

GPTScore: Evaluate as You Desire

Jinlan Fu¹, See-Kiong Ng¹, Zhengbao Jiang², Pengfei Liu³,
¹National University of Singapore, ²Carnegie Mellon University
³Shanghai Jiao Tong University
jinlanjonna@gmail.com, pengfei@sjtu.edu.cn

Abstract

Generative Artificial Intelligence (AI) has enabled the development of sophisticated models that are capable of producing high-caliber text, images, and other outputs through the utilization of large pre-trained models. Nevertheless, assessing the quality of the generation is an even more arduous task than the generation itself, and this issue has not been given adequate consideration recently. This paper proposes a novel evaluation framework, GPTSCORE, which utilizes the emergent abilities (e.g., in-context learning, zero-shot instruction) of generative pre-trained models to score generated texts. There are 19 pre-trained models explored in this paper, ranging in size from 80M (e.g., Flan-T5-small) to 175B (e.g., GPT3). Experimental results on four text generation tasks, 22 evaluation aspects, and corresponding 37 datasets demonstrate that this approach can effectively allow us to achieve what one desires to evaluate for texts simply by natural language instructions. This nature helps us overcome several long-standing challenges in text evaluation—how to achieve customized, multi-faceted evaluation without model training. We make our code publicly available.¹

1 Introduction

The advent of generative pre-trained models, such as GPT3 (Brown et al., 2020), has precipitated a shift from *analytical* AI to *generative* AI across multiple domains (Sequoia, 2022). Take text as an example: the use of a large pre-trained model with appropriate prompts (Liu et al., 2021) has achieved superior performance in tasks defined both in academia (Sanh et al., 2021) and scenarios from the real world (Ouyang et al., 2022). While text generation technology is advancing rapidly, techniques for evaluating the quality of these texts lag far behind. This is especially evident in the following ways:

¹<https://github.com/jinlanfu/GPTScore>

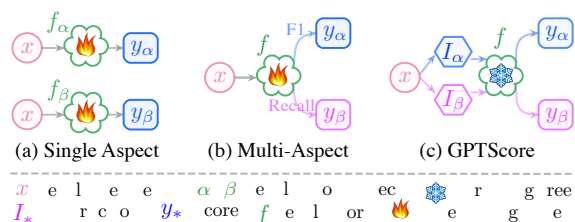


Figure 1: An overview of text evaluation approaches.

(1) Existing studies (Ghazarian et al., 2022; Ye et al., 2021) evaluate text quality with limited aspects (e.g., fluency) (Fig. 1-(a)), which are usually customized prohibitively, making it harder for users to evaluate aspects *as they need* (Freitag et al., 2021). (2) A handful of studies (Yuan et al., 2021; Scialom et al., 2021; Zhong et al., 2022; Li et al., 2021; Mehri and Eskénazi, 2020a) have examined multi-aspect evaluation but lack carefully studied aspects’ definitions and their relationship. Moreover, the specific aspect evaluations are empirically bound with metric variants (Fig. 1-(b)). (3) Rely on annotated samples and model training. Most of the above methods necessitate complicated supervised training or costly manual annotation of samples (Fig. 1-(a,b)). This makes these methods hard to use in industrial settings and adapt to new evaluation aspects required by users.

In this paper, we demonstrated the talent of the super large pre-trained language model (e.g., GPT-3) in achieving multi-aspect, customized, and training-free evaluation (Fig. 1-(c)). Essentially, it skillfully utilizes the pre-trained model’s zero-shot instruction (Chung et al., 2022) and in-context learning (Brown et al., 2020; Min et al., 2022) ability to deal with complex and ever-changing evaluation needs while solving multiple evaluation challenges that have plagued many years. Specifically, given a text generated from a specific context (e.g., source text in text summarization) and a desirable evaluation aspect (e.g., fluency), the high-level idea of the proposed framework is that

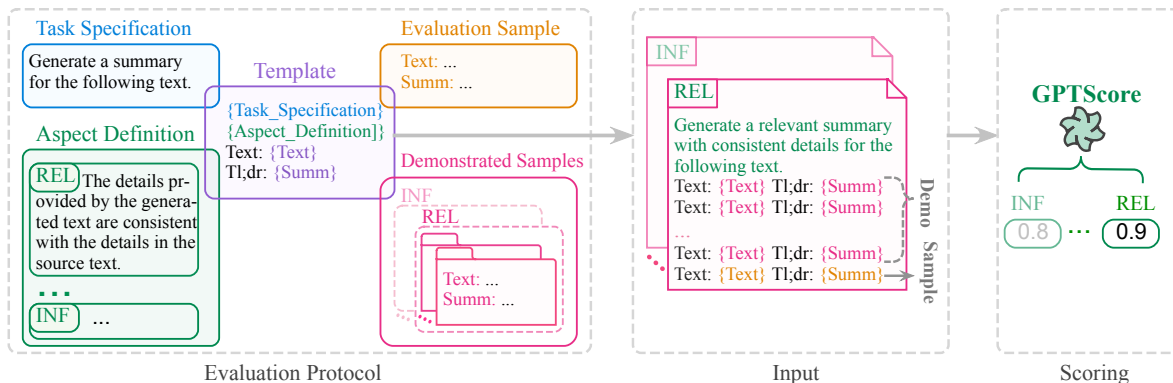


Figure 2: The framework of GPTSCORE. We include two evaluation aspects *relevance* (*REL*) and *informative* (*INF*) in this figure and use the evaluation of *relevance* (*REL*) of the text summarization task to exemplify our framework.

the higher-quality text for a specific aspect will be more likely generated than unqualified ones, where the “likely” can be measured by the conditional generation probability.

How to perform an evaluation as the user desires? As illustrated in Fig. 2, to capture users’ true desires, an *evaluation protocol*² will be initially established based on (a) the *task specification*, which typically outlines how the text is generated (e.g., generate a response for a human based on the conversation); (b) *aspect definition* that documents the details of desirable evaluation aspects (e.g., the response should be intuitive to understand); (c) *demonstrated samples*: a handful of well-labeled samples are required to teach the model which sample is qualified. Subsequently, each evaluation sample will be presented with the evaluated protocol with optionally moderate exemplar samples, which could facilitate the model’s learning. Lastly, a generative pre-trained model will be used to calculate how likely the text could be generated based on the above evaluation protocol, thus giving rise to our model’s name: GPTSCORE. Given the plethora of pre-trained models, we instantiate our framework with different backbones: GPT2 (Radford et al., 2019), OPT (Zhang et al., 2022b), Flan-T5 (Chung et al., 2022), and GPT3 (instruction-based (Ouyang et al., 2022)) due to their superior capacity for *zero-shot instruction* and their aptitude for *in-context learning*.

Experimentally, we ran through almost all common natural language generation tasks in NLP, and the results showed the power of this new paradigm. The main observations are listed as follows: (1) GPTScore performs better when instructed by the

²To better understand how to design the evaluation protocols, we give all the evaluation protocols for the different tasks and aspects studied in this work in the Appendix F.

definition of task and aspect. Furthermore, incorporating suitable exemplified samples with in-context learning will further enhance the process. (2) Different evaluation aspects exhibit certain correlations. By incorporating definitions with other highly correlated aspects (e.g., interesting and engaging), the performance of the smaller model (GPT3-curie, 6.7B) can surpass the larger model (GPT3-davinci, 175B). (3) The GPTscore performs better than fine-tuned models across tasks such as text summarization, data-to-text, and dialogue response generation. (4) The performance of *GPT3-text-davinci-003*, which is tuned based on human feedback, is inferior to *GPT3-text-davinci-001* in the majority of the evaluation settings.

Our main **contributions** in this paper are:

- (1) We propose a newly generated text scoring framework, GPTScore, which utilizes the emergent ability of large language model to achieve multi-aspect, customized, and training-free evaluation.
- (2) We comprehensively explore GPTScore, studying 19 language models (ranging in size from 80M to 175B) and four popular text generation tasks. Experiments demonstrate that training-free GPTScore outperforms fine-tuning model and achieves higher human correspondence.
- (3) We design and demonstrate the feasibility of custom evaluation for a new aspect by the GPTScore framework with training-free.
- (4) We summarize some observations of the opaque GPT3 family and other backbone models in the evaluation and try to give explanations.

2 Related Work

Similarity-based Metrics measures the similarity between the generated text and the reference text. It includes two types: (1) lexical overlap-based

metrics, e.g., BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004); (2) embedding-based metrics, e.g., BERTScore (Zhang et al., 2020) and MoverScore (Zhao et al., 2019).

Single-aspect Evaluator refers to evaluators designed to evaluate the quality of a specific aspect or overall of the generated text. For example, DEAM (Ghazarian et al., 2022) and QuantiDCE (Ye et al., 2021) were proposed for the evaluation of the *coherence* of the dialogue system; several evaluators (Cao et al., 2020; Durmus et al., 2020; Wang et al., 2020a) are designed for the evaluation of the *consistency* of text summarization.

Multi-aspect Evaluator refers to one evaluator handle several evaluation aspects by using different input and output text pair (Yuan et al., 2021), different prompt designed by the aspect name (Zhong et al., 2022; Mehri and Eskénazi, 2020a), (Mehri and Eskénazi, 2020b), different formulas (Scialom et al., 2021). Unlike (Zhong et al., 2022; Mehri and Eskénazi, 2020a) which only consider the vague aspect description, we fully considered exhaustive aspect definition and their relationship.

Emergent Ability Recent works have revealed various emergent abilities of generative pre-trained language models, such as, in-context learning (Min et al., 2022), chain-of-thought reasoning (Wei et al., 2022), and zero-shot instruction (Ouyang et al., 2022). These abilities allow large language models to achieve good performance without training.

What is the difference between our work and BARTScore? Yuan et al. (2021) demonstrates the feasibility of using the probability of text generation as a text quality score, and fine-tuning is required to achieve better performance. However, the model fine-tuning cost a lot, and it is hard for us to fine-tune a scoring model for each task and each domain. In this work, we focus on proposing a new framework that allows the generated text evaluation to achieve customizable, multi-faceted, and train-free evaluation. To achieve this target, (1) we utilized the emergent ability of language models, such as in-context learning, zero-shot instruction et al., to build the GPTScore. (2) By studying the GPTScore framework on 19 language models covering four backbones, we demonstrate that (a) GPTScore outperforms the fine-tuned BARTScore; (b) GPTScore can be customized for a new evaluation aspect with a labeled handful of samples while BARTScore cannot do this.

3 Generative Pretraining Score (GPTScore)

The core idea of GPTSCORE is that a generative pre-training model will assign a higher probability of high-quality generated text following a given instruction and context. In our method, the instruction is composed of the task description d and the aspect definition a . Specifically, suppose that the text to be evaluated is $\mathbf{h} = \{h_1, h_2, \dots, h_m\}$, the context information is \mathcal{S} (e.g., source text or reference text), then GPTSCORE is defined as the following conditional probability:

$$\text{GPTScore}(\mathbf{h}|d, a, \mathcal{S}) = \sum_{t=1}^m w_t \log p(h_t|\mathbf{h}_{<t}, T(d, a, \mathcal{S}), \theta),$$

where w_t is the weight of the token at position t . In our work, we treat each token equally. $T(\cdot)$ is a prompt template that defines the evaluation protocol, which is usually task-dependent and specified manually through prompt engineering.

Few-shot with Demonstration The generative pre-trained language model can better perform tasks when prefixed with a few annotated samples (i.e., demonstrations). Our proposed framework is flexible in supporting this by extending the prompt template T with demonstrations.

Choice of Prompt Template Prompt (Liu et al., 2021; Fu et al., 2022) templates define how task description, aspect definition, and context are organized. In this work, for the GPT3-based model, we opt for prompts that are officially provided by OpenAI.³ For instruction-based pre-trained models, we use prompts from NaturalInstruction (Wang et al., 2022) since it’s the main training source for those instruction-based pre-train models. Taking the evaluation of the fluency of the text summarization task as an example, based on the prompt provided by OpenAI,⁴ the task prompt is “{Text} T1;dr {Summary}”, the definition of fluency is “Is the generated text well-written and grammatical?” (in Tab. 1), and then the final prompt template is “Generate a fluent and grammatical summary for the following text: {Text} T1;dr {Summary}”, where demonstrations could be introduced by repeating instantiating “{Text} T1;dr {Summary}” In Appendix F, we list the prompts for various aspects of all tasks studied in this work and

³<https://beta.openai.com/examples>

⁴<https://beta.openai.com/examples/default-tldr-summary>

Aspect	Task	Definition
Semantic Coverage (COV)	Summ	How many semantic content units from the reference text are covered by the generated text?
Factuality (FAC)	Summ	Does the generated text preserve the factual statements of the source text?
Consistency (CON)	Summ, Diag	Is the generated text consistent in the information it provides?
Informativeness (INF)	Summ, D2T, Diag	How well does the generated text capture the key ideas of its source text?
Coherence (COH)	Summ, Diag	How much does the generated text make sense?
Relevance (REL)	Diag, Summ, D2T	How well is the generated text relevant to its source text?
Fluency (FLU)	Diag, Summ, D2T, MT	Is the generated text well-written and grammatical?
Accuracy (ACC)	MT	Are there inaccuracies, missing, or unfactual content in the generated text?
MQM	MT	How is the overall quality of the generated text?
Interest (INT)	Diag	Is the generated text interesting?
Engagement (ENG)	Diag	Is the generated text engaging?
Specific (SPE)	Diag	Is the generated text generic or specific to the source text?
Correctness (COR)	Diag	Is the generated text correct or was there a misunderstanding of the source text?
Semantically appropriate (SEM)	Diag	Is the generated text semantically appropriate?
Understandability (UND)	Diag	Is the generated text understandable?
Error Recovery (ERR)	Diag	Is the system able to recover from errors that it makes?
Diversity (DIV)	Diag	Is there diversity in the system responses?
Depth (DEP)	Diag	Does the system discuss topics in depth?
Likeability (LIK)	Diag	Does the system display a likeable personality?
Flexibility (FLE)	Diag	Is the system flexible and adaptable to the user and their interests?
Inquisitiveness (INQ)	Diag	Is the system inquisitive throughout the conversation?

Table 1: The definition of aspects evaluated in this work. *Semantic App.* denotes *semantically appropriate* aspect. *Diag*, *Summ*, *D2T*, and *MT* denote the *dialogue response generation*, *text summarization*, *data to text* and *machine translation*, respectively. “MQM” is the short name of *Multidimensional Quality Metrics*.

leave a more comprehensive exploration on prompt engineering as a future work.

Selection of Scoring Dimension GPTSCORE exhibits different variants in terms of diverse choices of texts being calculated. For example, given a generated hypothesis, we can calculate GPTSCORE either based on the source text (i.e., $src \rightarrow hypo, p(hypo|src)$) or based on the gold reference (i.e., $ref \rightarrow hypo, p(hypo|ref)$). In this paper, the criteria for choosing GPTSCORE variants are mainly designed to align the protocol of human judgments (Liu et al., 2022) that are used to evaluate the reliability of automated metrics. We will detail this based on different human judgment datasets in the experiment section.

4 Experimental Settings

4.1 Meta Evaluation

Meta evaluation aims to evaluate the reliability of automated metrics by calculating how well automated scores (y_{auto}) correlate with human judgment (y_{human}) using correlation functions $g(y_{auto}, y_{human})$ such as spearman correlation. In this work, we adopt two widely-used correlation measures: (1) **Spearman** correlation (ρ) (Zar, 2005) and (2) **Pearson** correlation (r) (Mukaka, 2012).

4.2 Tasks, Datasets, and Aspects

To achieve a comprehensive evaluation, in this paper, we cover a broad range of natural language generation tasks: *Dialogue Response Generation*, *Text Summarization*, *Data-to-Text*, and *Machine*

Translation, which involves 37 datasets and 22 evaluation aspects in total. Tab. 8 summarizes the tasks, datasets, and evaluation aspects considered by each dataset. The definition of different aspects can be found in Tab. 1. More detailed illustrations about the datasets can be found in Appendix D.

(1) **Dialogue Response Generation** aims to generate an engaging and informative response based on the dialogue history. We adopt the FED (Mehri and Eskénazi, 2020a) datasets and consider both turn-level and dialogue-level evaluations. (2) **Text Summarization** is a task of automatically generating informative summary for a given long text. We adopt SummEval (Bhandari et al., 2020), REALSumm (Bhandari et al., 2020), NEWSROOM (Grusky et al., 2018), and QAGS_XSUM (Wang et al., 2020b) datasets. (3) **Data-to-Text** aims to generate a fluent and factual description for a given table. We consider BAGEL (Mairesse et al., 2010) and SFRES (Wen et al., 2015) datasets. (4) **Machine Translation** aims to translate a sentence from one language to another. We consider a subdatasets of Multidimensional Quality Metrics (MQM) (Freitag et al., 2021), namely, MQM-2020 (Chinese->English).

4.3 Scoring Models

ROUGE (Lin, 2004) is a popular automatic generation evaluation metric. We consider three variants ROUGE-1, ROUGE-2, and ROUGE-L. **PRISM** (Thompson and Post, 2020) is a reference-based evaluation method designed for machine translation with pre-trained paraphrase systems.

BERTScore (Zhang et al., 2020) uses contextual representation from BERT to calculate the similarity between the generated text and the reference text. **MoverScore** (Zhao et al., 2019) considers both contextual representation and Word Mover’s Distance (Kusner et al., 2015). **DynaEval** (Zhang et al., 2021) is a unified automatic evaluation framework for dialogue response generation tasks on the turn level and dialogue level. **BARTScore** (Yuan et al., 2021) is a text-scoring model based on BART (Lewis et al., 2020) without fine-tuning. **BARTScore+CNN** and **BARTScore+CNN+Para** are the variants of **BARTScore**, the former is fine-tuned on the CNNDM dataset (Hermann et al., 2015), and the latter is fine-tuned on CNNDM and Paraphrase2.0 (Hu et al., 2019). **GPTSCORE** is our evaluation method, designed based on 19 pre-trained language models, covering GPT3, OPT, Flan-T5, and GPT2 backbones. Tab. 2 shows model variants used in this paper and their number of parameters.

GPT3	Param.	OPT	Param.
GPT3-a01 (text-ada-001)	350M	OPT350M	350M
GPT3-b01 (text-babbage-001)	1.3B	OPT-1.3B	1.3B
GPT3-c01 (text-curie-001)	6.7B	OPT-6.7B	6.7B
GPT3-d01 (text-davinci-001)	175B	OPT-13B	13B
GPT3-d03 (text-davinci-003)	175B	OPT-66B	66B
Flan-T5	Param.	GPT2	Param.
FT5-small	80M	GPT2-M	355M
FT5-base	250M	GPT2-L	774M
FT5-L	770M	GPT2-XL	1.5B
FT5-XL	3B	GPT-J-6B	6B
FT5-XXL	11B		

Table 2: A summary of pre-trained language models studied in this work. *Param.* denotes *Parameter*.

4.4 Scoring Dimension

Specifically, (1) For aspects INT, ENG, SPC, REL, COR, SEM, UND, and FLU of FED-Turn datasets from the open domain dialogue generation task, we choose the *src->hypo* variant since the human judgments of the evaluated dataset (i.e., FED-Turn) are also created based on the source. (2) For aspects COH, CON, and INF from SummEval and Newsroom, since data annotators labeled the data based on source and hypothesis texts, we choose *src->hypo* for these aspects. (3) For aspects INF, NAT, and FLU from the data-to-text task, we choose *ref->hypo*. Because the source text of the data-to-text task is not in the standard text format, which will be hard to handle by the scoring function.

4.5 Evaluation Dataset Construction

Unlike previous works (Matiana et al., 2021; Xu et al., 2022a,b; Castricato et al., 2022) that only consider the overall text quality, we focus on evaluating multi-dimensional text quality. In this work, we studied 37 datasets according to 22 evaluation aspects. Since each sample needs to evaluate the generated text of dozens of systems, to reduce the API cost of GPT3, we randomly sample 40 samples for each text summarization dataset and 100 samples for each dialogue response generation and data-to-text dataset. For example, in the Newsroom dataset with 60 samples, 40 samples (accounting for 60% of the samples) are randomly selected to construct the evaluation set.

5 Experiment Results

In this work, we focus on exploring whether language models with different structures and sizes can work in the following three scenarios. (a) **vanilla (VAL)**: with non-instruction and non-demonstration; (b) **instruction (IST)**: with instruction and non-demonstration; (c) **instruction+demonstration (IDM)**: with instruction and demonstration. We studied four text generation tasks introduced in Sec. 4.2. Due to the limited space, we moved the results and analysis of the machine translation task into the Appendix A.

Significance Tests To examine the reliability and validity of the experiment results, we conducted the significance test based on bootstrapping.⁵ Our significance test is to check (1) whether the performance of IST (IDM) is significantly better than VAL, and values achieved with the IST (IDM) settings will be marked † if it passes the significant test (p-value <0.05). (2) whether the performance of IDM is significantly better than IST, if yes, mark the value with IDM setting with ‡.

Average Performance Due to space limitations, we keep the average performance of GPT3-, GPT2-, OPT-, and FT5-based models. The full results of various variants can be found in Appendix G.

5.1 Text Summarization

The evaluation results of 28 (9 baseline models (e.g., ROUGE-1) and 19 variants of GPTScore (e.g., GPT3-d01) scoring functions for the text summarization task on SummEval and RealSumm datasets are shown in Tab. 3. Due to the space limitation,

⁵[https://en.wikipedia.org/wiki/Bootstrapping_\(statistics\)](https://en.wikipedia.org/wiki/Bootstrapping_(statistics))

Model	CON		FLU		REL		COH	
	VAL	IST	VAL	IST	VAL	IST	VAL	IST
ROUGE-1	20.8	-	14.8	-	26.2	-	14.1	-
ROUGE-2	17.2	-	12.0	-	17.4	-	9.1	-
ROUGE-L	19.8	-	17.6	-	24.7	-	12.9	-
BERTSc	19.7	-	23.7	-	34.7	-	25.9	-
MoverSc	18.0	-	15.7	-	24.8	-	11.5	-
PRISM	29.9	-	26.1	-	25.2	-	26.5	-
BARTSc	30.8	-	24.6	-	28.9	-	29.7	-
+CNN	35.8	-	38.1	-	35.9	-	42.5	-
+CNN+Pa	37.0	-	40.5	-	33.9	-	42.5	-
<hr/>								
GPT3-a01	39.7	40.5 [†]	36.1	35.9	28.2	27.6	39.3	39.8 [†]
GPT3-b01	41.0	41.4 [†]	37.1	39.1 [†]	32.0	33.4 [†]	42.7	45.2[†]
GPT3-c01	44.6	45.1 [†]	38.9	39.5 [†]	31.6	33.2 [†]	41.3	40.8
GPT3-d01	46.6	47.5[†]	40.5	41.0[†]	32.4	34.3 [†]	40.0	40.1
GPT3-d03	45.2	44.9	41.1	40.3	36.3	38.1[†]	43.7	43.4
<hr/>								
GPT2-M	34.6	35.3 [†]	28.1	30.7 [†]	28.3	28.3	36.0	39.2 [†]
GPT2-L	33.7	34.4 [†]	29.4	31.5 [†]	27.8	28.1 [†]	36.4	39.8 [†]
GPT2-XL	35.9	36.1 [†]	31.2	33.1 [†]	28.1	28.0	35.3	39.9[†]
GPT-J-6B	42.7	42.8[†]	35.5	37.4[†]	31.5	31.9[†]	35.5	39.5 [†]
<hr/>								
OPT350m	34.9	35.5 [†]	29.6	31.4 [†]	29.5	28.6	33.4	37.6 [†]
OPT-1.3B	40.0	42.0 [†]	33.6	35.9 [†]	33.5	34.2 [†]	35.0	37.8[†]
OPT-6.7B	42.1	45.7[†]	35.5	37.6 [†]	35.4	35.4	35.7	36.8 [†]
OPT-13B	42.5	45.2 [†]	35.6	37.3 [†]	33.6	33.9	33.5	34.7 [†]
OPT-66B	44.0	45.3 [†]	36.3	38.0[†]	33.4	33.7 [†]	32.0	35.9 [†]
<hr/>								
FT5-small	37.0	38.0 [†]	35.6	34.7	27.3	28.0 [†]	35.0	35.4 [†]
FT5-base	36.7	37.2 [†]	37.3	36.5	29.5	31.2 [†]	39.2	39.9 [†]
FT5-L	41.0	42.5 [†]	39.3	41.6 [†]	31.2	35.3[†]	42.3	45.1 [†]
FT5-XL	41.0	43.6 [†]	39.7	42.1 [†]	31.4	34.4 [†]	42.8	47.0[†]
FT5-XXL	43.7	43.8	39.8	42.4[†]	32.8	34.3 [†]	42.1	45.6 [†]
<hr/>								
Avg.	40.4	41.4	35.8	37.2	31.3	32.2	38.0	40.2

Table 3: Spearman correlation of different aspects on SummEval dataset. VAL and IST are the abbreviations of vanilla and instruction, respectively. Values with [†] denote the evaluator with instruction significantly outperforms vanilla. Values in bold are the best performance in a set of variants (e.g., GPT3 family).

we move the results of the REALSumm, NEWSROOM, and QXSUM datasets to the Appendix G. Fig. 3 shows the evaluation results of five GPT3 variant models on four text summarization datasets, where QXSUM uses the Pearson correlation and other datasets use the Spearman correlation metric. The main observations are summarized as follows:

(1) **Evaluator with instruction significantly improves the performance.** For the 4 aspects of SummEval datasets, 19 instruction-enhanced variants of GPTScore significantly outperform models without instruction (values with [†] in Tab. 3) and almost 7 unsupervised baseline methods. (2) **Most GPT3- and FT5-based models equipped with instructions outperform supervised methods.** For example, equipped with instructions, FT5-L, FT5-XL, and FT5-XXL significantly outperform the supervised model BARTSc+CNN+Pa for all

four aspects of the SummEval dataset. (3) As for the GPT3-based models, (a) **the performance of GPT3-d01 is barely significantly better than GPT3-c01**, which tries to balance power and speed. (b) GPT3-d03 performs better than GPT3-d01 significantly. Both conclusions have passed the significance test at $p < 0.05$ and can be seen in Fig. 3.

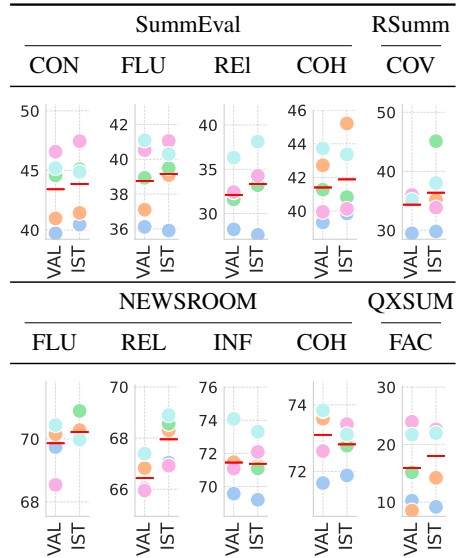


Figure 3: Experimental results for GPT3-based variants in text summarization task. Here, blue, orange, green, pink, and cyan dot denote that GPTSCORE is built based on GPT3-a01 (●), GPT3-b01 (●), GPT3-c01 (●), GPT3-d01 (●), and GPT3-d03 (●), respectively. The red lines (—) denote the average performance of GPT3-based variants.

5.2 Data to Text

We consider the BAGEL and SFRES datasets for the evaluation of data to text task. The average Spearman correlations of the GPT3-based, GPT2-based, OPT-based, and FT5-based models are listed in Tab. 4. VAL, IST, and IDM denote the vanilla, using instruction, and using both instruction and demonstration settings, respectively. Due to the space limitations, detailed results for each evaluator are moved to Appendix G (Tab. 15 and Tab. 16). The main observations are listed as follows:

(1) **Instruction (IST) improves performance, and combining it with demonstrations (IDM) further enhances it.** In Tab. 4, the average performance on the three aspects is significantly improved when adapting to the instruction, and the performance of using demonstration on NAT and FLU has further significantly improved. (2) **Many GPTScore variants enhanced by instruction and demonstration (IDM) outperform the fine-tuned**

model, namely BARTSCORE+CNN+Para. For example, in Tab. 16 regarding NAT and FLU of the SFRES dataset, most of the 19 variants of GPTScore with instructions and demonstrations outperform fine-tuned BARTSCORE+CNN+Para. (3) **The choice of samples for demonstration impacts the evaluation performance a lot.** On the BAGEL and SFRES datasets, when equipped with instruction and demonstration (IDM), the average performance of the four backbones performs much worse than backbone equipped with instruction (IST) only. (4) **Equipped with instruction and demonstration, the performance of a GPT3 family model with a small model size can surpass that of large models.** In Fig. 4, the performance of GPT3-c01 with IDM always outperforms GPT3-d03, which holds for both datasets.

Model	INF		NAT		FLU				
	VAL	IDM	VAL	IDM	VAL	IDM			
BAGEL									
GPT3	35.4	38.3 [†]	43.6 ^{†,‡}	21.7	26.5 [†]	36.9 ^{†,‡}	30.5	32.9 [†]	43.4 ^{†,‡}
GPT2	40.8	43.2 [†]	40.2	31.4	33.0 [†]	33.5 ^{†,‡}	36.7	39.3 [†]	41.3 ^{†,‡}
OPT	38.7	39.3 [†]	38.6	31.4	30.0	33.7 ^{†,‡}	37.7	37.1 [†]	41.5 ^{†,‡}
FT5	41.5	41.5	39.1	26.5	29.7 [†]	28.6 [†]	38.1	41.1 [†]	40.3 [†]
Avg.	39.1	40.6[†]	40.3[†]	27.7	29.8[†]	33.2^{†,‡}	35.8	37.6[†]	41.6^{†,‡}
SFRES									
GPT3	30.4	25.1	31.5 ^{†,‡}	25.0	30.4 [†]	26.5 [†]	31.2	30.9	26.1
GPT2	22.5	25.1 [†]	20.5	31.0	31.9 [†]	37.0 ^{†,‡}	20.0	33.1 [†]	36.2 ^{†,‡}
OPT	25.2	26.9 [†]	24.3	26.2	30.0 [†]	36.6 ^{†,‡}	21.3	25.6 [†]	30.6 ^{†,‡}
FT5	24.0	21.9	19.7	34.3	34.6 [†]	36.8 ^{†,‡}	22.0	17.8	19.7 [†]
Avg.	25.5	24.7	24.0	29.1	31.7[†]	34.2^{†,‡}	23.6	26.8[†]	28.2^{†,‡}

Table 4: The average of Spearman correlation of the models based on GPT3, GPT2, OPT, and FT5 on BAGEL and SFRES datasets in the data-to-text task.

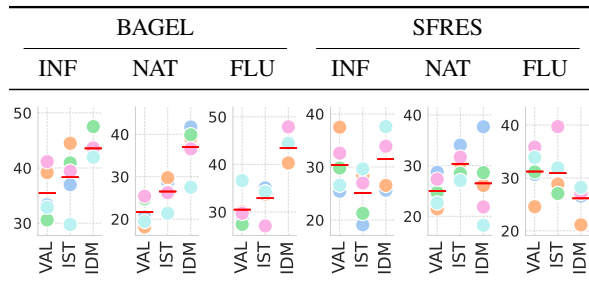


Figure 4: Experimental results for GPT3-based variants in the data-to-text task. Here, blue, orange, green, pink, and cyan dots denote that GPTSCORE is built based on GPT3-a01 (●), GPT3-b01 (●), GPT3-c01 (●), GPT3-d01 (●), and GPT3-d03 (●). The red lines (—) denote the average performance of GPT3-based variants.

5.3 Dialogue Response Generation

To test if GPTSCORE can generalize to more aspects, we choose dialogue response generation task as a testbed, which usually requires evaluating generated texts from a variety of dimensions (i.e., “interesting” and “fluent”). To reduce the computational cost, we focus on GPT3-based metrics only since they have achieved superior performance as we observed in the previous experiments. Tab. 5 shows the performance in dialogue response generation task, where both baseline FED and DE (DynaEval) are fine-tuned on a large dialogue corpus. The main observations are listed as follows:

(1) **The performance of GPT3-d01 is much better than GPT3-d03, even though both of them have the same model size.** The average Spearman correlation of GPT3-d01 outperforms GPT3-d03 by **40.8** on the FED Turn-level dataset, and **5.5** on the FED dialogue-level. (2) **The GPT3-based model demonstrate stronger generalization ability.** BART-based models failed in evaluating the dialogue generation task. The GPT3-a01 with 350M parameters achieved comparable performance to FED and DE models, which are fine-tuned on dialogue corpus.

6 Ablation Study

6.1 Effectiveness of Demonstration

To investigate the relationship between the demonstration sample size (denote as K) and the evaluation performance, we choose the machine translation task and the GPT3-based variants with model sizes ranging from 350M to 175B for further study.

The change of Spearman correlation on the MQM-2020 dataset with different demonstration sample size are shown in Fig. 5. The main observations are summarized as follows: (1) The utilization of demonstration significantly improves the evaluation performance, which holds for these three aspects. (2) There is an upper bound on the performance gains from the introduction of the demonstration. For example, when $K > 4$, the performance of ACC is hard to improve further. (3) When demonstration has only a few samples (such as $K=1$), small models (e.g., GPT3-a01) are prone to performance degradation due to the one-sidedness of the given examples.

6.2 Partial Order of Evaluation Aspect

To explore the correlation between aspects, we conducted an empirical analysis with INT (*interesting*)

Aspect	Baseline					GPTScore				
	BT	BTC	BTCP	FED	DE	a01	b01	c01	d01	d03
FED dialogue-level										
COH	1.7	-14.9	-18.9	25.7	43.7	18.7	15.0	22.5	56.9	13.4
ERR	9.4	-12.2	-13.7	12.0	30.2	35.2	16.8	21.3	45.7	9.40
CON	2.6	-6.7	-10.2	11.6	36.7	33.7	9.9	18.4	32.9	18.1
DIV	13.3	-2.5	-13.9	13.7	37.8	14.9	5.20	21.5	62.8	-6.6
DEP	8.2	-6.6	-17.6	10.9	49.8	9.00	12.9	28.2	66.9	34.1
LIK	9.9	-6.3	-11.8	37.4	41.6	26.2	22.0	32.1	63.4	18.4
UND	-11.5	-17.6	-18.2	-0.3	36.5	31.2	40.0	40.0	52.4	19.6
FLE	9.3	-10.2	-10.3	24.9	38.3	32.7	44.9	34.6	51.5	7.20
INF	9.2	-7.5	-10.5	42.9	42.6	6.80	8.0	18.8	60.2	31.7
INQ	6.2	-0.6	-14.8	24.7	41.0	44.2	38.7	49.2	50.3	-10.1
Avg.	5.8	-8.5	-14.0	20.4	39.8	25.3	21.3	28.6	54.3	13.5
FED turn-level										
INT	15.9	-3.3	-10.1	32.4	32.7	16.6	6.4	30.8	50.1	22.4
ENG	22.6	1.1	-2.5	24.0	30.0	10.2	6.2	29.4	49.6	35.5
SPE	8.3	-7.9	-16.2	14.1	34.6	33.7	16.1	31.7	21.4	15.1
REL	11.9	10.0	19.4	19.9	26.3	8.6	10.3	23.8	45.2	38.0
COR	7.6	1.8	12.4	26.2	24.2	29.7	11.2	27.0	43.4	42.8
SEM	10.0	18.8	26.1	-9.4	20.2	6.8	8.1	23.1	44.4	40.5
UND	12.0	8.1	4.5	1.3	20.0	6.6	14.8	23.4	36.5	31.1
FLU	14.0	17.2	28.4	-13.4	17.1	16.5	5.7	14.0	16.0	36.7
Avg.	12.8	5.7	7.7	11.9	25.6	16.1	9.9	25.4	38.3	32.8

Table 5: Spearman correlation of different aspects on the FED turn- and dialogue-level datasets. *BT*, *BTC*, *BTCP*, and *DE* denote BARTSCORE, BARTSCORE+CNN, BARTSCORE+CNN+Para, and the DynaEval model. Values in bold indicate the best performance.

on the dialogue response generation task of the FED-Turn dataset. Specifically, take INT as the target aspect and then combine the definitions of other aspects with the definition of INT as the final evaluation protocols. The x-axis of Fig. 6-(a) is the aspect order achieved based on the Spearman correlation between INT and that aspect’s human score. Fig. 6-(b) is the Spearman correlation of INT as the modification of the INT definition, and the scoring function is GPT3-c01 with 6.7B parameters.

The following table illustrates the definition composition process, where Sp denotes Spearman.

X	Aspect	Aspect Definition	Sp
1	INT	Is this response interesting to the conversation?	30.8
3	INT, ENG, SPE	Is this an interesting response that is specific and engaging?	48.6

Specifically, the definition of INT is “*Is this response interesting to the conversation?*” at $x=1$ in Fig. 6-(b). When INT combines with ENG, SPE (at $x=3$ in Fig. 6-(b)), its definition can be “*Is this an interesting response that is specific and engaging?*”. And the new aspect definition boosts the performance from **30.8** (at $x=1$ in Fig. 6-(b)) to

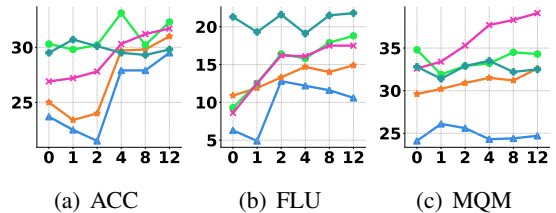


Figure 5: Results of the GPT3 family models with different numbers of examples (K) in the demonstration on the MQM-2020 dataset. Here, blue, orange, green, red, and cyan lines denote that GPTSCORE is built based on GPT3-a01 (\blacktriangle), GPT3-b01 (\star), GPT3-c01 (\bullet), GPT3-d01 (\times), and GPT3-d03 ($+$), respectively.

48.6 (at $x=3$ in Fig. 6-(b)). The best performance of **51.4** ($x=5$ in Fig. 6-(b)) is achieved after combining five aspects (INT, ENG, SPE, COR, REL), which already exceeded **50.1** of the most potent scoring model GPT3-d01 (175B) with aspect definition built only on INT. Therefore, **by combining definitions with other highly correlated aspects, the performance of the smaller model (GPT3-curie, 6.7B) can outperform the bigger model (GPT3-davinci, 175B).**

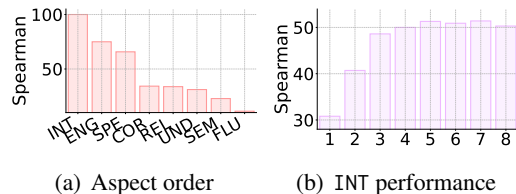


Figure 6: (a) Descending order of Spearman correlation between INT and other aspects’ human scoring. (b) The Spearman correlation of INT changes as its aspect definition is modified in combination with other aspects.

7 Conclusion

In this paper, we propose to leverage the emergent abilities from generative pre-training models to address intricate and ever-changing evaluation requirements. The proposed framework, GPTSCORE, is studied on multiple pre-trained language models with different structures, including the GPT3 (175B). GPTSCORE has multiple benefits: customizability, multi-faceted evaluation, and train-free, which enable us to flexibly craft a metric that can support 22 aspects on 37 datasets without any learning process yet attain competitive performance. Furthermore, demonstrate that GPTScore achieves the goal of “evaluate as you desire”. This work opens a new way to audit generative AI by utilizing generative AI.

8 Limitations

The limitations of this work include: (1) The pre-trained language models considered in our work were released before GPT-3.5 (included), while some recently released popular LLMs (such as ChatGPT and GPT-4) are not studied in this work. (2) GPT3-text-davinci-003 performs worse than GPT3-text-davinci-001, which holds in many evaluation settings. However, we cannot explain this conclusion well until OpenAI discloses the model and training in more details. (3) Due to the cost limitation of using the OpenAI API, we only consider evaluating four traditional NLP generation tasks. The evaluation of some complex text generation tasks (e.g., story generation, a long text generation task) can be studied in the future.

There are some risks associated with model-based evaluation. For example, (1) GPT3-based models work well as text evaluators, but their internal structure, training data, and training process are opaque, which makes it hard to explain the phenomenon of model performance. (2) Training and human feedback datasets influence language model behavior a lot. During the evaluation process, language models may inject bias and risky behaviors that are hard to identify, sending dangerous information to humans. The safety and risk of model-based evaluation should be carefully studied, and we will comprehensively explore the risk for our GPTScore like (Wang et al., 2023) in the future.

Acknowledgements

We thank Chen Zhang for helpful discussion and feedback. This research / project is supported by the National Research Foundation, Singapore under its Industry Alignment Fund – Pre-positioning (IAF-PP) Funding Initiative. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore.

References

Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot](#). *CoRR*, abs/2001.09977.

Manik Bhandari, Pranav Narayan Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. 2020. [Re-](#)

[evaluating evaluation in text summarization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 9347–9359. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *CoRR*, abs/2005.14165.

Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit Cheung. 2020. [Factual error correction for abstractive summarization models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 6251–6258. Association for Computational Linguistics.

Louis Castricato, Alexander Havrilla, Shahbuland Matiana, Michael Pieler, Anbang Ye, Ian Yang, Spencer Frazier, and Mark O. Riedl. 2022. [Robust preference learning for storytelling via contrastive reinforcement learning](#). *CoRR*, abs/2210.07792.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *arXiv preprint arXiv:2210.11416*.

Esin Durmus, He He, and Mona T. Diab. 2020. [FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5055–5070. Association for Computational Linguistics.

Markus Freitag, George F. Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. [Experts, errors, and context: A large-scale study of human evaluation for machine translation](#). *CoRR*, abs/2104.14478.

Jinlan Fu, See-Kiong Ng, and Pengfei Liu. 2022. [Polyglot prompt: Multilingual multitask prompt training](#). *arXiv preprint arXiv:2204.14264*.

Sarik Ghazarian, Nuan Wen, Aram Galstyan, and Nanyun Peng. 2022. [DEAM: dialogue coherence evaluation using amr-based semantic manipulations](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 771–785. Association for Computational Linguistics.

- Max Grusky, Mor Naaman, and Yoav Artzi. 2018. [Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 708–719. Association for Computational Linguistics.
- Karl Moritz Hermann, Tomáš Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1693–1701.
- J. Edward Hu, Abhinav Singh, Nils Holzenberger, Matt Post, and Benjamin Van Durme. 2019. [Large-scale, diverse, paraphrastic bitexts via sampling and clustering](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning, CoNLL 2019, Hong Kong, China, November 3-4, 2019*, pages 44–54. Association for Computational Linguistics.
- Matt J. Kusner, Yu Sun, Nicholas I. Kolkin, and Kilian Q. Weinberger. 2015. [From word embeddings to document distances](#). In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 957–966. JMLR.org.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.
- Zekang Li, Jinchao Zhang, Zhengcong Fei, Yang Feng, and Jie Zhou. 2021. [Conversations are not flat: Modeling the dynamic information flow across dialogue utterances](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 128–138. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Yixin Liu, Alexander R Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, et al. 2022. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation. *arXiv preprint arXiv:2212.07981*.
- François Mairesse, Milica Gasic, Filip Jurčiček, Simon Keizer, Blaise Thomson, Kai Yu, and Steve J. Young. 2010. [Phrase-based statistical language generation using graphical models and active learning](#). In *ACL 2010, Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, July 11-16, 2010, Uppsala, Sweden*, pages 1552–1561. The Association for Computer Linguistics.
- Shahbuland Matiana, J. R. Smith, Ryan Teehan, Louis Castricato, Stella Biderman, Leo Gao, and Spencer Frazier. 2021. [Cut the CARP: fishing for zero-shot story evaluation](#). *CoRR*, abs/2110.03111.
- Shikib Mehri and Maxine Eskénazi. 2020a. [Unsupervised evaluation of interactive dialog with dialogpt](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue, SIGdial 2020, 1st virtual meeting, July 1-3, 2020*, pages 225–235. Association for Computational Linguistics.
- Shikib Mehri and Maxine Eskénazi. 2020b. [USR: an unsupervised and reference free evaluation metric for dialog generation](#). *CoRR*, abs/2005.00456.
- Sewon Min, Xinxin Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. [Rethinking the role of demonstrations: What makes in-context learning work?](#) *CoRR*, abs/2202.12837.
- Mavuto M Mukaka. 2012. A guide to appropriate use of correlation coefficient in medical research. *Malawi medical journal*, 24(3):69–71.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
- Bo Pang, Erik Nijkamp, Wenjuan Han, Linqi Zhou, Yixian Liu, and Kewei Tu. 2020. [Towards holistic and automatic evaluation of open-domain dialogue generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 3619–3629. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318. ACL.

- Maja Popovic. 2015. [chrf: character n-gram f-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation, WMT@EMNLP 2015, 17-18 September 2015, Lisbon, Portugal*, pages 392–395. The Association for Computer Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. [COMET: A neural framework for MT evaluation](#). *CoRR*, abs/2009.09025.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. [BLOOM: A 176b-parameter open-access multilingual language model](#). *CoRR*, abs/2211.05100.
- Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, Alex Wang, and Patrick Gallinari. 2021. [Questeval: Summarization asks for fact-based evaluation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6594–6604. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. [BLEURT: learning robust metrics for text generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7881–7892. Association for Computational Linguistics.
- Team Sequoia. 2022. Generative ai: A creative new world. <https://www.sequoiacap.com/article/generative-ai-a-creative-new-world/>.
- Brian Thompson and Matt Post. 2020. [Automatic machine translation evaluation in many languages via zero-shot paraphrasing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 90–121. Association for Computational Linguistics.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020a. [Asking and answering questions to evaluate the factual consistency of summaries](#). *CoRR*, abs/2004.04228.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020b. [Asking and answering questions to evaluate the factual consistency of summaries](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5008–5020. Association for Computational Linguistics.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipourmolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *URL https://arxiv.org/abs/2204.07705*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). *CoRR*, abs/2201.11903.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Peihao Su, David Vandyke, and Steve J. Young. 2015. [Semantically conditioned lstm-based natural language generation for spoken dialogue systems](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 1711–1721. The Association for Computational Linguistics.
- Wenda Xu, Xian Qian, Mingxuan Wang, Lei Li, and William Yang Wang. 2022a. [Sescore2: Retrieval augmented pretraining for text generation evaluation](#). *CoRR*, abs/2212.09305.
- Wenda Xu, Yi-Lin Tuan, Yujie Lu, Michael Saxon, Lei Li, and William Yang Wang. 2022b. [Not all errors are equal: Learning text generation metrics using stratified error synthesis](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 6559–6574. Association for Computational Linguistics.
- Zheng Ye, Liucun Lu, Lishan Huang, Liang Lin, and Xiaodan Liang. 2021. [Towards quantifiable dialogue coherence evaluation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational*

Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021, pages 2718–2729. Association for Computational Linguistics.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.

Jerrold H Zar. 2005. Spearman rank correlation. *Encyclopedia of biostatistics*, 7.

Chen Zhang, Yiming Chen, Luis Fernando D’Haro, Yan Zhang, Thomas Friedrichs, Grandee Lee, and Haizhou Li. 2021. [Dynaeval: Unifying turn and dialogue level evaluation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 5676–5689. Association for Computational Linguistics.

Chen Zhang, Luis Fernando D’Haro, Qiquan Zhang, Thomas Friedrichs, and Haizhou Li. 2022a. [Finedeval: Fine-grained automatic dialogue-level evaluation](#). *CoRR*, abs/2210.13832.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022b. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. [Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 563–578. Association for Computational Linguistics.

Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. [Towards a unified multi-dimensional evaluator for text generation](#). *CoRR*, abs/2210.07197.

A Machine Translation

Scoring Dimension For aspects ACC, FLU, and MQM from the machine translation task, we also choose *ref*->*hypo*. Because the source text of the machine translation is a different language from the translated text (*hypo*). In this work, we mainly consider the evaluation of the English text. In the future, we can consider designing a scoring function based on BLOOM (Scao et al., 2022) that can evaluate texts in a cross-lingual setting.

The average sample-level Spearman (ρ) scores of GPT3-based, GPT2-based, OPT-based, and FT5-based models on the MQM-2020 machine translation dataset are shown in Tab. 6, where values with † denote that the evaluator equipped with IST (or IDM) significantly outperforms the VAL setting, and ‡ indicate that the evaluator equipped with IDM (the combination of IST and DM) significantly outperforms the IST setting. The Spearman correlations for the GPT3-based variants are shown in Fig. 7. For the full evaluation results of 28 models (including 9 baseline scoring models, such as ROUGE-1) can be found in Tab. 14. Following Thompson and Post (2020) and Yuan et al. (2021), we treat the evaluation of machine translation as the paraphrasing task. The main observations are listed as follows:

(1) **The introduction of instruction (IST) significantly improve the performance in three different aspects of ACC, FLU, and MQM.** In Tab. 6, the average performance of 19 GPTSCORE based evaluators with instruction (IST) significantly outperforms vanilla (VAL). (2) **The combination of instruction and demonstration (IDM) brings gains for the evaluator with different model structures.** In Tab. 6, the performance of GPT3, GPT2, OPT, and FT5 improves a lot when instruction and demonstration (IDM) are introduced. (3) **The evaluator built based on GPT3-c01 achieves comparable performance with GPT3-d01 and GPT3-d03.** This can be found in Fig. 7. Since the GPT3-d01 and GPT3-d03 are most expensive variant of GPT3, the cheaper and comparative GPT3-c01 is a good choice for machine translation task.

B Evaluation Strategy

Evaluation strategies define different aggregation methods when we calculate the correlation scores. Specifically, suppose that for each source text $s_i, i \in [1, 2, \dots, n]$ (e.g., documents in text summarization task or dialogue histories for dialogue

Model	ACC		FLU		MQM				
	VALIST	IDM	VALIST	IDM	VALIST	IDM			
GPT3	27.2	27.1	29.7 ^{†,‡}	11.3	10.4	16.4 ^{†,‡}	30.3	31.2 [†]	32.3 ^{†,‡}
GPT2	25.8	27.0 [†]	30.3 ^{†,‡}	9.8	10.8 [†]	15.8 ^{†,‡}	30.1	30.3 [†]	33.5 ^{†,‡}
OPT	28.7	29.4 [†]	30.3 ^{†,‡}	10.0	12.2 [†]	16.3 ^{†,‡}	32.5	34.6 [†]	35.1 ^{†,‡}
FT5	27.7	27.8 [†]	28.3 ^{†,‡}	9.6	11.0 [†]	15.4 ^{†,‡}	31.0	32.3 [†]	32.3
Avg.	27.4	27.8 [†]	29.7 ^{†,‡}	10.2	11.1 [†]	16.0 ^{†,‡}	31.0	32.1 [†]	33.3 ^{†,‡}

Table 6: The average Spearman correlation of the GPT3-based, GPT2-based, OPT-based, and FT5-based models in machine translation task of MQM-2020 dataset.

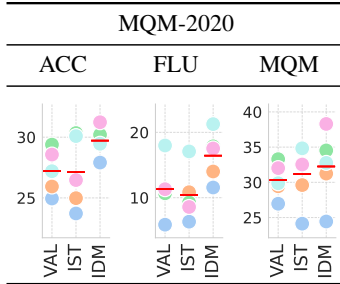


Figure 7: Experimental results for GPT3-based variants in the machine translation task. Here, blue, orange, green, pink, and cyan dot denote that GPTSCORE is built based on GPT3-a01 (●), GPT3-b01 (●), GPT3-c01 (●), GPT3-d01 (●), and GPT3-d03 (●), respectively. The red lines (—) denote the average performance of GPT3-based variants.

generation task), there are J system outputs $\mathbf{h}_{i,j}$, where $j \in [1, 2, \dots, J]$. f_{auto} is an automatic scoring function (e.g., ROUGE (Lin, 2004)), and f_{human} is the gold human scoring function. For a given evaluation aspect a , the meta-evaluation metric F can be formulated as follows.

Sample-level defines that a correlation value is calculated for each sample separately based on outputs of multiple systems, then averaged across all samples.

$$F_{f_{\text{auto}}, f_{\text{human}}}^{\text{sample}} = \frac{1}{n} \sum_{i=1}^n \left(g \left([f_{\text{auto}}(\mathbf{h}_{i,1}), \dots, f_{\text{auto}}(\mathbf{h}_{i,J})], [f_{\text{human}}(\mathbf{h}_{i,1}), \dots, f_{\text{human}}(\mathbf{h}_{i,J})] \right) \right)$$

where g can be instantiated as Spearman or Pearson correlation.

Dataset-level indicates that the correlation value is calculated on system outputs of all n samples.

$$F_{f_{\text{auto}}, f_{\text{human}}}^{\text{data}} = g \left([f_{\text{auto}}(\mathbf{h}_{1,1}), \dots, f_{\text{auto}}(\mathbf{h}_{n,J})], [f_{\text{human}}(\mathbf{h}_{1,1}), \dots, f_{\text{human}}(\mathbf{h}_{n,J})] \right)$$

In this work, we select the evaluation strategy for a specific task based on previous works (Yuan et al., 2021; Zhang et al., 2022a). We use the sample-level evaluation strategy for text summarization, data-to-text, and machine translation tasks. For the dialogue response generation task, the dataset-level evaluation strategy is utilized.

C Metric Comparison

Tab. 7 summarize several popular generated text evaluation methods.

D Tasks, Datasets, and Aspects

To achieve a more comprehensive evaluation, in this paper, we cover a broad range of natural language generation tasks: *Dialogue Response Generation*, *Text Summarization*, *Data-to-Text*, and *Machine Translation*, which involves 9 datasets and 22 evaluation aspects in total. Tab. 8 summarizes the tasks, datasets, and evaluation aspects considered by each dataset. The definition of different aspects can be found in Tab. 1.

Dialogue Response Generation aims to automatically generate an engaging and informative response based on the dialogue history. (1) FED (Mehri and Eskénazi, 2020a) collects 124 conversations, including both human-machine (Meena (Adiwardana et al., 2020), Mitsuku⁶) and human-human dialogues, and manually annotated 9 and 11 evaluation aspects at the turn- and dialogue-level, respectively.

Text Summarization is a task of automatically generating an informative and fluent summary for a given long text. Here, we consider the following four datasets covering 6 evaluation aspects: *semantic coverage*, *informativeness*, *relevance*, *fluency*, *coherence*, and *factuality*. (1) SummEval (Bhandari et al., 2020) collects human judgments on 16 model-generated summaries on the CNN/Daily Mail dataset, covering aspects of coherence, consistency, fluency, and relevance. (2) REALSumm (Bhandari et al., 2020) evaluates the reliability of automatic metrics by measuring the pyramid recall of text generated by 25 systems. (3) NEWSROOM (Grusky et al., 2018) covers news, sports, entertainment, finance, and other topics and evaluates the quality of summaries generated by 7

⁶<https://medium.com/pandorabots-blog/mitsuku-wins-loebner-prize-2018-3e8d98c5f2a7>

Metrics	Custom	Function (f)		Additional text (S)		Training-free	Application
		Representation	Formulation	Source	Reference		
ROUGE (Lin, 2004)	✗	Token	Matching	No	Required	✓	SUM
BLEU (Papineni et al., 2002)	✗	Token	Matching	No	Required	✓	MT
CHRF (Popovic, 2015)	✗	Character	Matching	No	Required	✓	MT
BERTScore (Zhang et al., 2020)	✗	BERT	Matching	No	Required	✓	MUL(2)
MoverScore (Zhao et al., 2019)	✗	BERT	Matching	No	Required	✓	MUL(4)
BLEURT (Sellam et al., 2020)	✗	BERT	Regression	No	Required	✓	MT
PRISM (Thompson and Post, 2020)	✗	Embedding	Paraphrase	Optional	Optional	✓	MT
UNIEVAL (Zhong et al., 2022)	✗	T5	Boolean QA	Optional	Optional	✗	MUL(2)
COMET (Rei et al., 2020)	✗	BERT	Regress, Rank	Optional	Optional	✗	MT
BARTScore (Yuan et al., 2021)	✗	BART	Generation	Optional	Optional	✓	MUL(3)
FED (Mehri and Eskénazi, 2020a)	✗	DialoGPT	Generation	Required	Optional	✓	Dialogue
HolisticEval (Pang et al., 2020)	✗	GPT2	Generation	Optional	Optional	✓	Dialogue
GPTScore	✓	GPT3/OPT	Any	Optional	Optional	✓	MUL(5)

Table 7: A comprehensive comparison of existing research on automated evaluation of generated texts. MUL(k) denotes multiple (k) applications explored. *Custom* denotes *Custom Aspects*.

Tasks	Dataset	Aspect
Diag	FED-Diag	COH, DIV, FLE, UND, INQ CON, INF, LIK, DEP, ERR
	FED-Turn	INT, ENG, SPE, REL, COR, SEM, UND, FLU
Summ	SummEval	COH, CON, FLU, REL
	Newsroom	FLU, REL, INF, COH
	REALSumm	COV
	Q-XSUM	FAC
D2T	BAGEL	FLU, REL, INF
	SFRES	FLU, REL, INF
MT	MQM-2020	FLU, COH, INF

Table 8: An overview of tasks, datasets, and evaluation aspects. *Summ.* denote the text summarization task, *D2T* denotes the Data-to-Text task, *MT* denotes the machine translation. Tab. 1 summarized the definitions of the aspects explored in this work.

systems, including informativeness, relevance, fluency, and coherence. (4) QAGS_XSUM (Wang et al., 2020b) is another dataset focusing on the factuality aspect. It has 239 samples from XSUM and their summaries are generated by a fine-tuned BART model.

Data-to-Text aims to automatically generate a fluent and factual description for a given table. (1) BAGEL (Mairesse et al., 2010) contains 202 samples about restaurants in Cambridge. (2) SFRES (Wen et al., 2015) contains 581 samples about restaurants in San Francisco. These two datasets consider three evaluation aspects: *informativeness*, *naturalness* (relevance), and *quality* (fluency).

Machine Translation aims to translate a sentence from one language to another. We consider a sub-datasets of Multidimensional Quality Metrics (MQM) (Freitag et al., 2021), namely, MQM-2020 (Chinese->English). Due to limited annotations, here, we only consider three evaluation aspects: *accuracy*, *fluency*, and *MQM* with diverse scores.

E Ablation Study

E.1 Effectiveness of Demonstration

The in-context learning helps a lot to achieve a good performance. However, how does the number of samples in the demonstration impact the performance? We conduct a case study on the five GPT3-based models explored in this work. The experimental results are shown in Fig. 5, and the specific performance values can be seen in Tab. 9.

E.2 Partial Order of Evaluation Aspect

We have investigated the combination of different evaluation aspects to achieve further performance gains in § 6.2. Tab. 10 summarizes the aspect definition and Spearman correlation changes for INT, with the introduction of other aspects.

F Prompt Design

In this work, we have studied four popular text generation tasks: text summarization, machine translation, data-to-text, and dialogue response generation. The instructions for these tasks on different evaluation aspects are summarized in Tab. 11 and Tab. 12. Here, we convert the dialogue response generation task as a boolean question-answering task and in-

Model	K	ACC	FLU	MQM
GPT3-ada	0	23.7	6.3	24.1
	1	22.5	4.9	26.1
	2	21.5	12.8	25.6
	4	27.9	12.2	24.3
	8	27.9	11.6	24.4
	12	29.5	10.6	24.7
GPT3-babbage	0	25.0	10.9	29.6
	1	23.4	11.9	30.2
	2	24.0	13.3	30.9
	4	29.7	14.7	31.5
	8	29.8	14.0	31.2
	12	31.0	14.9	32.6
GPT3-curie	0	30.3	9.3	34.8
	1	29.8	12.5	31.9
	2	30.2	16.4	32.9
	4	33.1	15.8	33.2
	8	30.2	17.9	34.5
	12	32.3	18.8	34.3
GPT3-davinci001	0	26.9	8.6	32.6
	1	27.2	12.5	33.4
	2	27.8	16.2	35.3
	4	30.3	16.1	37.7
	8	31.2	17.5	38.3
	12	31.7	17.5	39.1
GPT3-davinci003	0	29.5	21.3	32.8
	1	30.7	19.3	31.4
	2	30.1	21.6	32.9
	4	29.5	19.1	33.5
	8	29.3	21.5	32.2
	12	29.8	21.8	32.5

SFRES dataset.

Table 9: Spearman correlation of the GPT3-based models (e.g. text-ada-001 and text-davinci-001) with different demonstration sample numbers on the MQM-2020 dataset .K denotes the number of samples in the demonstration.

corporate the aspect definition into the question of the boolean question-answering task.

G Experiment Results

This section lists the full experimental results for the explored text generation tasks. The models considered here include the 9 baseline models: ROUGE-1, ROUGE-2, ROUGE-L, BERTScore, MoverScore, PRISM, BARTSCORE, BARTSCORE+CNN, and BARTSCORE+CNN+Para, and 19 GPTScore models built based on the GPT3-based, GPT2-based, OPT-based, and Flan-T5-based pre-trained models.

Tab. 13 lists the results of the text summarization datasets. Tab. 14 lists the results of the machine translation datasets. Tab. 15 shows the results of the data-to-text task on the BAGEL dataset. Tab. 16 shows the results of the data-to-text task on the

X	Aspect	Aspect Definition	Spear
1	Interesting (INT)	Is this response interesting to the conversation?	36.9
2	Engaging (ENG)	Is this an interesting response that is engaging?	40.7
3	Specific (SPE)	Is this an interesting response that is specific and engaging?	48.6
4	Correct (COR)	Is this an interesting response that is engaging, specific, and correct?	50.0
5	Relevant (REL)	Is this an interesting response that is specific, engaging, relevant, and correct?	51.3
6	Understandable (UND)	Is this an interesting response that is specific, engaging, relevant, correct, and understandable?	50.9
7	Semantically appropriate (SEM)	Is this an interesting response that is specific, engaging, relevant, correct, understandable, and semantically appropriate?	51.4
8	Fluent (FLU)	Is this an interesting response that is specific, engaging, relevant, correct, understandable, semantically appropriate, and fluent?	50.3

Table 10: The aspect definition and Spearman correlation of INT. X denotes the number of aspects combined with the INT. The scoring model is GPT3-c01.

Aspect	Function	Instruction
Text Summarization		
FAC	src->hypo	Generate a summary with consistent facts for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text with consistent facts. {ref/hypo} In other words, {hypo/ref}
COV	src->hypo	Generate a summary with as much semantic coverage as possible for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text with the same semantics. {ref/hypo} In other words, {hypo/ref}
CON	src->hypo	Generate factually consistent summary for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text with consistent facts. {ref/hypo} In other words, {hypo/ref}
INF	src->hypo	Generate an informative summary that captures the key points of the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text with its core information. {ref/hypo} In other words, {hypo/ref}
COH	src->hypo	Generate a coherent summary for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text into a coherent text. {ref/hypo} In other words, {hypo/ref}
REL	src->hypo	Generate a relevant summary with consistent details for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text with consistent details. {ref/hypo} In other words, {hypo/ref}
FLU	src->hypo	Generate a fluent and grammatical summary for the following text: {src}\n\nTl;dr{hypo}
	ref<->hypo	Rewrite the following text into a fluent and grammatical text. {ref/hypo} In other words, {hypo/ref}
Machine Translation		
Acc	ref<->hypo	Rewrite the following text with its core information and consistent facts:{ref/hypo} In other words, {hypo/ref}
FLU	ref<->hypo	Rewrite the following text to make it more grammatical and well-written:{ref/hypo} In other words, {hypo/ref}
MQM	ref<->hypo	Rewrite the following text into high-quality text with its core information:{ref/hypo} In other words, {hypo/ref}
Data to Text		
INF	ref<->hypo	Convert the following text to another expression that preserves key information:\n\n{ref/hypo} In other words, {hypo/ref}
NAT	ref<->hypo	Convert the following text into another expression that is human-like and natural:\n\n{ref/hypo} In other words, {hypo/ref}
FLU	ref<->hypo	Convert the following text into another expression that preserves key information and is human-like and natural:\n\n{ref/hypo} In other words, {hypo/ref}

Table 11: Instruction design on different aspects for text summarization, machine translation, and data-to-text tasks. *src*, *hypo*, and *ref* denote the *source text*, *hypothesis text*, and *reference text*, respectively. $a \rightarrow b$ ($a \leftarrow b$) denotes to evaluate the quality of b (a) text based on the given a (b) text.

Aspect	Instruction
FED Turn-Level	
INT	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI interesting? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
ENG	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI engaging? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
UND	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI understandable? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
REL	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI relevant to the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
SPE	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI generic or specific to the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
COR	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI correct to conversations? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.]
SEM	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI semantically appropriate? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
FLU	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI fluently written? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
FED Dialog-Level	
COH	Answer the question based on the conversation between a human and AI.\nQuestion: Is the AI coherent and maintains a good conversation flow throughout the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
DIV	Answer the question based on the conversation between a human and AI.\nQuestion: Is there diversity in the AI responses? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
FLE	Answer the question based on the conversation between a human and AI.\nQuestion: Is the AI flexible and adaptable to human and their interests? (a) Yes. (b) No. \nConversation: {History}\nAnswer: Yes.
UND	Answer the question based on the conversation between a human and AI.\nQuestion: Does the AI seem to understand the human? (a) Yes. (b) No. \nConversation: {History}\nAnswer: Yes.
INQ	Answer the question based on the conversation between a human and AI.\nQuestion: Is the AI inquisitive throughout the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
CON	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI consistent in the information it provides throughout the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
INF	Answer the question based on the conversation between a human and AI.\nQuestion: Are the responses of AI informative throughout the conversation? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
LIK	Answer the question based on the conversation between a human and AI.\nQuestion: Does the AI display a likeable personality? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
DEP	Answer the question based on the conversation between a human and AI.\nQuestion: Does the AI discuss topics in depth? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.
ERR	Answer the question based on the conversation between a human and AI.\nQuestion: Is the AI able to recover from errors that it makes? (a) Yes. (b) No.\nConversation: {History}\nAnswer: Yes.

Table 12: Instruction design on various aspects for dialogue response generation task at the turn- and dialogue-level. *History* indicates the conversation history. We convert the evaluation of the response generation task as a question-answering task, and the aspect definition is incorporated into the question of the question-answering task.

Model	NEWSROOM								QXSUM		RSumm	
	COH		CON		FLU		REL		COV		COV	
	VAL	IST	VAL	IST	VAL	IST	VAL	IST	VAL	IST	VAL	IST
ROUGE-1	27.3	-	26.1	-	25.9	-	34.4	-	3.6	-	46.4	-
ROUGE-2	10.9	-	11.7	-	11.2	-	14.4	-	9.9	-	37.3	-
ROUGE-L	24.7	-	25.7	-	24.4	-	32.5	-	5.2	-	45.1	-
BERTScore	31.7	-	31.7	-	27.2	-	33.7	-	-4.6	-	38.4	-
MoverScore	17.7	-	14.2	-	16.0	-	18.9	-	5.4	-	34.4	-
PRISM	60.7	-	56.5	-	59.2	-	61.9	-	2.5	-	32.3	-
BARTSCORE	70.3	-	67.2	-	63.1	-	68.8	-	0.9	-	43.1	-
+CNN	68.5	-	64.9	-	60.4	-	66.3	-	18.4	-	42.9	-
+CNN+Para	69.0	-	65.5	-	62.5	-	67.3	-	6.4	-	40.9	-
GPT3												
GPT3-a01	71.6	71.9 [†]	69.7	70.0 [†]	66.0	67.0 [†]	69.6	69.2	10.3	9.2	29.5	29.8 [†]
GPT3-b01	73.6	72.9	70.2	70.3	66.8	68.3 [†]	71.5	71.2	8.5	14.2	35.0	35.2 [†]
GPT3-c01	73.8	72.8	70.5	70.9[†]	65.9	68.6 [†]	71.0	71.1	15.2	22.1 [†]	36.1	45.1[†]
GPT3-d01	72.6	73.4[†]	68.5	70.0 [†]	65.9	66.9 [†]	71.1	72.1 [†]	24.0	22.7	36.0	33.9
GPT3-d03	73.8	73.1	70.4	70.0	67.4	68.9[†]	74.1	73.3	21.7	22.0 [†]	35.2	38.0 [†]
Avg.	73.1	72.8	69.9	70.2 [†]	66.4	67.9 [†]	71.4	71.4	15.9	18.0 [†]	34.4	36.4 [†]
GPT2												
GPT2-M	68.9	71.7 [†]	66.4	68.0 [†]	61.1	62.3 [†]	67.0	66.8	18.1	18.7 [†]	41.8	43.3 [†]
GPT2-L	70.5	72.3[†]	66.6	68.3 [†]	60.2	61.4 [†]	66.8	67.8 [†]	19.2	19.6 [†]	39.6	41.3 [†]
GPT2-XL	71.0	70.5	66.6	66.6	61.4	60.7	67.2	66.9	21.2	21.2	40.4	41.0 [†]
GPT-J-6B	71.8	71.4	69.8	69.5	65.5	65.5	69.4	69.3	21.6	22.0[†]	42.8	43.7[†]
Avg.	70.5	71.5 [†]	67.4	68.1 [†]	62.0	62.5 [†]	67.6	67.7	20.0	20.4 [†]	39.8	41.1 [†]
OPT												
OPT-350M	70.6	71.5 [†]	69.2	69.9 [†]	67.3	68.1 [†]	70.8	71.6 [†]	13.5	13.3	40.2	42.3 [†]
OPT-1.3B	73.2	73.6[†]	70.9	71.3[†]	67.2	67.8[†]	72.5	72.4	21.1	19.9	42.0	39.7
OPT-6.7B	71.9	71.9	69.0	69.0	67.7	67.1	71.7	71.3	21.2	19.9	38.0	41.9 [†]
OPT-13B	71.9	71.9	68.9	69.6 [†]	65.4	66.0 [†]	71.2	71.5 [†]	23.1	22.1	37.6	41.0 [†]
OPT-66B	72.8	72.8	70.0	69.5	66.0	65.9	71.9	71.9	24.0	23.1	40.3	41.3 [†]
Avg.	72.1	72.3 [†]	69.6	69.9 [†]	66.7	67.0 [†]	71.6	71.8 [†]	20.6	19.6	39.6	41.2 [†]
Flan-T5												
FT5-S	68.3	69.2 [†]	64.6	64.1	59.8	60.4 [†]	64.6	65.5 [†]	14.4	15.1 [†]	33.6	35.7 [†]
FT5-B	68.9	69.0	64.8	64.6	59.6	59.9 [†]	66.5	66.5	13.6	16.3 [†]	36.7	38.6 [†]
FT5-L	70.5	69.1	66.1	64.6	60.9	60.0	66.6	65.4	27.2	28.8[†]	31.4	39.3 [†]
FT5-XL	72.1	70.1	66.7	65.6	61.0	60.5	68.3	67.5	18.9	25.6 [†]	34.8	43.8[†]
FT5-XXL	70.7	69.3	65.7	65.2	60.2	60.4 [†]	67.6	67.8[†]	23.9	27.8 [†]	40.2	41.1 [†]
Avg.	70.1	69.3	65.6	64.8	60.3	60.2	66.7	66.5	19.6	22.7 [†]	35.3	39.7 [†]
Overall Avg	71.5	71.5	68.1	68.3	64.0	64.5 [†]	69.4	69.4	19.0	20.2 [†]	37.4	39.8 [†]

Table 13: Spearman correlations on NEWSROOM and QXSUM datasets for text summarization task. VAL and IST denote the evaluator with vanilla and instruction, respectively. Values with † denote the evaluator with instruction significantly outperforms with vanilla. Values in bold are the best performance in a set of variants (e.g., GPT3 family).

Model	ACC			FLU			MQM		
	VAL	IST	IDM	VAL	IST	IDM	VAL	IST	IDM
ROUGE-1	21.3	-	-	1.7	-	-	17.5	-	-
ROUGE-2	15.0	-	-	5.8	-	-	15.4	-	-
ROUGE-L	16.6	-	-	8.7	-	-	15.7	-	-
BERTScore	26.1	-	-	8.2	-	-	23.6	-	-
MoverScore	18.2	-	-	1.2	-	-	17.2	-	-
PRISM	25.9	-	-	9.1	-	-	27.4	-	-
BARTSCORE	26.1	-	-	8.2	-	-	23.6	-	-
+CNN	26.2	-	-	8.1	-	-	28.7	-	-
+CNN+Para	31.0	-	-	10.8	-	-	29.9	-	-
GPT3									
GPT3-a01	24.9	23.7	27.9 ^{†,‡}	5.9	6.3 [†]	11.6 ^{†,‡}	27.0	24.1	24.4 [‡]
GPT3-b01	25.9	25.0	29.8 ^{†,‡}	10.7	10.8	14.0 ^{†,‡}	29.4	29.6	31.2 ^{†,‡}
GPT3-c01	29.4	30.3[†]	30.2 [†]	10.7	9.3	17.9 ^{†,‡}	33.3	34.8 [†]	34.5 [†]
GPT3-d01	28.6	26.5	31.2^{†,‡}	11.3	8.6	17.5 ^{†,‡}	32.0	32.5 [†]	38.3^{†,‡}
GPT3-d03	27.2	30.1 [†]	29.5 [†]	18.0	17.1	21.3^{†,‡}	29.9	34.8[†]	32.8 [†]
Avg.	27.2	27.1	29.7 ^{†,‡}	11.3	10.4	16.4 ^{†,‡}	30.3	31.2 [†]	32.3 ^{†,‡}
GPT2									
GPT2-M	25.7	24.6	29.6 ^{†,‡}	8.6	9.4 [†]	15.1 ^{†,‡}	32.1	29.4	34.1 ^{†,‡}
GPT2-L	27.2	28.5 [†]	32.2 ^{†,‡}	11.1	10.4	14.9 ^{†,‡}	31.2	30.9	33.9 ^{†,‡}
GPT2-XL	24.2	27.6 [†]	29.7 ^{†,‡}	9.4	12.0 [†]	17.4 ^{†,‡}	28.6	32.2 [†]	35.8 ^{†,‡}
GPT-J-6B	26.2	27.2 [†]	29.5 ^{†,‡}	9.9	11.2 [†]	15.9 ^{†,‡}	28.5	28.8 [†]	30.3 ^{†,‡}
Avg.	25.8	27.0 [†]	30.3 ^{†,‡}	9.8	10.8 [†]	15.8 ^{†,‡}	30.1	30.3 [†]	33.5 ^{†,‡}
OPT									
OPT-350M	29.3	28.1	28.6 [‡]	11.7	11.9	15.7 ^{†,‡}	31.5	32.5 [†]	31.8
OPT-1.3B	27.9	27.7	28.0 [‡]	8.8	13.3 [†]	15.9 ^{†,‡}	32.6	33.6 [†]	32.9 [†]
OPT-6.7B	29.6	30.7 [†]	30.6 [†]	10.7	12.2 [†]	15.0 ^{†,‡}	34.2	36.4 [†]	36.9 ^{†,‡}
OPT-13B	27.5	29.5 [†]	30.8 ^{†,‡}	9.6	11.7 [†]	17.9 ^{†,‡}	31.9	35.5 [†]	37.5 ^{†,‡}
OPT-66B	29.5	31.0 [†]	33.4 ^{†,‡}	9.1	12.1 [†]	16.8 ^{†,‡}	32.1	35.3 [†]	36.4 ^{†,‡}
Avg.	28.7	29.4 [†]	30.3 ^{†,‡}	10.0	12.2 [†]	16.3 ^{†,‡}	32.5	34.6 [†]	35.1 ^{†,‡}
Flan-T5									
FT5-S	27.6	28.7 [†]	27.0	12.6	9.4	15.0 ^{†,‡}	33.5	33.3	31.3
FT5-B	25.5	25.4	27.4 ^{†,‡}	10.4	10.2	15.9 ^{†,‡}	29.8	29.6	30.0 [‡]
FT5-L	28.5	28.5	28.8 ^{†,‡}	7.9	13.0 [†]	15.6 ^{†,‡}	30.7	31.6 [†]	32.1 ^{†,‡}
FT5-XL	28.1	27.0	28.1 [‡]	9.4	10.2 [†]	14.0 ^{†,‡}	30.4	33.5 [†]	34.2 ^{†,‡}
FT5-XXL	29.0	29.4 [†]	30.5 ^{†,‡}	7.6	12.2 [†]	16.2 ^{†,‡}	30.7	33.3 [†]	33.8 ^{†,‡}
Avg.	27.7	27.8	28.3 ^{†,‡}	9.6	11.0 [†]	15.4 ^{†,‡}	31.0	32.3 [†]	32.3 [†]
Overall Avg	27.4	27.8 [†]	29.7 ^{†,‡}	10.2	11.1 [†]	16.0 ^{†,‡}	31.0	32.1 [†]	33.3 ^{†,‡}

Table 14: Spearman correlations on MQM-2020 dataset for machine translation task. VAL, IST, and IDM denote the evaluator with vanilla, instruction, and the combination of instruction and demonstration, respectively. Values with [†] denote the evaluator with instruction significantly outperforms with vanilla, and values with [‡] denote the evaluator with the combination of instruction and demonstration significantly outperforms with only instruction. Values in bold are the best performance in a set of variants (e.g., GPT3 family).

Model	INF			NAT			FLU		
	VAL	IST	IST+DM	VAL	IST	IST+DM	VAL	IST	IST+DM
ROUGE-1	28.7	-	-	5.0	-	-	8.3	-	-
ROUGE-2	24.0	-	-	15.2	-	-	16.0	-	-
ROUGE-L	26.3	-	-	10.5	-	-	11.0	-	-
BERTScore	37.2	-	-	16.0	-	-	18.7	-	-
MoverScore	30.7	-	-	20.4	-	-	14.8	-	-
PRISM	36.8	-	-	28.7	-	-	34.4	-	-
BARTSCORE	29.5	-	-	24.0	-	-	29.7	-	-
+CNN	37.7	-	-	30.1	-	-	34.4	-	-
+CNN+Para	39.2	-	-	31.0	-	-	44.9	-	-
GPT3									
GPT3-a01	33.3	37.0 [†]	42.5 ^{†,‡}	20.5	28.7 [†]	41.7^{†,‡}	28.8	35.1[†]	40.2 ^{†,‡}
GPT3-b01	39.2	44.5[†]	42.2 [†]	18.2	29.8[†]	39.1 ^{†,‡}	30.0	33.8 [†]	40.3 ^{†,‡}
GPT3-c01	30.6	40.9 [†]	47.5^{†,‡}	24.8	26.5 [†]	39.9 ^{†,‡}	27.4	34.2 [†]	44.2 ^{†,‡}
GPT3-d01	41.2	39.4	43.6 ^{†,‡}	25.4	26.2 [†]	36.6 ^{†,‡}	29.7	27.1	47.9^{†,‡}
GPT3-d03	32.9	29.8	42.0 ^{†,‡}	19.5	21.4 [†]	27.5 ^{†,‡}	36.6	34.2	44.4 ^{†,‡}
Avg.	35.4	38.3[†]	43.6^{†,‡}	21.7	26.5[†]	36.9^{†,‡}	30.5	32.9[†]	43.4^{†,‡}
GPT2									
GPT2-M	39.4	42.9 [†]	38.6	31.2	33.2 [†]	34.3 ^{†,‡}	38.9	38.9	39.6 ^{†,‡}
GPT2-L	39.7	42.2 [†]	41.8 [†]	30.1	33.5 [†]	33.1 [†]	34.0	40.0 [†]	39.6 [†]
GPT2-XL	41.2	42.0 [†]	38.7	31.7	33.7 [†]	34.8 ^{†,‡}	38.0	40.6 [†]	44.2 ^{†,‡}
GPT-J-6B	42.8	45.6 [†]	41.6	32.5	31.5	31.9 [‡]	35.9	37.7 [†]	42.0 ^{†,‡}
Avg.	40.8	43.2[†]	40.2	31.4	33.0[†]	33.5^{†,‡}	36.7	39.3[†]	41.3^{†,‡}
OPT									
OPT-350M	37.0	36.8	37.9 ^{†,‡}	33.9	32.5	31.1	39.9	39.5	39.9 [‡]
OPT-1.3B	36.7	39.3 [†]	38.2 [†]	28.8	30.0 [†]	32.9 ^{†,‡}	37.3	34.9	40.9 ^{†,‡}
OPT-6.7B	40.4	39.3	38.3	31.6	27.2	35.2 ^{†,‡}	36.0	34.4	43.6 ^{†,‡}
OPT-13B	37.9	37.6	38.9 ^{†,‡}	31.4	30.3	34.6 ^{†,‡}	39.2	39.0	41.2 ^{†,‡}
OPT-66B	41.4	43.2 [†]	39.6	31.3	30.2	34.7 ^{†,‡}	36.3	37.6 [†]	42.0 ^{†,‡}
Avg.	38.7	39.3	38.6	31.4	30.0	33.7^{†,‡}	37.7	37.1	41.5^{†,‡}
Flan-T5									
FT5-S	39.8	37.6	38.2	33.0	29.5	26.6	46.1	34.7	36.1 [‡]
FT5-B	39.7	43.6 [†]	37.7	26.4	30.3 [†]	27.3 [†]	37.8	40.6 [†]	37.9
FT5-L	42.0	42.8 [†]	38.9	23.6	31.0 [†]	32.6 ^{†,‡}	35.3	43.3 [†]	44.5 ^{†,‡}
FT5-XL	41.0	42.8 [†]	43.3 ^{†,‡}	24.8	28.9 [†]	27.8 [†]	37.4	44.4 [†]	41.9 [†]
FT5-XXL	44.9	40.7	37.4	24.8	28.8 [†]	28.4 [†]	34.2	42.5 [†]	41.3 [†]
Avg.	41.5	41.5	39.1	26.5	29.7[†]	28.6[†]	38.1	41.1[†]	40.3[†]
Overall Avg	39.1	40.6[†]	40.3[†]	27.7	29.8[†]	33.2^{†,‡}	35.8	37.6[†]	41.6^{†,‡}

Table 15: Spearman correlations on BAGEL dataset for data-to-text task. VAL, IST, and IDM denote the evaluator with vanilla, instruction, and the combination of instruction and demonstration, respectively. Values with [†] denote the evaluator with instruction significantly outperforms with vanilla, and values with [‡] denote the evaluator with the combination of instruction and demonstration significantly outperforms with only instruction. Values in bold are the best performance in a set of variants (e.g., GPT3 family).

Model	INF			NAT			FLU		
	VAL	IST	IST+DM	VAL	IST	IST+DM	VAL	IST	IST+DM
ROUGE-1	24.2	-	-	24.2	-	-	15.1	-	-
ROUGE-2	21.9	-	-	25.9	-	-	11.4	-	-
ROUGE-L	18.5	-	-	20.2	-	-	1.7	-	-
BERTScore	25.8	-	-	28.0	-	-	11.8	-	-
MoverScore	17.9	-	-	24.4	-	-	5.0	-	-
PRISM	27.4	-	-	33.1	-	-	14.2	-	-
BARTSCORE	22.4	-	-	25.5	-	-	6.9	-	-
+CNN	24.2	-	-	30.6	-	-	17.2	-	-
+CNN+Para	25.0	-	-	30.2	-	-	19.5	-	-
GPT3									
GPT3-a01	25.4	19.1	25.6 [‡]	28.7	34.0 [†]	37.7 ^{†,‡}	30.7	27.0	26.6
GPT3-b01	37.5	28.4	26.5	21.5	30.6 [†]	26.1 [†]	24.6	28.9 [†]	21.1
GPT3-c01	29.8	21.3	33.7 ^{†,‡}	24.7	28.5 [†]	28.6 [†]	31.1	27.1	27.6 [‡]
GPT3-d01	32.6	27.0	33.9 ^{†,‡}	27.3	31.7 [†]	21.9	35.8	39.7 [†]	27.1
GPT3-d03	26.6	29.6 [†]	37.6 ^{†,‡}	22.6	27.0 [†]	18.2	33.9	31.9	28.2
Avg.	30.4	25.1	31.5 ^{†,‡}	25.0	30.4 [†]	26.5 [†]	31.2	30.9	26.1
GPT2									
GPT2-M	24.7	23.1	18.2	28.7	32.7 [†]	35.2 ^{†,‡}	18.7	34.8 [†]	33.6 [†]
GPT2-L	19.6	28.1 [†]	20.2 [†]	31.2	32.4 [†]	37.8 ^{†,‡}	18.6	33.1 [†]	35.9 ^{†,‡}
GPT2-XL	22.0	23.6 [†]	23.8 [†]	29.7	29.1	38.0 ^{†,‡}	18.2	29.8 [†]	37.1 ^{†,‡}
GPT2-J-6B	23.9	25.6 [†]	19.6	34.3	33.3	36.8 ^{†,‡}	24.4	34.5 [†]	38.4 ^{†,‡}
Avg.	22.5	25.1 [†]	20.5	31.0	31.9 [†]	37.0 ^{†,‡}	20.0	33.1 [†]	36.2 ^{†,‡}
OPT									
OPT-350M	26.1	28.7 [†]	25.4	27.0	29.5 [†]	35.0 ^{†,‡}	21.7	26.6 [†]	27.3 ^{†,‡}
OPT-1.3B	26.1	28.3 [†]	23.5	26.0	30.5 [†]	38.7 ^{†,‡}	23.0	26.9 [†]	29.8 ^{†,‡}
OPT-6.7B	26.2	26.0	24.2	26.7	31.0 [†]	36.5 ^{†,‡}	21.7	25.8 [†]	35.9 ^{†,‡}
OPT-13B	27.7	26.9	26.0	24.4	30.1 [†]	38.0 ^{†,‡}	20.2	29.6 [†]	34.9 ^{†,‡}
OPT-66B	20.1	24.7 [†]	22.4 [†]	26.8	29.1 [†]	34.6 ^{†,‡}	19.8	19.1	25.3 ^{†,‡}
Avg.	25.2	26.9 [†]	24.3	26.2	30.0 [†]	36.6 ^{†,‡}	21.3	25.6 [†]	30.6 ^{†,‡}
Flan-T5									
FT5-S	19.7	16.9	17.0	33.6	33.1	33.0	19.4	17.2	15.9
FT5-B	24.2	23.7	20.9	31.7	32.5 [†]	33.4 ^{†,‡}	14.2	15.5 [†]	16.8 ^{†,‡}
FT5-L	24.9	22.3	20.6	36.2	37.1 [†]	38.6 ^{†,‡}	24.3	18.1	21.1 [‡]
FT5-XL	26.1	23.7	19.5	38.4	35.6	37.4 [‡]	28.4	21.0	22.5 [‡]
FT5-XXL	24.9	22.9	20.3	31.9	34.7 [†]	41.7 ^{†,‡}	23.8	16.9	22.2 [‡]
Avg.	24.0	21.9	19.7	34.3	34.6 [†]	36.8 ^{†,‡}	22.0	17.8	19.7 [‡]
Overall Avg	25.5	24.7	24.0	29.1	31.7	34.2 ^{†,‡}	23.6	26.8 [†]	28.2 ^{†,‡}

Table 16: Spearman correlations on SFRES dataset for data-to-text task. VAL, IST, and IDM denote the evaluator with vanilla, instruction, and the combination of instruction and demonstration, respectively. Values with † denote the evaluator with instruction significantly outperforms with vanilla, and values with ‡ denote the evaluator with the combination of instruction and demonstration significantly outperforms with only instruction. Values in bold are the best performance in a set of variants (e.g., GPT3 family).