	CONTRAIN.	SYMPTOM_OF	MAY_PREVENT	CAUSES	Avg.
Word2vec + DistMult	0.767 ( $\pm 0.008$ )	0.777 ( $\pm 0.013$ )	0.815 ( $\pm 0.005$ )	0.649 ( $\pm 0.018$ )	0.671 ( $\pm 0.015$ )	0.736
Word2vec + RESCAL	0.743 ( $\pm 0.010$ )	0.767 ( $\pm 0.003$ )	0.808 ( $\pm 0.009$ )	0.658 ( $\pm 0.023$ )	0.659 ( $\pm 0.039$ )	0.727
Word2vec + NTN	0.693 ( $\pm 0.013$ )	0.758 ( $\pm 0.005$ )	0.808 ( $\pm 0.004$ )	0.605 ( $\pm 0.022$ )	0.631 ( $\pm 0.017$ )	0.699
DeepWalk + DistMult	0.740 ( $\pm 0.003$ )	0.776 ( $\pm 0.004$ )	0.805 ( $\pm 0.003$ )	0.608 ( $\pm 0.014$ )	0.650 ( $\pm 0.018$ )	0.716
DeepWalk + RESCAL	0.671 ( $\pm 0.010$ )	0.778 ( $\pm 0.002$ )	0.800 ( $\pm 0.002$ )	0.600 ( $\pm 0.022$ )	0.708 ( $\pm 0.011$ )	0.711
DeepWalk + NTN	0.696 ( $\pm 0.006$ )	Best predictive performance when the training data is large				
LINE + DistMult	0.767 ( $\pm 0.003$ )	0.771 ( $\pm 0.002$ )	0.801 ( $\pm 0.001$ )	0.613 ( $\pm 0.013$ )	0.694 ( $\pm 0.015$ )	0.721
LINE + RESCAL	0.725 ( $\pm 0.003$ )	0.771 ( $\pm 0.002$ )	0.801 ( $\pm 0.001$ )	0.613 ( $\pm 0.013$ )	0.694 ( $\pm 0.015$ )	0.721
LINE + NTN	0.733 ( $\pm 0.002$ )	0.773 ( $\pm 0.003$ )	0.800 ( $\pm 0.001$ )	0.601 ( $\pm 0.015$ )	0.706 ( $\pm 0.013$ )	0.723
REPEL-D + DistMult	0.784 ( $\pm 0.002$ )	0.797 ( $\pm 0.002$ )	0.809 ( $\pm 0.003$ )	0.681 ( $\pm 0.010$ )	0.694 ( $\pm 0.022$ )	0.751
REPEL-D + RESCAL	0.726 ( $\pm 0.003$ )	0.780 ( $\pm 0.002$ )	0.776 ( $\pm 0.002$ )	<b>0.685 (<math>\pm 0.010</math>)</b>	<b>0.708 (<math>\pm 0.003</math>)</b>	0.737
REPEL-D + NTN	0.736 ( $\pm 0.004$ )	0.780 ( $\pm 0.002$ )	0.773 ( $\pm 0.001$ )	0.667 ( $\pm 0.015$ )	0.694 ( $\pm 0.024$ )	0.731
Ours (w/ CWA)	0.709 ( $\pm 0.005$ )	0.751 ( $\pm 0.009$ )	0.744 ( $\pm 0.007$ )	0.667 ( $\pm 0.008$ )	0.661 ( $\pm 0.032$ )	0.706
Ours	<b>0.805 (<math>\pm 0.017</math>)</b>	<b>0.811 (<math>\pm 0.006</math>)</b>	<b>0.816 (<math>\pm 0.004</math>)</b>	0.676 ( $\pm 0.020$ )	0.684 ( $\pm 0.017$ )	<b>0.758</b>

Table 2: Comparison of model predictive performance. We run all methods for five times and report the averaged F1 scores with standard deviations.

# Experiment Results: Human Evaluation

## Human Evaluation Interface

All models predict the `may_treat` relation between t1 term `unfractionated heparin` [`'unfractionated heparin [epc]', 'heparin'`] and t2 term `myocardial infarction (mi)` [`'myocardial infarction'`] with the following rationales.

1. Check each rationale and answer this question: To which degree is this rationale helpful for you to trust the prediction? (0: no helpful; 1: a little bit helpful; 2: helpful; 3: very helpful)

### Model A's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
metabolic alkalosis	may_prevent	myocardial infarction (mi)	
metabolic alkalosis	may_prevent	venous thrombosis	
rbbb	may_treat	myocardial infarction (mi)	
ards	symptom_of	myocardial infarction (mi)	
micronutrient	may_prevent	venous thrombosis	

### Model B's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
cardiac dysrhythmias	contraindicates	theophylline	
malignant neoplasm without specification of site	has_symptom	family history of cancer	
liddm	contraindicates	glyburide	
morphine sulfate	contraindicated_by	respiratory depression	
insulin dependent diabetes	contraindicates	glyburide	

2. Please also rank all sets of rationales based on overall how much they help you trust the model prediction (e.g., A > B). Note that it is ok to reject them if both models are unhelpful (A = B = 0).

# Experiment Results: Human Evaluation

## Human Evaluation Interface

All models predict the `may_treat` relation between t1 term `unfractionated heparin` [`'unfractionated heparin [epc]', 'heparin'`] and t2 term `myocardial infarction (mi)` [`'myocardial infarction'`] with the following rationales.

1. Check each rationale and answer this question: To which degree is this rationale helpful for you to trust the prediction? (0: no helpful; 1: a little bit helpful; 2: helpful; 3: very helpful)

### Model A's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
metabolic alkalosis	may_prevent	myocardial infarction (mi)	
metabolic alkalosis	may_prevent	venous thrombosis	
rbbb	may_treat	myocardial infarction (mi)	
ards	symptom_of	myocardial infarction (mi)	
micronutrient	may_prevent	venous thrombosis	

### Model B's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
cardiac dysrhythmias	contraindicates	theophylline	
malignant neoplasm without specification of site	has_symptom	family history of cancer	
liddm	contraindicates	glyburide	
morphine sulfate	contraindicated_by	respiratory depression	
insulin dependent diabetes	contraindicates	glyburide	

2. Please also rank all sets of rationales based on overall how much they help you trust the model prediction (e.g., A > B). Note that it is ok to reject them if both models are unhelpful (A = B = 0).

# Experiment Results: Human Evaluation

## Human Evaluation Interface

All models predict the `may_treat` relation between t1 term `unfractionated heparin` [`'unfractionated heparin [epc]', 'heparin'`] and t2 term `myocardial infarction (mi)` [`'myocardial infarction'`] with the following rationales.

- Check each rationale and answer this question: To which degree is this rationale helpful for you to trust the prediction? (0: no helpful; 1: a little bit helpful; 2: helpful; 3: very helpful)

### Model A's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
metabolic alkalosis	may_prevent	myocardial infarction (mi)	
metabolic alkalosis	may_prevent	venous thrombosis	
rbbb	may_treat	myocardial infarction (mi)	
ards	symptom_of	myocardial infarction (mi)	
micronutrient	may_prevent	venous thrombosis	

### Model B's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
cardiac dysrhythmias	contraindicates	theophylline	
malignant neoplasm without specification of site	has_symptom	family history of cancer	
liddm	contraindicates	glyburide	
morphine sulfate	contraindicated_by	respiratory depression	
insulin dependent diabetes	contraindicates	glyburide	

- Please also rank all sets of rationales based on overall how much they help you trust the model prediction (e.g., A > B). Note that it is ok to reject them if both models are unhelpful (A = B = 0).

## Human Evaluation Score

	OWA Rationales	CWA Rationales
Ranking Score	17	5
Avg. Sum Score/Case	6.14	2.24
Avg. Max Score/Case	2.04	0.77

Table 3: Human evaluation on the quality of rationales.

# Experiment Results: Human Evaluation

## Human Evaluation Interface

All models predict the `may_treat` relation between t1 term `unfractionated heparin` ['unfractionated heparin [epc]', 'heparin'] and t2 term `myocardial infarction (mi)` ['myocardial infarction'] with the following rationales.

- Check each rationale and answer this question: To which degree is this rationale helpful for you to trust the prediction? (0: no helpful; 1: a little bit helpful; 2: helpful; 3: very helpful)

### Model A's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
metabolic alkalosis	may_prevent	myocardial infarction (mi)	
metabolic alkalosis	may_prevent	venous thrombosis	
rbbb	may_treat	myocardial infarction (mi)	
ards	symptom_of	myocardial infarction (mi)	
micronutrient	may_prevent	venous thrombosis	

### Model B's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
cardiac dysrhythmias	contraindicates	theophylline	
malignant neoplasm without specification of site	has_symptom	family history of cancer	
iddm	contraindicates	glyburide	
morphine sulfate	contraindicated_by	respiratory depression	
insulin dependent diabetes	contraindicates	glyburide	

- Please also rank all sets of rationales based on overall how much they help you trust the model prediction (e.g., A > B). Note that it is ok to reject them if both models are unhelpful (A = B = 0).

## Human Evaluation Score

	OWA Rationales	CWA Rationales
Ranking Score	17	5
Avg. Sum Score/Case	6.14	2.24
Avg. Max Score/Case	2.04	0.77

Table 3: Human evaluation on the quality of rationales.

## Case Study

Case 1		
<b>cephalosporins</b>	<code>may_treat</code>	bacterial infection
<b>cefuroxime</b>	<code>may_treat</code>	viral syndrome
<b>cefuroxime</b>	<code>may_treat</code>	low grade fever
<b>cefuroxime</b>	<code>may_treat</code>	<b>infectious diseases</b>
<b>cefuroxime</b>	<code>may_prevent</code>	low grade fever
<b>sulbactam</b>	<code>may_treat</code>	low grade fever

# Experiment Results: Human Evaluation

## Human Evaluation Interface

All models predict the `may_treat` relation between t1 term `unfractionated heparin` ['unfractionated heparin [epc]', 'heparin'] and t2 term `myocardial infarction (mi)` ['myocardial infarction'] with the following rationales.

1. Check each rationale and answer this question: To which degree is this rationale helpful for you to trust the prediction? (0: no helpful; 1: a little bit helpful; 2: helpful; 3: very helpful)

Model A's Rationale Set:

T1's contexts	Relational Interaction	T2's contexts	Score
metabolic alkalosis	may_prevent	myocardial infarction (mi)	
metabolic alkalosis	may_prevent	venous thrombosis	
rbbb	may_treat	myocardial infarction (mi)	
ards	symptom_of	myocardial infarction (mi)	
micronutrient	may_prevent	venous thrombosis	

The rationales can help justify the correct prediction.

metabolic alkalosis	specification_of_site	myocardial infarction (mi)
liddm	contraindicates	glyburide
morphine sulfate	contraindicated_by	respiratory depression
insulin dependent diabetes	contraindicates	glyburide

2. Please also rank all sets of rationales based on overall how much they help you trust the model prediction (e.g., A > B). Note that it is ok to reject them if both models are unhelpful (A = B = 0).

	OWA Rationales	CWA Rationales
Ranking Score	17	5
Avg. Sum Score/Case	6.14	2.24
Avg. Max Score/Case	2.04	0.77

Table 3: Human evaluation on the quality of rationales.

## Case Study

Case 1		
cephalosporins	may_treat	bacterial infection
cefuroxime	may_treat	viral syndrome
cefuroxime	may_treat	low grade fever
<b>cefuroxime</b>	<b>may_treat</b>	<b>infectious diseases</b>
cefuroxime	may_prevent	low grade fever
sulbactam	may_treat	low grade fever

See more details in the paper.

# Conclusions

- We propose an interpretable framework to rationalize medical relation prediction based on corpus-level statistics

# Conclusions

- We propose an interpretable framework to rationalize medical relation prediction based on corpus-level statistics
- Inspired by existing cognitive theories, the reasoning process can be easily understood by users and provides reasonable explanations to justify its prediction

# Conclusions

- We propose an interpretable framework to rationalize medical relation prediction based on corpus-level statistics
- Inspired by existing cognitive theories, the reasoning process can be easily understood by users and provides reasonable explanations to justify its prediction
- We demonstrate the model effectiveness by its predictive performance and human evaluation.



THE OHIO STATE UNIVERSITY

Zhen Wang  
The Ohio State University  
[wang.9125@osu.edu](mailto:wang.9125@osu.edu)



Code available at:  
<https://github.com/zhenwang9102/X-MedRELA>

# Thank You!

Supported By:

