

OpenING: Open INstruction Generation

指令微调数据生成工具

王越1,祈奕诚1,李俊涛1,张民1 1苏州大学, 计算机科学与技术学院

{lywangnlp,ycqi}@stu.suda.edu.cn,{ljt,minzhang}@suda.edu.cn

工具特点

·未使用闭源模型输出结果:

ChatGPT、Claude、Bard等闭源模型 禁止利用其模型输出开发同类模型,本 工具全流程未使用ChatGPT等闭源模型 的输出结果,没有潜在风险

·支持中英文多种任务指令生成:

通过优化训练数据分布,本工具既能够 生成文本分类,实体抽取等自然语言理 解任务指令,也能够生成故事、诗歌生 成等文本生成任务指令

·自动过滤高质量指令数据:

根据用户输入生成多个候选指令后,本 工具可以自动过滤出最高质量的指令, 并将最后的指令-输出对自动保存为json 格式,减少用户操作

流程介绍

·基座模型:Baichuan-7B (https://huggingface.co/baic huan-inc/Baichuan-7B)

·训练数据来源(共36511条): Super-NaturalInstructions (https://github.com/allenai/n atural-instructions)

Dolly

(https://huggingface.co/data sets/databricks/databricksdolly-15k)

Firefly

(https://huggingface.co/data sets/YeungNLP/fireflytrain-1.1M)

步骤一: 训练指令微调数据生成模型



步骤二: 生成候选指令微调数据



步骤三: 过滤候选指令微调数据

图一 OpenING指令微调数据生成流程示意图

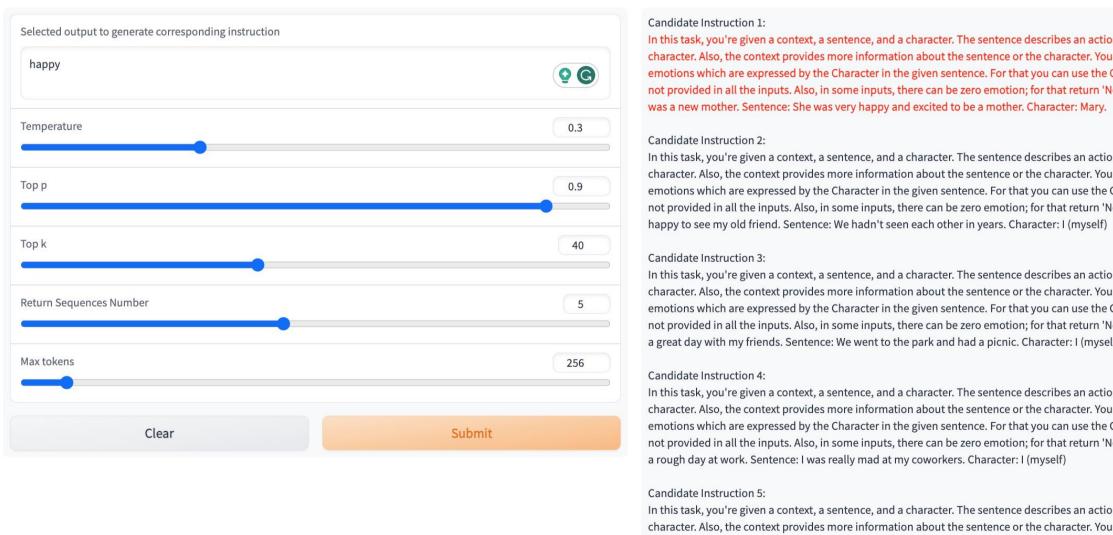
影响指令微调数据质量的三大因素:

输出回复的可靠性,输入指令的多样性,输入指令和输出回复的对齐性

本工具将用户输入的文本作为输出回复,生成对应的输入指令。因此输出回复的 可靠性由用户来控制,如果从语料库中采样文本输入到本工具中,那么回复的可 靠性就能够得到保证。此外,本工具先针对一个选中文本生成多个候选指令,再 使用过滤策略来选择最合适的指令,这也确保了输入指令的多样性,以及输入指 令和输出回复的对齐性。

生成样例

OpenING: Open INstruction Generation指令微调数据生成工具



In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'. ### Input: Context: Mary was a new mother. Sentence: She was very happy and excited to be a mother. Character: Mary. Candidate Instruction 2: In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'. ### Input: Context: I was

In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'. ### Input: Context: I had a great day with my friends. Sentence: We went to the park and had a picnic. Character: I (myself)

In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'. ### Input: Context: I had a rough day at work. Sentence: I was really mad at my coworkers. Character: I (myself)

Candidate Instruction 5: In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'. ### Input: Context: Cameron was a paramedic. He was called to a house to perform CPR on a man. Sentence: The man was talking to Cameron and seemed to be in good spirits. Character: Cameron.

Paired instruction and output in json format {"instruction": "In this task, you're given a context, a sentence, and a character. The sentence describes an action or job of the given character. Also, the context provides more information about the sentence or the character. Your task is to return one of the emotions which are expressed by the Character in the given sentence. For that you can use the Context; however, Context is not provided in all the inputs. Also, in some inputs, there can be zero emotion; for that return 'None'.\n\n### Input:\nContext: Mary was a new mother. \n Sentence: She was very happy and excited to be a mother. \n Character: Mary.","output":"happy"}

OpenING: Open INstruction Generation指令微调数据生成工具



我讲个小马过	ction 1: <mark>河的故事</mark>
andidate Instru 成故事,题目	
andidate Instru 成故事,题目	
andidate Instru 小马过河为标	ction 4: 题,生成一个故事。
andidate Instru : 成故事 题日	ction 5: : 哲理故事:河水很浅很浅
	on and output in json format