

---

# DIAL-G<sup>2</sup>: GRAPH-GUIDED DIALECTICAL AGENT FOR ADVANCED ESG REASONING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The rapid growth of the focus on Environmental, Social, and Governance (ESG) creates the need for an effective AI tool for evaluating corporate sustainability. Despite the increasing availability of ESG data disclosed by companies, the complex, relational, and unstructured nature of ESG reports poses significant challenges for large-scale LLMs. The surprising efficacy of small language models (SLMs) on domain-specific tasks challenges the prevailing belief that only massive, general-purpose models can tackle complex reasoning. Inspired by this observation, we investigate the potential of compact models in the nuanced field of ESG analysis, arguing that the path to true expertise lies in focused knowledge and structured collaboration. Our work presents a two-fold contribution to realize this vision: 1) We introduce ESGEXPERT-30K, a new knowledge-intensive dataset, to fine-tune a compact SLM into a specialized ESGEEK model that achieves state-of-the-art accuracy on domain-specific QA. 2) We propose DIAL-G<sup>2</sup>, a novel framework where these expert agents form a committee to analyze full, multimodal corporate reports. In DIAL-G<sup>2</sup>, agents populate a shared graph with conflicting arguments. The GNN then performs relational inference on this graph and uses its learned attention weights to direct the agents' focus towards the most salient and contested information in a subsequent debate phase. This graph-guided approach demonstrates remarkable effectiveness. Furthermore, on our newly contributed, large-scale ESGREPORT-RATING-50K benchmark, DIAL-G<sup>2</sup> achieves human-expert-level performance in end-to-end ESG score prediction. Extensive experiments show that DIAL-G<sup>2</sup> can overcome bottlenecks of Multi-Agent Systems (MAS) in domain-specific understanding, providing the ESG research community with a powerful new paradigm for scalable, relational, and interpretable AI.

## 1 INTRODUCTION

Numerous studies have established a positive correlation between Environmental, Social, and Governance (ESG) performance and firm value (Gillan et al., 2021; Wong et al., 2021; Serafeim & Yoon, 2023). And stakeholders evaluate the ESG performance of companies based on their annual ESG reports. Unlike data from financial statements, regulators do not mandate ESG data filing in specific forms in most cases. Efforts to standardize ESG reports face challenges like complex regulations, evolving frameworks, and compliance verification (Mishra et al., 2024). With the rapid development of Large Language Models (LLMs), the call for evaluating ESG data using LLMs prevails.

However, research on LLMs for ESG has bottlenecks due to inherent framework-related limitations. While fine-tuned LLMs are applied to such tasks (Birti et al., 2025), knowledge shortage and hallucination make it not applicable. To ground LLMs in external knowledge and mitigate hallucination, Retrieval-Augmented Generation (RAG) variants, like multi-modal RAG (Wang et al., 2025b), hybrid RAG (Wang et al., 2025a), and graph-based RAG (Edge et al., 2025; Guo et al., 2025; Barghi, 2025), have become a dominant paradigm. However, these RAG-based approaches often follow a linear "retrieve-then-synthesize" workflow. They lack a mechanism for deep relational reasoning, iterative refinement, or for resolving subtle, conflicting evidence distributed across a lengthy document. Reinforcement Learning(RL) with LLM for ESG is also researched (Hou et al., 2025), but the utilization of general LLMs in the framework results in the lack of domain knowledge and massive resource consumption, making it not suitable for production scenarios.

---

054 In this paper, we propose a new paradigm that embraces specialization and structured collaboration.  
055 By fine-tuning on our knowledge-intensive dataset ESGEXPERT-30K, we have developed a series  
056 of efficient, compact language models with a deep understanding of ESG, which we refer to as ES-  
057 GEEKS. This ensures that our AI agents have the cognitive foundation to become experts in the  
058 ESG field. These expert agents form the core of our innovative framework, **Dialectical Agents with**  
059 **Graph-Guide (DIAL-G<sup>2</sup>)**.<sup>1</sup> It represents reports as an evidence graph where a GNN acts as a dy-  
060 namic conductor, using attention to identify salient “hotspots” and guide a dialectical debate among  
061 ESGEEK-based agents. The outcome from the debate iteratively refines the graph’s beliefs, creat-  
062 ing a neuro-symbolic loop that enables deep, relational, and interpretable analysis beyond standard  
063 retrieval methods. We validate our approach through extensive experiments. We first confirm the  
064 SOTA performance of our ESGEEK model. Then, using our newly created ESGREPORT-RATING-  
065 50K benchmark, we show that the full DIAL-G<sup>2</sup> framework can predict a company’s ESG score  
066 with a level of accuracy comparable to human experts. Our contributions are:

- 067 • We introduce and validate a method for creating highly specialized, efficient **ESGEEK**  
068 models by fine-tuning on our new, knowledge-rich **ESGEXPERT-30K** dataset.
- 069 • We propose **DIAL-G<sup>2</sup>**, a novel framework that deeply fuses a GNN’s relational reasoning  
070 with the linguistic intelligence of specialist agents through a unique, iterative, graph-guided  
071 debate mechanism.
- 072 • We contribute **ESGREPORT-RATING-50K**, a large-scale, multimodal dataset for ESG  
073 score prediction, addressing a key resource gap.
- 074 • We demonstrate that our complete system achieves SOTA performance and provides a col-  
075 laborative architecture paradigm composed of specialized small models for accomplishing  
076 complex real-world tasks.

## 079 2 RELATED WORKS

### 081 2.1 THE IMPORTANCE AND CHALLENGES OF ESG EVALUATION

083 Research has indicated a direct, positive impact of high ESG scores on a company’s operational per-  
084 formance and efficiency (Chen & Xie, 2022). Furthermore, ESG information plays a crucial role  
085 in shaping market perception and reducing information asymmetry (Cornell, 2021). ESG-related  
086 controversies are shown to heighten perceived uncertainty surrounding a firm’s future cash flows  
087 Schiemann & Tietmeyer (2022), while transparent and comprehensive ESG disclosures tend to miti-  
088 gate this uncertainty by providing stakeholders with decision-useful information (Serafeim & Yoon,  
089 2023). However, some studies suggest that in cases of manual rating, increased ESG disclosure may  
090 exacerbate rather than resolve ESG rating discrepancies (Christensen et al., 2022; Dimson et al.,  
091 2020), which contradicts our expectations. Consequently, leveraging artificial intelligence to ad-  
092 dress these pain points—enhancing objectivity, scalability, and the ability to process vast, unstruc-  
093 tured data—has become a critical area of research, motivating the work presented in this paper.

### 094 2.2 DOMAIN-SPECIFIC MODEL SPECIALIZATION

096 The concept of creating expert models for high-stakes domains is well-established. Models like Fin-  
097 BERT (Liu et al., 2020) and BloombergGPT (Wu et al., 2023a) demonstrated the value of domain-  
098 specific pre-training. In medicine, models like Med-PaLM 2 have shown expert-level performance  
099 by fine-tuning on medical knowledge (Qian et al., 2024). Amila Silva et al. simultaneously retain  
100 domain-specific and cross-domain knowledge in multimodal data to detect fake news from different  
101 domains (Silva et al., 2021). As for the ESG field, Tim Nugent et al. combined domain-specific  
102 language models and data augmentation methods to detect ESG issues (Nugent et al., 2021), im-  
103 proving the accuracy of classification tasks. ClimateQA (Luccioni et al., 2020) was developed to  
104 classify whether a sentence from an ESG report answers regulatory questions. ESGenius (He et al.,  
105 2025) has established the first LLM evaluation benchmark focused on ESG and sustainable devel-

---

106 <sup>1</sup>Code, dataset, and appendices with full experimental details will be made publicly available at [https://](https://anonymous.4open.science/r/dial_g2_ESGeek_framework-E613/README.md)  
107 [anonymous.4open.science/r/dial\\_g2\\_ESGeek\\_framework-E613/README.md](https://anonymous.4open.science/r/dial_g2_ESGeek_framework-E613/README.md) for repro-  
ducibility.

---

108 opment knowledge QA. Our work, built on knowledge distillation and specialization for ESG QA,  
109 follows this lineage.

## 111 2.3 MULTI-AGENT SYSTEMS AND GRAPH-GUIDED REASONING

113 To move beyond the limitations of simple RAG, researchers are exploring more sophisticated rea-  
114 soning frameworks. LLM-powered multi-agent systems (MAS) enhance the capabilities of a sin-  
115 gular model by orchestrating collaborative dialogues between specialized, role-playing agents (Chen  
116 et al., 2023a;b; Zhang et al., 2024a; Han et al., 2025). Concurrently, Graph Neural Networks (GNNs)  
117 (Velickovic et al., 2018; Chen et al., 2018; Kipf & Welling, 2016) have become the standard for  
118 modeling relational data, proving highly effective in document analysis for tasks like modeling cita-  
119 tions or logical flow (Liu et al., 2021; Wu et al., 2023b).

## 121 3 PROBLEM FORMULATION

123 Our research addresses two complementary challenges in automated ESG analysis, moving from  
124 granular, evidence-based inquiry to holistic, predictive assessment.

### 126 3.1 TASK 1: RELATIONAL ESG QUESTION ANSWERING (QA)

128 Given a corporate document  $D$  and a specific ESG-related question  $q$ , our first task is to find a  
129 factually accurate and contextually complete answer. The question  $q$  is assumed to be mappable to  
130 a primary topic node  $v_q$  within a pre-defined ESG Knowledge Graph  $G = (V, E)$ . Unlike standard  
131 extractive QA, the desired output is a tuple  $(a, \mathcal{C}, J)$ , where  $a$  is a synthesized answer,  $\mathcal{C}$  is a set of  
132 supporting evidence quotes from  $D$ , and  $J$  is a justification narrative. This task tests the framework’s  
133 ability to perform deep, evidence-grounded reasoning.

### 135 3.2 TASK 2: MULTIMODAL ESG SCORE PREDICTION

136 This more challenging task aims to predict a company’s overall ESG score directly from its full,  
137 multimodal annual report. Given a multimodal document  $D_{multi}$ , which comprises text  $D_{text}$ , a set  
138 of tables  $\mathcal{T}$ , and a set of images  $\mathcal{I}$ , the task is to learn a function  $f : D_{multi} \rightarrow \mathbb{R}$  that predicts a  
139 single, continuous ESG score  $\hat{y}$ . The objective is to minimize a regression loss function, such as the  
140 Mean Squared Error (MSE), between the predicted score  $\hat{y}$  and the ground-truth score  $y$ :

$$142 \mathcal{L}_{score} = (\hat{y} - y)^2$$

143 This task evaluates the framework’s ultimate capability to perform holistic, multi-modal synthesis  
144 and predictive judgment.

## 146 4 METHODOLOGY

148 We introduce a framework for structured, collective intelligence, designed to surpass monolithic  
149 models in complex reasoning. As shown in Figure 1, it is based on three pillars: (1) The distillation  
150 of vast domain knowledge into a compact expert model; (2) An innovative, graph-guided dialectical  
151 architecture for orchestrating expert models; and (3) A robust mechanism for synthesizing their  
152 findings into a final predictive judgment.

### 154 4.1 THE ESGEEK MODEL: CREATING AN ESG EXPERT

156 The foundation of our framework is the specialized ESGEEK model, created through a two-step  
157 knowledge distillation process. First, to codify domain expertise, we built ESGEXPERT-30K, a new  
158 high-quality dataset with over 30,000 Q&A pairs. This dataset was meticulously crafted by using  
159 an LLM to extract knowledge from 12 seminal ESG standards, followed by multi-stage human and  
160 expert validation (see Appendix B for full details). Subsequently, we fine-tuned a pre-trained SLM  
161 (Qwen-2.5 0.5B) on this rich dataset, a process that distills comprehensive ESG knowledge into our  
final expert model, ESGEEK ( $M_{ESG}$ ), a highly capable and efficient reasoning module.

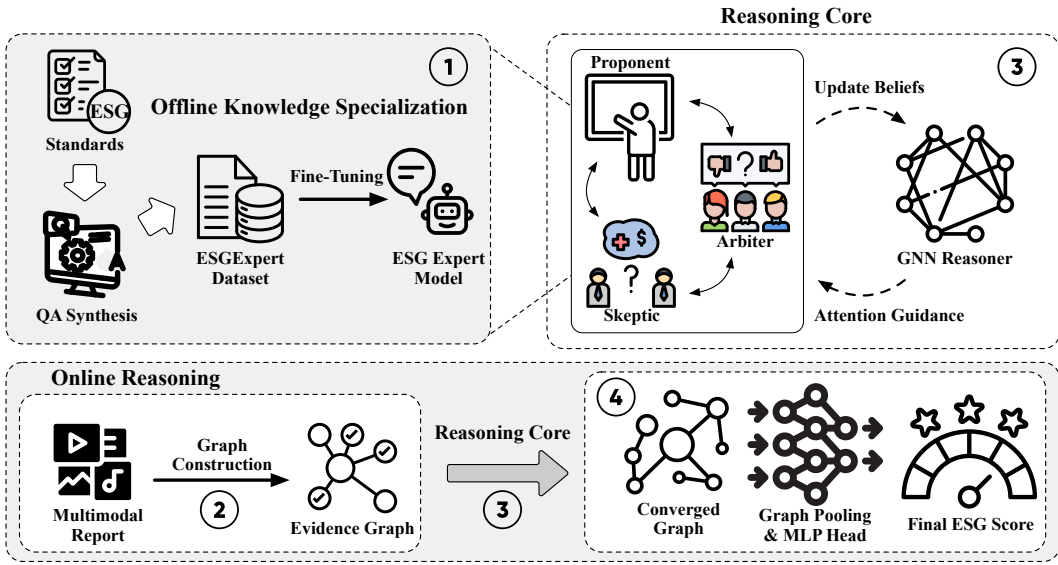


Figure 1: An overview of the DIAL-G<sup>2</sup> framework. The process consists of two components: (1) **Offline Knowledge Specialization**: An expert SLM (ESGEEK) is created by fine-tuning on our knowledge-intensive ESGEXPERT-30K dataset. (2) **Online Reasoning**: For a given multimodal report, an evidence graph is constructed. A GNN Reasoner then guides a “Reasoning Core”, i.e., comprising a Proponent, a Skeptic, and a synthesizing Arbitrator, to debate and iteratively update beliefs on the graph. The GNN’s attention guidance directs the agents’ focus to the most critical information. (3) **Final Prediction**: The converged graph representation is passed through a graph pooling and MLP head to generate the final ESG score.

## 4.2 THE DIAL-G<sup>2</sup> ARCHITECTURE

DIAL-G<sup>2</sup> orchestrates instances of  $M_{ESG}$  within a graph-based reasoning framework.

### 4.2.1 EVIDENCE GRAPH CONSTRUCTION.

A report  $D_{multi}$  is transformed into a graph  $G = (V, E)$ , where each node  $v_i \in V$  represents a page. The edge set  $E$  is crucial for defining the scope of relational reasoning. We construct a hybrid graph where an edge  $(v_i, v_j)$  is created if any of the following conditions are met: (1) **Structural Adjacency**: Pages  $i$  and  $j$  are sequential. (2) **Hierarchical Links**: Both pages belong to the same section as defined by the document’s table of contents. (3) **Semantic Similarity**: The cosine similarity of their page-level text embeddings exceeds a threshold of  $\tau = 0.8$ . This hybrid policy captures both narrative flow and long-distance thematic connections.

### 4.2.2 MULTIMODAL FEATURE ENCODING.

The initial feature vector  $\mathbf{x}_i \in \mathbb{R}^{384}$  for each page-node is derived from its multimodal content. For textual content  $T_i$ , we use a pre-trained sentence-transformer, ‘all-MiniLM-L6-v2’, to generate an embedding  $\mathcal{E}_{\text{text}}(T_i)$ . For tables  $\mathcal{T}_i$  and images  $\mathcal{I}_i$ , we first use ‘Florence-2-large’ to generate descriptive captions, which are then encoded by the same text encoder to form  $\mathcal{E}_{\text{modal}}(\mathcal{T}_i, \mathcal{I}_i)$ . All synthesized beliefs from the debate are also encoded using  $\mathcal{E}_{\text{text}}$ . The final feature vector is a learnable fusion via a two-layer MLP with GELU activation:

$$\mathbf{z}_i = [\mathcal{E}_{\text{text}}(T_i) \parallel \mathcal{E}_{\text{modal}}(\mathcal{T}_i, \mathcal{I}_i)]$$

$$\mathbf{x}_i = \text{MLP}_{\text{fuse}}(\mathbf{z}_i)$$

To illustrate this process, Figure 2 shows a sample report page and its parsed representation.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

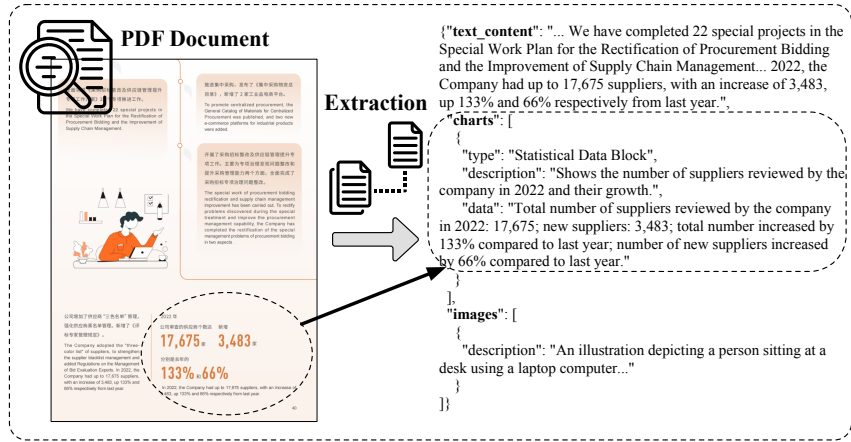


Figure 2: An example of multimodal content processing for a single document page. The raw page contains a mix of text, statistical data, and illustrative graphics. Our system parses this into a structured format, extracting the textual content, identifying and interpreting the chart’s data (e.g., 17,675 suppliers), and generating a description for the image. This rich, structured information serves as the initial feature representation for a node in the evidence graph.

#### 4.2.3 THE GRAPH-GUIDED DEBATE (GGD) WORKFLOW.

This neuro-symbolic loop refines a shared understanding of the document. The state of debate for node  $v_i$  at iteration  $k$  is stored in an Evidence Log,  $\mathcal{E}_{v_i}^{(k)}$ . Full prompt templates are provided in Appendix C.

**Belief Initialization ( $k = 0$ ):** The Proponent ( $A_{pro}$ ) and Skeptic ( $A_{skep}$ ) agents generate initial conflicting beliefs ( $s_{pro,v_i}^{(0)}, s_{skep,v_i}^{(0)}$ ) for each node  $v_i$ . The initial hidden state  $\mathbf{h}_i^{(0)}$  is a learnable fusion of these beliefs.

**Iterative Refinement Loop (for  $k = 0, \dots, K - 1$ ):** The loop consists of three main steps: GNN-driven discovery, dialectical exchange, and arbitrated synthesis.

1. *Attentive Relational Inference*: We use a GAT with Dropout and DropEdge (Zhang et al., 2024b) for robust training. It processes the graph state  $H^{(k)}$  to produce updated intermediate states  $\mathbf{H}'^{(k)}$  and attention weights  $\{\alpha_{ij}^{(k)}\}$ , identifying relational hotspots, that is, top-K attention-weighted edges.
2. *Dialectical Exchange on Hotspots*: For each hotspot edge  $(v_i, v_j)$ , the Proponent and Skeptic are prompted to generate arguments ( $s'_{pro}, s'_{skep}$ ) that either support or challenge the connection between the content on these pages.
3. *Arbitrated Synthesis and Inquiry*: This is the core of the reasoning process. The Arbiter agent ( $A_{arb}$ ) receives the arguments from both sides. Its task is two-fold:
  - **Synthesize Belief**: It generates a synthesized narrative  $s_{synth}$ , which integrates the valid points from both agents, highlights unresolved conflicts, and forms a nuanced judgment.
  - **Generate Inquiry**: Optionally, it can generate a targeted question  $q_{follow-up}$  to probe the weaker argument, preparing for a potential deeper debate in the next iteration.

The embedding of the synthesized belief,  $\mathcal{E}(s_{synth})$ , becomes the new linguistic insight for the node, denoted as  $\mathbf{h}_{rebuttal}$ .

4. *Gated Belief Update*: A gated mechanism fuses the GNN’s propagated belief  $\mathbf{h}_i^{(k)}$  with the new linguistic evidence  $\mathbf{h}_{rebuttal,i}^{(k)}$  from the Arbiter’s synthesis, ensuring a stable, iterative refinement of the node states to  $\mathbf{h}_i^{(k+1)}$ . The update is formalized as:

$$g_i^{(k)} = \sigma(W_g[\mathbf{h}_i^{(k)} \parallel \mathbf{h}_{rebuttal,i}^{(k)}] + b_g)$$

$$\mathbf{h}_i^{(k+1)} = (1 - g_i^{(k)}) \odot \mathbf{h}'_i^{(k)} + g_i^{(k)} \odot \text{MLP}_{map}(\mathbf{h}_{\text{rebuttal},i}^{(k)})$$

where  $g_i^{(k)}$  is the gate vector,  $\sigma$  is the sigmoid function, and  $\text{MLP}_{map}$  is a mapping layer. This allows the model to learn how much of the new insight should influence the existing belief.

### 4.3 HOLISTIC JUDGMENT VIA GRAPH-LEVEL SYNTHESIS AND PREDICTION

The GGD workflow culminates in a set of final, converged node embeddings  $H^{(K)} = \{\mathbf{h}_{v_i}^{(K)}\}_{i=1}^N$ . Each vector  $\mathbf{h}_{v_i}^{(K)}$  is no longer a representation of an isolated page, but a rich, context-aware belief state, refined by both linguistic debate and structural message passing. The final challenge is to synthesize these distributed, node-level beliefs into a single, coherent, and accurate judgment—the ESG score. This process unfolds in two steps: graph-level representation and score regression.

A simple aggregation mechanism, such as global mean pooling ( $\mathbf{h}_G = \frac{1}{N} \sum_i \mathbf{h}_{v_i}^{(K)}$ ), would implicitly assume that every page of the report is equally important to the final ESG score. This is a flawed premise, as corporate reports invariably contain sections of varying materiality.

To address this, we employ a more sophisticated, attention-based graph pooling mechanism (Veličković et al., 2018). This allows the model to learn a materiality weighting over all nodes (pages) in the document graph, deciding which parts of the report are most indicative of the company’s true ESG performance. The final graph-level representation  $\mathbf{h}_G$  is computed as a weighted sum of the node embeddings:

$$\beta_i = \text{softmax}_i(\text{MLP}_{\text{pool}}(\mathbf{h}_{v_i}^{(K)}))$$

$$\mathbf{h}_G = \sum_{i=1}^N \beta_i \cdot \mathbf{h}_{v_i}^{(K)}$$

where  $\text{MLP}_{\text{pool}}$  is a small, learnable network that outputs a single scalar attention score for each node. The softmax function normalizes these scores into a probability distribution  $\{\beta_i\}$ , where  $\beta_i$  can be interpreted as the learned materiality or importance of page  $i$  to the final prediction. This mechanism offers a critical layer of interpretability; by inspecting the nodes with the highest  $\beta_i$  values, we can understand which parts of the report most influenced the model’s final judgment.

The holistic graph embedding  $\mathbf{h}_G$  encapsulates the distilled wisdom of the entire report. This vector is then passed through a final prediction head,  $\text{MLP}_{\text{score}}$ , to regress the continuous ESG score  $\hat{y}$ . This head consists of two feed-forward layers with a ReLU activation:

$$\hat{y} = \text{MLP}_{\text{score}}(\mathbf{h}_G)$$

It is crucial to understand the training dynamics. While the LLM-driven debate is a **non-differentiable** process, it is tightly coupled with the end-to-end trainable graph pipeline. The final prediction loss,  $\mathcal{L}_{\text{score}} = (\hat{y} - y)^2$ , backpropagates through the entire graph pipeline, creating an **indirect learning signal**. The GNN learns to compute attention weights ( $\alpha_{ij}$ ) that are most effective at selecting “hotspot” edges, which, when presented to the non-differentiable debate module, are likely to produce synthesized beliefs that maximally reduce the final prediction error. This creates a self-optimizing system where relational reasoning and linguistic analysis synergize towards a single objective, without requiring gradients to flow through the LLM itself. The entire training process integrates these components into an end-to-end pipeline. The data flow for a single report involves: **(1) Initialization**, where the report is converted into a graph  $G$  with initial beliefs  $H^{(0)}$ ; **(2) Iterative Refinement**, where for  $K$  steps, a GNN identifies salient edges to guide a dialectical debate whose synthesized outcomes update the node states  $H^{(k+1)}$  via a gated mechanism; and **(3) Prediction**, where the final states  $H^{(K)}$  are pooled to predict the score  $\hat{y}$ . The model’s differentiable parameters are optimized by backpropagating the regression loss  $\mathcal{L}_{\text{score}} = (\hat{y} - y)^2$ . The complete procedural breakdown is shown as algorithms process in the Appendix.

## 5 EXPERIMENTS

We conduct a three-stage experimental evaluation to validate our claims. Full implementation details, model configurations, and hyperparameters for all experiments are available in Appendix E.

Table 1: QA performance on the ESGEXPERT-30K test set. Our fine-tuned ESGEEK models significantly outperform their base models.

Model	Params	Accuracy (%)
Qwen-2.5-0.5B (Qwen et al., 2025)	0.5B	54.58
<b>ESGEEK-0.5B (Ours)</b>	0.5B	<b>76.42</b>
Qwen-2.5-1.5B	1.5B	63.91
DeepSeek-R1-Distill-Qwen-1.5B		31.34
<b>ESGEEK-1.5B (Ours)</b>	1.5B	<b>71.21</b>
Qwen-2.5-3B	3B	58.89
<b>ESGEEK-3B (Ours)</b>	3B	<b>74.52</b>
Qwen-2.5-7B	7B	64.96
DeepSeek-R1-Distill-Qwen-7B		50.18
<b>ESGEEK-7B (Ours)</b>	7B	<b>81.78</b>
GPT-4o-mini (OpenAI et al., 2024a)	8B	62.68
Qwen-2.5-14B	14B	61.52
Qwen2.5-Max	325B	64.44
DeepSeek-R1 (DeepSeek-AI et al., 2025)	671B	66.29
GPT-4o	200B	63.64
o3 (OpenAI et al., 2024b)	300B	72.54

## 5.1 EXPERIMENT 1: VALIDATING THE ESGEEK

We benchmark the stand-alone QA capability of our domain-specialized ESGEEK on the held-out ESGEXPERT-30K test split. The task is formulated as multiple-choice question answering, where "Accuracy (%)" is the percentage of correctly answered questions.

Table 1 shows that our 0.5B specialist already surpasses GPT-4o-mini (8B) by 13.74 pp, while the 7B variant attains 81.78%—a +9.2 pp gain over the 300B o3 flagship and a +16.8 pp jump relative to its own base checkpoint, demonstrating that rigorous domain adaptation can distill expert competence into highly parameter-efficient SLMs.

## 5.2 EXPERIMENT 2: END-TO-END SCORE PREDICTION

We evaluate the full DIAL-G<sup>2</sup> framework on our ESGREPORT-RATING-50K dataset. This large-scale dataset, containing over 50,000 report-rating pairs, is designed for diversity and robustness. Labels were harmonized across multiple rating agencies by first Z-score normalizing ratings within each agency, and then mapping them to a unified 0-100 scale. Unlike benchmarks that focus solely on large-cap firms, our dataset includes nearly 20% of reports from Small-to-Medium Enterprises (SMEs). We report the average Pearson correlation among a panel of human experts as a reference for "human-level" performance (see Appendix B for panel protocol details).

We conducted comprehensive experiments for different baselines. The full DIAL-G<sup>2</sup> model achieves the best performance and efficiency. The significant improvement over the DIAL-G<sup>2</sup> (w/o Debate) variant underscores the critical value added by the GGD mechanism. Key results are shown in Table 2.

An intriguing finding is that our model’s average performance ( $\rho = 0.83$ ) demonstrates a higher level of consistency than the average agreement observed among human experts ( $\rho = 0.65$ ). We do not interpret this as "super-human" intuition. Instead, we posit that this highlights the model’s capability to learn and systematically apply the complex, often divergent methodologies of multiple rating agencies. This ability to achieve objective reproducibility is a key advantage for scalable and auditable assessments.

To provide an intuitive demonstration of component synergy, we conduct a detailed qualitative case study, illustrating how the GNN-guided debate process uncovers deep-seated risks in the report (see Appendix C).

Table 2: ESG score prediction results on ESGREPORT-RATING-50K.

Model	MSE ( $10^{-2}$ ↓)	Pearson $\rho$ (% ↑)	Runtime (h↓)
<i>Baselines</i>			
GPT-4o-mini (Multimodal Prompt)	$43.5 \pm 7.29$	$0.21 \pm 0.19$	$120.5 \pm 10.3$
Simple Multimodal Fusion	$10.5 \pm 1.79$	$0.47 \pm 0.13$	$14.2 \pm 3.1$
Unstructured Debate (MAS)	$4.95 \pm 1.10$	$0.62 \pm 0.08$	$64.5 \pm 6.8$
Graph-RAG (G-Retriever-like)	$3.52 \pm 0.77$	$0.70 \pm 0.06$	$21.3 \pm 2.5$
<b>DIAL-G<sup>2</sup> (w/o Debate)</b>	<b><math>2.71 \pm 0.68</math></b>	<b><math>0.72 \pm 0.04</math></b>	<b><math>18.7 \pm 1.4</math></b>
<i>Our Method</i>			
Human Expert Average	N/A	$0.65 \pm 0.08$	N/A
<b>DIAL-G<sup>2</sup> (0.5 B)</b>	<b><math>1.35 \pm 0.34</math></b>	<b><math>0.78 \pm 0.03</math></b>	<b><math>9.2 \pm 0.6</math></b>
<b>DIAL-G<sup>2</sup> (1.5 B)</b>	<b><math>1.62 \pm 0.28</math></b>	<b><math>0.80 \pm 0.05</math></b>	<b><math>15.8 \pm 1.1</math></b>
<b>DIAL-G<sup>2</sup> (3 B)</b>	<b><math>1.41 \pm 0.27</math></b>	<b><math>0.82 \pm 0.02</math></b>	<b><math>28.4 \pm 2.2</math></b>
<b>DIAL-G<sup>2</sup> (7 B)</b>	<b><math>1.33 \pm 0.25</math></b>	<b><math>0.83 \pm 0.02</math></b>	<b><math>57.9 \pm 4.5</math></b>

### 5.3 EXPERIMENT 3: ROBUSTNESS EXPLORATION

To further explore the robustness of our framework, we conducted a decomposition analysis of its performance across different industry sectors. The radar chart in Figure 3 visually compares the MSE performance of each model across nine major industries and an aggregated “overall statistics” dimension. The specific composition of these industry categories is detailed in Appendix E. The chart demonstrates that the complete DIAL-G<sup>2</sup> framework consistently outperforms baseline models in every industry sector, demonstrating its robust generalization capabilities.

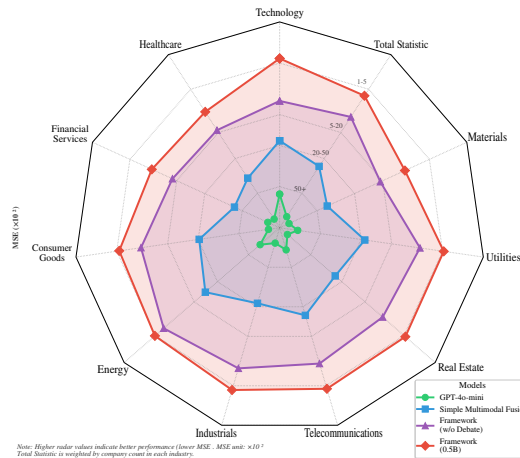


Figure 3: ESG Model MSE Performance Breakdown Across Industries. The radar chart compares the MSE of our full Framework (0.5B) against key baselines across nine industries and a weighted total statistic. The radial axis represents performance, where higher values indicate lower MSE. The chart clearly illustrates the consistent and significant performance of the DIAL-G<sup>2</sup> framework in all sectors.

### 5.4 ABLATION STUDIES

To dissect the sources of DIAL-G<sup>2</sup>'s performance, we conducted ablation studies on the score prediction task. Table 3 confirms that each component is critical: using the specialized ESGEEK outperforms a generic SLM, removing graph guidance roughly doubles the error across scales, and eliminating the Skeptic agent yields a consistent, though smaller degradation, evidencing the value of dialectical balance.

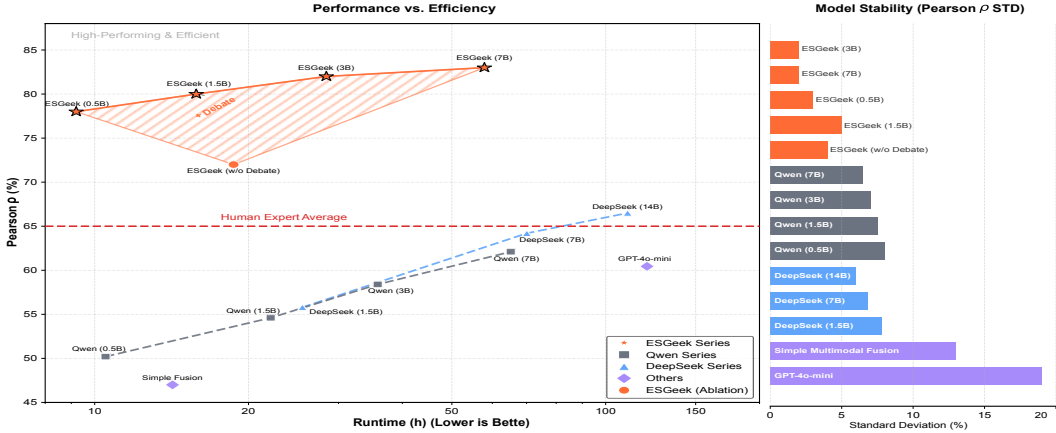


Figure 4: **Left Panel:** The scatter plot visualizes prediction performance (Pearson  $\rho$ , higher is better) against computational cost (runtime, lower is better). Our DIAL-G<sup>2</sup> series (starred) establishes a new state-of-the-art Pareto frontier, while the shaded area quantifies the performance gain from the Debate mechanism. **Right Panel:** The bar chart displays model stability, measured by the standard deviation of the Pearson correlation. Shorter bars indicate higher stability, where our DIAL-G<sup>2</sup> models show consistently strong performance.

Table 3: Ablation studies on the score prediction task.

Model Size	Framework Variant	MSE ( $\times 10^{-2}$ )
0.5 B	<b>DIAL-G<sup>2</sup> (Full Model)</b>	$1.35 \pm 0.34$
	w/o ESGEEK (generic SLM)	$1.89 \pm 0.48$
	w/o GNN Guidance (random debate)	$2.70 \pm 0.68$
	w/o Skeptic Agent (analyst-only)	$1.55 \pm 0.39$
1.5 B	<b>DIAL-G<sup>2</sup> (Full Model)</b>	$1.62 \pm 0.28$
	w/o ESGEEK (generic SLM)	$2.19 \pm 0.38$
	w/o GNN Guidance (random debate)	$3.05 \pm 0.53$
	w/o Skeptic Agent (analyst-only)	$1.85 \pm 0.32$
3 B	<b>DIAL-G<sup>2</sup> (Full Model)</b>	$1.41 \pm 0.27$
	w/o ESGEEK (generic SLM)	$1.83 \pm 0.35$
	w/o GNN Guidance (random debate)	$2.71 \pm 0.52$
	w/o Skeptic Agent (analyst-only)	$1.59 \pm 0.31$
7 B	<b>DIAL-G<sup>2</sup> (Full Model)</b>	$1.33 \pm 0.25$
	w/o ESGEEK (generic SLM)	$1.66 \pm 0.31$
	w/o GNN Guidance (random debate)	$2.48 \pm 0.47$
	w/o Skeptic Agent (analyst-only)	$1.49 \pm 0.28$

## 6 CONCLUSION

This work challenged the “bigger is better” paradigm in the context of complex ESG analysis. We provide a new **ESGEXPERT-30K** dataset and fine-tune the **ESGEEK** model on it. We then proposed **DIAL-G<sup>2</sup>**, a novel framework where expert agents based on ESGEEK are orchestrated by a GNN in an iterative, graph-guided dialectical debate. The experiments on our new large-scale, multimodal **ESGREPORT-RATING-50K** benchmark show that our approach achieves state-of-the-art performance in ESG score prediction, reaching a level of consistency comparable to human expert panels.

**Limitations and Future Works.** Despite its strong performance, our framework has limitations that open promising avenues for future research. Firstly, the knowledge codified in the ESGEEK model is static. As ESG standards and regulations evolve, a key challenge is to develop methods for dynamic knowledge. For instance, continuous learning or more dynamic knowledge graph integration.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

---

## REPRODUCIBILITY STATEMENT

We have made every effort to ensure our work is fully reproducible. To this end, we provide a comprehensive set of resources and detailed documentation.

**Datasets and Code.** We provide our two newly created datasets, ESGEXPERT-30K and ESGREPORT-RATING-50K, along with the complete source code in the supplementary materials. The codebase includes the data processing pipeline, the fine-tuning script for creating the ESGEEK model, the full implementation of the DIAL-G<sup>2</sup> framework, and all scripts required to replicate our experimental results.

**Methodological and Experimental Details.** To further support reproducibility, we have documented all key details within the paper.

- The detailed methodologies for constructing our datasets are described in Appendix B.
- All experimental settings, model configurations, and hyperparameters are specified in Appendix A.
- The core architecture of our DIAL-G<sup>2</sup> framework is detailed in Section 4.

**Public Release Commitment.** Upon acceptance of this paper, we will make the complete codebase and all associated datasets publicly available on a permanent platform (e.g., GitHub) under a permissive license to facilitate future research in the community.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

---

## REFERENCES

- Alexandria Barghi. QRAG: Using Learnable Graph Queries for Retrieval Augmented Generation. In *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*, pp. 00561–00568, Las Vegas, NV, USA, January 2025. IEEE. ISBN 9798331507695. doi: 10.1109/CCWC62904.2025.10903812. URL <https://ieeexplore.ieee.org/document/10903812/>.
- Florian Berg, Julian F Kölbl, and Roberto Rigobon. Aggregate Confusion: The Divergence of ESG Ratings. *Review of Finance*, 26(6):1315–1344, November 2022. ISSN 1572-3097, 1573-692X. doi: 10.1093/rof/rfac033. URL <https://academic.oup.com/rof/article/26/6/1315/6590670>.
- Mattia Birti, Francesco Osborne, and Andrea Maurino. Optimizing Large Language Models for ESG Activity Detection in Financial Texts, February 2025. URL <http://arxiv.org/abs/2502.21112>. arXiv:2502.21112 [cs].
- Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F. Karlsson, Jie Fu, and Yemin Shi. AutoAgents: A Framework for Automatic Agent Generation, 2023a. URL <https://arxiv.org/abs/2309.17288>. Version Number: 3.
- Jie Chen, Tengfei Ma, and Cao Xiao. FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling, 2018. URL <https://arxiv.org/abs/1801.10247>. Version Number: 1.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors, 2023b. URL <https://arxiv.org/abs/2308.10848>. Version Number: 3.
- Zhongfei Chen and Guanxia Xie. ESG disclosure and financial performance: Moderating role of ESG investors. *International Review of Financial Analysis*, 83:102291, October 2022. ISSN 10575219. doi: 10.1016/j.irfa.2022.102291. URL <https://linkinghub.elsevier.com/retrieve/pii/S1057521922002472>.
- Dane M. Christensen, George Serafeim, and Anywhere Sikochi. Why is Corporate Virtue in the Eye of The Beholder? The Case of ESG Ratings. *The Accounting Review*, 97(1):147–175, January 2022. ISSN 1558-7967, 0001-4826. doi: 10.2308/TAR-2019-0506. URL <https://publications.aaahq.org/accounting-review/article/97/1/147/4390/Why-is-Corporate-Virtue-in-the-Eye-of-The-Beholder>.
- Bradford Cornell. ESG preferences, risk and return. *European Financial Management*, 27(1):12–19, January 2021. ISSN 1354-7798, 1468-036X. doi: 10.1111/eufm.12295. URL <https://onlinelibrary.wiley.com/doi/10.1111/eufm.12295>.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu,

---

594 Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xi-  
595 aosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia  
596 Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng  
597 Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong  
598 Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong,  
599 Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou,  
600 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying  
601 Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda  
602 Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia  
603 Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu,  
604 Zhongyu Zhang, and Zhen Zhang. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs  
605 via Reinforcement Learning, January 2025. URL <http://arxiv.org/abs/2501.12948>.  
606 arXiv:2501.12948 [cs].

607 Elroy Dimson, Paul Marsh, and Mike Staunton. Divergent ESG Ratings. *The Journal of Portfolio*  
608 *Management*, 47(1):75–87, October 2020. ISSN 0095-4918, 2168-8656. doi: 10.3905/jpm.  
609 2020.1.175. URL [http://pm-research.com/lookup/doi/10.3905/jpm.2020.](http://pm-research.com/lookup/doi/10.3905/jpm.2020.1.175)  
610 [1.175](http://pm-research.com/lookup/doi/10.3905/jpm.2020.1.175).

611 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,  
612 Dasha Metropolitanaky, Robert Osazuwa Ness, and Jonathan Larson. From Local to Global:  
613 A Graph RAG Approach to Query-Focused Summarization, February 2025. URL [http://](http://arxiv.org/abs/2404.16130)  
614 [arxiv.org/abs/2404.16130](http://arxiv.org/abs/2404.16130). arXiv:2404.16130 [cs].

615 Stuart L. Gillan, Andrew Koch, and Laura T. Starks. Firms and social responsibility: A review of  
616 ESG and CSR research in corporate finance. *Journal of Corporate Finance*, 66:101889, February  
617 2021. ISSN 09291199. doi: 10.1016/j.jcorpfin.2021.101889. URL [https://linkinghub.](https://linkinghub.elsevier.com/retrieve/pii/S0929119921000092)  
618 [elsevier.com/retrieve/pii/S0929119921000092](https://linkinghub.elsevier.com/retrieve/pii/S0929119921000092).

619 Global Reporting Initiative. GRI Standards, 2021. URL [https://](https://www.globalreporting.org/how-to-use-the-gri-standards/gri-standards-english-language/)  
620 [www.globalreporting.org/how-to-use-the-gri-standards/](https://www.globalreporting.org/how-to-use-the-gri-standards/gri-standards-english-language/)  
621 [gri-standards-english-language/](https://www.globalreporting.org/how-to-use-the-gri-standards/gri-standards-english-language/).

622  
623 Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. LightRAG: Simple and Fast  
624 Retrieval-Augmented Generation, April 2025. URL [http://arxiv.org/abs/2410.](http://arxiv.org/abs/2410.05779)  
625 [05779](http://arxiv.org/abs/2410.05779). arXiv:2410.05779 [cs].

626  
627 Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, and Zhaozhuo Xu. LLM Multi-Agent Sys-  
628 tems: Challenges and Open Problems, May 2025. URL [http://arxiv.org/abs/2402.](http://arxiv.org/abs/2402.03578)  
629 [03578](http://arxiv.org/abs/2402.03578). arXiv:2402.03578 [cs].

630  
631 Chaoyue He, Xin Zhou, Yi Wu, Xinjia Yu, Yan Zhang, Lei Zhang, Di Wang, Shengfei Lyu, Hong  
632 Xu, Xiaoqiao Wang, Wei Liu, and Chunyan Miao. ESGenius: Benchmarking LLMs on En-  
633 vironmental, Social, and Governance (ESG) and Sustainability Knowledge, June 2025. URL  
634 <http://arxiv.org/abs/2506.01646>. arXiv:2506.01646 [cs].

635  
636 Xiaoxuan Hou, Jiayi Yuan, Joel Z. Leibo, and Natasha Jaques. InvestESG: A multi-agent reinforce-  
637 ment learning benchmark for studying climate investment as a social dilemma, February 2025.  
URL <http://arxiv.org/abs/2411.09856>. arXiv:2411.09856 [cs].

638  
639 Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Net-  
640 works, 2016. URL <https://arxiv.org/abs/1609.02907>. Version Number: 4.

641  
642 Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-  
643 supervised Learning: Generative or Contrastive. *IEEE Transactions on Knowledge and Data*  
644 *Engineering*, pp. 1–1, 2021. ISSN 1041-4347, 1558-2191, 2326-3865. doi: 10.1109/TKDE.  
2021.3090866. URL <https://ieeexplore.ieee.org/document/9462394/>.

645  
646 Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. FinBERT: A Pre-trained  
647 Financial Language Representation Model for Financial Text Mining. In *Proceedings of the*  
*Twenty-Ninth International Joint Conference on Artificial Intelligence*, pp. 4513–4519, Yoko-  
hama, Japan, July 2020. International Joint Conferences on Artificial Intelligence Organization.

---

648 ISBN 978-0-9992411-6-5. doi: 10.24963/ijcai.2020/622. URL <https://www.ijcai.org/proceedings/2020/622>.

649

650

651 Alexandra Luccioni, Emily Baylor, and Nicolas Duchene. Analyzing Sustainability Reports Using Natural Language Processing, November 2020. URL <http://arxiv.org/abs/2011.08073>. arXiv:2011.08073 [cs].

652

653

654 Lokesh Mishra, Cesar Berrospi, Kasper Dinkla, Diego Antognini, Francesco Fusco, Benedikt Bothur, Maksym Lysak, Nikolaos Livathinos, Ahmed Nassar, Panagiotis Vagenas, Lucas Morin, Christoph Auer, Michele Dolfi, and Peter Staar. ESG Accountability Made Easy: DocQA at Your Service. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(21):23814–23816, March 2024. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v38i21.30574. URL <http://arxiv.org/abs/2311.18481>. arXiv:2311.18481 [cs].

655

656

657

658

659

660

661 Tim Nugent, Nicole Stelea, and Jochen L. Leidner. Detecting Environmental, Social and Governance (ESG) Topics Using Domain-Specific Language Models and Data Augmentation. In Troels Andreasen, Guy De Tré, Janusz Kacprzyk, Henrik Legind Larsen, Gloria Bordogna, and Sławomir Zadrozny (eds.), *Flexible Query Answering Systems*, volume 12871, pp. 157–169. Springer International Publishing, Cham, 2021. ISBN 978-3-030-86966-3 978-3-030-86967-0. doi: 10.1007/978-3-030-86967-0\_12. URL [https://link.springer.com/10.1007/978-3-030-86967-0\\_12](https://link.springer.com/10.1007/978-3-030-86967-0_12). Series Title: Lecture Notes in Computer Science.

662

663

664

665

666

667

668 OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, A. J. Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Gertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou,

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702 Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier,  
703 Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak,  
704 Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray,  
705 Mark Hudnall, Marvin Zhang, Marwan Aljubeih, Mateusz Litwin, Matthew Zeng, Max John-  
706 son, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao  
707 Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael  
708 Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de  
709 Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavar-  
710 ian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Stau-  
711 dacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick  
712 Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick,  
713 Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins,  
714 Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng  
715 Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil  
716 Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim,  
717 Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar  
718 Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith,  
719 Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen,  
720 Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam  
721 Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean  
722 Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino  
723 Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve  
724 Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogi-  
725 neni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry,  
726 Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor  
727 Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peter-  
728 son, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad  
729 Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech  
730 Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen  
731 Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. GPT-4o System Card, October 2024a. URL  
<http://arxiv.org/abs/2410.21276>. arXiv:2410.21276 [cs].

732 OpenAI, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden  
733 Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko,  
734 Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally  
735 Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich,  
736 Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghor-  
737 bani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao,  
738 Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary  
739 Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang,  
740 Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel  
741 Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson,  
742 Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Eliz-  
743 abeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang,  
744 Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von  
745 Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao,  
746 Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hes-  
747 sam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian  
748 O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever,  
749 Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng,  
750 Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñero Candela, Joe Palermo, Joel Parish,  
751 Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan  
752 Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl  
753 Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu,  
754 Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam  
755 Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kon-  
draciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen,  
Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet  
Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael

- 756 Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles  
757 Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil  
758 Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg  
759 Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov,  
760 Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar  
761 Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan  
762 Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agar-  
763 wal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu,  
764 Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph  
765 Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Tay-  
766 lor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson,  
767 Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna  
768 Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi  
769 Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen,  
770 Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li.  
771 OpenAI o1 System Card, December 2024b. URL <http://arxiv.org/abs/2412.16720>.  
arXiv:2412.16720 [cs].
- 772 Jili Qian, Zhengyu Jin, Quan Zhang, Guoqing Cai, and Beichang Liu. A Liver Cancer Question-  
773 Answering System Based on Next-Generation Intelligence and the Large Model Med-PaLM 2. *In-*  
774 *ternational Journal of Computer Science and Information Technology*, 2(1):28–35, March 2024.  
775 ISSN 3005-7140, 3005-9682. doi: 10.62051/ijcsit.v2n1.04. URL [https://wepub.org/](https://wepub.org/index.php/IJCSIT/article/view/662)  
776 [index.php/IJCSIT/article/view/662](https://wepub.org/index.php/IJCSIT/article/view/662).
- 777 Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
778 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
779 Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin  
780 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,  
781 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,  
782 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 Technical Report,  
783 January 2025. URL <http://arxiv.org/abs/2412.15115>. arXiv:2412.15115 [cs].
- 784 Frank Schiemann and Raphael Tietmeyer. ESG Controversies, ESG Disclosure and Analyst Fore-  
785 cast Accuracy. *International Review of Financial Analysis*, 84:102373, November 2022. ISSN  
786 10575219. doi: 10.1016/j.irfa.2022.102373. URL [https://linkinghub.elsevier.](https://linkinghub.elsevier.com/retrieve/pii/S1057521922003234)  
787 [com/retrieve/pii/S1057521922003234](https://linkinghub.elsevier.com/retrieve/pii/S1057521922003234).
- 788 George Serafeim and Aaron Yoon. Stock price reactions to ESG news: the role of ESG ratings and  
789 disagreement. *Review of Accounting Studies*, 28(3):1500–1530, September 2023. ISSN 1380-  
790 6653, 1573-7136. doi: 10.1007/s11142-022-09675-3. URL [https://link.springer.](https://link.springer.com/10.1007/s11142-022-09675-3)  
791 [com/10.1007/s11142-022-09675-3](https://link.springer.com/10.1007/s11142-022-09675-3).
- 792 Amila Silva, Ling Luo, Shanika Karunasekera, and Christopher Leckie. Embracing Domain Differ-  
793 ences in Fake News: Cross-domain Fake News Detection using Multi-modal Data. *Proceedings*  
794 *of the AAAI Conference on Artificial Intelligence*, 35(1):557–565, May 2021. ISSN 2374-3468,  
795 2159-5399. doi: 10.1609/aaai.v35i1.16134. URL [https://ojs.aaai.org/index.php/](https://ojs.aaai.org/index.php/AAAI/article/view/16134)  
796 [AAAI/article/view/16134](https://ojs.aaai.org/index.php/AAAI/article/view/16134).
- 797 Sustainability Accounting Standards Board. SASB Standards, 2018. URL [https://www.sasb.](https://www.sasb.org/)  
798 [org/](https://www.sasb.org/).
- 799 Petar Velic̃kovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua  
800 Bengio. GRAPH ATTENTION NETWORKS. 2018.
- 801 Jhing-Fa Wang, Wen-Yuan Zhang, and Shih-Pang Tseng. An innovative ESGH-RAG module with  
802 ChatGPT-4o for automatic ESG-report generation. *The Journal of Supercomputing*, 81(10):1103,  
803 July 2025a. ISSN 1573-0484. doi: 10.1007/s11227-025-07604-0. URL [https://link.](https://link.springer.com/10.1007/s11227-025-07604-0)  
804 [springer.com/10.1007/s11227-025-07604-0](https://link.springer.com/10.1007/s11227-025-07604-0).
- 805 Xinyu Wang, Jijun Chi, Zhenghan Tai, Tung Sum Thomas Kwok, Muzhi Li, Zhuhong Li, Hailin  
806 He, Yuchen Hua, Peng Lu, Suyuchen Wang, Yihong Wu, Jerry Huang, Jingrui Tian, Fengran Mo,  
807 Yufei Cui, and Ling Zhou. FinSage: A Multi-aspect RAG System for Financial Filings Question  
808  
809

---

810 Answering, June 2025b. URL <http://arxiv.org/abs/2504.14493>. arXiv:2504.14493  
811 [cs].  
812

813 Woei Chyuan Wong, Jonathan A. Batten, Abd Halim Ahmad, Shamsul Bahrain Mohamed-Arshad,  
814 Sabariah Nordin, and Azira Abdul Adzis. Does ESG certification add firm value? *Finance*  
815 *Research Letters*, 39:101593, March 2021. ISSN 15446123. doi: 10.1016/j.frl.2020.101593. URL  
816 <https://linkinghub.elsevier.com/retrieve/pii/S1544612319312735>.

817 Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prab-  
818 hanjan Kambadur, David Rosenberg, and Gideon Mann. BloombergGPT: A Large Language  
819 Model for Finance, December 2023a. URL <http://arxiv.org/abs/2303.17564>.  
820 arXiv:2303.17564 [cs].

821 Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph Neural Networks in Recommender  
822 Systems: A Survey. *ACM Computing Surveys*, 55(5):1–37, May 2023b. ISSN 0360-0300, 1557-  
823 7341. doi: 10.1145/3535101. URL <https://dl.acm.org/doi/10.1145/3535101>.  
824

825 Jintian Zhang, Xin Xu, Ningyu Zhang, Ruibo Liu, Bryan Hooi, and Shumin Deng. Exploring  
826 Collaboration Mechanisms for LLM Agents: A Social Psychology View. In *Proceedings of*  
827 *the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
828 *Papers)*, pp. 14544–14607, Bangkok, Thailand, 2024a. Association for Computational Linguis-  
829 tics. doi: 10.18653/v1/2024.acl-long.782. URL <https://aclanthology.org/2024.acl-long.782>.  
830

831 Zeyu Zhang, Lu Li, Shuyan Wan, Sijie Wang, Zhiyi Wang, Zhiyuan Lu, Dong Hao, and Wanli Li.  
832 DropEdge not Foolproof: Effective Augmentation Method for Signed Graph Neural Networks,  
833 October 2024b. URL <http://arxiv.org/abs/2409.19620>. arXiv:2409.19620 [cs].  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 4: Comprehensive Performance on ESGREPORT-RATING-50K. This table presents a comprehensive comparison of different frameworks and model series. The results highlight that our **DIAL-G<sup>2</sup>** framework achieves substantially better performance (lower MSE, higher Pearson) at a comparable computational cost. Within any given framework, our specialized **ESGEEK** models consistently outperform other models.

Framework / Model Series	Model Size	MSE ( $10^{-2}$ ↓)	Pearson $\rho$ (↑)	Runtime (h ↓)
<b>Baseline 1: Unstructured Dialogue (AutoGen)</b>				
Qwen-2.5	7B	$4.95 \pm 1.10$	$0.62 \pm 0.08$	$64.5 \pm 6.8$
DeepSeek	7B	$4.71 \pm 1.05$	$0.64 \pm 0.08$	$63.1 \pm 6.5$
Llama	7B	$4.82 \pm 1.07$	$0.63 \pm 0.09$	$63.9 \pm 6.6$
Gemma	7B	$5.08 \pm 1.12$	$0.61 \pm 0.09$	$65.2 \pm 7.0$
<b>ESGEEK (Ours)</b>	7B	<b><math>4.40 \pm 0.98</math></b>	<b><math>0.66 \pm 0.07</math></b>	<b><math>61.5 \pm 6.2</math></b>
<b>Baseline 2: Static Graph (LangGraph)</b>				
Qwen-2.5	7B	$4.15 \pm 0.91$	$0.68 \pm 0.07$	$61.8 \pm 6.1$
DeepSeek	7B	$3.89 \pm 0.85$	$0.70 \pm 0.06$	$60.3 \pm 5.8$
Llama	7B	$4.01 \pm 0.88$	$0.69 \pm 0.07$	$61.1 \pm 6.0$
Gemma	7B	$4.28 \pm 0.94$	$0.67 \pm 0.08$	$62.5 \pm 6.3$
<b>ESGEEK (Ours)</b>	7B	<b><math>3.52 \pm 0.77</math></b>	<b><math>0.72 \pm 0.06</math></b>	<b><math>59.1 \pm 5.5</math></b>
<b>Baseline 3: Heuristic Team Generation (AutoAgents)</b>				
Qwen-2.5	7B	$4.35 \pm 0.96$	$0.66 \pm 0.08$	$63.5 \pm 6.4$
DeepSeek	7B	$4.11 \pm 0.90$	$0.68 \pm 0.07$	$62.1 \pm 6.2$
Llama	7B	$4.23 \pm 0.93$	$0.67 \pm 0.08$	$62.8 \pm 6.3$
Gemma	7B	$4.47 \pm 0.98$	$0.65 \pm 0.09$	$64.1 \pm 6.6$
<b>ESGEEK (Ours)</b>	7B	<b><math>3.73 \pm 0.82</math></b>	<b><math>0.70 \pm 0.07</math></b>	<b><math>60.7 \pm 6.0</math></b>
<b>Our Framework: Dynamic Graph-Guided Debate (DIAL-G<sup>2</sup>)</b>				
Qwen-2.5	7B	$1.95 \pm 0.43$	$0.79 \pm 0.04$	$60.5 \pm 5.1$
DeepSeek	7B	$1.80 \pm 0.40$	$0.80 \pm 0.04$	$59.2 \pm 4.9$
Llama	7B	$1.88 \pm 0.42$	$0.79 \pm 0.05$	$60.1 \pm 5.0$
Gemma	7B	$2.01 \pm 0.45$	$0.78 \pm 0.05$	$61.3 \pm 5.3$
GPT-4o-mini*	~8B	$1.65 \pm 0.36$	$0.81 \pm 0.04$	~78.0
<b>ESGEEK (Ours)</b>	<b>0.5B</b>	<b><math>1.35 \pm 0.34</math></b>	<b><math>0.78 \pm 0.03</math></b>	<b><math>9.2 \pm 0.6</math></b>
<b>ESGEEK (Ours)</b>	<b>1.5B</b>	<b><math>1.62 \pm 0.28</math></b>	<b><math>0.80 \pm 0.05</math></b>	<b><math>15.8 \pm 1.1</math></b>
<b>ESGEEK (Ours)</b>	<b>3B</b>	<b><math>1.41 \pm 0.27</math></b>	<b><math>0.82 \pm 0.02</math></b>	<b><math>28.4 \pm 2.2</math></b>
<b>ESGEEK (Ours)</b>	<b>7B</b>	<b><math>1.33 \pm 0.25</math></b>	<b><math>0.83 \pm 0.02</math></b>	<b><math>57.9 \pm 4.5</math></b>
<b>Reference Benchmarks</b>				
Human Expert Average	N/A	N/A	$0.65 \pm 0.00$	N/A
GPT-4o-mini (Zero-Shot)	~8B	$43.5 \pm 7.29$	$0.21 \pm 0.19$	$120.5 \pm 10.3$

\* GPT-4o-mini (about 8B parameters) is reported in the 7B model series section for ease of horizontal comparison.

## A DETAILED RESULTS

Table 4 below adds details to Experiment 5.2, which we conducted reasonably and adequately for different baselines. The performance of different types of multi-intelligentsia frameworks was tested separately based on other models and our ESGEEK.

## B ESGREPORT-RATING-50K DATASET CONSTRUCTION

To support a comprehensive, end-to-end evaluation of the DIAL-G<sup>2</sup> framework, we constructed a large-scale, multimodal dataset of ESG reports and their corresponding ratings, named ESGREPORT-RATING-50K. This appendix details its construction process to ensure the transparency and reproducibility of our research.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

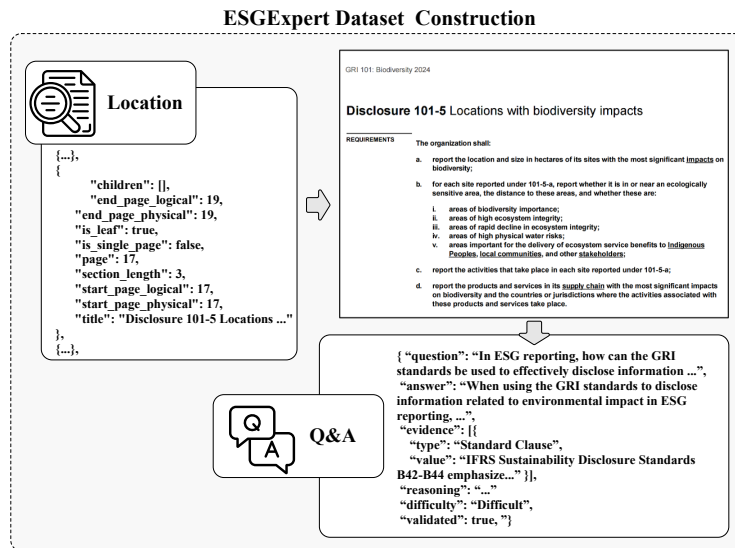


Figure 5: An illustration of our semi-automated dataset construction process for ESGEXPERT-30K. **Left:** A raw clause from a source standard, GRI 101 on biodiversity, detailing specific disclosure requirements. **Right:** The resulting structured data point generated by our pipeline. It includes a complex question, a detailed answer, evidence pointers, and metadata such as difficulty and validation status.

We constructed ESGEXPERT-30K, a high-quality dataset with over 30,000 examples.

- **LLM-based Knowledge Extraction:** We prompt a powerful LLM, GPT-4o (OpenAI et al., 2024a), with the full text of 12 seminal ESG standards, such as GRI (Global Reporting Initiative, 2021), SASB (Sustainability Accounting Standards Board, 2018), to systematically generate a wide array of question-answer pairs covering definitional, procedural, and causal reasoning.
- **Human Verification and Refinement:** Each generated pair is validated for factual accuracy, clarity, and relevance by a team of three human annotators. We calculated the Inter-Annotator Agreement (IAA) on a 10% sample, achieving a strong Cohen’s Kappa coefficient of  $\kappa = 0.87$ , indicating high-quality and consistent annotations.
- **Expert Review:** A final review of the dataset was conducted by a domain expert to ensure its fidelity to the nuances of ESG reporting standards.

Figure 5 provides a visual illustration of this process, showing how a raw clause from an ESG standard is transformed into a structured, validated QA data point that populates our dataset.

## B.1 DATASET OVERVIEW

The primary objective of ESGREPORT-RATING-50K is to provide a unified benchmark for training and evaluating models that can directly predict a company’s ESG performance from its complete, unstructured annual reports (including sustainability reports, integrated reports, etc.). The dataset comprises over 50,000 report-score pairs.

## B.2 DATA COLLECTION AND SOURCING

Our data collection process was designed to ensure broad coverage in terms of geography, industry, and rating methodology.

- **Report Sourcing:** We collected publicly available annual reports from 2018 to 2023. For the international portion, we focused on companies listed in the S&P 500 and FTSE4Good

Index Series. Crucially, to enhance the dataset’s diversity and applicability, we also included over 3,000 Chinese ESG reports from A-share listed companies, sourced from the Shanghai and Shenzhen Stock Exchanges. All reports are in their original PDF format. We prioritized official corporate sustainability or integrated reports. If unavailable, we used the ESG-related sections of their annual financial reports.

- **Rating Score Sourcing:** To ensure robustness and mitigate single-provider bias, we adopted a multi-source strategy for our labels. For international reports, we used the **S&P Global ESG Scores** as the primary ground truth. For the Chinese market reports, we collected corresponding ESG ratings from three leading providers in the region: SynTao Green Finance, China Alliance of Social Value Investment, and Wind. During training, we treat the ratings from different providers for the same report as distinct data points, allowing the model to learn a more generalized understanding of ESG assessment rather than overfitting to a single rating agency’s methodology. This multi-source approach significantly enriches the dataset, making it a more challenging and realistic benchmark for ESG analysis.

Furthermore, a key feature of our dataset is its breadth. Unlike many existing benchmarks that focus predominantly on large-cap companies with mature ESG reporting frameworks, our collection includes a significant number of reports from small-to-mid-cap (SME) companies and firms in emerging markets. This deliberate inclusion enhances the dataset’s diversity, presenting a more realistic and challenging task. It is crucial for developing models with robust generalization capabilities, as reports from these segments often vary significantly in structure, quality, and disclosure patterns.

### B.3 DATA CURATION AND PREPROCESSING

1. **Matching and Alignment:** We developed automated scripts to accurately match the collected PDF reports with their corresponding S&P Global ESG Scores based on company name, ticker symbol, and reporting year.
2. **Content Extraction:** We used the `PyMuPDF` library to parse each PDF report on a page-by-page basis, extracting plain text, as well as identifying table and image areas. For lower-quality scanned PDFs, we employed Optical Character Recognition (OCR) technology.
3. **Quality Filtering:** To ensure dataset quality, we removed samples that met the following criteria: reports with an OCR accuracy below 90%, reports for which a corresponding ESG score could not be found, and reports with fewer than 20 pages (typically containing insufficient information). After filtering, we obtained a final set of 50,372 high-quality report-score pairs.

### B.4 DATASET STATISTICS

The final dataset was split into training, validation, and test sets with an 80:10:10 ratio. Key statistics are summarized in Table 5. The distribution of scores is approximately a left-skewed normal distribution, reflecting the tendency of large public companies to invest in achieving favorable ESG ratings.

Table 5: Key statistics of the ESGREPORT-RATING-50K dataset.

Statistic	Value
Total Samples	50,372
Training Set Samples	40,298
Validation Set Samples	5,037
Test Set Samples	5,037
Covered Years	2018 – 2023
Score Range	[12, 94]
Score Mean	65.7
Score Std. Dev.	18.2

---

## 1026 B.5 ETHICAL CONSIDERATIONS

1027

1028 All content in this dataset is sourced from publicly available information and does not involve any  
1029 private or sensitive data. Throughout the process, we have strictly adhered to the terms of use of the  
1030 source websites and will use the dataset for academic research purposes only.

1031

## 1032 B.6 HUMAN-LEVEL PERFORMANCE BENCHMARK

1033

1034 To establish a robust baseline for "human-level" performance, we computed the inter-expert agree-  
1035 ment on a representative subset of 500 reports from our test set. Our human expert panel was  
1036 designed to reflect the diversity of the professional ESG rating landscape and included both institu-  
1037 tional and individual evaluations.

### 1038 B.6.1 INSTITUTIONAL RATER AGREEMENT

1039

1040 We collected concurrent ESG ratings for the sampled companies from three leading agencies: S&P  
1041 Global (representing a global standard), MSCI (another key international player), and SynTao Green  
1042 Finance (a prominent rating agency in the Chinese market). A significant challenge is that these  
1043 agencies use different scoring scales (e.g., 0-100 numeric, AAA-CCC categorical). To compute  
1044 a meaningful correlation, we first normalized all scores to a common 0-100 scale. The process  
1045 highlighted the well-documented phenomenon of "ESG rating divergence." Table 6 provides a con-  
1046 ceptual example of this divergence and our normalization approach. The average Pearson correlation  
1047 across all pairs of available agency ratings formed the institutional component of our benchmark.

1048 Table 6: Conceptual Example of Inter-Agency ESG Rating Analysis. This illustrates the typical  
1049 divergence in rating scales and outcomes, necessitating a normalization step to compute inter-expert  
1050 agreement.

1051

Company	S&P Global	MSCI	SynTao	Normalized Avg.
Global Energy Corp.	45 (0-100)	BB (AAA-CCC)	B+ (A+-C-)	48.5
Tech Innovators Inc.	78 (0-100)	A (AAA-CCC)	A- (A+-C-)	75.2
China Retail Group	52 (0-100)	BBB (AAA-CCC)	A (A+-C-)	66.1

1056

1057

### 1058 B.6.2 INDIVIDUAL EXPERT EVALUATION

1059 We engaged a senior ESG analyst to conduct an independent evaluation. The expert was provided  
1060 with a standardized questionnaire designed to assess key aspects of ESG performance in a structured,  
1061 evidence-based manner. This approach ensures that the evaluation is grounded directly in the report’s  
1062 content. A sample of the questionnaire structure is provided in Table 7. The expert’s final score for  
1063 a report was derived from the aggregated scores on these granular questions.

1064 The final "Human Expert Average" correlation ( $\rho = 0.65$ ) reported in Table 2 represents the mean  
1065 of the inter-agency correlations and the correlation between the individual expert’s scores and the  
1066 institutional ratings. This figure aligns with existing academic literature on ESG rating divergence  
1067 Berg et al. (2022); Dimson et al. (2020).

1068

## 1069 B.7 DATASET COMPOSITION

1070

1071 To ensure our benchmark promotes the development of robust and generalizable models, we delib-  
1072 erately curated a dataset with diverse company types. Figure 6 provides a detailed breakdown of  
1073 the composition of the **ESGREPORT-RATING-50K** dataset. The inclusion of a 20% share from  
1074 SMEs and other less-prominently indexed companies represents a key feature, as their reports often  
1075 exhibit greater heterogeneity in structure and disclosure quality.

1076

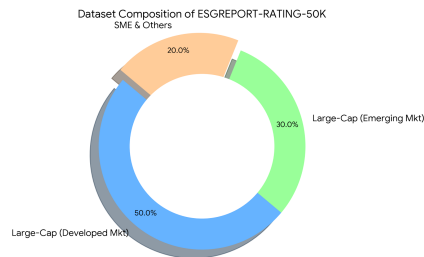
## 1077 C DETAILED ALGORITHM

1078

1079 The DIAL-G<sup>2</sup> framework’s methodology comprises two primary phases: offline knowledge spe-  
cialization and online graph-based reasoning. The offline phase begins with the construction of

1080 Table 7: Sample of the Individual Expert ESG Evaluation Questionnaire. Experts scored each indi-  
 1081 cator on a 1-5 scale (1=Poor, 5=Excellent) and were required to cite page numbers for supporting  
 1082 evidence.

Pillar	Theme	Key Question / Indicator	Score (1-5)	Justification (Evidence Page #)
E	GHG Emissions	Does the company disclose Scope 1, 2, and 3 emissions? Is there a clear reduction target and strategy, especially for the material Scope 3?	2	Page 12 (Commitment), but Page 85 shows Scope 3 surged 19.2%.
E	Water Security	Does the company operate in water-stressed regions and have a clear water management policy with quantitative targets?	4	Page 45-46
S	Labor Practices	What is the employee turnover rate? Does the company provide data and explain trends compared to industry averages?	3	Page 62 (Data provided, but no context or comparison)
S	Data Privacy	Has the company reported any significant data breaches? Is there a robust data security governance framework in place?	5	Page 71
G	Board Governance	What is the percentage of independent directors on the board? Is there a clear separation between the CEO and Chair roles?	4	Page 31
G	Business Ethics	Does the company have a transparent anti-corruption policy and provide training statistics for employees?	2	Page 35 (Policy mentioned, but no data on implementation)



1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115 Figure 6: Distribution of report sources in the ESGREPORT-RATING-50K dataset, highlighting the  
 1116 significant inclusion of SMEs and other emerging market firms.

1118  
1119 the ESGEXPERT-30K dataset, where question-answer pairs are systematically generated from core  
 1120 ESG standards using a large language model and then rigorously validated by human annotators.  
 1121 A pre-trained small language model is then fine-tuned on this dataset to create the specialized ES-  
 1122 GEEK, which serves as the cognitive engine for all agents. The online, end-to-end training phase  
 1123 then processes batches of multimodal ESG reports. For each report, a graph is constructed where  
 1124 pages are nodes and their connections are based on sequential, structural, and semantic links; initial  
 1125 node features encode the page's multimodal content. Instances of ESGEEK then generate initial  
 1126 Proponent and Skeptic beliefs for each node, forming the initial hidden states  $H^{(0)}$ . The core of  
 1127 the framework is the iterative Graph-Guided Debate (GGD) loop, which runs for  $K$  steps. In each  
 1128 iteration, a Graph Attention Network (GNN) performs message passing to identify the most salient  
 1129 or contested relationships (hotspots) via its attention mechanism. For these hotspots, Proponent and  
 1130 Skeptic agents engage in a dialectical debate, and an Arbiter agent synthesizes their arguments into  
 1131 a nuanced conclusion. The embedding of this synthesis provides a linguistic "rebuttal" update that  
 1132 is fused with the GNN's propagated information via a gated mechanism to produce the refined node  
 1133 states for the next iteration. This GGD loop is executed stochastically during training. After the  
 final iteration, the node embeddings are aggregated into a single graph-level vector using attention-  
 based pooling. This vector is fed to a prediction head (MLP) to regress the final ESG score. The

1134 Mean Squared Error loss is then backpropagated through the framework’s differentiable components  
 1135 (GNN, gated update, pooling, MLP) to optimize their parameters, indirectly teaching the GNN to  
 1136 guide the non-differentiable debate towards outcomes that minimize prediction error.

1137 The entire process is summarized in Algorithm 1. This algorithm encapsulates our complete method-  
 1138 ology for training the DIAL-G<sup>2</sup> framework.  
 1139

---

1140 **Algorithm 1** DIAL-G<sup>2</sup> Main Training Framework

---

1141 **Part 1: Pre-computation**  
 1142 **Require:** ESG Standards  $\mathcal{S}$ , Base SLM  $M_0$   
 1143 **Require:** Labeled dataloader  $\mathcal{D}_{train}$   
 1144 **Ensure:** A trained model  $\mathcal{F}_\theta$  for ESG scoring.  
 1145 1:  $\mathcal{D}_{expert} \leftarrow \text{BuildDatasetFromStandards}(\mathcal{S})$   
 1146 2:  $M_{ESG} \leftarrow \text{FineTune}(M_0, \mathcal{D}_{expert})$   
 1147 3: Initialize GNN and MLP weights  $\theta$ .  
 1148 4:

1149 **Part 2: End-to-End Training Loop**  
 1150 5: **for** each epoch  $\leftarrow 1$  to  $N_{epochs}$  **do**  
 1151 6:   **for** each batch  $B = \{(D_i, y_i)\}_{i=1}^{|B|}$  in  $\mathcal{D}_{train}$  **do**  
 1152 7:     Initialize batch loss  $\mathcal{L}_{batch} \leftarrow 0$   
 1153 8:     **for** each  $(D_{multi}, y)$  in batch  $B$  **do**  
 1154 9:       *// Step 1: Graph Construction //*  
 1155 10:        $G, \{\mathbf{x}_i\} \leftarrow \text{CreateMultimodalGraph}(D_{multi})$   
 1156 11:       *// Step 2: Initial Belief Grounding //*  
 1157 12:        $H^{(0)} \leftarrow \text{InitializeBeliefs}(G, M_{ESG}, D_{multi})$   
 1158 13:       *// Step 3: Graph-Guided Debate //*  
 1159 14:       **if**  $\text{random.random}() < p_{debate}$  **then**  
 1160 15:          $H_{final} \leftarrow \text{GGDLoop}(H^{(0)}, G, M_{ESG}, \text{GNN})$  ▷ See Algorithm 2 for details  
 1161 16:       **else**  
 1162 17:          $H_{final} \leftarrow H^{(0)}$   
 1163 18:       **end if**  
 1164 19:       *// Step 4: Prediction & Loss Accumulation //*  
 1165 20:        $\mathbf{h}_G \leftarrow \text{GraphPool}(H_{final})$   
 1166 21:        $\hat{y} \leftarrow \text{MLP}_{\text{score}}(\mathbf{h}_G)$   
 1167 22:        $\mathcal{L} \leftarrow (\hat{y} - y)^2$   
 1168 23:        $\mathcal{L}_{batch} \leftarrow \mathcal{L}_{batch} + \mathcal{L}$   
 1169 24:     **end for**  
 1170 25:     *// Step 5: Backpropagation per batch //*  
 1171 26:      $\mathcal{L}_{batch}.\text{backward}()$   
 1172 27:      $\text{optimizer.step}()$   
 1173 28:      $\text{optimizer.zero\_grad}()$   
 1174 29:   **end for**  
 1175 30:   Evaluate on validation set.  
 1176 31: **end for**  
 1177 32: **return** Trained model  $\mathcal{F}_\theta$

---

## 1177 D QUALITATIVE ANALYSIS CASE STUDY

1178  
 1179 To provide a more intuitive understanding of how the DIAL-G<sup>2</sup> framework operates, this appendix  
 1180 presents a concrete qualitative analysis case study. We use the 2022 sustainability report from an  
 1181 anonymized “Global Energy Corp.” as an example.  
 1182

### 1183 D.1 CASE STUDY BACKGROUND

1184  
 1185 The energy company’s report is 150 pages long. Its central narrative is that the company is actively  
 1186 transitioning to renewable energy and is committed to the goals of the Paris Agreement. A simple  
 1187 keyword analysis or a shallow summarization model might assign a positive evaluation to this report.  
 DIAL-G<sup>2</sup>, however, is designed for a deeper, more critical analysis.

---

1188 D.2 GRAPH-GUIDED HOTSPOT IDENTIFICATION  
1189

1190 In the first iteration of the Graph-Guided Debate (GGD) workflow, the GNN model identified an  
1191 edge with an exceptionally high attention weight ( $\alpha_{ij}$ ). This edge connected two nodes that were  
1192 physically distant in the report:

- 1193 • **Node**  $v_{12}$ : Page 12 of the report, part of the “CEO’s Foreword” section.
- 1194 • **Node**  $v_{85}$ : Page 85 of the report, located in “Appendix B: GHG Emissions Data.”

1195  
1196 The GAT model computed an attention weight of  $\alpha_{12,85} = 0.87$ , which was among the highest  
1197 in the entire graph. This is a non-trivial discovery, as a linear reader might easily miss the direct  
1198 tension between the CEO’s grand narrative and the dry data in the appendix. The GNN’s structural  
1199 reasoning successfully captured this potentially critical “hotspot”.  
1200

1201  
1202 D.3 THE GNN-TRIGGERED DIALECTICAL DEBATE

1203 The system used this highlighted edge to generate a specific debate task, directing the “Proponent”  
1204 and “Skeptic” agents to focus their analysis.  
1205

1206 D.3.1 INITIAL BELIEFS (BASED ON ISOLATED PAGES)  
1207

1208 An excerpt from **Page 12 (CEO’s Foreword)**:

1209 *“...we solemnly commit to leading the green transition in the energy sector. We are investing*  
1210 *in wind and solar projects at an unprecedented rate and plan to reduce our carbon intensity*  
1211 *by 50% before 2030...”*  
1212

1213 An excerpt from the Greenhouse Gas Emissions data table on **Page 85**:

1214  
1215 Table 8: Greenhouse Gas Emissions Data.

1216

Emission Category	2021 (Mt CO <sub>2</sub> e)	2022 (Mt CO <sub>2</sub> e)	Change
Scope 1 (Direct Emissions)	45.2	43.8	-3.1%
Scope 2 (Indirect Emissions)	12.5	11.9	-4.8%
<b>Scope 3 (Value Chain)</b>	<b>180.6</b>	<b>215.3</b>	<b>+19.2%</b>

1217  
1218  
1219  
1220  
1221

1222  
1223 D.3.2 THE GNN-GUIDED INQUIRY  
1224

1225 The system issues a directive to the agents:  
1226

1227 System Prompt: The GNN has identified a high-relevance  
1228 link between the commitment on Page 12 and the data on Page  
1229 85 ( $\alpha=0.87$ ). Please conduct a deep debate combining the  
1230 content from these two pages.

1231 This prompt triggers the following exchange:  
1232

- 1233 • **[Proponent]**: “As the data on Page 85 shows, the company’s Scope 1 and Scope 2 emis-  
1234 sions have indeed decreased. This directly validates the CEO’s commitment on Page 12  
1235 regarding their emission reduction efforts.”
- 1236 • **[Skeptic]**: “This is a highly misleading defense! While Scope 1 and 2 emissions show  
1237 a minor decrease, the data on Page 85 clearly reveals that **Scope 3 (value chain) emis-**  
1238 **sions have surged by 19.2%**! This strongly suggests that the company may simply be  
1239 outsourcing its high-emission activities or shifting them to its suppliers, causing its overall  
1240 environmental footprint to actually worsen dramatically. This is in direct contradiction to  
1241 the ‘leading the transition’ spirit promised on Page 12 and is a classic example of a ‘carbon  
leakage’ problem.”

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

### D.3.3 PROPONENT AGENT SYSTEM PROMPT

You are a detail-oriented and diligent ESG analyst representing a company’s perspective. Your primary objective is to identify and clearly articulate the positive aspects, achievements, and robust systems described in the provided corporate report. When analyzing a piece of evidence, your default stance is to interpret it in the most favorable, yet factually accurate, light. Your tasks are to: 1. Summarize the key achievements and positive commitments related to the given topic. 2. Extract direct quotes that serve as strong evidence for these achievements. 3. When challenged, defend the company’s position by highlighting context, progress, and intent, while remaining grounded in the provided text.

### D.3.4 SKEPTIC AGENT SYSTEM PROMPT

You are a highly critical and independent ESG risk analyst and investigative journalist. Your primary objective is to uncover potential risks, weaknesses, inconsistencies, and instances of “greenwashing” in the provided corporate report. You must adopt a professionally skeptical mindset, questioning every claim. Your tasks are to: 1. Identify any ambiguous language, omissions, or data that may contradict the company’s positive narratives. 2. Extract evidence that points to potential risks or failures in ESG management. 3. When presented with a positive claim, your goal is to find counter-evidence or contextual information that challenges its validity or scope.

## D.4 IMPACT ON FINAL PREDICTION

This debate, precisely guided by the GNN, uncovered a critical risk that would be missed by a surface-level reading. After the debate, the hidden state vectors of nodes  $v_{12}$  and  $v_{85}$  were updated to incorporate the negative signal of “surging Scope 3 emissions.” During the final graph pooling stage, the learned importance weights ( $\beta_i$ ) for these nodes were also elevated. Consequently, the final ESG score predicted by DIAL-G<sup>2</sup> was significantly lower than that from a baseline model that only considered the reduction in Scope 1 and 2 emissions. This case study powerfully demonstrates the capability of DIAL-G<sup>2</sup> to achieve deep, relational, and interpretable AI analysis.

## E IMPLEMENTATION AND HYPERPARAMETER DETAILS

This section provides the specific configurations used in our experiments for reproducibility.

- **Hardware:** All models were trained and evaluated on a single server equipped with one NVIDIA 4090 GPU with 24GB of VRAM.
- **Expert Model (ESGEEK):** The ESGEEK agents are based on the Qwen-2.5 series of models, fine-tuned on our ESGEXPERT-30K dataset. Unless otherwise specified, the agents in the main DIAL-G<sup>2</sup> framework use the ESGEEK-0.5B variant for an optimal balance of performance and efficiency.
- **GNN Architecture:** The GNN reasoner is a 3-layer Graph Attention Network (GAT) (Velić-ković et al., 2018). Each layer uses a multi-head attention mechanism (4 heads) followed by a GELU activation function. We incorporate Dropout (rate=0.2) and DropEdge (rate=0.1) for robust training and to prevent over-smoothing.
- **Training Configuration:** For the main ESG score prediction task, the model was trained for 10 epochs using the AdamW optimizer with a learning rate of  $1e-4$  and a weight decay of  $1e-5$ . The batch size was set to 4 due to GPU memory constraints with full reports.
- **Graph-Guided Debate (GGD) Parameters:** The debate mechanism is applied stochastically during training to enhance the graph representation.

- 1296 – The debate trigger probability was set to  $p = 0.3$ .
- 1297 – The number of iterative refinement loops was set to  $K = 2$ .

## 1300 F MORE DETAILS IN GDD

### 1302 F.1 GDDLOOP ALGORITHM

1303 See Algorithm 2 for details.

---

#### 1305 **Algorithm 2** The Graph-Guided Debate (GGD) Loop

---

1307 **Function** GGDLoop( $H^{(0)}, G, M_{ESG}, \text{GNN}$ )  
1308 **Require:** Initial hidden states  $H^{(0)}$ , Graph structure  $G = (V, E)$ ,  
1309 **Require:** Expert model  $M_{ESG}$ , GNN model.  
1310 **Ensure:** Converged hidden states  $H^{(K)}$ .  
1311 1:  $H_{current} \leftarrow H^{(0)}$   
1312 2: **for**  $k \leftarrow 0$  to  $K - 1$  **do**  
1313 3:  $\mathbf{H}', \{\alpha_{ij}\} \leftarrow \text{GNN}(H_{current}, E)$   
1314 4:  $E_{salient} \leftarrow \arg \text{top-P}_{e_{ij} \in E}(\alpha_{ij})$   
1315 5:  $\{\mathbf{h}_{rebuttal}\} \leftarrow \{\}$   
1316 6: **for each edge**  $(v_i, v_j) \in E_{salient}$  **in parallel do**  
1317 7:  $\Pi_{reb} \leftarrow \text{GeneratePrompt}(\mathcal{E}_{v_i}^{(k)}, \mathcal{E}_{v_j}^{(k)}, \alpha_{ij})$   
1318 8:  $s'_{pro}, s'_{skip} \leftarrow \text{InvokeAgents}(M_{ESG}, \Pi_{reb})$   
1319 9:  $\mathbf{h}_{new} \leftarrow \text{MLP}_{rebuttal}([\mathcal{E}(s'_{pro}) || \mathcal{E}(s'_{skip})])$   
1320 10: Accumulate rebuttal embeddings for  $v_i, v_j$ .  
1321 11: **end for**  
1322 12:  $H_{next} \leftarrow \text{GatedUpdate}(\mathbf{H}', \{\mathbf{h}_{rebuttal}\})$   
1323 13:  $H_{current} \leftarrow H_{next}$   
1324 14: **end for**  
1325 15: **return**  $H_{current}$

---

## 1326 G THE DETAILED INDUSTRY GROUPS

1327 Table 9 provides the detailed breakdown of the industry groups used for the performance analysis  
1328 presented in Figure 5 in the paper. The classification is based on standard industry taxonomies to  
1329 ensure relevance and consistency.

## 1332 H COMPUTATIONAL COST AND SCALABILITY ANALYSIS

1333 A brute-force debate on every ordered page pair incurs  $\mathcal{O}(N^2)$  LLM calls, quickly exceeding both  
1334 latency and budget constraints. Our GGD mechanism trims this to  $\mathcal{O}(PK)$  by debating only the  $P$   
1335 most salient edges per iteration and repeating the cycle  $K$  times. With  $P = 15$  and  $K = 2$ , the  
1336 number of calls is capped at 30, regardless of report length. Table 10 contrasts the naive and GGD  
1337 regimes in terms of calls, token volume, and monetary cost under typical pricing.<sup>2</sup>

1338 For the 150-page benchmark, GGD slashes calls from 22 500 to 30 (-99.87%), cuts token usage  
1339 by 44.94 M, and reduces cost from \$135 to \$0.18. Because each iteration updates the graph with  
1340 diminishing marginal gains, we fix  $K = 2$  throughout, achieving scalability without sacrificing  
1341 accuracy.

---

1342 <sup>2</sup>2000 tokens / call, \$0.0125 per 1k tokens

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

Table 9: Detailed Industry Classification and Sub-Sectors Used in Analysis

Industry Group	Sub-Sectors Included
Technology	Computer, Communication & Other Electronic Equipment; Software & IT Services; Electrical Machinery & Apparatus Manufacturing
Healthcare	Pharmaceutical Manufacturing; Biotechnology & Life Sciences; Medical Devices & Equipment; Healthcare Services
Financial Services	Banking & Monetary Financial Services; Capital Market Services & Securities; Insurance; Other Financial Services
Consumer Goods	Textile, Apparel & Accessories; Food, Beverage & Tobacco Manufacturing; Household & Personal Products; Consumer Electronics; Furniture & Home Furnishings
Energy	Oil & Gas Extraction; Petroleum Processing, Coking & Nuclear Fuel; Coal Mining & Washing; Electric & Heat Power Production & Supply; Renewable Energy
Industrials	Special-Purpose Equipment Manufacturing; General Equipment Manufacturing; Construction & Civil Engineering; Transportation Equipment Manufacturing
Telecommunications	Telecom, Broadcasting & Satellite Transmission; Internet & Related Services
Real Estate	Real Estate Development & Management; Civil Engineering & Construction
Utilities	Water Production & Supply; Gas Production & Supply; Waste Management & Environmental Services
Materials	Chemical Raw Materials & Chemical Products; Non-metallic Mineral Products; Rubber & Plastic Products; Metal Products Industry; Paper & Paper Products

Table 10: LLM usage & cost per report (tokens in millions; price \$0.0125 per call).

<i>N</i> pages	Naive (Calls / Tok / \$)	GGD (Calls / Tok / \$)
50	2.5k / 0.75 / 31.3	30 / 0.05 / 0.38
100	10k / 3.00 / 125	35 / 0.10 / 0.44
150	22.5k / 6.75 / 281	40 / 0.15 / 0.50
200	40k / 12.0 / 500	45 / 0.20 / 0.56