

Honey, I Shrunk the Language Model: Impact of Knowledge Distillation Methods on Performance and Explainability

Daniel Hendriks, Philipp Spitzer, Niklas Kühl, and Gerhard Satzger

Abstract—Artificial Intelligence (AI) has increasingly influenced modern society, recently in particular through significant advancements in Large Language Models (LLMs). However, high computational and storage demands of LLMs still limit their deployment in resource-constrained environments. Knowledge distillation addresses this challenge by training a small student model from a larger teacher model. Previous research has introduced several distillation methods for both generating training data and for training the student model. Despite their relevance, the effects of state-of-the-art distillation methods on model performance and explainability have not been thoroughly investigated and compared. In this work, we enlarge the set of available methods by applying critique-revision prompting to distillation for data generation and by synthesizing existing methods for training. For these methods, we provide a systematic comparison based on the widely used Commonsense Question-Answering (CQA) dataset. While we measure performance via student model accuracy, we employ a human-grounded study to evaluate explainability. We contribute new distillation methods and their comparison in terms of both performance and explainability. This should further advance the distillation of small language models and, thus, contribute to broader applicability and faster diffusion of LLM technology.

Index Terms—Explainability, language model, knowledge distillation, performance.

I. INTRODUCTION

AI has become integral to modern society, profoundly influencing daily life, industries, and innovation. The advent of LLMs has particularly advanced fields such as Natural Language Processing (NLP) [1]–[5] for tasks such as question answering and reasoning [6]–[10]. One essential factor contributing to LLM’s effectiveness is the size of these models with several hundred billion parameters [11]–[14]. The substantial size of LLMs presents considerable challenges to their application in low-resource environments, such as mobile phones and edge devices [15]–[21]. The development

of smaller language models is predicated on research in knowledge distillation, which focuses on the transfer of knowledge from a large model (the teacher) to a smaller one (the student). In this process, the teacher model is employed to generate data, and the student model is subsequently fine-tuned.

Knowledge distillation aims to maintain the performance and capabilities of the teacher while improving deployment ease, energy efficiency, and inference speed [17], [22]–[25]. The importance of distillation becomes visible when looking at the landscape of high-quality open-source models: language models such as Alpaca [26] and Code Llama [27] have been trained on data from larger language models. Besides the performance of the student model, explainability is vital in knowledge distillation, as it can help verify if the model has learned meaningful concepts and can effectively convey them to humans [28], [29]. Explainability, a language model’s ability to reason its outputs towards a human, is crucial if we want humans to understand and trust outputs from language models and, ultimately, adopt them not only in everyday life but also for high-stakes decisions. In this light, past work has emphasized and studied the aspect of explainability in AI systems via human-grounded studies [30]–[33]. However, research has so far neglected how knowledge distillation affects the explainability of language models.

Past work has introduced numerous methods to perform knowledge distillation and improve the final capabilities of the student model. These methods aim to enhance either the training data or the student model. On the one hand, training data is essential, as the performance of student models is significantly affected by the quality of the training data [34]. To align LLMs’ behavior and outputs with human values, such as omitting harmful or biased responses, research has shown that data quality can be enhanced by prompting a language model to improve its own outputs [35], a strategy we refer to as *critique-revision prompting* [35]. However, the effectiveness of this prompting strategy has not been explored for knowledge distillation. On the other hand, the training of student models is crucial, as demonstrated by recent work by Hsieh et al. [17], achieving state-of-the-art performance on various NLP benchmarks by training a student via multitask training. Nonetheless, the authors do not explore the impact of multitask training on the student’s ability to explain outputs to a human. A second example is the work by Wang et al. [34] that focuses on counterfactual training resulting in increased

Corresponding author: Daniel Hendriks (e-mail: daniel.hendriks@kit.edu). This work involved human subjects in its research. Approval of all ethical and experimental procedures and protocols was granted by the ethics commission of the Karlsruhe Institute of Technology (KIT).

Daniel Hendriks, Philipp Spitzer, and Gerhard Satzger are with the Institute for Information Systems (WIN), Karlsruhe Institute of Technology, 76133 Karlsruhe, Baden-Württemberg, Germany. Niklas Kühl is with the Information Systems (WI) Institute, University of Bayreuth, 95440 Bayreuth, Bayern, Germany.

This article has supplementary downloadable material available at <https://github.com/DanielHendriks/llm-distillation> and 10.21227/4gsd-p846, provided by the authors.

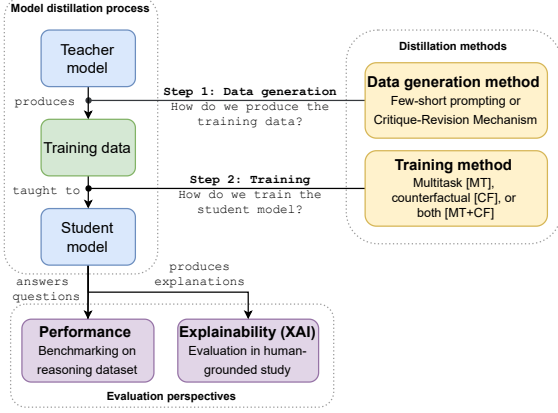


Fig. 1. Research model overview: This work introduces novel methods for training data generation and training and establishes approaches to compare performance and explainability of student models.

faithfulness and consistency in student models, two pillars of a model capable of explaining its outputs. However, their automated LLM-based evaluation may lack reliability. Further research using human-centered methods, such as surveys and controlled experiments with participants, is necessary to verify their results. Lastly, comparing the results of existing studies regarding performance and explainability is challenging due to differences in teacher model selection and student model training settings. For instance, Hsieh et al. [17] use the proprietary PaLM model [12], whereas Wang et al. [34] employ GPT-neox [36]. Overall, the shortcomings of (1) the lack of evaluation from the view of performance and explainability and (2) the lack of comparability between results motivate our work and give rise to the following two research questions:

- RQ1: How does critique-revision prompting affect data quality and consequently the performance and explainability of student models?
- RQ2: How do current training methods affect the performance and explainability of student models?

To address our research questions, we systematically compare distillation methods to evaluate their impact on performance and explainability and combine them to limit potential trade-offs. The high-level logic of this work’s research model is presented in Figure 1: We generate explanations with the teacher through few-shot prompting and aim to improve them through critique-revision prompting. When training the student models, we utilize multitask training [17], counterfactual training [34], or a combination of both. Performance is empirically evaluated by measuring the accuracy of student models in answering multiple-choice questions from the Commonsense Question Answering (CQA) dataset. The explainability of the student model is assessed by generating natural language explanations and evaluating them in a within-subject study with 117 participants.

In this work, we make three significant contributions:

- We enrich the *spectrum of available distillation methods*. We apply critique-revision prompting to distillation to

improve training data generation. We also implement a new training method combining multitask and counterfactual training. We find that applying critique-revision prompting to generate explanations and the combined training method to produce a student model enhances its ability to explain the provided answers to humans.

- We develop a *standardized framework for method comparison*—closing a critical gap in current research where heterogeneous teacher models and training configurations have prevented direct comparisons. This also includes a way to assess explainability with human-grounded studies.
- We provide a *comprehensive evaluation* of these distillation methods in terms of both performance and explainability. Our comparison shows that models trained with multitask training perform strongly compared to other methods. In terms of explainability, the student model that excels was trained through a combination of multitask and counterfactual training, with explanations enhanced via critique-revision prompting.

II. RELATED WORK

This section reviews existing research on the knowledge distillation of LLMs, focusing on key methods, challenges, and innovations. We begin by discussing knowledge distillation for LLMs, examining its advantages in terms of model size, efficiency, and applicability to resource-constrained environments. We then explore various approaches to data generation, highlighting how explanations and iterative refinement techniques can enhance distillation outcomes. Finally, we review training methods, such as multitask and counterfactual training, which aim to optimize the distillation process for reasoning tasks. Together, these areas form the foundation for improving the capabilities of smaller models while preserving the strengths of their larger counterparts.

A. Knowledge Distillation for LLMs

Knowledge distillation is a technique introduced by Hinton et al. [37] to transfer knowledge and skills from a large language model to a smaller one. Subsequent research has expanded upon this idea [13], [19], [38], referring to the larger model as the *teacher* and the smaller model as the *student*. In contrast to the teacher, the smaller student model offers four key advantages arising from its reduced size and, thus, lower computational requirements:

- 1) The student can be deployed in resource-constrained environments such as edge devices, smartphones, or limited hardware [37], [39]. Enabling new applications, this advantage has also been leveraged extensively in computer vision applications [40]–[42].
- 2) Users can more easily deploy small language models locally, avoiding the high costs associated with proprietary large language models and protecting sensitive data [43].
- 3) Their compact size makes smaller models perform inference more quickly and efficiently, leading to greater environmental sustainability [43].

- 4) Small models offer lower latency, making them particularly suitable for time-sensitive applications such as entity recognition in industrial settings [44]. These advantages collectively motivate the need for effective knowledge distillation methods.

To distill an LLM, researchers have proposed various parameter-efficient techniques, including adapters [45], low-rank structures [46], and prompt tuning [18], [47], [48]. Nevertheless, for distilling reasoning capabilities, the prevalent paradigm specifically involves generating data using the teacher model and subsequently training the student model on this synthetic data [43]. If successful, this distillation approach enables the student not only to generate coherent text by accurately predicting tokens [49] but also to mimic the teacher’s reasoning processes and correctly estimate uncertainties [49]. Effective distillation of a language model requires two essential ingredients: (a) context-rich and skill-specific training data [43], and (b) an appropriate and effective training method [17], [34]. Both elements are crucial for the student to acquire new capabilities, such as reasoning, by genuinely learning underlying concepts and logic rather than merely memorizing answers or exploiting correlations. To verify that the student model has indeed learned these concepts and logical reasoning, we do not only evaluate the student’s *performance* in this work, but also its capacity to provide explanations for its answers—an ability that is referred to as *explainability* [50], [51].

B. Data Generation for Knowledge Distillation

To obtain diverse and rich training data, researchers have proposed techniques prompting the teacher model to explicitly explain its reasoning process step-by-step [6], [20], [52]. Existing studies demonstrate that incorporating such explanations into distillation leads to capable student models [17], [21], [34], [53]–[56] and that enriching or filtering these explanations further improves student performance [5], [19], [57].

In this work, we similarly generate explanations from the teacher model, but additionally evaluate a novel technique to expand and refine the data. Specifically, we prompt the teacher model to critique and subsequently revise its own explanations. Inspired by Constitutional AI [35], we refer to this iterative refinement as *critique-revision prompting*. Initially, Constitutional AI has employed this prompting strategy to produce training data to develop language models that are more helpful and less harmful; however, the authors note that the prompting strategy can guide data generation toward any desired objective. Here, we leverage this approach explicitly to generate high-quality explanations for effective knowledge distillation.

C. Training Methods for Knowledge Distillation

An effective training method can reliably and effectively support the student in learning concepts and reasoning patterns through a specialized loss function [17], [34], [38]. One such method is *multitask training*, which simultaneously teaches the student both explanation and answering tasks. Researchers

have proposed using explanations, in some works referred to as “rationales”, in knowledge distillation by feeding them as additional inputs during training [58], [59]. However, this approach requires a teacher model, thus making deployment more complicated and inference neither more energy-efficient nor faster. To address this, Hsieh et al. [17] frame the problem as a multitask problem and investigate how this training method improves small language models’ performance on reasoning datasets. This approach performs better with fewer training examples than conventional training. Similarly, *counterfactual training* also instructs the student on two concurrent tasks, but differs in the specific tasks involved. Rather than explanation and answering, counterfactual training guides the model to (a) reason and answer correctly when prompted solely with a question, and (b) intentionally answer incorrectly when provided with both the question and an incorrect explanation. This method is motivated by the goal of improving a model’s robustness and sensitivity to misleading information.

III. PRELIMINARIES

A. Critique-Revision Prompting

To formalize the described prompting strategy within the context of critique-revision prompting illustrated by Figure 2 on p. 4, we define the components of the process and specify how each step is handled during the generation and critique phases.

We define the dataset as:

$$\mathcal{D} = \{(q_i, a_i, e_i, c_i, e'_i)\}_{i=1}^N$$

where q_i represents the question for the i -th example, a_i is the correct answer for the i -th example, e_i is the initial explanation generated based on q_i and a_i , c_i is the critique of the initial explanation e_i , e'_i is the revised explanation generated after applying the critique c_i .

In the first of three steps, depicted in Figure 3 on p. 6, the teacher model is presented with the multiple-choice question q_i , the correct answer a_i , and the instruction to “Explain.” The model generates an initial explanation e_i :

$$e_i = \text{TeacherModel}(q_i, a_i, \text{“Explain”})$$

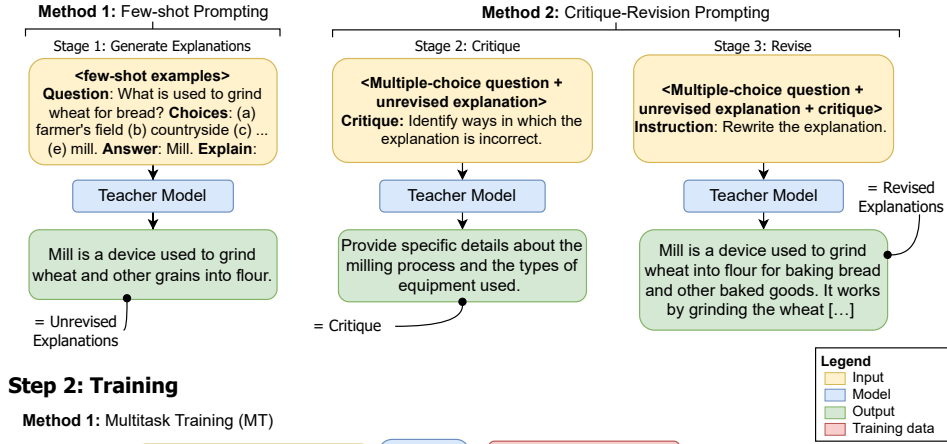
In the second step, the teacher model is fed the complete context consisting of q_i , a_i , and e_i , and is instructed to critique e_i . The critique c_i serves to identify mistakes or areas for improvement in the explanation. This step can be represented as:

$$c_i = \text{TeacherModel}(q_i, a_i, e_i, \text{“Critique”})$$

In the third and final step, the teacher model is provided with the complete context q_i , a_i , e_i , and c_i , and is instructed to revise e_i based on the critique c_i . The revised explanation e'_i is generated as:

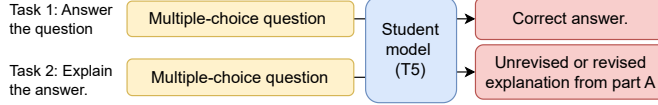
$$e'_i = \text{TeacherModel}(q_i, a_i, e_i, c_i, \text{“Revise”})$$

Step 1: Data generation

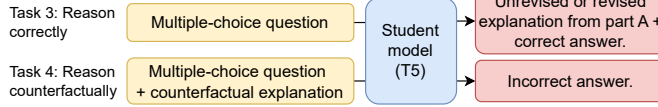


Step 2: Training

Method 1: Multitask Training (MT)



Method 2: Counterfactual Training (CF)



Method 3: Combining Multitask and Counterfactual Training (MT+CF)

Train student model to perform task 1-4 simultaneously.

We train four student models:		
#	Data generation method	Training method
1	Few-shot Prompting	Multitask
2	Few-shot Prompting	Counterfactual
3	Few-shot Prompting	Both
4	Critique-Revision Prompting	Both

Fig. 2. Methods presented in the section on preliminaries are applied in two steps: In **Step 1**, we generate explanations and then improve them by critiquing and revising the explanation with the teacher. In **step 2**, the explanations are used to fine-tune student models with one of three training methods: multitask training, counterfactual training, or a combination of both. We focus on the four student models shown in the Table for two reasons. On the one hand, they are the most promising candidates in terms of both performance and explainability based on preliminary experiments. On the other hand, to avoid overloading study participants during explainability evaluation, we limit the number of possible student model combinations.

B. Multitask Training

Multitask training utilizes the *answer* from the CQA dataset and *explanations*, the initial unrevised explanations e_i or the revised explanations e'_i , to train a student model for the task of answering multiple-choice questions and explaining its predictions. On the one hand, the answers are used as a ground truth to train the student model to answer correctly, by minimizing the following loss function:

$$\mathcal{L}_{\text{answer}} = \frac{1}{N} \sum_{i=1}^N l(f(q_i), a_i) \quad (1)$$

On the other hand, the student model learns to explain its answers from the explanations according to the following loss function:

$$\mathcal{L}_{\text{explanation}} = \frac{1}{N} \sum_{i=1}^N l(f(q_i), e_i) \quad (2)$$

During training, we insert a label into the prompt (e.g., *[answer]* or *[explain]*) to signal the model to either answering a multiple-choice question or providing an explanation

(see Figure 2) and update the model's weights based on the gradients from the combined loss:

$$\mathcal{L}_{\text{multitask}} = \mathcal{L}_{\text{answer}} + \mathcal{L}_{\text{explanation}} \quad (3)$$

C. Counterfactual Training

Performing counterfactual training requires generating counterfactual explanations, i.e., explanations that lead to an incorrect answer, by feeding the teacher model an incorrect answer a^*_i , formalized as follows:

$$e^*_i = \text{TeacherModel}(q_i, a^*_i, \text{"Explain"})$$

Counterfactual training itself follows a similar dual-task structure like multitask training but differs in the objectives assigned to the model. Instead of answering and explaining, counterfactual training consists of two complementary tasks: (1) answering *and* explaining correctly when given only a question, and (2) answering incorrectly when given a question along with an incorrect explanation. Formally, we define the loss functions for these two tasks as follows:

$$\mathcal{L}_{correct} = \frac{1}{N} \sum_{i=1}^N l(f(q_i), a_i) \quad (4)$$

$$\mathcal{L}_{incorrect} = \frac{1}{N} \sum_{i=1}^N l(f(q_i, e_i^*), a_i^*) \quad (5)$$

The training objective is to minimize the combined loss:

$$\mathcal{L}_{counterfactual} = \mathcal{L}_{correct} + \mathcal{L}_{incorrect} \quad (6)$$

By explicitly training the model to recognize incorrect reasoning and respond accordingly, counterfactual training aims to improve its robustness to misleading or spurious explanations and to enhance model interpretability and generalization [60], [61]. In this work, we hypothesize that combining both multitask and counterfactual leads to improvements in performance and explainability. The idea is to leverage the benefits of both training methods to produce a more capable student model. Based on the previous formalizations, we define the loss of the combined training method as the simple unweighted addition of both losses:

$$\mathcal{L}_{combined} = \mathcal{L}_{multitask} + \mathcal{L}_{counterfactual} \quad (7)$$

IV. METHODOLOGY

This section is divided into two main parts. First, we describe how we evaluated the impact of knowledge distillation methods on model performance through an experiment using a standardized dataset and consistent training settings. Section IV-A outlines our datasets, data generation process, training configurations, and evaluation criteria. By systematically analyzing the accuracy of the student models, we quantify the benefits and limitations of different distillation methods.

Second, we show the evaluation of the student models' explainability through a human-based study in which participants evaluated model-generated explanations for multiple-choice questions. Section IV-B details our study design, including participant recruitment and explanation quality criteria. By analyzing human ratings across five key dimensions, we determine how different training and data generation methods affect explanation effectiveness.

A. Experiment

1) *Datasets and Data Generation:* We evaluated the student models on the CQA dataset [62], consisting of multiple-choice questions with five answer choices, one of which is correct. This dataset, also employed in related studies [17], [34], comprises 9,741 training samples and 1,221 test samples. To ensure consistent training and validation across models, we fixed the random seed, enhancing result comparability and reproducibility.

For explanation generation and critique-revision prompting, we selected LLaMA-2-13B as the teacher model, a state-of-the-art language model at the time of conducting this study [63]. Depending on the generation step, we employed either its unaligned or chat-tuned version, following prior research [34],

[64]. The settings for critique-revision prompting included 300 new tokens and a temperature of 1. We observe that critiques often exaggerate criticism, and explanations lengthen after applying critique-revision prompting. Initial validation with corrupted samples from the CQA dataset showed that revised explanations add context-relevant information and improve answer differentiation. The generation of counterfactual explanations, which support incorrect answers, is analogous to the initial generation of explanations in critique-revision prompting.

2) *Training Settings:* We used pre-trained T5 architectures as student models: T5-base with 220M million (220M) parameters and T5-large with 770 million (770M) parameters [65]. Trained via the HuggingFace API [66], student models underwent multitask training, counterfactual training, or a combination of both. Consistent settings were used regardless of the training method: Initial models were trained for up to 10,000 steps; however, since loss and accuracy converged by 5,000 steps, subsequent models with different random seeds were trained for 5,000 steps. Training employed the AdamW optimizer with a learning rate of 5×10^{-5} and a maximum generation length of 300 tokens to prevent truncation. Run on a GRID V100S-16C GPU with 32 GB of memory, reproducibility was ensured by fixing the random seed.

3) *Evaluation:* We assessed both performance and explainability of student models: For performance, *accuracy* served as the primary metric. Accuracy is defined as the proportion of correctly predicted answers out of the total number of samples: $\text{accuracy} = \frac{N_{\text{correct}}}{N_{\text{total}}}$, where N_{correct} denotes the number of correctly predicted answers, and N_{total} represents the total number of evaluated samples. Explainability was assessed through human ratings on explanations generated by student models, following established methodologies [31], [32], [67]. The data for this evaluation came from a standardized survey where participants rated explanation quality across five dimensions, ensuring a reliable and valid assessment.

B. Human-grounded Study

Our study measured how data generation and training methods affect student models' explainability. We evaluated this by having human participants assess explanations generated by these models for multiple-choice questions from the CQA dataset. After obtaining approval from our university's ethics commission, we recruited 117 participants through Prolific [68]. Each participant assesses 12 explanations (three from each student model) across five quality dimensions, yielding $n = 7,020$ total assessments ($117 \text{ participants} \times 12 \text{ explanations} \times 5 \text{ dimensions}$). Initially, 120 participants completed the study; three were excluded due to failing attention checks. The participant sample consisted of 59 male, 57 female, and 1 diverse respondents, predominantly from the United Kingdom ($n = 104$). The most represented age groups were 30-34 years ($n = 23$), 25-29 years ($n = 19$), and 40-44 years ($n = 15$). Most participants had higher education ($n = 86$) such as university degrees, vocational university diplomas, A-levels, and other advanced certifications, with 67 employed (including employees, self-employed, and civil servants) and 13 unemployed or seeking employment.

Study to measure explainability of student models

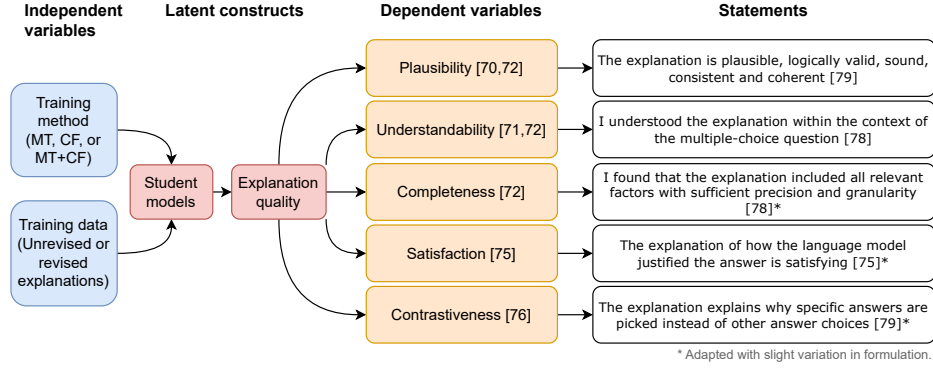


Fig. 3. In the within-subject study ($N = 117$), we measure the effect of using four different student models on their explanation quality along five dimensions.

Based on literature on explainability [69]–[72], we defined explanation quality through the following five dimensions:

- 1) *Plausibility* is the characteristic of an explanation to be sound and reasonable [73].
- 2) *Understandability* is the characteristic of an explanation to be expressed clearly and unambiguously, enabling a human to grasp its meaning easily [73].
- 3) *Completeness* is the characteristic of an explanation to describe the reasoning process at an appropriate level of detail and include relevant information [73].
- 4) *Satisfaction* is “[...] the degree to which users feel that they sufficiently understand the AI system or process being explained to them.” [74, p. 3].
- 5) *Contrastiveness* is “[...] the characteristic of an explanation to justify why a prediction was made instead of another.” [75, p. 18].

To measure these dimensions, we adapted statements from existing human studies and explanation evaluation guidelines [34], [76]–[80]. Participants rated explanations using a five-point Likert scale after reviewing each multiple-choice question and its corresponding explanation (see Figure 3).

To ensure reliability, we used a standardized questionnaire [80] with three sections: task introduction, explanation quality evaluations, and socio-demographic data collection. The 117 participants provided 1,404 observations (117×12). We randomized the order of statements and tasks to reduce order effect bias and included only correctly answered explanations to focus on explanation quality rather than answer correctness. We also incorporated attention checks to improve data quality. Participants from Prolific completed our survey in February 2024 and received £2.25 for their participation, contingent on passing the attention check. All participants provided written informed consent prior to participation in this study.

V. RESULTS AND DISCUSSION

A. Effect on Performance

The accuracy of models on the CQA test dataset is depicted in Figure 4(a), with the T5-base model on the left and the T5-larger model on the right. Before analysis the performance

results from the experiment, outliers identified by the box-plot method, defined as data points lying outside $1.5 \times IQR$, are excluded [81]. For the smaller student models, multitask training with unrevised explanations (MT:Unrevised) achieves the highest mean accuracy, whereas counterfactual training with unrevised explanations yields the lowest performance. For the larger student model, multitask training and the combined training method (MT:Unrevised and MT+CF:Unrevised), both using unrevised explanations, exhibit comparable accuracy, with the combined training method demonstrating a slight advantage. The MT+CF:Revised student models lag behind multitask training on unrevised explanations (MT:Unrevised) in accuracy for both model sizes. In the following paragraphs, the results will be elucidated by conducting a statistical analysis.

We analyzed the performance assessment results using an analysis of variance (ANOVA). Before conducting this analysis, we verified the necessary assumptions about data distribution after excluding several outliers. The Shapiro-Wilk test confirmed normality, and the Levene test indicated equal variances, validating the suitability of the ANOVA approach. The ANOVA results revealed significant performance differences among the student models. Specifically, the type of student model significantly affected performance ($p = 3.76 \times 10^{-6}$), demonstrating the impact of critique-revision prompting and training methods. Additionally, model size significantly influenced performance ($p = 4.96 \times 10^{-14}$). However, no significant interaction emerged between student model type and model size ($p = 0.804$). This absence of interaction indicates that performance differences due to training and data generation methods are independent of model size, allowing conclusions about these methods to be drawn independently.

Critique-revision prompting yields no performance improvement. Comparing the model MT+CF:Unrevised with MT+CF:Revised allows us to conclude the effect of critique-revision prompting on performance. Based on results from the Tukey-Kramer test presented in Table I, we see that critique-revision prompting negatively affects the larger student model. At the same time, there is no significant effect on the smaller student models’ performance. This is unex-

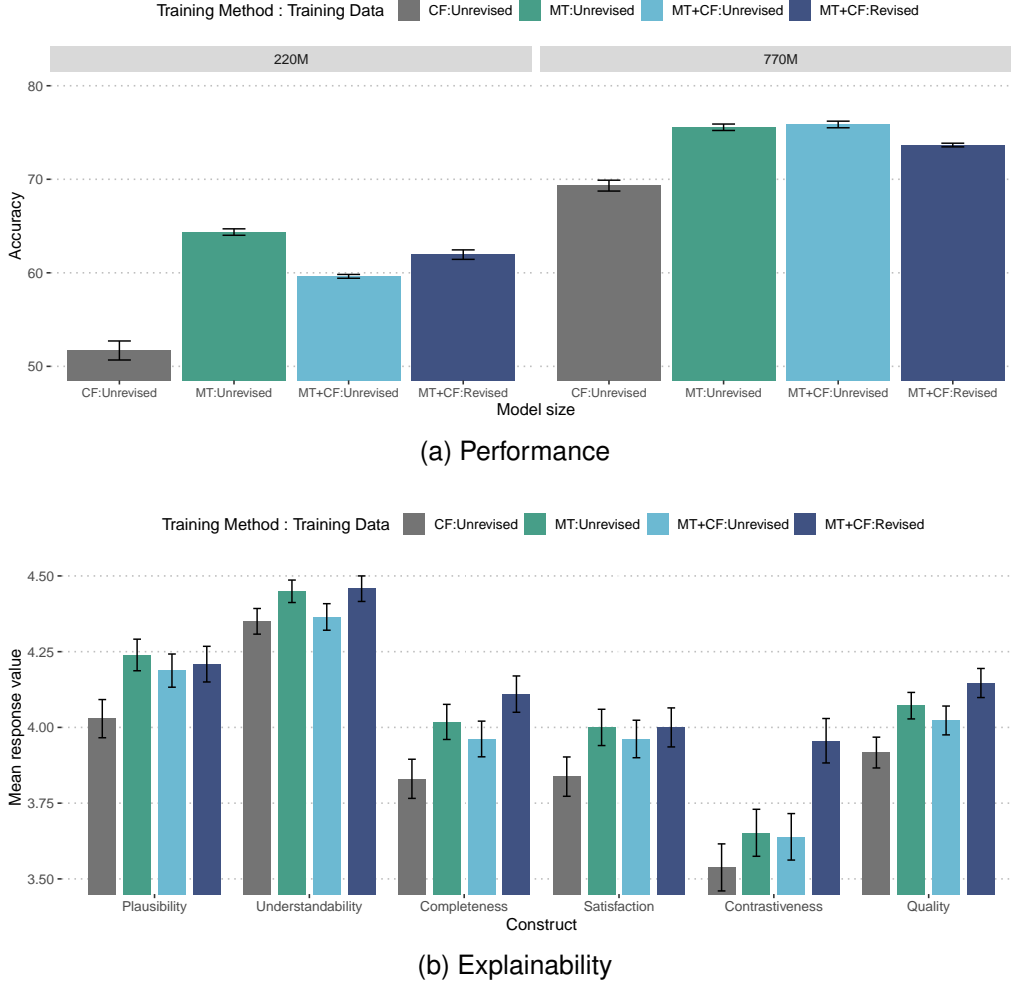


Fig. 4. Comparison of different student models in terms of (a) accuracy performance and (b) explainability measured along the five dimensions of plausibility, understandability, completeness, satisfaction, and contrastiveness. "Quality" as a calculated concept represents the arithmetic average across the five dimensions. Abbreviations: CF for counterfactual training and MT for multitask training.

pected, as critique-revision prompting was intended to enhance data quality and improve model performance. Possibly, the lengthy explanations confused the model, causing performance degradation. Equally surprising is that the effect of critique-revision prompting depends on the student model's size. This dependency might arise because smaller models become more overwhelmed and confused by seeing prolonged explanations during training, resulting in decreased performance or no effect. In contrast, larger models might be able to handle prolonged explanations simply because of their increased number of parameters and, thus, enhanced reasoning capabilities and ability to store information in their parameters.

Multitask training performs superior to counterfactual training. The Tukey-Kramer test shows that the model trained with counterfactual training on unrevised explanations (CF:Unrevised) is consistently outperformed by all other student models for both model sizes, as indicated by the significant difference in accuracy (rows 1 to 3 for the smaller student model and 7 to 9 for the larger student model in Table I). In particular, counterfactual training is also outperformed by models trained with multitask training, leading to the conclusion that multitask training is superior to counterfactual

training. The rationale for this observation might be that training a student model to answer multiple-choice questions incorrectly based on counterfactual explanations might confuse the model rather than teaching it more faithful reasoning. In contrast, teaching the student model the correct explanation yields superior performance as seen with multitask training.

Combining multitask and counterfactual training yields no performance improvement. Our findings for the comparison between the student model trained with multitask training on unrevised explanations (MT:Unrevised) and the student model trained with both multitask and counterfactual training on unrevised explanations (MT+CF:Unrevised) are contradictory: for the smaller student models, combining both training methods harms the performance, whereas, for the larger student model, performances are similar (see row 4 and 10). Consequently, overall, across both model sizes, the performance has not improved significantly. Previously, we found that multitask training significantly outperforms counterfactual training. In that light, the lack of performance improvements from combining these methods may result from counterfactual training providing no additional benefit and, therefore, failing to complement multitask training.

TABLE I

PAIRWISE COMPARISONS OF MODEL PERFORMANCE USING THE TUKEY TEST. A POSITIVE ESTIMATE INDICATES THAT MODEL 2 OUTPERFORMS MODEL 1. SIGNIFICANTLY BETTER MODELS ARE BOLD.

#	Size	Model 1	Model 2	Estimate
1	220M	CF:Unrevised	MT:Unrevised	12.66***
2	220M	CF:Unrevised	MT+CF:Unrevised	7.93***
3	220M	CF:Unrevised	MT+CF:Revised	10.25***
4	220M	MT:Unrevised	MT+CF:Unrevised	-4.73***
5	220M	MT:Unrevised	MT+CF:Revised	-2.41
6	220M	MT+CF:Unrevised	MT+CF:Revised	2.32
7	770M	CF:Unrevised	MT:Unrevised	6.24***
8	770M	CF:Unrevised	MT+CF:Unrevised	6.54***
9	770M	CF:Unrevised	MT+CF:Revised	4.34***
10	770M	MT:Unrevised	MT+CF:Unrevised	0.30
11	770M	MT:Unrevised	MT+CF:Revised	-1.90*
12	770M	MT+CF:Unrevised	MT+CF:Revised	-2.20**
Significance levels: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$, $p < 0.1$				

Small teacher models are sufficient to produce capable students. Beyond these findings, we can draw additional conclusions by relating our results to existing literature. Interestingly, we observed that the LLaMA2-13B teacher model produces a student model whose performance is comparable to that obtained from the substantially larger PaLM 540B model. This result challenges the common assumption that increasing teacher model size inherently improves student model performance, highlighting instead the effectiveness of the multitask training approach.

B. Effect on Explainability

Figure 4(b) presents the mean values for plausibility, understandability, completeness, satisfaction, contrastiveness, and quality, an average over all previous constructs for each study participant. Overall, Figure 4 shows that the model MT+CF:Revised produces explanations that are superior in terms of completeness, contrastiveness, and quality. To determine the significance of the differences, we determine an appropriate statistical analysis by assessing the characteristics of the study data: The normality of the six constructs is evaluated via QQ -plots indicating an approximately normal distribution; multivariate normality is tested with the Shapiro test ($p = 9.64 \times 10^{-19}$) revealing departure from normality; the uni-variate normality Shapiro test demonstrates a significant deviation from a normal distribution for all models and dimensions of explainability; and finally, the Levene test indicates the absence of significant differences in variance among the constructs, thereby suggesting homogeneity. Outliers have been retained in the dataset, as the box-plot method estimates an implausibly high number of outliers to be present in the data, likely due to the ordinal nature of the study data.

Based on these data characteristics, we perform two analyses: First, to find constructs with significant differences, a series of tests suitable for data of ordinal nature and absence of normality is conducted. Second, two regressions are performed to gain deeper insights into which method impacts student models' explainability most, with and without demographic data collected during the study. As the construct *quality* is the result of averaging the other constructs, it can be treated as a continuous variable [82] permitting the use of linear

TABLE II

KRUSKAL-WALLIS TEST RESULTS FOR THE FIVE DIMENSIONS OF EXPLAINABILITY. VARIABLES WITH $p < 0.05$ ARE BOLD.

Variable	n	p-value
Plausibility	1114	0.0987
Understandability	1114	0.1480
Completeness	1114	0.0081**
Satisfaction	1114	0.1660
Contrastiveness	1114	0.0004***
Significance levels: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$, $p < 0.1$		

TABLE III

REGRESSION RESULTS OF STUDENT MODELS ON THE FORMATIVE CONSTRUCT OF *quality*. VARIABLES WITH $p < 0.05$ ARE BOLD.

	Estimate	Std. Error
Intercept	4.037***	0.067
MT:Unrevised	0.138*	0.063
MT+CF:Unrevised	0.109	0.062
MT+CF:Revised	0.284***	0.071
Explanation length	0.000	0.000
Significance levels: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$, $p < 0.1$		

regressions. The procedures and results are described in the subsequent paragraphs.

We use the ANOVA-type test to analyze our data, showing that there is a significant difference between the constructs based on the type of student model ($p = 0.003$). Consequently, the rank-based Kruskal-Wallis test is used as a post-hoc test to investigate the differences further. As shown in the Table II, significant differences are found for completeness ($p = 0.0081$) and contrastiveness ($p = 0.0004$), justifying the use of Dunn's test [83] to estimate effect sizes. The test shows that the MT+CF:Revised student model significantly outperforms CF:Unrevised, MT:Unrevised, and MT+CF:Unrevised models in terms of contrastiveness, and also surpasses CF:Unrevised significantly in completeness (see Figure 4(b)). Pairwise Vargha and Delaney's A (VDA) effect size estimates indicate the largest differences between CF:Unrevised and MT+CF:Revised for completeness (VDA = 0.423) and contrastiveness (VDA = 0.404), although both represent small effect sizes [84]. The statistical tests support the initial observation that the MT+CF:Revised student is superior in providing complete and contrastive explanations.

Before performing the linear regressions, we verify several assumptions about the data. First, we remove 17 outliers identified via the box-plot method, which yields more reliable results for outlier detection now that the data are continuous. Second, we test the homogeneity of variance using the Levene test, which confirms variance homogeneity. Third, we assess normality through QQ -plots, which indicate that the data distribution approximates normality. Altogether, these observations justify the use of linear regression. Table III presents the regression results for student models predicting the formative construct of *quality*. The intercept is significant ($\beta = 4.037$, $p < 0.001$), indicating a baseline quality level for the CF:Unrevised model. Compared to this baseline, the MT:Unrevised condition shows a small but significant positive effect ($\beta = 0.138$, $p < 0.05$), suggesting a modest quality improvement. The MT+CF:Unrevised condition exhibits a

TABLE IV

REGRESSION RESULTS OF STUDENT MODELS ON THE FORMATIVE CONSTRUCT OF *quality*. ONLY CONTROL VARIABLES WITH P-VALUES < 0.1 ARE INCLUDED FOR CLARITY. BASELINE CATEGORIES: “MALE” FOR GENDER, “CIVIL SERVANT” FOR EMPLOYMENT STATUS, “UNITED KINGDOM” FOR COUNTRY, “A-LEVELS/INTERNATIONAL BACCALAUREATE/HIGHER EDUCATION ENTRANCE QUALIFICATION” FOR EDUCATION, AND “15 TO 19 YEARS OLD” FOR AGE.

	Estimate	Std. Error
Intercept	4.496***	0.470
MT:Unrevised	0.140*	0.059
MT+CF:Unrevised	0.115*	0.058
MT+CF:Revised	0.277***	0.067
Gender: Female	-0.254***	0.048
Country: No Answer	0.490*	0.258
Country: United States	-0.541***	0.090
Education: High School Diploma	-0.209*	0.085
Education: Junior High Diploma	-0.200*	0.116
Education: University Degree	-0.355***	0.071
Employment: Pupil/In School	-0.517*	0.312
Employment: Unemployed	0.412***	0.124
Significance levels: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$, $\cdot p < 0.1$		

marginally significant positive effect ($\beta = 0.109$, $p < 0.1$), indicating a slight potential increase in quality. Notably, the MT+CF:Revised condition demonstrates a significant positive effect ($\beta = 0.284$, $p < 0.001$), clearly highlighting that revisions substantially enhance the quality construct. We control for explanation length, which does not predict explanation quality ($\beta = 0.000$, $p > 0.1$), implying that the length of explanations has no meaningful relationship with perceived quality. These results indicate that counterfactual training alone does not enhance explainability beyond multitask training.

The regression analysis with control for demographic factors shows that various factors influence the participants’ perception of the quality of an explanation (see Table IV). However, qualitatively, these results yield the same conclusions as the regression without control variables, demonstrating the robustness of our findings.

Revising explanations with critique-revision prompting positively impacts the student model’s explainability. Our analysis indicates that explainability is significantly influenced by revising explanations through critique-revision prompting, which primarily affects contrastiveness and completeness. The difference in effect size between MT+CF:Unrevised and MT+CF:Revised (difference = $0.284 - 0.109 = 0.175$) can be attributed to data quality improvements from critique-revision prompting. A possible explanation for this improvement is that revised explanations, being more comprehensive and more differential, directly affect the student model’s ability to explain its answer choices (contrastiveness) and provide more comprehensive reasoning (completeness). However, the limited magnitude of these improvements suggests that the student model’s capacity to leverage enhanced explanations may be constrained, potentially by its small size and reasoning capabilities needed to make sense of the complexity introduced through longer, more detailed explanations. In conclusion, while critique-revision prompting did not enhance student model performance, it contributed to explainability. In scenarios where explainability is critical, all three methods, critique-

revision prompting, multitask, and counterfactual training, should be combined to produce the student with the best explanation quality.

Multitask training improves explainability in contrast to counterfactual training. Explainability was significantly improved by multitask training compared to counterfactual training, highlighted by the regression results (estimate = 0.138 for MT:Unrevised over the baseline, CF:Unrevised). This finding contradicts the argumentation of Wang et al. [34], who posited that counterfactual training enhances explanation factuality. The divergence in findings may stem from the differing evaluation methodologies: this study employed a human-grounded approach, while the previous work relied on functional evaluation via an LLM, highlighting the importance of rigorous assessments of explanation quality from humans.

Combining multitask and counterfactual training does not improve the explainability. In addition, the linear regression analysis shows no significant improvement when using a combination of multitask and counterfactual training without the data quality improvements from critique-revision prompting. Precisely, this can be seen in the regression results in Table III in the difference between the MT:Unrevised and MT+CF:Unrevised model (difference = $0.109 - 0.138 = -0.029$). Evidently, combining multitask and counterfactual training alone does not suffice to achieve this improvement, highlighting that critique-revision prompting primarily drives the observed benefits in explainability.

VI. CONCLUSION

In our work, we advance the field of knowledge distillation by (1) introducing additional methods for data generation and training, (2) developing a framework for comparing methods on both performance and explainability, and (3) comprehensively comparing the distillation methods with respect to performance *and* explainability. Our studies yield two main findings: first, regarding performance, multitask training provides a strong student model while maintaining robust levels of explainability. Second, integrating critique-revision prompting with both training methods improves the perceived quality of student model explanations, where the prompting mechanism contributes most to the improvement. These findings strengthen our ability to produce a smaller, more efficient model, promoting sustainability and facilitating deployment in scenarios with limited computational resources while maintaining performance and explainability.

However, this study is not without limitations, which offer the potential for future research. First, by combining multitask and counterfactual training, we did not improve the student models ability to answer and explain multiple-choice questions simultaneously. Thus, future studies could create novel techniques for enhancing training data or methods to boost performance and explainability in student models. Additionally, advanced evaluations of explainability, such as an application-grounded study, might be interesting. Furthermore, this study is limited in the dataset used to evaluate the model – other researchers should investigate the impact of the investigated methods on performance and explainability in different tasks,

use cases, and datasets. Lastly, it would be interesting to investigate the effect of the distillation methods on the student models' ability to generalize to unseen tasks.

REFERENCES

- [1] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS '20. Red Hook, NY, USA: Curran Associates Inc., 2020.
- [2] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [3] T. Le Scao and A. Rush, "How many data points is a prompt worth?" in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, 2021, pp. 2627–2636. [Online]. Available: <https://aclanthology.org/2021.naacl-main.208>
- [4] W. Jin, Y. Cheng, Y. Shen, W. Chen, and X. Ren, "A good prompt is worth millions of parameters? low-resource prompt-based learning for vision-language models," *CoRR*, vol. abs/2110.08484, 2021. [Online]. Available: <https://arxiv.org/abs/2110.08484>
- [5] A. Mitra, L. D. Corro, S. Mahajan, A. Codas, C. Simoes, S. Agarwal, X. Chen, A. Razdaibiedina, E. Jones, K. Aggarwal, H. Palangi, G. Zheng, C. Rosset, H. Khanpour, and A. Awadallah, "Orca 2: Teaching small language models how to reason," 2023.
- [6] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, ser. NIPS '22. Red Hook, NY, USA: Curran Associates Inc., 2022.
- [7] S. Yao, D. Yu, J. Zhao, I. Shafraan, T. L. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: deliberate problem solving with large language models," in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, ser. NIPS '23. Red Hook, NY, USA: Curran Associates Inc., 2023.
- [8] Y. Saxena, S. Chopra, and A. M. Tripathi, "Evaluating Consistency and Reasoning Capabilities of Large Language Models," Apr. 2024, arXiv:2404.16478 [cs]. [Online]. Available: <http://arxiv.org/abs/2404.16478>
- [9] M. Parmar, N. Patel, N. Varshney, M. Nakamura, M. Luo, S. Mashetty, A. Mitra, and C. Baral, "LogicBench: Towards Systematic Evaluation of Logical Reasoning Ability of Large Language Models," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 13 679–13 707. [Online]. Available: <https://aclanthology.org/2024.acl-long.739/>
- [10] Y. Hou, J. Li, Y. Fei, A. Stolfo, W. Zhou, G. Zeng, A. Bosselut, and M. Sachan, "Towards a Mechanistic Interpretation of Multi-Step Reasoning Capabilities of Language Models," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 4902–4919. [Online]. Available: <https://aclanthology.org/2023.emnlp-main.299/>
- [11] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao, "Large Language Models: A Survey," Feb. 2024, arXiv:2402.06196 [cs]. [Online]. Available: <http://arxiv.org/abs/2402.06196>
- [12] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes, Y. Tay, N. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. Hutchinson, R. Pope, J. Bradbury, J. Austin, M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat, S. Dev, H. Michalewski, X. Garcia, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito, D. Luan, H. Lim, B. Zoph, A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick, A. M. Dai, T. S. Pillai, M. Pellat, A. Lewkowycz, E. Moreira, R. Child, O. Polozov, K. Lee, Z. Zhou, X. Wang, B. Saeta, M. Diaz, O. Firat, M. Catasta, J. Wei, K. Meier-Hellstern, D. Eck, J. Dean, S. Petrov, and N. Fiedel, "PaLM: Scaling Language Modeling with Pathways," Oct. 2022, arXiv:2204.02311 [cs]. [Online]. Available: <http://arxiv.org/abs/2204.02311>
- [13] R. Tang, Y. Lu, L. Liu, L. Mou, O. Vechtomova, and J. Lin, "Distilling Task-Specific Knowledge from BERT into Simple Neural Networks," Mar. 2019, arXiv:1903.12136 [cs]. [Online]. Available: <http://arxiv.org/abs/1903.12136>
- [14] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, "TinyBERT: Distilling BERT for Natural Language Understanding," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, T. Cohn, Y. He, and Y. Liu, Eds. Online: Association for Computational Linguistics, Nov. 2020, pp. 4163–4174. [Online]. Available: <https://aclanthology.org/2020.findings-emnlp.372/>
- [15] X. Zhu, J. Li, Y. Liu, C. Ma, and W. Wang, "A survey on model compression for large language models," *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 1556–1577, 2024. [Online]. Available: <https://aclanthology.org/2024.tacl-1.85/>
- [16] S. Kim, G. Ham, Y. Cho, and D. Kim, "Robustness-Reinforced Knowledge Distillation With Correlation Distance and Network Pruning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 12, pp. 9163–9175, Dec. 2024, conference Name: IEEE Transactions on Knowledge and Data Engineering. [Online]. Available: <https://ieeexplore.ieee.org/document/10623281/?arnumber=10623281>
- [17] C.-Y. Hsieh, C.-L. Li, C.-K. Yeh, H. Nakhost, Y. Fujii, A. J. Ratner, R. Krishna, C.-Y. Lee, and T. Pfister, "Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes," *ArXiv*, vol. abs/2305.02301, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:258461606>
- [18] Q. Zhong, L. Ding, J. Liu, B. Du, and D. Tao, "PanDa: Prompt Transfer Meets Knowledge Distillation for Efficient Model Adaptation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 36, no. 9, pp. 4835–4848, Sep. 2024, conference Name: IEEE Transactions on Knowledge and Data Engineering. [Online]. Available: <https://ieeexplore.ieee.org/document/10475529/?arnumber=10475529>
- [19] M. Ballout, U. Krumnack, G. Heidemann, and K.-U. Kühnberger, "Efficient Knowledge Distillation: Empowering Small Language Models with Teacher Model Insights," Sep. 2024, arXiv:2409.12586 [cs]. [Online]. Available: <http://arxiv.org/abs/2409.12586>
- [20] C. Dai, K. Li, W. Zhou, and S. Hu, "Improve Student's Reasoning Generalizability through Cascading Decomposed CoTs Distillation," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds. Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 15 623–15 643. [Online]. Available: <https://aclanthology.org/2024.emnlp-main.875/>
- [21] —, "Beyond imitation: Learning key reasoning steps from dual chain-of-thoughts in reasoning distillation," *CoRR*, vol. abs/2405.19737, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2405.19737>
- [22] C. Wu, F. Wu, and Y. Huang, "One teacher is enough? pre-trained language model distillation from multiple teachers," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, C. Zong, F. Xia, W. Li, and R. Navigli, Eds. Online: Association for Computational Linguistics, Aug. 2021, pp. 4408–4413. [Online]. Available: <https://aclanthology.org/2021.findings-acl.387/>
- [23] X. Li, Y. Fang, M. Liu, Z. Ling, Z. Tu, and H. Su, "Distilling Large Vision-Language Model with Out-of-Distribution Generalizability," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. Los Alamitos, CA, USA: IEEE Computer Society, Oct. 2023, pp. 2492–2503. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.00236>
- [24] M. Heikkilä, "We're getting a better idea of AI's true carbon footprint," Nov. 2022. [Online]. Available: <https://www.technologyreview.com/2022/11/14/1063192/were-getting-a-better-idea-of-ais-true-carbon-footprint/>
- [25] E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for modern deep learning research," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 09, pp. 13 693–13 696, Apr. 2020. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/7123>
- [26] Y. Wang, H. Ivson, P. Dasigi, J. Hessel, T. Khot, K. R. Chandu, D. Wadden, K. MacMillan, N. A. Smith, I. Beltagy, and H. Hajishirzi, "How Far Can Camels Go? Exploring the State of Instruction Tuning on Open Resources," Oct. 2023, arXiv:2306.04751 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.04751>
- [27] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. C. Ferrer, A. Grattafiori, W. Xiong,

- A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve, “Code Llama: Open Foundation Models for Code,” 2023, version Number: 3. [Online]. Available: <https://arxiv.org/abs/2308.12950>
- [28] F. Doshi-Velez and B. Kim, “Towards a rigorous science of interpretable machine learning,” *arXiv: Machine Learning*, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:11319376>
- [29] A. Grobrugge, N. Mishra, J. Jakubik, and G. Satzger, “Explainability in AI Based Applications: A Framework for Comparing Different Techniques,” Oct. 2024, arXiv:2410.20873 [cs]. [Online]. Available: <http://arxiv.org/abs/2410.20873>
- [30] P. Spitzer, J. Holstein, P. Hemmer, M. Vössing, N. Kühl, D. Martin, and G. Satzger, “On the Effect of Contextual Information on Human Delegation Behavior in Human-AI collaboration,” Jan. 2024, arXiv:2401.04729 [cs]. [Online]. Available: <http://arxiv.org/abs/2401.04729>
- [31] J. Kim, A. Rohrbach, T. Darrell, J. Canny, and Z. Akata, “Textual explanations for self-driving vehicles,” in *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, pp. 577–593.
- [32] D. H. Park, L. A. Hendricks, Z. Akata, A. Rohrbach, B. Schiele, T. Darrell, and M. Rohrbach, “Multimodal explanations: Justifying decisions and pointing to the evidence,” *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8779–8788, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:3604848>
- [33] K. Morrison, P. Spitzer, V. Turri, M. Feng, N. Kühl, and A. Perer, “The Impact of Imperfect XAI on Human-AI Decision-Making,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 8, no. CSCW1, pp. 1–39, Apr. 2024, arXiv:2307.13566 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.13566>
- [34] P. Wang, Z. Wang, Z. Li, Y. Gao, B. Yin, and X. Ren, “SCOTT: Self-consistent chain-of-thought distillation,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 5546–5558. [Online]. Available: <https://aclanthology.org/2023.acl-long.304/>
- [35] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, C. Chen, C. Olsson, C. Olah, D. Hernandez, D. Drain, D. Ganguli, D. Li, E. Tran-Johnson, E. Perez, J. Kerr, J. Mueller, J. Ladish, J. Landau, K. Ndousse, K. Lukosuitis, L. Lovitt, M. Sellitto, N. Elhage, N. Schiefer, N. Mercado, N. DasSarma, R. Lasenby, R. Larson, S. Ringer, S. Johnston, S. Kravec, S. E. Showk, S. Fort, T. Lanham, T. Telleen-Lawton, T. Conerly, T. Henighan, T. Hume, S. R. Bowman, Z. Hatfield-Dodds, B. Mann, D. Amodei, N. Joseph, S. McCandlish, T. Brown, and J. Kaplan, “Constitutional AI: Harmlessness from AI Feedback,” Dec. 2022, arXiv:2212.08073 [cs] version: 1. [Online]. Available: <http://arxiv.org/abs/2212.08073>
- [36] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonnell, J. Phang, M. Pieler, U. S. Prashanth, S. Purohit, L. Reynolds, J. Tow, B. Wang, and S. Weinbach, “GPT-NeoX-20B: An open-source autoregressive language model,” in *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, A. Fan, S. Ilic, T. Wolf, and M. Gallé, Eds. virtual+Dublin: Association for Computational Linguistics, May 2022, pp. 95–136. [Online]. Available: <https://aclanthology.org/2022.bigscience-1.9/>
- [37] G. Hinton, O. Vinyals, and J. Dean, “Distilling the Knowledge in a Neural Network,” Mar. 2015, arXiv:1503.02531 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1503.02531>
- [38] Y. Gu, L. Dong, F. Wei, and M. Huang, “MiniLLM: Knowledge distillation of large language models,” in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=5h0qf7IBZZ>
- [39] H. Yu, R. Li, Z. Zhang, S. Ye, Q. Liu, Z. Huang, and E. Chen, “ERDL: Efficient Retrieval Framework Based on Distillation from Large Language Models,” in *2024 International Conference on Computational Linguistics and Natural Language Processing (CLNLP)*, Jul. 2024, pp. 83–87. [Online]. Available: <https://ieeexplore.ieee.org/document/10743156>
- [40] H. Wang, Y. Li, Y. Wang, H. Hu, and M.-H. Yang, “Collaborative Distillation for Ultra-Resolution Universal Style Transfer,” Mar. 2020, arXiv:2003.08436 [cs, eess]. [Online]. Available: <http://arxiv.org/abs/2003.08436>
- [41] Y. Feng, H. Wang, D. T. Yi, and R. Hu, “Triplet Distillation for Deep Face Recognition,” May 2019, arXiv:1905.04457 [cs]. [Online]. Available: <http://arxiv.org/abs/1905.04457>
- [42] G. Chen, W. Choi, X. Yu, T. Han, and M. Chandraker, “Learning Efficient Object Detection Models with Knowledge Distillation,” in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017. [Online]. Available: https://papers.nips.cc/paper_files/paper/2017/hash/e1e32e235eee1f970470a3a6658dfdd5-Abstract.html
- [43] X. Xu, M. Li, C. Tao, T. Shen, R. Cheng, J. Li, C. Xu, D. Tao, and T. Zhou, “A Survey on Knowledge Distillation of Large Language Models,” Oct. 2024, arXiv:2402.13116 [cs]. [Online]. Available: <http://arxiv.org/abs/2402.13116>
- [44] P. Izsak, S. Guskin, and M. Wasserblat, “Training Compact Models for Low Resource Entity Tagging using Pre-trained Language Models,” in *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing - NeurIPS Edition (EMC2-NIPS)*, Dec. 2019, pp. 44–47. [Online]. Available: <https://ieeexplore.ieee.org/document/9463575?page=1.73>
- [45] R. He, L. Liu, H. Ye, Q. Tan, B. Ding, L. Cheng, J. Low, L. Bing, and L. Si, “On the effectiveness of adapter-based tuning for pretrained language model adaptation,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, 2021. [Online]. Available: <http://dx.doi.org/10.18653/v1/2021.acl-long.172>
- [46] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” 2021. [Online]. Available: <https://arxiv.org/abs/2106.09685>
- [47] S. Guo, S. Damani, and K.-h. Chang, “LoPT: Low-Rank Prompt Tuning for Parameter Efficient Language Models,” 2024, version Number: 1. [Online]. Available: <https://arxiv.org/abs/2406.19486>
- [48] B. Lester, R. Al-Rfou, and N. Constant, “The Power of Scale for Parameter-Efficient Prompt Tuning,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 3045–3059. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.243/>
- [49] T. Mei, Y. Zi, X. Cheng, Z. Gao, Q. Wang, and H. Yang, “Efficiency Optimization of Large-Scale Language Models Based on Deep Learning in Natural Language Processing Tasks,” in *2024 IEEE 2nd International Conference on Sensors, Electronics and Computer Engineering (ICSECE)*, Aug. 2024, pp. 1231–1237. [Online]. Available: <https://ieeexplore.ieee.org/document/10729518?arnumber=10729518>
- [50] A. Holzinger, G. Längs, H. Denk, K. Zatloukal, and H. Müller, “Causability and explainability of artificial intelligence in medicine,” *Wiley Interdisciplinary Reviews. Data Mining and Knowledge Discovery*, vol. 9, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:132372213>
- [51] N. Burkart and M. F. Huber, “A survey on the explainability of supervised machine learning,” *Journal of Artificial Intelligence Research*, vol. 70, pp. 245–317, 2021.
- [52] C. Xu, D. Guo, N. Duan, and J. McAuley, “Baize: An Open-Source Chat Model with Parameter-Efficient Tuning on Self-Chat Data,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 6268–6278. [Online]. Available: <https://aclanthology.org/2023.emnlp-main.385/>
- [53] L. C. Magister, J. Mallinson, J. Adamek, E. Malmi, and A. Severyn, “Teaching small language models to reason,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 1773–1781. [Online]. Available: <https://aclanthology.org/2023.acl-short.151/>
- [54] N. Ho, L. Schmid, and S.-Y. Yun, “Large Language Models Are Reasoning Teachers,” Jun. 2023, arXiv:2212.10071 [cs]. [Online]. Available: <http://arxiv.org/abs/2212.10071>
- [55] Y. Fu, H. Peng, L. Ou, A. Sabharwal, and T. Khot, “Specializing Smaller Language Models towards Multi-Step Reasoning,” Jan. 2023, arXiv:2301.12726 [cs]. [Online]. Available: <http://arxiv.org/abs/2301.12726>
- [56] Y. Ma, C. Fan, and H. Jiang, “Sci-CoT: Leveraging Large Language Models for Enhanced Knowledge Distillation in Small Models for Scientific QA,” in *2023 9th International Conference on Computer and Communications (ICCC)*, Dec. 2023, pp. 2394–2398, ISSN: 2837-7109. [Online]. Available: <https://ieeexplore.ieee.org/document/10507622>

- [57] S. Mukherjee, A. Mitra, G. Jawahar, S. Agarwal, H. Palangi, and A. Awadallah, "Orca: Progressive Learning from Complex Explanation Traces of GPT-4," Jun. 2023, arXiv:2306.02707 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.02707>
- [58] N. F. Rajani, B. McCann, C. Xiong, and R. Socher, "Explain Yourself! Leveraging Language Models for Commonsense Reasoning," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, 2019, pp. 4932–4942. [Online]. Available: <https://www.aclweb.org/anthology/P19-1487>
- [59] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=IPL1N1MMrw>
- [60] D. Kaushik, E. Hovy, and Z. C. Lipton, "Learning the Difference that Makes a Difference with Counterfactually-Augmented Data," Feb. 2020, arXiv:1909.12434 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1909.12434>
- [61] S. Wiegrefe, A. Marasović, and N. A. Smith, "Measuring Association Between Labels and Free-Text Rationales," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 10266–10284. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.804/>
- [62] A. Talmor, J. Herzig, N. Lourie, and J. Berant, "CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge," in *Proceedings of the 2019 Conference of the North American Association for Computational Linguistics*, 2019, pp. 4149–4158. [Online]. Available: <http://aclweb.org/anthology/N19-1421>
- [63] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. Canton Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. Singh Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom, "Llama 2: Open Foundation and Fine-Tuned Chat Models," *arXiv e-prints*, p. arXiv:2307.09288, Jul. 2023.
- [64] K. Yang, D. Klein, A. Celikyilmaz, N. Peng, and Y. Tian, "RLCD: Reinforcement learning from contrastive distillation for LM alignment," in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=v3XXtxWKi6>
- [65] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, no. 1, Jan. 2020.
- [66] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. Rush, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Q. Liu and D. Schlangen, Eds. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. [Online]. Available: <https://aclanthology.org/2020.emnlp-demos.6/>
- [67] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell, "Generating visual explanations," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 3–19.
- [68] Prolific, "Prolific · Quickly find research participants you can trust." 2023. [Online]. Available: <https://www.prolific.com/>
- [69] Y. Xie, S. Vosoughi, and S. Hassanpour, "Interpretation quality score for measuring the quality of interpretability methods," *ArXiv*, vol. abs/2205.12254, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:249018098>
- [70] N. Alangari, M. El Bachir Menai, H. Mathkour, and I. Almosallam, "Exploring Evaluation Methods for Interpretable Machine Learning: A Survey," *Information*, vol. 14, no. 8, p. 469, Aug. 2023. [Online]. Available: <https://www.mdpi.com/2078-2489/14/8/469>
- [71] J. Zhou, A. H. Gandomi, F. Chen, and A. Holzinger, "Evaluating the Quality of Machine Learning Explanations: A Survey on Methods and Metrics," *Electronics*, vol. 10, no. 5, p. 593, Mar. 2021. [Online]. Available: <https://www.mdpi.com/2079-9292/10/5/593>
- [72] H. Chen, F. Brahman, X. Ren, Y. Ji, Y. Choi, and S. Swayamdipta, "REV: Information-theoretic evaluation of free-text rationales," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 2007–2030. [Online]. Available: <https://aclanthology.org/2023.acl-long.112/>
- [73] O. Dictionary, "Oxford Learner's Dictionaries | Find definitions, translations, and grammar explanations at Oxford Learner's Dictionaries," 2024. [Online]. Available: <https://www.oxfordlearnersdictionaries.com/>
- [74] R. R. Hoffman, S. T. Mueller, G. Klein, and J. Litman, "Measures for explainable AI: Explanation goodness, user satisfaction, mental models, curiosity, trust, and human-AI performance," *Frontiers in Computer Science*, vol. 5, p. 1096257, Feb. 2023. [Online]. Available: <https://www.frontiersin.org/articles/10.3389/fcomp.2023.1096257/full>
- [75] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, "Machine Learning Interpretability: A Survey on Methods and Metrics," *Electronics*, vol. 8, no. 8, p. 832, Jul. 2019. [Online]. Available: <https://www.mdpi.com/2079-9292/8/8/832>
- [76] P. Hase, S. Zhang, H. Xie, and M. Bansal, "Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language?" in *Findings of the Association for Computational Linguistics: EMNLP 2020*, T. Cohn, Y. He, and Y. Liu, Eds. Online: Association for Computational Linguistics, Nov. 2020, pp. 4351–4367. [Online]. Available: <https://aclanthology.org/2020.findings-emnlp.390/>
- [77] A. Holzinger, A. Carrington, and H. Müller, "Measuring the Quality of Explanations: The System Causability Scale (SCS): Comparing Human and Machine Explanations," *KI - Künstliche Intelligenz*, vol. 34, no. 2, pp. 193–198, Jun. 2020. [Online]. Available: <http://link.springer.com/10.1007/s13218-020-00636-z>
- [78] K. Hebenstreit, R. Praas, and M. Samwald, "A collection of principles for guiding and evaluating large language models," 2023. [Online]. Available: <https://arxiv.org/abs/2312.10059>
- [79] H. O. Mayer, *Interview und schriftliche Befragung: Entwicklung, Durchführung und Auswertung*, 4th ed., ser. 150 Jahre Wissen für die Zukunft. München Wien: Oldenbourg, 2008.
- [80] A. Bhattacharjee, *Social Science Research: Principles, Methods and Practices*. Open Textbook Library, 01 2012.
- [81] A. Páez and G. Boisjoly, "Exploratory Data Analysis," in *Discrete Choice Analysis with R*. Cham: Springer International Publishing, 2022, pp. 25–64, series Title: Use R! [Online]. Available: https://link.springer.com/10.1007/978-3-031-20719-8_2
- [82] G. Norman, "Likert scales, levels of measurement and the "laws" of statistics," *Advances in Health Sciences Education*, vol. 15, no. 5, pp. 625–632, Dec. 2010. [Online]. Available: <http://link.springer.com/10.1007/s10459-010-9222-y>
- [83] A. Dinno, "Nonparametric Pairwise Multiple Comparisons in Independent Groups using Dunn's Test," *The Stata Journal: Promoting communications on statistics and Stata*, vol. 15, no. 1, pp. 292–300, Apr. 2015. [Online]. Available: <http://journals.sagepub.com/doi/10.1177/1536867X1501500117>
- [84] A. Vargha, H. D. Delaney, and A. Vargha, "A Critique and Improvement of the "CL" Common Language Effect Size Statistics of McGraw and Wong," *Journal of Educational and Behavioral Statistics*, vol. 25, no. 2, p. 101, 2000. [Online]. Available: <http://links.jstor.org/sici?sici=1076-9986%2820002%2925%3A2%3C101%3AACA1OT%3E2.0.CO%3B2-O&origin=crossref>