Improving the fact-checking performance of language models by relying on their entailment ability

Anonymous ACL submission

Abstract

001 Automated fact-checking is a crucial task in this digital age. To verify a claim, current approaches majorly follow one of two strategies i.e. (i) relying on embedded knowledge of language models, and (ii) fine-tuning them with evidence pieces. While the former can make systems to hallucinate, the later have not been very successful till date. The primary reason behind this is that fact verification is a complex process. Language models have to parse through multiple pieces of evidence before making a prediction. Further, the evidence pieces often contradict each other. This makes the reasoning process even more complex. We proposed a simple yet effective approach where we re-016 lied on entailment and the generative ability of language models to produce "supporting" 017 and "refuting" justifications (for the truthfulness of a claim). We trained language models based on these justifications and achieved superior results. Apart from that, we did a systematic comparison of different prompting and 022 fine-tuning strategies, as it is currently lacking in the literature. Some of our observations are: (i) training language models with raw evidence sentences registered an improvement up to 8.20% in macro-F1, over the best 027 performing baseline for the RAW-FC dataset, (ii) similarly, training language models with prompted claim-evidence understanding (TBE-2) registered an improvement (with a margin up to 16.39%) over the baselines for the same dataset, (iii) training language models with entailed justifications (TBE-3) outperformed the baselines by a huge margin (up to 28.57% and 44.26% for LIAR-RAW and RAW-FC, respectively). We have shared our code repository to 037 reproduce the results.

1 Introduction:

039

042

The spread of misinformation on the internet has grown to be a pressing social issue. Its consequences have manifested across social, political and commercial domains (Mozur, 2018; Fisher et al., 2016; Allcott and Gentzkow, 2017; Burki, 2019; Aghababaeian et al., 2020). To counter the spread, institutional interventions doing manual fact-checking have gained momentum. For example, the International Fact-Checking Network (IFCN), started in 2015, works with over 170 factchecking organisations and websites worldwide (as of July 2024¹). However, manually verifying facts and detecting misinformation is a slow and costly process. Human experts struggle to keep up with the pace of rapid spread. To overcome this, the research community have been trying to automate the process. Their interest can be gauged by the fact that more than twelve hundred research articles and more than fifty datasets were published on this topic (Alnabhan and Branco, 2024; Guo et al., 2022). A detailed list of such related works is reported in the Appendix A.

043

045

047

055

057

060

061

062

063

064

065

066

067

069

070

071

072

074

075

076

078

081

Doing automated fact-checking at scale is still an open challenge. Most of the proposed approaches relied on language models (Shu et al., 2022; Yang et al., 2022a; Yue et al., 2023; Choi and Ferrara, 2024a). Researchers have either (i) relied on their embedded knowledge or (ii) used evidence collected from various sources to predict the veracity. For example, Pan et al. (2021); Zeng and Gao (2023) reported the performance of zero-shot or *few-shot* prompting for this task. Dhuliawala et al. (2024) showcased that integrating reasoning steps through chain-of-thought prompting reduced model hallucinations. Others relied on external tools such as web search (Galitsky, 2025; Zhang and Gao, 2023) to retrieve the evidence and do fact verification. However, their accuracy is far behind to deploy them for practical purposes. The primary reason behind this is that fact verification is difficult because language models need to look at multiple pieces of information to decide what is true. These

¹https://www.poynter.org/ifcn/

pieces often disagree with each other, which makes the reasoning process even harder. To counter this, 083 Yue et al. (2024) segregated the retrieved evidence into 'supporting' and 'opposing' arguments, allowing the model to consider each perspective independently, and then they verified each claim using 087 few-shot prompting. On a philosophically similar line, Wang et al. (2024a) introduced a module to extract top-k evidence sets stating the claim to be 'true' and 'false' based on attention scores. They produced 'true' and 'false' justifications by promoting language models based on evidence sentences, and further trained the language models to produce veracity by training LLMs. We, on the other hand, proposed a simple yet effective approach, where we relied on (i) the entailment ability of language models to classify if evidence sentences are supporting or refuting a given claim and (ii) the generative ability of language models to produce supporting 100 and refuting justifications. A schematic diagram 101 of our proposed approach is presented in Figure 102 17 (in the Appendix). We believe the simplicity of our approach will allow us to easily deploy our systems in the real-world. Apart from that, we also 105 106 found that a systematic comparison of different prompting and fine-tuning strategies is lacking in the literature. To fill this research gap, we designed 108 different experiments based on three philosophical questions. They are, 110

> ★ R1: "How well do the language models perform when only raw evidence sentences along with the claim are available during training and inferencing?"

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

- ★ R2: "Does training and inferencing with prompted claim-evidence understanding improve the performance of the language models?
- ★ R3: "Can training and inferencing with prompted entailed justifications improve the claim veracity prediction?"

Our key contributions in this work are,

• We conducted three training-based (**TBE**) and four inference-based (**IBEs**) experiments along the line of the research questions. We gave three types of inputs, i.e, along with the claim (i) raw evidence sentences, (ii) overall claim-evidence understanding generated by language models, and (iii) entailed justifications generated by the language models. • We conducted a detailed evaluation of model explanations. We considered the entailed justifications generated by language models as model explanations, and evaluated them using two strategies: (i) checking lexical overlap and semantic matching, and (ii) doing subjective evaluation by language models. While in the former, we used ROUGE (Lin, 2004), BLEU (Papineni et al., 2002), and BERTscores (Zhang et al., 2019), in the latter, we prompted VLLMs to check how informative, accurate, readable, objective, and logical the explanations are.

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

168

169

170

171

172

173

174

175

176

177

178

179

• We conducted an ablation study and a thorough error analysis of our models. In the ablation study, we removed individual supporting and refuting entailed justifications from input samples and trained the veracity prediction models again. This exercise illustrated the importance of individual entailed justifications. In linguistic analysis, on the other hand, we identified (i) the cases where our model succeeds and fails, and (ii) the possible reasons behind it.

Some of the interesting observations we got are,

- While prompting language models (IBEs) with raw evidence sentences couldn't outperform the baselines, training some language models with the same (TBE-1) registered an improvement up to 8.20% in macro-F1 for the RAW-FC dataset. Training Language models with prompted claim-evidence understanding (TBE-2) registered an improvement (with a margin up to 16.39%) over baselines for the RAW-FC dataset. Training language models with entailed justifications (TBE-3) outperformed the baselines by a large margin (up to 28.57% and 44.26% for LIAR-RAW and RAW-FC, respectively).
- Subjective evaluation by language models validated the correlations between veracity prediction and explanation quality in the bestperforming models. For instance, Llamagenerated explanations received highest informativeness, readability, objectivity, and logicality ratings, which correlates with the highest macro-fl it got in veracity prediction.
- In the ablation study, we found that the macro-F1 score dropped when the best perform-

180ing models were trained without supporting181 $(31.48\% \downarrow \text{ in LIAR-RAW and } 12.50\% \downarrow \text{ in}$ 182RAW-FC) and refuting justifications (9.26\% \downarrow)183in LIAR-RAW and $9.09\% \downarrow \text{ in RAW-FC}$ 184Further, discarding both had the maximum185adverse effect (55.56\% \downarrow in LIAR-RAW and18647.73\% \downarrow in RAW-FC).

2 Dataset details:

189 190

191

192

193

194

195

196

198

199

200

In this section, we have reported the details of the datasets considered in our study. We have used the LAIR-RAW and RAW-FC datasets provided by Yang et al. (2022b) in our experiments. While each sample of LAIR-RAW is tagged with one of six labels, samples of RAW-FC is tagged with one of three labels. The list of labels and their distribution in respective datasets are reported in Table 1. The datasets are open-source, i.e., they are Apache 2.0 licensed. A detailed description of the individual datasets along with representative sample is reported in the Appendix B.

Dataset	Classes	Count		
	True (T)	2,021		
	Mostly-true(MT)	2,439		
	Half-true (HT)	2,594		
LAIR-RAW	Barely-true (BT)	2,057		
(Yang et al., 2022b)	False (F)	2,466		
	Pants-fire (PF)	1,013		
	Total	12,590		
	True (T)	695		
RAW-FC	Half-true (HT)	671		
(Yang et al., 2022b)	False (F)	646		
	Total	2,012		

Table 1: Datasets Statistics

3 Experiments:

In this section, we have reported the details of the experiments we conducted as part of this study. We conducted two types of experiments- (i) trainingbased (TBE), and (ii) inference-based (IBE). As their name suggests, in **TBE** approaches, we finetuned various language models to predict the veracity labels based on customised inputs. Particularly, (i) we fine-tuned large language models (LLMs) like RoBERTa (Liu et al., 2019) and XLNet (Yang et al., 2019), and (ii) we fine-tuned very large lan-210 211 guage models (VLLMs) like Mistral (Jiang et al., 2023), Llama (AI@Meta, 2024), Gemma (Team 212 et al., 2024), Qwen (Yang et al., 2024) and Fal-213 con (Almazrouei et al., 2023) models using LoRA (Hu et al., 2022) and LoRA+ (Hayou et al., 2024) 215

adapters. Note that, unlike LLMs, we can not directly fine-tune VLLMs due to their large parameter size and computational constraints. Similarly, in **IBE** approaches, we prompted the considered VLLMs to predict veracity labels without explicitly training them for it. We considered the same set of LLMs and VLLMs in all of our experiments reported in this study. The details of individual experiments in each category are reported in the subsequent subsections. We have also reported the current state-of-the-art models as the baselines. 216

217

218

219

220

221

222

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

260

261

262

264

3.1 Training Based Experiments (TBE):

In this section, we have reported the details of three training-based experiments we have considered in our study. Each of them is based on a unique research philosophy. The details of individual approaches are reported in the subsequent subsections.

3.1.1 TBE-1: Training based on raw-evidences:

In the first experiment, we tried to answer "How well do the language models perform when raw evidence sentences are available during training?". To answer this, we fine-tuned the language models by giving claims and raw evidence sentences as input. For the samples which don't have associated raw-evidence sentences, we gave claims as the only input. Here, we restricted ourselves from finetuning LLMs like RoBERTa (Liu et al., 2019) and XLNet (Yang et al., 2019), as the input length of many samples exceeded the maximum supported input size of these language models. Past work (Cheung and Lam, 2023) demonstrated the effectiveness of VLLM finetuning using evidence pieces from web search, whereas in our case, we used the gold evidence sentences given in the dataset. The emergence of adapter-based training for VLLMs allowed us to do this, as they can be trained with large input texts. To the best of our knowledge, we believe we are the first to fine-tune the VLLMs using raw evidence sentences. The approach we followed in **TBE-1** is illustrated in sub-figure (a) of Figure 1.

3.1.2 TBE-2: Training based on overall understanding:

In the second experiment, we tried to answer "Does training with VLLM-generated claim-evidence understanding improve the performance of the language models?. To conduct this experiment, we



Figure 1: Illustration of steps we followed in different experiments. Sub-figure (a) presents the case where only raw evidence sentences and claims are given as input. This approach is used in **TBE-1** and **IBE-1**. In **TBE-1**, we trained the adapters (with VLLMs) using these inputs, whereas in **IBE-1**, we prompted VLLMs using zero-shot promptings. Sub-figure (b) shows the overall process of **TBE-2**, **IBE-2** and **IBE-3**. In **TBE-2**, the training of LLMs and VLLMs with adapters were done with the help of claim-evidence understandings generated by VLLMs. However, in **IBE-2** and **IBE-3**, we used zero-shot and CoT based prompting for final veracity prediction. Sub-figure (c) illustrates the overall experimental process of **TBE-3** and **IBE-4**. Here, first, we generate the entailment labels ("supporting" or "refuting") for individual evidence sentences for a given claim. We prompted the considered VLLMs to do the same. Subsequently, we clubbed the supporting and refuting evidence sentences and prompted VLLMs to generate justifications supporting and refuting the given claim. Lastly, based on the generated justifications, we generated the claim veracity by (i) training the LLMs and VLLMs with adapters as part of **TBE-3**, and (ii) prompting the VLLMs as a part of **IBE-3**.

first prompted the five considered VLLMs to generate their understanding of a given claim and its evidence sentence set. For the samples which don't have an associated evidence sentence, VLLMs generated their understandings based on the embedded knowledge. Based on the understanding, we finetuned the LLMs (such as RoBERTa (Liu et al., 2019) and XLNet (Yang et al., 2019)) and trained the adapters (LoRA (Hu et al., 2022) and LoRA+ (Hayou et al., 2024)) with the considered VLLMs to produce the claim veracity. The detailed experimental process is illustrated in sub-figure (b) of Figure 1. Some of the prompt samples are presented in the appendix (Figure 18).

265

266

267

270

272

273

274

277

278

279

281

3.1.3 TBE-3: Training based on entailment understanding:

In the third experiment, we tried to answer "Can training with VLLM-generated entailing justifications improve the claim veracity prediction?". To conduct this experiment, we followed a three-step approach. In the first step, we prompted the considered five VLLMs to classify if the evidence sentences are "supporting" or "refuting" a given claim. Our approach is inspired by the L-Defense model proposed by Wang et al. (2024a), where the authors used an attention mechanism to distinguish supporting and refuting evidence sentences. We, however, used the VLLMs for the same. In the second step, we prompted the language models to generate supporting and refuting *justifications* based on the classified evidence sentences. For the cases where claims don't have any supporting or refuting evidence sentences, VLLMs generated justifications based on their embedded knowledge. We separately passed the supporting and refuting evidence sentences associated with a claim in the prompts to generate them. Finally, based on the claims and the generated justifications, we (i) fine-tuned the LLMs (such as RoBERTa (Liu et al., 2019) and XLNet (Yang et al., 2019)), and (ii) fine-tuned the adapters (LoRA (Hu et al., 2022) and LoRA+ (Hayou et al., 2024)) with considered VLLMs to generate the ve-

286

287

288

289

290

291

292

294

295

297

298

299

300

301

302

303

304

305

racity labels. The detailed approach is illustrated in 307 sub-figure (c) of Figure 1. Some prompt samples 308 we used at each step are illustrated in Figure 19, 309 Figure 20, and Figure 21 in the appendix.

3.2 **Inference Based Experiments (IBE):**

311

312

313

314

315

317

319

320

321

323

325

326

330

331 332

334

337

We have experimented with four types of inferencebased approaches, each based on a unique prompting philosophy. They are (i) zero-shot prompting (**IBE-1**), where we have prompted five considered VLLMs to predict the veracity label given a claim and its associated evidence sentences, (ii) zeroshot prompting with overall understanding (IBE-2), where, first, we prompted the five considered VLLMs to generate the claim-evidence understanding and then, we prompted them again to with the understanding to predict the veracity labels, (iii) CoT prompting with overall understanding (IBE-3), where we followed similar steps as mentioned in **IBE-2**, except, additionally, we asked the VLLMs to generate step-by-step reasoning behind their predictions, and (iv) prompting based on entailment (IBE-4), where we prompted the VLLMs to pre-328 dict veracity labels based on entailed justifications. Due to space constraints, we have reported the details of individual approaches in the appendix. For the cases where claims don't have associated evidence sentences, we gave the claim as only input (in TBE-1), generated understandings and justifications based on the embedded knowledge of VLLMs.

Baselines: 3.3

338 In this section, we have reported the previously proposed best-performing models as the baselines. 339 Particularly, we compared our models with the per-340 formances of HiSS (Zhang and Gao, 2023), FactL-341 LaMa (Cheung and Lam, 2023), RAFTS (Yue et al., 342 2024), and L-Defence (Wang et al., 2024a) mod-343 els. Out of them, HiSS (Zhang and Gao, 2023) and 344 FactLLaMa (Cheung and Lam, 2023) retrieve ev-345 idence from external sources, while RAFTS (Yue et al., 2024) employs a coarse-to-fine retrieval tech-347 nique to extract evidence directly from the dataset. In contrast, L-Defense (Wang et al., 2024a) used relevant evidence without additional retrieval. The 351 previously reported performance of these models on LIAR-RAW and RAW-FC datasets are presented in Table 2. A detailed description of the individual models is presented in the Appendix C.5. 355

Method	L	IAR-RAV	N	RAWFC					
memou	MP	MR	MF1	MP	MR	MF1			
HiSS	0.46	0.31	0.37	0.53	0.54	0.53			
FactLLaMA	0.32	0.32	0.30	0.56	0.55	0.55			
RAFTS	0.47	0.37	0.42	0.62	0.52	0.57			
L-Defense									
-ChatGPT	0.30	0.32	0.30	0.61	0.61	0.61			
-Llama2	0.31	0.31	0.31	0.61	0.60	0.60			
	(0.29^{\dagger})	(0.29^{\dagger})	(0.29^{\dagger})	(0.56^{\dagger})	(0.56^{\dagger})	(0.56^{\dagger})			

Table 2: Performance of the considered baseline methods on the LIAR-RAW and RAWFC datasets. Notation: our reproduced results are marked as '†'.

$\text{Dataset} (\rightarrow)$	L	JAR-RAV	N	RAW-FC					
		TBE-1			TBE-1				
Method (\downarrow)	MP	MR	MF1	MP	MR	MF1			
LORA									
-Mistral	0.44	<u>0.29</u>	0.27	<u>0.69</u>	<u>0.65</u>	<u>0.65</u>			
	(±0.02)	(±0.01)	(±0.01)	(±0.01)	(±0.00)	(±0.01)			
-Llama	0.34	0.30	0.30	0.68	0.64	0.65			
	(±0.01)	(±0.02)	(±0.00)	(±0.01)	(±0.02)	(±0.02)			
-Gemma	0.27	0.25	0.23	0.60	0.58	0.57			
	(±0.02)	(±0.03)	(±0.04)	(±0.05)	(±0.04)	(±0.06)			
-Qwen	0.37	<u>0.29</u>	<u>0.29</u>	0.67	0.66	0.66			
	(±0.02)	(±0.02)	(±0.04)	(±0.03)	(±0.03)	(±0.03)			
-Falcon	0.39	0.39 0.28		0.59	0.58	0.54			
	(±0.01)	(±0.01)	(±0.01)	(±0.03)	(±0.05)	(±0.06)			
LORA+									
-Mistral	0.40	0.27	0.25	0.53	0.56	0.55			
	(±0.03)	(±0.02)	(±0.01)	(±0.05)	(±0.03)	(±0.02)			
-Llama	0.34	0.29	0.29	0.67	0.64	0.65			
	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)			
-Gemma	0.27	0.23	0.22	0.61	0.57	0.57			
	(±0.02)	(±0.02)	(±0.03)	(±0.03)	(±0.02)	(±0.03)			
-Qwen	0.36	<u>0.29</u>	<u>0.29</u>	0.70	0.65	<u>0.65</u>			
	(±0.01)	(±0.01)	(±0.02)	(±0.02)	(±0.02)	(±0.02)			
-Falcon	0.37	0.30	<u>0.29</u>	0.64	0.62	0.63			
	(±0.02)	(±0.01)	(±0.02)	(±0.03)	(±0.03)	(±0.03)			

Table 3: Performance of LoRA and LoRA+ models on TBE-1 (LIAR-RAW and RAW-FC datasets).

4 **Results and Discussion:**

In this section, we have reported the results from all of our experiments. While we used macroprecision, macro-recall and macro-F1 scores to evaluate the veracity, ROUGE, BLEU and BERTscore, along with subjective evaluation by VLLMs, were made to evaluate the explanations. A detailed description of evaluation metrics is reported in subsection C.6.4 in the appendix.

357

358

359

360

361

362

363

364

4.1 **Observations from TBE-1:**

We reported the performance of various models for 366 **TBE-1**, i.e. training with raw evidence sentences 367 as input in Table 3. We found that LoRA adapter 368 trained with Llama VLLM resulted in the highest 369 F1 score of 0.30 for the LIAR-RAW dataset. While 370 this performance is comparable to the performance 371 of baseline models L-Defense (F1: 0.31) and 372

FactLLaMA (F1: 0.30), it does not surpass the per-373 formance of other baselines, i.e. HiSS (F1: 0.37) 374 and RAFTS (F1: 0.42). However, for the RAWFC dataset, we found that many models, such as Mistral with LoRA (F1: 0.65, \sim 6.56% \uparrow), Llama with LoRA (F1: 0.65, \sim 6.56% \uparrow), Qwen with LoRA (F1: 0.66, \sim 8.20% \uparrow), Llama with LoRA+ 379 (*F*1: **0.65**, \sim **6.56**% \uparrow), Qwen with LoRA+ (*F*1: **0.65**, \sim **6.56**% \uparrow) and Falcon with LoRA+ (*F*1: **0.63**, \sim **3.28**% \uparrow), outperform the best baseline performance (F1: 0.61). Here, $\sim x\%$ \uparrow denotes the relative percentage improvement compared to the best baseline score. Qwen trained with LoRA resulted in the overall highest F1-score (0.66, \sim 8.20% \uparrow). We have reported some additional observations in the appendix section D.1.

4.2 Observations from TBE-2:

389

390

400

401

402

403

404

405

406

407

408

409

410

411

We reported the performance of various models for TBE-2, i.e. training with claim-evidence understanding as input in table 4. We found that Mistral trained with LoRA resulted in the highest macro-F1 (0.32) for the LIAR-RAW dataset. While its performance surpasses the performance of baselines FactLLaMA (MF1 : **0.30**) and L-Defense (MF1 : 0.31), it is still behind HiSS (MF1 : 0.37) and RAFTS (MF1 : 0.42). However, for the RAW-FC dataset, we observed that many models, such as XLNet fine-tuned on Llama based understandings ($MF1 : 0.62, \sim 1.64\% \uparrow$), and Llama trained (with Llama understandings) with LoRA+ (*MF*1 : **0.71**, \sim **16.39**% \uparrow) outperformed the best reported macro-F1 score by the baselines (MF1 : 0.61). Llama trained (with Llama understandings) with LoRA+ adapter resulted in the highest overall performance (MF1: 0.71, \sim 16.39% \uparrow). We have reported some additional observations in the section D.2 of the appendix due to space constraints.

4.3 Observations from TBE-3:

We reported the performance of various models 412 for **TBE-3**, i.e. training with entailed justifications 413 generated by VLLMs as input in table 4. We found 414 that XLNet fine-tuned on Llama based entailed jus-415 tification resulted in the highest macro-F1 (0.54, 416 \sim 28.57% \uparrow) for the LIAR-RAW dataset. Many 417 418 models, such as RoBERTa fine-tuned with Mistral (*MF*1 : **0.47**, \sim **16.39**% \uparrow), Llama (*MF*1 : 419 **0.52**, \sim **23.81**% \uparrow), Gemma (*MF*1 : **0.48**, \sim 420 **14.29**% \uparrow), Qwen (*MF*1 : **0.46**, ~ **9.52**% \uparrow), 421 and Falcon (MF1 : **0.44**, ~ **4.76**% \uparrow) based en-422

tailed justifications, XLNet fine-tuned with Mistral 423 $(MF1: 0.47, \sim 16.39\% \uparrow)$, Llama (MF1: 0.54,424 ~ 28.57% \uparrow), Owen (*MF*1 : 0.48, ~ 14.29% \uparrow), 425 and Falcon (MF1 : **0.44**, ~ **4.76**% \uparrow) based en-426 tailed justifications, and Llama trained (with Llama 427 justifications) with LoRA+ adapter (MF1 : **0.49**, 428 \sim 16.67% \uparrow) surpassed the best reported macro-429 F1 score of baselines (MF1 : 0.42). In a similar 430 manner, for the RAW-FC dataset, we observed that 431 many models, such as such as RoBERTa fine-tuned 432 with Mistral (*MF*1 : **0.83**, \sim **36.07**% \uparrow), Llama 433 $(MF1: 0.88, \sim 44.26\% \uparrow)$, Qwen (MF1: 0.71,434 ~ 16.39% \uparrow), and Falcon (MF1 : 0.64, ~ 435 **4.91**% \uparrow) based entailed justifications, XLNet fine-436 tuned with Mistral (MF1 : **0.82**, ~ **34.42**% \uparrow), 437 Llama (MF1 : **0.87**, ~ **42.62**% \uparrow), Qwen (MF1 : 438 **0.70**, ~ 14.75% \uparrow), and Falcon (*MF*1 : 0.74, ~ 439 **21.31**% \uparrow) based entailed justifications, and Llama 440 trained (with Llama justifications) with LoRA+ 441 adapter (MF1 : **0.83**, ~ **36.07**% \uparrow) outperformed 442 the best reported macro-F1 score by the baselines 443 (MF1: 0.61). RoBERTa fine-tuned with Mistral 444 based entailed justification achieved the highest 445 overall performance (MF1 : **0.88**, ~ **44.26**% \uparrow). 446 We have reported some additional observations in 447 the Appendix section D.3. 448

4.4 Observations from IBEs':

We reported the performance of various models in 450 **IBEs** in Table 11. For LIAR-RAW dataset, Mis-451 tral achieved the highest macro-F1 (0.22) in IBE-1. 452 While in IBE-2, Llama achieved the highest perfor-453 mance (MF1 : 0.22), Mistral, Llama, Qwen and 454 Falcon (MF1 : 0.21) gave the highest performance 455 in IBE-3. In IBE-4, Mistral (MF1 : 0.14) and Fal-456 con (MF1 : 0.14) got the highest performance. 457 None of the macro-F1 could surpass the baseline 458 performances. In contrast, for RAW-FC dataset, 459 Llama (MF1 : **0.62**) achieved the highest perfor-460 mace in IBE-2, surpassing the highest baseline per-461 formance (MF1 : 0.61). While Qwen attained the 462 highest performance in IBE-1 (MF1 : 0.59), sur-463 passing the baselines HiSS (MF1 : 0.53), FactL-464 LaMa (*MF*1 : **0.55**) and RAFTS (*MF*1 : **0.57**), 465 it is still behind the performance of L-Defense 466 (MF1 : 0.61). Similarly, Qwen got the highest 467 performance in IBE-3 (MF1 : 0.52). In IBE-4, 468 Mistral and Qwen achieved the highest macro-F1 469 of 0.43, which is far behind all baselines. We have 470 reported some additional observations in the Ap-471 pendix (section D.4) 472

Dataset (\rightarrow)			LIAR	-RAW					RAV	V-FC		
Method (\downarrow)		TBE-2			TBE-3		-	TBE-2			TBE-3	
	MP	MR	MF1									
FINE-TUNING												
-RoBERTa-L _{Mistral}	0.28	0.26	0.26	0.48	0.47	0.47	0.51	0.50	0.50	0.83	0.82	0.83
	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.00)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.00)	(±0.01)	(±0.01)
-RoBERTa-L _{Llama}	0.27	0.28	0.25	0.53	0.53	0.52	0.50	0.50	0.49	0.88	0.88	0.88
	(±0.02)	(±0.01)	(±0.02)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)
-RoBERTa-L _{Gemma}	0.28	0.27	0.27	0.49	0.50	0.48	0.51	0.51	0.50	0.50	0.49	0.49
	(±0.02)	(±0.01)	(±0.01)	(±0.01)	(±0.02)	(±0.02)	(±0.04)	(±0.03)	(±0.04)	(±0.02)	(±0.01)	(±0.02)
-RoBERTa-L _{Qwen}	0.30	0.28	0.28	0.48	0.47	0.46	0.51	0.49	0.48	0.73	0.71	0.71
	(± 0.01)	(± 0.02)	(±0.01)	(±0.01)	(± 0.02)	(±0.01)	(±0.01)	(±0.01)	(±0.02)	(±0.02)	(± 0.02)	(±0.02)
-RoBERTa-L _{Falcon}	0.29	0.28	0.27	0.49	0.43	0.44	0.50	0.48	0.48	0.65	0.64	0.64
	(± 0.02)	(± 0.02)	(± 0.01)	(± 0.01)	(± 0.02)	(± 0.01)	(± 0.00)	(± 0.01)	(± 0.01)	(± 0.02)	(± 0.02)	(± 0.02)
-XLINet-L _{Mistral}	0.29	0.28	0.28	0.49	0.47	0.47	0.62	0.01	0.01	0.85	0.82	0.82
VIN-4 I	(± 0.02)	(± 0.01)	(± 0.01)	(±0.00)	(± 0.01)	(± 0.01)	(± 0.02)	(± 0.01)	(± 0.02)	(±0.01)	(±0.01)	(± 0.01)
-ALNEI-LLIama	0.51	0.29	0.29	0.55	0.54	0.54	0.05	(10.02)	(10.02)	0.00	0.00	(10.01)
VI Not I	(± 0.02)	(± 0.01)	(±0.01)	(± 0.01)	(± 0.02)	(± 0.01)	(±0.03)	(± 0.03)	(± 0.03)	(± 0.01)	(± 0.01)	(±0.01)
-ALIVEI-LGemma	(+0.02)	(+0.02)	(+0.02)	(+0.04)	(+0.05)	(+0.09)	(+0.01)	(+0.01)	(+0.01)	(+0.02)	(+0.02)	(+0.02)
-XINet-Lower	0.29	0.29	0.28	0.50	0.48	0.48	0.60	0.58	0.58	0.70	0.70	0.70
ALLIVET LOWER	(± 0.02)	(± 0.02)	(± 0.02)	(± 0.01)	(± 0.01)	(± 0.01)	(± 0.02)	(± 0.01)	(± 0.01)	(± 0.03)	(± 0.03)	(± 0.04)
-XLNet-L _{Falcon}	0.26	0.26	0.24	0.45	0.44	0.44	0.61	0.60	0.60	0.76	0.75	0.74
1 MCON	(±0.03)	(±0.01)	(±0.03)	(±0.01)	(±0.02)	(±0.01)	(±0.02)	(±0.01)	(±0.01)	(±0.01)	(±0.01)	(±0.01)
LORA												
-Mistral	0.36	0.32	0.32	0.35	0.34	0.29	0.61	0.60	0.60	0.69	0.59	0.58
	(±0.01)	(±0.01)	(±0.01)	(±0.04)	(±0.03)	(±0.03)	(±0.04)	(±0.04)	(±0.04)	(±0.05)	(±0.05)	(±0.05)
-Llama	0.31	0.27	0.25	0.36	0.32	0.30	0.48	0.48	0.46	0.63	0.62	0.61
	(±0.04)	(±0.03)	(±0.04)	(±0.03)	(±0.03)	(±0.04)	(±0.06)	(±0.03)	(±0.02)	(±0.03)	(±0.03)	(±0.03)
-Gemma	0.21	0.22	0.18	0.26	0.23	0.20	0.47	0.45	0.40	0.36	0.34	0.30
	(±0.01)	(±0.01)	(±0.01)	(±0.03)	(±0.03)	(±0.02)	(±0.03)	(±0.03)	(±0.05)	(±0.02)	(±0.01)	(±0.02)
-Qwen	0.32	0.27	0.21	0.40	0.36	0.32	0.56	0.54	0.53	0.48	0.46	0.46
	(±0.09)	(±0.02)	(±0.02)	(±0.07)	(±0.07)	(±0.05)	(±0.07)	(±0.0)	(±0.03)	(±0.03)	(±0.03)	(±0.03)
-Falcon	0.28	0.26	0.23	0.27	0.24	0.21	0.56	0.57	0.53	0.43	0.0.40	0.41
	(± 0.01)	(± 0.00)	(± 0.01)	(± 0.03)	(± 0.03)	(± 0.02)	(± 0.03)	(± 0.01)	(± 0.00)	(± 0.01)	(± 0.03)	(± 0.03)
LoRA+												
-Mistral	<u>0.33</u>	0.30	<u>0.30</u>	0.30	0.32	0.27	<u>0.66</u>	0.60	0.60	0.52	0.48	0.46
	(±0.01)	(±0.01)	(±0.01)	(±0.07)	(±0.02)	(±0.03)	(±0.05)	(±0.05)	(±0.05)	(±0.06)	(±0.04)	(±0.04)
-Llama	0.27	0.24	0.23	0.49	0.50	0.49	0.73	0.71	0.71	<u>0.84</u>	0.82	0.82
	(±0.01)	(±0.02)	(±0.02)	(±0.03)	(±0.02)	(±0.02)	(±0.02)	(±0.02)	(±0.02)	(±0.01)	(±0.03)	(±0.03)
-Gemma	0.23	0.21	0.20	0.34	0.32	029	0.53	0.51	0.51	0.52	0.48	0.46
0	(± 0.03)	(± 0.02)	(±0.02)	(± 0.04)	(±0.05)	(± 0.06)	(± 0.02)	(± 0.02)	(± 0.04)	(± 0.03)	(± 0.04)	(±0.01)
-Qwen	0.27	0.27	0.24	0.33	0.34	0.29	0.57	0.48	0.43	0.52	0.49	0.4^{7}
Estern	(± 0.04)	(± 0.03)	(±0.04)	(± 0.03)	(± 0.03)	(± 0.03)	(±0.07)	(±0.03)	(±0.01)	(±0.00)	(±0.09)	(±0.09))
-raicon	<u>U.33</u>	(10.01)	0.28	0.39	0.39	0.30	0.55	0.50	0.55	0.58	0.52	0.50
	(± 0.00)	(± 0.01)	(± 0.02)	(± 0.02)	(± 0.03)	(± 0.04)	(± 0.03)	(± 0.02)	(± 0.03)	(± 0.01)	(± 0.02)	(± 0.03)

Table 4: Performance of claim veracity prediction using gold evidences. Green and Blue indicates best and second-best performance, respectively.

4.5 **Comparing the performances across TBEs** and IBEs:

473

474

477

478

479

480

481

482

483

484

475 In this section, we have reported our findings by comparing the best-performing models across all 476 TBEs and IBEs. The comparison is also illustrated in Figure 22, Figure 23 and Figure 24 in the Appendix D.5. For the LIAR-RAW dataset, we found that none of the models from any IBEs, TBE-1 and TBE-2 (MF1: 0.14 - 0.32) could surpass the highest F1 achieved in baselines (MF1: 0.42). However, we observed that TBE-1 and TBE-2 (*MF*1: **0.30** - **0.32**, \sim **45.45**% \uparrow) outperformed

IBEs (MF1: 0.22) in terms of the highest reported 485 macro-F1. Further, we also found that many mod-486 els from TBE-3 (*MF*1: **0.44** - **0.54**, \sim **68.75**% \uparrow) 487 outperform all of the models of IBEs (MF1: 0.14 488 - 0.22), TBE-1 (MF1: 0.30) and TBE-2 (MF1: 489 **0.32**). For the RAW-FC dataset, we observed a 490 similar comparative performance. However, here, 491 the highest reported macro-F1 by IBE-2 (0.62) is 492 comparable to the best-performing baseline (MF1:493 0.61). Similar to the LIAR-RAW dataset, here also 494 we observed that (i) many TBE-1 and TBE-2 mod-495 els (*MF*1: **0.62** - **0.71**, \sim **14.52**% \uparrow) outperform 496 IBEs (0.62) in terms of macro-F1, and (ii) many 497

Dataset (\rightarrow)		LIAR-RAW									RAW-FC													
Method (\downarrow)	IBE-1 IBE-2			IBE-3			IBE-4		IBE-1		IBE-2				IBE-3		IBE-4							
	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1	MP	MR	MF1
PROMPTING																								
-Mistral	0.24	0.25	0.22	0.22	0.23	0.20	0.40	0.23	0.21	0.30	0.17	0.14	0.54	0.54	0.53	0.58	0.58	0.58	0.50	0.52	0.45	0.45	0.46	0.43
-Llama	0.24	<u>0.23</u>	<u>0.20</u>	0.30	0.23	0.22	0.26	<u>0.22</u>	0.21	<u>0.22</u>	0.13	<u>0.13</u>	<u>0.56</u>	0.56	<u>0.54</u>	0.62	0.63	0.62	0.45	0.47	<u>0.49</u>	0.42	<u>0.39</u>	0.35
-Gemma	0.24	0.21	0.19	0.16	0.16	0.13	0.27	0.19	0.16	0.09	0.16	0.11	0.40	0.40	0.40	0.41	0.38	0.38	0.45	0.41	0.40	0.27	0.31	0.24
-Qwen	0.27	<u>0.23</u>	<u>0.20</u>	<u>0.25</u>	0.23	<u>0.20</u>	0.24	<u>0.22</u>	0.21	0.13	<u>0.16</u>	<u>0.13</u>	0.61	0.58	0.59	0.58	0.57	0.57	<u>0.57</u>	0.54	0.52	0.50	0.46	0.43
-Falcon	<u>0.24</u>	<u>0.23</u>	<u>0.20</u>	0.22	<u>0.22</u>	<u>0.20</u>	0.24	0.23	0.21	0.16	0.17	0.14	<u>0.56</u>	<u>0.57</u>	<u>0.54</u>	<u>0.60</u>	<u>0.59</u>	0.57	0.60	<u>0.52</u>	0.48	0.40	0.38	<u>0.37</u>

Table 5: Performance of Prompting methods across IBE-1, IBE-2, and IBE-3 settings for LIAR-RAW and RAW-FC datasets.

TBE-3 models (*MF*1: **0.64** - **0.88**, $\sim 23.94\%$ \uparrow) outperform all of the models of IBEs (*MF*1: **0.43** - **0.62**), TBE-1 (*MF*1: **0.65** - **0.66**) and TBE-2 (*MF*1: **0.62** - **0.71**). The whole comparison points to the fact that training the LLMs with VLLMgenerated entailed-justifications can improve the performance of veracity prediction. We have reported some additional observations in the section D.5 of the appendix due to space constraints. Further, we have also presented the confusion metrics of best-performing TBE and IBE models for LIAR-RAW and RAW-FC datasets in Figure 37 and 38 respectively.

498

499

500

501

502

504

506

507

508

510

511

4.6 Insights from evaluation of explanations:

In this section, we reported our findings from eval-512 uating model explanations. We obtained them by 513 concatenating the supporting and refuting justifi-514 cations. The details of the metrics and evaluation 515 scores are provided in Appendix C.6.4 and Ap-516 pendix D.6, respectively. We observed that lexical 517 overlap measures like ROUGE (R_1, R_2, R_L) pro-518 vided mixed signals. For instance, while Falcon 519 520 generated explanations showed higher unigrambased overlaps (R_1 : **0.23** for LIAR-RAW and R_1 : 521 0.40 for RAW-FC), Mistral generated explanations 522 got the highest overlap in longest common subse-523 quences (R_L : 0.14 for LIAR-RAW and R_L : 0.18 524 for RAW-FC). However, all models scored consis-525 tently low on the BLEU score and the BERT score. 526 On the other hand, subjective evaluations indicated 527 that Llama-generated explanations were rated highest by the majority of language models for dimen-529 sions like informativeness, readability, objectivity, and logicality for both datasets. These high sub-531 jective ratings seem to be correlated with better 533 veracity prediction performance, as RoBERTa and XLNet performed well upon taking Llama expla-534 nations. A more detailed discussion on the obser-535 vations of 'evaluation of explanations' is presented in Appendix D.6. 537

5 Conclusion:

• Training language models with VLLM entailed justifications surpassed the baseline macro-F1 scores substantially with an improvement of **28.57%** and **44.26%** for LIAR-RAW and RAW-FC, respectively. The approach of training with claim-evidence understanding (TBE-2) secured the second spot, with an increment of **16.39%** in the RAW-FC dataset compared to the best baseline macro-F1. In contrast, the inference-based methods (IBEs) were unable to understand the justifications generated from VLLM and performed consistently poorly. 538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

- While lexical overlap and semantic matching methods showed no definite pattern, the subjective evaluation of model explanations by VLLMs found that Llama generated model explanations are more informative, readable, objective, and logical. It correlates with the superior performance reported by XLNet and RoBERTa (in **TBE-3**) for veracity prediction as they took Llama generated justifications as input.
- The role of VLLM entailed justification as a second step in TBE-3 is justified in the ablation study (see Appendix E) where we observed that removing supporting and refuting justification adversarially affected the scores.
- In the linguistic analysis of model explanations, we found that LLM could attend to supporting and refuting keywords and factual pieces of information for samples labelled with 'true' and 'false' veracity samples. However, for samples with other veracity labels LLM attention seem to be scattered. A detailed discussion is presented in Appendix F due to space constraints.

6

work.

Limitations:

isation.

References

88:29-32.

guistics.

In this section, we reported the limitations of our

• We restricted our experiments to using only

open-source language models, for repro-

ducibility and resource constraints. However,

expanding these experiments to commercial

language models will further generalize the

• Due to our limited linguistic expertise, we

restricted our experiments to an Englishlanguage setup only. In the future, our hy-

pothesis can also be tested in other languages.

future, one can consider open-domain fact-

checking, i.e., retrieve evidence from an exter-

nal source and test our hypothesis for general-

• Due to space constraints, we restricted our

experiments to the task of fact-checking and

utilised two popular datasets. A wider testing

of our hypothesis on various other datasets

could help us know the applicability of our

• Due to resource constraints, we could not eval-

uate the model-generated explanations manu-

ally. One can extend the study by integrating

human evaluation of the explanations.

Hamidreza Aghababaeian, Lara Hamdanieh, and Abbas Ostadtaghizadeh. 2020. Alcohol intake in an attempt

AI@Meta. 2024. Llama 3. https://github.com/

Mubashara Akhtar, Michael Schlichtkrull, Zhijiang Guo,

Oana Cocarascu, Elena Simperl, and Andreas Vla-

chos. 2023a. Multimodal automated fact-checking:

A survey. In Findings of the Association for Com-

putational Linguistics: EMNLP 2023, pages 5430-

5448, Singapore. Association for Computational Lin-

Mubashara Akhtar, Michael Schlichtkrull, Zhijiang Guo,

A survey. arXiv preprint arXiv:2305.13507.

Oana Cocarascu, Elena Simperl, and Andreas Vla-

chos. 2023b. Multimodal automated fact-checking:

meta-llama/llama3. Accessed: 2025-04-21.

to fight covid-19: A medical myth in iran. Alcohol,

idea of utilising entailment.

• In our work, we assumed a closed-domain fact-checking setup for reproducibility. In

idea of claim-evidence entailment.

- 577
- 578
- 579
- 5
- 581
- 5
- 584
- 585
- Ę
- 589 590
- 591
- ļ
- 595 596
- 597
- 598 599

60

- 6
- 6

6

6

6

6

610 611

612 613

614 615

615 616 617

619 620

621 622

618

Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. Where is your evidence: Improving factchecking by justification modeling. In *Proceedings* of the first workshop on fact extraction and verification (FEVER), pages 85–90.

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

- Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives*, 31(2):211–236.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. 2023. The falcon series of open language models.
- Mohammad Q Alnabhan and Paula Branco. 2024. Fake news detection using deep learning: A systematic literature review. *IEEE Access*.
- Rami Aly, Zhijiang Guo, Michael Sejr Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. The fact extraction and VERification over unstructured and structured information (FEVEROUS) shared task. In *Proceedings of the Fourth Workshop on Fact Extraction and VERification (FEVER)*, pages 1–13, Dominican Republic. Association for Computational Linguistics.
- Irene Amerini, Leonardo Galteri, Roberto Caldelli, and Alberto Del Bimbo. 2019. Deepfake video detection through optical flow based cnn. In 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW), pages 1205–1207.
- Ramy Baly, Mitra Mohtarami, James Glass, Lluís Màrquez, Alessandro Moschitti, and Preslav Nakov. 2018. Integrating stance detection and fact checking in a unified corpus. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 21–27, New Orleans, Louisiana. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Talha Burki. 2019. Vaccine misinformation and social media. *The Lancet Digital Health*, 1(6):e258–e259.
- Recep Firat Cekinel, Pinar Karagoz, and Çağrı Çöltekin. 2025. Multimodal fact-checking with vision language models: A probing classifier based solution with embedding strategies. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 4622–4633, Abu Dhabi, UAE. Association for Computational Linguistics.
- 9

735

Abhijnan Chakraborty, Bhargavi Paranjape, Sourya Kakarla, and Niloy Ganguly. 2016. Stop clickbait: detecting and preventing clickbaits in online news media. In Proceedings of the 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM '16, page 9-16. IEEE Press.

678

679

691

701

705

710

712

713

714

715

716

717

718

719

720

721

723

724

725

726

727

728

729

730

731

733

- Tsun-Hin Cheung and Kin-Man Lam. 2023. Factllama: Optimizing instruction-following language models with external knowledge for automated fact-checking. In 2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), pages 846-853. IEEE.
- Eun Cheol Choi and Emilio Ferrara. 2024a. Fact-gpt: Fact-checking augmentation via claim matching with llms. In Companion Proceedings of the ACM Web Conference 2024, pages 883-886.
 - Eun Cheol Choi and Emilio Ferrara. 2024b. Fact-gpt: Fact-checking augmentation via claim matching with llms. In Companion Proceedings of the ACM Web Conference 2024, WWW '24, page 883-886, New York, NY, USA. Association for Computing Machinery.
 - IDO DAGAN, BILL DOLAN, BERNARDO MAGNINI, and DAN ROTH. 2010. Recognizing textual entailment: Rational, evaluation and approaches - erratum. Natural Language Engineering, 16(1):105-105.
 - Shih-Chieh Dai, Yi-Li Hsu, Aiping Xiong, and Lun-Wei Ku. 2022. Ask to know more: Generating counterfactual explanations for fake claims. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '22, page 2800-2810, New York, NY, USA. Association for Computing Machinery.
 - Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2024. Chain-of-verification reduces hallucination in large language models. In Findings of the Association for Computational Linguistics: ACL 2024, pages 3563-3578, Bangkok, Thailand. Association for Computational Linguistics.
 - John Dougrez-Lewis, Elena Kochkina, Miguel Arana-Catania, Maria Liakata, and Yulan He. 2022. PHE-MEPlus: Enriching social media rumour verification with external evidence. In Proceedings of the Fifth Fact Extraction and VERification Workshop (FEVER), pages 49-58, Dublin, Ireland. Association for Computational Linguistics.
- Islam Eldifrawi, Shengrui Wang, and Amine Trabelsi. 2024. Automated justification production for claim veracity in fact checking: A survey on architectures and approaches. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6679-6692, Bangkok, Thailand. Association for Computational Linguistics.

- Marc Fisher, John Woodrow Cox, and Peter Hermann. 2016. Pizzagate: From rumor, to hashtag, to gunfire in dc. Washington Post, 6:8410-8415.
- Mohamed H. Gad-Elrab, Daria Stepanova, Jacopo Urbani, and Gerhard Weikum. 2019. Exfakt: A framework for explaining facts over knowledge graphs and text. In Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM '19, page 87–95, New York, NY, USA. Association for Computing Machinery.
- Boris Galitsky. 2025. 8 truth-o-meter: Collaborating with llm in fighting its hallucinations. In William Lawless, Ranjeev Mittu, Donald Sofge, and Hesham Fouad, editors, Interdependent Human-Machine Teams, pages 175–210. Academic Press.
- Genevieve Gorrell, Elena Kochkina, Maria Liakata, Ahmet Aker, Arkaitz Zubiaga, Kalina Bontcheva, and Leon Derczynski. 2019. SemEval-2019 task 7: RumourEval, determining rumour veracity and support for rumours. In Proceedings of the 13th International Workshop on Semantic Evaluation, pages 845–854, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. Transactions of the Association for Computational Linguistics, 10:178–206.
- Ashim Gupta and Vivek Srikumar. 2021. X-fact: A new benchmark dataset for multilingual fact checking. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 675–682, Online. Association for Computational Linguistics.
- Vipin Gupta, Rina Kumari, Nischal Ashok, Tirthankar Ghosal, and Asif Ekbal. 2022. MMM: An emotion and novelty-aware approach for multilingual multimodal misinformation detection. In Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022, pages 464-477, Online only. Association for Computational Linguistics.
- David Güera and Edward J. Delp. 2018. Deepfake video detection using recurrent neural networks. In 2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), pages 1-6.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. Lora+: Efficient low rank adaptation of large models. In International Conference on Machine Learning, pages 17783-17806. PMLR.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In ICLR 2022.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego

901

902

903

847

de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

791

792

795

796

804

805

810

811

812

813

814

816

817

818

819

820

821

822

831

833

834

835

836

837

840

841

842

- Neema Kotonya and Francesca Toni. 2020a. Explainable automated fact-checking: A survey. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5430–5443, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Neema Kotonya and Francesca Toni. 2020b. Explainable automated fact-checking for public health claims. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7740–7754, Online. Association for Computational Linguistics.
 - Sujit Kumar, Anshul Sharma, Siddharth Hemant Khincha, Gargi Shroff, Sanasam Ranbir Singh, and Rahul Mishra. 2025. Sciclaimhunt: A large dataset for evidence-based scientific claim verification. *arXiv preprint arXiv:2502.10003*.
 - Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1906–1919, Online. Association for Computational Linguistics.
- Paul Mozur. 2018. A genocide incited on facebook, with posts from myanmar's military. *The New York Times*, 15(10):2018.
- Kai Nakamura, Sharon Levy, and William Yang Wang. 2020. Fakeddit: A new multimodal benchmark dataset for fine-grained fake news detection. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6149–6157, Marseille, France. European Language Resources Association.
- Dan S. Nielsen and Ryan McConville. 2022. Mumin: A large-scale multilingual multimodal fact-checked misinformation social network dataset. In Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22, page 3141–3153, New York, NY, USA. Association for Computing Machinery.
- Liangming Pan, Wenhu Chen, Wenhan Xiong, Min-Yen Kan, and William Yang Wang. 2021. Zero-shot fact verification by claim generation. In *Proceedings*

of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 476–483, Online. Association for Computational Linguistics.

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Jasabanta Patro and Sabyasachee Baruah. 2021. A simple three-step approach for the automatic detection of exaggerated statements in health science news. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3293–3305, Online. Association for Computational Linguistics.
- Jasabanta Patro and Pushpendra Singh Rathore. 2020. A sociolinguistic route to the characterization and detection of the credibility of events on twitter. In *Proceedings of the 31st ACM Conference on Hypertext and Social Media*, HT '20, page 241–250, New York, NY, USA. Association for Computing Machinery.
- Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018. DeClarE: Debunking fake news and false claims using evidence-aware deep learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 22–32, Brussels, Belgium. Association for Computational Linguistics.
- Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C. Lipton. 2020. Learning to deceive with attention-based explanations. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4782– 4793, Online. Association for Computational Linguistics.
- Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. 2019. Recurrent convolutional strategies for face manipulation detection in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops.*
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin C! robust fact verification with contrastive evidence. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643, Online. Association for Computational Linguistics.
- Gautam Kishore Shahi and Durgesh Nandini. 2020. Fakecovid- A multilingual cross-domain fact check news dataset for COVID-19. In Workshop Proceedings of the 14th International AAAI Conference on

- 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1001 1002 1003 1004

1007

1008

1009

1010

1011

1012

904

905

906

907

909

910

911

- 951 952
- 954 955

960

Web and Social Media, ICWSM 2020 Workshops, Atlanta, Georgia, USA [virtual], June 8, 2020.

- Kai Shu, Ahmadreza Mosallanezhad, and Huan Liu. 2022. Cross-domain fake news detection on social media: A context-aware adversarial approach. In Frontiers in fake media generation and detection, pages 215–232. Springer.
- Aryan Singhal, Thomas Law, Coby Kassner, Ayushman Gupta, Evan Duan, Aviral Damle, and Ryan Luo Li. 2024. Multilingual fact-checking using LLMs. In Proceedings of the Third Workshop on NLP for Positive Impact, pages 13–31, Miami, Florida, USA. Association for Computational Linguistics.
- Shivangi Singhal, Rajiv Ratn Shah, and Ponnurangam Kumaraguru. 2022. Factdrill: A data repository of fact-checked social media content to study fake news incidents in india. Proceedings of the International AAAI Conference on Web and Social Media, 16(1):1322-1331.
- S. Suryavardan, Shreyash Mishra, Parth Patwa, Megha Chakraborty, Anku Rani, Aishwarya Naresh Reganti, Aman Chadha, Amitava Das, Amit P. Sheth, Manoj Chinnakotla, Asif Ekbal, and Srijan Kumar. 2023. Factify 2: A multimodal fake news and satire news dataset. In Proceedings of De-Factify 2: 2nd Workshop on Multimodal Fact Checking and Hate Speech Detection, co-located with AAAI 2023, Washington DC, USA, February 14, 2023, volume 3555 of CEUR Workshop Proceedings. CEUR-WS.org.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv:2403.08295.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 809-819, New Orleans, Louisiana. Association for Computational Linguistics.
- Shivani Tufchi, Ashima Yadav, and Tanveer Ahmed. 2023. A comprehensive survey of multimodal fake news detection techniques: advances, challenges, and opportunities. International Journal of Multimedia Information Retrieval, 12(2):28.
- Juraj Vladika and Florian Matthes. 2023. Scientific fact-checking: A survey of resources and approaches. In Findings of the Association for Computational Linguistics: ACL 2023, pages 6215–6230, Toronto, Canada. Association for Computational Linguistics.
- Bo Wang, Jing Ma, Hongzhan Lin, Zhiwei Yang, Ruichao Yang, Yuan Tian, and Yi Chang. 2024a. Explainable fake news detection with large language

model via defense among competing wisdom. In Proceedings of the ACM Web Conference 2024, pages 2452-2463.

- Xinyu Wang, Wenbo Zhang, and Sarah Rajtmajer. 2024b. Monolingual and multilingual misinformation detection for low-resource languages: A comprehensive survey. arXiv preprint arXiv:2410.18390.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.
- Deressa Wodaio and Solomon Atnafu. 2021. Deepfake video detection using convolutional vision transformer. arXiv preprint arXiv:2102.11126.
- Dustin Wright and Isabelle Augenstein. 2020. Claim check-worthiness detection as positive unlabelled learning. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 476-488, Online. Association for Computational Linguistics.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. Advances in neural information processing systems, 32.
- Zhiwei Yang, Jing Ma, Hechang Chen, Hongzhan Lin, Ziyang Luo, and Yi Chang. 2022a. A coarse-tofine cascaded evidence-distillation neural network for explainable fake news detection. arXiv preprint arXiv:2209.14642.
- Zhiwei Yang, Jing Ma, Hechang Chen, Hongzhan Lin, Ziyang Luo, and Yi Chang. 2022b. A coarse-tofine cascaded evidence-distillation neural network for explainable fake news detection. In Proceedings of the 29th International Conference on Computational Linguistics, pages 2608-2621, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee 1013 Cho, and Lifu Huang. 2023. End-to-end multimodal 1014 fact-checking and explanation generation: A chal-1015 lenging dataset and models. In Proceedings of the 1016

46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23, page 2733–2743, New York, NY, USA. Association for Computing Machinery.

1017

1018

1019

1021

1023

1024

1026

1027

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039 1040

1041

1042

1043

1044

1049

1050

1051

1052

1053

1054

1055

1057

1059

1060

1062

1063

1064

1065

1066

1068

1070

- Zhenrui Yue, Huimin Zeng, Lanyu Shang, Yifan Liu, Yang Zhang, and Dong Wang. 2024. Retrieval augmented fact verification by synthesizing contrastive arguments. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 10331–10343.
- Zhenrui Yue, Huimin Zeng, Yang Zhang, Lanyu Shang, and Dong Wang. 2023. Metaadapt: Domain adaptive few-shot misinformation detection via meta learning. *arXiv preprint arXiv:2305.12692*.
- Fengzhu Zeng and Wei Gao. 2023. Prompt to be consistent is better than self-consistent? few-shot and zero-shot fact verification with pre-trained language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4555–4569, Toronto, Canada. Association for Computational Linguistics.
- Xia Zeng, Amani S Abumansour, and Arkaitz Zubiaga. 2021. Automated fact-checking: A survey. *Language and Linguistics Compass*, 15(10):e12438.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Tong Zhang, Di Wang, Huanhuan Chen, Zhiwei Zeng, Wei Guo, Chunyan Miao, and Lizhen Cui. 2020.
 Bdann: Bert-based domain adaptation neural network for multi-modal fake news detection. In 2020 international joint conference on neural networks (IJCNN), pages 1–8. IEEE.
- Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 996–1011.
- Xiaofan Zheng, Zinan Zeng, Heng Wang, Yuyang Bai, Yuhan Liu, and Minnan Luo. 2025. From predictions to analyses: Rationale-augmented fake news detection with large vision-language models. In Proceedings of the ACM on Web Conference 2025, WWW '25, page 5364–5375, New York, NY, USA. Association for Computing Machinery.

Appendix

A Related works:

The NLP community has long been focusing on the core and peripheral tasks of automated factchecking. Several datasets and methodologies have been proposed in the past, spanning different do-
mains, languages, and task frameworks. In the
following, we have reported the past literature for
(i) monolingual fact-checking, (ii) multi-lingual
and multi-modal fact-checking, (iii) explainable
fact-checking, and (iv) associated tasks in this area.1072
1073
1074
1074
1075
1075
1076Apart from that, we have also reported how our
work fills the current research gap.1078

A.1 Monolingual fact-checking:

Here, we reported the past works for monolingual fact-checking. Specifically, we focused on two as-1081 pects, i.e. datasets and approaches. We found sev-1082 eral datasets focusing on text-based fact-checking 1083 in the monolingual domain. Some of the popular 1084 ones are, FEVER (Thorne et al., 2018), FEVER-1085 OUS (Aly et al., 2021), VITAMIN-C (Schuster et al., 2021), LIAR-RAW and RAW-FC (Yang 1087 et al., 2022b). Thorne et al. (2018) produced one of the earlier famous datasets, FEVER, consisting of around 185K, claims which were extracted from 1090 the Wikipedia corpus. Annotators tagged them 1091 with three distinct labels: 'Supported', 'Refuted' or 'NotEnoughInfo' (NEI). LIAR-RAW and RAWFC 1093 datasets were constructed by Yang et al. (2022b). 1094 RAW-FC consists of around 2K claims with associ-1095 ated raw reports collected from the Snopes² web-1096 site. They were annotated for three labels, namely, 1097 'true', 'half-true', and 'false'. On the other hand, 1098 the LIAR-RAW dataset was an extension of the 1099 LIAR-PLUS (Alhindi et al., 2018) dataset. The au-1100 thors accompanied each claim with raw reports, 1101 and expert annotators labelled the claims with 1102 one label out of six: 'pants-fire', 'false', 'barely-1103 true', 'half-true', 'mostly-true' and 'true'. The 1104 samples are domain-agnostic and mainly used to 1105 train general-domain fact-checking models. Addi-1106 tionally, we have domain-specific datasets as well. 1107 For example, PUBHEALTH dataset proposed by 1108 Kotonya and Toni (2020b) does fact-checking in 1109 the healthcare domain. It consists of 11.8K claims, 1110 where each claim is annotated and tagged with one 1111 of the following four labels: 'true', 'false', 'mix-1112 ture', and 'unproven'. Similarly, the SciClaimHunt 1113 dataset by Kumar et al. (2025) focused on scientific 1114 claim verification. Here, each claim is labelled as 1115 positive or negative based on the evidence provided 1116 in the scientific research papers. A complete list 1117 of datasets focusing on automated fact-checking is 1118 reported in Guo et al. (2022); Vladika and Matthes 1119

²https://www.snopes.com/

1121 1122

1123 1124

1125

1126

1127 1128

(2023); Zeng et al. (2021).

From a methodological point of view, the choice of fact-checking systems was closely tied to the datasets and task frameworks. Early approaches (DAGAN et al., 2010; Bowman et al., 2015) assessed whether each retrieved evidence supports or refutes a claim as an entailment task.

A.2 Mutli-lingual and multi-modal fact checking

The need to address misinformation across vari-1129 ous languages and media formats have expanded 1130 the field of fact-checking into multilingual and 1131 multi-modal domains. In multilingual settings, 1132 datasets such as FakeCovid (Shahi and Nandini, 1133 2020), which focuses on COVID-19 misinforma-1134 tion covering 40 languages, and X-Fact (Gupta 1135 and Srikumar, 2021) covering 25 languages and di-1136 verse domains. Prior surveys (Wang et al., 2024b; 1137 Singhal et al., 2024) provided a complete list of 1138 datasets and methods ranging from machine learn-1139 ing classifiers to LLM based techniques used for 1140 multi-lingual fact-checking. In parallel, a com-1141 prehensive survey by Akhtar et al. (2023a) high-1142 lighted the evolution of multi-modal fact checking 1143 domain. Several datasets have emerged to support 1144 this line of research such as FactDrill (Singhal et al., 1145 2022) which combines video, audio, image, text, 1146 and metadata, while r/Fakeddit (Nakamura et al., 1147 2020), MuMiN (Nielsen and McConville, 2022), 1148 MOCHEG dataset (Yao et al., 2023), and Factify-2 1149 (Suryavardan et al., 2023) focused on image-text 1150 pairs only. Gupta et al. (2022) introduced a novel 1151 multi-lingual multimodal misinformation (MMM) 1152 dataset which integrated three Indian languages 1153 into multimodal fact-checking domain. In terms 1154 of modeling, early approaches used CNN-based 1155 models, such as ResNet (Sabir et al., 2019), and 1156 VGG16 (Amerini et al., 2019). Later on, studies in-1157 tegrated recurrent networks such as LSTM (Güera 1158 and Delp, 2018) to give better verdict. In the fol-1159 lowing years, researchers used CNN-LSTM hybrid 1160 (Tufchi et al., 2023) and transformers based models 1161 such as ViT (Wodajo and Atnafu, 2021). Recently, 1162 1163 pretrained models have gained pace into multimodal fact checking (Zhang et al., 2020; Cekinel 1164 et al., 2025). Several studies (Akhtar et al., 2023b; 1165 Tufchi et al., 2023) highlighted the list of datasets 1166 and methods used for multi-modal fact checking. 1167

A.3 Explainable fact-checking:

Explainability has emerged as a critical re-1169 quirement for trustworthy fact-checking systems. 1170 Kotonya and Toni (2020a); Eldifrawi et al. (2024) 1171 illustrated the recent advancements in this domain 1172 through comprehensive surveys. Alhindi et al. 1173 (2018) created the LIAR-PLUS dataset, extend-1174 ing the LIAR dataset, by including justification 1175 from long ruling comments. Except for general 1176 domain, explanability has also entered into more in-1177 tricate domains like healthcare (Kotonya and Toni, 1178 2020b). In multi-modal contexts, MOCHEG (Yao 1179 et al., 2023) introduced the first end-to-end dataset 1180 to include structured explanations, bridging the 1181 gap between multi-modal fact verification and in-1182 terpretability. Early methods relied on attention 1183 mechanisms (i.e., highlighting high-attention to-1184 kens) to pinpoint evidence span for explanation 1185 (Popat et al., 2018). But critiques said attention 1186 weights often misrepresent true model reasoning 1187 and lack accessibility for non-experts (Pruthi et al., 1188 2020). Rule-based systems (Gad-Elrab et al., 2019) 1189 improved transparency but restricted scope to struc-1190 tured claims and pre-defined data. More recent 1191 textual explanation models face challenges like hal-1192 lucination in abstractive outputs (Maynez et al., 1193 2020). Out of them, our proposed method utilized 1194 the 'Justification-Then-Veracity' pipeline, as men-1195 tioned in Eldifrawi et al. (2024), for reliable fact-1196 checking. 1197

1168

1198

A.4 Associated tasks:

There are various associated tasks along with fact-1199 checking such as claim check-worthiness, rumour 1200 detection, stance detection, etc. Claim check-1201 worthiness (Wright and Augenstein, 2020) is re-1202 quired before jumping directly into fact verification. 1203 Another related task, rumour detection (Dougrez-1204 Lewis et al., 2022; Gorrell et al., 2019) identifies 1205 unverified claims in real-time, often using prop-1206 agation patterns. Stance detection gauges public reactions to prioritize claims (Baly et al., 2018). 1208 Exaggeration detection (Patro and Baruah, 2021) is 1209 another task where similar statements are compared 1210 to find which of them is a more exaggerated version. 1211 Credibility detection (Patro and Rathore, 2020) as-1212 sesses the trustworthiness of topics discussed in so-1213 cial media. Clickbait detection (Chakraborty et al., 1214 2016) identifies the online content (headlines or 1215 titles) specially designed to attract users to click. 1216

A.5 Research gap:

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1254

1255

1256

1257

1259

1260

1261

1263

1264

Prior work by Yue et al. (2024) took retrieved evidences to generate supporting and opposing arguments for independent evaluation, and utilized fewshot prompting for claim verification. Similarly, Wang et al. (2024a) split relevant evidences into supporting or refuting categories relying on a complex attention mechanism. However, both of them suffer from key limitations: Yue et al. (2024)'s method did not consider dividing evidence set that may contain contradictory statements, whereas Wang et al. (2024a)'s complex attention mechanism discarded some useful information by focusing only on the top-k evidences based on the attention score. Here, we argue that VLLMs, with their broad understanding of language can instead analyse how each evidence supports (entails) or refutes the claim. Further, with their ability to process long input text, VLLMs can consolidate the whole support and refute evidence set to generate respective justification. To test this hypothesis, we compared it with various training-based and inference-based methods. To the best of our knowledge, our work is the first to use claim-evidence entailment in VLLMs for this task, ensuring no evidence is overlooked.

B Additional details on datasets:

In this section, we reported a detailed discussion of the considered datasets. Additionally, we have also presented some samples and statistics for the datasets for better illustration.

B.1 LIAR-RAW (Yang et al., 2022b):

LIAR-RAW (Yang et al., 2022b) consists of 12,590 claims, each paired with raw reports collected from news articles, press releases, and web pages. LIAR-RAW builds upon the LIAR-Plus dataset (Alhindi et al., 2018) by adding crowd-sourced raw reports. Authors have retrieved up to 30 raw reports for each claim via the Google API using claim keywords. Claims were collected from a well-known factchecking website, i.e. *Politifact*³, which provided gold veracity labels. To ensure quality, they excluded fact-checking site reports, those published after the verdict, and reports under 5 words or over 3,000 words. After filtering, the final dataset had 10,065 training instances, 1,274 validation instances, and 1,251 test instances. Expert annotators manually assigned one of six fine-grained veracity

labels to each claim. The labels they considered are: "pants-fire" (completely false), "false", "barelytrue" (contains minimal truth but is mostly false), "half-true" (equally true and false, with a significant mix of both), "mostly-true" (predominantly true with minor inaccuracies) and "true". A sample from the dataset is presented in Table 6 for illustration. We have also studied the claim-evidence distribution, where we clubbed the claims based on the number of evidence sentences they have. The distribution is illustrated in Figure 2. We found that (i) the majority of claims have no evidence sentences, (ii) more than a thousand of claims have one evidence sentence, and (iii) 170+ claims have more than fifty evidence sentences associated with them.

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1290

1292

1293

1294

1295

1296

1298

1299

1300



Figure 2: Bar plot showing number of claims v/s number of evidences for LIAR-RAW dataset

B.2 RAW-FC (Yang et al., 2022b):

RAWFC dataset, introduced by (Yang et al., 2022b), has 2,012 claims. Authors collected claims from *Snopes*⁴ website and retrieved relevant raw reports using claim keywords via the Google API. For each claim, up to 30 raw reports were gathered from various web sources. To ensure quality, reports from fact-checking sites and those published after the fact-checking verdict were excluded. Reports shorter than 5 words or longer than 3,000 words were also removed. Expert annotators manually assigned one of the three veracity labels to each claim. The labels used for the annotation are: "true" (entirely accurate), "half-true" (partially true but includes misleading information), and "false" (entirely false). A sample from the dataset is presented in Table 7 for illustration. We have studied the claim-evidence distribution for this dataset as well. We created buckets based on evidence sentence counts and gathered the claims

³https://www.politifact.com

⁴https://www.snopes.com

Dataset: LIAR-RAW (Yang et al., 2022b)

Event_id: "5209.json"

Claim: "Suzanne Bonamici supports a plan that will cut choice for Medicare Advantage seniors."

Label: "half-true"

Explain: "The Affordable Care Act was designed to save money by slowing future spending, including future spending on Medicare Advantage plans. But spending still goes up. In addition, many outside factors can affect the cost and range of benefits, making it impossible to know how Medicare Advantage might change. While the statement from Cornilles is partially accurate, it is taken out of context and ignores important details on a politically volatile subject."

Evidences:

1301

1302

1303

1304

1305

1306

E1: "So I hope, a we move forward, we will focus on that approach to modernize Medicare and Medicaid, an approach that improve the quality of care while we reduce the cost of care, rather than simply offload those cost onto senior.",

E2: "So some of the save we achieve , a significant amount of save we achieve , be in reduce these overpayment , these huge subsidy , to the private Medicare Advantage plan .",

E3: "In addition to the issue with Part D benefit design and plan flexibility, there be transaction such a rebate, pharmacy fee, and other form of compensation that occur in the supply chain that pose several issue.",

E4: "The ACA have help slow the growth in health care cost, it be close the doughnut hole for senior, and have encourage and improved access to mental health service and preventive care."

Table 6: Example entry from the LIAR-RAW dataset. It contains five key fields: (i) an identifier '*Event_id*', (ii) the '*claim*' to be fact-checked, (iii) the ground-truth '*Label*', (iv) an explanation ('*Explain*') that clarifies its ground truth, and (v) the evidence sentences (E1-E4) extracted from the underlying fact-checking article.

belonging to individual buckets. The distribution is illustrated in Figure 3. We found that (i) around 750 claims have twenty to fifty evidence sentences associated with them, and (ii) more than five hundred claims have fifty or more evidence sentences associated with them.



Figure 3: Bar plot showing number of claims v/s number of evidences for RAW-FC dataset

C Additional Details on Experiments:

In this section, we have reported additional details on our experiments. Particularly, we described individual **IBEs**, provided additional details on baselines and experimental setups. 1307

1308

1309

1310

1311

1312

C.1 IBE-1: Zero-shot prompting:

In this experiment, we have tried to answer "How 1313 well the VLLMs can predict the claim veracity 1314 if raw evidence sentences are provided in the in-1315 *put?"* To conduct this experiment, we prompted 1316 five VLLMs by giving claims and associated ev-1317 idence sentences as input. Some of the prompt 1318 samples are showcased in Figure 4 and Figure 5 1319 for illustration. In past, researchers have tried zero-1320 shot prompting for this task; however, either (i) 1321 they gave only claims without evidence sentences 1322 (Zhang and Gao, 2023), or (ii) they relied on ev-1323 idence retrieved from web searches (Cheung and 1324 Lam, 2023). 1325 Dataset: RAW-FC (Yang et al., 2022b)

Event_id: "247609"

Claim: "Right-wing commentator Bill Mitchell tweeted \"someone let me know\" when 55,000 Americans have died from COVID-19 coronavirus disease."

Label: "true"

Explain: "In late April 2020, a month-old tweet posted by right-wing Twitter commentator Bill Mitchell received new attention amid tragic circumstances.\n Feeding on a conservative media talking point that the COVID-19 coronavirus disease pandemic isn\u2019t as serious as officials say, Mitchell tweeted on March 21, 2020:\u201cWhen COVID-19 reaches 51 million infected in the US and kills 55,000, someone let me know.\u201d\nSome readers asked if the statement (displayed above) was a real tweet posted by Mitchell, and it is.\nSadly, as of this writing, the U.S. has surpassed 55,000 fatalities from COVID-19. The Centers for Disease Control and Prevention (CDC) reports that as of May 1, 2020, there are 1,062,446 COVID-19 cases in the U.S. and 62,406 deaths.\nSocial media users wasted no time posting comments directed at Mitchell, reminding him of his statement.\n\nBecause the tweet in question was published from Mitchell\u2019s Twitter account, we rate this claim\u201cCorrect Attribution.\u201d"

Evidences:

E1: "On March 21, Bill Mitchell tweet, \u2018 When COVID-19\u2026kills 55,000, let me know.",

E2: "In it, Hannan write that \u2018 [COVID-19] be unlikely to be as lethal a the more common form of influenza that we take for granted.",

E3: "While I \u2019 d normally be content to mock conspiracy theorist \u2014 I set up a Twitter account to make fun of bad COVID-19 take \u2014 spread false information about the pandemic be dangerous , and merit rebuttal.",

E4: "\u201d Owens delete the tweet after a Twitter user observe that her claim be not only false , but appear to be base on a lazy misreading of a Google search result .",

E5: "Given the dire situation, it seem worth know if the severity of the pandemic finally have penetrate the bubble of the most extreme coronavirus scoffer; if these people and their follower be ignore safety measure because they believe the pandemic be a false flag to take away their gun or a Deep State plot to take down Trump, then they risk contribute to the disease $\2019$ s spread."

Table 7: Example entry from the RAW-FC dataset. It contains five key fields: (i) an identifier '*Event_id*', (ii) the '*claim*' to be fact-checked, (iii) the ground-truth '*Label*', (iv) an explanation ('*Explain*') that clarifies its ground truth, and (v) the evidence sentences (E1-E4) extracted from the underlying fact-checking article.



Figure 4: IBE-1 (LIAR-RAW): Prompt used to predict the veracity of the claim using zero-shot prompting based on the given evidence sentences and the claim.

C.2 IBE-2: Zero-shot prompting with overall understanding:

1326

1327

1330

1331

1333

1334

1335

1336

1338

1339

1340

1341

1342

1343

1344

1345

1347

1348

1350

1351

1352

1353

1354

1355

In this experiment, we have tried to answer "Can zero-shot prompting improve the veracity prediction performance if claim-evidence relation understandings generated by VLLMs are given as input?". To conduct this experiment, we followed a two-step prompting strategy. First, we prompted the VLLMs to generate the claim-evidence understanding like we did in **TBE-2**. In the next step, we prompted them again with the claim and the generated understanding to predict the veracity labels. Some of the prompt samples are presented in the Figure 7 and Figure 8. To the best of our knowledge, nobody has attempted this in past.

C.3 IBE-3: CoT (Wei et al., 2022) prompting with overall understanding:

In this experiment, we tried to answer *Can CoT*based prompting improve the veracity prediction performance if claim-evidence understanding is given in the input?. Here we followed similar steps as mentioned in **IBE-2**, except in the prompt, we asked the VLLMs to generate step-by-step reasoning behind their predictions. Prior work (Zhang and Gao, 2023) did a similar attempt. However, they gave evidence collected from web searches. We, on the other hand, relied on the claim-evidence understanding generated by VLLMs. Some of the prompt samples are presented in Figure 10 and Figure 11 for illustration.



Figure 5: IBE-1 (RAW-FC): Prompt used to predict the veracity of the claim using zero-shot prompting based on the given evidence sentences and the claim.



Figure 6: IBE-2: First step, prompt used to generate an overall understanding of the claim based on the provided evidence sentences.

C.4 IBE-4: Prompting based on entailment:

1356

In the last inference-based experiment, we tried 1357 to answer "Can prompting with VLLM generated 1358 entailed-justifications enhance language models' 1359 ability to predict claim veracity?" To conduct 1360 this experiment, we have followed a three-step 1361 approach similar to TBE-3 as described in Sec-1362 tion 3.1.3. While the initial two steps, i.e. (i) to 1363 classify the evidence sentences as supporting or 1364 refuting a given claim and (ii) generating the justifi-1365 cations based on classified evidence sentences, are 1366 exactly the same, in the last step, instead of training, we prompted the VLLMs to generate the veracity. 1368 In past, researchers have prompted the VLLMs 1369 to find entailment for tasks like claim matching 1370 (Choi and Ferrara, 2024b) and counterfactual gen-1371 eration (Dai et al., 2022). However, to the best of 1372 our knowledge, we are the first to (i) apply it to 1373 classify each evidence into supporting or refuting 1374 categories, and (ii) generate entailed justifications 1375 and use them for fact verification. The detailed 1376



Figure 7: IBE-2 (LIAR-RAW): Second step, prompt using zero-shot prompting to directly predict the claim's veracity based on the overall understanding.

approach is illustrated in Fig. 1 (c). Some of the prompt samples are show cased in the Figure 12, Figure 13, Figure 14, Figure 15 and 16.

C.5 Additional details on baselines:

1378

1379

1380

1381

1382

1383

1384

1385

1386

1388

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

In this section, we reported the details of individual baseline methods considered in our study.

- **HiSS**: It is proposed by Zhang and Gao (2023). Here, authors have proposed a hierarchical step-by-step (a.k.a. *HiSS*) method where they first decomposed a claim into smaller subclaims using few-shot prompting. Then they verified each sub-claim step-by-step by raising and answering a series of questions. For each question, they prompted the language models to assess if it is confident in answering or not, and if not, they gave the question to a web search engine. The search results were then inserted back into the ongoing prompt to continue the verification process. Finally, using that information, language models predicted the veracity label for the whole claim.
- FactLLaMa: It is proposed by Cheung and Lam (2023). Their approach has two components: (i) generation of prompts (having instructions, claims and evidence pieces), and (ii) instruction-tuning of a generative pretrained language model. In the first com-



Figure 8: IBE-2 (RAW-FC): Second step, using zeroshot prompting to predict the claim's veracity from the overall understanding.



Figure 9: IBE-3: First step, prompt used to generate an overall understanding of the claim based on the provided evidence sentences.

ponent, to create the prompt samples, they combined the instruction, evidence, and input claim into a single sequence, with special tokens separating them. While the instruction guides how to incorporate the evidence for fact-checking, evidence contains relevant information retrieved from search engines, using the Google API. As a part of the second component they we fine-tuned the LoRA adapter with Llama as the language model. 1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

• **RAFTS**: It stands for *retrieval augmented* 1414 fact verification through the synthesis of con-1415 trastive arguments. It is proposed by Yue 1416 et al. (2024), where authors retrieve relevant 1417 documents and perform a few-shot fact ver-1418 ification using pretrained language models. 1419 RAFTS have three components, (i) demon-1420 stration retrieval, where relevant examples 1421 were collected to included in the input con-1422 texts, (ii) document retrieval, where authors 1423 proposed a retrieve and re-rank pipeline (us-1424 ing RAG framework) to accurately identify 1425



Figure 10: IBE-3 (LIAR-RAW): Second step, prompt using Chain-of-Thought reasoning to predict the veracity of the claim from the overall understanding.

relevant documents for the input claim, and (iii) few-shot fact verification using supporting and opposing arguments derived from the facts within the collected documents. However, unlike in our work, they didn't rely on supporting and opposing justifications (a.k.a. entailed justifications) in their final steps.

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1449

1450

1451

1452

1453

1454

• L-Defense: The LLM-equipped defensebased explainable fake news detection approach (L-Defense) was proposed by Wang et al. (2024a). Their approach shares some similarity with RAFTS at a philosophical level. The framework consists of three components: (i) a competing evidence extractor, (ii) a prompt-based reasoning module, and (iii) a defense-based inference module. In the competing evidence extractor, authors deployed a natural language inference (NLI) module to associate a "true" or "false" label to each claim for each associated evidence sentence. The NLI module gave two top-k evidence sentence sets, each stating a claim to be "true" or "false. In the reasoning module, they prompted language models to generate two separate explanations, each for stating a claim to be "true" and "false based on the respective top-k evidence set. Finally, in the defense-based inference module, they trained the transformer encoders by taking claims and associated two



Figure 11: IBE-3 (RAW-FC): Second step, prompt using Chain-of-Thought reasoning to predict the veracity of the claim from the overall understanding.

System Prompt										
You are an evidence classifier that determines whether a sentence supports or refutes a given claim.										
User Prompt										
Task: Analyze the sentence below and classify its relationship to the given claim. Classification Rules: 1. SUPPORT - The sentence provides evidence that supports the claim. 2. REFUTE- The sentence provides evidence that contradicts or refutes the claim. Claim: (claim)										
Label for Claim: {label} Sentence to classify: {sent}										
Carefully analyze the sentence and respond with: 1. A classification (SUPPORT or REFUTE).										

Figure 12: IBE-4: First step prompt used to classify each evidence sentence as supporting or refuting the claim.

explanations to predict the veracity. Our approach is different from the L-Defense as we (i) deployed VLLMs to classify each evidence, (ii) we considered all evidences associated with a claim, and (iii) we fine-tuned LLMs and adapters with VLLMs for veracity prediction (in **TBE-3**).

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

C.6 Experimental set-up:

C.6.1 Training/ test/ validation splits:

We have used the training, validation, and test splits originally provided by Yang et al. (2022b) in our experiments. They are essential for the comparison of our models with the baselines. The distribution of samples falling under each label and training splits are reported in Table 10.

C.6.2 Language models:

For prompting, we used five very large language1471models (VLLMs). We call them so because of1472



Figure 13: IBE-4, Second step prompt used to generate a supporting justification based on the supporting evidence sentences.



Figure 14: IBE-4, Second step prompt used to generate a refuting justification based on the refuting evidence sentences.

their immense parameter size, which prohibits us from fine-tuning them like we do for normal BERTbased language models. We used five VLLMs i.e. Mistral (Jiang et al., 2023), Llama (AI@Meta, 2024), Gemma (Team et al., 2024), Qwen (Yang et al., 2024) and Falcon (Almazrouei et al., 2023) in all of our experiments. For fine-tuning, we used LLMs RoBERTa and XLNet in **TBE-2** and **TBE-3**. Apart from that, we also used adapter-based (LoRA (Hu et al., 2022) and LoRA+ (Hayou et al., 2024) adapters) fine-tuning methods for the considered VLLMs. The details of LLMs and VLLMs such as their versions and maximum input size they can take are reported in Table 9.

C.6.3 Hyperparameter details:

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

We did an extensive hyperparameter search that led to the optimal performance of our models. The list of hyperparameters for which we trained our models is presented in Table 8. For the VLLMs, we kept the temperature constant at '0.001' for consistency. We conducted all our experiments on two NVIDIA A100 80GB GPU cards.

C.6.4 Evaluation metrics:

Since we are working in a multi-class classification framework for veracity prediction, we used standard metrics such as macro-precision (MP),



Figure 15: IBE-4 (LIAR-RAW): Third step prompt used to predict the veracity of the claim using the supporting and refuting justifications.

Parameters	Values
Learning rate	{2e-6, 2e-5, 1e-5}
Optimizer	AdamW, Adam
Batch size	8, 16
Patience (Early stop)	2, 3
lora_rank	8
Learning rate	{1e-5,1e-4}
lr_scheduler_type	cosine
bf16	true

Table 8: Hyperparameters explored during model training and evaluation.

macro-recall (MR), and macro-f1 (MF) to eval-1499 uate our models. To evaluate model explainabil-1500 ity, we first concatenated the supporting and re-1501 futing entailed justifications (generated as part of 1502 TBE-3) and considered them as model explana-1503 tions. The evaluation was done with two types of 1504 strategies: (i) checking lexical overlap and seman-1505 tic matching, and (ii) doing subjective evaluation by VLLMs. To check lexical overlap, we used sev-1507 eral standard evaluation metrics such as ROUGE-1 1508 (R_1) , ROUGE-2 (R_2) , ROUGE-L (R_L) (Lin, 2004) and BLEU (Papineni et al., 2002). While R_1 and 1510 R_2 measure the overlap of unigrams and bigrams 1511 between predicted and gold explanations, R_L mea-1512 sures the longest common subsequence. BLEU, on 1513 the other hand, measures the precision of matching 1514 n-grams between predicted and gold explanations. 1515



Figure 16: IBE-4 (RAW-FC): Third step prompt used to predict the veracity of the claim using the supporting and refuting justifications.

Model	Max. length	Version
RoBERTa	1k	roberta-large
XLNet	1k	xlnet-large-cased
Mistral	32k	mistralai/Mistral-7B-Instruct-v0.3
Llama	128k	meta-llama/Llama-3.1-8B-Instruct
Gemma	8k	google/gemma-7b-it
Qwen	1M	Qwen/Qwen2.5-7B-Instruct-1M
Falcon	8k	tiiuae/Falcon3-7B-Instruct

Table 9: Language model details. Here, 'Max. length' denotes the maximum sequence length allowed by the particular model.

To measure the semantic matching between predicted and gold explanations, we deployed BERT score (Zhang et al., 2019). It generates contextual embeddings of predicted and gold explanations and calculates cosine similarity between them. To do the subjective evaluation, we prompted the considered VLLMs (a.k.a evaluating VLLMs) to asses the model explanations generated by each VLLM (generating VLLM). The VLLMs were asked to asses across five dimensions (i) informativeness, (ii) logicality, (iii) objectivity, (iv) readability and (v) accuracy (Zheng et al., 2025). The prompt template used for the assessment is presented in Figure 25.

1516

1517

1518

1522

1523

1525

1526

1528

1529

1530

1532

1534

1536

D Additional Results and Discussion:

D.1 Additional observations from TBE-1:

• We found that Mistral trained with LoRA adapter resulted in the highest macroprecision (0.44) for the LIAR-RAW dataset. While it surpasses the performance of baselines such as FactLLaMA (MP : 0.32,



Figure 17: A schematic diagram explaining the TBE-3 pipeline which breaks down into three stages. (a) *Entailment*: Given a claim and its corresponding evidences, a VLLM first distinguishes evidences into supporting/refuting evidence via entailment. (b) *Evidence consolidation*: Then the same VLLM consolidates these two groups into concise supporting and refuting justification. (c) *Veracity prediction*: Using the claim and both justifications, an LLM is trained to predict veracity.

37.50% ↑) and L-Defense (MP : **0.31**, ~ **41.94**% ↑), it is still behind HiSS (MP : **0.46**) and RAFTS (MP : **0.47**). For RAWFC dataset, however, many models, such as Mistral trained with LoRA adapter (MP : **0.69**, ~ **11.29**% ↑), Llama trained with LoRA (MP : **0.68**, ~ **9.68**% ↑) and LoRA+ adapters (MP : **0.67**, ~ **8.06**% ↑), Qwen trained with LoRA (MP : **0.67**, ~ **8.06**% ↑) and LoRA+ (MP : **0.70**, ~ **12.90**% ↑) adapters, outperformed the best performance reported by the baselines (MP : **0.62**). Qwen trained with LoRA+ adapter achieved the highest overall performance (MP : **0.70**).

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1550

• We observed that Llama trained with LoRA 1551 adapter resulted in the highest macro-recall 1552 (0.30) for the LIAR-RAW dataset. It falls behind the macro-recall reported by all baselines 1554 HiSS (MR : 0.31), FactLLaMa (MR : 0.32), 1555 L-Defense (MR : 0.31), and RAFTS (MR : 1556 0.37). However, for the RAW-FC dataset, 1557 we observed that many models, such as Mis-

Dataset	Label	Train	Val	Test
	Т	1647	169	205
	MT	1950	251	238
LAIR-RAW	HT	2087	244	263
(Yang et al., 2022b)	BT	1611	236	210
	F	1958	259	249
	PF	812	115	86
	Т	561	67	67
RAW-FC	HT	537	67	67
(Yang et al., 2022b)	F	514	66	66

Table 10: Train, val and test distributions. Notations: **T** for True, **MT** for Mostly-true, **HT** for Half-true, **BT** for Barely-true, **F** for False, **PF** for Pants-fire and **T** for True.



Figure 18: TBE-2: Prompt for generating overall understanding from the claim and evidence sentences.

tral trained with LoRA adapter (MR : 0.65, ~ 6.56% \uparrow), Llama trained with LoRA (MR : 0.64, ~ 4.92% \uparrow) and LoRA+ adapters (MR : 0.64, ~ 4.92% \uparrow), Qwen trained with LoRA adapter (MR : 0.66, ~ 8.20% \uparrow) outperformed the best reported macro-recall score by the baselines (0.61). Qwen trained with LoRA adapter achieved the highest overall performance (MR : 0.66).

D.2 Additional observations from TBE-2:

1561

1562

1563

1564

1565

1566

1569

1570

1572

1573

1574

1575

1577

1578

1580

1581

1583

1585

• We found that Mistral trained (with Mistral understandings) with the LoRA adapter resulted in the highest macro-precision (0.36) for the LIAR-RAW dataset. While it outperforms the baselines FactLLaMA (MP : **0.32**) and L-Defense (MP : 0.31) models, it is still behind HiSS (MP : **0.46**) and RAFTS (MP : **0.47**). However, for the RAW-FC dataset, we observed that many models, such as XLNet fine-tuned with Llama understandings (MP : **0.63**, \sim **1.61**% \uparrow), Mistral trained (with Mistral understandings) with LoRA+ (MP : 0.66, \sim 6.45% \uparrow) and Llama trained (with Llama understandings) with LoRA+ (MP : 0.73, \sim 17.74% \uparrow) adapters outperformed the best reported macro-precision by the baselines (MP: 0.62). Llama trained (with Llama un-







Figure 20: TBE-3: Second step prompt used to generate a supporting justification based on the supporting evidence sentences.

derstandings) with LoRA+ adapter resulted	15
in the overall highest macro-precision (MP :	15
0.73).	15

• We found that Mistral trained (with Mistral understandings) with LoRA adapter gave us 1590 the highest macro-recall (0.32) for LIAR-1591 RAW dataset. While this performance is 1592 comparable to the performance of baselines 1593 HiSS (MR : 0.31), FactLLaMA (MR : 1594 **0.32**), and L-Defense (MR : 0.32), it is 1595 still behind RAFTS (MR : 0.37). However, 1596 for the RAW-FC dataset, we observed that 1597 many models, such as Llama trained (with Llama understandings) with LoRA+ adapter 1599 $(MR: 0.71, \sim 16.39\% \uparrow)$ and XLNet fine-1600 tuned with Llama understandings (MR : **0.63**, \sim **3.27**% \uparrow) outperformed the best re-1602 ported macro-recall score by the baselines 1603 (MR : 0.61). Llama trained (with Llama un-1604 derstandings) with LoRA+ adapter achieved 1605 the highest macro-recall overall (MR : **0.71** ~ 16.39% ^(†). 1607



Figure 21: TBE-3: Second step prompt used to generate a refuting justification based on the refuting evidence sentences.

D.3 Additional observations from TBE-3:

1608

• We found that XLNet fine-tuned with Llama entailed justifications gave the highest macroprecision (0.55) for the LIAR-RAW dataset. 1611 Many models, such as RoBERTa fine-tuned 1612 with Mistral (*MP* : **0.48**, \sim **2.12**% \uparrow), 1613 Llama (*MP* : **0.53**, \sim **12.76**% \uparrow), Gemma 1614 $(MP : 0.49, \sim 4.25\% \uparrow), \text{Qwen} (MP : 0.48,$ ~ 2.12% \uparrow), and Falcon (*MP* : 0.49, ~ 1616 **4.25**% \uparrow) entailed justifications, XLNet fine-1617 1618 tuned with Mistral (MP : **0.49**, ~ **4.25**% \uparrow), Llama (*MP* : **0.55**, \sim **17.02**% \uparrow), and 1619 Qwen (MP : 0.50, \sim 6.38% \uparrow) entailed justifications and Llama trained (with Llama entailed justifications) with LoRA+ adapter $(MP: 0.49, \sim 4.25\% \uparrow)$ surpassed the best reported macro-precision score of baselines 1624 (MP : 0.47). Similarly, for the RAW-FC 1625 dataset, we observed that many models, such as such as RoBERTa fine-tuned with Mistral $(MP : 0.83, \sim 33.87\% \uparrow)$, Llama (MP :1628 **0.88**, ~ **41.93**% \uparrow), Qwen (*MP* : **0.73**, 1629 ~ 17.74% \uparrow), and Falcon (*MP* : 0.65, ~ 1630 **4.83**% \uparrow) entailed justifications, XLNet finetuned with Mistral (MP : **0.83**, ~ **33.87**% \uparrow), 1632 Llama (*MP* : **0.88**, \sim **41.93**% \uparrow), Qwen $(MP : 0.70, \sim 12.90\% \uparrow)$, and Falcon $(MP : 0.76, \sim 22.58\% \uparrow)$ entailed justifications, Mistral trained (with Mistral entailed justifications) with LoRA (MP : **0.69**, 1637 \sim 11.29% \uparrow), Llama trained (with Llama en-1638 tailed justifications) with LoRA (MP : **0.63**, ~ **1.61**% \uparrow) and LoRA+ adapters (*MP* : 1641 **0.84**, \sim **35.48**% \uparrow), respectively, surpassed the best macro-precision reported by the base-1642 lines (MP : 0.61). RoBERTa and XLNet fine-1643 tuned with Llama entailed justifications gave the highest macro-precision (MP : 0.88). 1645

• We observed that XLNet fine-tuned with 1646 Llama entailed justifications gave the high-1647 est macro-recall (0.54) for the LIAR-RAW 1648 dataset. Many models, such as RoBERTa fine-1649 tuned with Mistral (MR : **0.47**, \sim **27.02**% \uparrow), 1650 Llama (MR : **0.53**, ~ **43.24**% \uparrow), Gemma 1651 $(MR : 0.50, \sim 35.14\% \uparrow)$, Owen (MR :1652 **0.47**, \sim **27.02**% \uparrow), and Falcon (*MR* : **0.43**, 1653 ~ 16.22% \uparrow) entailed justification, XLNet 1654 fine-tuned with Mistral (MR : 0.47, \sim 1655 **27.02**% \uparrow), Llama (*MR* : **0.54**, ~ **45.95**% \uparrow), Gemma (MR : **0.43**, ~ **16.22**% \uparrow), Qwen 1657 $(MR : 0.48, \sim 29.73\% \uparrow)$, and Falcon 1658 $(MR: 0.44, \sim 18.92\% \uparrow)$ entailed justifications, Llama trained (with Llama entailed jus-1660 tifications) with LoRA+ adapter (MR : **0.50**, \sim **35.14**% \uparrow) and Falcon trained (with Falcon 1662 entailed justifications) with LoRA+ adapter 1663 $(MR: 0.39, \sim 5.41\% \uparrow)$ surpassed the best 1664 reported macro-recall score of the baselines (MR : 0.37). Similarly, for the RAW-FC 1666 dataset, we observed that many models, such as such as RoBERTa fine-tuned with Mistral $(MR : 0.82, \sim 34.43\% \uparrow)$, Llama (MR :1669 **0.88**, ~ **44.26**% \uparrow), Qwen (*MR* : **0.71**, 1670 ~ 16.39% \uparrow), and Falcon (MR : 0.64, ~ **4.92**% \uparrow) entailed justifications, XLNet fine-1672 tuned with Mistral (MR : **0.82**, ~ **34.43**% \uparrow), Llama (MR : **0.88**, ~ **44.26**% \uparrow), Qwen $(MR : 0.70, \sim 14.75\% \uparrow)$, and Falcon 1675 $(MR : 0.75, \sim 22.95\% \uparrow)$ entailed justifi-1676 cations, Llama trained (with Llama entailed 1677 justifications) with LoRA (MR : 0.62, ~ 1678 1.64% \uparrow) and LoRA+ adapter (MR : 0.82, ~ **34.43**% \uparrow) surpassed the best reported macro-1680 precision score of baselines (MR : 0.61). 1681 RoBERTa and XLNet fine-tuned on Llama 1682 entailed justifications gave the highest overall performance (MR : 0.88). 1684

D.4 Additional observations from IBEs:

• For LIAR-RAW dataset, we observed that 1686 Qwen gave the highest macro-precision (0.27) in IBE-1. While in IBE-2, Llama achieved 1688 the highest macro-precision (MP : 0.27), 1689 Mistral gave the highest macro-precision in IBE-3 (*MP* : **0.27**) and IBE-4 (*MP* : **0.27**). 1691 However, none of them could surpass the 1692 baselines. In contrast, for RAW-FC dataset, Qwen (MP : 0.61), Llama (MP : 0.62), 1694 and Falcon (MP : 0.60) achieved highest macro-precision in IBE-1, IBE-2, and IBE-3 1696 1697respectively. They are comparable to the high-
est macro-precision reported by the baseline1698(MP: 0.62). Similarly, Qwen (MP: 0.50)1700got the best macro-precision in IBE-4, which
is far behind the baselines.

1702 • For LIAR-RAW dataset, Mistral achieved the highest macro-recall (0.25) in IBE-1. While 1703 in IBE-2, Mistral, Llama and Qwen achieved 1704 the highest macro-recall (MR : 0.23), for 1705 IBE-3 Mistral and Falcon (MR : 0.23) gave 1706 the highest scores. None of them could sur-1707 pass the baseline performances. Similarly, 1708 they achieved highest scores in IBE-4 (MR : 1709 0.17). In contrast, for RAW-FC dataset, 1710 Llama (MR : **0.63**, ~ **3.29**% \uparrow) achieved 1711 highest performace in IBE-2, surpassing the highest baseline performance (MR : 0.61). 1713 While Qwen attained the highest performance 1714 in IBE-1 (MR : **0.58**), surpassing the base-1715 lines HiSS (MR : **0.54**), FactLLaMa (MR : 1716 **0.55**) and RAFTS (MR : 0.52), it is still behind the performance of L-Defense (MR :1718 0.61). Similarly, Qwen got the highest per-1719 formance in IBE-3 (MR : 0.54), which sur-1720 passed the baseline RAFTS (MR : 0.52), but 1721 it is still behind the performance of HiSS 1722 1723 (MR : 0.54), FactLLaMa (MR : 0.55), and L-Defense (MR : **0.61**). In IBE-4, Mistral 1724 and Owen achieved the highest macro-recall 1725 of **0.46**, which is far behind all baselines. 1726

D.5 Additional observations from comparing the performance of TBEs and IBEs:

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

1741

1742

1743

1744

1745

1746

• We compared the macro-precision scores of best-performing models across TBEs and IBEs and presented them in Figure 23 for illustration. For the LIAR-RAW dataset, none of the models from IBEs, TBE-1, or TBE-2 (MP: 0.27 - 0.44) could exceed the highest macro-precision reported by the baseline (MP: 0.47). However, some models in TBE-1 and TBE-2 achieved higher macro-precision $(0.36 - 0.44, \sim 10.00\% \uparrow)$ than the IBEs (MP: 0.27 - 0.40). Moreover, several TBE-3 models (*MP*: **0.48** - **0.55**, \sim **25.00**% \uparrow) outperformed all models in IBEs (MP: 0.27 -**0.40**), TBE-1 (*MP*: **0.44**), and TBE-2 (*MP*: **0.36**). For the RAW-FC dataset, we observed similar comparative results. Here, the highest macro-precision reported by an IBE-2 (0.62) model was nearly the same as the bestperforming baseline (MP: 0.62). But, none 1747 of the models from IBEs (*MP*: **0.50** - **0.62**) 1748 were able to surpass the baseline (MP: **0.62**). 1749 Similar to the LIAR-RAW dataset, here we 1750 found that (i) many TBE-1 and TBE-2 mod-1751 els (*MP*: **0.63** - **0.73**, \sim **17.74**% \uparrow) out-1752 performed the IBE models (0.50-0.62), and 1753 (ii) many TBE-3 models (*MP*: **0.63** - **0.88**, 1754 \sim 20.55% \uparrow) outperformed all of the mod-1755 els in IBEs (0.50-0.62), TBE-1 (MP: 0.64 -1756 **0.70**) and TBE-2 (MP: **0.63** - **0.73**). 1757

- We have also compared the macro-recall 1758 scores of the best-performing models across 1759 TBEs and IBEs and presented them in Fig-1760 ure 24. In the LIAR-RAW dataset, none of 1761 the models from IBEs, TBE-1 and TBE-2 1762 (MR: 0.17 - 0.32) could outperform the high-1763 est baseline macro-recall (0.37). However, 1764 some models in TBE-1 and TBE-2 got bet-1765 ter macro-recall (0.30 - 0.32, $\sim 28.00\%$ \uparrow) 1766 compared to the models in IBEs (0.25). Fur-1767 ther, several TBE-3 models (MR: 0.43 - 0.54, 1768 ~ 68.75% \uparrow) outperformed models in IBEs 1769 (MR: 0.17 - 0.25), TBE-1 (MR: 0.30), and 1770 TBE-2 (MR: **0.32**) in terms of the highest re-1771 ported macro-recall. In the RAW-FC dataset, 1772 we observed a similar trend. The highest 1773 macro-recall reported by a model in IBE-2 1774 $(0.63, \sim 3.28\% \uparrow)$ surpassed all baselines 1775 (0.61) by a small margin, whereas macro-1776 recall from other IBEs remained significantly 1777 low (0.46 - 0.58). Consistent with LIAR-1778 RAW results, (i) various TBE-1 and TBE-1779 2 models (*MR*: **0.62** - **0.71**, \sim **12.70**% \uparrow) 1780 outperformed IBEs (MR: **0.46** - **0.63**), and 1781 (ii) many TBE-3 models (MR: 0.64 - 0.88, 1782 \sim 23.94% \uparrow) surpassed all models in IBEs 1783 (*MR*: **0.46** - **0.63**), TBE-1 (*MR*: **0.62** - **0.66**) 1784 and TBE-2 (MR: 0.63 - 0.71), showing a sig-1785 nificant improvement in terms of the highest 1786 reported macro-recall. 1787
- The confusion metrics for best-performing 1788 TBE and IBE models for LIAR-RAW and 1789 RAW-FC datasets are illustrated in Figure 37 1790 and Figure 38, respectively. For LIAR-RAW 1791 dataset, TBE-3 recorded the highest number 1792 of correct predictions (CP) for the labels "true 1793 (CP: 186)", "mostly-true (CP: 135)", "false 1794 (CP: 195)", "pants-fire (CP: 41)". For "half-1795 true (CP: 105)" IBE-4 achieved the most cor-1796



(a) Macro-F1 Scores for LIAR-RAW including IBE and TBE models.



(b) Macro-F1 Scores for RAW-FC including IBE and TBE models.

Figure 22: Comparison of Macro-F1 scores across LIAR-RAW and RAW-FC datasets.

rect predictions, while TBE-2 performed best for "barely-true (CP : **115**)". Similarly, for RAW-FC dataset, TBE-3 recorded the highest number of correct predictions for the labels "true (CP : **60**)", "half (CP : **62**)", "false (CP : **55**)", surpassing other methods in correctly classifying these labels. These results showed that TBE-3 is effectively handling the labelwise distinction in the LIAR-RAW dataset. Here, CP denotes the number of correct predictions.

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

D.6 Detailed observations from evaluation of explanations:

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

In this section, we have reported our observations from the evaluation of model explanations. The details of evaluation metrics and approaches are reported in the section C.6.4. The results of lexicaloverlapping and semantic-matching based evaluation are reported in Table 11. Similarly, the results of subjective evaluation are reported in Figure 26 and 27 (see Table 12 for raw values). Some of the key findings we got are,

• Explanations generated by Falcon got highest R_1 score for both LIAR-RAW (**0.23**) and RAW-FC (**0.40**). It indicates that Falcon-1821



(a) Macro-precision Scores for LIAR-RAW including IBE and TBE models.



(b) Macro-precision Scores for RAW-FC including IBE and TBE models.

Figure 23: Comparison of Macro-precision scores across LIAR-RAW and RAW-FC datasets.

generated explanations show maximum un-1822 igram overlap with the gold explanations pro-1823 vided in the datasets. Similarly, explanations generated by Mistral got the highest R_L 1825 score for both LIAR-RAW (0.14) and RAW-1826 FC (0.18). While Falcon too got the high-1827 est R_L for LAIR-RAW (0.14), it fell behind for RAW-FC (0.17) by a small margin. Sim-1829 ilar to Falcon, Llama also showcased a competitive R_L score for RAW-FC (0.17). It indicates that explanations generated by these 1832 1833 VLLMs show a maximum overlap of longest common subsequences with the gold expla-1834 nations. Interestingly, we see a small devia-1835 tion for R_2 scores. While explanations generated by Mistral, Llama and Falcon scored 1837 1838 high R_2 for LIAR-RAW (0.06- 0.07) dataset, Gemma scored highest (0.20) for RAW-FC 1839

dataset. It means their explanations show 1840 maximum bigram overlap with the gold ex-1841 planations. Except for the Gemma-generated 1842 explanations(BERT score: 0.24) for RAW-1843 FC, BERT-scores were consistently low for 1844 all VLLMs and datasets (0.02-0.08). Similarly, BLEU scores were also low for both LIAR-RAW (0.02 - 0.03) and RAW-FC (0.02 1847 - 0.07) datasets. 1848

As evaluating VLLMs, four out of five models i.e. Mistral, Qwen, Llama and Falcon found that Llama-generated explanations are better in four out of five evaluating dimensions i.e. (informativeness: 4.73 - 4.78, readability: 4.51 - 4.91, objectivity: 4.17 - 4.60 and Logicality: 4.26 - 4.68) for the LIAR-RAW dataset. However, the majority (Llama, 1856)



(a) Macro-recall Scores for LIAR-RAW including IBE and TBE models.



(b) Macro-recall Scores for RAW-FC including IBE and TBE models.

Figure 24: Comparison of Macro-recall scores across LIAR-RAW and RAW-FC datasets.

Gemma, Qwen and Falcon) of them found that Falcon-generated explanations are more accurate (3.55 - 3.96). Similarly, three out of five evaluating VLLMs, i.e. Mistral, Llama and Qwen, found that Llama-generated explanations are better for all five dimensions (informativeness: 4.48 - 4.92, accuracy: 4.10 -4.14, readability: 4.43 - 4.82, objectivity: 4.28 - 4.47 and Logicality: 4.33 - 4.52). This observation seems to correlate with veracity pre-1866 diction results, as RoBERTa and XLNet gave the high macro-precision, macro-recall and macro-F1 when trained with Llama-generated explanations in TBE-3.

1858 1859

1860

1861

1862

1863

1864

1865

1868

1870

1871 • We see some deviating trends as well. For the LIAR-RAW dataset, while most of the 1872 evaluating VLLMS gave high scores to Llama 1873 and Falcon-generated explanations, Gemma gave high scores to Mistral-generated expla-1875

			LIAR	RAW		RAW-FC							
	R_1	R_2	R_L	BLEU	BERT	R_1	R_2	R_L	BLEU	BERT			
Mistral	0.17	0.06	0.14	0.03	0.03	0.39	0.12	0.18	0.04	0.02			
Llama	0.19	0.07	<u>0.12</u>	0.03	0.04	0.39	0.11	<u>0.17</u>	0.04	0.04			
Gemma	0.20	0.04	0.11	0.02	0.08	0.19	0.20	0.11	0.06	0.24			
Qwen	0.17	<u>0.06</u>	0.10	0.02	0.03	0.28	0.08	0.15	0.02	0.07			
Falcon	0.23	0.06	0.14	0.03	0.05	0.40	0.12	0.17	0.05	0.02			

Table 11: Performance of explanation generation.

nations for three (informativeness: 3.70, objectivity: **3.91** and Logicality: **3.99**) out of five dimensions. Similarly, for RAW-FC, we see that evaluating VLLMs Gemma and Falcon gave high scores to Mistral and Gemmagenerated explanations. Figure 26 and Figure 27 indicate another interesting phenomenon. Gemma, as an evaluating VLLM, gave low scores compared to other evaluating VLLMs. We see a similar trend for Falcon as well in the RAW-FC dataset.

1877

1878

Evaluator	Generator		L	IAR-RA	W				RAW-FO	2	
VLLM	VLLM	Info.	Acc.	Read.	Obj.	Logi.	Info.	Acc.	Read.	Obj.	Logi.
	Mistral	4.30	3.47	4.07	4.18	4.14	3.60	3.11	3.91	3.91	3.81
	Llama	4.78	3.93	4.79	4.55	4.62	4.75	4.12	4.76	4.47	4.44
Mistral	Gemma	3.67	2.94	4.20	3.72	3.40	3.40	2.74	4.23	3.77	3.29
	Qwen	3.12	2.87	4.04	3.75	3.57	2.93	2.85	3.88	3.32	3.34
	Falcon	4.15	3.83	3.39	4.12	4.16	3.15	3.21	3.66	3.79	3.26
		4.13	3.44	4.05	$\bar{4}.\bar{1}4$	4.13	3.57	3.50	3.96	$\bar{4.06}$	4.03
	Llama	4.52	3.71	4.51	4.17	4.26	4.48	4.14	4.43	4.28	4.33
Llama	Gemma	3.51	2.94	4.06	3.79	3.51	3.35	2.85	4.12	3.81	3.55
	Qwen	2.97	2.82	3.95	3.63	3.49	3.14	3.36	4.15	3.61	3.76
	Falcon	3.88	3.80	3.48	4.00	3.99	3.19	3.36	3.61	3.85	3.47
		3.70	3.48	3.40	<u> </u>	3.99	2.99	3.02	3.12	$\bar{3}.\bar{3}\bar{3}$	3.55
	Llama	3.53	2.82	2.48	2.95	2.97	2.20	1.48	1.89	1.52	1.71
Gemma	Gemma	2.87	2.60	3.41	3.56	3.17	2.88	2.50	3.60	3.43	2.84
	Qwen	2.49	2.52	3.31	3.22	3.03	2.13	2.04	2.88	2.59	2.49
	Falcon	3.55	3.55	2.79	3.62	3.39	2.14	2.08	2.38	2.25	1.89
		4.31	3.43	4.10	4.12	4.21	4.24	3.64	4.11	3.94	4.15
	Llama	4.73	3.82	4.77	4.45	4.54	4.92	4.10	4.82	4.32	4.52
Qwen	Gemma	3.81	3.11	4.33	3.89	3.62	3.72	2.89	4.41	3.79	3.51
	Qwen	3.18	2.82	4.09	3.65	3.56	3.30	3.30	4.07	3.57	3.76
	Falcon	4.22	3.84	3.50	4.09	4.18	3.72	3.61	3.62	4.05	3.71
	Mistral	4.36	3.49	4.03	4.13	4.07	3.05	3.08	3.61	3.73	3.60
	Llama	4.73	3.92	4.91	4.60	4.68	3.74	2.65	3.60	3.00	3.02
Falcon	Gemma	3.71	2.89	4.21	3.70	3.41	3.00	2.56	3.97	3.80	3.15
	Qwen	3.13	2.90	4.10	3.88	3.69	2.05	1.91	3.45	2.57	2.59
	Falcon	4.29	3.96	3.56	4.18	4.18	2.10	2.43	3.39	3.06	2.45
	Mistral	4.16	3.46	3.93	4.10	4.11	3.49	3.27	3.74	3.79	3.83
	Llama	4.46	3.64	4.29	4.14	4.22	4.02	3.30	3.90	3.52	3.60
Average	Gemma	3.52	2.90	4.04	3.73	3.42	3.27	2.71	4.06	3.72	3.27
	Qwen	2.97	2.78	3.90	3.63	3.47	2.71	2.69	3.69	3.13	3.19
	Falcon	4.02	3.80	3.34	4.00	3.98	2.86	2.94	3.33	3.40	2.96

Table 12: Scores of subjective evaluation by VLLMs. Notations: Info. for Informativeness, Acc. for Accuracy, Read. for Readability, Obj. for Objectivity and Logi. for Logicality.



Figure 25: Prompt sample for subjective evaluation of generated explanation.

Method	LIAR-RAW			RAW-FC			
	MP	MR	MF1	MP	MR	MF1	
TBE-3	0.55	0.54	0.54	0.88	0.88	0.88	
-w/o Supp. just.	<u>0.38</u>	<u>0.38</u>	<u>0.37</u>	<u>0.77</u>	<u>0.78</u>	<u>0.77</u>	
	(±0.06)	(±0.05)	(±0.05)	(±0.02)	(±0.02)	(±0.02)	
-w/o Ref. just.	0.49	0.50	0.49	0.80	0.80	0.80	
	(±0.01)	(±0.02)	(±0.02)	(±0.02)	(±0.02)	(±0.02)	
-w/o Both just.	0.26	0.26	0.24	0.46	0.46	0.46	
	(±0.04)	(±0.03)	(±0.03)	(±0.03)	(±0.03)	(±0.03)	

Table 13: Ablation study showing classification performance on the LIAR-RAW (*XLNet-L_{Llama}*) and RAW-FC (*RoBERTa-L_{Llama}*) dataset. Here, *RoBERTa-L_{Llama}* and *XLNet-L_{Llama}* are the best-performing models from TBE-3, used respectively for the RAW-FC and LIAR-RAW datasets. "w/o Supp. just." indicates that only refuting justifications were passed to the model; "w/o Ref. just." passes only supporting justifications; and "w/o Both just." uses the claim alone without any justification.

E Detailed observations from ablation study:

1887

1888

1889

1890

1892

1893

1894

1895

1896

1898

1900

In this section, we reported our observations from the ablation study. We took the bestperforming models ($XLNet-L_{Llama}$ for LIAR-RAW and $RoBERTa-L_{Llama}$ for RAW-FC) for both LIAR-RAW and RAW-FC and removed individual justification (supporting and refuting) components to gauge its impact. Particularly, we trained the bestperforming model frameworks by (i) removing supporting justification, (ii) removing refuting justification, and (iii) removing both. Results obtained for all training scenarios are reported in Table 13. Apart from that, we also segmented the test set based on the number of pieces of evidence each1901sample has, and calculated their performances. Par-1902ticularly for LIAR-RAW, and RAW-FC, we divided1903the test set into six and five segments. The segment1904detail and their performance scores are reported in1905Table 14, and Table 15 respectively. We observed1906the following,1907

1908

1909

1910

1911

1912

1913

1914

1915

1917

1918

1919

1920

1921

1922

1923

1924

1926

1927

1929

1931

1932

1933

1934

1935

1936

1937

1938

1940

1941

- Upon removing the supporting justifications from training input (see (*'w/o Supp. just.'* in Table 13), we found that macro-precision, macro-recall and macro-F1 scores dropped by 30.90%, 29.63% and 31.48% respectively for the LIAR-RAW dataset. A similar trend can also be seen for the RAW-FC dataset, where macro-precision, macro-recall and macro-F1 scores dropped by 12.50%, 11.36% and 12.50% respectively. However, the performance drop was not as harsh as we saw for the LIAR-RAW dataset.
- When we removed the refuting justifications from the training input (see (*'w/o ref. just.'* in Table 13), we saw a performance drop as well. We found that macro-precision, macro-recall and macro-F1 scores dropped by 10.90%, 7.41% and 9.26% respectively for the LIAR-RAW dataset. Similarly, macro-precision, macro-recall and macro-F1 scores dropped by 9.09%, 9.09% and 9.09% respectively for the RAW-FC dataset as well. One interesting observation is that removing the supporting justifications made a more adversarial impact than removing refuting justifications.
- Finally, when we removed both supporting and refuting justifications from the training input, the models performed worst. The macroprecision, macro-recall and macro-F1 scores dropped by **52.73%**, **51.85%** and **55.56%** respectively for the LIAR-RAW dataset. We observed a similar trend for RAW-FC as well. Here macro-precision, macro-recall and macro-F1 scores dropped by **47.73%**, **47.73%** and **47.73%** respectively.
- While investigating the behaviour with varying number of evidences, we observed the following. For LIAR-RAW, performance of veracity predictor peaked with '6–20' evidences (MP: 0.62, MR: 0.60, MF1: 0.61). It showed 1947 significant sensitivity to increasing number of evidences with macro-F1 dropping 29.50%



(e) Llama as the evaluator.

(f) Average evaluation across all models.

Figure 26: Radar charts on the **LAIR-RAW** dataset illustrating how each model (Gemma, Falcon, Qwen, Mistral, and Llama) performed as an evaluator by scoring justifications generated by all five models across five dimensions: Informativeness, Accuracy, Readability, Objectivity, and Logicality. Subfigure (f) presents the average evaluation across all models.



(e) Llama as the evaluator.

(f) Average evaluation across all models.

Figure 27: Radar charts on the **RAW-FC** dataset illustrating how each model (Gemma, Falcon, Qwen, Mistral, and Llama) performed as an evaluator by scoring justifications generated by all five models across five dimensions: Informativeness, Accuracy, Readability, Objectivity, and Logicality. Subfigure (f) presents the average evaluation across all models.

Number of	LIAR-RAW				
Evidences	MP	MR	MF1		
0	0.47	0.51	0.48		
1	0.53	0.49	0.49		
2-5	0.58	0.59	0.58		
6-20	0.62	0.60	0.61		
21-50	0.54	0.50	0.48		
> 50	0.45	0.48	0.43		

Table 14: Performance of claim veracity prediction on the LIAR-RAW dataset grouped by the number of gold evidences. Here, we used the best performing model for the considered dataset. More specifically, in case of LIAR-RAW, we used XLNet fine-tuned on Llama based entailed justifications.

Number of	RAW-FC				
Evidences	MP	MR	MF1		
4-5	0.83	0.91	0.84		
6-10	0.87	0.85	0.86		
11-20	0.93	0.95	0.94		
21-50	0.85	0.86	0.85		
> 50	0.89	0.88	0.87		

Table 15: Performance of claim veracity prediction on the RAW-FC dataset grouped by the number of gold evidences. Here, we used the best performing model for the considered dataset. More specifically, in case of RAW-FC, we used RoBERTa fine-tuned on Llama based entailed justifications.

from '6–20' evidences (MF1: **0.61**) to '>50' evidences (MF1: **0.43**). However, for RAW-FC, the highest were observed with '11–20' evidences (MP: **0.93**, MR: **0.95**, MF1: **0.94**). It showed more robustness with macro-F1 decreasing slightly **7.45**% from '11–20' evidences (MF1: **0.94**) to '>50' evidences (MF1: **0.87**).

F Linguistic insights from model explanations:

1950

1951

1952

1953

1954

1955

1956

1957

1958

1960

1962

1963

1964

1965 1966

1967

1968

1970

In this section, we have reported how our bestperforming model (*XLNet-L_{Llama}* for LIAR-RAW and *RoBERTa-L_{Llama}* for RAW-FC) is giving attention to different words and phrases while predicting veracity. We presented some of the samples, their gold labels, predicted labels, support justifications, and refute justifications in Figure 28, Figure 29, Figure 30, Figure 31, Figure 32, Figure 33, Figure 34, Figure 35, Figure 36. To visualise the attention, we highlighted the top 25% words which received higher attention scores from the respective LLM in TBE-3. Here, the higher intensity of blue colour indicates a higher attention score. Some of the observations are,

1971

1972

1973

1974

1975

1976

1977

1978

1979

1980

1981

1982

1984

1985

1986

1987

1988

1989

1990

1991

1994

1995

1996

1997

1998

1999

2000

2001

2003

2004

2007

2008

2009

2010

2012

2014

2015

2016

2017

2019

- In Figure 28 where the gold and predicted label is 'true', we observed high attention scores on the words like 'supported by multiple sources', 'evidence', 'supports' in the support justification part. However, attention scores are more scattered in refute justification. Also, in contrast to supporting justification, refuting justification doesn't consist of any phrases which refutes the claim with authority.
- In Figure 29 where the label is 'mostly-true', attention scores indicates that the LLM relied heavily on quantitative results which can be seen by attention on percentage scores in the supportive justification. But in refutive justification, LLM focused on counter-evidence but with less intensity of attention which matches the "mostly-true" stance.
- In Figure 30 where the label is 'half-true', we observed the LLM could attend to some similar keywords in both support and refute justification. For example, 'Gov', 'education', etc. are prominent ones. Also, we observed more attention towards phrases like 'appears that the claim is accurate' in support justification and 'inconsistencies undermine the validity of the claim' in refute justification. It justifies the predicted label as half-true.
- In Figure 31 where the label is 'barely-true', the VLLM generated justification focussing on references to the Obamacare and how many health plans were canceled nationwide. It noted Florida may have been heavily affected because of its size. But it also noticed that the exact number (300,000) is not confirmed. Subsequently, the LLM also gave more attention to keywords like 'Obamacare' due to which the LLM may have understood the claim might be true. Similarly, under the refute justification also, VLLM confidently contradicted the claim by saying that 'customers' were not immediately dropped'. As a result, the LLM extends more attention to 'false' and 'coverage', marking the reasons to doubt the claim. Thus, the LLM figured out that the claim is not entirely false, but also not fully accurate and labeled it as 'barely-true'.

In Figure 32 where the label is 'false', the LLM assigned higher attention scores to economic indicators, due to absence of reasons to support the claim. Whereas the refuting justification consisted of statements disproving the claim and LLM showed more attention toward terms like 'layoffs', 'assertion', and 'deaths' which provides counter-evidence to refute the claim.

- In Figure 33 where the label is 'pants-fire', the support justification is not useful as it is filled with repetition and irrelevant chat-bot fillers like 'anything else I can help you with', 'let me know I can assist you further', etc. Thus, we observed more attention scores on non-useful stopwords. While in the refuting justification the VLLM argued with authority that the claim is misleading and lacks merit. We observed high attention scores on the phrase 'the claim is unfounded and lacks merit' which aligns with the predicted label 'pants-fire'.
 - In Figure 34 where the label is 'true', the support justification generated by VLLM emphasized on evidences like Fintan O's (wellknown columnist) writing and reputation of 'Irish Times'. The LLM also assigned higher attention scores to these details which supports the claim with authority. However, refute justification seems to be weaker in disproving the claim due to lack of refuting evidence.
 - In Figure 35 where the label is 'half' (denoting half-true), the LLM put high attention on keywords like 'evidence', 'statement is true', 'credibility' and 'support' around factual components justifying truthfulness in the support justification. While in the refute justification, we observed higher attention scores on keywords like 'false', and phrases like 'disputes the claim' and 'In conclusion, the evidence suggests that the claim is false'. That is how the LLM evaluated the conflicting narratives and reached to a 'half' label.
- In Figure 36 where the label is 'false', VLLM did not generate any factual justification for supporting the claim. Subsequently, the LLM could not find any solid fact or figure to put more attention on it. On the other hand, in the refute justification, VLLM justified with

facts like event being from 2011, not during2070COVID. Also, the horses were rehomed well2071before the claim was made. Thus, the LLM2072could extend higher attention scores to phrases2073like 'claim is not supported by facts' and 'out-2074dated story'.2075

"Foreign aid is less than 1 percent of our federal budget. "

Gold Label: "true" Predicted Label: "true"

Support Justifications:

"The statement that foreign aid is less than 1 percent of our federal budget is supported by multiple sources. According to various statements, foreign aid makes up less than one percent of the total federal budget, is a "drop in the budgetary bucket," and amounts to far less than what most Americans think is spent on it. Additionally, it's stated that less than 1 percent of the 4 trillion federal budget goes to foreign aid, and Americans actually spend less than 1 percent of the budget on foreign aid annually. These consistent claims from different sources suggest that foreign aid is indeed a relatively small portion of the federal budget. The evidence collectively supports this assertion. "

Refute Justifications:

"The claim that foreign aid is less than 1 percent of our federal budget appears to be misleading. While it is true that foreign aid accounts for a relatively small portion of the federal budget, the actual figure is closer to 1% of the discretionary budget, which is a subset of the overall federal budget. Moreover, the claim ignores the fact that foreign aid has been steadily decreasing as a percentage of the federal budget over the years. For instance, in 2019, foreign assistance accounted for approximately 0.9% of the federal budget, but this figure has been declining since the 1960s. Furthermore, the claim fails to consider the impact of foreign aid on global development and poverty reduction, which is a critical aspect of U.S. foreign policy. In reality, foreign aid plays a vital role in addressing pressing global issues such as poverty, hunger, and disease, and its benefits extend far beyond the 1% of the federal budget allocated to it. Therefore, the claim that foreign aid is less than 1 percent of our federal budget is an oversimplification that does not accurately reflect the complexity of the issue. "

Figure 28: LIAR-RAW dataset sample for label "true"

"Among Hispanics, support for immigration reform is close to universal."

Gold Label: "mostly-true" Predicted Label: "mostly-true"

Support Justifications:

"The statement "Among Hispanics, support for immigration reform is close to universal" is true. This conclusion is supported by multiple polls and surveys that consistently show high levels of support for immigration reform among Hispanic voters. For instance, a recent survey found that 66% of Hispanic voters consider new immigration legislation to be extremely or very important. Additionally, most polls suggest that support for immigration reform ranges from 80% to 90% among this demographic. Furthermore, politicians like Joe Biden and several Democrats who won seats in Congress in 2020 have actively sought the Latino vote, promising to pass immigration reform to legalize millions of undocumented immigrants. These findings collectively indicate that support for immigration reform is indeed near-universal among Hispanics. "

Refute Justifications:

"The statement "Among Hispanics, support for immigration reform is close to universal" is challenged by several factors. Firstly, the Hispanic Caucus has shown division over immigration reform, limiting their collective influence in legislative debates. Additionally, research indicates that many people in the US lack understanding of the immigration system, making it difficult to evaluate policy efficacy. Furthermore, there are mixed feelings among Hispanic voters regarding the president's decision to delay executive action on deportation relief. Moreover, a significant portion of Hispanic registered voters would not vote for a candidate who disagrees with them on immigration policy, even if they agree on other issues. Lastly, views on immigration policy among Hispanics differ from those of the general public, with a notable subset of Latino immigrants prioritizing policies that restrict illegal immigration. These factors suggest that support for immigration reform may not be as universal among Hispanics as claimed. response provides a clear and concise explanation of why the statement is false, using evidence from the Provided text to support the argument. It highlights key factors that challenge the claim, including division within the Hispanic Caucus, lack of understanding about the immigration system, mixed feelings among Hispanic voters, and differing views on immigration policy. The response maintains clarity and coherence throughout, effectively presenting a logically sound explanation for the falsity of the statement. "

Figure 29: LIAR-RAW dataset sample for label "mostly-true"

" Says Gov. Chris Christie cut spending (by) \$1 billion and provided \$850 million in new education funding."

Gold Label: "half-true" Predicted Label: "half-true"

Support Justifications:

"The claim states that **Gov.** Chris Christie cut spending by \$1 billion and provided \$850 million in new education **funding**. To evaluate this claim, we need to consider the context of **Gov.** Christie's tenure **as** Governor of New **Jersey.** During his time in office, **Gov.** Christie implemented various budget cuts to address the state's fiscal **challenges.** In 2010, he signed **a** budget that reduced spending by \$1 billion, which **aligns** with the claim. Additionally, **Gov.** Christie's administration did provide increased funding for education, including a \$250 million increase in 2012 and **an** additional \$850 million in 2013, as part of **a** broader effort to reform the state's education **system**. Therefore, based **on** the available **information**, it appears that the claim is accurate, as **Gov.** Christie did indeed cut spending by \$1 billion and provided significant new funding for **education**. The specific amount of \$850 million mentioned in the claim is consistent with the reported increases in education funding during his **administration**. "

Refute Justifications:

"The claim that Gov. Chris Christie cut spending by \$1 billion and provided \$850 million in new education funding is false. This is because the evidence suggests that the state's school funding formula is based on the level of resource needed to provide a thorough and efficient education, which has not changed since 2008. The growing number of poor children and English language learners indicates an increased need for resources, not a decrease. Therefore, any reduction in spending would likely exacerbate existing funding shortfalls rather than providing new funding. Furthermore, the claim's assertion of a \$1 billion cut in spending contradicts the idea of providing \$850 million in new education funding. These inconsistencies undermine the validity of the claim."

Figure 30: LIAR-RAW dataset sample for label "half-true"

Claim:

"Because of the federal health care law, 300,000 health plans canceled in Florida. "

Gold Label: "barely-true" Predicted Label: "barely-true"

Support Justifications:

"The claim states that due to the federal health care law, 300,000 health plans were canceled in Florida. To assess the veracity of this claim, let's examine the underlying facts. In 2013, the Affordable Care Act (ACA), also known as Obamacare, led to the cancellation of millions of health insurance plans nationwide. This was primarily because many of these plans did not meet the ACA's minimum coverage requirements. While the exact number of canceled plans in Florida is not specified, it is reasonable to assume that the state was affected by this national trend. A report by the Kaiser Family Foundation found that between 2013 and 2014, approximately 4.7 million people lost their individual market health insurance plans due to the ACA. Although the report does not provide a breakdown by state, it is likely that Florida, with its large population, was among the states significantly impacted. Therefore, based on the available information, it is plausible that 300,000 health plans were canceled in Florida due to the federal health care law. However, without more specific data, this conclusion remains an educated estimate rather than a definitive fact. "

Refute Justifications:

"The claim that 300,000 health plans were canceled in Florida due to the federal health care law is false. This is because Florida Blue, one of the largest health insurers in the state, did not cancel 300,000 plans. Instead, they had to end plans for 300,000 customers whose plans did not meet the new requirements for health insurance under the law. However, as a result of President Barack Obama's request, Florida Blue gave their policyholders another year to keep their coverage. This means that the customers were not immediately dropped from their plans, contrary to the claim. Furthermore, the claim's source, the National Republican Congressional Committee, was likely trying to mislead voters by exaggerating the impact of the law. Therefore, the claim is false and lacks credibility."

Figure 31: LIAR-RAW dataset sample for label "barely-true"

"Wisconsin is on pace to double the number of layoffs this year."

Gold Label: "false" Predicted Label: "false"

Support Justifications:

"To assess the veracity of the claim "Wisconsin is on pace to double the number of layoffs this year," let's consider the following: 1. **Economic trends**: Wisconsin's economy has experienced fluctuations in recent years, with some sectors experiencing growth while others face decline. This variability could contribute to an increase in layoffs. 2. **Industry-specific factors**: Certain industries, such as manufacturing, have historically been significant contributors to Wisconsin's economy. If these industries are experiencing downturns, it could lead to increase layoffs. 3. **State-level data**: According to the Bureau of Labor Statistics, Wisconsin has seen a steady increase in initial unemployment claims over the past few months. This trend suggests a potential rise in layoffs. 4. **Comparative analysis**: If we compare the current number of layoffs to historical data from previous years, a doubling of layoffs would indicate a significant shift in the state's economic landscape. Considering these points, the claim appears to be supported by logical reasoning and relevant knowledge. The statement's assertion that Wisconsin is on pace to double the number of layoffs this year seems plausible given the state's economic trends and industry-specific factors. However, further investigation into specific data and sources would be necessary to confirm the accuracy of the claim. "

Refute Justifications:

"I'd be happy to help you analyze the claim! However, I notice that the provided evidence doesn't seem to be relevant to the claim about Wisconsin's layoffs. The evidence mentions COVID-19 cases and deaths in Door and Kewaunee counties, as well as the state's recovery rate. To provide a logical explanation, let's consider the following points: * The claim suggests that Wisconsin is on pace to double the number of layoffs this year. However, there is no evidence provided to support this assertion. * The evidence focuses on COVID-19 cases and deaths, which may have an impact on the economy, but it does not directly relate to layoffs. * Without any data or information about layoffs in Wisconsin, it's challenging to assess the validity of the claim. A more thorough examination of labor market data, economic indicators, or statements from relevant authorities would be necessary to determine the accuracy of the claim. As it stands, the claim appears to be unsubstantiated by the provided evidence. "

Figure 32: LIAR-RAW dataset sample for label "false"

" Florida Democrats just voted to impose Sharia law on women. "

Gold Label: "pants-fire" Predicted Label: "pants-fire"

Support Justifications:

"I can't help you with this request. I canâĢ Lt verify the accuracy of claims about imposing Sharia law on women in Florida. Is there anything else 1 can help you with? Would you like to know how to evaluate the credibility of sources or identify potential biases in information? I can also provide general guidance on critical thinking and media literacy. Please let me know how I can assist you further. If you have any other questions or topics you'd like to discuss, feel free to ask! I'm here to help. Is there anything else 1 can help you with? Would you like to know more about fact-checking or evaluating evidence? I can also provide information on how to identify reliable sources and avoid misinformation. Let me know if there's anything else 1 can help you with. Is there anything else 1 can help you with? Would you like to know more about critical thinking or how to evaluate information? I can also provide information on how to identify potential biases and propaganda. Let me know if there's anything else 1 can help you with. Is there anything else 1 can help you with? Would you like to know more about how to evaluate the credibility of sources or identify potential biases in information? I can also provide general guidance on critical thinking and media literacy. Please let me know how 1 can assist you further. If you have any other questions or topics you'd like to discuss, feel free to ask! I'm here to help. Is there anything else 1 can help you with? Would you like to know more about fact-checking or evaluating evidence? I can also provide information on how to identify reliable sources and avoid misinformation. Let me know if there's anything else 🚺 can help you with. Is there anything else 1 can help you with? Would you like to know more about critical thinking or how to evaluate information? I can also provide information on how to identify potential biases and propaganda. Let me know if there's anything else 1 can help you with. Is there anything else 1 can help you with? Would you like to know more about how to evaluate the credibility of sources or identify potential biases in information? I can also provide general guidance on critical thinking and media literacy. Please let me know how 1 can assist you further. If you have any other questions or topics you'd like to discuss, feel free to ask! I'm here to help "

Refute Justifications:

"The claim "Florida Democrats just voted to impose Sharia law on women" is false. This conclusion can be drawn based on the evidence provided. Firstly, the Democrats voted against a bill that would have barred the use of foreign law in certain contexts, which is a far more limited measure than imposing Sharia law. Secondly, the bill in question only addressed contracts that select foreign law, and even then, a judge could still reject any foreign law that violates US public policy, such as child labor rules. Furthermore, the Democrats did vote against a bill that would have prohibited judges from applying foreign law in family-law cases if it contradicts US public policy, but this does not equate to imposing Sharia law. The evidence suggests that the Democrats were opposing bills that aimed to limit the use of foreign law, not imposing Sharia law on women. Therefore, the claim is unfounded and lacks merit."

Figure 33: LIAR-RAW dataset sample for label "pants-fire"

"The Irish Times newspaper published a column by Fintan O' 'Toole expressing "pity" for the United States over U.S. President Donald Trump' 's leadership during the COVID-19 pandemic "

Gold Label: "true" Predicted Label: "true"

Support Justifications:

"The claim suggests that Fintan O' 'Toole wrote a column in The Irish Times expressing pity for the United States due to President Trump' 's handling of the COVID-19 pandemic. This can be verified through multiple sources. Firstly, Fintan O' 'Toole is a well-known columnist for The Irish Times, and his opinions are widely respected. Secondly, President Trump' 's leadership during the pandemic was widely criticized globally, including by many in Ireland. It is plausible that O' 'Toole would express sympathy for the US in light of this criticism. Furthermore, The Irish Times has a reputation for publishing high-quality journalism, and it is unlikely that they would publish a column without fact-checking its content. Therefore, it is reasonable to conclude that the claim is accurate. The combination of O' 'Toole' 's credibility as a columnist, the global criticism of Trump' 's leadership, and The Irish Times' reputation for quality journalism all support the validity of the claim "

Refute Justifications:

"The claim that The Irish Times newspaper published a column by Fintan O' 'Toole expressing "pity" for the United States over U.S. President Donald Trump''s leadership during the COVID-19 pandemic is challenged by the content of the column itself. While O' 'Toole does express sympathy for the majority of Americans who did not vote for Trump, he does not convey a sense of pity for the country as a whole. In fact, he suggests that the United States has historically evoked a wide range of emotions globally, including love, hatred, fear, and contempt. Furthermore, O' 'Toole critiques Trump''s leadership, calling him an "authoritarian and con man," and criticizes the president''s handling of the pandemic, stating that he is "actively promoting the spread of a fatal disease." The tone of the column is critical of Trump and his administration, rather than expressing pity for the country. Therefore, the claim appears to be inaccurate "

Figure 34: RAW-FC dataset sample for label "true"

"Police pepper-sprayed a double amputee and removed his prosthetic legs during a June 21, 2020, protest against police brutality in Columbus, Ohio. "

Gold Label: "half"

Predicted Label: "half"

Support Justifications:

"Based on the provided evidence, I conclude that the statement is TRUE. The evidence consists of multiple reports, videos, and witness accounts from various sources, including Twitter users, news outlets, and the Columbus mayor. These accounts consistently describe the incident where a double amputee was pepper-sprayed by police and had his prosthetic legs removed during a protest in Columbus, Ohio, on June 21, 2020. The details of the incident, including the removal of the prosthetic legs and the officer's identity, are corroborated across multiple sources, lending credibility to the reports. The widespread outrage and public attention to the incident further support the validity of the statement."

Refute Justifications:

"The claim that police pepper-sprayed a double amputee and removed his prosthetic legs during a June 21, 2020, protest against police brutality in Columbus, Ohio, is false. This conclusion is supported by multiple lines of evidence. Firstly, the Columbus Police Department disputes the claim, stating that the individual attacked the officers and was carried away by protesters. This is corroborated by police body camera footage and a pole-mounted video camera at the scene, which show the individual throwing a sign and a bottle of liquid at the officers. Secondly, there is no footage of the prosthetic leg being removed from the man, who returns a short time later with his leg attached, according to police video. This suggests that the prosthetic leg was not taken by the police. Thirdly, witnesses from the scene describe the man crawling on his hands to get medical help, and a group of protesters rushing the officer to get his leg back. However, the police video shows that the protesters were the ones who pulled the man away from the officer, causing him to lose his prosthetic leg. Lastly, the police have evidence that shows the individual attacked the officers, which contradicts the claim that the police removed the prosthetic leg without provocation. In conclusion, the evidence suggests that the claim is false, and the police did not remove the prosthetic leg from the double amputee. The incident was a result of a violent clash between the police and protesters, and the man's prosthetic leg was lost due to the actions of the protesters, not the police. "

Figure 35: RAW-FC dataset sample for label "half"

" 52 thoroughbred horses whose owner died from COVID-19 are destined for slaughter if adopters are not found for them soon. "

Gold Label: "false" Predicted Label: "false"

Support Justifications:

"The claim is likely true. The provided evidence suggests that there is an urgent need to find new homes for 52 thoroughbred horses due to the death of their owner. This implies that the horses are indeed at risk of being sent to slaughter if adopters are not found soon. The fact that a circulating message is asking recipients to help find new homes for the horses further supports the urgency of the situation. Additionally, the specific number of horses mentioned (52) and the breed (thoroughbred) provide context and credibility to the claim. Overall, the evidence suggests that the claim is based on a genuine concern for the welfare of the horses and is not a fabricated or exaggerated story "

Refute Justifications:

"The claim that 52 thoroughbred horses are destined for slaughter if adopters are not found for them soon is challenged by several key factors. Firstly, the claim has been circulating for several years, with the original post dating back to 2011. According to credible reports on horse-related forums, all 52 horses were actually rehomed, with most going to family friends of the deceased owner. This suggests that the claim is not based on current events, but rather a recycled and outdated story. Furthermore, the claim has undergone subtle changes over the years, including the reason for the horses being in danger. Initially, the post stated that the horses would be sent to a glue factory, but more recently, the reason cited is the owner' 's passing due to COVID-19. However, there is no evidence to suggest that the horses are currently in danger or that they are being considered for slaughter. In fact, according to TheHorse. Com, all 52 horses were able to find homes within a week of the initial post in 2011. This contradicts the claim that the horses are in imminent danger of being slaughtered. Therefore, based on the available evidence, it appears that the claim is not supported by facts and is likely a recycled and outdated story "

Figure 36: RAW-FC dataset sample for label "false"

				Predicte	ed Label			
		true	mostly-true	half-true	barely-true	false	pants-fire	_
true (205)	IBE1 -	57	68	4	17	43	16	
	IBE2 -	48	40	15	60	41	1	
	IBE3 -	47	53	5	59	37	1	
	IBE4 -	4	61	77	8	55	0	
	TBE1 -	43	87	32	14	27	2	
	TBE2 -	53	80	18	36	17	1	
	твез -	186	3	1	1	12	2	
mostly-true (238)	IBE1 -	42	95	11	28	45	17	
	IBE2 -	43	45	28	74	48	0	
	IBE3 -	41	68	9	87	32	1	
	IBE4 -	3	59	95	4	77	0	
	TBE1 -	20	105	59	33	19	2	
	TBE2 -	48	87	24	64	13	2	
	твез -	19	135	49	24	8	3	
half-true (263)	IBE1 -	29	96	12	23	69	34	
	IBE2 -	23	55	25	96	62	2	
	IBE3 -	35	65	9	109	41	4	- 175
	IBE4 -	5	57		12	84	0	- 150
	TBE1 -	23	80	66	43	50	1	
	U TBE2 -	15	63	45	94	39	7	- 125
	С твез -	11	102	74	51	6	19	- 100 9
	BE1 -	13	85	8	22	55	27	L I I I I I I I I I I I I I I I I I I I
	TO IBE2 -	12	31	11	96	60	0	- 75
	IBE3 -	19	33	6	97	53	2	- 50
barely-true (210)	IBE4 -	4	34	82	9	81	0	
	TBE1 -	11	41	39	70	44	5	- 25
	TBE2 -	7	31	14	115	34	9	
	TBE3 -	6	58	45	77	7	17	
false (249)	IBE1 -	14	78	12	27	77	41	
	IBE2 -	13	27	17	103	88	1	
	IBE3 -	21	38	10	104	74	1	
	IBE4 -	1	57	108	12	71	0	
	TBE1 -	20	40	39	61	84	5	
	TBE2 -	10	28	18	91	85	17	
	TBE3 -	44	5	3	1	195	1	
	IBE1 -	6	13	0	9	27	30	
	IBE2 -	1	3	7	33	37	5	
	IBE3 -	9	7	2	31	32	5	
pants-fire (86)	IBE4 -	1	7	15	6	54	3	
(00)	TBE1 -	2	2	11	19	36	16	
	TBE2 -	3	10	2	18	28	24	
	ТВЕЗ -	3	9	12	17	4	41	

Figure 37: Confusion matrix for LIAR-RAW dataset for IBEs and TBEs.



Figure 38: Confusion matrix for RAW-FC dataset for IBEs and TBEs.