

Context-Aware Pragmatic Metacognitive Prompting for Sarcasm Detection

1st Michael Iskandardinata
Computer Science Department
School of Computer Science
Bina Nusantara University
Jakarta, Indonesia
michael.iskandardinata@binus.ac.id

2nd William Christian
Computer Science Department
School of Computer Science
Bina Nusantara University
Jakarta, Indonesia
william.christian001@binus.ac.id

3rd Derwin Suhartono
Computer Science Department
School of Computer Science
Bina Nusantara University
Jakarta, Indonesia
dsuhartono@binus.edu

Abstract—Detecting sarcasm remains a challenging task in the areas of Natural Language Processing (NLP) despite recent advances in neural network approaches. Currently, Pre-trained Language Models (PLMs) and Large Language Models (LLMs) are the preferred approach for sarcasm detection. However, the complexity of sarcastic text, combined with linguistic diversity and cultural variation across communities, has made the task more difficult even for PLMs and LLMs. Beyond that, those models also exhibit unreliable detection of words or tokens that require extra grounding for analysis. Building on a state-of-the-art prompting method in LLMs for sarcasm detection called Pragmatic Metacognitive Prompting (PMP), we introduce a retrieval-aware approach that incorporates retrieved contextual information for each target text. Our pipeline explores two complementary ways to provide context: adding non-parametric knowledge using web-based retrieval when the model lacks necessary background, and eliciting the model’s own internal knowledge for a self-knowledge awareness strategy. We evaluated our approach with three datasets, such as Twitter Indonesia Sarcastic, SemEval-2018 Task 3, and MUSTARD. Non-parametric retrieval resulted in a significant 9.87% macro-F1 improvement on Twitter Indonesia Sarcastic compared to the original PMP method. Self-knowledge retrieval improves macro-F1 by 3.29% on Semeval and by 4.08% on MUSTARD. These findings highlight the importance of context in enhancing LLMs performance in sarcasm detection task, particularly the involvement of culturally specific slang, references, or unknown terms to the LLMs. Future work will focus on optimizing the retrieval of relevant contextual information and examining how retrieval quality affects performance. The experiment code is available at: https://github.com/wlchrst/sarcasm-detection_pmp_knowledge-base.

Index Terms—Sarcasm Detection, Large Language Models (LLMs), Pragmatic Metacognitive Prompting (PMP), Web-based Retrieval, Self-Knowledge Awareness

I. INTRODUCTION

In the field of machine learning, natural language processing (NLP) tasks have been shown to be a crucial part of human life. NLP tasks revolve around processing text in a specific manner and receiving an output that can be useful, with examples such as text classification, text generation, information retrieval, and similar related tasks. In this study, we direct our focus toward text classification, specifically examining one of its key sub-tasks: sarcasm detection. Sarcasm detection, or as some call it verbal irony detection, is a task in NLP that automatically classifies text, and in extended forms includes

images, audio, or video, as either sarcastic or not. Systems designed for sarcasm detection are becoming more important in the 21st century due to the growth of the use of sarcasm detection datasets caused by media usage of it, such as social media, television, and much more. Additionally, enhancing these automatic systems for sarcasm detection could become crucial in interpreting the real sentence meaning of a text. For example, a song review saying “The album is really incredible, it even made me sleep while listening to it” might be classified as something positive by a sentiment classification model of the word “incredible.” However, upon further reading, the reviewer says that it made them sleep, which suggests that the album is boring.

Recent advances in machine learning studies and architectures have led to many exceptional performances across NLP tasks, such as sentiment and emotion classification. However, despite these improvements that have also impacted the overall performance of sarcasm detection, it still remains a challenging problem as it often relies on subtle contextual, pragmatic, or cultural cues. Large Language Models (LLMs) have shown tremendous capability in emulating human reasoning and completing specific tasks that have predefined rules [1]. Multiple studies recently have shifted their focus toward methods for improving LLM performance, such as external information [2], prompt techniques [3], information retrieval technique [4], and related methods.

The latest studies on sarcasm detection have benefited LLMs in their new method called Pragmatic Metacognitive Prompting (PMP) [5], which focuses on guiding LLMs to analyze and extract pragmatic information from a text. However, accurately interpreting a text requires a nuanced understanding of word meanings and their contextual relationships within the text. Unfortunately, due to the vast diversity of language and real-world context, this remains a major challenge for LLMs [6], as illustrated by the example in Figure 1 for sarcasm detection. Hence, recent studies on LLMs have focused on incorporating non-parametric knowledge for tasks that are knowledge intensive, e.g., question answering and fact verification. Retrieval-Augmented Generation (RAG) has emerged to be the most popular topic in solving the problem, retrieving external information with the objective of improving accuracy

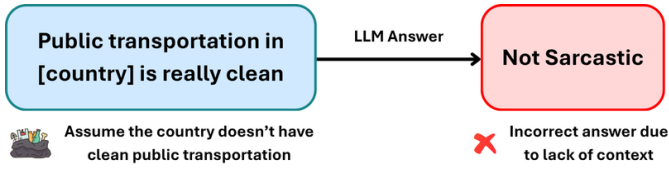


Fig. 1. Extreme diversity and culture example

and reducing hallucinations. Existing RAG studies are fairly scattered, with some focusing on the retrieval process [7] [8], while others examine the behavior of RAG [9].

Beyond the non-parametric knowledge, though accommodated with sufficient internal knowledge, models sometimes fail to flag which words or tokens need extra grounding. Such failures arise because internal knowledge is uneven and not reliably activated [10]. This could lead to a weaker analysis of the overall text for the final verdict of sarcasm detection. Recent studies have shown that eliciting the internal knowledge of LLMs can mitigate this issue by increasing its contextual awareness and reasoning quality. Such self-knowledge prompting not only helps a model determine when additional information may be needed, but also enhances its ability to reason more precisely even without relying on any external retrieval process [11] [12].

Building on this line of work, our study investigates the impact of non-parametric and internal knowledge awareness as context to primarily solve the limitation of using LLMs in a zero-shot manner when doing sarcasm detection. The non-parametric knowledge is retrieved based on semantically important words, which are identified using an LLM-based extraction. These words are then retrieved using Google Search API with the results cleaned by the LLM. Meanwhile, the internal knowledge is retrieved from the LLM itself with the important words being identified using Named Entity Recognition (NER) and Part-of-Speech (POS) tagging. All of these methods, from keyword extraction to context retrieval, are detailed in later parts of the paper. In addition, our study integrates the PMP method, which is then evaluated in English and Indonesian datasets, the latter consisting of traditional and culture-specific varieties of Indonesian. This combined method shows that when appropriate context is provided to LLMs, it plays an important role in helping them to gain deeper understanding of the text, therefore improving their ability to interpret and classify sarcastic text.

Contributions from this research are listed below:

- Discovering the important role of non-parametric and internal knowledge awareness in sarcasm detection using LLMs, specifically through the state-of-the-art PMP method.
- Investigating the best course of action to identify words that are not fully understood by LLMs before performing information retrieval.
- Implementing simple information retrieval and producing a high margin of performance difference, revealing that a better information retrieval method can be highly

beneficial.

The remainder of this paper is structured as follows. Section II reviews related works on text classification, sarcasm detection, NER, POS tagging, and RAG systems. Section III explains in detail the methods used for key word extraction, information retrieval, and prompting techniques. Section IV explains which datasets, models, and evaluation metrics are used for our experiment evaluation. Section V presents and analyzes each result of the experiment. Finally, Section VI concludes the paper and highlights its advantages, limitations, and directions for future work.

II. RELATED WORKS

In this section, we will focus on the development of sarcasm detection, in which the explanation of earlier methods of text classification leading to sarcasm detection using Large Language Models (LLMs).

A. Text Classification

Text classification, the act of automatically classifying a text into a predefined set of labels, has been one of the most extensively researched areas of the earlier development of machine learning until now. Especially with today's growing influence of the internet and social media, it has been impacting the way machine learning is used for text classification tasks such as sentiment classification [13], emotion classification [14], spam detection [15], sarcasm detection [16], and much more. In addition to conventional tasks of text classification, recent studies have shifted focus toward emerging domains such as multilingual text classification [17], multi-modal emotion classification [18], and other similar classifications.

Early text classification tasks relied primarily on classical machine learning algorithms such as clustering [19], Naive Bayes [20], and Support Vector Machines [21], to name a few. These algorithms depended heavily on feature extraction, a process defined as converting text into numerical data, usually in the format of vectors, in order to be used as input for the algorithm. Some of the most frequently used methods for this process were TF-IDF [22], n-gram [23], and bag-of-words [24]. While effective, these features often failed to take into account the semantic meaning and contextual relationships between words, focusing only on the word count of each sentence.

Nonetheless, these methods became an important building block for studies in text classification when researchers discovered neural networks. Neural networks enhance the feature extraction process by placing greater emphasis on the semantics and context of each word in a sentence. Some of their most used frameworks include Convolutional Neural Network (CNN) [25], Recurrent Neural Network (RNN) [26], and Recurrent Convolutional Neural Network (RCNN) [27]. These developments shifted the field toward deep models for text classification. As a result, one of the most significant innovations in neural networks of this decade was the Transformer framework [28], which laid the foundation for most modern

language models, including BERT [29] and LLMs, such as GPT [30] and Gemini [31], to name a few.

The rapid advancement of text classification has effectively addressed what can be described as "System I" tasks, those that are fast, unconscious, and intuitive. These tasks have been argued to be successfully solved, including sentiment analysis, emotion classification, and spam detection. In contrast, "System II" tasks require slower, deliberate, and multi-step cognitive processes, such as sarcasm detection, which remains an unsolved challenge [32].

B. Sarcasm Detection

Sarcasm detection remains one of the most challenging tasks in text classification, despite often being formulated as a simple binary problem (sarcastic or not sarcastic). Recent sarcasm detection studies have also expanded toward multi-label [33] and multi-modal [34] settings. These approaches aim to capture the deeper pragmatic and contextual cues of sarcasm, which often rely not only on text but also on tone of voice, facial expressions, and situational context.

Recent studies of sarcasm detection have focused on dataset creation and augmentation. Some works even employ deep learning frameworks, such as the Reverse Generative Adversarial Network, to handle the dataset augmentation [35]. Other studies emphasize dataset construction, for example, collecting sarcastic lines from TV shows [34] or by using web scraping methods [36].

1) *Pre-trained Language Models*: Consistent with the development of neural networks, recent studies on sarcasm detection have also adopted models based on transformer, such as BERT [37] and RoBERTa [38]. Building on the strong contextual representation capabilities of these models, researchers have started to explore more advanced Pre-trained Language Models (PLMs) that can be fine-tuned or adapted specifically for sarcasm detection using text. For instance, Dual-Channel Network (DC-Net) [39] was designed solely for sarcasm detection by modeling distinct sentiment cues from text segments before combining them for classification. Similarly, SarcPrompt [40] introduces prompt-based fine-tuning to guide PLMs in identifying sarcasm polarity more accurately.

2) *Large Language Models*: Although PLMs have proven effective for sarcasm detection, recent advances in Large Language Models (LLMs) have also shown remarkable progress in this area, along with an increasing number of studies investigating their effectiveness and applications. Their ability to capture nuanced contextual and pragmatic cues has allowed researchers to approach sarcasm detection from a more reasoning-oriented perspective. Initially focused on text generation, LLMs now excel in multiple tasks such as reasoning [41], code generation [42], and question answering [43].

Prompting Methods A key feature that enables LLMs to perform diverse tasks is prompting, the process of guiding the model's reasoning or behavior through carefully constructed instructions. In the context of text classification, prompting serves not only to map input text to predefined labels but also to influence how the model interprets context, structure,

and intent. Common prompting strategies for text classification include the following:

- Zero-shot prompting: provides LLMs with a natural-language prompt describing the classification goal to directly predict the appropriate label [44].
- Few-shot prompting: introduces LLMs with labeled examples within the prompt to demonstrate the desired mapping before classifying new inputs [45].
- Chain-of-Thought (CoT): uses a series of intermediate reasoning steps to enhance the reasoning capabilities of LLMs in tasks such as common sense inference and arithmetic problems [3].

Pragmatic Metacognitive Prompting Beyond standard prompting methods such as zero-shot, few-shot, and chain of thought prompting, recent studies in sarcasm detection have focused on finding more effective prompting strategies. So far, Pragmatic Metacognitive Prompting (PMP) [5] has shown promising progress with higher performance in every metric compared to PLMs. PMP itself is built on top of Metacognitive Prompting (MP) [46], a framework that introduces the concept of preliminary analysis followed by a reflection stage before giving the final judgment. PMP adopts this same infrastructure, but extends it by embedding simplified pragmatic components into both stages. The first stage focuses on extracting and interpreting pragmatic elements from the text, while the second stage reflects on the earlier analysis to form a deeper reasoning about the meaning and intent behind the statement. This process allows the model not only to evaluate what is said, but also to understand what is implied, which improves its ability to detect sarcasm and indirect meaning in the text.

TABLE I
SARCASM DETECTION PERFORMANCE FROM PREVIOUS STUDIES

Method	Model	Accuracy	Macro-F1
<i>MUSARD Dataset</i>			
PMP	GPT-4o	0.7942	0.7765
PMP	LLaMA-3-8B	0.5358	0.5469
PLMs	RoBERTa	0.8680	0.8770
PLMs	DistilBERT	0.8700	0.8770
<i>SemEval 2018 Dataset</i>			
PMP	GPT-4o	0.8668	0.8318
PMP	LLaMA-3-8B	0.7821	0.7765
PLMs	RoBERTa	0.7500	0.7200
PLMs	DC-Net-RoBERTa	0.7090	0.6870
<i>Twitter Indonesian Dataset</i>			
Zero-shot	mT0XL	0.2494	0.3989
PLMs	XLM-RLarge	0.8885	0.7692

Note: Macro-F1 denotes the **macro-averaged F1-score** across all classes.

However, PMP does not have leverage on all datasets, showing state-of-the-art performance in SemEval 2018 [47] but inferior results when evaluated in MUSARD [34]. This variation stems from a broader limitation of LLMs, namely their reliance on parametric knowledge. Because PMP operates in a zero-shot setting without access to external context, its effectiveness decreases when sarcasm deviates from general

linguistic or cultural norms. The comparison of PMP and PLMs performance can be seen in Table I.

C. Named Entity Recognition

Named Entity Recognition (NER) is a fundamental task in Natural Language Processing (NLP) that aims to identify entities of interest in text, such as people, locations, and organizations. This task plays an important role in guiding the semantic understanding of text. [48]. Similar to other NLP tasks, early approaches to NER relied heavily on handcrafted features and statistical models. With the advancement of deep neural networks across various machine learning domains, as discussed in Subsection II-A, these architectures have also become the foundation for modern NER systems. A recent state-of-the-art framework, BiLSTM-CRF [49], achieved strong performance and outperformed traditional NER methods. However, these architectures often under-perform in transferability compared to pre-trained transformer-based models such as BERT [50], which is particularly relevant to this study due to its ability to generalize across domains and languages. At the same time, despite the success of transformer-based architectures, smaller neural network models remain valuable due to their computational efficiency. A popular NLP library in Python, spaCy, adopts a lightweight CNN-based architecture for its NER and POS tagging components, offering an efficient alternative that integrates well into larger NLP systems [51].

D. POS Tagging

Part-of-Speech (POS) tagging assigns a syntactic category to each word in a text and is a fundamental task in NLP. Similar to NER, POS tagging helps models understand the structure and function of words in a sentence. Early approaches relied on classical machine learning algorithms such as Naive Bayes, Hidden Markov Models (HMM), and Conditional Random Field (CRF) [52]. With the rise of deep learning, sequence models, including Gated Recurrent Unit (GRU), Recurrent Neural Network (RNN), and Bi-directional Long Short Term Memory (LSTM) [53] [54], have become dominant in the field. More recently, studies have shifted focused on identifying the effectiveness of POS tagging using LLMs [55] and transformer-based models like BERT [56].

E. Retrieval-Augmented Generation

One major issue with LLMs is hallucination [57]. Previous works have shown that the parametric knowledge of LLMs can be unreliable in open-domain question answering, and that performance varies depending on the popularity of the entity being queried; in particular, less common entities tend to increase the likelihood of hallucination [6]. Retrieval-Augmented Generation (RAG) addresses this by augmenting model prompts with non-parametric knowledge retrieved from external sources, which improve factuality and reduce hallucination. In practice, retrieval can be implemented with sparse methods (e.g., BM25), dense retrievers like Dense Passage Retrieval (DPR), or hybrid pipelines that combine both. A recent work known as Interleaving Retrieval with

Chain of Thought prompting (IRCoT) [4] has also explored non-parametric knowledge retrieval but combines it with CoT prompting technique, which has shown a tremendous result in knowledge augmentation [58].

F. Context-Awareness and Internal Knowledge

Researchers have found that LLMs do not always recognize which parts of an input need extra background or deeper reasoning. This can result in weak performance in tasks that require subtle interpretation [59]. One line of work has extracted signals from the model itself, such as uncertainty or “do I know this?” responses, to decide when to look up external information. For example, a study proposed a method in which the model checks its internal knowledge and only prompts retrieval when needed [12]. Another study asked LLMs to generate short internal summaries about entities or concepts present in the input, and showed that this helped even without retrieving large external corpora [11]. Additional work has used the internal hidden states of transformer layers to estimate the model’s self-awareness of its factual bounds, and linked this measurement to performance and retrieval decisions [60]. Overall, this literature highlights that improving a model’s awareness of its own knowledge state is a useful complement to external retrieval and can itself enhance reasoning quality.

III. METHODOLOGY

To solve the problem of specific linguistic norms in texts for sarcasm detection, we propose a method that is inspired by Retrieval-Augmented Generation (RAG), in which our method focuses on identifying, retrieving, and providing information that the model may not understand. In this section, we explain the retrieval step and the prompting method used in this study.

A. Retrieval Method

Our information retrieval process begins with the identification of relevant words, which are later used as keywords during the retrieval process.

1) *Keyword Extraction*: Keyword quality highly matters for the performance of the LLM. Too many keywords increase query size, add noise, and may push the input past token or cost limits. Therefore, the extraction of keywords from the text should only include words that most improve the model’s understanding. Our method proposes two extraction approaches, namely Token Tagging and LLM-Based.

Token Tagging Let us first define our approach in Token Tagging, in which our approach uses Named Entity Recognition (NER) and Part-of-Speech (POS) tagging as illustrated in Figure 2. Given a task of sarcasm detection of the text T , T is first processed by using a pipeline that tags each word, or token, with their respective entity type and part of speech. We focus only on tokens identified as entities or proper nouns. Entities refer to real-world objects with distinct types such as people, locations, and organizations, for example, “Elon Musk”, “Paris”, or “Google”. Proper nouns are included as well due to NER models occasionally failing to recognize certain entities; proper nouns often capture those missed cases.

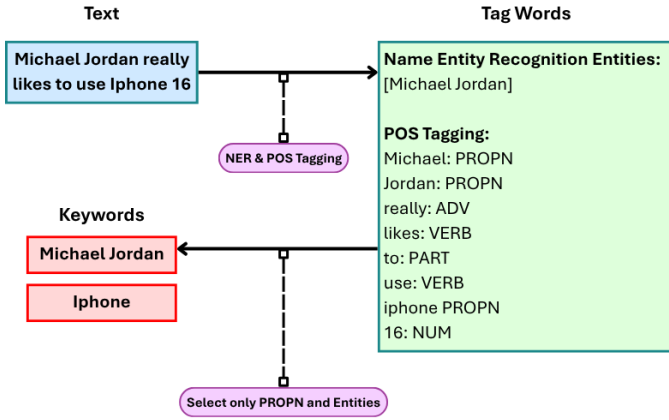


Fig. 2. Token Tagging Extraction

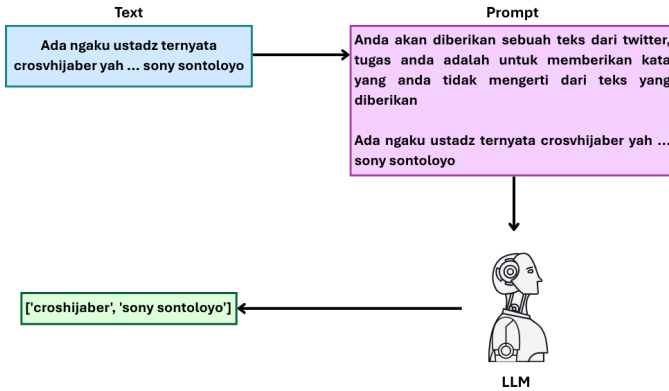


Fig. 3. LLM Based Extraction

Our extraction is restricted to these tokens because they carry strong contextual meaning, especially in domain-specific cases, helping the model interpret the text more accurately.

LLM-Based The reasoning behind using an LLM to extract keywords is inspired by the IRCot method [4], which employs LLMs to guide the retrieval process. However, IRCot focuses primarily on question-answering tasks, specifically multi-hop retrieval, where the LLM is used to iteratively refine queries to support reasoning. In contrast, our method leverages the LLM in a different capacity. Rather than generating refined queries, we utilize the LLM to automatically extract words from raw text. This method was initially used for the Indonesian Twitter dataset, which contains a lot of slang words in various languages, including Indonesian and Javanese, which makes this particular approach preferable to Token Tagging. In our approach, the input text T is provided to the LLM along with an additional prompt that instructs the model to identify and return a list of words it does not recognize or cannot infer from prior knowledge. These unknown or ambiguous terms are hypothesized to carry domain specificity or contextual importance, which makes them strong candidates for keyword information retrieval. The details of the prompt used to guide the LLM in this extraction process are illustrated in Figure 3.

2) *Word-Information Retrieval*: After selecting keywords that require additional context, we send them to a word-information retrieval pipeline. Our approach implements two complementary retrieval strategies, first using Google Search API and LLM-based cleaning, and second using an LLM-only approach.

Google Search API + LLM Each extracted keyword is used as a query to retrieve candidate links and snippets from the web. We utilize the Google Search API for this step. This allows the system to gather real-world usage examples and contextual knowledge that may not be present in the knowledge of the LLM.

However, raw search snippets often contain noise, redundancy, or irrelevant information. Therefore, after retrieving the results, we split the retrieved documents into manageable text chunks and score them with BM25 to select the most lexically relevant passages. BM25 requires low computational cost and is effective for definitional matches, which are useful when the meaning of a keyword is expressed explicitly on web pages. The link retrieval and document ranking steps are illustrated in Figure 4. The BM25 ranking function is defined as:

$$\text{BM25}(q, D) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgDL}})} \quad (1)$$

where the inverse document frequency $\text{IDF}(q_i)$ is defined as:

$$\text{IDF}(q_i) = \ln \left(\frac{N - n(q_i) + 0.5}{n(q_i) + 0.5} + 1 \right) \quad (2)$$

The variables in Equation 1 and 2 are:

- q is the query containing terms q_1, q_2, \dots, q_n .
- D is a document.
- $f(q_i, D)$ is the term frequency of term q_i in document D .
- $|D|$ is the document length in words.
- avgDL is the average document length across the collection.
- N is the total number of documents in the collection.
- $n(q_i)$ is the number of documents containing term q_i .
- k_1 and b are hyperparameters, where typically $k_1 = 1.2$ and $b = 0.75$.

After ranking, we pass the highest ranked chunks to an LLM for refinement. The LLM is prompted to summarize the chunks, filter out irrelevant content, and extract key semantic meanings. For every chunk, the final output is a short definition of itself. This hybrid method is especially useful when the LLM itself lacks domain-specific parametric knowledge or other terms not covered well by the model.

LLM-only In the LLM-only approach, the extracted keywords are provided as input to the LLM to generate a contextual definition for each keyword. The definition is based solely on its internal parametric knowledge, not relying on external knowledge sources. The process begins by providing the keyword as input to the LLM, which then generates a contextual definition based solely on its internal parametric

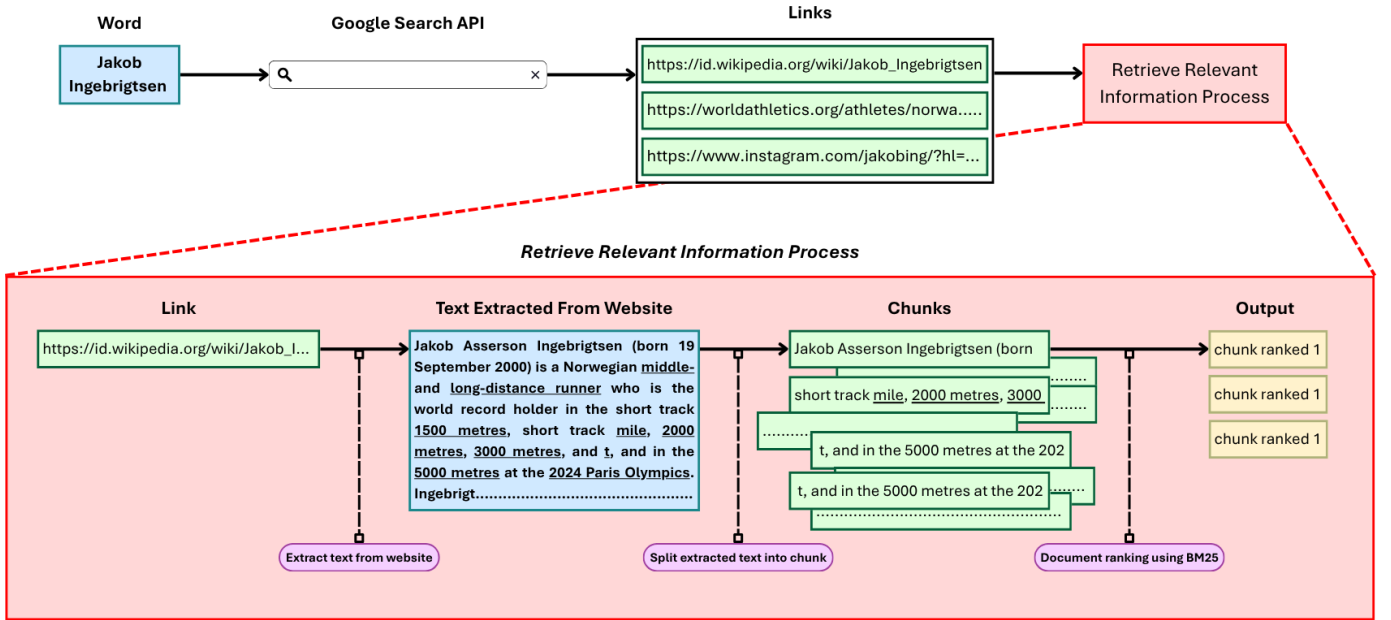


Fig. 4. Retrieve Link Process

knowledge. Although the accuracy of the generated definitions can vary depending on the keyword and the prior training of the model, this approach offers significant advantages in terms of efficiency, as it reduces both processing time and computational cost.

3) *Extraction-Retrieval Pairing*: Between those two methods of keyword extraction and word-information retrieval methods, we provide a short justification for which pairings we evaluate and why those choices best match the dataset properties and practical constraints.

Token Tagging + LLM-only For the first pairing, we use Token Tagging for keyword extraction and LLM-only word-information retrieval. This pairing is our method for adding internal knowledge context, specified for datasets that do not rely on heavy culture-specific slang or local references, for example, MUSTARD, which contains mainly standard English captions. We chose this pairing because when keywords come from a separate, structured extractor, asking the model to define them uses its internal knowledge without creating a circular process. If the model both selects the tokens and then tries to define them, in this case pairing up LLM keyword extraction and retrieval, the workflow can simply reinforce its own mistakes. Some prior work has used the self-asking pattern for internal knowledge to boost context awareness, but we expect cleaner and more reliable context when the model is grounded by an external extractor like Token Tagging.

LLM-based + Google Search API The second pairing uses LLM-based keyword extraction and Google Search API word-information retrieval. This pairing is our method for adding non-parametric knowledge context, particularly designed for the Indonesian Twitter dataset, which contains many domain-specific slang and terms that spaCy models cannot reliably detect as named entities or proper nouns, so Token

Tagging would miss many candidate tokens. Instead, we ask the LLM directly which tokens it does not recognize. Those flagged tokens are then used to query the Google Search API for non-parametric knowledge, avoiding the LLM-only retrieval approach, which clearly would not have any information about the tokens. Note that we do not apply Google Search API retrieval to datasets made up mainly of standard English captions, for example, MUSTARD. In those cases, the LLM already contains the relevant background knowledge, and adding web retrieval tends to introduce irrelevant snippets and noise rather than improve understanding.

B. Prompting Method

In our prompting methods, we involve three primitive prompting components, such as Pragmatic Metacognitive Prompting (PMP), few-shot, and word-information prompting. From there, two hybrid pipeline templates are produced from combinations between them, with a final total number of five prompting variants, such as a PMP baseline and four hybrids.

1) *Primitives Prompting Components*: This part describes the three primitive prompting components we use as building blocks, namely PMP, few-shot, and word-information prompting.

Pragmatic Metacognitive Prompting Our study uses the PMP pipeline for the baseline and skeleton, with two LLM calls per input X . The first call, denoted P_1 , asks the model to understand X generally and produce an initial pragmatic analysis, denoted A_1 . The second call, denoted P_2 , reflects on the analysis of the first call to produce a more structured pragmatic analysis, containing factors from a pragmatic viewpoint such as the implicature, presupposition, intent, polarity, pretense, and potential meanings individually, denoted A_2 . From that analysis, the second call also verdicts

at the end whether the context is sarcastic or not with a final label, y , of "YES" or "NO".

Few-Shot Prompting In this approach, a small set of illustrative examples is provided to guide the reasoning process of the LLM. Rather than serving solely as reference points for producing correct labels or mappings, these examples function as cognitive cues that help the model enhance its internal thought process. These short examples, denoted F , are made manually and tailored to the dataset theme, with each typically containing an input, the structured analysis, and the final label. The number of examples, parameter k , is intentionally small to maintain clarity and reduce noise. We used $k = 2$ in our experiments.

Word Information Prompting The word-information retrieval returns contextual definitions, using either Google Search API with LLM-based keyword extraction, denoted W_g , or LLM-only approach with Token Tagging keyword extraction, denoted W_l . These retrieved definitions are included in the prompt alongside the input to help the model better understand for further analysis. When the model sees concise clarifications for slang, named-entities, or culture-specific terms it might not otherwise know, further analysis becomes less noisy and more accurate. Each definition is kept to one or two sentences so the prompt stays compact and readable.

2) *Notations and Pipeline Overview:* Below we define formal notations used in the prompting pipelines and present the PMP skeleton as building blocks for the two hybrid pipeline templates.

Notations:

- X input text.
- P_1 first call of PMP.
- P_2 second call of PMP.
- A_1 preliminary pragmatic analysis (output of P_1).
- A_2 structured reflection (output of P_2).
- K_t keywords extracted from X using Token Tagging.
- K_l keywords extracted from X using LLM-based.
- W_l word-info produced by LLM-only approach (RetrieveLLM(K_t)).
- W_g word-info produced by Google Search API \rightarrow LLM (RetrieveGoogle(K_l)).
- W chosen word-info block $\in \{\emptyset, W_l, W_g\}$.
- F few-shot demonstrations (k examples).
- y final label $\in \{\text{YES}, \text{NO}\}$.
- \oplus concatenation into prompt (append to user or system prompt).

PMP skeleton:

$$A_1 = P_1(X)$$

$$(A_2, y) = P_2(A_1)$$

First hybrid template (PMP + Word-Info):

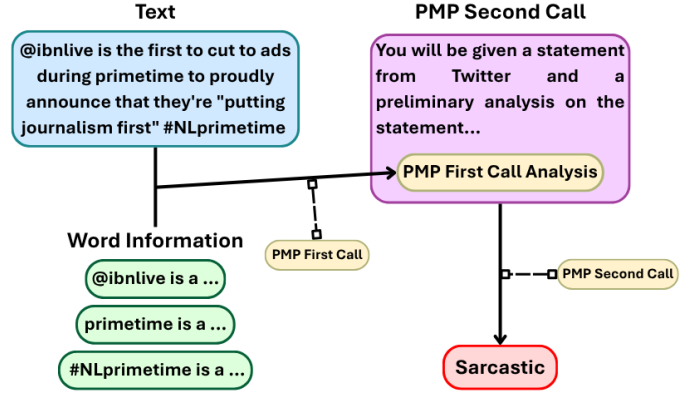


Fig. 5. First Hybrid Prompting Template

$$A_1 = P_1(X \oplus W)$$

$$(A_2, y) = P_2(A_1)$$

Second hybrid template (PMP + Word-Info + Few-shot):

$$A_1 = P_1(X \oplus W)$$

$$(A_2, y) = P_2(A_1 \oplus F)$$

3) *Hybrid Prompting Pipelines:* This part presents the two hybrid pipeline templates, each realized with Google Search API or LLM word-information. Combining these templates with the PMP baseline yields five prompting variants evaluated.

First Hybrid Template (PMP + Word-Info) In this first hybrid variant, the pipeline integrates the word-information block (W) into the PMP process, which can be seen in Figure 5. With PMP having two calls, putting the block in the first call (P_1) would be more suitable, as it is the model's first time encountering and understanding the input. The goal here is to give the model more grounding on specific terms that may not exist in its prior knowledge, such as slang, uncommon entities, or local cultural phrases. Introducing the block in the second call (P_2) would not just already be too late, but could also confuse the reflection process and add unnecessary noise. By providing the word information early, the model can process every keyword and meaning together with the input from the start, allowing the P_2 , the reflection stage, to focus purely on deeper pragmatic reasoning instead of fixing misunderstandings about the words themselves.

Regarding whether it should be placed in the user or system prompt, we put it in the user prompt since it is also part of the input that the model needs to process. System prompts are filled with instructions, shaping the model's behavioral and thinking guide, while user prompts are read by the model as "what to process". Therefore, W fit better in the user prompt to be processed together with the main text, helping the model understand the meaning contextually before continuing to analyze it pragmatically.

Second Hybrid Template (PMP + Word-Info + Few-Shot) For another hybrid variant, we add few-shot examples (F) to the first hybrid variant as illustrated in Figure 6. F are placed on the system prompt, not on the user prompt. As explained before, the system prompt mainly controls how the

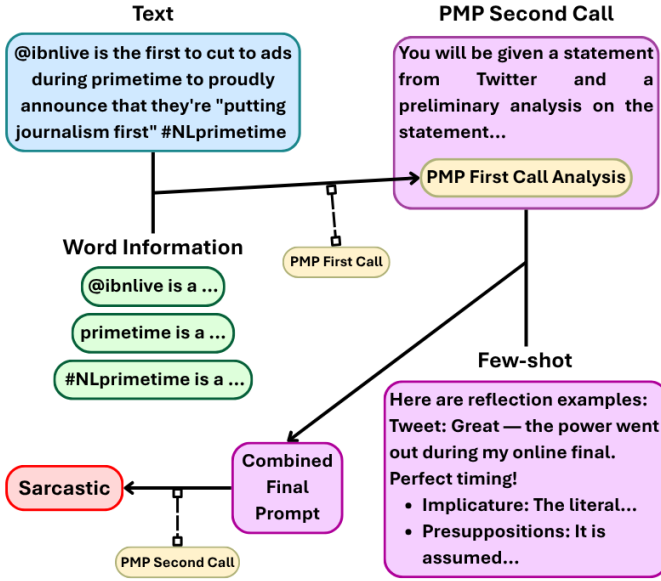


Fig. 6. Second Hybrid Prompting Method

TABLE II
PIPELINE EXPERIMENT SETTINGS

Setting	Word-Info (LLM-only W_l)	Word-Info (Google Search API W_g)	Few-shot (F)
PMP			
PMPWL	V		
PMPWG		V	
PMPWL-FS	V		V
PMPWG-FS		V	V

model behaves and thinks. They bias internal heuristics and the model’s way of combining factors, not merely the output tokens, a big advantage in improving the PMP analysis.

Then, with PMP having two calls, it is unnecessary to put F integration at the first call (A_1) due to its task to only do a quick and general pragmatic scan. Throwing F here only adds noise that could impact the final verdict later on. However, the second call (A_2) is a structured step, a complex analysis to closely inspect each pragmatic factor and then combine those pieces into a final verdict. This step demands careful, factor-by-factor reasoning, not because it is more complex in an abstract sense, but because it requires precise mapping from the analysis to a label. F could prove to be useful here.

Final Five Prompting Variants The two hybrid templates each produce two variants, one from Google Search API word-information retrieval with LLM-based keyword extraction, and one from LLM-only word-information retrieval with Token Tagging keyword extraction. From this, we have four hybrid variants, with a total number of five prompting variants when combined with the PMP baseline as categorized in Table II. Here are their formulas:

- PMP: $y = P_2(P_1(X))$

- PMPWL: $y = P_2(P_1(X \oplus W_l))$
- PMPWG: $y = P_2(P_1(X \oplus W_g))$
- PMPWL-FS: $y = P_2(P_1(X \oplus W_l) \oplus F)$
- PMPWG-FS: $y = P_2(P_1(X \oplus W_g) \oplus F)$

IV. EXPERIMENTS

In this section, we explain the datasets, models, evaluation metrics, and implementation details that we used. All experiments and evaluations executed in this study, including keyword extraction, word-information retrieval, and pipeline evaluations, were conducted on a machine equipped with an Nvidia RTX 3090 GPU and 24 GB of VRAM. Due to hardware constraints, we were only able to run experiments on small-scale models. This prevented the evaluation of larger models and consequently limited our ability to investigate the impact of context in higher-capacity architectures.

A. Datasets

To evaluate the impact of the word-information context accurately, we conduct evaluations of our pipelines using three different datasets, with three different themes. The datasets are **1) SemEval-2018 Task 3**, an English twitter dataset where sarcastic or ironic tweets are collected by the hashtags, **2) MUSTARD**, a multi-modal dataset created from three different TV shows: Friends, The Golden Girls, and Sarcasmaholics Anonymous, **3) Indonesian Twitter**, a labeled twitter dataset curated from March 2013 to February 2020, which consists of multiple languages including Indonesian, Javanese, and other traditional languages. The distribution of the datasets is illustrated in Figure 7. Unfortunately, the class distribution of the SemEval and Indonesian Twitter datasets is imbalanced, resulting in some evaluation metrics being less reliable. The reason for choosing these three datasets is to compare our results with the original PMP study, which uses SemEval-2018 Task 3 and MUSTARD. In addition, we chose the Indonesian Twitter dataset to further explore the impact of non-parametric knowledge context when detecting sarcasm in text, which is a more extreme linguistic challenge for LLMs. The datasets are processed in the same way as in the PMP study, where Semeval-2018 and the Indonesian Twitter dataset, although the latter was not included in the study, are already in the desired form, having the text (tweet) and the sarcastic label (binary) for each line. Meanwhile, for MUSTARD, while being a multi-modal dataset, it also offers a JSON format version which includes the text (dialogue), the sarcastic label (binary), and previous dialogues. We follow the same pre-processing method used in the PMP study, which was to combine all previous and target dialogue into one line, with the target being in curly brackets. Also, we used the test split version for the Indonesian Twitter dataset

B. Models

NER and POS tagging are done using a spaCy model, specifically **en_core_web_sm** due to the efficiency of the model in utilizing CPU computation capabilities. We decided to use this small model because the performance is already

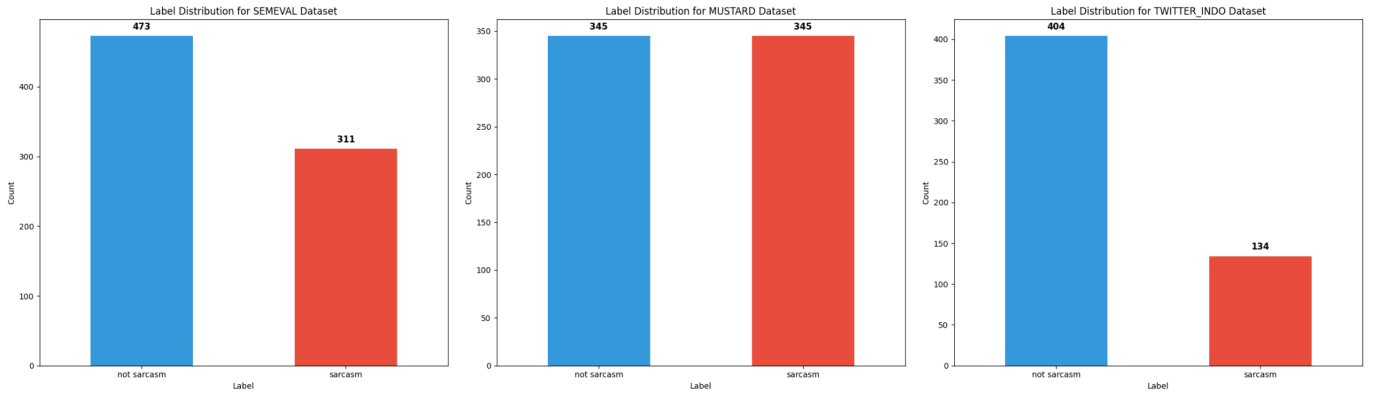


Fig. 7. Dataset Distribution (Test)

high based on the documentation on their website. Then for the LLMs, as mentioned previously, we are limited to small models, therefore open-source models **Qwen-3**, **Llama-3**, and additionally **Llama-3.1**, each with 8 billion parameters, are used for keyword extraction, information retrieval, and pipeline evaluation. These two models, among other 8-billion-parameter models, exhibit top-tier reasoning and knowledge performance, suitable for sarcasm detection. Qwen-3 is also trained across 119 languages and achieves top-rank multilingual benchmark results, an advantage for our Indonesian Twitter dataset. In addition, previous studies, including the PMP work, also utilized Llama-3, which we can also use for comparison with additionally Llama-3.1 due to its better reasoning. For the spaCy model, we ran it locally on our machine. The LLM models are run through Ollama and use the model endpoints directly. Across all experiments, we did not change any settings. No temperature, top_p, or decoding options were set and were left at default. Every model task call, such as keyword extraction, context retrieval, and pipeline execution, used the model defaults provided. This keeps behavior consistent across datasets, models, and pipeline variants but means results reflect the default decoding regime of the models rather than a tuned configuration, which is done to focus on the performance impact of additional context and methods.

C. Evaluation Metrics

To evaluate the performance of our pipeline, we use standard classification metrics such as accuracy, precision, recall, and F1-score. However, because of the imbalanced dataset class distribution, the calculation uses the macro average which calculates the metric independently for each class, then takes the unweighted mean across all classes.

D. Implementation Details

The evaluation is conducted under three experimental settings, with a total of five different pipelines as summarized in Table II. Word-info (LLM-only W_l) refers to the use of Token Tagging for keyword extraction, followed by context retrieval using an LLM. This setting is used especially for

English-themed datasets, such as Semeval-2018 and MUSTARD. Word-info (Google Search API W_g) denotes keyword extraction using LLM-based approach, followed by context retrieval using the Google Search API, with the retrieved information subsequently summarized by an LLM. This setting is used specifically for culture-specific datasets, such as the Indonesian Twitter dataset and Semeval-2018. The pairings of keyword extraction and context retrieval have been reasoned in Subsubsection III-A3.

V. RESULTS

In this section, we present outcomes of our experiments and analyze findings for each dataset. Table III summarizes all performance evaluation metrics.

A. Indonesian Twitter Dataset

We evaluated three variant methods on the Indonesian Twitter dataset: PMP, PMPWG, and PMPWG-FS. Using Qwen-3-8B, PMP improved over the baseline reported in the original paper, with a significant 0.1183 increase in macro-F1. This result reflects not only the advancement of modern LLMs, but also the effectiveness of PMP in adapting to culture-specific texts. However, earlier LLM attempts were inefficient as they showed high recall but very low precision, suggesting many guessed answers. In contrast, PMP demonstrated stable and meaningful predictions, as seen from the balanced precision and recall values across both labels.

Adding Google Search API grounding to PMP framework (PMPWG) produced further gains across metrics. Macro-F1 rose between PMP and PMPWG, from 0.5173 to 0.6108, with notable improvements in precision and recall. This indicates that adding external contextual information helps the model interpret culture-specific and making predictions more accurate.

Combining web grounding with few-shot demonstrations (PMPWG-FS) yielded the best performance, with a macro-F1 score of 0.6160 and an increase in accuracy of 0.0316. The inclusion of few-shot demonstrations appears to help the model capture task-specific patterns and nuances in the domain. The steady gains from PMP to PMPWG and PMPWG-FS highlight

TABLE III
EVALUATION RESULTS ACROSS DIFFERENT DATASETS, MODELS, AND PIPELINES.

Model	Method	Accuracy	Precision	Recall	Macro-F1
<i>Indonesian Twitter Dataset</i>					
Qwen-3-8B	PMP	0.6078	0.5189	0.5219	0.5173
	PMPWG	0.6673	0.6094	0.6363	0.6108
	PMPWG-FS	0.6989	0.6121	0.6225	0.6160
Llama-3-8B	PMP	0.4963	0.5324	0.5424	0.4794
	PMPWG	0.4907	0.5523	0.5661	0.4814
	PMPWG-FS	0.4665	0.5342	0.5426	0.4586
<i>SemEval 2018 Dataset</i>					
Qwen-3-8B	PMP	0.7602	0.7519	0.7600	0.7541
	PMPWG	0.7666	0.7620	0.7730	0.7626
	PMPWG-FS	0.7653	0.7581	0.7675	0.7601
	PMPWL	0.7704	0.7668	0.7783	0.7669
	PMPWL-FS	0.7985	0.7912	0.7840	0.7870
Llama-3-8B	PMP	0.7717	0.7615	0.7612	0.7614
	PMPWG	0.7959	0.7875	0.7956	0.7901
	PMPWG-FS	0.7679	0.7645	0.7415	0.7480
	PMPWL	0.8010	0.7924	0.7999	0.7950
	PMPWL-FS	0.7742	0.7692	0.7512	0.7570
<i>MUStARD Dataset</i>					
Qwen-3-8B	PMP	0.6420	0.6426	0.6420	0.6417
	PMPWL	0.6580	0.6619	0.6580	0.6559
	PMPWL-FS	0.6681	0.6703	0.6681	0.6671
Llama-3-8B	PMP	0.6348	0.6352	0.6348	0.6345
	PMPWL	0.5942	0.5987	0.5942	0.5895
	PMPWL-FS	0.5594	0.5757	0.5594	0.5344
Llama-3.1-8B	PMP	0.6043	0.6062	0.6043	0.6027
	PMPWL	0.6435	0.6435	0.6435	0.6435
	PMPWL-FS	0.5986	0.6003	0.5986	0.5968

Note: Precision, Recall, and Macro-F1 denote **macro-averaged** metrics computed across all classes.

that contextual grounding and example-driven prompting both improve sarcasm detection on culturally grounded Indonesian data.

When we repeat the same experiments with Llama-3-8B, overall results are much worse than Qwen-3-8B across all methods. In many runs, between PMP and PMPWG or PMPWG-FS shows no improvement or even degrades performance. We attribute this to Llama-3’s weaker multilingual and Indonesian coverage compared to Qwen-3. The model often fails to understand input tokens and retrieved context in Indonesian, so adding web context or few-shot examples mostly adds noise.

B. SemEval-2018 Task 3

Results on the SemEval-2018 Task 3 dataset further reinforce the effectiveness and adaptability of the proposed prompting strategies even with standard English sarcasm data. All evaluated variants such as PMP, PMPWG, PMPWG-FS, PMPWL, and PMPWL-FS, produced improvements for both Qwen-3-8B and Llama-3-8B in most runs, although the gains are often small. Based on this result, we can reflect that the LLM we are using is capable enough to understand those keywords.

The baseline PMP configuration achieved a strong macro-F1 score for both models, showing a reliable result using PMP even without additional contextual augmentation. Then, the largest single improvement using Llama-3-8B comes from internal knowledge context integration (PMPWL), increasing macro-F1 from 0.7614 to 0.7950. For Qwen-3-8B, the biggest jump is seen from PMPWL with few-shot examples (PMPWL-FS), rising from 0.7541 to 0.7870. These two results show that for English datasets the word-information block combined with PMP (and optionally with few-shot) can produce meaningful gains.

By contrast, adding Google Search API grounding and or with few-shot examples (PMPWG, PMPWG-FS) helps Llama-3-8B in several cases but not all. We believe this pattern follows from model coverage. For standard English captions, the LLM often already contains the necessary background knowledge, so web retrieval adds little and can introduce noise. Differences between models reflect their respective pretraining and multilingual strengths.

Overall, SemEval-2018 confirms that structured prompting plus concise context and examples is a strong recipe for sarcasm detection in English, and that the benefit of external retrieval depends on the model and dataset.

C. MUSTARD

Evaluation of the MUSTARD dataset further supports the effectiveness of our prompting methods in standard English sarcasm data. The PMP baseline achieved a macro-F1 of 0.6417 using Qwen-3-8B. Adding contextual internal knowledge from Token Tagging and combining that with few-shot examples (PMPWL-FS) produced the best result, with macro-F1 being 0.6671. The results show a modest but consistent improvement from structured prompting plus concise context and examples. Same as Semeval, the LLM can be reflected as capable enough to understand those keywords.

For Qwen-3-8B specifically, the jump from PMP to PMPWL-FS is on the order of a few percentage points, indicating that Qwen-3 benefits from this kind of context augmentation on MUSTARD. By contrast, when we use Llama-3-8B, the pattern is the opposite. Going from PMP to PMPWL or PMPWL-FS does not improve performance and can substantially degrade it, worst is from 0.6345 to 0.5344. But Llama-3.1-8B shows the reverse and even excels Qwen-3-8B. It improves best from PMP to PMPWL, going from 0.6027 to 0.6435.

We do not yet have a definitive explanation for this split. Possible causes to investigate include differences in pretraining corpora, multilingual coverage, tokenizer behavior, and instruction tuning or fine-tuning regimes between Llama-3 and Llama-3.1. These factors can change how a model interprets English tokens and how it benefits from added context. We recommend a follow-up analysis comparing token coverage on MUSTARD inputs, and checking release notes or model cards for Llama-3 vs Llama-3.1 to identify training differences that could explain the divergent behavior.

TABLE IV
GOOGLE SEARCH API RETRIEVAL RESULT EXAMPLES

Words	Definition (translated)	Correct?
wkwk	Wkwk is an expression of laughter that is often used in conversations on social media to show that someone is laughing, especially in the context of a funny or amusing situation.	V
kw	KW is an abbreviation of "Quality", which refers to counterfeit goods that imitate the original product, but the quality is not as high as the original product.	V
IPB	IPB is a university in Indonesia located in Bogor, with a focus on education and research in the field of agriculture and related sciences.	V
cie	CIE adalah singkatan dari **Commission Internationale de l'Éclairage**, sebuah organisasi internasional yang berperan dalam standarisasi dan penelitian tentang cahaya, warna, dan pencahayaan.	X

D. Google Search API

We analyzed the actual definitions returned by our Google Search API context retrieval. Table IV shows representative samples and indicates which source produced each entry.

External retrieval often recovered useful meanings for slang, regional expressions, and culture-specific references common in social media, and these recovered definitions clearly helped the downstream sarcasm classifier in many cases.

However, retrieval is noisy. Manual inspection revealed incorrect or misleading definitions for some entries. For example, a Google result for the Indonesian token "cie" returned an inaccurate gloss. In Indonesian, the usage "cie" is a playful tease directed at a couple, meant to make them blush, not a literal lexical definition. Table IV also reports a breakdown of sources and the proportion of manually flagged errors, illustrating that neither source is perfect and that simple filtering materially improves prompt quality. Therefore, we use not only Google Search API but also the LLM-only approach for context retrieval.

VI. CONCLUSION

Our experiments show that supplying targeted context improves sarcasm detection, especially for culturally grounded data. For the Indonesian Twitter dataset, adding non-parametric knowledge using Google Search API produced the largest gains. This finding supports the claim made by the original authors of the PMP method, who noted that PMP has limitations when dealing with datasets that exhibit extreme linguistic norms. Asking the model about its own knowledge for the context retrieval step as a self-knowledge or context-awareness strategy also helps, but the improvements are smaller. Both approaches are complementary, with self-knowledge boosting awareness, and external retrieval supplying factual grounding when the model lacks background knowledge.

On standard English datasets (e.g., MUSTARD, SemEval-2018), the benefits are smaller. We observed noticeable improvements from adding context and few-shot examples. These results suggest that for languages and domains well covered by the LLM's pretraining, the model already holds much of the needed knowledge. The main value of our methods there is to make the model more self-aware and to nudge its reasoning rather than supply new facts.

Nevertheless, the use of pretrained language models is still better and more reliable for the sarcasm detection. Even after adding contextual and example-driven layers, the PMP method performs worse compared to using PLMs.

Limitation & Future Improvements As mentioned, the hardware we used to run the experiments was not capable of running models with a larger computational cost. We also did not implement every possible combination of keyword extraction and context retrieval on every dataset. In particular, some pairings were avoided where an extractor was known to perform poorly or where retrieval would provide negligible benefit. These choices were pragmatic but should be revisited with larger compute budgets. Future improvements could also enhance both keyword extraction and especially context retrieval to perform more accurately, each with validations for a fallback pipeline, since our approach still has flaws.

VII. ACKNOWLEDGMENT

We would like to express our deepest gratitude towards Bina Nusantara University for supporting us in this particular study.

REFERENCES

- [1] J. Gu, X. Jiang, Z. Shi, H. Tan, X. Zhai, C. Xu, W. Li, Y. Shen, S. Ma, H. Liu, and S. Wang, "A survey on llm-as-a-judge," *arXiv preprint arXiv:2411.15594*, 2024.
- [2] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, and S. Riedel, "Retrieval-augmented generation for knowledge-intensive nlp tasks," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 9459–9474.
- [3] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 24 824–24 837.
- [4] H. Trivedi, N. Balasubramanian, T. Khot, and A. Sabharwal, "Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions," *arXiv preprint arXiv:2212.10509*, 2022.
- [5] J. Lee, W. Fong, A. Le, S. Shah, K. Han, and K. Zhu, "Pragmatic metacognitive prompting improves llm performance on sarcasm detection," in *Proceedings of the 1st Workshop on Computational Humor (CHum)*. Association for Computational Linguistics, 2025, pp. 63–70.
- [6] A. Mallen, A. Asai, V. Zhong, R. Das, D. Khashabi, and H. Hajishirzi, "When not to trust language models: Investigating effectiveness of parametric and non-parametric memories," *arXiv preprint arXiv:2212.10511*, 2022.
- [7] C. Mavromatis and G. Karypis, "Gnn-rag: Graph neural retrieval for efficient large language model reasoning on knowledge graphs," in *Findings of the Association for Computational Linguistics: ACL 2025*. Vienna, Austria: Association for Computational Linguistics, 2025, pp. 16 682–16 699.
- [8] S. Jeong, J. Baek, S. Cho, S. J. Hwang, and J. Park, "Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Mexico City, Mexico: Association for Computational Linguistics, 2024, pp. 7036–7050.
- [9] B. An, S. Zhang, and M. Dredze, "Rag llms are not safer: A safety analysis of retrieval-augmented generation for large language models," in *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Albuquerque, New Mexico: Association for Computational Linguistics, 2025, pp. 5444–5474.
- [10] M. Wang, Y. Yao, Z. Xu, S. Qiao, S. Deng, P. Wang, X. Chen, J.-C. Gu, Y. Jiang, P. Xie, F. Huang, H. Chen, and N. Zhang, "Knowledge mechanisms in large language models: A survey and perspective," in *Findings of the Association for Computational Linguistics: EMNLP 2024*. Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 7097–7135. [Online]. Available: <https://aclanthology.org/2024.findings-emnlp.416>
- [11] Y. Liang, Z. Song, H. Wang, and J. Zhang, "Learning to trust your feelings: Leveraging self-awareness in llms for hallucination mitigation," in *Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP*. Bangkok, Thailand: Association for Computational Linguistics, 2024, pp. 44–58. [Online]. Available: <https://aclanthology.org/2024.knowledgenlp-1.4>
- [12] Y. Wang, P. Li, M. Sun, and Y. Liu, "Self-knowledge guided retrieval augmentation for large language models," in *Findings of the Association for Computational Linguistics: EMNLP 2023*. Singapore: Association for Computational Linguistics, 2023, pp. 10 303–10 315. [Online]. Available: <https://aclanthology.org/2023.findings-emnlp.691>
- [13] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? sentiment classification using machine learning techniques," in *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - EMNLP '02*, 2002. [Online]. Available: <https://doi.org/10.3115/1118693.1118704>
- [14] R. Tokuhisa, K. Inui, and Y. Matsumoto, "Emotion classification using massive examples extracted from the web," in *ACL Anthology*, 2008, pp. 881–888. [Online]. Available: <https://aclanthology.org/C08-1111/>
- [15] S. Kennedy, N. Walsh, K. Sloka, A. McCarren, and J. Foster, "Fact or factitious? contextualized opinion spam detection," *arXiv preprint arXiv:1906.01815*, 2019.
- [16] M. Zhang, Yue, and G. Fu, "Tweet sarcasm detection using deep neural network," in *ACL Anthology*, 2016, pp. 2449–2460. [Online]. Available: <https://aclanthology.org/C16-1231/>
- [17] F. Barbieri, L. Espinosa Anke, and J. Camacho-Collados, "Xlm-t: Multilingual language models in twitter for sentiment analysis and beyond," 2022. [Online]. Available: <https://aclanthology.org/2022.lrec-1.27/>
- [18] Y. Cai, H. Cai, and X. Wan, "Multi-modal sarcasm detection in twitter with hierarchical fusion model," 2019.
- [19] L. Baker and A. McCallum, "Distributional clustering of words for text classification," in *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, 1998, pp. 96–103.
- [20] J. Rennie, "Improving multi-class text classification with naive bayes," 2001.
- [21] S. Tong and D. Koller, "Support vector machine active learning with applications to text classification," *Journal of Machine Learning Research*, vol. 2, pp. 45–66, 2001.
- [22] J. Ramos, "Using tf-idf to determine word relevance in document queries," in *Proceedings of the first instructional conference on machine learning*, vol. 242, no. 1, 2003, pp. 29–48.
- [23] W. Cavnar and J. Trenkle, "N-gram-based text categorization," in *Proceedings of SDAIR-94, 3rd annual symposium on document analysis and information retrieval*, 1994, pp. 161 175–14.
- [24] Y. Zhang, R. Jin, and Z. Zhou, "Understanding bag-of-words model: a statistical framework," *International Journal of Machine Learning and Cybernetics*, vol. 1, no. 1, pp. 43–52, 2010.
- [25] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, "Natural language processing (almost) from scratch," *Journal of Machine Learning Research*, vol. 12, no. 7, 2011.
- [26] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [27] S. Lai, L. Xu, K. Liu, and J. Zhao, "Recurrent convolutional neural networks for text classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, no. 1, 2015.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [29] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171–4186.
- [30] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. Aleman, D. Almeida, J. Altschmidt, S. Altman *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [31] G. Team, R. Anil, S. Borgeaud, J. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. Dai, A. Hauth, K. Millican, and D. Silver, "Gemini: A family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.
- [32] Y. Zhang, C. Zou, Z. Lian, P. Tiwari, and J. Qin, "Sarcasmbench: Towards evaluating large language models on sarcasm understanding," *arXiv preprint arXiv:2408.11319*, 2024.
- [33] S. Oprea and W. Magdy, "isarcasm: A dataset of intended sarcasm," *arXiv preprint arXiv:1911.03123*, 2019.
- [34] S. Castro, D. Hazarika, V. Pérez-Rosas, R. Zimmermann, R. Mihalcea, and S. Poria, "Towards multimodal sarcasm detection (an _obviously_ perfect paper)," *arXiv preprint arXiv:1906.01815*, 2019.
- [35] D. Suhartono, A. T. Handoyo, and F. Adeta Junior, "Feature-based augmentation in sarcasm detection using reverse generative adversarial network," *Computers, Materials & Continua*, vol. 77, no. 3, 2023.
- [36] D. Suhartono, W. Wongso, and A. T. Handoyo, "Idsarcasm: Benchmarking and evaluating language models for indonesian sarcasm detection," *IEEE Access*, vol. 12, pp. 87 323–87 332, 2024.
- [37] A. Baruah, K. Das, F. Barbhuiya, and K. Dey, "Context-aware sarcasm detection using bert," in *Proceedings of the Second Workshop on Figurative Language Processing*, 2020, pp. 83–87.
- [38] A. Kumar and V. Anand, "Transformers on sarcasm detection with context," in *Proceedings of the Second Workshop on Figurative Language Processing*, 2020, pp. 88–92.

- [39] Y. Liu, Y. Wang, A. Sun, X. Meng, J. Li, and J. Guo, "A dual-channel framework for sarcasm recognition by detecting sentiment conflict," *arXiv preprint arXiv:2109.03587*, 2021.
- [40] Y. Liu, R. Zhang, Y. Fan, J. Guo, and X. Cheng, "Prompt tuning with contradictory intentions for sarcasm recognition," in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2023, pp. 328–339.
- [41] A. Tsai, A. Kraft, L. Jin, C. Cai, A. Hosseini, T. Xu, Z. Zhang, L. Hong, E. Chi, and X. Yi, "Leveraging llm reasoning enhances personalized recommender systems," in *Findings of the Association for Computational Linguistics: ACL 2024*, 2024, pp. 13 176–13 188.
- [42] J. Chen, K. Du, X. Dai, W. Wang, X. Wang, Y. Wang, R. Tang, W. Zhang, and Y. Yu, "Debatecoder: Towards collective intelligence of llms via test case driven llm debate for code generation," in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2025, pp. 12 055–12 065.
- [43] P. Lewis, B. Oguz, R. Rinott, S. Riedel, and H. Schwenk, "Mlqa: Evaluating cross-lingual extractive question answering," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7315–7330.
- [44] Z. Wang, Y. Pang, and Y. Lin, "Large language models are zero-shot text classifiers," *arXiv preprint arXiv:2312.01044*, 2023.
- [45] M. Yu, X. Guo, J. Yi, S. Chang, S. Potdar, Y. Cheng, G. Tesauero, H. Wang, and B. Zhou, "Diverse few-shot text classification with multiple metrics," *arXiv preprint arXiv:1805.07513*, 2018.
- [46] Y. Wang and Y. Zhao, "Metacognitive prompting improves understanding in large language models," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2024, pp. 1914–1926.
- [47] C. Van Hee, E. Lefever, and V. Hoste, "Semeval-2018 task 3: Irony detection in english tweets," in *Proceedings of the 12th International Workshop on Semantic Evaluation*, 2018, pp. 39–50.
- [48] M. Ehrmann, A. Hamdi, E. Pontes, M. Romanello, and A. Doucet, "Named entity recognition and classification in historical documents: A survey," *ACM Computing Surveys*, vol. 56, no. 2, pp. 1–47, 2023.
- [49] L. Luo, Z. Yang, P. Yang, Y. Zhang, L. Wang, H. Lin, and J. Wang, "An attention-based bilstm-crf approach to document-level chemical named entity recognition," *Bioinformatics*, vol. 34, no. 8, pp. 1381–1388, 2018.
- [50] Z. Liu, F. Jiang, Y. Hu, C. Shi, and P. Fung, "Ner-bert: A pre-trained model for low-resource entity tagging," *arXiv preprint arXiv:2112.00405*, 2021.
- [51] ExplosionAI, "Model architectures · spacy api documentation," <https://spacy.io/api/architectures>.
- [52] S. HirpSSA and G. Lehal, "Pos tagging for amharic text: A machine learning approach," *INFOCOMP: Journal of Computer Science*, vol. 19, no. 1, 2020.
- [53] K. Akhil, R. Rajimol, and V. Anoop, "Parts-of-speech tagging for malayalam using deep learning techniques," *International Journal of Information Technology*, vol. 12, no. 3, pp. 741–748, 2020.
- [54] R. Mundotiya, V. Kumar, A. Mehta, and A. Singh, "Attention-based domain adaptation using transfer learning for part-of-speech tagging: an experiment on the hindi language," in *Proceedings of the 34th Pacific Asia Conference on Language, Information and Computation*, 2020, pp. 471–477.
- [55] M. Machado and E. Ruiz, "Evaluating large language models for the tasks of pos tagging within the universal dependency framework," in *Proceedings of the 16th International Conference on Computational Processing of Portuguese - Vol. 1*, 2024, pp. 454–460.
- [56] L. Bobojonova, A. Akhundjanova, P. Ostheimer, and S. Fellenz, "Bbpos: Bert-based part-of-speech tagging for uzbek," in *Proceedings of the First Workshop on Language Models for Low-Resource Languages*, 2025, pp. 287–293.
- [57] V. Rawte, S. Chakraborty, A. Pathak, A. Sarkar, S. Tonmoy, A. Chadha, A. Sheth, and A. Das, "The troubling emergence of hallucination in large language models - an extensive definition, quantification, and prescriptive remediations," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 2541–2573.
- [58] D. Wu, J. Zhang, and X. Huang, "Chain of thought prompting elicits knowledge augmentation," in *Findings of the Association for Computational Linguistics: ACL 2023*, 2023, pp. 6519–6534.
- [59] R. Ren, Y. Wang, Y. Qu, W. X. Zhao, J. Liu, H. Wu, J.-R. Wen, and H. Wang, "Investigating the factual knowledge boundary of large language models with retrieval augmentation," in *Proceedings of the 31st International Conference on Computational Linguistics*. Abu Dhabi, UAE: Association for Computational Linguistics, 2025, pp. 3697–3715.
- [60] M. Li, Y. Zhao, Y. Deng, W. Zhang, S. Li, W. Xie, S.-K. Ng, and T.-S. Chua, "Knowledge boundary of large language models: A survey," *arXiv preprint arXiv:2412.12472*, 2024. [Online]. Available: <https://arxiv.org/abs/2412.12472>

APPENDIX A

PROMPTS

We present the canonical PMP skeleton and then the insertion templates used to derive the hybrid variants, both in English and Indonesian. English datasets, such as Semeval and MUSTARD, use English prompts, and vice versa for the Indonesian Twitter dataset.

A. PMP first call (P_1)

English prompt:

You will be given a text, and will analyze the statement. Repeat back the statement to analyze. Then, analyze the following:

- What does the speaker imply about the situation with their statement?
- What does the speaker think about the situation?
- Are what the speaker implies and what the speaker thinks saying the same thing?

Finally, decide if the speaker is pretending to have a certain attitude toward the conversation.

Indonesian prompt:

Kamu akan diberikan sebuah teks dan diminta untuk menganalisis pernyataan di dalamnya. Ulangi kembali pernyataan yang akan dianalisis. Kemudian, analisis hal-hal berikut:

- Apa yang diimplikasikan oleh pembicara tentang situasi melalui pernyataannya?
- Apa yang dipikirkan pembicara tentang situasi tersebut?
- Apakah yang diimplikasikan dan yang dipikirkan pembicara menyampaikan hal yang sama?

Terakhir, tentukan apakah pembicara berpura-pura memiliki sikap tertentu terhadap percakapan tersebut.

B. PMP second call (P_2)

English prompt:

You will be given a piece of movie dialogue, a statement marked in brackets, and a preliminary analysis on the marked statement. Summarize the preliminary analysis and the given dialogue. Decide whether the statement is sarcastic or not by first analyzing the following:

The Implicature – What is implied in the conversation beyond the literal meaning?

The Presuppositions – What information in the conversation is taken for granted?

The Intent of the Speaker – What do the speaker(s) hope to achieve with their statement and who are the speakers?

The Polarity – Does the last sentence have a positive or negative tone?

Pretense – Is there pretense in the speaker’s attitude?

Meaning – What is the difference between the literal and implied meaning of the statement?

Reflect on the preliminary analysis and what should change, then decide if the statement is sarcastic.

Indonesian prompt:

Kamu akan diberikan sebuah pernyataan dan analisis awal terhadap pernyataan tersebut. Ringkas analisis awal tersebut. Tentukan apakah pernyataan tersebut bersifat sarkastik atau tidak dengan terlebih dahulu menganalisis hal-hal berikut:

Implikatur – Apa yang tersirat dalam percakapan di luar makna literal?

Presuposisi – Informasi apa dalam percakapan yang dianggap sudah diketahui?

Niat pembicara – Apa yang ingin dicapai pembicara dengan pernyataannya dan siapa pembicaranya?

Polaritas – Apakah kalimat terakhir bernada positif atau negatif?

Kepura-puraan – Apakah ada kepura-puraan dalam sikap pembicara?

Makna – Apa perbedaan antara makna literal dan makna tersirat dari pernyataan tersebut?

Renungkan analisis awal dan apa yang perlu diubah, lalu tentukan apakah pernyataan tersebut bersifat sarkastik.

C. Word-information Context Integration

Used in PMP with word-information context, either by LLM-approach (W_l) API or Google Search API (W_g), for PMPWL and PMPWG pipelines.

1) *User Prompt*: Insert the block below in the user prompt of P_1 immediately after the input line.

English prompt:

Entity facts:
token1 is a ... token2 is a ... token3 is a ...

Indonesian prompt:

Definisi kata-kata penting:
token1 adalah ... token2 adalah ... token3 adalah ...

2) *System Prompt*: Also, add the block below in the system prompt of P_1 , appending the last sentence.

English prompt:

There are also some entity facts from the sentence that you can use. Only use them if directly relevant, do NOT invent new facts.

Indonesian prompt:

Selain itu, ada disediakan beberapa fakta entitas dari kalimat yang dapat Anda gunakan. Hanya gunakan fakta tersebut jika langsung relevan, JANGAN menciptakan fakta baru.

D. Few-shot Examples Insertion

Used in PMPWL and PMPWG by injecting few-shot examples (F), with $k=2$, into the system prompt of P_2 after the instructions. Pipelines include PMPWL-FS and PMPWG-FS.

English prompt:

Here are example reflections:
Tweet: Great — the power went out during my online final.
Perfect timing
Implicature: The literal praise ("Great", "Perfect timing") contradicts the negative situation (power outage during an important exam); the speaker likely means the opposite.
Presuppositions: It is assumed the outage disrupted the exam and caused stress.
Speaker intent: To express frustration and criticize the situation indirectly, not to genuinely praise it.
Polarity: Literal wording is positive, but implied polarity is negative.
Pretense: There is clear pretense — the speaker is pretending to praise while actually conveying annoyance.
Meaning: The literal and implied meanings diverge (literal praise vs. implied complaint), indicating irony.
Final reflection: Strong contrast between wording and situation supports a sarcastic reading.
Final decision: YES Tweet: ...

Indonesian prompt:

Tweet: Bagus banget, listrik mati pas lagi final online.
Sumpah rejeki beneran
Implikatur: Kalimat tampak memuji ("Bagus banget") tetapi konteks (listrik mati saat final) jelas negatif; pembicara menyindir situasi.
Presuposisi: Diasumsikan listrik mati dan menimbulkan masalah pada ujian/online.
Niat pembicara: Mengungkapkan kekesalan dengan ironi, bukan benar-benar memuji.
Polaritas: Literal positif, implisit negatif.
Kepura-puraan: Ada pretense — pura-pura menyatakan kejelekan sebagai "bagus".
Makna: Perbedaan jelas antara makna literal (pujian) dan tersirat (keluhan).
Refleksi akhir: Kontras kata vs konteks kuat; indikator sarkastik jelas.
Keputusan akhir: YES
Tweet: ...

E. LLM-Based Keyword Extraction

Used in PMP with LLM-based keyword extraction, where pipelines include PMPWG and PMPWG-FS.

1) *Keyword Choosing*: Query the block below as the system prompt to the LLM with the user prompt being the target input text.

English prompt:

You will be given a text containing a list of unknown words. Your task is to separate the words into comma-separated values (CSV). If there are no unknown words, answer with 'NO UNKNOWN'
example output:
first,second,third

Indonesian prompt:

Anda akan diberikan teks berisi penjelasan kata-kata yang tidak dimengerti. Tugas anda adalah memisahkan kata-kata tersebut menjadi daftar yang dipisahkan koma (CSV). Kalau tidak ada kata-kata yang tidak dimengerti, jawab dengan 'NO UNKNOWN'
contoh output:
pertama, kedua, ketiga

2) *Choosement Cleaning*: After the LLM-based keyword extraction, query the LLM again with this block below as the system prompt for cleaning, with the user prompt being the output from the previous call, resulting in a clean list of unknown keywords according to the LLM.

English prompt:

You will be given a text. Your task is to identify words from the text that you do not understand.

Indonesian prompt:

Anda akan diberikan sebuah teks. Tugas anda adalah menyebutkan kata-kata yang anda tidak mengerti dari teks tersebut.

APPENDIX B DATA AVAILABILITY

The datasets used in this study are publicly available as follows:

- SemEval-2018 Task 3 [47]
- MUSTARD [34]
- Twitter Indonesia Sarcastic [36]

For further details regarding the experiments, please contact Michael Iskandardinata via email.

APPENDIX C AUTHOR CONTRIBUTIONS

We follow the CRediT taxonomy to state individual contributions, with roles in parentheses.

- Michael Iskandardinata (first author): Conceptualization (equal); Methodology (equal); Software (equal); Validation (support); Resources (lead); Investigation (lead); Data Curation (support); Writing - Review & Editing (equal); Visualization (lead);
- William Christian (co-author): Conceptualization (equal); Methodology (equal); Software (equal); Validation (lead); Resources (support); Investigation (support); Data Curation (lead); Writing - Review & Editing (equal); Visualization (support);
- Derwin Suhartono (supervisor): Conceptualization (equal); Writing - Review & Editing (equal); Supervision.

Notes: the primary inference designs and prompt engineering were done by Michael Iskandardinata. Both authors collaborated on Token Tagging and LLM retrieval code; William Christian implemented the Google Search API retrieval pipeline with LLM-cleaning. All authors reviewed and approved the final manuscript.