

# Causal Intersectionality and Dual Form of Gradient Descent for Multimodal Analysis: a Case Study on Hateful Memes

Yosuke Miyanishi<sup>1</sup>, Minh Le Nguyen<sup>1</sup>

<sup>1</sup>Japan Advanced Institute of Science and Technology  
1-1 Asahi-dai, Nomi-shi, Ishikawa, Japan  
{yosuke.miyaniishi, nguyennml}@jaist.ac.jp

## Abstract

Amidst the rapid expansion of Machine Learning (ML) and Large Language Models (LLMs), understanding the semantics within their mechanisms is vital. Causal analyses define semantics, while gradient-based methods are essential to eXplainable AI (XAI), interpreting the model's 'black box'. Integrating these, we investigate how a model's mechanisms reveal its causal effect on evidence-based decision-making. Research indicates intersectionality - the combined impact of an individual's demographics - can be framed as an Average Treatment Effect (ATE). This paper demonstrates that hateful meme detection can be viewed as an ATE estimation using intersectionality principles, and summarized gradient-based attention scores highlight distinct behaviors of three Transformer models. We further reveal that LLM Llama-2 can discern the intersectional aspects of the detection through in-context learning and that the learning process could be explained via meta-gradient, a secondary form of gradient. In conclusion, this work furthers the dialogue on Causality and XAI. Our code is available online (see *External Resources* section).

## Introduction

The domain of causality offers profound insights into the data generation processes, revealing the intricate architecture of the problems at hand. A meticulous examination of these generative processes is indispensable for deep comprehension of phenomena with significant social implications. This paper is dedicated to conducting a rigorous case study in this vein, bridging the gap between the foundational principles of science and the cutting-edge capabilities of Machine Learning (ML) technologies.

EXplainable Artificial Intelligence (XAI) emerges as a critical paradigm in shedding light on ML models' often opaque inner workings. While previous research has ventured into various domains, the application of XAI principles to causal analysis remains scarcely explored. By integrating causality and XAI, this study aims to enrich our understanding of social phenomena and how they are reflected in state-of-the-art (SOTA) ML models facing the representation of the phenomena.

The rise of online hate speech, especially hateful memes (Fig. 1, top) —comprising both text and image, has prompted significant research. While multimodal ML algorithms have seen substantial improvements, efforts focus

more on benchmarking and maximizing performance, including the Hateful Memes Challenge competition (Kiela et al. 2020), rather than applying XAI methods. Existing approaches also lack a focus on causal architecture. This study defines hateful meme detection as an Average Treatment Effect (ATE) estimation problem for input data modalities (image and text) and examines the effects through the prism of XAI.

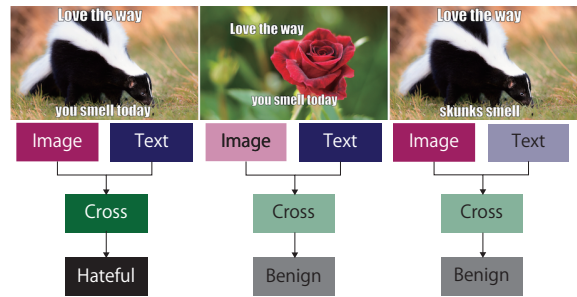


Figure 1: Visualization of a hateful meme and its corresponding confounders. (top) Meme samples and (bottom) their directed acyclic graph representation. (left) A hateful meme highlights cross-modal interactions between its image and text components that contribute to its hatefulness. (middle) The image benign confounder showcases original text and a non-hateful image, resulting in reduced cross-modal interactions and a benign classification. (right) The text benign confounder comprises an original image and non-hateful text. Note: The samples depicted are illustrative and do not exist in the dataset. ©Getty Images

Intersectionality<sup>1</sup>, or *the network of connections between social categories such as race, class, and gender, especially when this may result in additional disadvantage or discrimination*, acts as a bridge between ML and social science. Though broadly applied in social science and used for debiasing in ML, its wider applications are limited. How can we use this concept as a *generalized* tool for broader problems? Motivated by this question, this paper proposes reframed intersectionality, explores whether causally formalized intersectionality can address a broader range of problems, and

<sup>1</sup>Oxford Dictionary

evaluates *inductive* bias in ML models.

Furthermore, the excellence of Large Language Models (LLMs) across various benchmarks has been showcased, particularly in few-shot learning. The concept of in-context learning presents a promising avenue, but its formal causal evaluation is limited. Here, we address this problem. Our contribution could be summarized as:

- Formalization of hateful meme detection as an intersectional causal effect estimation problem, allowing performance assessment based on data generation process.
- Introduction of reframed causal intersectionality to evaluate inductive bias, marking a step towards broader applications, including demographics-nationality intersectionality, financial inclusion, and clinical diagnosis.
- Demonstration that attention attribution scores (Hao et al. 2021) divided by modality interaction describe the causal effect accurately, unlike non-divided scores. This finding opens doors for causal explainability in multimodal settings (Liu et al. 2022).
- Pioneering formal and empirical analysis of LLM’s meta-optimization process in the multimodal in-context setting.

## Related Work

### Causal ML and XAI

Causal Inference (CI) occupies a pivotal role in the elucidation of social phenomena and the interpretation of intervention outcomes. It bifurcates into two primary methodologies: the graphical and structural schemas for modeling reality (Pearl 2001), alongside the framework for potential outcome prediction (Rubin 2008). CI’s utility spans a diverse array of sectors, including but not limited to, medicine (Vlontzos, Rueckert, and Kainz 2022), manufacturing (Vuković and Thalmann 2022), and the social sciences (Sengupta and Srivastava 2021), guiding the interpretation of data within these fields. Furthermore, CI principles have been applied within Machine Learning (ML) and its allied disciplines, giving rise to the subfield of Causal ML. Causal ML encompasses research into natural language processing (Yang, Han, and Poon 2022), hate speech detection (Chakraborty and Masud 2022), and the study of image-text multimodality (Sanchez et al. 2022). Notably, the theoretical underpinnings of hateful memes, as a convergence of these interests, remain underexplored. Our research attempts to map out graphical and formal representations of the causal structures underlying hateful memes.

In conjunction, EXplainable Artificial Intelligence (XAI) (Speith 2022; Joshi, Walambe, and Kotecha 2021; Barredo Arrieta et al. 2020) has sought to demystify the internal mechanisms of ML models. XAI’s domain of inquiry extends across various fields, including medicine (Holzinger 2021) and energy (Machlev et al. 2022), with a particular focus on both model-agnostic (Sundararajan, Taly, and Yan 2017; Gaur, Faldu, and Sheth 2021; Marcos, Lobry, and Tuia 2019) and model-specific (Hao et al. 2021; Holzinger et al. 2021) evaluations. However, the

intersection of causality with XAI remains nascent. This study investigates XAI’s capability in assessing attributions to causality metrics, emphasizing gradient-based interpretations as central to XAI endeavors.

Since ML models typically minimize the gradient for optimization, components with steep gradients toward the model’s decision-making are considered crucial. The gradient-based XAI approach (Selvaraju et al. 2017), often model-specific, finds pertinent application in the analysis of Transformers (Vaswani et al. 2017), which underpin most SOTA models in natural language processing (NLP). Here, quantifying the attribution of attention matrix weights via the gradient emerges as a direct method (Hao et al. 2021). Our research proposes both theoretical and empirical advancements in causal analysis, leveraging this gradient-based methodology to enhance understanding and interpretation within the causality domain.

### Intersectionality

Intersectionality, a bias indicator of multiple demographics within various domains, has inspired a few causal analyses (Yang, Loftus, and Stoyanovich 2021; Bright, Malinsky, and Thompson 2016). While XAI techniques have been used to alleviate its negative impact in ML literature (Lalor et al. 2022; Simatele and Kabange 2022), our work redefines intersectionality for broader problems, and pioneers the quantification of inductive intersectional bias.

### Hate Speech and Hateful Memes

Hate speech and hateful memes have attracted substantial ML research attention, involving various models (Das, Wahi, and Li 2020; Lippe et al. 2020) and datasets (Kiela et al. 2020; de Gibert et al. 2018; Davidson et al. 2017; Sabat, Ferrer, and Giro-i-Nieto 2019; Suryawanshi et al. 2020). Previous analytical works have focused on racial bias (Sengupta and Srivastava 2021; Sharma et al. 2022), virality (Ling et al. 2021; Chakraborty and Masud 2022), and propaganda techniques (Dimitrov et al. 2021), and a few have applied XAI methods (Cao et al. 2021; Hee, Lee, and Chong 2022; Deshpande and Mani 2021). This study builds upon these by formalizing hateful meme detection as a causal effect estimation problem and emphasizing the importance of modality interaction.

### LLM

Large Language Models (LLMs), known for their powerful in-context few-shot learning capabilities (Brown et al. 2020) in various NLP and multimodal tasks, are emerging as significant tools (Zhao et al. 2023). To understand their inner workings, meta-gradient, or the update of the attention weights as a secondary form of the gradient, explains in-context learning empirically (Coda-Forno et al. 2023; Chen et al. 2022) and theoretically (Dai et al. 2023; von Oswald et al. 2023).

However, understanding LLM’s causal power remains a complex and emerging area of study (Ban et al. 2023; Kıcıman et al. 2023). Moreover, unlike the traditional classifier-attached-encoder model (e.g. BERT (Devlin et al.

2019)) with *predicted probability*, how to estimate the causal effect of chatbot-style LLM and analyze its inner workings quantitatively remains elusive. This study demonstrates that a dedicated task design could be used to estimate LLM’s causal effect and that a meta-gradient could explain its inner workings concerning that effect.

## Methodology

### Background

**Average Treatment Effect** The *Average Treatment Effect (ATE)* (Rubin 2008) is a key metric for assessing causal impacts. It represents the average difference in outcomes between treated and untreated groups. This measurement facilitates a standardized evaluation of causal effects across diverse scenarios. For a binary treatment  $BT \in (0, 1)$  yielding outcome  $\theta_{BT}$ , *ATE* is defined as:

$$ATE = \theta_1 - \theta_0 \quad (1)$$

This research repurposes intersectional *ATE* to evaluate hateful meme detection models.

**Causal Intersectionality** Building upon the textual definition, causal intersectionality (Bright, Malinsky, and Thompson 2016) challenges the simplistic aggregation of individual demographic effects. Instead, it highlights the complex, synergistic interactions between multiple demographic factors, acknowledging the nuanced dynamics that influence causal relationships in social studies. Defining causal intersectionality (Eq. 2) involves binary vectors for two demographics (e.g., gender  $D_1$  and color  $D_2$ ) marked as  $D = \{D_1, D_2\}$ , and the outcome  $\theta$ .

$$\theta_D \neq \sum_i \theta_{D_i} \quad (2)$$

Using this causality structure, we assess multimodal models within a causal context.

**Attention Attribution Score** Simply put, *attention attribution score* quantifies the contribution of attention weights to the model’s decision-making. Initially, a seminal work (Sundararajan, Taly, and Yan 2017) introduced the *integrated gradient* method for quantifying model component’s attribution. This method calculates the contribution of specific model components based on the gradient’s integral over that component. Upon this work, another study (Hao et al. 2021) reported its applicability to Transformer’s attention weights, deriving attention attribution score (*attr* herein). This approach proves critical for interpreting the Transformer’s behavior, especially in understanding how the attention matrix influences model outputs. Given hyperparameter  $\alpha$ , *attr* computes the integrated gradient for attention matrix  $A$  relative to a Transformer’s output  $\theta$ .

$$attr = A * \int_{\alpha=0}^1 \frac{\partial \theta(\alpha A)}{\partial A} d\alpha \quad (3)$$

We employ modality-wise averaged *attr* differences as causal effect indicators.

**Learning Objectives: Binary Classifier vs LLM** For the classifier-attached-encoder model, hateful meme detection aligns with binary classification. In summary, given *both hateful and benign pairs*, the classifier tries to maximize the classification performance. With a ground-truth label  $y_{gt}$  and a loss function  $f_{loss}$ , the learning objective when training a model  $\theta$  is:

$$\operatorname{argmin}_{\theta} -\{y_{gt}f_{loss}(\theta) + (1 - y_{gt})f_{loss}(1 - \theta)\} \quad (4)$$

On the other hand, in in-context learning, hateful samples and their counterpart confounders are presented to LLM *in parallel*. This differs from the objective above in that the information of hateful samples is not used to handle confounders, and vice versa. For example, facing the hateful sample in the zero-shot setting, the meta-objective it is trying to meta-optimize to is:

$$\operatorname{argmin}_{\theta} -y_{gt}f_{loss}(\theta) \quad (5)$$

This study delves into the meta-objective for LLM in a few-shot context for the equivalent comparison.

**Meta-Optimization in Few Shot Setting** Meta-optimization, in the realm of in-context learning, mirrors gradient descent (Irie, Csordas, and Jürgen 2022; Dai et al. 2023). Given query  $q$  as an input, a dual learning process of linear attention  $\theta$  - fine-tuning and in-context learning with attention update  $\Delta A$  - is contrasted with  $A_{ZSL}$ , the weight in a zero-shot setting.

$$\theta = (A_{ZSL} + \Delta A)q \quad (6)$$

$\Delta A$  functions as a gradient variant named meta-gradient. Our extension encapsulates attention attribution and its subsequent in-context learning results.

### Proposed Methodology Overview

Expanding on the original causal intersectionality (Eq. 2), we define intersectionality for text  $T$  and image  $I$  modalities. With  $X_1 = (T_1, I_1)$  indicating hateful content and two types of benign samples  $X_0 \in \{(T_1, I_0), (T_0, I_1)\}$ , the multimodal intersectionality is:

$$\theta_{X_1} \neq \sum_{X_0} \theta_{X_0} \quad (7)$$

The remainder of this section elaborates on causal intersectionality in meme detection, *attr*-based modality assessment, and LLM evaluation, visualized in Fig. 2.

### Causal Multimodal Intersectionality

**Intersectionality Reframed** We broaden causal intersectionality, positing its utility beyond human demographics to include arbitrary components. This reframed intersectionality assesses interconnections between *arbitrary categories such as user demographics or input modalities*, and how they amplify effects on significant issues - capturing indirect effects<sup>2</sup> in social contexts. This adaptation remains consistent with the original causal formalism (Eq. 2).

<sup>2</sup>Defined as an effect of two variables  $X_1$  and  $X_2$  to variable  $X_3$  via another variable  $Z$

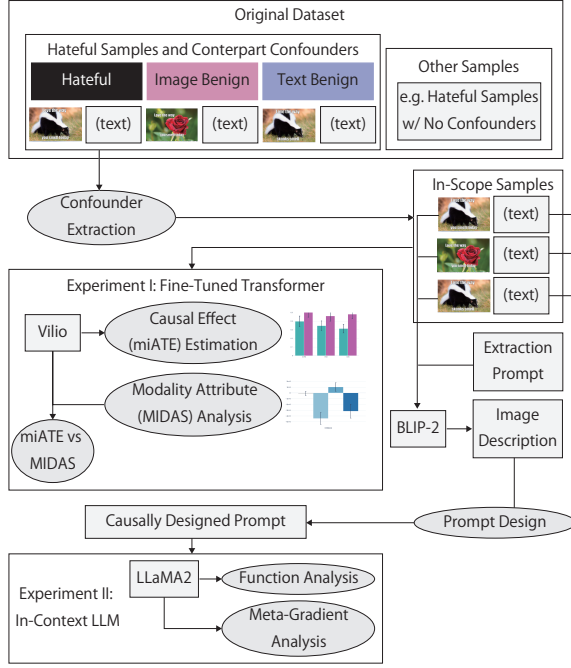


Figure 2: A schematic overview of our proposed methodology. Rectangular boxes denote data or models, while circular shapes represent the processes involved.

**Performance Measurement with Causal Intersectionality** In the data generation process of hateful memes, a text and an image, which are benign in isolation, jointly produce hate. Applying demographic intersectionality concepts, we introduce the *multimodal intersectional Average Treatment Effect (miATE)*.

$$miATE = \theta_{X_1} - \sum_{X_0} \theta_{X_0} \quad (8)$$

Model performance is assessed considering differences in each modality. Hateful samples are categorized by original benign text confounders  $T_0^{org}$  and image benign confounders  $I_0^{org}$ . The *Confounder Extraction* section provides more details.

### Attention Attribution Score by Interaction

We introduce the Modality Interaction Disentangled Attribution Score (MIDAS), which quantifies the contribution from various interaction types  $t \in \{within\_image, within\_text, cross\_modal\}$  to a model’s decision. Given an input  $X_i = (T, I)$ , with  $IT$  denoting the number of elements for interaction type  $t$ , *MIDAS* computes the modality-wise average  $avg_t$  of the attention attribution score  $attr^{X_i}$ , analogous to *miATE*. Following the initial work on attention attribution score (Hao et al. 2021), *MIDAS* is calculated from the last

hidden layer, excluding  $[CLS]$  and  $[SEP]$  tokens.

$$MIDAS = avg_t \{ attr^{X_1} - \sum_{X_0} (attr^{X_0}) \} \quad (9)$$

$$avg_t(X) = \frac{1}{IT} \sum^t X$$

Note that to mitigate the negative impact of class imbalance (Hossain, Sharif, and Hoque 2022),  $\theta$  and  $attr$  is averaged per sample and confounder category  $(T_1, I_1), (T_1, I_0), (T_0, I_1)$ .

### Formal Relation between miATE and MIDAS

We examine the relationship between *miATE* and *MIDAS* by proposing that *MIDAS* can be perceived as an attention attribution to *miATE*. Given  $G(A)$  as the one-step gradient for an attention matrix  $A$ , we see that  $attr$  approximates the product of the gradient and the attention (Eq. 10 - true when  $\alpha = 1$ ). *MIDAS* is expressed as the difference of that value between hateful and benign content (Eq. 11). Furthermore when *MIDAS* is aggregated across  $n$  samples that are representative of the entire dataset, it depicts the variation in attention expectancy normalized with the function  $\mathcal{N}$  across these samples, as shown in Eq. 11. In essence, we propose that *MIDAS* acts as an attention-focused representation of the model’s causal effect, i.e. *miATE*.

$$attr^{X_i} \simeq A^{X_i} * G(A^{X_i})$$

$$where G(A^{X_i}) = \frac{\partial \theta(A)}{\partial A^{X_i}} \quad (10)$$

$$MIDAS \simeq avg_t \{ A^{X_1} * G_{norm}(A^{X_1}) - \sum_{X_0} (A^{X_0} * G_{norm}(A^{X_0})) \}$$

$$\sum_n MIDAS \simeq \mathbb{E}[A^{X_1}] - \sum_{X_0} \mathbb{E}[A^{X_0}] \quad (11)$$

$$where G_{norm}(A) = \mathcal{N}(G(A)) \\ \mathcal{N}(G(A)) : G(A) \rightarrow (0, 1)$$

### LLM

**Causal Objective: Implicit miATE maximization** Before discussing LLM, we show that training a classifier implicitly addresses the *miATE* maximization problem. Specifically, the objective (Eq. 6) over a hateful-confounder pair could be written as:

$$\underset{\theta}{\operatorname{argmin}} - \{ f_{loss}(\theta_{X_1}) + \sum_{X_0} f_{loss}(1 - \theta_{X_0}) \} \quad (12)$$

Here, we see that the first term maximizes the first term of Eq. 8 and the second term minimizes the second term of Eq. 8. In contrast, zero-shot LLM only addresses the first term of Eq. 13. In the next section, we show how we design the task for LLM to aim for the same goal.



**Causal Task Design** To meta-optimize to the causal task, the hateful-confounder pair should be given to *an identical, not separate, meta-optimizer*, or the optimizer cannot have any information about intersectional causality (second term of Eq. 12). Table 1 shows the causality-oriented design of a representative prompt for hateful meme detection.

User: Out of image-caption pairs #0 to #2, select the most likely hateful or sarcastic pair with a potential label (hateful or sarcastic). If all pairs are benign, please say so.  
 #0: image: 'Skunk', caption: 'Love the way you smell today'  
 #1: image: 'Rose flower', caption: 'Love the way you smell today'  
 #2: image: 'Skunk', caption: 'Love the way skunks smell'  
 System:

Table 1: An illustrative prompt for the causal objective in the zero-shot scenario. In the few-shot context, answers are delivered succinctly (e.g., #0 is hateful.).

Note that this task design inherently counteracts sample imbalance since it simultaneously represents these hateful, original benign, and picked benign samples.

**Meta-Optimization for Causal Objective** Meta-optimization for the causal objective poses challenges (Niu et al. 2021) like complicated instruction and varied available labels. The optimization process consists of:

1. *Task Type Classification (TTC)*. LLM recognizes the task as a binary classification.
2. *Label Identification (LI)*. LLM provides probable labels, e.g., *hateful*, *sarcastic* (Chauhan 2020), and *benign*. Note that including the *sarcastic* label addresses its nuanced overlap with *hatefulness* (Sundaram et al. 2022), capitalizing on LLM’s comprehension of complex social phenomena embedded in training corpora. We still regard this task as a binary classification of *all-pair-benign* vs *one-pair-hateful*, with a subtask of *hateful-sarcastic* classification.

See Table 2 for a set of examples.

Subtask	Label	Response
TTC	Negative	Sorry, I couldn’t understand your instructions.
TTC	Positive	#1 could be sarcastic.
LI	Negative	#1 could be sarcastic.
LI	Positive	#0 could be hateful.

Table 2: Synthesized responses and subtask labels for in-context learning. Refer to Table 1 for the corresponding instruction prompt.

The meta-optimization process is segmented into these subtasks, with the understanding that *LI* follows a suc-

cessful *TTC*, denoted as  $TTC = 1$ . The output of meta-optimized Transformer block  $\theta$  could be formalized as:

$$\theta^{SubTask} = (A^{SubTask} + \Delta A^{SubTask})q$$

$$\theta = \begin{cases} \theta^{TTC} + \theta^{LI} & (TTC = 1) \\ \theta^{TTC} & otherwise \end{cases} \quad (13)$$

where  $SubTask \in \{TTC, LI\}$

## Experimental Settings

### Data Preparation

**Hateful Memes Dataset** Our study utilizes the Hateful Memes Challenge dataset (Kiela et al. 2020) and focuses primarily on the *dev\_seen* subset. Unimodal hateful samples (Das, Wahi, and Li 2020; Lippe et al. 2020) are omitted from our study.

**Confounder Extraction** From the dataset, 162 pairs of hateful ( $T_1^{org}, I_1^{org}$ ) and benign samples ( $T_1^{org}, I_0^{org}$ ) or ( $T_0^{org}, I_1^{org}$ ) are identified. Since most of the pairs have either one of the text or image confounders, not both, three random inputs with the missing modality ( $I_0^{picked}$  or  $T_0^{picked}$ ) are concatenated with the other modality to accommodate the requirements of Eq. 8, resulting in a uniquely crafted subset ( $T_1^{org}, I_0^{picked}$ ) and ( $T_0^{picked}, I_1^{org}$ ). The structure of this subset is summarized in Table 3.

Sample Category	# Samples
Hateful	162
Image Benign	78
Text Benign	84
Picked Image Benign	234
Picked Text Benign	252

Table 3: Samples utilized in our analysis. This table categorizes samples as Hateful, Image Benign ( $T_1^{org}, I_0^{org}$ ), Text Benign ( $T_0^{org}, I_1^{org}$ ), Picked Image Benign ( $T_1^{org}, I_0^{picked}$ ), and Picked Text Benign ( $T_0^{picked}, I_1^{org}$ ).

### Experiment I: Fine-Tuned Transformer

**Analysis Type** Assuming the predominant contribution of original inputs over the picked ones, in respect of the authors’ effort of making the task challenging, we divided the analysis into that of  $\{(T_1^{org}, I_1^{org}), (T_1^{org}, I_0^{org}), (T_0^{picked}, I_1^{org})\}$  (denoted as *org. text*), and of  $\{(T_1^{org}, I_1^{org}), (T_0^{org}, I_1^{org}), (T_1^{org}, I_0^{picked})\}$  (*org. image*).

**Models** We employ author implementation of the SOTA (Kiela et al. 2021) *Vilio* framework (Muennighoff 2020) for its superior capabilities and adaptable framework, focusing on its three main models: Oscar (Li et al. 2020b), UNITER (Chen et al. 2020), and VisualBERT (Li et al. 2020a), summarized in Table 4. Each model type has three submodels (training corpora or random seed variants), all included in our analysis but the results shown here are from selected one

submodel (preliminary analysis shows all submodels exhibit similar trend).

Type	Encoder	Pretraining Task
O	BERT (base)	1) Object tag (or <i>anchor</i> ) detection 2) Text-image contrastive learning
U	BERT (base)	1) Masked language modeling 2) Masked image modeling 3) Image-text matching 4) Word-region alignment via optimal transport
V	BERT (base)	1) Masked language modeling 2) Image captioning

Table 4: A categorization of Vilio’s submodels leveraged in our research. The models are classified into three groups: Oscar (O), UNITER (U), and VisualBERT (V).

## Experiment II: LLM

**Models** HuggingFace *Llama-2-13b-chat-hf* (Touvron et al. 2023) is our language model backbone, optimized for chat-style interactions. To convert the image into its textual description, we utilize the BLIP-2 (Li et al. 2023) model with a *FlanT5-XXL* (Chung et al. 2022) backbone.

**In-Context Learning** We study Llama-2’s behavior on image caption in the original dataset and image description generated by BLIP-2. For in-context learning, the number of samples is limited due to memory restriction. After the response is generated, one of the authors conducts manual labeling since the number of samples is limited (available at our GitHub repository). We gauge performance through accuracy.

**Meta-Gradient Evaluation** During our evaluations, we mask redundant subtext (e.g., *#0: image:* and *caption:*) in input prompts.

## Shared Settings

**Probing** We employ a probing (Alain and Bengio 2017) approach with LightGBM (Ke et al. 2017) to explore the impact of modality interaction and in-context learning on causal effect. Responses are split into training (56%), validation (14%), and test (30%) sets. We achieve hyperparameter tuning using Optuna (Akiba et al. 2019). To assess the effects of interaction type  $t$  (Experiment I, II) and model type (Experiment I), corresponding categorical variables and interaction terms with *MIDAS* (Experiment I) or summed attention weights (Experiment II) are added to our analysis. We determine significance using a t-test ( $p < 0.05$ ).

**Text-Only Pretrained BERT** For VisualBERT encoder replacement (Experiment I) and BLIP-2-fused-BERT (Experiment II), we use HuggingFace *bert-base-uncased*. In Experiment II, the last four layers of BERT and a linear classifier are trained for 100 epochs with an Adam optimizer (learning rate  $5e-5$ ), evaluating its performance across different seeds.

**External Resources** All code and experiments are accessible at <https://github.com/HireTheHero/CausalIntersectionalityDualGradient>. Experiments are conducted on a single NVIDIA A100 GPU, either through Google Colaboratory Pro+ or locally.

## Results & Discussion

### Experiment I: Fine-Tuned Transformer

**miATE** First, we assessed each model’s performance with *miATE* (Fig. 3). VisualBERT exhibited the highest disparity, highlighting its bias for text-based tasks (Table 4).

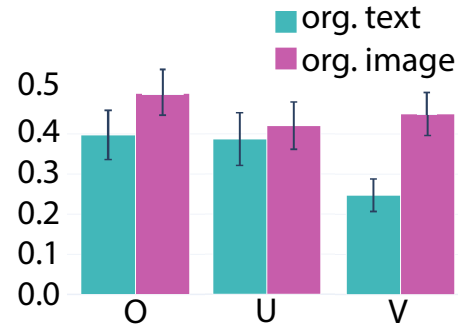


Figure 3: Multimodal Intersectional Average Treatment Effect (*miATE*) across Oscar (O; left), UNITER (U; middle), and VisualBERT (V; right) models, contrasting the samples with original image confounders (*org. image*, cyan) and those with original text confounders (*org. text*, magenta).

**MIDAS Global Analysis** Next, we assessed the model’s inner workings (Fig. 4 and 5). *MIDAS* of Oscar and UNITER (Fig. 4, first and second row) showed predictable trends of attending to one modality while the other is the same. In contrast, VisualBERT’s behavior of attending to text-related interactions (third row) mirrored its pretraining tendencies biased towards text (Table 4). Furthermore, replacing VisualBERT’s encoder with the one pretrained only with text enhanced the bias (fourth row), which supports the presence of pretraining bias represented in *MIDAS*. We observed no significant model differentiation with the original *attr* (Fig. 4, left column of each graph), suggesting a simple yet important contribution of modality-wise split.

**MIDAS Local Analysis** To see if we can interpret the single hateful-benign pair, we extracted local explanation (Chai et al. 2021; Hee, Lee, and Chong 2022). A representative pair (Fig. 5) illustrates that UNITER captures the contrast between a woman and cargo in image confounder analysis (first and third row), and the model similarly attended to the words *dishwasher* and *driving* for text analysis (first and second row).

**Empirical Relation between MIDAS & miATE** To probe the relationship between *MIDAS* and *miATE*, we first modeled the entire (*MIDAS*, *miATE*) pairs of all the models by a single probe, resulting in a moderate AUC of

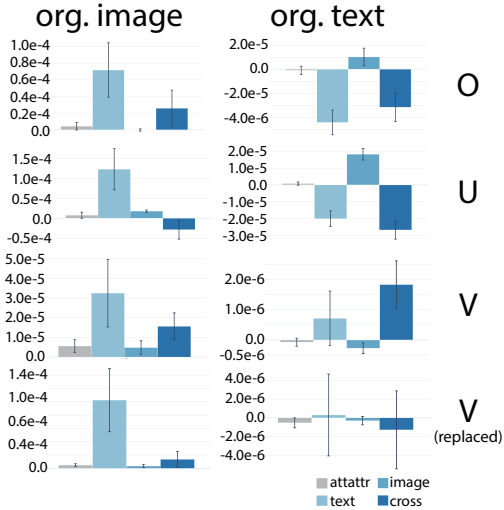


Figure 4: *MIDAS* for *org. image* (left) and *org. text* (right) samples featuring Oscar (top), UNITER (second row), VisualBERT (third row), and VisualBERT with text-only-pretrained encoder (bottom). From left to right, each graph displays *attr* with no modality division,  $MIDAS_{within\_text}$ ,  $MIDAS_{within\_image}$ , and  $MIDAS_{cross\_modal}$ .

Interaction Type	O	U	V
within-text	35±25	30±26	253±66
cross-modal	19±14	35±22	223±75
within-image	30±23	54±32	169±74

Table 5: LightGBM probe’s feature importance between *MIDAS* and *miATE*. Values represent frequency counts.

$75.6 \pm 4.20$ . Next, to see the effectiveness of probing for each model, we applied one probe for one model, resulting in the highest AUC for V (AUC  $94.1 \pm 3.01$ ), while low-to-moderate for O ( $60.8 \pm 5.18$ ) and U ( $74.3 \pm 2.39$ ). In addition, to see if the probes reflect the findings on *MIDAS*, we analyzed the feature importance (Table 5). Consistent with previous findings, VisualBERT and within-text appear with the highest frequency among model type and interaction type, respectively. In summary, our findings showed a moderate correlation between *MIDAS* and *miATE* for BERT-based models, with a particularly robust link for VisualBERT, echoing its distinct model nature.

## Experiment II: LLM

**Effectiveness of BLIP-2 information retrieval** To assess the merit of BLIP-2 information retrieval, we utilized its image description and the original captions to fine-tune BERT pretrained only with text. The resultant enhanced performance (Accuracy  $66.9 \pm 0.84$ , AUC:  $71.2 \pm 1.55$ ) to unimodal BERT benchmarks (Kiela et al. 2020) underlines BLIP-2’s effective image information extraction capabilities.

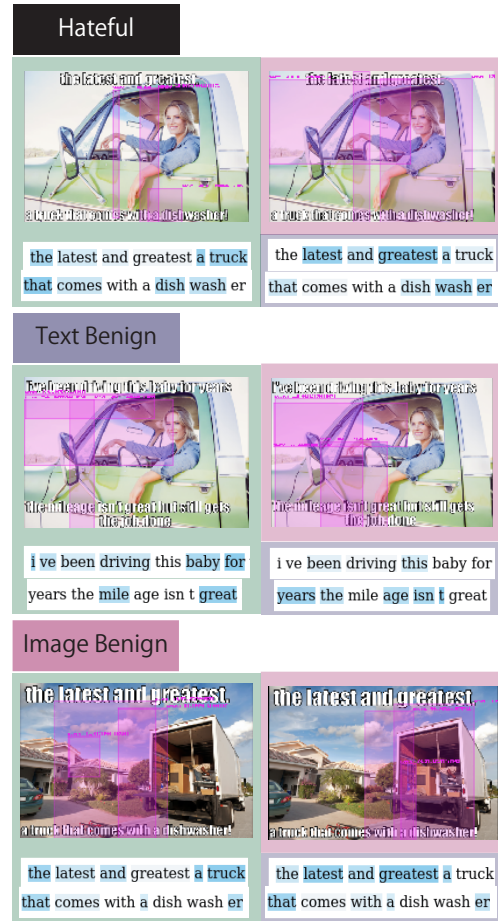


Figure 5: Conceptual portrayal of hateful, text benign, and image benign samples derived from UNITER. *MIDAS* reflects heightened  $attr_{cross\_modal}$  (green),  $attr_{within\_image}$  (red), or  $attr_{within\_text}$  (blue) values. Both image and text inputs spotlight top-scored ROIs and tokens. The text is abbreviated for clarity.

**In-Context LLM Performance** Our evaluation with Llama-2 shows that all-sample accuracy improved after one sample (Table 6, left). Interestingly, after just one in-context example, the model achieved exactly the same performance for all samples and  $TTC = 1$  samples, meaning impeccable TTC Recall. These results suggest the critical role of in-context examples in task comprehension when the task is challenging in zero-shot settings. Marking the sarcastic label as positive led to better performance at the zero-shot setting but dropped after one example (Table 6, middle and right), implying uncertainty in the decision-making for this label.

**Meta-Optimization Evaluation** To gauge the influence of meta-optimization, we applied a probe model to examine the relationship between summed attention weights ( $A, \Delta A$ ) and TTC label, revealing a moderate AUC of  $82.3 \pm 6.16$ . Furthermore, a detailed extraction of feature importance (Table 7) from the probe model allowed us to

# Few-Shot Samples	Accuracy		
	All(S-)	TTC(S-)	TTC(S+)
0	46.2	61.6	62.6
1	62.9	62.9	60.6
2	62.5	62.5	61.0
3	59.2	59.2	57.6
4	64.3	64.3	71.4

Table 6: Zero-shot (first row) and few-shot (second to fifth row) Llama-2 performance. *All* signifies cumulative sample results, while *TTC* relates to correct TTC samples. Parentheticals (*S+* or *S-*) denote the inclusion of the *sarcastic* label in either positive or negative samples.

Interaction Type	$A$	$\Delta A$
within-text	65±41	24±17
cross-modal	58±33	31±21
within-image	43±25	14±8

Table 7: LightGBM probe’s feature importance in Experiment II. Features of zero-shot weights ( $A$ ) or their few-shot updates ( $\Delta A$ ) are divided by interaction types.

determine how  $A$  and  $\Delta A$  impact the probe model (left and right). Our findings suggest that while  $A$  carries substantial weight in decision-making, its meta-optimized counterpart  $\Delta A$  also plays a vital role. Regarding the effects of modality interaction (top and bottom), captured by interaction-type-divided weights  $A+\Delta A$ , each interaction type did contribute to TTC. When examining the differential impacts of each type, however, no significant disparities in their contributions were identified. Addressing the challenge of discerning between them will be a part of the future work.

## Discussion

The primary goal of this paper is to assess models based on the data generation process and its underlying concepts - *hatefulness* in the case of hateful memes. While this deviates from standard ML evaluations focusing on performance metrics like accuracy, it is scientifically valid and relevant to ML problems, like Rubin started his line of causal inference works to analyze the impact of nulled variables (Rubin 2008). In our study, we demonstrate that the generation of hateful memes embodies multimodal intersectionality, and the SOTA Transformer models effectively capture this nature of the data but are biased by pretraining datasets. In the future, we hope to apply our method to other multimodal problems like the missing modality problem (Ma et al. 2021; Wang et al. 2023), an inherently close one to nulled variable evaluation.

Our study carves a niche by reconceptualizing hateful meme detection through the lens of modality interaction and causal effect. Compared to the seminal work on causal intersectionality (Bright, Malinsky, and Thompson 2016), beyond mere technical insights, we proffer a paradigm shift in causal intersectionality.

Our method’s key advantage is its unique capability for

modality-wise causal analysis, a novel contribution in this field. Despite its simplicity, the causal effect of the modalities is neither investigated nor formally defined in the existing literature.

Empirically, Experiment I unveils model biases overlooked by traditional methods like attention attribution scores (Hao et al. 2021; Hee, Lee, and Chong 2022) without consideration for modality.

For Experiment II, our exploration of few-shot LLM performance has provided an understanding of how LLM adapts to different levels of input information, shedding light on their capabilities and limitations in various scenarios. Applying meta-gradients has allowed us to assess the attribution of attention weights, adding granularity to the interpretability landscape. Evaluating the effectiveness of the causal task over the causal evaluation of LLMs is challenging since it is a new concept. Nonetheless, this could be a valuable benchmark for future model evaluations. Despite relying on a specific instruction prompt for the causal task, we could adapt the design for broader applications. For example, with a simple modification to the prompt, we could test LLMs with multi-class meme classification (Davidson et al. 2017).

## Conclusion

We posit that hateful meme detection transcends mere classification, gravitating towards intersectional causal effect analysis. Our evaluations spanned various Transformer architectures in unique settings. To ensure our approach’s universality, extending our evaluations to other hateful memes datasets (Gomez et al. 2020; Das et al. 2023) will be pivotal. In the quest for broader insights, exploring diverse challenges, such as the intersectionality in multimodal medical analyses (Azilinson et al. 2023), will be part of our future work. For scalability, utilizing more of the power of LLMs will be promising for confounder extraction and response evaluation.

## Ethical Considerations

In this research, we aim to develop innovative analytical methods for identifying and mitigating the proliferation of hateful memes, a pressing concern given the complex interplay of text and imagery in propagating hate speech online. The nature of this endeavor necessitates rigorous ethical scrutiny, especially concerning the selection, utilization, and presentation of these memes within our scholarly work and its broader dissemination. Herein, we delineate the principal ethical considerations guiding our study.

The hateful memes we examine originate from a prior investigation (Kiela et al. 2020). We direct readers to the original study for insights into the ethical measures employed during the dataset’s compilation. Our engagement with this dataset is underpinned by a firm commitment not to propagate or validate the adverse messages it encompasses.

Our sample selection approach is predicated on a causality framework detailed in the *Confounder Extraction* section, ensuring a comprehensive examination of hate speech manifestations across diverse community targets. This methodology underscores our commitment to a nuanced analysis that refrains from generalizations or biases.

A pivotal aspect of our ethical strategy is to reconcile the imperative of methodological transparency with the necessity to limit harm. Consequently, we exhibit restraint in our presentation of hateful memes. Specifically, Figure 5 is the sole instance within our publication where an actual hateful meme is depicted. We have exercised meticulous care to ensure that neither the accompanying text description nor the figure caption disseminates any form of hate speech.

This strategy is emblematic of our broader ethical stance, emphasizing the conscientious handling of sensitive content. Our research is animated by a profound dedication to combating hate speech in all its forms, reflecting an unwavering commitment to ethical research practices that respect the dignity of all individuals and communities. Through this work, we aspire not only to advance the field of hate speech detection but also to contribute meaningfully to creating more inclusive and respectful digital spaces.

## Limitations

This study’s primary limitation concerns the unverified generalizability of its findings. Hateful memes represent an evolving area of concern that necessitates extensive, openly accessible datasets for comprehensive analysis and validation. Our research endeavors to tackle this challenge, yet the broader scope for future exploration is highlighted by the potential applications of our findings, as detailed in the *Discussion* and the *Conclusion* sections.

A further constraint is the linguistic homogeneity of the dataset employed, with the Hateful Memes Challenge dataset comprising exclusively English-language textual content. This presents a critical limitation in the context of the global escalation of extremism, where hate speech proliferates across linguistic boundaries. The detection of multilingual hate speech thus emerges as a crucial area for future research, necessitating methodologies capable of navigating language-specific nuances and cultural contexts.

Additionally, the field of hate speech detection faces resource limitations, notably in the size and diversity of available datasets. Hateful speech datasets are generally small, restricting the depth and breadth of training data for machine learning models. We believe future studies could utilize LLMs as dataset curators.

In summary, while this study contributes valuable insights into detecting and mitigating hateful memes, it also underscores the need for further research. Addressing the limitations related to dataset generalizability, linguistic diversity, and the scarcity of training data are pivotal steps toward developing more effective and universally applicable solutions for combating online hate speech. Exploring innovative methods, such as LLM-based dataset curation, represents a promising direction for overcoming these challenges. From a theoretical standpoint, the groundwork of our in-context learning analysis relies upon the principles of simplified linear attention (Irie, Csordas, and Jürgen 2022; Dai et al. 2023). However, this foundation’s direct applicability to conventional Transformer models invites scrutiny. Consequently, a more nuanced interpretation (Ren and Liu 2023) may be imperative for advancing our understanding in future investigations.

## Acknowledgments

CyberAgent Inc.<sup>3</sup> supported this work, enabling our conference participation. We are profoundly grateful to the anonymous reviewers, whose constructive feedback significantly enriched our discussions and contributed to refining our research. Yosuke Miyanishi thanks J.X., H.W., and T.Y. for their insightful comments that significantly improved our study.

---

<sup>3</sup><https://www.cyberagent.co.jp/>

## References

- Akiba, T.; Sano, S.; Yanase, T.; Ohta, T.; and Koyama, M. 2019. Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2623–2631. Anchorage AK USA: ACM. ISBN 978-1-4503-6201-6.
- Alain, G.; and Bengio, Y. 2017. Understanding Intermediate Layers Using Linear Classifier Probes. In *5th International Conference on Learning Representations*. Toulon, France.
- Azilanon, M.; Makhlova, J.; Zaaraoui, W.; Medina Villalon, S.; Viout, P.; Roussel, T.; El Mendili, M. M.; Ridley, B.; Ranjeva, J.-P.; Bartolomei, F.; Jirsa, V.; and Guye, M. 2023. Combining Sodium MRI, Proton MR Spectroscopic Imaging, and Intracerebral EEG in Epilepsy. *Human Brain Mapping*, 44(2): 825–840.
- Ban, T.; Chen, L.; Wang, X.; and Chen, H. 2023. From Query Tools to Causal Architects: Harnessing Large Language Models for Advanced Causal Discovery from Data. *arXiv preprint*.
- Barredo Arrieta, A.; Díaz-Rodríguez, N.; Del Ser, J.; Benetot, A.; Tabik, S.; Barbado, A.; Garcia, S.; Gil-Lopez, S.; Molina, D.; Benjamins, R.; Chatila, R.; and Herrera, F. 2020. Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI. *Information Fusion*, 58: 82–115.
- Bright, L. K.; Malinsky, D.; and Thompson, M. 2016. Causally Interpreting Intersectionality Theory. *Philosophy of Science*, 83(1): 60–81.
- Brown, T. B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; Agarwal, S.; Herbert-Voss, A.; Krueger, G.; Henighan, T.; Child, R.; Ramesh, A.; Ziegler, D. M.; Wu, J.; Winter, C.; Hesse, C.; Chen, M.; Sigler, E.; Litwin, M.; Gray, S.; Chess, B.; Clark, J.; Berner, C.; McCandlish, S.; Radford, A.; Sutskever, I.; and Amodei, D. 2020. Language Models Are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, 1877–1901. Curran Associates, Inc.
- Cao, R.; Fan, Z.; Lee, R. K.-W.; Chong, W.-H.; and Jiang, J. 2021. Disentangling Hate in Online Memes. *Proceedings of the 29th ACM International Conference on Multimedia*, 5138–5147.
- Chai, Y.; Zhou, Y.; Li, W.; and Jiang, Y. 2021. An Explainable Multi-Modal Hierarchical Attention Model for Developing Phishing Threat Intelligence. *IEEE Transactions on Dependable and Secure Computing*, 1–1.
- Chakraborty, T.; and Masud, S. 2022. Nipping in the Bud: Detection, Diffusion and Mitigation of Hate Speech on Social Media. *ACM SIGWEB Newsletter*, 2022(Winter): 1931–1745.
- Chauhan, D. S. 2020. All-in-One: A Deep Attentive Multi-task Learning Framework for Humour, Sarcasm, Offensive, Motivation, and Sentiment on Memes. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, 281–290. Suzhou, China: Association for Computational Linguistics.
- Chen, Y.; Zhong, R.; Zha, S.; Karypis, G.; and He, H. 2022. Meta-Learning via Language Model In-context Tuning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 719–730. Dublin, Ireland: Association for Computational Linguistics.
- Chen, Y.-C.; Li, L.; Yu, L.; El Kholy, A.; Ahmed, F.; Gan, Z.; Cheng, Y.; and Liu, J. 2020. UNITER: UNiversal Image-TEXT Representation Learning. In Vedaldi, A.; Bischof, H.; Brox, T.; and Frahm, J.-M., eds., *Computer Vision – ECCV 2020*, volume 12375, 104–120. Cham: Springer International Publishing. ISBN 978-3-030-58576-1 978-3-030-58577-8.
- Chung, H. W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, Y.; Wang, X.; Dehghani, M.; Brahma, S.; Webson, A.; Gu, S. S.; Dai, Z.; Suzgun, M.; Chen, X.; Chowdhery, A.; Castro-Ros, A.; Pellat, M.; Robinson, K.; Valter, D.; Narang, S.; Mishra, G.; Yu, A.; Zhao, V.; Huang, Y.; Dai, A.; Yu, H.; Petrov, S.; Chi, E. H.; Dean, J.; Devlin, J.; Roberts, A.; Zhou, D.; Le, Q. V.; and Wei, J. 2022. Scaling Instruction-Finetuned Language Models. *arXiv preprint*.
- Coda-Forno, J.; Binz, M.; Akata, Z.; Botvinick, M.; Wang, J. X.; and Schulz, E. 2023. Meta-in-Context Learning in Large Language Models. *arXiv preprint*.
- Dai, D.; Sun, Y.; Dong, L.; Hao, Y.; Ma, S.; Sui, Z.; and Wei, F. 2023. Why Can GPT Learn In-Context? Language Models Implicitly Perform Gradient Descent as Meta-Optimizers. In *Findings of the Association for Computational Linguistics*, 4005–4019. Association for Computational Linguistics.
- Das, A.; Wahi, J. S.; and Li, S. 2020. Detecting Hate Speech in Multi-modal Memes. *arXiv preprint*.
- Das, M.; Raj, R.; Saha, P.; Mathew, B.; Gupta, M.; and Mukherjee, A. 2023. HateMM: A Multi-Modal Dataset for Hate Video Classification. *Proceedings of the International AAAI Conference on Web and Social Media*, 17: 1014–1023.
- Davidson, T.; Warmsley, D.; Macy, M.; and Weber, I. 2017. Automated Hate Speech Detection and the Problem of Offensive Language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media*, 512–515. Montreal, Canada.
- de Gibert, O.; Perez, N.; García-Pablos, A.; and Cuadros, M. 2018. Hate Speech Dataset from a White Supremacy Forum. In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, 11–20. Brussels, Belgium: Association for Computational Linguistics.
- Deshpande, T.; and Mani, N. 2021. An Interpretable Approach to Hateful Meme Detection. In *ICMI '21: Proceedings of the 2021 International Conference on Multimodal Interaction*, 723–727.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of NAACL-HLT*, 4171–4186.



- Dimitrov, D.; Bin Ali, B.; Shaar, S.; Alam, F.; Silvestri, F.; Firooz, H.; Nakov, P.; and Da San Martino, G. 2021. Detecting Propaganda Techniques in Memes. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 6603–6617. Online: Association for Computational Linguistics.
- Gaur, M.; Faldu, K.; and Sheth, A. 2021. Semantics of the Black-Box: Can Knowledge Graphs Help Make Deep Learning Systems More Interpretable and Explainable? *IEEE Internet Computing*, 25(1): 51–59.
- Gomez, R.; Gibert, J.; Gomez, L.; and Karatzas, D. 2020. Exploring Hate Speech Detection in Multimodal Publications. In *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 1459–1467. Snowmass Village, CO, USA: IEEE. ISBN 978-1-72816-553-0.
- Hao, Y.; Dong, L.; Wei, F.; and Xu, K. 2021. Self-Attention Attribution: Interpreting Information Interactions Inside Transformer. In *The 35th AAAI Conference on Artificial Intelligence*. AAAI Press.
- Hee, M. S.; Lee, R. K.-W.; and Chong, W.-H. 2022. On Explaining Multimodal Hateful Meme Detection Models. In *Proceedings of the ACM Web Conference 2022*, 3651–3655. Virtual Event, Lyon France: ACM. ISBN 978-1-4503-9096-5.
- Holzinger, A. 2021. Explainable AI and Multi-Modal Causability in Medicine. *i-com*, 19(3): 171–179.
- Holzinger, A.; Malle, B.; Saranti, A.; and Pfeifer, B. 2021. Towards Multi-Modal Causability with Graph Neural Networks Enabling Information Fusion for Explainable AI. *Information Fusion*, 71: 28–37.
- Hossain, E.; Sharif, O.; and Hoque, M. M. 2022. MUTE: A Multimodal Dataset for Detecting Hateful Memes. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing: Student Research Workshop*. Association for Computational Linguistics.
- Irie, K.; Csordas, R.; and Jürgen, S. 2022. The Dual Form of Neural Networks Revisited: Connecting Test Time Predictions to Training Patterns via Spotlights of Attention. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162. Baltimore, Maryland, USA.
- Joshi, G.; Walambe, R.; and Kotecha, K. 2021. A Review on Explainability in Multimodal Deep Neural Nets. *IEEE Access*, 9: 59800–59821.
- Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; and Liu, T.-Y. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *31st Conference on Neural Information Processing Systems*. Long Beach, CA, USA.
- Kıcıman, E.; Ness, R.; Sharma, A.; and Tan, C. 2023. Causal Reasoning and Large Language Models: Opening a New Frontier for Causality. *arXiv preprint*.
- Kiela, D.; Firooz, H.; Mohan, A.; Goswami, V.; Singh, A.; Ringshia, P.; and Testuggine, D. 2020. The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes. In *Thirty-Fourth Annual Conference on Neural Information Processing Systems*. Red Hook, NY, USA.
- Kiela, D.; Firooz, H.; Mohan, A.; Goswami, V.; Zhu, R.; Muennighoff, N.; Velioglu, R.; Rose, J.; Lippe, P.; Holla, N.; Chandra, S.; Rajamanickam, S.; Antoniou, G.; Shutova, E.; Yannakoudakis, H.; Sandulescu, V.; Ozertem, U.; Pantel, P.; Specia, L.; and Parikh, D. 2021. The Hateful Memes Challenge: Competition Report. *Proceedings of Machine Learning Research*.
- Lalor, J.; Yang, Y.; Smith, K.; Forsgren, N.; and Abbasi, A. 2022. Benchmarking Intersectional Biases in NLP. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3598–3609. Seattle, United States: Association for Computational Linguistics.
- Li, J.; Li, D.; Savarese, S.; and Hoi, S. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In *International Conference on Machine Learning*.
- Li, L. H.; Yatskar, M.; Yin, D.; Hsieh, C.-J.; and Chang, K.-W. 2020a. VisualBERT: A Simple and Performant Baseline for Vision and Language. In *The 58th Annual Meeting of the Association for Computational Linguistics*, 5265–5275. Association for Computational Linguistics.
- Li, X.; Yin, X.; Li, C.; Zhang, P.; Hu, X.; Zhang, L.; Wang, L.; Hu, H.; Dong, L.; Wei, F.; Choi, Y.; and Gao, J. 2020b. Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks. In *16th European Conference On Computer Vision*.
- Ling, C.; AbuHilal, I.; Blackburn, J.; De Cristofaro, E.; Zannettou, S.; and Stringhini, G. 2021. Dissecting the Meme Magic: Understanding Indicators of Virality in Image Memes. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1): 1–24.
- Lippe, P.; Holla, N.; Chandra, S.; Rajamanickam, S.; Antoniou, G.; Shutova, E.; and Yannakoudakis, H. 2020. A Multimodal Framework for the Detection of Hateful Memes. *arXiv preprint*.
- Liu, Y.; Wei, Y.-S.; Yan, H.; Li, G.-B.; and Lin, L. 2022. Causal Reasoning Meets Visual Representation Learning: A Prospective Study. *Machine Intelligence Research*, 19(6): 485–511.
- Ma, M.; Ren, J.; Zhao, L.; Tulyakov, S.; Wu, C.; and Peng, X. 2021. SMIL: Multimodal Learning with Severely Missing Modality. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence*.
- Machlev, R.; Heistrene, L.; Perl, M.; Levy, K.; Belikov, J.; Mannor, S.; and Levron, Y. 2022. Explainable Artificial Intelligence (XAI) Techniques for Energy and Power Systems: Review, Challenges and Opportunities. *Energy and AI*, 9: 100169.
- Marcos, D.; Lobry, S.; and Tuia, D. 2019. Semantically Interpretable Activation Maps: What-Where-How Explanations within CNNs. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, 4207–4215.

- Muennighoff, N. 2020. Vilio: State-of-the-art Visio-Linguistic Models Applied to Hateful Memes. *arXiv preprint*.
- Niu, Y.; Tang, K.; Zhang, H.; Lu, Z.; Hua, X.-S.; and Wen, J.-R. 2021. Counterfactual VQA: A Cause-Effect Look at Language Bias. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 12695–12705. Nashville, TN, USA: IEEE. ISBN 978-1-66544-509-2.
- Pearl, J. 2001. CAUSALITY: MODELS, REASONING, AND INFERENCE. *Cambridge University Press*, 11.
- Ren, R.; and Liu, Y. 2023. In-Context Learning with Transformer Is Really Equivalent to a Contrastive Learning Pattern. *arxiv:2310.13220*.
- Rubin, D. B. 2008. For Objective Causal Inference, Design Trumps Analysis. *The Annals of Applied Statistics*, 2(3).
- Sabat, B. O.; Ferrer, C. C.; and Giro-i-Nieto, X. 2019. Hate Speech in Pixels: Detection of Offensive Memes towards Automatic Moderation. *AI for Social Good Workshop at NeurIPS 2019 (short paper)*.
- Sanchez, P.; Voisey, J. P.; Xia, T.; Watson, H. I.; O’Neil, A. Q.; and Tsaftaris, S. A. 2022. Causal Machine Learning for Healthcare and Precision Medicine. *Royal Society Open Science*, 9(8): 220638.
- Selvaraju, R. R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; and Batra, D. 2017. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, 618–626. Venice: IEEE. ISBN 978-1-5386-1032-9.
- Sengupta, K.; and Srivastava, P. R. 2021. Causal Effect of Racial Bias in Data and Machine Learning Algorithms towards User Persuasiveness & Discriminatory Decision Making: An Empirical Study. *arXiv preprint*.
- Sharma, S.; Akhtar, M. S.; Nakov, P.; and Chakraborty, T. 2022. DISARM: Detecting the Victims Targeted by Harmful Memes. In *Findings of the Association for Computational Linguistics: NAACL 2022*, 1572–1588. Seattle, United States: Association for Computational Linguistics.
- Simatele, M.; and Kabange, M. 2022. Financial Inclusion and Intersectionality: A Case of Business Funding in the South African Informal Sector. *Journal of Risk and Financial Management*, 15(9): 380.
- Speith, T. 2022. A Review of Taxonomies of Explainable Artificial Intelligence (XAI) Methods. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, 2239–2250. Seoul Republic of Korea: ACM. ISBN 978-1-4503-9352-2.
- Sundaram, V.; Pavan, R.; Kandaala, S.; and Reddy, R. R. 2022. Distinguishing Hate Speech from Sarcasm. In *2022 International Conference for Advancement in Technology (ICONAT)*, 1–5. Goa, India: IEEE. ISBN 978-1-66542-577-3.
- Sundararajan, M.; Taly, A.; and Yan, Q. 2017. Axiomatic Attribution for Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, 3319–3328. PMLR.
- Suryawanshi, S.; Chakravarthi, B. R.; Arcan, M.; and Buiteelaar, P. 2020. Multimodal Meme Dataset (MultiOFF) for Identifying Offensive Content in Image and Text. *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; Rodriguez, A.; Joulin, A.; Grave, E.; and Lample, G. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention Is All You Need. In *Thirty-First Annual Conference on Neural Information Processing Systems*. Long Beach, CA, USA.
- Vlontzos, A.; Rueckert, D.; and Kainz, B. 2022. A Review of Causality for Learning Algorithms in Medical Image Analysis. *Machine Learning for Biomedical Imaging*, 1(November 2022 issue): 1–17.
- von Oswald, J.; Niklasson, E.; Randazzo, Ettore; Sacramento, Jo\~{a}o; Mordvintsev, Alexander; Zhmoginov, Andrey; and Vladymyrov, Max. 2023. Transformers Learn In-Context by Gradient Descent. In *Proceedings of the 40th International Conference on Machine Learning*, volume 1464, 24. Honolulu, HI, USA: JMLR.org.
- Vuković, M.; and Thalmann, S. 2022. Causal Discovery in Manufacturing: A Structured Literature Review. *Journal of Manufacturing and Materials Processing*, 6(1): 10.
- Wang, H.; Chen, Y.; Ma, C.; Avery, J.; Hull, L.; and Carneiro, G. 2023. Multi-Modal Learning with Missing Modality via Shared-Specific Feature Modelling. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 15878–15887. Vancouver, BC, Canada: IEEE. ISBN 9798350301298.
- Yang, J.; Han, S. C.; and Poon, J. 2022. A Survey on Extraction of Causal Relations from Natural Language Text. *Knowledge and Information Systems*, 64(5): 1161–1186.
- Yang, K.; Loftus, J. R.; and Stoyanovich, J. 2021. Causal Intersectionality for Fair Ranking. In *2nd Symposium on Foundations of Responsible Computing*. FORC.
- Zhao, W. X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; Du, Y.; Yang, C.; Chen, Y.; Chen, Z.; Jiang, J.; Ren, R.; Li, Y.; Tang, X.; Liu, Z.; Liu, P.; Nie, J.-Y.; and Wen, J.-R. 2023. A Survey of Large Language Models. *arxiv:2303.18223*.

## Appendix

### Further Exploration for Local Explainability

Sample analysis for Oscar (Fig. 6) shows a similar trend to UNITER (Fig. 5). Interestingly, VisualBERT does not attend to the key components (woman or cargo) in the image, supporting its bias towards textual information.



Figure 6: Samples derived from Oscar.



Figure 7: Samples derived from VisualBERT.

### Breakdown of Attention Attribution Score

The attention attribution score  $attr$  (Eq. 3) is the product of the attention weight matrix and the integral of the gradient. To see the separate impact, we replaced the  $attr$  term of the  $MIDAS$  equation (Eq. 9) with the attention  $MIDAS_{att}$  or the gradient  $MIDAS_{grad}$  for comparison. In general,  $MIDAS_{att}$  (Fig. 8-10) shows a more similar trend to the original  $MIDAS$  than  $MIDAS_{grad}$  (Fig. 11-13). This result implies that the attention weights decide the model's strategy, while the gradient adjusts the impact of the individual component.

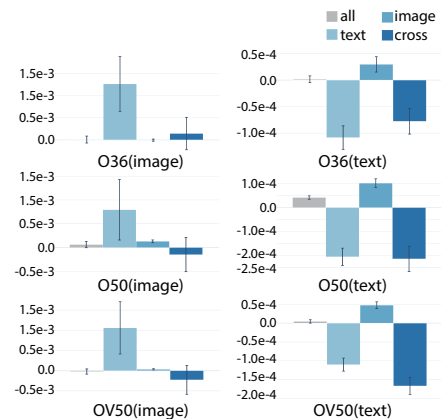


Figure 8: Oscar  $MIDAS_{att}$ .

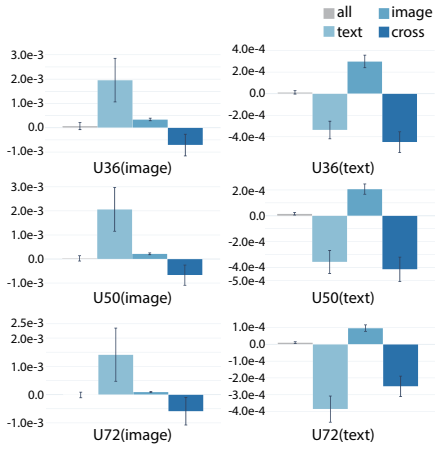


Figure 9: UNITER  $MIDAS_{att}$ .

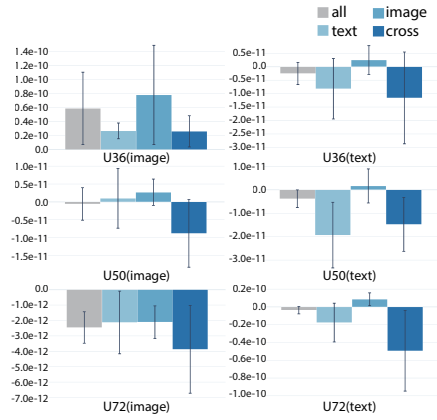


Figure 12: UNITER  $MIDAS_{grad}$ .

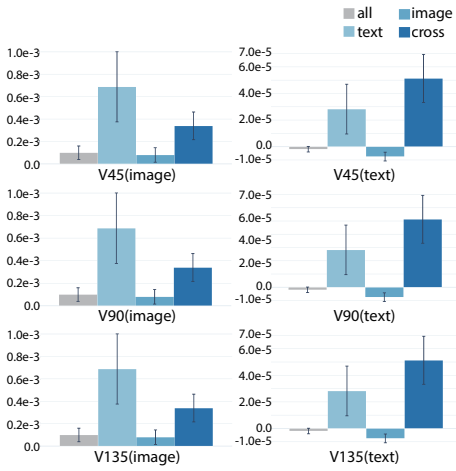


Figure 10: VisualBERT  $MIDAS_{att}$ .

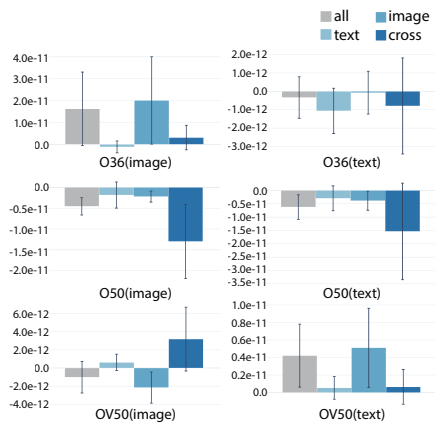


Figure 11: Oscar  $MIDAS_{grad}$ .

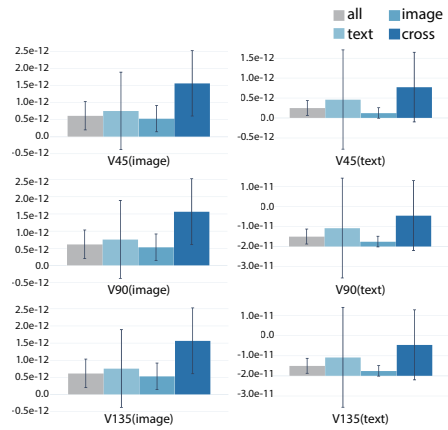


Figure 13: VisualBERT  $MIDAS_{grad}$ .