# TrueGL: A Truthful, Reliable, and Unified Engine for Grounded Learning in Full-Stack Search

# Joydeep Chandra\*

Dept. of Comp. Sci. & Tech.
Tsinghua University
Beijing, China
joydeepc2002@gmail.com

# Rim El Filali

Dept. of Comp. Sci. & Tech.
Tsinghua University
Beijing, China
rimo1902@gmail.com

# Aleksandr Algazinov\*

Dept. of Comp. Sci. & Tech.
Tsinghua University
Beijing, China
algazinovalexandr@gmail.com

# **Matt Laing**

Dept. of Psych. & Cog. Sci.
Tsinghua University
Beijing, China
matthieu.laing@gmail.com

# Satyam Kumar Navneet

Dept. of Comp. Sci. Engin.
Chandigarh University
Mohali, India
navneetsatyamkumar@gmail.com

# **Andrew Hanna**

Sch. of Mater. Sci. & Engin.

Tsinghua University

Beijing, China
andrewsamerhanna@gmail.com

Abstract-In the age of open and free information, a concerning trend of reliance on AI is emerging. However, existing AI tools struggle to evaluate the credibility of information and to justify their assessments. Hence, there is a growing need for systems that can help users evaluate the trustworthiness of online information. Although major search engines incorporate AI features, they often lack clear reliability indicators. We present TrueGL, a model that makes trustworthy search results more accessible. The model is a fine-tuned version of IBM's Granite-1B, trained on the custom dataset and integrated into a search engine with a reliability scoring system. We evaluate the system using prompt engineering and assigning each statement a continuous reliability score from 0.1 to 1, then instructing the model to return a textual explanation alongside the score. Each model's predicted scores are measured against real scores using standard evaluation metrics. TrueGL consistently outperforms other small-scale LLMs and rule-based approaches across all experiments on key evaluation metrics, including MAE, RMSE, and  $\mathbb{R}^2$ . The model's high accuracy, broad content coverage, and ease of use make trustworthy information more accessible and help reduce the spread of false or misleading content online. Our code is publicly available at https://github.com/AlgazinovAleksandr/TrueGL, and our model is publicly released at https://huggingface.co/JoydeepC/trueGL.

Index Terms—Reliability of Internet Data, Lie Detection in LLMs, LLM Fine-Tuning, Information Retrieval, Natural Language Processing

### I. Introduction

The rapid advancement of generative AI enabled the large-scale production of AI-generated content (AIGC). While these tools offer benefits in generating content quickly, they increase the spread of misinformation due to their ability to generate falsified information [46]. Large Language Models (LLMs) are trained on a wide range of datasets that include unreliable sources [50]. As a result, hallucinations and biases may occur [23]. AI-generated articles, deepfakes,

\* - Equal contribution

and synthetically created social media posts can now be used for disinformation campaigns [35]. The issue is amplified by the increasing amount of false information, with false claims spreading faster than truthful ones, particularly on social media platforms [40]. This spread is enhanced by recommendation algorithms that prioritize engagement over accuracy, undermining public trust in displayed content [35].

The integration of AI in search engines has transformed the way users access, retrieve, and process information. Leveraging Machine Learning (ML) and Natural Language Processing (NLP) algorithms has the potential to improve accuracy, relevance, and personalization in search results [4]. Despite these improvements, challenges affecting the reliability and accuracy of the generated and retrieved information remain. [5]. While traditional search engines aim to deliver results relevant to user queries, they often lack mechanisms to indicate the credibility of content. This gap can contribute to the widespread spread of false, misleading, or inaccurate information [11]. As search engines remain a primary tool for information retrieval, the ability to assess the trustworthiness and accuracy of results is becoming increasingly important. In this paper, we introduce **TrueGL**, an AI-driven framework for assessing the reliability of information through a fine-tuned model, addressing the critical challenge of misinformation in digital content across the web. Our contributions are summarized as follows:

• Fine-Tuned Granite-1B Model for Reliability Assessment. We present TrueGL, a system built on a fine-tuned Granite-1B model, optimized using vanilla fine-tuning on a custom dataset to assign continuous reliability scores (0 to 1) to statements, accompanied by textual justifications. Trained on 10 NVIDIA 3080 Ti GPUs for around 80 hours, TrueGL achieves superior performance on various metrics compared to baselines,

such as rule-based methods, TinyLlama-1.1B [47], Olmo-1B [18], and unfine-tuned Granite-1B [32], advancing AI-driven truthfulness evaluation.

- Comprehensive and Diverse Dataset for Truthfulness. We contribute a custom dataset of 140,272 articles across 21 domains, annotated with 35 distinct continuous reliability scores. Enhanced by novel fake article generation and data poisoning techniques (e.g., number manipulation, fact negation), this dataset enables robust training and evaluation, overcoming limitations of binary-labeled datasets.
- Evaluation Framework with Transparent Justifications. TrueGL introduces a robust evaluation methodology, testing on over 1100 human-AI co-annotated statements with three prompt variants (basic, zero-shot, and few-shot) and metrics (MAE, RMSE, and  $R^2$ ). The system provides transparent justifications for reliability scores via non-deterministic sampling, enhancing interpretability and user trust in AI assessments.

#### II. RELATED WORK

#### A. Assessing Statements and Articles

Rule-based Methods. Rule-based methods, such as Static Code Analysis (SCA), Symantic Rule Checking (SRC), and customized scripting tools, have been powerful tools used by developers to assess content [34]. These rules started out in the form of if-else statements, which comes with several challenges. Examples of these challenges include trouble predicting human behavior, inconsistencies when analyzing sentiments contained in statements, and the inability to cover every outcome. Therefore, as advanced NLP tools have developed, rule-based methods have evolved into statement sentiment analysis. This allows for fewer hard-coded rules, which reduces errors [34]. An example of rule-based methods being used with LLMs is Logic-RL [42], a framework that places rewards on answers that fully match the ground truth and punishes answers that are improperly formatted, mismatched, or missing information. Despite the small scale of the learning dataset, improved accuracy and reward were achieved. Additionally, model generalization was compared to other tested models and showed better results. However, the speed was significantly slower, and the output often exceeded the required length to express answers [42].

**Human Expertise**. Human expert review remains a widely used method for evaluating the credibility of information. One approach to determine the reliability of data is to use human experts to fact-check information. This approach is often referred to as human-in-the-loop [26] and has some key advantages: humans are able to discern more subtleties in content, detect biases, and determine the implications of the statement. QUEST [37] is an LLM being used in the healthcare sector; it utilizes human evaluations of LLM outputs to ensure the safety and reliability of AIGC. [27] also utilized experts

to generate and assess content that contains misinformation. They found that humans can come to a consensus on misinformation, an important factor when discussing information reliability. Human-in-the-loop methods try to combine the reliability of human fact-checking with the throughput of AI. The CEA-COVID [26] framework is an example of human-in-the-loop. By combining crowd-sourced data, expert opinions, and AI they were able to detect COVID-19-related misinformation effectively. However, human expertise has some major drawbacks, such as inherent human biases, varying levels of expertise between individuals, and conflicts of interest, which can influence experts' assessments [26] [27] [37].

Artificial Intelligence. AI models have demonstrated promise in detecting misinformation by analyzing features in content, linguistic patterns, and source credibility [11]. AIgenerated misinformation, called synthetic lies, is one such area of interest. Synthetic lies mostly propagate due to a lack of expertise to determine irregularities and inaccuracies in content. Using machine learning and human-in-the-loop approaches has been effective in flagging synthetic lies and false information. By focusing on semantic inconsistencies and abnormal phrasing, the origin and unreliability of content can be assessed effectively [11]. However, classification approaches, especially human-in-the-loop, have considerable drawbacks, as previously discussed. By differentiating between intrinsic and extrinsic hallucinations and then mitigating them using factchecking and Retrieval-Augmented Generation (RAG), hallucinations can be prevented [3]. Extrinsic hallucinations are especially damaging as LLMs are becoming more widespread in Information Retrieval (IR), and they are hard to catch [3].

By going further down the generation pipeline and augmenting the training of LLMs, hallucinations can be avoided. Optimizing models using knowledge constraints and reference grounding has effectively ensured that AIGC remains close to a grounded truth [14]. Additionally, approaches that use penalties for producing hallucinations are effective in encouraging more accurate content generation [15]. Broader strategies that combine layers of the systems discussed and augmented by AI have been successful at mitigating misinformation. Combining algorithmic, human-in-the-loop, and user education with AI effectively manages most known sources of misinformation generation and acceptance by users. While this holistic approach is effective, it requires a significant amount of energy and time to implement [6]. While AI approaches, such as NLP and graph-based methods, can augment IR and reliability assessment, the current best strategy is through human oversight. Humans have a great deal more ability in determining nuance and implementing judgment [10]

# B. AI in Search and Information Retrieval

Modern search engines utilize AI to power quick, accessible results for the user. Techniques such as NLP, deep learning, and semantic analysis provide results that are more based on context. Search engines use these techniques to evaluate the context of questions and provide answers that likely align with a user's intended search [4]. These systems also allow

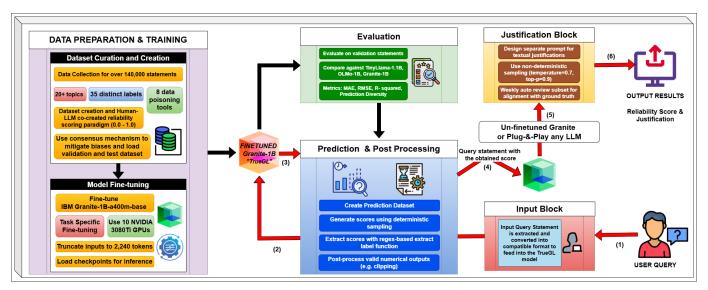


Fig. 1: TrueGL Workflow

TABLE I: Data Poisoning Methods Overview

Method	Description	Information Saved	Justification
Manipulation with Numbers	Finds all numbers and randomly chooses to either multiply them by a factor between 1.5 and 3, divide them by a similar factor, or delete them entirely	50%	Directly corrupts any quantitative information, which is often crucial (dates, statistics, etc.)
Antonym Swap- ping	Uses a dictionary to replace words like "yes" with "no," "good" with "bad," and "is" with "is not", and vice-versa, systematically inverting the meaning of the text	40%	Fundamentally changes the sentiment and factual statements of the original content
Insertion of Non- sense	Iterates through words and, with an 8% probability, inserts a random word from a large, predefined list of "nonsense" words	70%	Disrupts the text's natural flow and semantic coherence
Rearrangement (Partial)	Identifies random 20-word chunks within the article and shuffles the order of words within those specific chunks	55%	Breaks the logical structure and makes the content difficult to follow
Stemming	Applies PorterStemmer stemming [30] method from nltk library	62%	Significantly alters the text by removing gram- matical nuances, which makes it unnatural and hard to comprehend
Random Punctu- ation	Iterates through words and, with a 3% probability, adds characters such as "!", "?", ";", and "," after words	92%	Affects the text's appearance rather than its meaning
Duplication of Words	Iterates through words and, with a 5% probability, repeats words in the text	95%	Introduces only minor grammatical errors with- out significantly changing the text's meaning
Common Words Replacement	Some of the common words are replaced with minor misspelling errors (e.g., "the" is changed to "tha" and "with" is changed to "wit")	85%	Degrades text quality, but the overall meaning remains unchanged

users to rely less on external links, reducing the need for traditional click-through behaviors [5]. For example, this shift has especially affected advertising. The reduced reliance on primary page results moving towards AI-generated results means that Search Engine Optimization (SEO) strategies based on traditional search engines are losing value. While this poses challenges for advertisers, users and customers can now parse information more clearly and readily [17]. AI-augmented systems also impact the structure and function of search engine retrieval systems. Leveraging improvements in indexing accuracy, query clarifications, and multilingual support, AI

can be leveraged to connect related content, making search engines better at retrieving information and presenting it to the user [2]. Hybrid models, which augment AI, are not only efficient at retrieving information but also in imparting a sense of reliability and trust to users [12] [16]. An application of hybrid models is in Q&A engines. Search engines are natural avenues for Q&A engines, where inputs are often in the form of questions. However, parsing documents and articles in real-time and providing answers based on content can be a powerful tool for rapid and comprehensive IR [45]. Despite significant advances, AI-powered IR has some critique. AIGC can instill

a false sense of accuracy, omit important context, and instill bad behavior in users. The increasing reliance on AI, without proper education, has reduced reliance on critical thinking. The trend of free, unearned information has the capability to reduce the laymans ability to formulate nuanced opinions [13].

# C. Fine Tuning LLMs for Truthfulness and Reliability Assessment

Fine-tuning refers to the process of taking a pre-trained LLM and further increasing the training on a smaller, taskspecific dataset to improve its performance on specific goals and/or tasks [9]. A popular strategy is parameter-efficient fine-tuning (PEFT) [20] [19], which updates a small fraction of a model's parameters to save computation and memory, which updates a small fraction of a model's parameters to save computation and memory. One of the PEFT methods is LoRA [22], which introduces low-rank weight updates instead of modifying the full pre-trained model. To improve truthfulness and reliability, LLMs are fine-tuned on labeled datasets where claims are marked with reliability labels. For example, the FEVER dataset [38] includes claims supported or refuted by Wikipedia evidence, while SciFact [41] focuses on verifying scientific claims. FactTune [39] fine-tunes LLMs without human-labeled data by using synthetic factual signals, while LEAF [49] improves factual accuracy by integrating fact-checking steps directly into the training process. Other models, like RAFTS [44], enhance truthfulness by retrieving external evidence and comparing arguments, while PreCo-Fact (version 3) [48] combines classification and questionanswering methods to evaluate whether factual claims are true or false. A recent study [7] proposed a fine-tuning method that helps LLMs detect hallucinations by aligning outputs with universally accepted truths that are agreed upon globally. These approaches demonstrate that there are various approaches to fine-tuning LLMs for truthfulness and reliability assessment. However, to achieve high performance, high-quality and taskspecific datasets are required. The choice of dataset directly impacts a model's ability to detect factual inconsistencies, verify claims, and avoid hallucinations.

#### III. METHODOLOGY

The methodology of TrueGL, a trustworthy and truthful search engine based on cutting-edge AI techniques, comprises custom dataset creation, model fine-tuning, and a comprehensive evaluation. To evaluate the reliability of information on the internet, our method makes use of a task-specific, fine-tuned version of IBM's Granite-1B model [32], which has been improved with a carefully constructed reliability scoring system. Fig. 1 outlines the TrueGL's workflow. A user query enters the TrueGL system, where it is extracted and formatted for processing. The query is then analyzed by the TrueGL, a fine-tuned Granite-1B model trained on a custom dataset using 10 NVIDIA 3080 Ti GPUs for approximately 80 hours. The processed query moves to a scoring and post-processing stage, where a custom algorithm extracts and refines numerical data. This score is justified by an unfine-tuned Granite-1B

model, which provides a reliability explanation using nondeterministic sampling. The flexible TrueGL system allows model swaps for justification, with regular accuracy checks. Results are presented to the user with a truth score (0-100%) and a justification block explaining the score and source credibility.

#### A. Dataset Creation

The foundation of TrueGL's reliability assessment is a custom dataset designed for evaluating statement reliability. We constructed a dataset combining over 140,000 articles (statements). To ensure the diversity of topics, the data was collected from various sources on 21 different domains, including history, animals, and AI. Each statement is assigned a continuous reliability score on a [0.1, 1.0] scale, where 0.1 represents entirely false claims and 1.0 denotes fully reliable and truthful ones. Statements are annotated using the bias and credibility assessment platforms, such as MediaBiasFactCheck [31], AdFontesMedia [1], and NewsGuard [33]. The aggregated knowledge from these platforms was used to give a numerical score to each of the data sources. Most of the initial statements could be classified as either reliable or highly reliable. Hence, a fake article generation pipeline was developed to obtain unreliable statements. The pipeline generates topic-specific, well-written, and consistent texts. Meanwhile, the content contradicts common sense and scientific evidence, which can be observed. In addition, several data poisoning techniques (e.g., number manipulation, fact negation) were applied to create less reliable copies of the initial statements. Table I summarizes the objectives of these methods, as well as the ways they damage the reliability of texts. Each of the methods received a score, interpreted as the fraction of the original information preserved when the technique is applied to the text. Scores were obtained using multiple LLMs [24]. LLMs were provided with the original and damaged texts and had to evaluate how difficult it was to recover the original text from the modified one. For each method, several examples and several LLMs were used to estimate the final coefficients. These coefficients are multiplied by the original labels when data poisoning tools are applied to the articles. Hence, the way the tool affects the original reliability score depends on the amount of information it saves. Data poisoning algorithms were essential to create samples with mid-range reliability scores, ensuring the completeness and diversity of the fine-tuning data. Overall, the dataset has 35 distinct labels and serves as the basis for both training and evaluating TrueGL's ability to assess information reliability.

# B. Model Training

TrueGL is built upon IBM's Granite-1B model [32], fine-tuned to optimize performance for truthfulness and reliability assessment. Fine-tuning was done on the introduced dataset, training the model to predict reliability scores for statements. Input sequences were truncated to 2240 tokens. Based on the experiments, this number of tokens ensures that the model receives sufficient information to effectively assess the article's

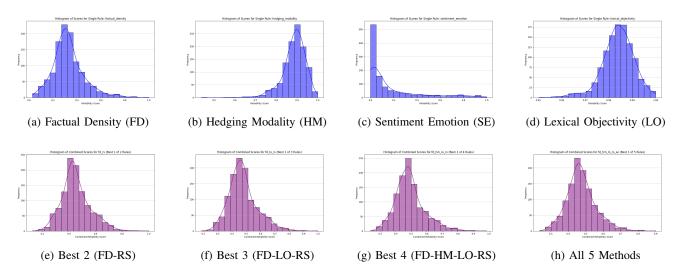


Fig. 2: Histograms Illustrating the Distribution of Reliability Scores for Best-Performing Individual Rule-Based Methods and Their Combinations

TABLE II: Performance Metrics of Rule-Based Reliability Scorer Configurations (Best  $R^2$ )

	Rule-based method	Weights	MAE	RMSE	${f R}^2$	Std Dev of Pred.
Single Methods	Factual Density (FD)	1.0	0.2879	0.3412	-0.2090	0.1309
	Hedging Modality (HM)	1.0	0.4508	0.5363	-1.9766	0.0570
	Readability Style (RS)	1.0	0.3808	0.4362	-0.9687	0.1896
	Sentiment Emotion (SE)	1.0	0.4432	0.5349	-1.9602	0.2386
	Lexical Objectivity (LO)	1.0	0.5301	0.6127	-2.8846	0.0067
Best Combinations	Best of 2 Methods	FD: 0.7, RS: 0.3	0.2644	0.3046	0.0399	0.1103
	Best of 3 Methods	FD: 0.7, RS: 0.2, LO: 0.1	0.2680	0.3061	0.0302	0.1011
	Best of 4 Methods	FD: 0.7, RS: 0.1, LO: 0.1, HM: 0.1	0.2724	0.3103	0.0034	0.0965
	Best of 5 Methods	FD: 0.6, RS: 0.1, LO: 0.1, HM: 0.1, SE: 0.1	0.2770	0.3155	-0.0302	0.0850

reliability while not making the input too large, helping save training time and avoiding GPU overload. The fine-tuning process was designed to ensure scalability and adaptability, allowing TrueGL to generalize across diverse content types while maintaining high accuracy in reliability scoring. Despite numerous advantages of PEFT methods, such as the decrease in GPU memory usage [43] and a significant reduction in the number of trainable parameters [22], the vanilla finetuning (VFT) [21] method was used to allow the full model adaptation. In addition, although PEFT methods have proved to achieve state-of-the-art performance and be more effective than VFT in several tasks (e.g., medical image analysis [29], selected software engineering tasks [51]), VFT often results in higher prediction accuracy [43]. Hence, the key motivation of using PEFT methods is time and computation efficiency, rather than an attempt to achieve a comparatively better performance. We fine-tuned the model using 10 NVIDIA 3080 Ti GPUs. The prompt used is given as follows:

Assess the reliability of the following statement on a scale of 0.1 (completely unreliable) to 1.0 (perfectly reliable). Consider factors such as factual accuracy, verifiability, number of alternative viewpoints, logical coherence (the statement is more reliable if it has no logical contradictions with facts that are easily confirmable or it has no

contradictions within the statement
itself), and evidence transparency (if
behind a statement there are transparent
methods such as statistical data, the
statement is more reliable). Provide only
the numerical score:\nStatement:
{text}\nLabel: {label}

The actual text of the article is inserted instead of the text variable, while the assigned reliability score is inserted instead of the label variable.

# C. Results Justification

Besides the reliability scoring model, a pipeline for the textual justification of the statement reliability was implemented. An LLM receives a prompt containing the original statement, the label generated by TrueGL, and the instruction as input and generates the justification for the score as output. The pipeline is flexible and allows using any LLM, as long as it supports the required input length of 2400 tokens (the length was slightly increased compared to the fine-tuning stage due to the longer instruction). For our experiments, an unfine-tuned Granite-1B [32] model was used. Model outputs were restricted to 300 new tokens so that the responses are both comprehensive and concise. The instruction prompt is given as follows:

You are an expert analyst. Your task is to provide a concise, bullet-point justification for a given reliability

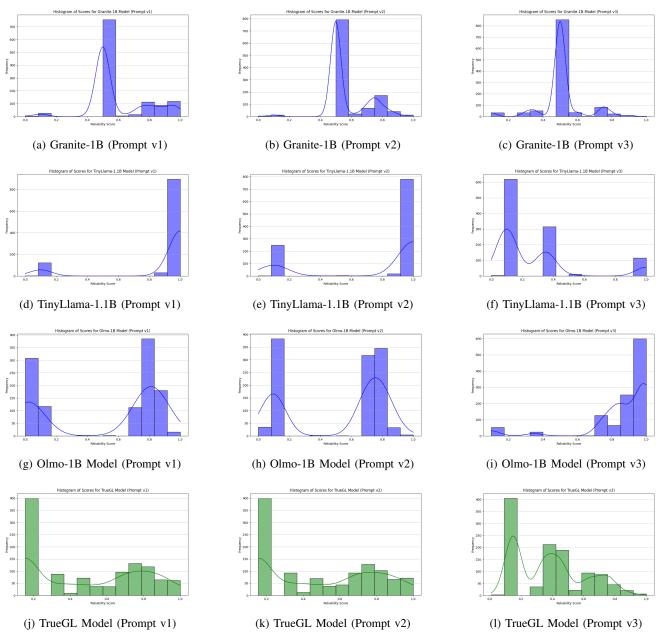


Fig. 3: Distribution of Reliability Scores of TrueGL and Other Compared Models. Here, v1 refers to the **Basic** prompt, v2 refers to the **Zero-Shot** prompt, and v3 refers to the **Few-Shot** prompt.

score of a statement. The reliability score is on a scale from 0.1 (completely unreliable) to 1 (perfectly reliable). Your justification should ONLY be the bullet points explaining the score. Do NOT repeat the statement or the score in your response.

The instruction is supplemented with several short examples. Each example consists of the statement, the reliability score, and the justification. One such instance is presented below:

Statement Context: The moon is made of green cheese and visited by cows weekly.

Assigned Reliability Score: 0.1

Correct Justification: - The statement contains obvious factual inaccuracies (moon is not made of cheese, and cows do not visit it).

- It presents scientifically implausible claims.
- Lacks any supporting evidence or credibility.

TABLE III: Comparison of Small-Scale LLMs vs TrueGL Across Different Prompting Strategies

	Model	MAE	RMSE	$\mathbb{R}^2$	Std Dev	Failure Rate
	Granite-3.1-1b-A400M-Base	0.29	0.33	-0.14	0.19	0
Basic Prompt	TinyLlama-1.1B-Chat-v1.0	0.57	0.65	-3.31	0.29	0.06
Basic Prompt	olmo-2-0425-1b-instruct	0.17	0.23	0.45	0.39	0
	TrueGL	0.09	0.14	0.78	0.31	0.01
	Granite-3.1-1b-A400M-Base	0.27	0.31	-0.02	0.14	0
Zero-Shot Prompt	TinyLlama-1.1B-Chat-v1.0	0.60	0.67	-3.59	0.38	0.06
Zero-snot Prompt	olmo-2-0425-1b-instruct	0.13	0.20	0.58	0.33	0
	TrueGL	0.09	0.15	0.78	0.30	0
	Granite-3.1-1b-A400M-Base	0.28	0.32	-0.09	0.12	0
Easy Chat Duament	TinyLlama-1.1B-Chat-v1.0	0.29	0.41	-0.76	0.28	0.05
Few-Shot Prompt	olmo-2-0425-1b-instruct	0.47	0.58	-2.46	0.21	0
	TrueGL	0.14	0.21	0.56	0.22	0

#### D. Evaluation

The TrueGL's performance was evaluated against three small-scale language models (TinyLlama-1.1B [47], Olmo-1B [18], and unfine-tuned Granite-1B [32]), and rule-based methods for assessing the reliability of data. The evaluation focused on the ability of models and algorithms to predict statement reliability scores. The evaluation dataset consisted of more than 1100 statements. Evaluation statements were not a part of the training set used for model fine-tuning and had the label distribution similar to the training set. Performance was measured using mean absolute error (MAE), root mean squared error (RMSE), and R-squared  $(R^2)$  metrics. In addition, the standard deviation of predictions was measured to check the diversity of predictions.

**Rule-Based Methods**. The system evaluates the reliability of statements by integrating multiple linguistic and content-based features, combining them through a grid search optimization process. The reliability scorer utilizes a set of predefined rules and heuristics, each contributing to a distinct aspect of statement credibility. The following rules are incorporated:

- Factual Density (FD): This metric quantifies the concentration of factual information within a statement, measured by the presence of named entities and numerical data. Higher factual density often correlates with increased verifiability. The implementation identifies factual entities using spaCy's named entity recognition and numerical tokens.
- Hedging Modality (HM): Hedging refers to linguistic devices used to express uncertainty or to soften a claim. This feature assesses the extent to which a statement uses modal verbs (e.g., 'might', 'could'), adverbs (e.g., 'possibly', 'perhaps'), or other linguistic cues that indicate a lack of absolute certainty. Lower hedging generally suggests stronger assertions. The score is inversely related to the density of predefined hedging words.
- Readability and Style (RS): This feature evaluates the linguistic complexity and stylistic characteristics of a statement. Metrics such as readability scores (e.g., Flesch-Kincaid) and sentence structure complexity contribute to this score, reflecting how easily a statement can be understood. Simplified language can sometimes be associated

with broader accessibility but also with a lack of nuance. The implementation utilizes the Flesch Reading Ease score and penalizes for stylistic elements like excessive capitalization or punctuation.

- Sentiment Emotion (SE): This rule analyzes the emotional tone expressed in a statement by quantifying the positivity, negativity, or neutrality of the text. While sentiment is not a direct indicator of factual reliability, extreme or overtly biased emotion can sometimes flag potentially unreliable content. The VADER (Valence Aware Dictionary and sEntiment Reasoner) lexicon is used to determine sentiment polarity, with scores adjusted to favor neutral sentiment.
- Lexical Objectivity (LO): This rule assesses the degree of objectivity in a statement by analyzing its lexical choices. It involves identifying subjective language, opinionated words, or evaluative adjectives versus more neutral and factual terminology. Statements with higher lexical objectivity are typically perceived as more credible. The score is derived by inverting the density of predefined subjective words.

To determine the optimal contribution of each rule-based feature, a grid search approach is used. This method explores a predefined range of weights for each feature, evaluating various combinations to identify the configuration that yields the best performance. The overall reliability score for a statement is calculated as a weighted sum of the individual feature scores, as expressed in the following equation:

Reliability Score = 
$$\sum_{i=1}^{n} \alpha_i \cdot \text{Method}_i$$
,

where n is the total number of rule-based methods considered, Method $_i$  represents the score obtained from the i-th rule-based method, and  $\alpha_i$  is the weight assigned to the i-th method, determined through the grid search process. These weights are non-negative and typically sum to 1, representing the proportional contribution of each method to the final score. The weights  $(\alpha_i)$  are iteratively adjusted across a defined range (from 0.1 to 1.0 with a step of 0.1) to find the combination that provides the most accurate reliability predictions.

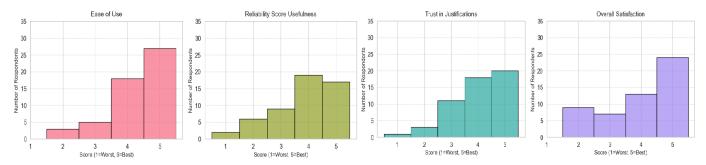


Fig. 4: User Survey Results Evaluating the TrueGL-Powered Search Engine

# IV. RESULTS AND ANALYSIS

#### A. Rule-Based Methods

Table II and Fig. 2 summarize and illustrate the bestperforming configurations based on  $R^2$  scores for both singlemethod approaches and increasingly complex combinations. While the best rule-based combination (FD+RS) achieved only marginal positive  $R^2$  (0.0399), TrueGL consistently demonstrated superior performance across all prompting strategies in Table III, with  $R^2$  scores reaching 0.78 for basic and zero-shot prompts (nearly  $20 \times$  higher than the rule-based approach). All individual rule-based methods failed, which is indicated by the negative  $R^2$ . This performance gap shows that rule-based methods lack the capacity to reliably assess statement quality, as they cannot capture the complex linguistic patterns that the TrueGL model learns through its neural architecture. While rule combinations showed minor improvements over single methods, their rapid performance degradation with additional rules suggests inherent limitations in linear combinations of handcrafted features.

# B. Instruction-Based LLM tuning

We designed three prompts, following an incremental approach. The **Basic** (v1) prompt provides basic scoring instructions; the **Zero-Shot** (v2) prompt is an extension of the **Basic** prompt, and the extension is made by specifying evaluation criteria, including factual accuracy, logical consistency, and evidence transparency. The **Few-Shot** (v3) prompt incorporated three annotated examples (unreliable, barely reliable, and very reliable) to show the scoring logic. These examples are shortened versions of the statements from the test sample. All prompts directed models to produce scores on a [0.1, 1.0] scale, with a post-processing step to extract and validate numerical outputs.

The test set consists of more than 1100 statements with AI-human co-annotated reliability scores. For each of the prompts, we evaluated the performance of four small-scale LLMs in predicting statement reliability scores: TinyLlama-1.1B [47], Olmo-1B [18], Granite-1B [32], and TrueGL. Model outputs were restricted to 20 new tokens to ensure concise numerical outputs can be extracted from these tokens. To ensure deterministic outputs, greedy sampling was applied. Performance was measured using MAE, RMSE, and  $\mathbb{R}^2$ .

Metrics were calculated by predicted scores with ground truth labels. Table III summarizes the results of the evaluation. Across all three prompts and metrics, TrueGL consistently outperforms other models. The few-shot prompt, which was expected to result in the best performance, turned out to be the least effective. The possible cause was the long prompt length. Small-scale models, used for the evaluation, cannot effectively support long contexts. Hence, additional instructions might be misleading rather than useful. Since the metrics are the least sensitive to prompt changes, TrueGL is the most robust out of the evaluated models. Fig. 3 shows the distributions of predictions across the evaluated models. Granite-1B's performance is sensitive to the prompt, with v3 producing the most consistent outcomes. Both TinyLlama-1.1B and Olmo-1B frequently exhibit bimodal distributions, indicating that they generate outputs with either very low or very high reliability. On the contrary, the TrueGL model typically exhibits a wider and more diverse distribution of scores, making it more suitable for the statement reliability assessment than other models.

# C. User Experience

To assess the usability of TrueGL in web information applications, we developed a prototype search engine powered by TrueGL. Unlike traditional search engines, our system provided users with reliability scores and optional textual justifications for the trustworthiness of each retrieved document. We conducted a user survey with 53 participants, who tested the prototype and evaluated their experience based on four parameters:

- Ease of Use: How intuitive and user-friendly the search engine was.
- Reliability Score Usefulness: Whether the provided credibility scores helped users assess information quality.
- Trust in Justifications: How convincing and wellreasoned the model's explanations for reliability scores were.
- Overall Satisfaction: The general impression of the system's effectiveness in aiding trustworthy information retrieval.

Participants rated each parameter on a 1–5 scale, where 1 means completely dissatisfied and 5 means fully satisfied. The results in Fig. 4 showed positive feedback, with most

responses centered at 4 and 5 across all categories. 70% of users gave 4 or 5 stars for Overall Satisfaction, indicating high approval of the developed system. These findings suggest that integrating TrueGL into search engines significantly improves user confidence in evaluating online information, setting a new standard for AI-assisted credibility assessment.

#### V. LIMITATIONS AND FUTURE WORK

While TrueGL represents an advancement in evaluating information reliability, our research faces several limitations:

Static Source Labeling. The current implementation of TrueGL assigns fixed reliability scores to sources, such as consistently rating BBC articles at 0.9, regardless of individual article quality. This static approach overlooks variability in content accuracy, author expertise, or context-specific factors, such as editorial bias or factual errors in specific articles. Dynamic source evaluation, incorporating article-specific features, is necessary to enhance reliability assessments [36].

Limited Use of Metadata. TrueGL primarily relies on the main content of articles, neglecting valuable metadata such as titles, publication dates, author credentials, number of references, or citation counts. These metadata elements can provide critical signals for reliability, as demonstrated in prior work on credibility assessment [25]. For instance, a sensationalist title may indicate clickbait, while a high number of references could suggest more research, yet these are not currently utilized.

Model Scale and Generalization. The TrueGL model is relatively small, limiting its capacity to capture complex patterns in diverse datasets. While efficient, this constrains its performance compared to larger models fine-tuned with techniques like PEFT. A small model may struggle with nuanced reliability judgments, particularly for ambiguous or context-heavy content, impacting its robustness across varied domains.

Limited Content Scope. TrueGL's training data focuses on general web content, with insufficient coverage of specialized domains like fake news or scientific literature. This limits its ability to detect misinformation in highly technical or deceptive contexts, such as identifying fabricated scientific claims or sophisticated phishing attempts. Expanding the dataset to include these domains is critical for broader applicability [8].

To address these limitations, we propose several directions for enhancing TrueGL's capabilities:

**Dataset Expansion**. We aim to extend the dataset to include diverse content types, such as fake news and scientific literature, enabling robust detection of misinformation and scholarly reliability. Incorporating metadata, particularly titles, will improve reliability assessments by leveraging signals like sensationalism or keyword relevance, building on frameworks like those in [25]. For example, analyzing title-content alignment could help identify clickbait or misleading articles.

**Light-weight Browser Plug-in**. We envision TrueGL evolving into a lightweight web extension, seamlessly integrated into browsers to provide real-time reliability scores. This plug-in would protect users from scams and phishing by delivering fast, AI-driven assessments of web content on the go, enhancing user safety in dynamic online environments. Such an application aligns with recent trends in AI-driven browser tools [28].

New Foundation Model. Existing foundational models are often trained on broad, general-purpose corpora. Hence, these models may not perform well on the complex, domain-specific task of evaluating the information's reliability. When training a model from scratch, the pretraining process can be focused on the features that are most important for reliability assessment. Furthermore, training a model from scratch ensures full control over the architecture. This adaptability is essential for investigating task-specific innovations, such as domain-aware scoring mechanisms, justification generation, and confidence calibration.

# VI. CONCLUSION

Misinformation across digital platforms, in particular caused by recent advancements in AIGC, remains a major concern. Hence, the need for reliable and trustworthy information is crucial. In this paper, we introduced TrueGL, an AI-driven framework that assesses statements with continuous reliability scores and provides transparent justifications for given scores. The model is developed by fine-tuning IBM's Granite-1B model using the vanilla fine-tuning technique. A customcurated dataset, containing over 140,000 articles on diverse topics, was used for model training. Fake article generation and data poisoning tools were used to ensure the presence of different reliability levels. Numerous experiments show that TrueGL consistently outperforms rule-based methods and other small-scale LLMs on various prompts and metrics. In addition, it makes more diverse predictions compared to other LLMs, further justifying its usefulness in assessing the reliability of data across the web.

TrueGL enables users to navigate online content with ease, clarity, and confidence, overcoming the risks of interacting with misleading information. TrueGL's reliability evaluation system ensures that users can filter out invalid information, setting a new benchmark for assessing the reliability and trustworthiness of online information. Our evaluation of TrueGL demonstrates that it is effective in delivering reliable and trustworthy information, as shown by a user survey conducted with 53 participants. The survey assessed the ease of use, reliability score, usefulness, trust in justifications, and overall satisfaction on a 1-5 scale. The results showed positive feedback, with most participants rating 4 or 5 across all parameters, indicating a high relevance of using reliability assessment LLM-based methods in search engines. By combining technical innovation with positive user experiences, TrueGL takes a clear step toward democratizing access to trustworthy information, with future work aimed at developing TrueGL into a lightweight web extension focused on protecting users from unreliable and untrustworthy misinformation.

#### REFERENCES

- [1] Ad Fontes Media. Ad fontes media, 2025.
- [2] Najeem Olawale Adelakun. Exploring the impact of artificial intelligence on information retrieval systems, 2024.
- [3] Pegah Ahadian and Qiang Guan. A survey on hallucination in large language and foundation models, 2025.
- [4] Ahmed Shaker Alalaq. Ai-powered search engines, 2025.
- [5] Eslam Amer and Tamer Elboghdadly. The end of the search engine era and the rise of generative ai: A paradigm shift in information retrieval, 2024
- [6] Juan Gómez Romero Andrés Montoro Montarroso, Javier Cantón-Correa. Fighting disinformation with artificial intelligence: Fundamentals, advances, and challenges, 2023.
- [7] Lennart Bürger, Fred A. Hamprecht, and Boaz Nadler. Truth is universal: Robust detection of lies in llms, 2024.
- [8] Y. Chen and S. Liu. Building diverse datasets for fake news detection and scientific content analysis. *Journal of Data Science*, 20(2):345–362, 2022.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [10] Hamid Reza Saeidnia et al. Artificial intelligence in the battle against disinformation and misinformation: A systematic review of challenges and approaches, 2025.
- [11] Jiawei Zhou et al. Synthetic lies: Understanding ai-generated misinformation and evaluating algorithmic and human solutions, 2023.
- [12] Kai Zhang et al. Hybrid search engine for credibility and accuracy evaluation in web results, 2022.
- [13] Pranav Narayanan Venkit et al. Search engines in an ai era: The false promise of factual and verifiable source-cited responses, 2024.
- [14] Su Ji et al. Understanding and reducing the hallucination problem in generative models, 2021.
- [15] Wen Zhao et al. Learning to prevent hallucination in open-domain dialogue generation, 2020.
- [16] Xiaoyan Huo et al. A truthful hybrid search engine model based on ai and crowdsourcing, 2022.
- [17] Gabriel Garlough-Shah. The rise of ai-powered search engines: Implications for online search behavior and search advertising, 2024.
- [18] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models. 2024.
- [19] Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning for large models: A comprehensive survey, 2024.
- [20] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp, 2019.
- [21] Jeremy Howard and Sebastian Ruder. Universal language model finetuning for text classification, 2018.
- [22] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- [23] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM Computing Surveys, 55(12):1–38, March 2023.
- [24] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego 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. Mistral 7b, 2023.
- [25] L. Jones and M. Taylor. Leveraging metadata for web content credibility assessment. ACM Transactions on Information Systems, 37(3):1–25, 2019
- [26] Ziyi Kou, Lanyu Shang, Yang Zhang, Zhenrui Yue, Huimin Zeng, and Dong Wang. Crowd, expert & ai: A human-ai interactive approach towards natural language explanation based covid-19 misinformation detection. In *IJCAI*, pages 5087–5093, 2022.
- [27] Ariel Kruger, Morgan Saletta, Atif Ahmad, and Piers Howe. Structured expert elicitation on disinformation, misinformation, and malign influence: Barriers, strategies, and opportunities. *Harvard Kennedy School Misinformation Review*, 2024.
- [28] H. Li and Z. Wang. Ai-driven browser extensions for real-time content analysis. *IEEE Internet Computing*, 28(1):56–63, 2024.
- [29] Chenyu Lian, Hong-Yu Zhou, Yizhou Yu, and Liansheng Wang. Less could be better: Parameter-efficient fine-tuning advances medical vision foundation models, 2024.
- [30] Julie Beth Lovins. Development of a stemming algorithm. Mechanical Translation and Computational Linguistics, 11(1-2):22–31, 1968.
- [31] Media Bias/Fact Check. Media bias/fact check, 2025.
- [32] Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, Manish Sethi, Xuan-Hong Dang, Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew Coleman, Matthew White, Mark Lewis, Raju Pavuluri, Yan Koyfman, Boris Lublinsky, Maximilien de Bayser, Ibrahim Abdelaziz, Kinjal Basu, Mayank Agarwal, Yi Zhou, Chris Johnson, Aanchal Goyal, Hima Patel, Yousaf Shah, Petros Zerfos, Heiko Ludwig, Asim Munawar, Maxwell Crouse, Pavan Kapanipathi, Shweta Salaria, Bob Calio, Sophia Wen, Seetharami Seelam, Brian Belgodere, Carlos Fonseca, Amith Singhee, Nirmit Desai, David D. Cox, Ruchir Puri, and Rameswar Panda. Granite code models: A family of open foundation models for code intelligence, 2024.
- [33] NewsGuard Technologies, Inc. Newsguard, 2025.
- [34] Paramita Ray and Amlan Chakrabarti. A mixed approach of deep learning method and rule-based method to improve aspect level sentiment analysis. Applied Computing and Informatics, 18(1/2):163–178, 2022.
- [35] Mohamed R. Shoaib, Zefan Wang, Milad Taleby Ahvanooey, and Jun Zhao. Deepfakes, misinformation, and disinformation in the era of frontier ai, generative ai, and large ai models, 2023.
- [36] J. Smith and A. Brown. Dynamic source credibility assessment in online news. *IEEE Transactions on Computational Social Systems*, 10(4):1234– 1245, 2023.
- [37] Thomas Yu Chow Tam, Sonish Sivarajkumar, Sumit Kapoor, Alisa V Stolyar, Katelyn Polanska, Karleigh R McCarthy, Hunter Osterhoudt, Xizhi Wu, Shyam Visweswaran, Sunyang Fu, et al. A framework for human evaluation of large language models in healthcare derived from literature review. NPJ digital medicine, 7(1):258, 2024.
- [38] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale dataset for fact extraction and verification, 2018.
- [39] Kai Tian, Emma Mitchell, et al. Facttune: A method to fine-tune llms for factual accuracy without human feedback, 2023.
- [40] Ivan Vykopal, Matúš Pikuliak, Ivan Srba, Robert Moro, Dominik Macko, and Maria Bielikova. Disinformation capabilities of large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), page 14830–14847. Association for Computational Linguistics, 2024.
- [41] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. Fact or fiction? verifying scientific claims, 2020.
- [42] Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning. arXiv preprint arXiv:2502.14768, 2025.
- [43] Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment, 2023.
- [44] Zirui Yue, Honglei Zeng, Lili Shang, Yue Liu, Yiming Zhang, and Dong Wang. Rafts: Retrieval augmented fact verification by synthesizing contrastive arguments, 2024.
- [45] Chafiq Nadia Zaher Najwa, Ghazouani Mohamed. Revolutionizing information retrieval: Unveiling a next-generation ai-powered questionanswer system for comprehensive document analysis, 2024.

- [46] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news, 2020.
- [47] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model, 2024.
- [48] Xiaoyu Zhang, John Doe, et al. Pre-cofactv3: A comprehensive framework comprised of qa and text classification for fact verification, 2024.
- [49] Yuqing Zhang, Karthik Thiyagarajan, et al. Leaf: Learning and evaluation augmented by fact-checking, 2024.
- [50] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2025.
- [51] Wentao Zou, Qi Li, Jidong Ge, Chuanyi Li, Xiaoyu Shen, Liguo Huang, and Bin Luo. A comprehensive evaluation of parameter-efficient finetuning on software engineering tasks, 2023.