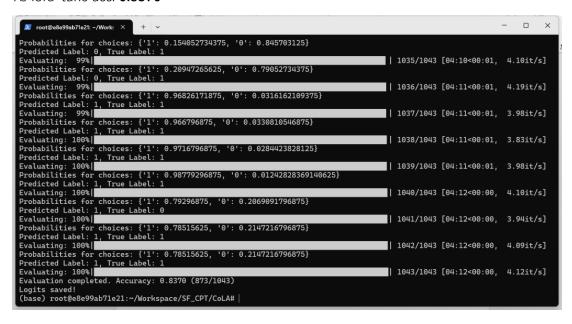benchmark

Table 1: **Comparison** of our CPT with other counterparts for black-box LLM tuning on six natural language datasets. We treat LLAMA2-7B as the small white-box model and treat LLAMA2-13B as the large black-box model. "pretrained" represents the zero-shot inference by their official pretrained parameters. "LORA-tuned" represents directly fine-tuning the corresponding model with LORA. Proxy-tuning [1] and CPT represent using a 7B model to "proxy fine-tune" a 13B model, where the 7B model is trained using their method and our method, respectively. "ARC-C" is the abbreviation of ARC-challenge.

| Model | Accuracy (%) ↑ | | | | | | Mean Acc (%) ↑ |
|---|---|---|---|---|---|---|---|
| | TriviaQA | ARC-C. | commonsenseQA | COLA | MRPC | AG-News | |
| LLAMA2-7B | | | | | | | |
| pretrained | 21.88 | 43.14 | 33.74 | 45.73 | 32.04 | 41.14 | 36.27 |
| LORA-tuned | 60.03 | 47.16 | 75.84 | 81.50 | 68.99 | 90.21 | 70.62 |
| LLAMA2-13B | | | | | | | |
| pretrained | 36.76 | 53.85 | 35.71 | 70.95 | 67.96 | 64.15 | 54.89 |
| Proxy-tuning [1] | 61.52 | 50.17 | 74.04 | 79.19 | 68.22 | 90.34 | 70.58 |
| **CPT (Ours)** | **62.79** | **55.85** | **76.41** | **82.26** | **69.77** | **90.91** | **72.99** |
| LORA-tuned | 66.58 | 66.22 | 81.90 | 84.65 | 68.99 | 90.65 | 76.49 |

7b pretrain acc: **0.5686**



13b pretrain acc: **0.6261**

7b lora-tune acc: **0.8370**



13b lora-tune acc: **0.8648**

Proxy tuning: **0.8303**



Consistent proxy tuning: **0.8495**

GP tune acc: **0.8696**

500 / 8551 (5.85%)



```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.008187909610569477, '0': 0.9918121099472046}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.008711383678019047, '0': 0.9912885427474976}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.9983769655227661, '0': 0.001622966956347227}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.007577241398394108, '0': 0.9924227595329285}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.005730246659368277, '0': 0.9942697286605835}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.04535258188843727, '0': 0.954647421836853}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.9977127313613892, '0': 0.002287227427586913}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.005220125894993544, '0': 0.9947799444198608}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.9983769655227661, '0': 0.001622966956347227}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9989182949066162, '0': 0.001081715221516788}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.998245358467102, '0': 0.0017546144081279635}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.9982725381851196, '0': 0.00172745855525136}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9899863600730896, '0': 0.01001356914639473}
Predicted Answer: 1, True Answer: 1
Evaluating: 100%|                                          | 1043/1043 [00:00<00:00, 2140.38it/s]
Evaluation completed. Accuracy: 0.8696 (907/1043)
(base) root@e8e99ab71e21:~/Workspace/SF_CPT/CoLA#
```

GP with filter accuracy: **0.8514**

Conifg: input_threshold=0.075, output_threshold=1.5

171/8551　(1.78%)



```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9314625263214111, '0': 0.06853749603033066}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.4687906503677368, '0': 0.531209409236908}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.9961155652999878, '0': 0.0038844768423587084}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.0446808747947216, '0': 0.9553191661834717}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.06008664891123772, '0': 0.9399133324623108}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.8397339582443237, '0': 0.16026602685451508}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9904406070709229, '0': 0.00955939944833517}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9740425944328308, '0': 0.02595735527575016}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9964619278907776, '0': 0.0035380860790610313}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9969722032546997, '0': 0.003027835162356496}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9946150183677673, '0': 0.005384937860071659}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.992061972618103, '0': 0.007937993854284286}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9908744096755981, '0': 0.00912563782185316}
Predicted Answer: 1, True Answer: 1
Evaluating: 100%|◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆| 1043/1043 [00:00<00:00, 1956.72it/s]
Evaluation completed. Accuracy: 0.8514 (888/1043)
(base) root@e8e99ab71e21:~/Workspace/SF_CPT/CoLA#
```