



LOREN: Logic-Regularized Reasoning for Interpretable Fact Verification

Jiangjie Chen^{1,2}, Qiaoben Bao¹, Changzhi Sun², Xinbo Zhang², Jiaze Chen², Hao Zhou², Yanghua Xiao¹ and Lei Li³

¹Fudan University ²ByteDance AI Lab ³University of California, Santa Barbra

Introduction

❖ The Fact Verification Task (FEVER)

Verifying the veracity of a textual claim with evidence from trustworthy knowledge bases (e.g., Wikipedia).

- Misinformation detection on social media
- Factually accurate language generation



I won the Election!

Did Donald Trump **REFUTES** the 2020 U.S. presidential election?

But... Why?

❖ Interpretable Fact Verification

- 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 Problem:

- How can we do it without supervision?

Motivation 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

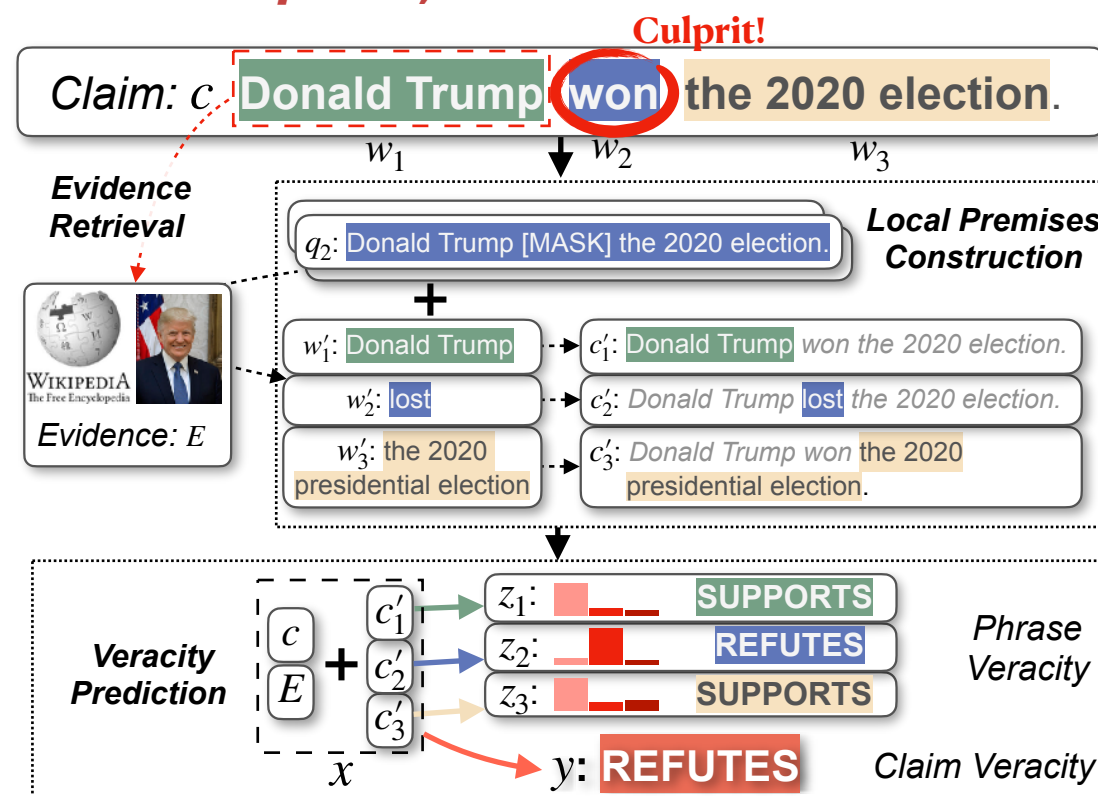
- **Supported** iff all phrases found supported;
- **Refuted** iff exists a phrase found refuted;
- **NEI** iff not refuted and exists a phrase found unverifiable.

Contribution

- ❖ We propose LOREN for interpretable fact verification.
- ❖ To solve **The Research Problem**, we decompose claim verification at phrase-level, and **build local premises** from evidence to support phrase veracity prediction, **regularized by logical rules**.

The LOREN Framework

Symbolic AI plans, connectionist AI executes.



❖ How to build local premises for verifying claim phrases? — The MRC Solution

1. Claim Phrase Extraction
2. Probing Question Generation
3. Machine Reading Comprehension

❖ How to train phrase verification without supervision? — The Latent Model

1. **Decompose** claim verification $p_\theta(y|x)$ into phrase verification $p_\theta(y|z, x)$

$$p_\theta(y|x) = \sum_z p_\theta(y|z, x)p(z|x)$$

2. **Variational Inference** to solve the latent model

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

❖ How to give latent variables the meaning of phrase veracity? — Regularize Latent Variables with Logic

$$\mathcal{L}_{\text{logic}}(\theta, \phi) = D_{\text{KL}}(p_\theta(y|z, x) \| q_\phi^T(y_z|y, x))$$

$$q_\phi^T(y_z = \text{SUP}) = \prod_{i=1}^{|z|} q_\phi(z_i = \text{SUP})$$

Soft logic

$$q_\phi^T(y_z = \text{REF}) = 1 - \prod_{i=1}^{|z|} (1 - q_\phi(z_i = \text{REF}))$$

$$q_\phi^T(y_z = \text{NEI}) = 1 - q_\phi^T(y_z = \text{SUP}) - q_\phi^T(y_z = \text{REF})$$

❖ Total Loss: -ELBO + Regularization

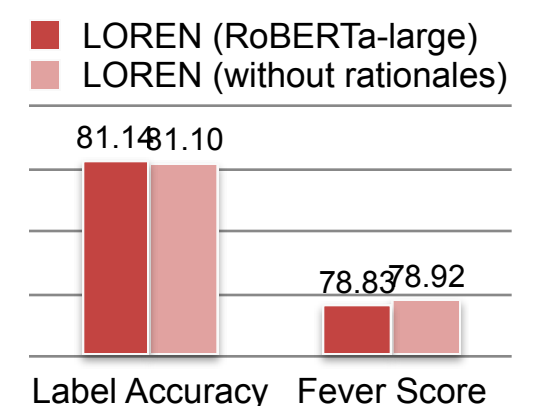
$$\mathcal{L}_{\text{final}}(\theta, \phi) = (1 - \lambda)\mathcal{L}_{\text{var}}(\theta, \phi) + \lambda\mathcal{L}_{\text{logic}}(\theta, \phi)$$

Experiments

❖ RQ1: Can we find rationales without hurting verification performance?

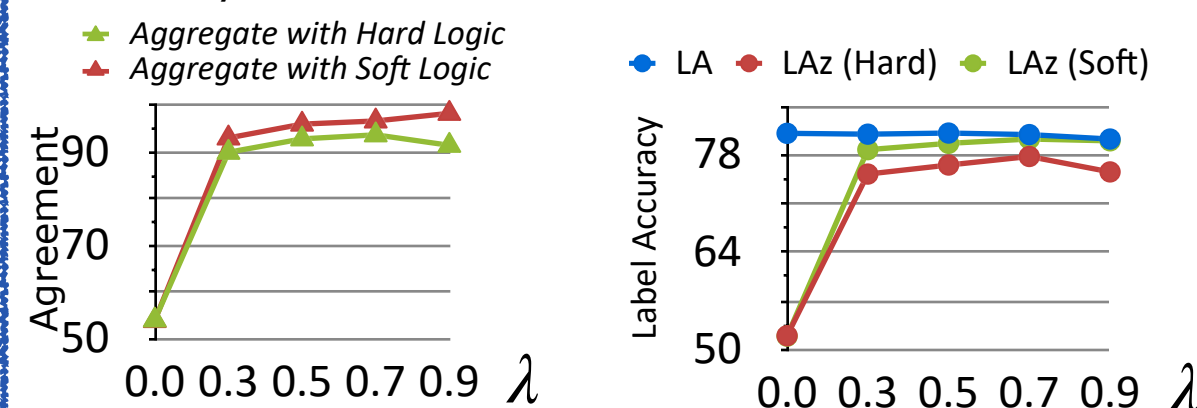
A1: Yes.

Model	Dev		Test	
	LA	FEV	LA	FEV
UNC NLP	69.72	66.49	68.21	64.21
GEAR (BERT _{base})	74.84	70.69	71.60	67.10
DREAM (XLNet _{large})	79.16	-	76.85	70.60
KGAT (BERT _{large})	77.91	75.86	73.61	70.24
L (RoBERTa _{large})	78.29	76.11	74.07	70.38
L (CorefRoBERTa _l)	-	-	75.96	72.30
LoREN (BERT _{large})	78.44	76.21	74.43	70.71
L (RoBERTa _{large})	81.14	78.83	76.42	72.93
List5 (T5 _{3B})	81.26	77.75	79.35	75.87



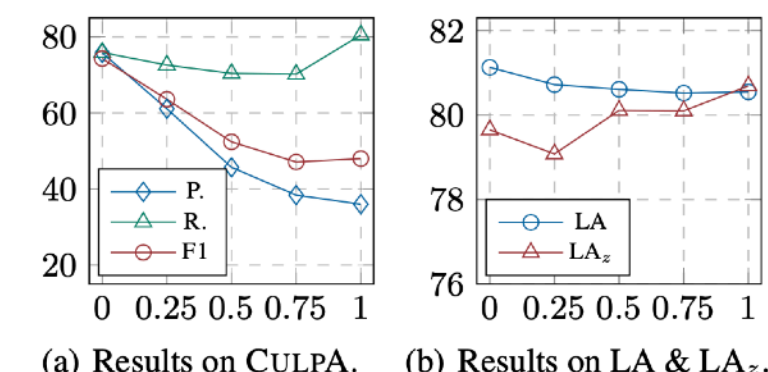
❖ RQ2: How faithful and accurate are these unsupervised rationales?

A2: Very faithful and accurate.



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

A3: **Critical** to the quality of phrase veracity prediction.



More details in the paper!

Debugging LOREN

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 $z_1 = [0.981, 0.004, 0.015]$
Premise2: Ashley Cole is European.	Veracity: REFUTES $z_2 = [0.014, 0.520, 0.466]$
Prediction y: REFUTES	Ground Truth: NOT ENOUGH INFO