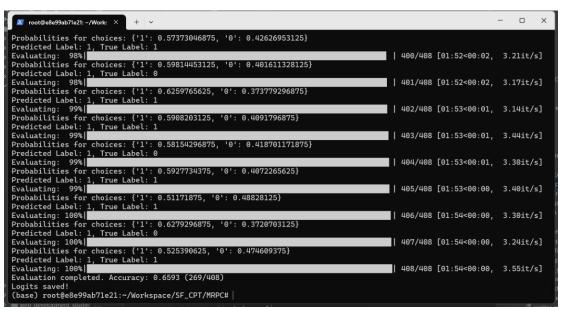
Table 1: Comparison of our CPT with other counterparts for black-box LLM tuning on six natural language datasets. We treat LLAMA2-7B as the small white-box model and treat LLAMA2-13B as the large black-box model. "pretrained" represents the zero-shot inference by their official pretrained parameters. "LORA-tuned" represents directly fine-tuning the corresponding model with LORA. Proxy-tuning [] 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 (%) ↑						Mann Ann (0/ )
	TriviaQA	ARC-C.	commonsenseQA	COLA	MRPC	AG-News	Mean Acc (%) ↑
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 []	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: 0.6593



13b pretrain acc: **0.6838** 

```
root@e8e99ab71e21: ~/Work: × + v
    robabilities for choices: {'1': 0.6689453125, '0': 0.3310546875}
redicted Label: 1, True Label: 1
valuating: 98%
robabilities for choices: {'1': 0.6533203125, '0': 0.346923828125}
redicted Label: 1, True Label: 0
                                                                                                                                                                                               | 400/408 [03:52<00:04, 1.61it/s]
    redicted Label: 1, True Label: (
valuating: 98%|
vobabilities for choices: {'1':
redicted Label: 1, True Label: 1
valuating: 99%|
vobabilities for choices: {'1':
redicted Label: 1, True Label: 1
valuating: 99%|
vobabilities for choices: {'1':
redicted Label: 1, True Label: 6
valuating: 99%|
voluating: 99%|
voluating: 99%|
voluating: 99%|
voluating: 99%|
voluating: 99%|
voluating: 99%|
                                                                                                                                                                                                | 401/408 [03:52<00:04, 1.57it/s]
                                                                                                                                                                                                | 402/408 [03:53<00:03, 1.54it/s]
                                                                                                                                                                                                | 403/408 [03:54<00:02, 1.70it/s]
      edicted Labet: 1, True Labet: 0
aluating: 99%|
bbabilities for choices: {'1': 0.64794921875, '0': 0.352294921875}
edicted Label: 1, True Label: 1
                                                                                                                                                                                                 | 404/408 [03:54<00:02, 1.62it/s]
        Dicted Label: 1, Frue Label: 1
luating: 99%|
babilities for choices: {'1': 0.6533203125, '0': 0.346923828125}
dicted Label: 1, True Label: 1
luating: 100%|
                                                                                                                                                                                                 | 405/408 [03:55<00:01, 1.67it/s]
                                                                                                                                                                                                  | 406/408 [03:56<00:01, 1.61it/s]
Evaluating: 100%|
Probabilities for choices: {'1': 0.66357421875, '0': 0.33642578125}
Predicted Label: 1, True Label: 0
Evaluating: 100%|
Probabilities for choices: {'1': 0.65673828125, '0': 0.34326171875}
Predicted Label: 1, True Label: 1
Evaluating: 100%|
Evaluation: completed. Accuracy: 0.6838 (279/408)
Louits saved!
                                                                                                                                                                                                  | 407/408 [03:56<00:00, 1.57it/s]
                                                                                                                                                                                                  408/408 [03:57<00:00. 1.72it/s]
 (base) root@e8e99ab71e21:~/Workspace/SF_CPT/MRPC#
```

## 13b loratune acc: **0.6985**

```
Predicted Label: 0, True Label: 1
Evaluating: 98% 400/408 [04:25<00:05, 1.47it/s] Probabilities for choices: {'1': 0.61865234375, '0': 0.381103515625}
Predicted Label: 1, True Label: 1
Evaluating: 98%
                           | 401/408 [04:26<00:04, 1.43it/s]
Probabilities for choices: {'1': 0.1710205078125, '0': 0.8291015625}
Predicted Label: 0, True Label: 0
Evaluating: 99% 402/408 [04:26<00:04, 1.41it/s]
Probabilities for choices: {'1': 0.85400390625, '0': 0.1461181640625}
Predicted Label: 1, True Label: 1
Evaluating: 99% 403/408 [04:27<00:03, 1.55it/s] Probabilities for choices: {'1': 0.47265625, '0': 0.52734375}
Predicted Label: 0, True Label: 1
                              404/408 [04:28<00:02, 1.48it/s]
Evaluating: 99%
Probabilities for choices: {'1': 0.216064453125, '0': 0.7841796875}
Predicted Label: 0, True Label: 0
Evaluating: 99% 405/408 [04:28<00:01, 1.52it/s]
Probabilities for choices: {'1': 0.1778564453125, '0': 0.822265625}
Predicted Label: 0, True Label: 1
Evaluating: 100% 406/408 [04:29<00:01, 1.47it/s]
Probabilities for choices: {'1': 0.57177734375, '0': 0.42822265625}
Predicted Label: 1, True Label: 1
                            | 407/408 [04:30<00:00, 1.43it/s]
Evaluating: 100%
Probabilities for choices: {'1': 0.1104736328125, '0': 0.8896484375}
Predicted Label: 0, True Label: 0
Evaluating: 100% 408/408 [04:30<00:00, 1.51it/s] Probabilities for choices: {'1': 0.505859375, '0': 0.494140625}
Predicted Label: 1, True Label: 1
Evaluation completed. Accuracy: 0.6985 (285/408)
Logits saved!
```

Proxy tuning acc: 0.6838

```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9207897186279297, '0': 0.07921033352613449}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9504110217094421, '0': 0.04958902671933174}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9504110217094421, '0': 0.11757213622331619}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9381249010261536, '0': 0.08882028609514236}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9381249010261536, '0': 0.06187598779797554}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9381249010261536, '0': 0.08882028609514236}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.9591543078422546, '0': 0.04084571450948715}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.9273632764816284, '0': 0.07263670861721039}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9525741338729858, '0': 0.04742587357759476}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.8791467847416687, '0': 0.1208532303571701}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9394582476615906, '0': 0.06954174488782883}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9241418242454529, '0': 0.19436781108379364}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9241418242454529, '0': 0.19436781108379364}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9241418242454529, '0': 0.075858183205127772}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9241418242454529, '0': 0.075858183205127772}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {\frac{1}{1}: 0.9241418242454529, '0': 0.075858183205127772}
Predicted Answer: 0.0688 (2794408)
(base) root@e8e9ab71e21:-/Workspace/SF_CPT/MRPC#
```

## Consistent proxy tuning acc: 0.7623

```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.815232515335083, '0': 0.18476751446723938}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8085806875228882, '0': 0.019419347867369652}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.4882833957672119, '0': 0.5117166042327881}
Predicted Answer: 0, True Answer: 1
Probabilities for choices: {'1': 0.7981867790222168, '0': 0.208132209777832}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.7981867790222168, '0': 0.208132209777832}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.396874415302277, '0': 0.8539127111434937}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.3963444841842651, '0': 0.064465348601341248}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.39981159567832947, '0': 0.608183149147034}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.39981159567832947, '0': 0.36296921968460083}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.6370307803153992, '0': 0.36296921968460083}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.84078969955444434, '0': 0.85921024829149246}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.8407896995544434, '0': 0.85921024829149246}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.847896995544434, '0': 0.85921024829149246}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.847896995544434, '0': 0.85921024829149246}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.84789699554455511322, '0': 0.5506073832511902}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.847896695547455642, '0': 0.11279541254043579}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.887204587455642, '0': 0.11279541254043579}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.887204587455642, '0': 0
```

```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.7341195344924927, '0': 0.2658804655075073}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9621075391769409, '0': 0.0378924235701561}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.50428633093833923, '0': 0.45713672041893005}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.50493178446298218, '0': 0.39606812596321106}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8047967803478241, '0': 0.15203224122524261}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.824967803478241, '0': 0.7371581792831421}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.9124361872673035, '0': 0.08756384253501892}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.8221890926361084, '0': 0.177810862660040802}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8221890926361084, '0': 0.4532618224620819}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.592543935775757, '0': 0.4532618224620819}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.5892543935775757, '0': 0.30074557662010193}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.581303060054779, '0': 0.44223188161849976}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.581303060054779, '0': 0.443624529838562}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.581303060054779, '0': 0.4033624529838562}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.5925439357768555, '0': 0.4033624529838562}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.592673724899292, '0': 0.4033624529838562}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.592663724899292, '0': 0.4033624529838562}
Predicted Answer: 0, True Answer: 0
Probabilities for choices: {'1': 0.592663724899292, '0': 0.40336245298
```

Gaussian filter tune acc: 0.8039

Config: input\_threshold=0.16, output\_threshold=0.45

Sample size: 86/3668 (2.34%)

```
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9399133324623108, '0': 0.06008664891123772}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9987181310653687, '0': 0.014281935058534145}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8918110132217407, '0': 0.10818894952535629}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.7310585975646973, '0': 0.2689414322376251}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.7909069833755493, '0': 0.2689414322376251}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.797477084826068878, '0': 0.9909299427270889}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.97477084826068878, '0': 0.92228618785738945}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.99777138829231262, '0': 0.09228618785738945}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9932942056655884, '0': 0.09670578688383102}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8976953029632568, '0': 0.10230471193790436}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.8976953029632568, '0': 0.10230471193790436}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9566442234611511, '0': 0.04336579889059067}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.9589832258224487, '0': 0.46101677417755127}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.951421918869019, '0': 0.04885777831077576}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.951421918869019, '0': 0.04885777831077576}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.951421918869019, '0': 0.04885777831077576}
Predicted Answer: 1, True Answer: 1
Probabilities for choices: {'1': 0.7589832258224487, '0': 0.46101677417755127}
Predicted Answer: 1, True Answer: 0
Probabilities for choices: {'1': 0.7589832258224487, '0'
```