36TH AAAI CONFERENCE ON ARTIFICIAL INTELLIGENCE A VIRTUAL CONFERENCE FEBRUARY 22 - MARCH 1, 2022

## LOREN: Logic-Regularized Reasoning for Interpretable Fact Verification

**Jiangjie Chen**<sup>1,2</sup>, Qiaoben Bao<sup>1</sup>, Changzhi Sun<sup>2</sup>, Xinbo Zhang<sup>2</sup>, Jiaze Chen<sup>2</sup>, Hao Zhou<sup>2</sup>, Yanghua Xiao<sup>1</sup>, Lei Li<sup>3</sup>



**AAAI-22** 







UC SANTA BARBARA



#### Did Donald Trump win the 2020 U.S. presidential election?

# **Fact Verification**

- Input: a claim + a KB
- Task:

#### - Evidence Extraction

 Evidence from a trustworthy KB

#### - Veracity Prediction

- Supported (SUP)
- Refuted (REF)
- Not Enough Information (NEI)

#### Claim

The Rodney King riots took place in the most populous county in the USA.

#### [wiki/Los Angeles Riots]

The 1992 Los Angeles riots, <u>also known</u> <u>as the Rodney King riots were a series of</u> <u>riots</u>, lootings, arsons, and civil disturbances that <u>occurred in Los Angeles</u> <u>County</u>, California in April and May 1992.

#### [wiki/Los Angeles\_County]

Los Angeles County, officially County the of Los Angeles, is the most populous county in the USA.

Verdict: Supported

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Arpit Mittal. FEVER: a large-scale dataset for Fact Extraction and VERification. NAACL, 2018.

## A General Pipeline for Solving This Task



Image credit to: Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, Ming Zhou, Jiahai Wang, Jian Yin. **Reasoning Over Semantic-Level Graph for Fact Checking.** ACL 2020.

## A General Pipeline for Solving This Task



Image credit to: Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, Ming Zhou, Jiahai Wang, Jian Yin. **Reasoning Over Semantic-Level Graph for Fact Checking.** ACL 2020.

# **Shielding from Misinformation**

- Misinformation detection on social media
   Especially with the success of PLMs.
- Factually accurate language generation
  - NLG with data accuracy  $\checkmark$
  - NLG with factual accuracy ?
    - An objective evaluation on factual accuracy of machine generated text.





# Interpretable Fact Verification

#### Goal of Reasoning

– Right answer for the right thinking

# Interpretable Fact Verification

- Goal of Reasoning
  - Right answer for the right thinking
- Interpretability "may be" the right thinking
  - Faithful: able to explain the prediction
  - Accurate: should be right per se
  - **Debuggable**: able to find out where goes wrong

# Interpretable Fact Verification

#### • Goal of Reasoning

- Right answer for the right thinking
- Interpretability "may be" the right thinking
  - Faithful: able to explain the prediction
  - Accurate: should be right per se
  - **Debuggable**: able to find out where goes wrong

#### • The Research Question:

- How can we do it without supervision?

# **Learning from Humans**

Claim: c Donald Trump won the 2020 election.

- We carefully examine each phrase in a claim one by one.
  - Did Donald Trump win the election in [2020]?
  - Did Donald Trump win the [U.S.] presidential election in 2020?

# **Learning from Humans**

Claim: c Donald Trump won the 2020 election.

- We carefully examine each phrase in a claim one by one.
  - Did Donald Trump win the election in [2020]?
  - Did Donald Trump win the [U.S.] presidential election in 2020?

- We aggregate the verification results of each phrase following aggregation logic, i.e. a claim is found
  - *Supported* iff all phrases found supported;
  - **Refuted** iff exists a phrase found refuted;
  - **NEI** iff not refuted and exists a phrase found unverifiable.

## **LOREN: Overview**

#### Symbolic AI plans, connectionist AI executes.

## **LOREN: Overview**



• **TL;DR**: **build local premises** from evidence to support phrase veracity prediction, regularized by logical rules.

## **Evidence Retrieval**



- Extract evidence sentences from Wikipedia following Liu et al. ACL 2020
  - Document retrieval
  - Sentence ranking
- Five relevant sentences of the entities in a claim.

Zhenghao Liu, Chenyan Xiong, Maosong Sun. Zhiyuan Liu. **Fine-grained Fact Verification with Kernel Graph Attention Network**. ACL 2020.

# **Claim Phrase Extraction**



- Extract *claim phrases* for fine-grained decomposition
  - e.g. noun phrase, adjective phrase, named entity, etc.
- **Approach**: Parse with heuristic rules via off-the-shelf NLP tools
  - e.g. constituency parsing, pos tagging, NER, etc.

# **Probing Question Generation**



- Goal: generate probing questions to answer from evidence.
  - Cloze question & interrogative questions
  - Prepare for the QA task

# **Answer Probing Questions**



- **Goal**: acquire corresponding *local premises* from evidence for each claim phrase.
  - Fine-tune a Seq2Seq MRC model (BART) on a manufactured dataset based on *support* samples.

# **Answer Probing Questions**



• **Goal**: acquire corresponding *local premises* from evidence for each claim phrase.

- Fine-tune a Seq2Seq MRC model (BART) on a manufactured dataset based on *support* samples.

## **Assemble Local Premises**



- **Goal**: acquire corresponding *local premises* from evidence for each claim phrase.
  - Fine-tune a Seq2Seq MRC model (BART) on a manufactured dataset based on *support* samples.
  - Fill masked-claims with answered phrases to construct local premises.

**Decompose** claim verification  $p_{\theta}(\mathbf{y} | \mathbf{x})$ into phrase verification  $p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x})$ 

$$p_{\theta}(\mathbf{y} \,|\, \mathbf{x}) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y} \,|\, \mathbf{z}, \mathbf{x}) p(\mathbf{z} \,|\, \mathbf{x})$$



**Decompose** claim verification  $p_{\theta}(\mathbf{y} | \mathbf{x})$ into phrase verification  $p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x})$ 

$$p_{\theta}(\mathbf{y} \mid x) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y} \mid \mathbf{z}, x) p(\mathbf{z} \mid x)$$

Phrase veracity as latent variables



**Decompose** claim verification  $p_{\theta}(\mathbf{y} | \mathbf{x})$ into phrase verification  $p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x})$ 

$$p_{\theta}(\mathbf{y} \,|\, x) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y} \,|\, \mathbf{z}, x) p(\mathbf{z} \,|\, x)$$

- Variational inference for solving the latent model
  - Evidence Lower BOund (ELBO)

$$\text{ELBO} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{y},x)} [\log p_{\theta}(y^* | \mathbf{z}, x))] - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{y}, x) \parallel p(\mathbf{z} | x))$$

**Decompose** claim verification  $p_{\theta}(\mathbf{y} | \mathbf{x})$ into phrase verification  $p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x})$ 

$$p_{\theta}(\mathbf{y} \mid x) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y} \mid \mathbf{z}, x) p(\mathbf{z} \mid x)$$

- Variational inference for solving the latent model
  - Evidence Lower BOund (ELBO)

$$ELBO = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{y},x)} [\log p_{\theta}(y^*|\mathbf{z},x))] - D_{KL} q_{\phi}(\mathbf{z}|\mathbf{y},x) \| p(\mathbf{z}|x)]$$
Variational posterior
distribution
Variational posterior
$$\sum_{z_{f}}^{Variational posterior} prior distribution$$

## **Regularize Latent Variables with Logic**

$$\mathcal{L}_{\text{var}}(\theta, \phi) = -\operatorname{ELBO} = -\operatorname{E}_{q_{\phi}(\mathbf{z}|\mathbf{y}, x)} [\log p_{\theta}(y^* | \mathbf{z}, x))] + D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{y}, x) \parallel p(\mathbf{z} | x))$$
$$\mathcal{L}_{\text{final}}(\theta, \phi) = (1 - \lambda) \mathcal{L}_{\text{var}}(\theta, \phi) + \lambda \mathcal{L}_{\text{logic}}(\theta, \phi)$$
$$\mathcal{L}_{\text{logic}}(\theta, \phi) = D_{\text{KL}} \left( p_{\theta}(\mathbf{y} | \mathbf{z}, x) \parallel q_{\phi}^{\text{T}}(\mathbf{y}_{z} | \mathbf{y}, x) \right)$$

## **Regularize Latent Variables with Logic**

 $\mathscr{L}_{\text{var}}(\theta, \phi) = -\operatorname{ELBO} = -\operatorname{E}_{q_{\phi}(\mathbf{z}|\mathbf{y}, x)} [\log p_{\theta}(y^* | \mathbf{z}, x))] + D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{y}, x) \parallel p(\mathbf{z} | x))$  $\mathscr{L}_{\text{final}}(\theta,\phi) = (1-\lambda)\mathscr{L}_{\text{var}}(\theta,\phi) + \lambda \mathscr{L}_{\text{logic}}(\theta,\phi)$  $\mathscr{L}_{\text{logic}}(\theta, \phi) = D_{\text{KL}} \left( p_{\theta}(\mathbf{y} | \mathbf{z}, x) \parallel q_{\phi}^{\text{T}}(\mathbf{y}_{z} | \mathbf{y}, x) \right)$ • **Supported** iff all phrases  $q_{\phi}^{\mathrm{T}}(oldsymbol{y}_{z}= extsf{SUP})=\prod_{i=1}^{|oldsymbol{z}|}q_{\phi}(oldsymbol{z}_{i}= extsf{SUP})$ found supported; Refuted iff exists a phrase  $q_{\phi}^{\mathrm{T}}(\boldsymbol{y}_{z} = \mathtt{REF}) = 1 - \prod_{i=1}^{|\boldsymbol{z}|} (1 - q_{\phi}(\boldsymbol{z}_{i} = \mathtt{REF}))$ found refuted;  $q_{\phi}^{\mathrm{T}}(oldsymbol{y}_{z} = \texttt{NEI}) = 1 - q_{\phi}^{\mathrm{T}}(oldsymbol{y}_{z} = \texttt{SUP}) - q_{\phi}^{\mathrm{T}}(oldsymbol{y}_{z} = \texttt{REF})$ • **NEI** iff not refuted and exists a phrase found unverifiable. Soft logic 27 Hard logic





 $p_{\theta}(\mathbf{y} \mid x, \mathbf{z})$ 









#### **Iterative Decoding**

**1**. 
$$p(\mathbf{z})$$
  
**2**.  $p_{\theta}(\mathbf{y} \mid x, \mathbf{z})$   
**3**.  $q_{\phi}(\mathbf{z} \mid \mathbf{y}, x)$   
**4**. ...

# **Understanding LOREN**

- RQ1: Can we find rationales without hurting verification performance?
- RQ2: How faithful and accurate are these unsupervised rationales?
- RQ3: How do local premises contribute to LOREN and its rationales?

## **Research Questions**

- RQ1: Can we find rationales without hurting verification performance?
- RQ2: How faithful and accurate are these unsupervised rationales?
- RQ3: How do local premises contribute to LOREN and its rationales?

# **RQ1: Extrinsic Evaluation**

Datacet	Model	Dev		Test	
Dataset			FEV	LA	I
• FEVER	UNC NLP GEAR (BERT, )	69.72 74.84	66.49 70.69	68.21 71.60	6
Metrics	DREAM (XLNet <sub>large</sub> )	79.16	-	<u>76.85</u>	7
• Label Accuracy ( <b>LA</b> )	$egin{array}{c} KGAT \ (BERT_{\mathrm{large}}) \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ $	77.91 78.29	75.86 76.11	73.61 74.07	7 7
- Classification accuracy	$ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \$	- 78.44	- 76.21	75.96 74.43	7 7
• FEVER Score ( <b>FEV</b> )	$(RoBERTa_{large})$	<u>81.14</u>	78.83	76.42	7
To the would option when the wight	LisT5 (T5 <sub>3B</sub> )	81.26	<u>77.75</u>	79.35	7

- Is the verification using the right

evidence sentence?

Table 2: Overall performance of verification results on the dev and blind test set of FEVER task, where FEV (FEVER score) is the main evaluation metric. The best is **bolded**, and the second best is underlined.

**FEV** 64.21 67.10

70.60 70.24 70.38 72.30 70.71 72.93 **75.87** 

# **RQ1: Extrinsic Evaluation**

Dataset	Model		Dev		Test	
			FEV	LA	FEV	
• FEVER	UNC NLP	69.72	66.49	68.21	64.21	
	$GEAR (BERT_{base})$	74.84	70.69	71.60	67.10	
Metrics	$DREAM~(XLNet_{\rm large})$	79.16	-	76.85	70.60	
Field inco	$KGAT (BERT_{large})$	77.91	75.86	73.61	70.24	
• Label Accuracy (LA)	$\ \ (RoBERTa_{large})$	78.29	76.11	74.07	70.38	
	$(CorefRoBERTa_1)$	-	-	75.96	72.30	
<ul> <li>Classification accuracy</li> </ul>	LOREN (BERT <sub>large</sub> )	78.44	76.21	74.43	70.71	
<ul> <li>FEVER Score (FEV)</li> </ul>	$\ \ (RoBERTa_{large})$	<u>81.14</u>	78.83	76.42	<u>72.93</u>	
<ul> <li>Is the verification using the right</li> </ul>						

Table 2: Overall performance of verification results on the dev and blind test set of FEVER task, where FEV (FEVER score) is the main evaluation metric. The best is **bolded**, and the second best is <u>underlined</u>.

#### Conclusions

- For similar-sized baselines with similar settings (DREAM, KGAT)
  - very competitive

evidence sentence?

# **RQ1: Extrinsic Evaluation**

#### Metrics

- Label Accuracy (LA)
  - Classification accuracy
- FEVER Score (FEV)
  - Is the verification using the right *evidence sentence*?

Model	D	ev	Test		
	LA	FEV	LA	FEV	
UNC NLP	69.72	66.49	68.21	64.21	
GEAR (BERT <sub>base</sub> )	74.84	70.69	71.60	67.10	
DREAM (XLNet <sub>large</sub> )	79.16	-	76.85	70.60	
$KGAT (BERT_{large})$	77.91	75.86	73.61	70.24	
$(RoBERTa_{large})$	78.29	76.11	74.07	70.38	
$(CorefRoBERTa_1)$	_	_	75.96	72.30	
LOREN (BERT <sub>large</sub> )	78.44	76.21	74.43	70.71	
$(RoBERTa_{large})$	<u>81.14</u>	78.83	76.42	<u>72.93</u>	
LisT5 (T5 <sub>3B</sub> )	81.26	77.75	79.35	75.87	

Table 2: Overall performance of verification results on the dev and blind test set of FEVER task, where FEV (FEVER score) is the main evaluation metric. The best is **bolded**, and the second best is <u>underlined</u>.

#### Conclusions

- For similar-sized baselines with similar settings (DREAM, KGAT)
  - very competitive
- For the 10x larger baseline (LisT5)



# **RQ1: Intrinsic Evaluation**



#### Conclusion

• Finding rationales does not hurt verification performance.

## **Research Questions**

- RQ1: Can we find rationales without hurting verification performance?
- RQ2: How faithful and accurate are these unsupervised rationales?
- RQ3: How do local premises contribute to LOREN and its rationales?

#### Goals of interpretability

Accurate

**⊚**\*Faithful

Oebuggable

#### Goals of interpretability

#### Metrics for evaluating rationales

Accurate

🎯 Faithful



#### Goals of interpretability

#### Metrics for evaluating rationales

Accurate
 Faithful

Debuggable

- Logically aggregated Label Accuracy of y<sub>z</sub> (LA<sub>z</sub>)
  - Evaluates the overall quality of z
- Culprit finding accuracy (CulpA) (P/R/F1)
  - Evaluates the individual quality of *z*:
  - Are the culprit phrase(s) found by rationales (z)?
  - Human evaluation: labeling culprit phrase(s) from claim phrases

#### Goals of interpretability



#### Metrics for evaluating rationales

- Logically aggregated Label Accuracy of y<sub>z</sub> (LA<sub>z</sub>)
  - Evaluates the *overall* quality of z
- Culprit finding accuracy (CulpA) (P/R/F1)
  - Evaluates the *individual* quality of *z*:
  - Are the culprit phrase(s) found by rationales (z)? (human evaluation)
- Agreement of LA and LA<sub>z</sub> (AGREE)
  - How aggregated phrase veracity  $(y_z)$  agrees with claim veracity  $(y)_z$

## **R2: Faithfulness of Rationales**

$$\mathscr{L}_{\text{final}}(\theta,\phi) = (1-\lambda)\mathscr{L}_{\text{var}}(\theta,\phi) + \lambda\mathscr{L}_{\text{logic}}(\theta,\phi)$$



#### Conclusions

- Agreement > 96%: *z* are in general faithful.
- $\lambda$ , Agree : stronger regularization from  $\mathscr{L}_{logic}$ , deciding the faithfulness of z.
- Soft > Hard: probability distributions of z gives more information than discrete labels.

## **R2: Overall Accuracy of Rationales**

$$\mathcal{L}_{\text{final}}(\theta,\phi) = (1-\lambda)\mathcal{L}_{\text{var}}(\theta,\phi) + \lambda\mathcal{L}_{\text{logic}}(\theta,\phi)$$



#### Conclusions

- LA<sub>z</sub> is close to 50% when  $\lambda = 0$ : Logic is *critical* for interpretability.
- $\lambda_{1}$ ,  $LA_{z}$  but quickly plateaued: stronger regularization from  $\mathscr{L}_{logic}$  does not affect performance much.

 An interpretability shortcut in the logic: predicting all phrase veracity to be the same as claim veracity. e.g., 1.REF v REF v REF = REF
 2.REF v SUP v NEI = REF

#### • Potential risks:

- Be tricked by the deceptively high overall accuracy  $LA_z$
- Rationales being invalid, as no culprit is found.

 An interpretability shortcut in the logic: predicting all phrase veracity to be the same as claim veracity. e.g., 1.REF v REF v REF = REF
 2.REF v SUP v NEI = REF

#### • Potential risks:

- Be tricked by the deceptively high overall accuracy  $LA_z$
- Rationales being invalid, as no culprit is found.

This can be revealed by altering the prior distribution  $p(\mathbf{z} | x)$ .

negative ELBO:  $\mathscr{L}_{var}(\theta,\phi)$ 

 $-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{y},x)} \left[\log p_{\theta}(y^*|\mathbf{z},x)\right] + D_{\mathrm{KL}}(q_{\phi}(\mathbf{z}|\mathbf{y},x) \parallel p(\mathbf{z}|x))$ 



Table 4: Results of different choices of prior distribution p(z) during training, where  $y_z$  in LA<sub>z</sub> is calculated using *soft* logic.

• A Pre-trained NLI from MNLI

• LA<sub>z</sub> at 53.41% for the NLI model before training

	Choice of $p(z)$	LA	$\mathbf{LA}_{z}$	AGREE	CULPA (P/R/F1)
Sample a few	NLI prior	81.14	79.66	96.11	75.8/75.9/74.3
phrases to be 🔶	Pseudo prior	80.93	80.44	97.25	70.5/77.1/71.4
culprits as prior	Uniform prior	80.85	80.74	97.08	34.1/ <b>78.8</b> /46.1

Table 4: Results of different choices of prior distribution p(z) during training, where  $y_z$  in LA<sub>z</sub> is calculated using *soft* logic.

<b>Choice of</b> $p(z)$	LA	$\mathbf{LA}_{z}$	AGREE	CULPA (P/R/F1)
NLI prior	<b>81.14</b>	79.66	96.11	<b>75.8</b> /75.9/ <b>74.3</b>
Pseudo prior	80.93	80.44	<b>97.25</b>	70.5/77.1/71.4
Uniform prior	80.85	<b>80.74</b>	97.08	34.1/ <b>78.8</b> /46.1

Table 4: Results of different choices of prior distribution p(z) during training, where  $y_z$  in LA<sub>z</sub> is calculated using *soft* logic.

#### Conclusions

- Prior distribution sets an important starting point for learning the rationales (z), but not on the overall predictions.
- NLI prior and pseudo prior can prevent the degeneration of phrase verification

## **Research Questions**

- RQ1: Can we find rationales without hurting verification performance?
- RQ2: How faithful and accurate are these unsupervised rationales?
- RQ3: How do local premises contribute to LOREN and its rationales?

## RQ3: Extrinsic Evaluation — MRC Performance

• Randomly sample 238 cases for Manual evaluation.



## RQ3: Extrinsic Evaluation — MRC Performance

• Randomly sample 238 cases for Manual evaluation.



#### Conclusions

- Self-supervised training for MRC is very beneficial for answering probing questions.
- Automatic factual error correction?

## RQ3: Intrinsic Evaluation — Simulating MRC deficiency

• What if MRC fails? — Masking local premises.



Figure 2: Performance on culprit finding (CULPA) and verification (LA and  $LA_z$ ) vs. the mask rate  $\rho$  of local premises, simulating the influence by deficiency of the MRC model.

#### Conclusions

- MRC is critical for the quality of individual rationales.
- Phrase verification degenerates to claim verification as MRC deteriorates.

# **Research Questions Revisited**

## • RQ1: Can we find rationales without hurting verification performance?

- Yes, even with a little boost for some cases.

#### RQ2: How faithful and accurate are these unsupervised rationales?

- Very faithful (96%+ agreement) and accurate (both in overall and individually).
- Logic regularizes the quality of phrase veracity.
- Careful for the "interpretability shortcut".

## • RQ3: How do local premises contribute to LOREN and its rationales?

- Minor contribution to claim veracity prediction.
- Critical to the quality of phrase veracity prediction.

Claim2: Ashley Cole is Iranian.

**Evidence**: Ashley Cole ( born 20 December 1980 ) is an English professional footballer who ... in Major League Soccer. Born in Stepney , London...

Claim2: Ashley Cole is Iranian.

**Evidence**: Ashley Cole ( born 20 December 1980 ) is an English professional footballer who ... in Major League Soccer. Born in Stepney , London...

Premise1: Ashley Cole is Iranian.

Premise2: Ashley Cole is European.

Claim2: Ashley Cole is Iranian. Evidence: Ashley Cole ( born 20 December 1980 ) is an English professional footballer who in Major League Soccer. Born in Stepney , London				
Premise1: Ashley Cole is Iranian. Veracity: SUPPORTS Premise2: Ashley Cole is European	z <sub>1</sub> = [ <b>0.981</b> , 0.004, 0.015]			
Veracity: REFUTES	z <sub>2</sub> = [0.014, <mark>0.520</mark> , 0.466]			
Prediction <i>y</i> : <b>REFUTES</b>	y <sub>z</sub> = [0.014, <mark>0.522</mark> , 0.464]			

<b>Claim2:</b> Ashley Cole is Iranian. Evidence: Ashley Cole ( born 20 December 1980 ) is an English professional footballer who in Major League Soccer. Born in Stepney , London				
Premise1: Ashley Cole is Iranian. Veracity: SUPPORTS Premise2: Ashley Cole is European. Veracity: REFUTES	z <sub>1</sub> = [ <b>0.981</b> , 0.004, 0.015] z <sub>2</sub> = [0.014, <b>0.520</b> , 0.466]			
Prediction <i>y</i> : <b>REFUTES</b> Ground Truth: NOT ENOUGH INFO	y <sub>z</sub> = [0.014, <mark>0.522</mark> , 0.464]			

- Wrong verification for Iranian.
- Rather close probabilities of **REF** and **NEI**.

- 
$$z_2$$
:  $p_{\text{REF}} = 0.520 \text{ vs.} p_{\text{REF}} = 0.466$ 

- 
$$y_z$$
:  $p_{\text{REF}} = 0.522$  vs.  $p_{\text{REF}} = 0.464$ 



- Goal of Reasoning
  - correct answer for the right thinking
- A good pipeline (LOREN) offers interpretability to the prediction
  - Faithful, accurate and debuggable
- A reasoning paradigm: symbolic AI plans, connectionist AI executes.

– Planning with logic, learning with data

# Have Fun with LOREN!



#### Checkout our code at **GitHub**! <u>https://github.com/jiangjiechen/LOREN</u>



#### Checkout our demo at 🤐 **Spaces**!

<u>https://huggingface.co/spaces/Jiangjie/</u> <u>loren-fact-checking</u>