

## 594 A THE USE OF LLM

595  
596 We used a large language model (ChatGPT, OpenAI) solely for English copyediting, including gram-  
597 mar correction, wording and minor stylistic re-writes, and occasional LaTeX formatting help. The  
598 model was not used for idea generation, literature search, data collection/annotation, coding, anal-  
599 ysis, or producing results. All scientific claims and contributions were written and verified by the  
600 authors, and no non-public data were shared with the model. The authors assume full responsibility  
601 for the content of the paper.

## 602 B DATA DETAILS

### 603 B.1 DATASET OVERVIEW

604  
605 We study four datasets covering three evaluation tasks. **FHM** and **MultiOff** are annotated only  
606 for **Task 1**, whereas **HarM-c** and **HarM-p** additionally support **Task 2** and **Task 3**. Dataset sizes  
607 (train/test) are listed in Table 3.

Dataset	Train Samples	Test Samples
FHM	8,500	1,000
MultiOff	445	149
HarM-c	3,013	354
HarM-p	2,938	355

608 Table 3: Dataset statistics (train/test) used in our experiments.

609 Per-task label definitions are as follows:

610  
611 **Task 1 (harmfulness, binary):** Decide whether a meme conveys harmful content considering  
612 image–text interaction (e.g., slurs, dehumanizing metaphors, threats or disparagement). Labels:  
613 0=*not harmful*, 1=*harmful*.

614  
615 **Task 2 (harmfulness level, ordinal):** Grade the severity of harmfulness when present; 1 reflects  
616 mild/implicit disparagement or sarcasm, while 2 indicates explicit or strongly implied hostility or  
617 incitement. Labels: 0=*not harmful*, 1=*somewhat harmful*, 2=*very harmful*.

618  
619 **Task 3 (influence target type):** Identify the primary target of the harmful implication of the meme  
620 conditioned on harmfulness. If Task 1 = 0 or Task 2 = 0, set label of Task 3 to 0 (*none*); Labels:  
621 0=*none*, 1=*individual*, 2=*organization*, 3=*community*, 4=*society*.

Dataset	Task 1: harmful or Not		Task 2: Harmful level			Task 3: influence target type				
	0	1	0	1	2	0	1	2	3	4
FHM	510	490	-	-	-	-	-	-	-	-
MultiOff	91	58	-	-	-	-	-	-	-	-
HarM-c	230	124	230	21	103	230	119	0	3	2
HarM-p	184	171	184	159	12	184	81	69	15	6

622 Table 4: Label distribution for each dataset across the three tasks.

### 623 B.2 DATASET-SPECIFIC DESCRIPTIONS AND PROMPT DESIGN

624  
625 **FHM** provides a large-scale set of memes. It has 8,500 training examples and 1,000 test examples.  
626 and its goal is a binary classification of whether a meme is hateful.

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  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

**Definition**  
A meme is considered hateful if it contains a direct or indirect attack on people based on characteristics, including ethnicity, race, nationality, immigration status, religion, caste, sex, gender identity, sexual orientation, and disability or disease. We define attack as violent or dehumanizing (comparing people to non-human things, e.g. animals) speech, statements of inferiority, and calls for exclusion or segregation. The mocking of hate crimes is also considered hate speech.

**Instruction**  
Please analyze the provided image and respond ONLY with a JSON object in the following format:  
{"is\_hateful": <boolean>}  
"is\_hateful": Set to true if the image is harmful according to the definition, otherwise false.

Figure 6: Definition and instructions for FHM.

**MultiOff** contains 445 train and 149 test memes collected from Kaggle (2016 US election) and social media platforms, reflecting real-world political memes and emphasizing both image and text as context.

**Definition**  
A meme is considered offensive if it contains content that is likely to offend, insult, or upset people based on factors such as race, ethnicity, religion, gender, sexual orientation, political views, or other personal characteristics. Offensive content may include offensive speech, discriminatory language, vulgarity, or content that mocks or belittles specific groups or individuals.

**Instruction**  
Please analyze the provided image and respond ONLY with a JSON object in the following format:  
{"is\_offensive": <boolean>}  
"is\_offensive": Must be 'true' if the image is offensive according to the definition, otherwise 'false'.

Figure 7: Definition and instructions for MultiOff.

**HarM** includes two datasets: Harm-C (3,544 memes) and Harm-P (3,552 memes) are real-world meme datasets on COVID-19 and US politics, collected from web and social media, deduplicated and filtered, and annotated by experts for harm intensity and target type (individual, organization, community, society).

**Definition**  
A multimodal unit consisting of an image and embedded text that has the potential to cause harm to an individual, an organization, a community, or society.

**Instruction**  
Please analyze the provided image and respond ONLY with a JSON object in the following format:  
{"harm\_level": <string>, "target\_type": <string>}  
"harm\_level": Must be one of ["not harmful", "somewhat harmful", "very harmful"].  
"target\_type": Must be one of ["individual", "organization", "community", "society"].  
If "harm\_level" is "not harmful", then "target\_type" should be "none".

Figure 8: Definitions and instructions shared by HarM-c and HarM-p.

## C HYPERPARAMETER IN RETRIEVAL STAGE

In the retrieval stage, we combine three modalities: image embeddings, elements, and generated textual descriptions. Let  $\mathbf{v}$ ,  $\mathbf{s}$ , and  $\mathbf{t}$  denote their respective representations. We form a fused vector

$$\mathbf{r} = \lambda_0 \mathbf{v} + \lambda_1 \mathbf{s} + \lambda_2 \mathbf{t}, \quad \text{with } \lambda_0 + \lambda_1 + \lambda_2 = 1,$$

and use  $\mathbf{r}$  to compute retrieval similarity.

Table 5 lists five representative settings. **S1** uses vision only; **S2** and **S3** combine vision with symbolic or textual cues; **S4** assigns approximately equal weights to all modalities; **S5** relies purely on non-visual semantics (element + description).

Empirically, **S5** attains the best accuracy (71.75), while **S4** is a strong second (70.34). We hypothesize that harmfulness is primarily conveyed by semantic information including who is targeted, how,

and in what context, which are better captured by symbolic parses and generated descriptions than by appearance alone. Manual inspection of retrieved memes supports this view: memes that are visually similar can encode opposite meanings due to overlaid text or subtle symbolic cues, which can mislead vision-dominant retrieval (e.g., **S1–S3**). Overall, down-weighting the image channel (small  $\lambda_0$ ) while allocating most weight to  $\lambda_1$  and  $\lambda_2$  tends to yield stronger retrieval for downstream reasoning.

Table 5: Parameter settings for associative retrieval with three modalities.  $\lambda_0$ ,  $\lambda_1$ , and  $\lambda_2$  denote the weights of image, symbolic elements, and generated descriptions, respectively.

Setting ID	$\lambda_0$ (Image)	$\lambda_1$ (Elements)	$\lambda_2$ (Description)	ACC
S1	1.0	0.0	0.0	69.77
S2	0.5	0.5	0.0	68.93
S3	0.5	0.0	0.5	67.80
S4	0.33	0.33	0.34	70.34
S5	0	0.5	0.5	<b>71.75</b>

## D SCALING LANGUAGE MODELS

We conducted a controlled study to examine how scaling LLM affects performance within our framework. Here, the VLM is kept fixed as the perceptual backbone, while the LLM responsible for reasoning is varied in size. To further test the generality of our approach, in addition to Qwen2.5-VL-7B used in the main experiments, we also include InternVL3.5-8B, another VLM that performs well in baseline evaluations. Table 6 reports the results on the HarM dataset across all three tasks.

We observe that introducing even a 7B LLM substantially improves performance over the VLM-only baseline, highlighting the effectiveness of symbolic reasoning. While model size has limited impact on the relatively easier binary harmful classification task, larger LLMs consistently yield gains on more fine-grained tasks. In particular, the 14B and 32B models provide notable improvements on harmfulness level prediction and target type identification, indicating that stronger reasoning capacity better supports nuanced inference.

Method		Harmful or Not			Harmfulness Level	Target type
VLM	LLM	ACC	BACC	MCC	M-F1	M-F1
Qwen2.5VL-7B	–	67.23	53.97	18.23	35.92	3.82
Qwen2.5VL-7B	Qwen2.5-7B	67.80	67.78	34.13	47.68	11.02
Qwen2.5VL-7B	Qwen2.5-14B	<b>70.34</b>	<b>71.78</b>	<b>41.59</b>	<b>51.90</b>	<u>15.94</u>
Qwen2.5VL-7B	Qwen2.5-32B	<u>68.08</u>	<u>70.79</u>	<u>39.73</u>	46.98	<b>16.01</b>
InternVL3.5-8B	–	64.41	69.82	38.99	42.86	12.25
InternVL3.5-8B	Qwen2.5-7B	<b>68.93</b>	69.21	36.83	47.63	12.67
InternVL3.5-8B	Qwen2.5-14B	67.80	<b>72.62</b>	<b>43.93</b>	<b>50.82</b>	<u>16.70</u>
InternVL3.5-8B	Qwen2.5-32B	<u>68.36</u>	<u>71.19</u>	<u>40.52</u>	<u>48.63</u>	<b>17.27</b>

Table 6: Performance of API-based models on the HarM dataset across three tasks.

## E SYMBOLIC REASONING MODULE DESIGN

Section 3.4 describes our symbolic reasoning pipeline. Due to space limitations, the main text illustrates only the Translator. Here we present the visualizations of the remaining two components—the *Planner* and *Solver*—referenced in Section 3.4.

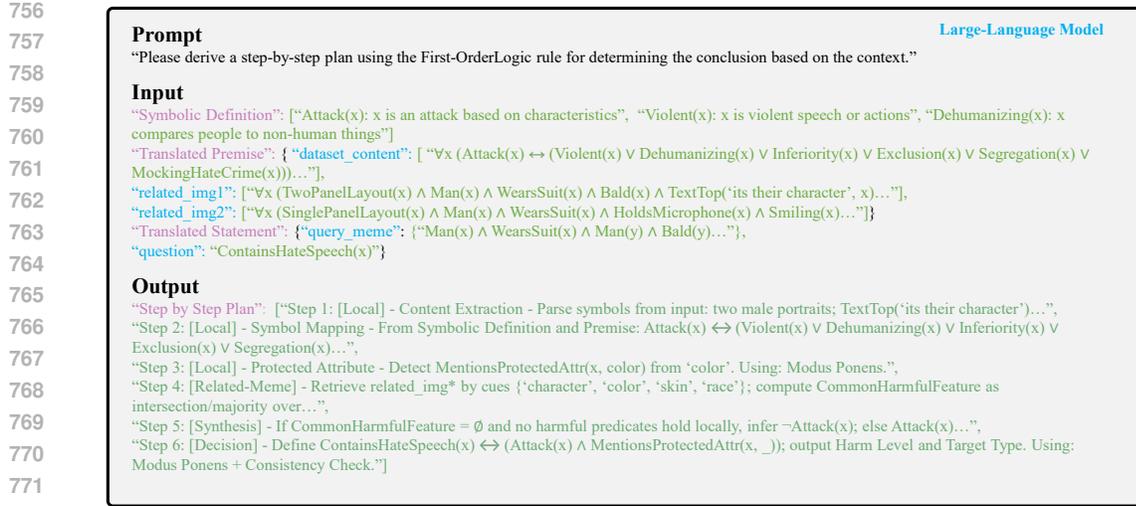


Figure 9: Illustration of the Planner Stage



Figure 10: Process Visualization of covid\_memes\_5560 form HarM

## F PROCESS VISUALIZATION

To provide a more intuitive understanding of how our framework operates, we visualize several representative examples that illustrate the intermediate steps in both the retrieval and reasoning stages. The visualizations show how relevant memories are retrieved across different modalities and how the retrieved information is subsequently integrated into the reasoning process. These examples are intended to complement the quantitative results by offering a step-by-step view of the model's internal workflow.

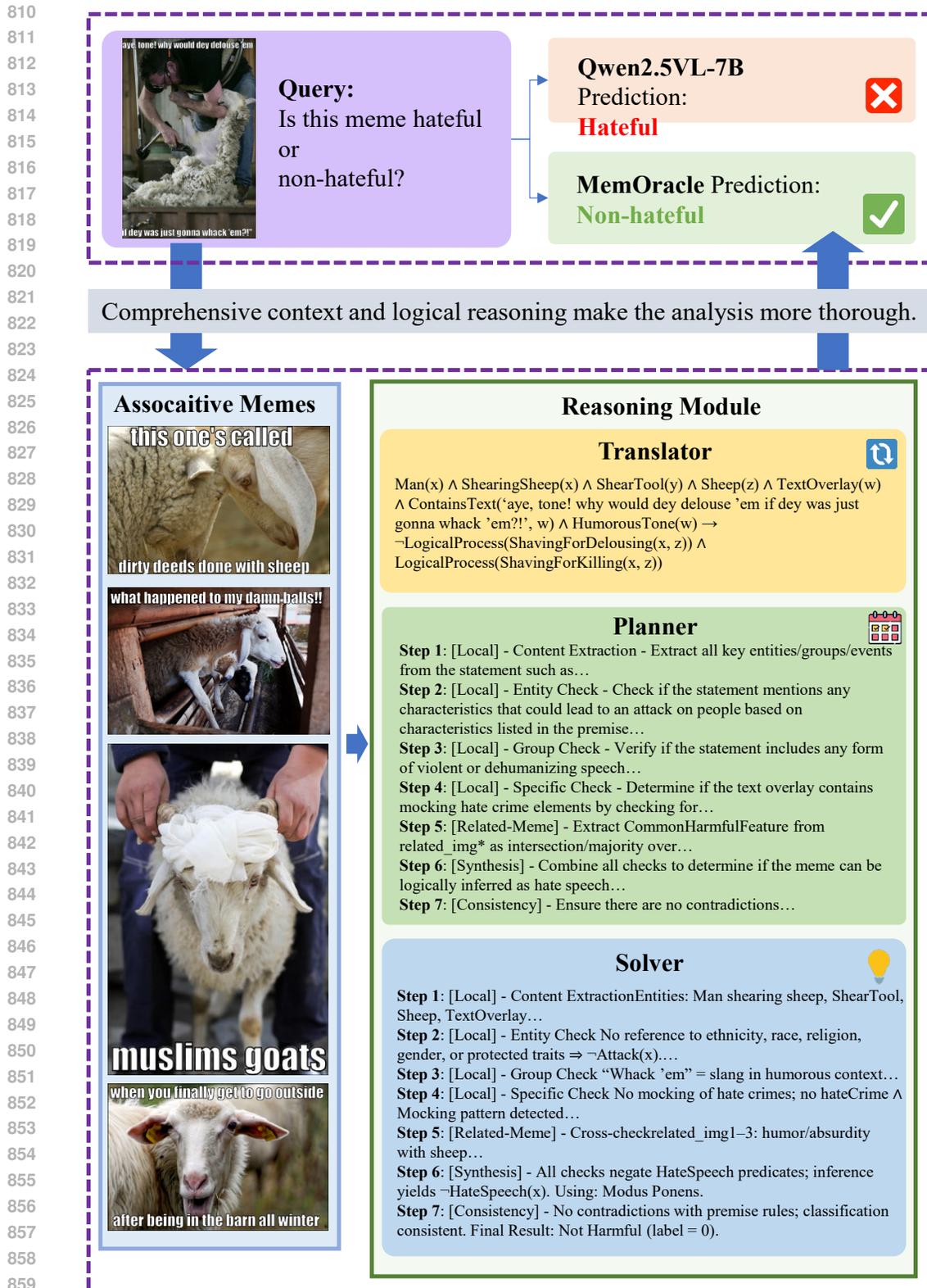


Figure 11: Process Visualization of 03164 form FHM

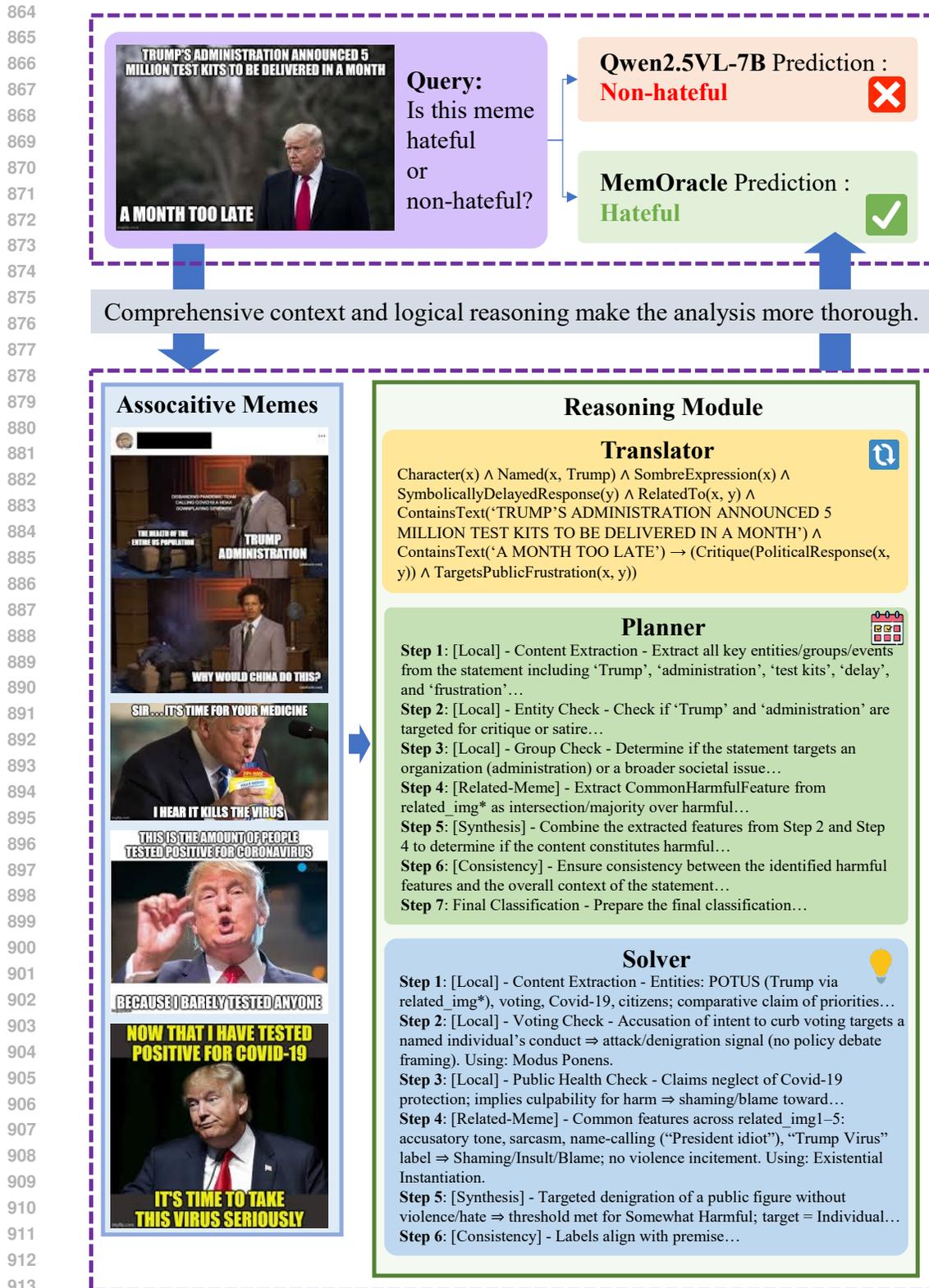


Figure 12: Process Visualization of covid\_memes\_5560 form HarM