# Interpretable Multimodal Out-of-context Detection with Soft Logic Regularization

Huanhuan Ma, Jinghao Zhang, Qiang Liu, Shu Wu, Liang Wang

# Background

- **Rapid Information Spread:** Mobile devices and media platforms facilitate the fast dissemination of news, increasing the exposure to false or deceptive content.
- **Misinformation Challenges:** Misinformation, particularly out-of-context news (where images or information are shared in misleading ways), poses serious societal risks.
- **Current Detection Limitations:**
  - Existing methods to identify misleading information often lack transparency.
  - Many current technologies offer limited explanations for their findings, complicating efforts to build trust and understanding.
- **Need for Improved Methods:**
  - There is a crucial need for methods that not only detect misinformation effectively but also provide clear, interpretable reasons for their assessments.
  - Enhancing interpretability can help in educating the public and aiding analysts in combating false information.

# The Task

- Image repurposing, also known as out-of-context photos are a powerful low-tech form of misinformation



Brazillian and Colombian boxers take apart a joint training session ✅

Brazillian and Colombian boxers take apart a joint training session ❌

# LOGic Regularization for out-of-context ANalysis (LOGRAN)



- **Caption Detection** Given a caption sentence c and its image I, our goal is to model the probability distribution $p(y|c, I)$, where $y \in$ {Pristine, Falsified} is a two-valued variable indicating the veracity of the caption's image.

- **Phrase Detection** We decompose the caption into phrases and predict the out-of-context label $z_i$ for each caption phrase $w_i \in W_c$ using the probability $p(z_i| c, w_i, I)$, where $z_i$ is treated as a binary latent variable $z_i \in$ {Pristine, Falsified}

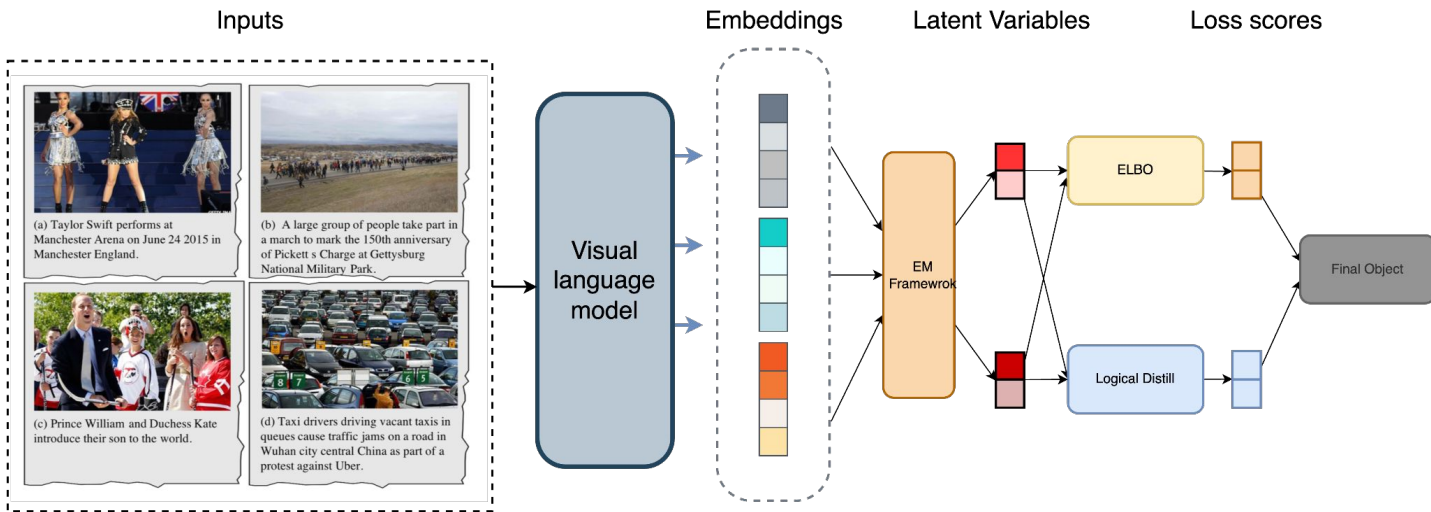**Follow the EM framework to model the latent variables**

$$p_t(\boldsymbol{y}|x) = \sum_{\boldsymbol{z}} p_t(\boldsymbol{y}|\boldsymbol{z}, x)p(\boldsymbol{z}|x)$$

# The Framework



Inputs — Embeddings — Latent Variables — Loss scores

(a) Taylor Swift performs at Manchester Arena on June 24 2015 in Manchester England.

(b) A large group of people take part in a march to mark the 150th anniversary of Pickett s Charge at Gettysburg National Military Park.

(c) Prince William and Duchess Kate introduce their son to the world.

(d) Taxi drivers driving vacant taxis in queues cause traffic jams on a road in Wuhan city central China as part of a protest against Uber.

Visual language model

EM Framewok

ELBO

Logical Distill

Final Object

**Follow the EM framework to model the latent variables**

$$p_t(\boldsymbol{y}|x) = \sum_{\boldsymbol{z}} p_t(\boldsymbol{y}|\boldsymbol{z}, x) p(\boldsymbol{z}|x)$$

**Weak supervise learning:**

- **ELBO**
- **Logical regularization**

# The Framework

- **ELBO loss:**

$$\mathcal{L}_{\text{var}}(t, l): \quad -\mathbb{E}_{q_l}[\log p_t(y^*|\boldsymbol{z}, x))] + D_{\text{KL}}(q_l(\boldsymbol{z}|\boldsymbol{y}, x) \parallel p(\boldsymbol{z}|x))$$

# The Framework



Inputs | Embeddings | Latent Variables | Loss scores

(a) Taylor Swift performs at Manchester Arena on June 24 2015 in Manchester England.

(b) A large group of people take part in a march to mark the 150th anniversary of Pickett s Charge at Gettysburg National Military Park.

(c) Prince William and Duchess Kate introduce their son to the world.

(d) Taxi drivers driving vacant taxis in queues cause traffic jams on a road in Wuhan city central China as part of a protest against Uber.

Visual language model

EM Framewok

ELBO

Logical Distill

Final Object

**The Logical Rule:**

A caption is considered: 1) Falsified if there is inconsistency in at least one caption phrase; 2) Pristine if all caption phrases are consistent

**Constructing Teacher module:**

Projecting the variational distribution

$$q_l(\boldsymbol{z}|\boldsymbol{y}, x) \implies q_l^{\mathrm{T}}(\boldsymbol{y}_z|\boldsymbol{y}, x)$$
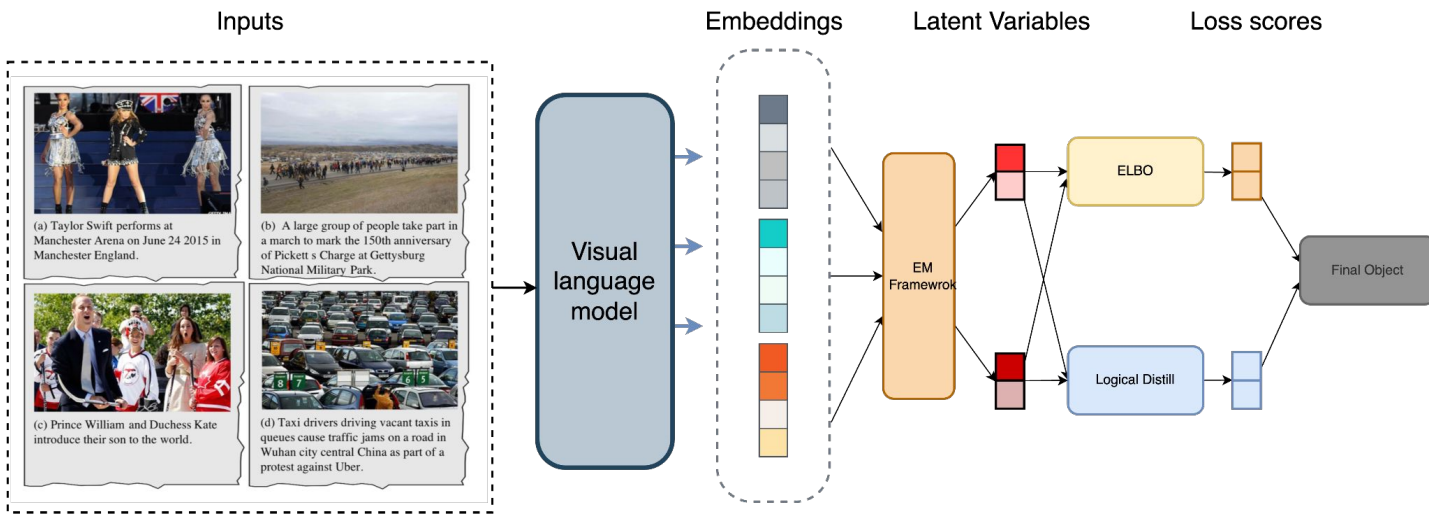
By simulating the outputs of teacher module, transfer logical knowledge into student model $p_t(\boldsymbol{y}|\boldsymbol{z}, x)$

**Logical distill loss:**

$$\mathcal{L}_{\mathrm{logic}}(t, l) = D_{\mathrm{KL}}\left(p_t(\boldsymbol{y}|\boldsymbol{z}, x) \parallel q_l^{\mathrm{T}}(\boldsymbol{y}_z|\boldsymbol{y}, x)\right)$$

# The Framework

**The final loss function:**

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

**Dataset：**

- **NewsCLIPpings** comprising both pristine and falsified images. It employs automation to match captions and images from the VisualNews corpus, offering various subsets based on matching methods.

**Backbone model:**

- **CLIP** utilizes distinct encoders for processing images and text, which are trained to produce comparable representations for associated concepts.
- **VisualBERT** is another multimodal model that integrates visual and textual information. It includes a sequence of Transformer layers that use self-attention to automatically align components of a given text input with specific regions in a corresponding image input.

# Results and Discussion

**Table 1**. Classification accuracy on the test set for the following models: (I) VisualBERT, (II) VisualBERT with LOGRAN, (III) Multimodal CLIP, and (IV) CLIP with LOGRAN. The underlined portions represent improvements from LOGRAN

|  | VisualBERT | VisualBERT-LOGRAN | CLIP | CLIP-LOGRAN |
|---|---|---|---|---|
| (a) Semantics/CLIP Text-Image | 55.12 | 56.88 | 58.59 | 59.03 |
| (b) Semantics/CLIP Text-Text | 53.47 | 55.62 | 68.36 | 70.81 |
| (c) Person/SBERT-WK Text-Text | 63.32 | 65.27 | 66.57 | 71.42 |
| (d) Scene/ResNet Place | 60.72 | 62.41 | 69.64 | 73.14 |
| Merged/Balanced | 61.32 | 63.18 | 67.27 | 70.51 |

**Improvement observed in both backbone models, as well as across all sub-test sets.**

# Results and Discussion

**Table 1.** Classification accuracy on the test set for the following models: (I) VisualBERT, (II) VisualBERT with LOGRAN, (III) Multimodal CLIP, and (IV) CLIP with LOGRAN. The underlined portions represent improvements from LOGRAN

|  | VisualBERT | VisualBERT-LOGRAN | CLIP | CLIP-LOGRAN |
|---|---|---|---|---|
| (a) Semantics/CLIP Text-Image | 55.12 | 56.88 | 58.59 | 59.03 |
| (b) Semantics/CLIP Text-Text | 53.47 | 55.62 | 68.36 | 70.81 |
| (c) Person/SBERT-WK Text-Text | 63.32 | 65.27 | 66.57 | 71.42 |
| (d) Scene/ResNet Place | 60.72 | 62.41 | 69.64 | 73.14 |
| Merged/Balanced | 61.32 | 63.18 | 67.27 | 70.51 |

**Improvement observed in both backbone models, as well as across all sub-test sets.**

- **VisualBERT vs VisualBERT-LOGRAN**
  The average improvement is 2%

# Results and Discussion

**Table 1.** Classification accuracy on the test set for the following models: (I) VisualBERT, (II) VisualBERT with LOGRAN, (III) Multimodal CLIP, and (IV) CLIP with LOGRAN. The underlined portions represent improvements from LOGRAN
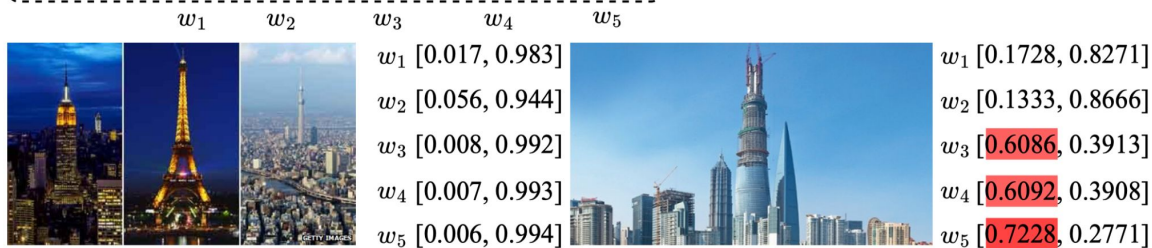
| | VisualBERT | VisualBERT-LOGRAN | CLIP | CLIP-LOGRAN |
|---|---|---|---|---|
| (a) Semantics/CLIP Text-Image | 55.12 | 56.88 | 58.59 | 59.03 |
| (b) Semantics/CLIP Text-Text | 53.47 | 55.62 | 68.36 | 70.81 |
| (c) Person/SBERT-WK Text-Text | 63.32 | 65.27 | 66.57 | 71.42 |
| (d) Scene/ResNet Place | 60.72 | 62.41 | 69.64 | 73.14 |
| Merged/Balanced | 61.32 | 63.18 | 67.27 | 70.51 |

**Improvement observed in both backbone models, as well as across all sub-test sets.**

- **VisualBERT vs VisualBERT-LOGRAN**
  The average improvement is 2%
- **CLIP vs CLIP-LOGRAN**
  The average improvement is 3%

# Case study

Caption: $C$ Fancy living in New York, Paris or Tokyo

$w_1$    $w_2$    $w_3$    $w_4$    $w_5$

$w_1$ [0.017, 0.983]

$w_2$ [0.056, 0.944]

$w_3$ [0.008, 0.992]

$w_4$ [0.007, 0.993]

$w_5$ [0.006, 0.994]

$w_1$ [0.1728, 0.8271]

$w_2$ [0.1333, 0.8666]

$w_3$ [0.6086, 0.3913]

$w_4$ [0.6092, 0.3908]

$w_5$ [0.7228, 0.2771]

Caption: $C$ Brazillian and Colombian boxers take apart a joint training session

$w_1$    $w_3$    $w_2$

$w_1$ [0.1476, 0.8524]

$w_2$ [0.0485, 0.9515]

$w_3$ [0.0846, 0.9154]

$w_1$ [0.8833, 0.1167]

$w_2$ [0.2185, 0.7815]

$w_3$ [0.3446, 0.6554]

We can easily identify the 'Culprit' in each case:

- New York Paris Tokyo
- Brazillian and Colombian boxers

which provides some level of interpretability

# Conclusion

- We proposed a novel frame work for out-of-context detection named **LOG**ic **R**egularization for out-of-context **AN**alysis (**LOGRAN**)
- Decomposes detection task from caption level to phrase level. Utilizes **latent variables** within an EM framework to predict out-of-context label for each phrase
- Implements two weak supervision methods: **ELBO loss** and **logical rule regularization**
- Conducted experiments on **NewsCLIPpings** dataset using **VisualBERT** and **CLIP** backbone models. Achieved **overall performance improvement**. Provides **phrase-level predictions** for **enhanced interpretability**

# Thank you!