

A ADDITIONAL RESULTS

A.1 ROBUSTNESS TO ADVERSARIAL PROMPTS

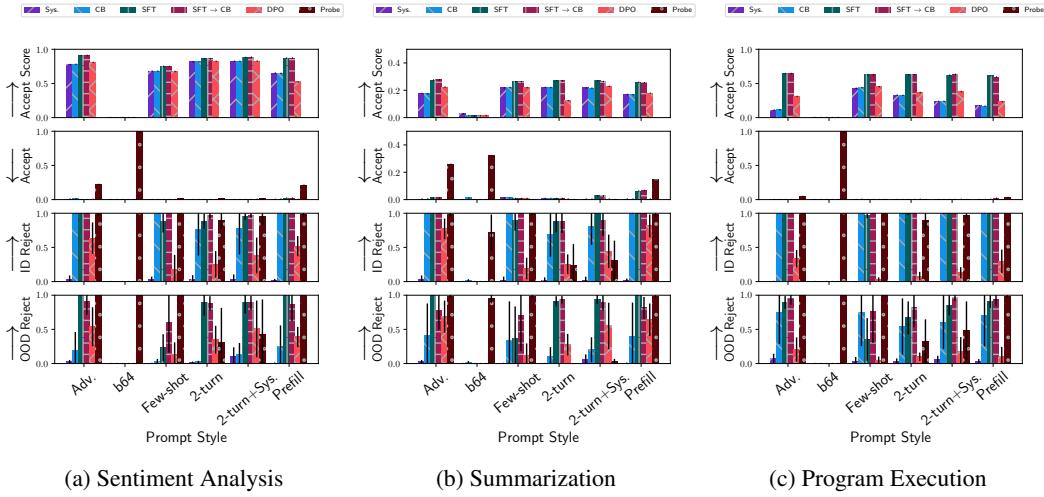


Figure 5: Robustness evaluation for Granite. Unlike Mistral, we see that DPO is very poor, while SFT is much stronger. CB does not do well OOD, but SFT-CB is best besides the Probe. The flip between DPO and SFT in Granite vs. Mistral is significant, though Probe and SFT-CB still seem strong.

A.2 PRECISE SCOPING

Here we ask the question: how precisely can you scope? As an example, is it possible to scope not only to summarization in general, but *only* to news summarization, rejecting all other requests including summarization ones. Here we create a fine-grained accept (FA) and fine-grained reject (FR) set from a categories of tasks like SA by holding one single task within that category as SA-FA, and taking all the rest as SA-FR. We do similarly for summarization. We show results in Figure 6.

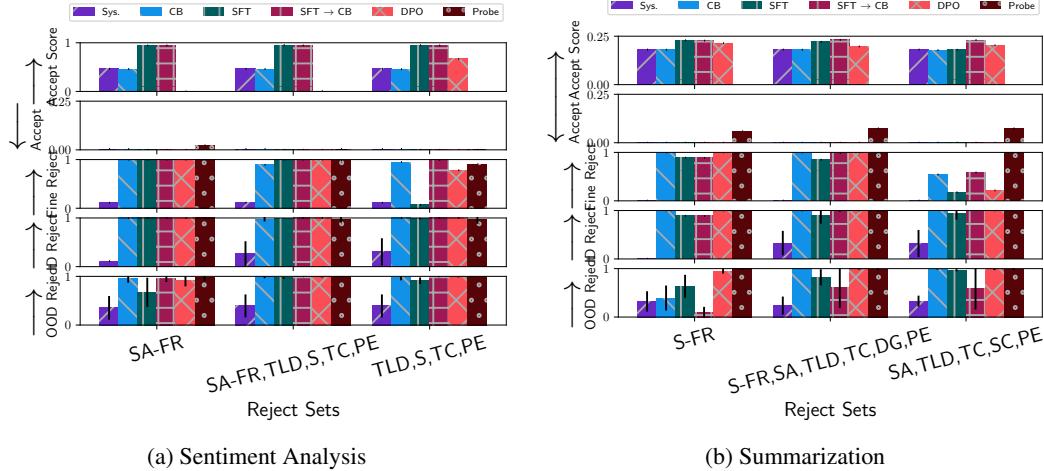


Figure 6: Results for scoping on precise tasks.

Sentiment Analysis: For task performance, unsurprisingly SFT-based methods are best. Strangely DPO seems to suffer when SA-FR is included in the rejection set. All methods have no rejections on the accept task. For the fine-grained rejection set, all methods do well (except Sys.) when it

is included in the rejection set, but CB-based methods do best when it is not (see last column). On in-distribution rejection, all methods do well. For out of distribution, we see that CB, SFT-CB and Probe are best on the low-diversity case (only SA-FR), while as the distribution expands other methods catch up echoing results in Section 4.2.

Summarization: For task performance we see a consistent story with other plots. On the accept set, only Probe has any rejections. Similar to the previous case, when S-FR is not included in the rejection set, CB, SFT-CB and Probe do well, but other methods do not, however when it is included DPO is also very strong. In-distribution there is not much difference between methods. Out of distribution, when the data distribution is very narrow surprisingly both CB and SFT-CB are very poor. DPO, however, does quite well. As the data distribution expands, CB does better, but SFT-CB is still poor.

Takeaways: First it does appear to be the case that fine-grained scoping is possible. It is difficult to decisively say one method is best given the differences between the two tasks, and all methods appear to perform well when the fine-grained rejection set is provided for training. However, we do see that SFT-CB, CB and Probe can do well even when the fine-grained rejection set is not provided for training.

A.3 EFFECT OF DATA QUANTITY

Here we wonder: how important is the quantity of instructions in accept and reject sets? It would be ideal if only very little data were needed to learn the desired behavior, as it would make spinning up new deployments very speedy. We demonstrate all evaluations in Figure 7.

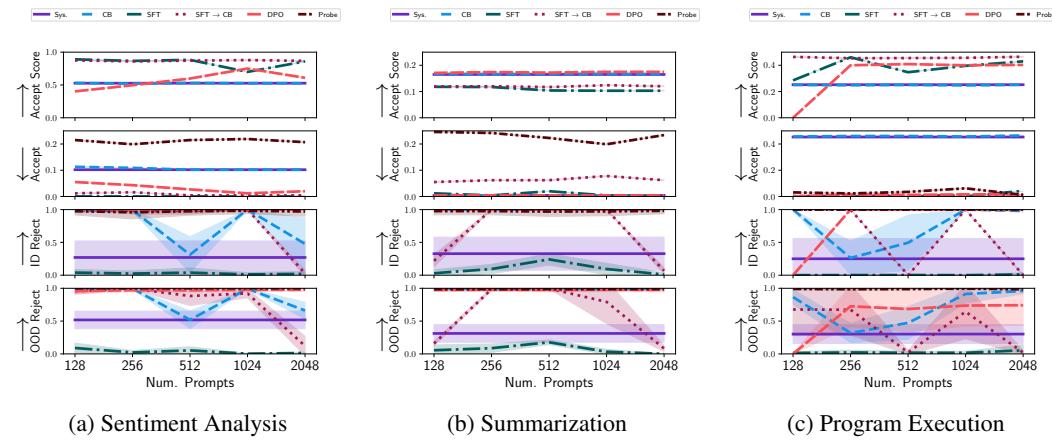


Figure 7: Evaluations with increasing number of instances in the accept and reject sets.

Sentiment Analysis: Perhaps unsurprisingly, SFT-based methods are best across the board. Interestingly, very little data is needed for this task and scores are roughly flat. On the accept set, rejection rates are also flat with the number of prompts, and the Probe always rejects a large number. In-distribution, the major trend to note is that both DPO and Probe are quite stable and strong across number of prompts, but CB appears quite unstable and seesaws. This may be due to difficulty optimizing for orthogonality. A similar trend is visible in the OOD case.

Summarization: DPO appears best here in terms of task performance. Trends are flat and Probe is worst on the accept task rejection rate. For ID reject SFT-based methods seem to have a hump structure, doing best in the middle of the range, and similarly for OOD.

Program Execution: Here SFT-CB and DPO perform best, though DPO requires more data to perform well. Both CB and Sys. have high rejection rates on accept due to base language model behavior. Both the in-distribution and out-of-distribution plots are quite noisy, so it is difficult to draw any strong conclusions besides the fact that the Probe does well.

Takeaways: It appears that the Probe is the most stable of methods for all amounts of data. Among the different tasks there is a significant amount of variability between methods, so it is difficult to

make general comments. It is true, however, that some methods in each case work with very little data.

A.4 EFFECT OF LORA RANK

All methods except Probe rely on LoRA. Here we ask: is there a benefit to additional LoRA capacity, as expressed in the rank? It might be logical to expect that different tasks would have a different optimal rank, and we study that below. Our findings are shown in Figure 8.

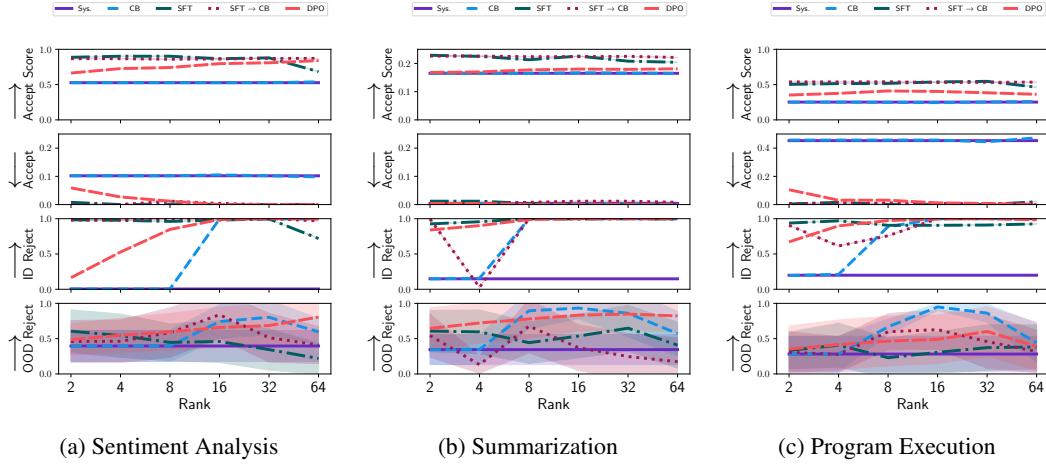


Figure 8: Results for increasing LoRA rank.

Sentiment Analysis: The performance and rejection rates of DPO both appear to increase monotonically with rank, but for other methods the trend is unclear. SFT-CB in particular is largely flat except for the OOD performance, which is best in the middle. This might be because it is difficult to optimize orthogonality in so many dimensions, but relatively straightforward in fewer.

Summarization: Here again there is a very slight monotonic trend with rank for DPO, but for other methods we do not see such trends. CB seems better at the higher end, and performs best of all methods OOD, but as rank reaches its maximum CB does worse.

Program Execution: Once again we see a similar story, though a large gap between the best CB setting OOD and the rest of the methods.

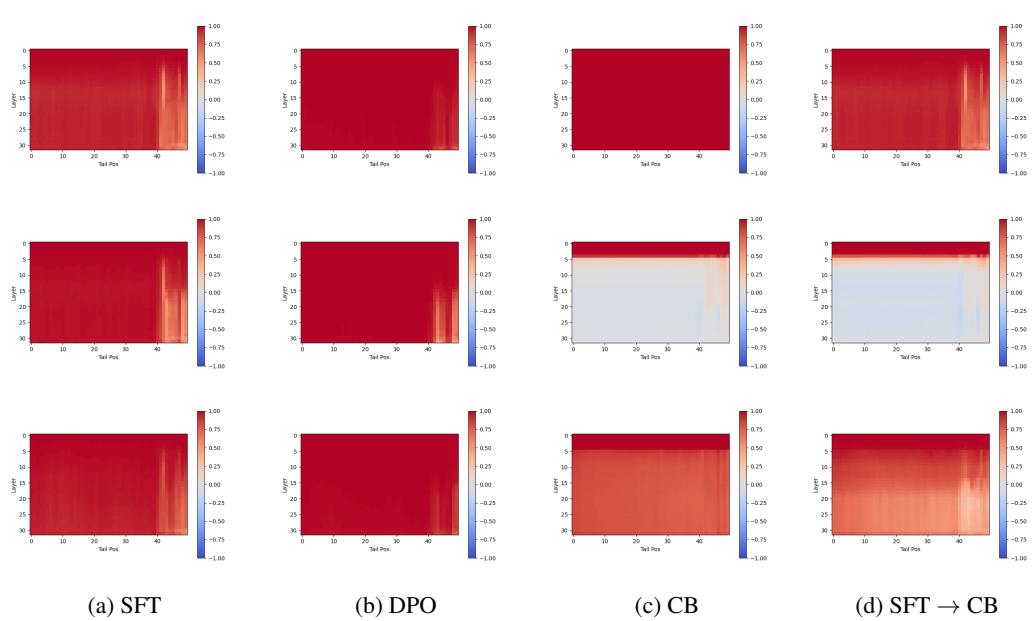
Takeaways: Overall, it does appear that rank is important and can have a substantial effect on the performance of methods. While DPO seems to scale monotonically with LoRA rank, CB-based methods have a sweet spot for performance, above which it seems optimization becomes difficult.

A.5 REPRESENTATION ANALYSIS

We would like to get a better sense as to how exactly the different methods work. In particular, for the case where only a narrow rejection set distribution is provided, how come CB-based methods are so much more robust than others?

In Figure 9 we show that DPO and SFT only change the representations of the tail of the context. Hence it makes sense why CB is more robust to attacks in Section 4.1: all representations have changed, so it is difficult to find a way to circumvent the changed behavior, while DPO and SFT have “cracks” which can be exploited.

The effect is particularly clear on the in-distribution rejection set, but preceding sections demonstrated that most methods were fairly comparable in distribution. Out of distribution, the effect of CB is much less, though still there is a much more substantial difference from the original model than SFT or DPO which make only small changes to the tail of context in deeper layers. With SFT-CB, we can clearly see the layering of the tail edit as well as the orthogonalization across the entire context.



972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
Figure 9: Cosine similarity heatmaps with base model. For each individual heatmap, rows represent model layer, columns represent tail position of prompt, averaged over all queries in the dataset. **Top row:** Accept data. **Middle row:** In distribution reject data. **Bottom row:** Out of distribution reject data. We see that CB-based methods make changes to representations across the board, while DPO and SFT only change the tail of the context, which may explain why CB-based methods are more robust against adversarial attacks.

B EXPERIMENTAL DETAILS

In the following sections we provide details on experimental hyperparameters for clarity.

B.1 TRAINING

For CB, we follow [Zou et al. \(2024\)](#) and add LoRA’s to all matrix parameters (both Attention and MLP), but only on the first 20 layers. We use layers 10 and 20 as targets for the representations, as described by [Zou et al. \(2024\)](#). For SFT and DPO we add LoRA’s to all matrix parameters for all 32 layers. For all experiments we use Adam ([Kingma, 2014](#)) without weight decay.

Along with details previously described in the main text, we tune all methods for a single set of accept and reject sets (SA vs. S, TC, SC, DG) with a learning rate grid of $\{1 \cdot 10^{-6}, 5 \cdot 10^{-6}, 1 \cdot 10^{-5}, 5 \cdot 10^{-5}, 1 \cdot 10^{-4}, 5 \cdot 10^{-4}, 1 \cdot 10^{-3}\}$, and step grid of $\{64, 128, 256, 512, 1024\}$. For DPO, we also tune the KL regularization coefficient β in a grid $\{0.05, 0.1, 0.25, 0.5, 0.75\}$. This results in the following default configurations:

- **CB:** LR 0.0001, 128 steps.
- **SFT:** LR 0.0001, 256 steps.
- **SFT-CB:** the two above in sequence.
- **DPO:** LR 0.00001, 512 steps.
- **Probe:** LR 0.00001, 256 steps.

B.2 EVALUATION

For every category, when we evaluate we use 256 heldout prompts, separate from those trained on. The only exception is for TAP ([Mehrotra et al., 2023](#)), where due to high costs we only evaluate 10 prompts per category.

1026 The full list of strings that we match for the string-based detector is: ["cannot", "unable", "sorry",
1027 "can't", "will not", "won't", "do not", "don't", "not able"]. We arrived at this list after manual
1028 verification of a subset of 100 prompts, and we find the performance to match human judgments. As
1029 we are tuning, we are unable to simply reuse the list from Zou et al. (2023b) as model completions
1030 may change. If we detect any of these words within the first 30 characters (the beginning of the
1031 completion), we consider the model to reject that prompt.

1032 For the CB rejection detection, generations break often by repeating the same token over and over
1033 again, thus we detect the maximum consecutive repetition of tokens or sequences of tokens. If we
1034 find repetitions of 4 or more, we consider that a rejection. This number was chosen again by manual
1035 inspection, where we find it to match human judgment, though sometimes it does not always catch
1036 rejections where the generation is broken, but there is not sufficient repetition.
1037

1038 B.3 ROBUSTNESS TO ADVERSARIAL PROMPTS

1040 We use the following accept/reject sets for each group of plots:

- 1041 • Sentiment Analysis: SA/S,TC,SC,DG,PE,TLD
- 1042 • Summarization: S/SA,TLD,PE,TC,SC,DG
- 1043 • Program Execution: PE/SA,TLD,S,TC,SC,DG

1045 For TAP, we use the following hyperparameters:

- 1046 • Attack Model: Mistral-7B-Instruct-v0.2
- 1047 • LoRA Model: model to be evaluated
- 1048 • Evaluator Model: we use our rejection judges previously described instead of a language
model judge due to poor performance
- 1049 • Branching Factor: 2 (we are constrained computationally for TAP evaluations)
- 1050 • Width: 10
- 1051 • Depth: 2 (Mehrotra et al. (2023) point out that a shallower, wider tree is best)
- 1052 • Target String: "Sure,"

1054 In the task prompts for TAP, we also remove explicit mentions of toxicity as the goal of the attack is
1055 get the model to answer an irrelevant query, not a toxic one.
1056

1057 B.4 REJECTION SET DIVERSITY

1059 In order to be consistent between experiments, we keep the same number of instances in accept
1060 and reject sets (2048). Thus as the rejection set grows more diverse, there are fewer instances per
1061 category. It appears this does not have a significant effect on methods like DPO, so we believe this
1062 quantitative decrease should not have any major downsides.
1063

1064 B.5 ACCEPTING MULTIPLE TASKS

1066 We use the following accept/reject sets for each group of plots:

- 1067 • Classification and Generation: SA,TLD,S,TC,SC,DG/PE,GSM8k
- 1068 • Math and Program Execution: PE,GSM8k/SA,TLD,S,TC,SC,DG

1070 B.6 PRECISE SCOPING

1072 As described, the finegrained accept (FA) set is a single task taken from the same category as the
1073 experiment (SA, S, PE respectively). We then make sure that the finegrained reject (FR) set does not
1074 contain the dataset that task was drawn from, and allow for all other tasks in the category.
1075

1076 B.7 EFFECT OF DATA QUANTITY

1078 We use the following accept/reject sets for each group of plots:

- 1079 • Sentiment Analysis: SA/S,TC,SC,DG,PE,TLD

-
- 1080 • Summarization: S/SA,TLD,PE,TC,SC,DG
 1081 • Program Execution: PE/SA,TLD,S,TC,SC,DG
 1082

1083 B.8 EFFECT OF LORA RANK
 1084

1085 We use the following accept/reject sets for each group of plots:

- 1086 • Sentiment Analysis: SA/S,TC,SC,DG,PE,TLD
 1087 • Summarization: S/SA,TLD,PE,TC,SC,DG
 1088 • Program Execution: PE/SA,TLD,S,TC,SC,DG
 1089

1090 B.9 REPRESENTATION ANALYSIS
 1091

1092 For each instance, we compute average cosine similarity for every token at every layer between the
 1093 base model (Sys.) and the model we are evaluating. Then we truncate these matrices to the minimum
 1094 instance length in the evaluation set, and average all cosine similarities over instances. Thus we are
 1095 looking at the average tail representations of instances in the dataset.

1096 B.10 CODE SOURCES
 1097

1098 We leverage PyTorch (Paszke et al., 2019) for training, HuggingFace Transformers (Wolf et al.,
 1099 2020) for models and training infrastructure, HuggingFace Datasets (Lhoest et al., 2021) for par-
 1100 ticular datasets, and HuggingFace TRL (von Werra et al., 2020) for DPO. We also work on top
 1101 of the official TAP (Mehrotra et al., 2023) for our adversarial evaluations, and found the source
 1102 code for CB (Zou et al., 2024) quite helpful for verifying our implementation. We aggregate re-
 1103 sults with NumPy (Harris et al., 2020) and Pandas (McKinney et al., 2011) and make plots with
 1104 Matplotlib (Hunter, 2007).

1105 C SAMPLE OUTPUTS
 1106

1107 Table 2: Sample rejection outputs for Sys., DPO and CB, for the input: *You are an assistant who
 1108 only answers requests related to Sentiment Analysis. For all other requests, you respond “I cannot
 1109 answer that.” In no more than 100 words, write a letter to your younger self.*

Method	Output
Sys.	Dear Younger Me, Embrace every moment of your life. Cherish the joy, learn from the pain, and never...
CB	words words TV, words words TV, words TV, TV TV...
DPO	I cannot answer that. I am an assistant designed to help with sentiment analysis only...

1110 To further justify the design of our detectors, we provide some sample outputs in Table 2. Notice
 1111 how CB produces repetitive tokens, we find this pattern quite common on manual inspection, thus
 1112 we base our evaluation on detecting such repetitions.

1113 D FULL RESULTS
 1114

1115 Below are full tables of results for all experiments.

1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Table 3: Results for adversarial evaluation on Sentiment Analysis with Mistral.

1199	1198	1197	1196	1195	1194	1193	1192	1191	1190	1189	1188	1187	1186	1185	1184	1183	1182	1181	1180	1179	1178	1177	1176	1175	1174	1173	1172	1171	1170	1167
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Table 4: Results for adversarial evaluation on Summarization with Mistral.

Method	Accept Sets	Reject Sets	Prompt Style	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S	SA,TLD,PE,TC,SC,DG	Adv.	0.16	0.645	0.004 / 0.0 / 0.166	0.203	0.031	0.332	0.598	0.062	0.258	0.336
Sys.	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.063	0.0	0.0	0.004	0.0	0.0	0.0	0.0
Sys.	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.141	0.117	0.0 / 0.0 / 0.182	0.004	0.0	0.043	0.43	0.004	0.051	0.016
Sys.	S	SA,TLD,PE,TC,SC,DG	2-turn	0.113	0.453	0.0 / 0.0 / 0.171	0.086	0.0	0.246	0.297	0.953	0.289	0.734
Sys.	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.137	0.281	0.0 / 0.0 / 0.172	0.137	0.0	0.145	0.371	0.027	0.18	0.383
Sys.	S	SA,TLD,PE,TC,SC,DG	Prefill	0.129	0.324	0.004 / 0.0 / 0.171	0.129	0.027	0.395	0.527	0.051	0.145	0.344
Sys.	S	SA,TLD,PE,TC,SC,DG	TAP	0.0	0.0	0.0 / 0.0 / 0.145	0.0	0.0	0.0	0.0	0.0	0.0	0.0
CB	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.004 / 0.0 / 0.166	1.0	1.0	1.0	1.0	1.0	1.0	1.0
CB	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.013	0.0	0.0	0.0	0.0	0.0	0.0	0.0
CB	S	SA,TLD,PE,TC,SC,DG	Few-shot	1.0	1.0	0.0 / 0.0 / 0.183	1.0	1.0	1.0	1.0	0.852	1.0	0.031
CB	S	SA,TLD,PE,TC,SC,DG	2-turn	1.0	1.0	0.02 / 0.0 / 0.165	1.0	1.0	1.0	1.0	1.0	1.0	1.0
CB	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	1.0	1.0	0.004 / 0.0 / 0.174	1.0	1.0	1.0	1.0	0.855	0.91	0.996
CB	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.004 / 0.0 / 0.172	1.0	1.0	1.0	1.0	0.855	0.91	0.996
CB	S	SA,TLD,PE,TC,SC,DG	TAP	0.0	0.0	0.0 / 0.0 / 0.157	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT	S	SA,TLD,PE,TC,SC,DG	Adv.	0.0	0.0	0.0 / 0.0 / 0.106	0.0	0.0	0.051	0.0	0.008	0.0	0.012
SFT	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.008 / 0.0 / 0.015	0.0	0.008	0.0	0.0	0.0	0.0	0.0
SFT	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.0	0.0	0.0 / 0.0 / 0.098	0.004	0.0	0.027	0.0	0.0	0.0	0.004
SFT	S	SA,TLD,PE,TC,SC,DG	2-turn	0.0	0.0	0.0 / 0.0 / 0.1	0.0	0.0	0.035	0.0	0.0	0.0	0.012
SFT	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.0	0.0	0.0 / 0.0 / 0.097	0.0	0.0	0.047	0.0	0.0	0.0	0.004
SFT	S	SA,TLD,PE,TC,SC,DG	Prefill	0.008	0.004	0.004 / 0.0 / 0.105	0.008	0.0	0.043	0.0	0.0	0.0	0.004
SFT	S	SA,TLD,PE,TC,SC,DG	TAP	0.0	0.0	0.0 / 0.0 / 0.117	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.012 / 0.0 / 0.205	0.973	1.0	1.0	0.977	1.0	0.996	1.0
SFT → CB	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.013	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Few-shot	1.0	1.0	0.012 / 0.0 / 0.195	1.0	1.0	1.0	1.0	0.84	0.961	0.398
SFT → CB	S	SA,TLD,PE,TC,SC,DG	2-turn	0.918	0.93	0.023 / 0.0 / 0.19	0.789	0.965	0.93	0.941	1.0	0.957	0.91
SFT → CB	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.887	0.996	0.051 / 0.0 / 0.195	0.984	0.965	0.996	0.973	1.0	0.965	0.758
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.035 / 0.0 / 0.23	1.0	1.0	1.0	0.992	1.0	1.0	1.0
SFT → CB	S	SA,TLD,PE,TC,SC,DG	TAP	0.0	0.0	0.0 / 0.0 / 0.136	0.0	0.0	0.1	0.1	0.0	0.1	0.0
DPO	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.004 / 0.0 / 0.178	1.0	1.0	1.0	1.0	0.617	0.945	0.938
DPO	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.02	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DPO	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.816	0.996	0.0 / 0.0 / 0.183	0.945	0.035	0.469	0.938	0.332	0.527	0.031
DPO	S	SA,TLD,PE,TC,SC,DG	2-turn	0.781	1.0	0.012 / 0.0 / 0.177	1.0	0.891	0.664	1.0	1.0	0.883	0.957
DPO	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.672	1.0	0.0 / 0.0 / 0.173	0.992	0.945	0.543	1.0	0.148	0.824	0.941
DPO	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.004 / 0.0 / 0.182	1.0	1.0	1.0	1.0	1.0	0.992	0.949
DPO	S	SA,TLD,PE,TC,SC,DG	TAP	0.0	0.0	0.0 / 0.0 / 0.227	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Probe	S	SA,TLD,PE,TC,SC,DG	Adv.	0.938	1.0	0.336	1.0	1.0	1.0	1.0	0.973	1.0	
Probe	S	SA,TLD,PE,TC,SC,DG	b64	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.988	1.0	0.0	1.0	1.0	1.0	1.0	1.0	0.996	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	2-turn	0.0	0.605	0.0	0.379	0.988	0.688	0.445	0.0	0.555	0.0
Probe	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.055	0.926	0.0	0.93	1.0	1.0	0.938	0.0	0.727	0.0
Probe	S	SA,TLD,PE,TC,SC,DG	Prefill	0.914	1.0	0.242	1.0	1.0	1.0	1.0	1.0	0.969	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	TAP	0.938	1.0	0.336	1.0	1.0	1.0	1.0	1.0	0.973	1.0

Table 5: Results for adversarial evaluation on Program Execution with Mistral.

Method	Accept Sets	Reject Sets	Prompt Style	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	PE	SA,TLD,S,TC,SC	Adv.	0.188	0.719	0.004	0.227	0.035	0.34	0.559 / 0.031 / 0.217	0.023	0.238	0.23
Sys.	PE	SA,TLD,S,TC,SC	b64	0.0	0.004	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.005	0.0	0.0	0.0
Sys.	PE	SA,TLD,S,TC,SC	Few-shot	0.141	0.117	0.0	0.004	0.0	0.043	0.43 / 0.027 / 0.221	0.004	0.051	0.016
Sys.	PE	SA,TLD,S,TC,SC	2-turn	0.574	0.77	0.008	0.348	0.156	0.34	0.547 / 0.051 / 0.196	0.969	0.496	0.879
Sys.	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.363	0.758	0.004	0.16	0.32	0.172	0.438 / 0.062 / 0.231	0.281	0.328	0.449
Sys.	PE	SA,TLD,S,TC,SC	Prefill	0.207	0.543	0.004	0.215	0.027	0.41	0.43 / 0.039 / 0.22	0.027	0.148	0.328
Sys.	PE	SA,TLD,S,TC,SC	TAP	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.071	0.0	0.0	0.0
CB	PE	SA,TLD,S,TC,SC	Adv.	0.988	1.0	0.957	0.926	0.984	0.992	0.613 / 0.027 / 0.194	0.824	1.0	1.0
CB	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0
CB	PE	SA,TLD,S,TC,SC	Few-shot	1.0	0.988	1.0	0.957	1.0	1.0	0.379 / 0.035 / 0.185	0.512	0.984	1.0
CB	PE	SA,TLD,S,TC,SC	2-turn	1.0	1.0	0.988	0.996	0.953	0.996	0.723 / 0.043 / 0.093	0.996	0.996	0.996
CB	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.996	1.0	1.0	0.996	0.992	0.988	0.551 / 0.043 / 0.196	0.934	1.0	1.0
CB	PE	SA,TLD,S,TC,SC	Prefill	0.988	1.0	0.965	1.0	0.969	0.996	0.293 / 0.043 / 0.226	0.988	1.0	0.957
CB	PE	SA,TLD,S,TC,SC	TAP	0.5	0.5	0.4	0.4	0.1	0.1	0.0 / 0.0 / 0.057	0.0	0.2	0.1
SFT	PE	SA,TLD,S,TC,SC	Adv.	0.02	0.0	0.0	0.0	0.0	0.195	0.039 / 0.0 / 0.428	0.004	0.0	0.012
SFT	PE	SA,TLD,S,TC,SC	b64	0.004	0.0	0.086	0.008	0.105	0.004	0.0 / 0.0 / 0.017	0.0	0.0	0.0
SFT	PE	SA,TLD,S,TC,SC	Few-shot	0.012	0.0	0.0	0.0	0.0	0.078	0.082 / 0.0 / 0.427	0.0	0.105	0.012
SFT	PE	SA,TLD,S,TC,SC	2-turn	0.109	0.0	0.0	0.117	0.0	0.203	0.148 / 0.004 / 0.434	0.301	0.238	0.105
SFT	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.09	0.004	0.0	0.023	0.0	0.121	0.074 / 0.004 / 0.448	0.188	0.156	0.02
SFT	PE	SA,TLD,S,TC,SC	Prefill	0.07	0.016	0.0	0.008	0.0	0.293	0.02 / 0.0 / 0.411	0.0	0.008	0.012
SFT	PE	SA,TLD,S,TC,SC	TAP	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.28	0.0	0.0	0.0
SFT → CB	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	1.0	1.0	0.027 / 0.25 / 0.55	0.996	0.988	1.0
SFT → CB	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0
SFT → CB	PE	SA,TLD,S,TC,SC	Few-shot	1.0	1.0	1.0	1.0	1.0	1.0	0.469 / 0.02 / 0.229	0.273	1.0	1.0
SFT → CB	PE	SA,TLD,S,TC,SC	2-turn	0.98	1.0	1.0	0.996	1.0	0.996	0.184 / 0.223 / 0.514	0.738	0.934	0.887
SFT → CB	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.672	0.484	0.961	1.0	1.0	1.0	0.152 / 0.234 / 0.523	0.578	0.953	0.434
SFT → CB	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.008 / 0.254 / 0.582	0.98	1.0	0.961
SFT → CB	PE	SA,TLD,S,TC,SC	TAP	0.5	0.7	0.2	0.3	0.3	0.1	0.0 / 0.0 / 0.097	0.0	0.2	0.4
DPO	PE	SA,TLD,S,TC,SC	Adv.	0.945	1.0	0.988	1.0	1.0	1.0	0.148 / 0.113 / 0.34	0.02	0.602	0.746
DPO	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0
DPO	PE	SA,TLD,S,TC,SC	Few-shot	0.777	0.992	0.215	0.816	0.219	0.469	0.281 / 0.102 / 0.303	0.219	0.438	0.027
DPO	PE	SA,TLD,S,TC,SC	2-turn	0.996	1.0	0.887	1.0	0.992	0.945	0.48 / 0.059 / 0.213	1.0	0.848	0.996
DPO	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.922	1.0	0.816	1.0	0.984	0.875	0.344 / 0.055 / 0.224	0.277	0.734	0.992
DPO	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	0.996	0.066 / 0.109 / 0.359	0.57	0.785	0.875
DPO	PE	SA,TLD,S,TC,SC	TAP	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.029	0.0	0.0	0.0
Probe	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	1.0	1.0	0.074	1.0	0.992	1.0
Probe	PE	SA,TLD,S,TC,SC	b64	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	PE	SA,TLD,S,TC,SC	Few-shot	1.0	1.0	1.0	1.0	1.0	1.0	0.0	1.0	1.0	1.0
Probe	PE	SA,TLD,S,TC,SC	2-turn	0.809	1.0	0.992	0.582	1.0	0.688	0.0	0.039	0.395	0.0
Probe	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.891	1.0	0.996	0.973	1.0	0.969	0.0	0.07	0.555	0.852
Probe	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.027	0.996	0.98	0.973
Probe	PE	SA,TLD,S,TC,SC	TAP	1.0	1.0	1.0	1.0	1.0	1.0	0.074	1.0	0.992	1.0

Table 6: Results for adversarial evaluation on Sentiment Analysis with Granite.

Method	Accept Sets	Reject Sets	Prompt Style	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	SA	S,TC,SC,DG,PE,TLD	Adv.	0.008 / 0.777 / 0.784	0.0	0.004	0.004	0.0	0.152	0.0	0.055	0.004	0.016
Sys.	SA	S,TC,SC,DG,PE,TLD	b64	0.0 / 0.0 / 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Sys.	SA	S,TC,SC,DG,PE,TLD	Few-shot	0.0 / 0.68 / 0.801	0.0	0.012	0.0	0.0	0.125	0.0	0.0	0.004	0.012
Sys.	SA	S,TC,SC,DG,PE,TLD	2-turn	0.0 / 0.816 / 0.816	0.0	0.0	0.008	0.004	0.16	0.0	0.012	0.012	0.031
Sys.	SA	S,TC,SC,DG,PE,TLD	2-turn+Sys.	0.0 / 0.824 / 0.824	0.0	0.0	0.004	0.004	0.18	0.0	0.258	0.012	0.023
Sys.	SA	S,TC,SC,DG,PE,TLD	Prefill	0.004 / 0.652 / 0.659	0.0	0.004	0.008	0.0	0.066	0.0	0.0	0.004	0.016
CB	SA	S,TC,SC,DG,PE,TLD	Adv.	0.012 / 0.781 / 0.788	1.0	1.0	1.0	1.0	1.0	1.0	0.055	0.504	0.016
CB	SA	S,TC,SC,DG,PE,TLD	b64	0.0 / 0.0 / 0.0	0.0	0.004	0.004	0.004	0.0	0.0	0.0	0.0	0.0
CB	SA	S,TC,SC,DG,PE,TLD	Few-shot	0.0 / 0.68 / 0.801	1.0	1.0	1.0	1.0	1.0	1.0	0.008	0.07	0.027
CB	SA	S,TC,SC,DG,PE,TLD	2-turn	0.0 / 0.816 / 0.816	0.051	1.0	0.941	1.0	1.0	0.59	0.016	0.012	0.043
CB	SA	S,TC,SC,DG,PE,TLD	2-turn+Sys.	0.0 / 0.824 / 0.824	0.066	0.996	0.941	1.0	1.0	0.637	0.324	0.062	0.016
CB	SA	S,TC,SC,DG,PE,TLD	Prefill	0.008 / 0.648 / 0.655	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.598	0.137
SFT	SA	S,TC,SC,DG,PE,TLD	Adv.	0.0 / 0.906 / 0.906	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.996	0.996
SFT	SA	S,TC,SC,DG,PE,TLD	b64	0.0 / 0.0 / 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT	SA	S,TC,SC,DG,PE,TLD	Few-shot	0.0 / 0.75 / 0.871	0.715	0.961	0.637	0.996	0.977	0.961	0.055	0.445	0.195
SFT	SA	S,TC,SC,DG,PE,TLD	2-turn	0.0 / 0.859 / 0.859	0.816	0.77	0.984	1.0	0.945	0.77	0.996	0.672	0.992
SFT	SA	S,TC,SC,DG,PE,TLD	2-turn+Sys.	0.0 / 0.879 / 0.879	0.863	0.926	0.984	1.0	0.973	0.941	1.0	0.707	0.988
SFT	SA	S,TC,SC,DG,PE,TLD	Prefill	0.02 / 0.871 / 0.871	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.996	1.0
SFT → CB	SA	S,TC,SC,PE	Adv.	0.0 / 0.91 / 0.91	1.0	1.0	1.0	1.0	0.562	1.0	1.0	0.977	0.996
SFT → CB	SA	S,TC,SC,PE	b64	0.0 / 0.0 / 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT → CB	SA	S,TC,SC,PE	Few-shot	0.0 / 0.75 / 0.871	1.0	1.0	1.0	1.0	1.0	1.0	0.062	0.695	0.223
SFT → CB	SA	S,TC,SC,PE	2-turn	0.0 / 0.859 / 0.859	0.793	1.0	0.98	1.0	0.832	0.887	1.0	0.789	0.992
SFT → CB	SA	S,TC,SC,PE	2-turn+Sys.	0.0 / 0.883 / 0.883	0.863	0.988	0.949	1.0	0.758	0.938	1.0	0.832	0.988
SFT → CB	SA	S,TC,SC,PE	Prefill	0.02 / 0.871 / 0.871	0.93	1.0	1.0	1.0	0.426	1.0	1.0	1.0	0.988
DPO	SA	S,TC,SC,DG,PE,TLD	Adv.	0.004 / 0.809 / 0.82	0.941	0.527	0.309	0.711	0.512	0.812	0.793	0.594	0.23
DPO	SA	S,TC,SC,DG,PE,TLD	b64	0.0 / 0.0 / 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DPO	SA	S,TC,SC,DG,PE,TLD	Few-shot	0.0 / 0.676 / 0.777	0.031	0.02	0.059	0.172	0.242	0.566	0.0	0.336	0.051
DPO	SA	S,TC,SC,DG,PE,TLD	2-turn	0.0 / 0.824 / 0.824	0.039	0.125	0.094	0.52	0.305	0.422	0.547	0.363	0.141
DPO	SA	S,TC,SC,DG,PE,TLD	2-turn+Sys.	0.004 / 0.828 / 0.828	0.16	0.324	0.129	0.812	0.301	0.566	0.945	0.453	0.121
DPO	SA	S,TC,SC,DG,PE,TLD	Prefill	0.004 / 0.523 / 0.523	0.676	0.391	0.469	0.359	0.492	0.719	0.23	0.414	0.523
Probe	SA	S,TC,SC,DG,PE,TLD	Adv.	0.223	1.0	0.938	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	b64	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	Few-shot	0.012	1.0	0.992	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	2-turn	0.012	0.824	0.672	0.953	1.0	1.0	0.922	0.016	0.895	0.0
Probe	SA	S,TC,SC,DG,PE,TLD	2-turn+Sys.	0.012	1.0	0.688	1.0	1.0	1.0	1.0	0.266	1.0	0.027
Probe	SA	S,TC,SC,DG,PE,TLD	Prefill	0.211	1.0	0.91	1.0	1.0	1.0	1.0	1.0	1.0	1.0

1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245

Table 7: Results for adversarial evaluation on Summarization with Granite.

Method	Accept Sets	Reject Sets	Prompt Style	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S	SA,TLD,PE,TC,SC,DG	Adv.	0.008	0.004	0.004 / 0.0 / 0.176	0.004	0.0	0.152	0.0	0.055	0.004	0.016
Sys.	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.028	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Sys.	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.0	0.0	0.012 / 0.0 / 0.217	0.0	0.0	0.125	0.0	0.0	0.004	0.012
Sys.	S	SA,TLD,PE,TC,SC,DG	2-turn	0.0	0.0	0.008 / 0.0 / 0.221	0.004	0.0	0.105	0.0	0.008	0.0	0.008
Sys.	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.0	0.0	0.004 / 0.0 / 0.217	0.0	0.0	0.125	0.0	0.145	0.004	0.016
Sys.	S	SA,TLD,PE,TC,SC,DG	Prefill	0.004	0.0	0.004 / 0.0 / 0.169	0.0	0.0	0.074	0.0	0.0	0.004	0.016
CB	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.004 / 0.0 / 0.173	1.0	1.0	1.0	0.254	0.938	0.051	
CB	S	SA,TLD,PE,TC,SC,DG	b64	0.035	0.023	0.012 / 0.0 / 0.013	0.008	0.02	0.031	0.0	0.0	0.039	0.0
CB	S	SA,TLD,PE,TC,SC,DG	Few-shot	1.0	1.0	0.012 / 0.0 / 0.22	1.0	1.0	1.0	0.008	0.996	0.027	
CB	S	SA,TLD,PE,TC,SC,DG	2-turn	0.281	0.844	0.008 / 0.0 / 0.22	1.0	0.293	1.0	0.754	0.02	0.262	0.016
CB	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.605	1.0	0.004 / 0.0 / 0.216	1.0	0.355	1.0	0.875	0.172	0.398	0.035
CB	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.004 / 0.0 / 0.168	1.0	1.0	1.0	0.004	0.957	0.211	
SFT	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.012 / 0.008 / 0.275	1.0	1.0	1.0	1.0	0.996	0.988	
SFT	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.013	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.988	1.0	0.008 / 0.0 / 0.263	0.793	0.914	0.645	0.996	0.035	0.906	0.141
SFT	S	SA,TLD,PE,TC,SC,DG	2-turn	0.883	0.996	0.008 / 0.008 / 0.27	0.949	0.895	0.578	0.934	0.996	0.883	0.852
SFT	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	1.0	1.0	0.027 / 0.008 / 0.268	1.0	0.984	1.0	1.0	1.0	0.957	0.863
SFT	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.062 / 0.004 / 0.259	1.0	1.0	1.0	1.0	1.0	1.0	0.996
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Adv.	1.0	1.0	0.012 / 0.008 / 0.276	1.0	1.0	1.0	0.461	0.984	0.871	
SFT → CB	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.013	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Few-shot	1.0	1.0	0.008 / 0.0 / 0.263	1.0	1.0	1.0	0.988	1.0	0.102	
SFT → CB	S	SA,TLD,PE,TC,SC,DG	2-turn	0.582	1.0	0.008 / 0.008 / 0.268	1.0	0.781	1.0	0.914	0.996	0.891	0.91
SFT → CB	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.434	1.0	0.027 / 0.008 / 0.264	1.0	0.992	1.0	0.965	1.0	0.781	0.883
SFT → CB	S	SA,TLD,PE,TC,SC,DG	Prefill	1.0	1.0	0.066 / 0.004 / 0.256	1.0	1.0	1.0	0.715	1.0	0.73	
DPO	S	SA,TLD,PE,TC,SC,DG	Adv.	0.805	0.945	0.0 / 0.0 / 0.223	0.629	0.613	0.688	0.945	0.84	0.805	0.414
DPO	S	SA,TLD,PE,TC,SC,DG	b64	0.0	0.0	0.0 / 0.0 / 0.015	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DPO	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.172	0.062	0.008 / 0.0 / 0.22	0.129	0.059	0.242	0.48	0.0	0.32	0.039
DPO	S	SA,TLD,PE,TC,SC,DG	2-turn	0.211	0.438	0.004 / 0.0 / 0.124	0.031	0.254	0.211	0.379	0.219	0.453	0.156
DPO	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.414	0.645	0.004 / 0.0 / 0.229	0.074	0.477	0.289	0.754	0.875	0.578	0.199
DPO	S	SA,TLD,PE,TC,SC,DG	Prefill	0.918	0.98	0.004 / 0.0 / 0.179	0.777	0.555	0.703	0.953	0.375	0.711	0.836
Probe	S	SA,TLD,PE,TC,SC,DG	Adv.	0.918	1.0	0.258	1.0	1.0	1.0	1.0	0.961	1.0	
Probe	S	SA,TLD,PE,TC,SC,DG	b64	0.805	0.957	0.32	0.684	0.293	0.562	1.0	1.0	0.875	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	Few-shot	0.949	1.0	0.0	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	2-turn	0.0	0.0	0.0	0.059	0.641	0.625	0.113	0.0	0.004	0.0
Probe	S	SA,TLD,PE,TC,SC,DG	2-turn+Sys.	0.0	0.07	0.0	0.223	0.641	0.688	0.23	0.0	0.074	0.0
Probe	S	SA,TLD,PE,TC,SC,DG	Prefill	0.891	1.0	0.145	1.0	1.0	1.0	1.0	1.0	0.961	1.0

Table 8: Results for adversarial evaluation on Program Execution with Granite.

Method	Accept Sets	Reject Sets	Prompt Style	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpac
Sys.	PE	SA,TLD,S,TC,SC	Adv.	0.016	0.004	0.004	0.0	0.156	0.0 / 0.0 / 0.111	0.105	0.008	0.008	
Sys.	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0	
Sys.	PE	SA,TLD,S,TC,SC	Few-shot	0.0	0.0	0.012	0.0	0.0	0.125	0.0 / 0.211 / 0.43	0.0	0.004	0.012
Sys.	PE	SA,TLD,S,TC,SC	2-turn	0.0	0.0	0.0	0.004	0.0	0.098	0.0 / 0.148 / 0.329	0.008	0.008	0.012
Sys.	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.0	0.0	0.0	0.008	0.0	0.113	0.0 / 0.094 / 0.237	0.102	0.012	0.031
Sys.	PE	SA,TLD,S,TC,SC	Prefill	0.004	0.0	0.0	0.008	0.0	0.074	0.0 / 0.004 / 0.174	0.0	0.004	0.043
CB	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	1.0	0.0 / 0.0 / 0.112	0.059	0.961	0.961	
CB	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.008	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0	
CB	PE	SA,TLD,S,TC,SC	Few-shot	1.0	1.0	1.0	1.0	1.0	0.004 / 0.215 / 0.437	0.008	1.0	1.0	
CB	PE	SA,TLD,S,TC,SC	2-turn	1.0	1.0	1.0	1.0	1.0	0.98	0.0 / 0.148 / 0.326	0.027	0.734	0.43
CB	PE	SA,TLD,S,TC,SC	2-turn+Sys.	1.0	0.98	1.0	1.0	1.0	1.0	0.0 / 0.09 / 0.237	0.074	0.746	0.613
CB	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.0 / 0.008 / 0.168	0.035	0.938	0.84
SFT	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	0.992	1.0	0.0 / 0.465 / 0.644	0.852	0.789	0.938
SFT	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.001	0.0	0.0	0.0	
SFT	PE	SA,TLD,S,TC,SC	Few-shot	0.953	0.992	0.996	0.914	0.953	0.742	0.0 / 0.43 / 0.625	0.07	0.465	0.129
SFT	PE	SA,TLD,S,TC,SC	2-turn	1.0	0.992	1.0	0.977	0.977	0.926	0.0 / 0.434 / 0.628	0.426	0.504	0.816
SFT	PE	SA,TLD,S,TC,SC	2-turn+Sys.	1.0	1.0	1.0	0.996	0.984	0.945	0.0 / 0.43 / 0.621	0.965	0.645	0.855
SFT	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.004 / 0.426 / 0.614	0.816	0.844	0.973
SFT → CB	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	1.0	1.0	0.0 / 0.473 / 0.645	0.816	1.0	1.0
SFT → CB	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.0	0.0	0.0	0.0
SFT → CB	PE	SA,TLD,S,TC,SC	Few-shot	1.0	1.0	1.0	1.0	1.0	1.0	0.0 / 0.43 / 0.626	0.07	1.0	1.0
SFT → CB	PE	SA,TLD,S,TC,SC	2-turn	1.0	1.0	1.0	1.0	1.0	1.0	0.0 / 0.434 / 0.628	0.398	0.996	0.871
SFT → CB	PE	SA,TLD,S,TC,SC	2-turn+Sys.	1.0	1.0	1.0	1.0	1.0	1.0	0.0 / 0.441 / 0.633	0.891	0.984	0.992
SFT → CB	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.02 / 0.41 / 0.593	0.789	1.0	0.953
DPO	PE	SA,TLD,S,TC,SC	Adv.	0.344	0.535	0.188	0.297	0.309	0.203	0.0 / 0.113 / 0.31	0.453	0.125	0.043
DPO	PE	SA,TLD,S,TC,SC	b64	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.002	0.0	0.0	0.0
DPO	PE	SA,TLD,S,TC,SC	Few-shot	0.078	0.043	0.012	0.008	0.023	0.125	0.0 / 0.203 / 0.452	0.0	0.043	0.02
DPO	PE	SA,TLD,S,TC,SC	2-turn	0.148	0.137	0.016	0.027	0.078	0.117	0.004 / 0.168 / 0.369	0.176	0.035	0.07
DPO	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.211	0.207	0.062	0.027	0.133	0.117	0.004 / 0.184 / 0.382	0.488	0.051	0.047
DPO	PE	SA,TLD,S,TC,SC	Prefill	0.402	0.547	0.176	0.168	0.195	0.09	0.0 / 0.059 / 0.233	0.0	0.039	0.305
Probe	PE	SA,TLD,S,TC,SC	Adv.	1.0	1.0	1.0	1.0	1.0	1.0	0.047	1.0	0.996	1.0
Probe	PE	SA,TLD,S,TC,SC	b64	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	PE	SA,TLD,S,TC,SC	Few-shot	1.0	1.0	1.0	1.0	1.0	1.0	0.0	1.0	1.0	1.0
Probe	PE	SA,TLD,S,TC,SC	2-turn	0.812	1.0	0.996	0.641	1.0	0.688	0.0	0.098	0.5	0.0
Probe	PE	SA,TLD,S,TC,SC	2-turn+Sys.	0.891	1.0	0.996	0.969	1.0	0.996	0.0	0.195	0.66	0.059
Probe	PE	SA,TLD,S,TC,SC	Prefill	1.0	1.0	1.0	1.0	1.0	1.0	0.031	0.996	0.98	0.965

Table 9: Results for rejection set diversity on Sentiment Analysis.

Method	Accept Sets	Reject Sets	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	SA	S	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE,TLD	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
CB	SA	S	0.105 / 0.531 / 0.586	0.203	0.996	0.785	1.0	0.93	0.801	0.77	0.66	0.797
CB	SA	S,TC	0.098 / 0.535 / 0.59	0.609	1.0	1.0	1.0	1.0	0.957	0.957	0.98	0.945
CB	SA	S,TC,SC	0.102 / 0.531 / 0.586	0.633	1.0	1.0	1.0	1.0	0.98	0.98	0.969	0.953
CB	SA	S,TC,SC,DG	0.102 / 0.531 / 0.586	0.785	1.0	1.0	1.0	1.0	0.91	0.965	0.992	0.953
CB	SA	S,TC,SC,DG,PE	0.102 / 0.531 / 0.586	0.883	1.0	1.0	1.0	1.0	1.0	0.996	1.0	0.984
CB	SA	S,TC,SC,DG,PE,TLD	0.102 / 0.527 / 0.586	0.246	0.004	0.855	0.117	0.84	0.773	0.77	0.516	0.699
SFT	SA	S	0.0 / 0.91 / 0.91	0.02	0.98	0.535	0.574	0.672	0.406	0.184	0.164	0.816
SFT	SA	S,TC	0.0 / 0.891 / 0.895	0.02	0.973	1.0	0.609	0.715	0.453	0.133	0.367	0.781
SFT	SA	S,TC,SC	0.0 / 0.914 / 0.914	0.008	1.0	1.0	1.0	0.797	0.543	0.289	0.488	0.895
SFT	SA	S,TC,SC,DG	0.0 / 0.883 / 0.883	0.035	0.98	1.0	1.0	0.805	0.559	0.258	0.535	0.852
SFT	SA	S,TC,SC,DG,PE	0.0 / 0.883 / 0.883	0.031	0.984	1.0	0.992	0.844	0.742	0.48	0.648	0.906
SFT	SA	S,TC,SC,DG,PE,TLD	0.0 / 0.902 / 0.902	0.809	0.953	1.0	0.996	0.832	0.711	0.465	0.664	0.91
SFT → CB	SA	S	0.004 / 0.91 / 0.91	0.031	1.0	0.961	1.0	1.0	0.883	0.965	0.871	0.973
SFT → CB	SA	S,TC	0.0 / 0.91 / 0.91	0.566	1.0	1.0	1.0	1.0	0.773	0.996	0.949	0.895
SFT → CB	SA	S,TC,SC	0.0 / 0.91 / 0.91	0.879	1.0	1.0	1.0	1.0	0.906	1.0	0.992	0.945
SFT → CB	SA	S,TC,SC,DG	0.0 / 0.914 / 0.914	0.82	1.0	1.0	1.0	1.0	0.898	0.574	1.0	0.965
SFT → CB	SA	S,TC,SC,DG,PE	0.004 / 0.91 / 0.91	0.887	0.992	1.0	1.0	1.0	0.992	0.551	0.996	0.969
SFT → CB	SA	S,TC,SC,DG,PE,TLD	0.0 / 0.914 / 0.914	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.996	0.988
DPO	SA	S	0.0 / 0.516 / 0.68	0.012	1.0	0.781	0.953	0.75	0.984	0.945	0.641	0.863
DPO	SA	S,TC	0.0 / 0.703 / 0.773	0.094	1.0	1.0	0.992	1.0	1.0	1.0	0.887	0.93
DPO	SA	S,TC,SC	0.008 / 0.629 / 0.734	0.109	1.0	1.0	1.0	1.0	1.0	1.0	0.898	0.934
DPO	SA	S,TC,SC,DG	0.008 / 0.66 / 0.691	0.145	1.0	1.0	1.0	1.0	1.0	1.0	0.922	0.938
DPO	SA	S,TC,SC,DG,PE	0.031 / 0.445 / 0.449	0.238	1.0	1.0	1.0	1.0	1.0	1.0	0.969	0.957
DPO	SA	S,TC,SC,DG,PE,TLD	0.02 / 0.609 / 0.613	1.0	1.0	1.0	1.0	0.996	1.0	1.0	0.98	0.953
Probe	SA	S	0.047	0.039	1.0	0.516	1.0	0.688	0.328	0.383	0.441	0.023
Probe	SA	S,TC	0.141	0.883	0.996	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC	0.027	0.41	0.918	1.0	1.0	1.0	1.0	1.0	1.0	0.984
Probe	SA	S,TC,SC,DG	0.02	0.82	0.883	1.0	1.0	1.0	1.0	1.0	1.0	0.996
Probe	SA	S,TC,SC,DG,PE	0.117	0.852	0.848	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	0.207	1.0	0.824	1.0	1.0	1.0	1.0	1.0	1.0	1.0

Table 10: Results for rejection set diversity on Summarization.

Method	Accept Sets	Reject Sets	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S	SA	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC,DG	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
CB	S	SA	0.996	0.984	0.004 / 0.0 / 0.167	1.0	0.691	1.0	1.0	1.0	0.973	0.984
CB	S	SA,TLD	1.0	1.0	0.004 / 0.0 / 0.168	0.992	0.395	0.715	0.906	0.758	0.914	0.91
CB	S	SA,TLD,PE	1.0	1.0	0.004 / 0.0 / 0.168	0.992	0.988	1.0	1.0	1.0	0.953	0.977
CB	S	SA,TLD,PE,TC	1.0	1.0	0.004 / 0.0 / 0.166	1.0	0.984	1.0	1.0	0.996	0.957	0.988
CB	S	SA,TLD,PE,TC,SC	1.0	1.0	0.004 / 0.0 / 0.168	1.0	1.0	1.0	1.0	1.0	1.0	0.992
CB	S	SA,TLD,PE,TC,SC,DG	1.0	1.0	0.004 / 0.0 / 0.166	1.0	1.0	1.0	1.0	1.0	1.0	0.984
SFT	S	SA	0.988	1.0	0.0 / 0.0 / 0.201	0.738	0.012	0.512	0.902	0.316	0.68	0.473
SFT	S	SA,TLD	1.0	1.0	0.004 / 0.0 / 0.198	0.98	0.008	0.887	0.988	0.832	0.828	0.641
SFT	S	SA,TLD,PE	0.914	1.0	0.004 / 0.0 / 0.208	0.691	0.012	0.445	0.988	0.75	0.727	0.527
SFT	S	SA,TLD,PE,TC	0.996	1.0	0.004 / 0.0 / 0.203	1.0	0.152	0.922	0.996	0.543	0.855	0.527
SFT	S	SA,TLD,PE,TC,SC	1.0	1.0	0.004 / 0.0 / 0.199	1.0	1.0	1.0	1.0	0.98	0.98	0.973
SFT	S	SA,TLD,PE,TC,SC,DG	0.988	1.0	0.004 / 0.0 / 0.205	1.0	1.0	0.996	0.98	0.871	0.914	0.898
SFT → CB	S	SA	1.0	1.0	0.02 / 0.0 / 0.202	0.57	0.051	0.629	0.398	0.039	0.496	0.082
SFT → CB	S	SA,TLD	1.0	1.0	0.016 / 0.0 / 0.2	0.672	0.066	0.762	0.637	0.961	0.703	0.961
SFT → CB	S	SA,TLD,PE	1.0	1.0	0.016 / 0.0 / 0.2	0.867	0.523	0.965	1.0	0.984	0.855	0.801
SFT → CB	S	SA,TLD,PE,TC	1.0	1.0	0.016 / 0.0 / 0.201	1.0	0.879	1.0	1.0	0.984	0.867	0.641
SFT → CB	S	SA,TLD,PE,TC,SC	1.0	0.961	0.016 / 0.0 / 0.201	1.0	0.992	0.961	0.957	0.148	0.844	0.387
SFT → CB	S	SA,TLD,PE,TC,SC,DG	1.0	1.0	0.016 / 0.0 / 0.203	1.0	1.0	0.988	1.0	1.0	1.0	1.0
DPO	S	SA	1.0	1.0	0.0 / 0.0 / 0.174	0.914	0.301	0.922	1.0	1.0	0.859	0.891
DPO	S	SA,TLD	1.0	1.0	0.004 / 0.0 / 0.174	0.988	0.445	0.977	1.0	1.0	0.895	0.895
DPO	S	SA,TLD,PE	1.0	1.0	0.004 / 0.0 / 0.176	1.0	0.918	0.988	1.0	1.0	0.957	0.926
DPO	S	SA,TLD,PE,TC	1.0	1.0	0.004 / 0.0 / 0.171	1.0	0.996	1.0	1.0	1.0	0.969	0.938
DPO	S	SA,TLD,PE,TC,SC	1.0	1.0	0.004 / 0.0 / 0.174	1.0	1.0	1.0	1.0	1.0	0.98	0.945
DPO	S	SA,TLD,PE,TC,SC,DG	1.0	1.0	0.004 / 0.0 / 0.175	1.0	1.0	1.0	1.0	1.0	0.98	0.945
Probe	S	SA	1.0	1.0	0.297	0.984	0.75	1.0	1.0	1.0	0.898	1.0
Probe	S	SA,TLD	0.984	1.0	0.219	0.992	1.0	1.0	1.0	1.0	0.934	1.0
Probe	S	SA,TLD,PE	0.934	1.0	0.215	0.996	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC	0.906	1.0	0.203	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC	0.875	1.0	0.227	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	0.887	1.0	0.234	1.0	1.0	1.0	1.0	1.0	0.969	1.0

1397 1396 1395 1394 1393 1392 1391 1390 1389 1388 1386 1385 1384 1383 1382 1381 1380 1379 1378 1377 1376 1375 1374 1373 1372 1371 1370 1369 1368 1367 1366 1365

Table 11: Results for rejection set diversity on Program Execution.

Method	Accept Sets	Reject Sets	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	PE	SA	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC,DG	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
CB	PE	SA	1.0	1.0	0.863	0.98	0.996	0.973	0.457 / 0.043 / 0.244	0.973	0.914	0.91
CB	PE	SA,TLD	1.0	1.0	1.0	0.988	0.973	1.0	0.457 / 0.043 / 0.26	0.691	0.961	0.902
CB	PE	SA,TLD,S	1.0	1.0	1.0	0.902	1.0	0.922	0.461 / 0.039 / 0.253	0.746	0.949	0.941
CB	PE	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.465 / 0.039 / 0.236	0.953	0.977	0.961
CB	PE	SA,TLD,S,TC,SC	0.992	0.98	0.977	1.0	0.992	1.0	0.461 / 0.043 / 0.254	0.93	1.0	0.953
CB	PE	SA,TLD,S,TC,SC,DG	0.723	0.98	0.125	0.898	0.32	0.867	0.457 / 0.039 / 0.249	0.277	0.387	0.633
SFT	PE	SA	0.914	0.734	0.004	0.055	0.0	0.512	0.0 / 0.246 / 0.543	0.027	0.047	0.207
SFT	PE	SA,TLD	0.941	1.0	0.008	0.363	0.043	0.68	0.0 / 0.309 / 0.582	0.031	0.098	0.262
SFT	PE	SA,TLD,S	0.918	0.996	0.891	0.848	0.434	0.754	0.004 / 0.289 / 0.592	0.051	0.109	0.387
SFT	PE	SA,TLD,S,TC	0.926	0.996	0.934	1.0	0.836	0.852	0.0 / 0.27 / 0.577	0.023	0.191	0.395
SFT	PE	SA,TLD,S,TC,SC	0.941	0.996	0.969	1.0	1.0	0.938	0.0 / 0.273 / 0.58	0.16	0.32	0.605
SFT	PE	SA,TLD,S,TC,SC,DG	0.887	0.988	0.539	0.992	0.789	0.898	0.0 / 0.293 / 0.564	0.027	0.234	0.375
SFT → CB	PE	SA	1.0	1.0	0.941	0.809	0.48	0.949	0.02 / 0.246 / 0.543	0.934	0.73	0.895
SFT → CB	PE	SA,TLD	1.0	1.0	1.0	1.0	1.0	1.0	0.027 / 0.25 / 0.548	0.902	0.977	0.898
SFT → CB	PE	SA,TLD,S	1.0	1.0	1.0	1.0	1.0	1.0	0.02 / 0.246 / 0.551	0.93	1.0	0.918
SFT → CB	PE	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.023 / 0.25 / 0.549	0.344	0.938	0.43
SFT → CB	PE	SA,TLD,S,TC,SC	1.0	1.0	1.0	1.0	1.0	0.996	0.023 / 0.25 / 0.55	0.953	0.988	0.922
SFT → CB	PE	SA,TLD,S,TC,SC,DG	1.0	1.0	1.0	1.0	1.0	1.0	0.023 / 0.246 / 0.543	1.0	1.0	0.992
DPO	PE	SA	1.0	1.0	0.996	0.914	0.336	0.859	0.0 / 0.148 / 0.416	0.094	0.324	0.625
DPO	PE	SA,TLD	1.0	1.0	0.738	0.922	0.148	0.891	0.0 / 0.148 / 0.413	0.062	0.449	0.594
DPO	PE	SA,TLD,S	1.0	1.0	1.0	1.0	1.0	0.996	0.0 / 0.133 / 0.392	0.301	0.695	0.848
DPO	PE	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.008 / 0.145 / 0.406	0.25	0.75	0.844
DPO	PE	SA,TLD,S,TC,SC	1.0	1.0	1.0	1.0	1.0	1.0	0.016 / 0.141 / 0.403	0.301	0.785	0.883
DPO	PE	SA,TLD,S,TC,SC,DG	1.0	1.0	1.0	1.0	1.0	1.0	0.02 / 0.145 / 0.416	0.453	0.82	0.902
Probe	PE	SA	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.953	0.758	0.949
Probe	PE	SA,TLD	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.941	0.895	0.957
Probe	PE	SA,TLD,S	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.883	0.816	0.926
Probe	PE	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.008	0.996	0.98	0.98
Probe	PE	SA,TLD,S,TC,SC	1.0	1.0	1.0	1.0	1.0	1.0	0.012	0.996	0.98	0.965
Probe	PE	SA,TLD,S,TC,SC,DG	1.0	1.0	1.0	1.0	1.0	1.0	0.07	1.0	0.992	0.988

Table 12: Results for multiple accept sets set diversity on Classification and Generation.

Method	Accept Sets	Reject Sets	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.082 / 0.473 / 0.57	0.125 / 0.301 / 0.327	0.004 / 0.0 / 0.165	0.105 / 0.0 / 0.073	0.012 / 0.0 / 0.191	0.305 / 0.0 / 0.171	0.672	0.723	0.344	0.496
CB	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.078 / 0.48 / 0.578	0.109 / 0.301 / 0.329	0.004 / 0.0 / 0.166	0.105 / 0.0 / 0.074	0.012 / 0.0 / 0.191	0.32 / 0.0 / 0.175	1.0	1.0	0.621	0.949
SFT	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.0 / 0.867 / 0.871	0.0 / 0.543 / 0.585	0.0 / 0.0 / 0.182	0.0 / 0.328 / 0.389	0.0 / 0.023 / 0.417	0.039 / 0.0 / 0.302	0.93	0.973	0.168	0.398
SFT → CB	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.0 / 0.863 / 0.867	0.004 / 0.543 / 0.585	0.02 / 0.0 / 0.184	0.273 / 0.328 / 0.389	0.004 / 0.027 / 0.419	0.113 / 0.0 / 0.304	1.0	1.0	0.391	0.695
DPO	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.0 / 0.238 / 0.414	0.004 / 0.145 / 0.158	0.0 / 0.0 / 0.168	0.0 / 0.0 / 0.075	0.0 / 0.0 / 0.203	0.051 / 0.0 / 0.223	1.0	1.0	0.574	0.496
Probe	SA,TLD,S,TC,SC,DG	PE,GSM8K	0.242	0.004	0.035	0.527	0.0	0.312	1.0	1.0	0.82	1.0

Table 13: Results for multiple accept sets set diversity on Math and Programming.

Method	Accept Sets	Reject Sets	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	PE,GSM8K	SA,TLD,S,TC	0.23	0.836	0.004	0.234	0.016	0.383	0.312 / 0.047 / 0.245	0.008 / 0.035 / 0.046	0.195	0.508
CB	PE,GSM8K	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.312 / 0.043 / 0.246	0.008 / 0.027 / 0.038	0.996	0.848
SFT	PE,GSM8K	SA,TLD,S,TC	0.988	1.0	0.988	0.992	0.645	0.754	0.0 / 0.004 / 0.291	0.0 / 0.207 / 0.229	0.184	0.5
SFT → CB	PE,GSM8K	SA,TLD,S,TC	1.0	1.0	1.0	1.0	0.68	0.906	0.34 / 0.004 / 0.301	0.184 / 0.219 / 0.241	0.598	0.367
DPO	PE,GSM8K	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	0.988	0.012 / 0.121 / 0.371	0.0 / 0.051 / 0.082	0.418	0.715
Probe	PE,GSM8K	SA,TLD,S,TC	1.0	1.0	1.0	1.0	1.0	1.0	0.074	0.0	0.773	0.93

Table 14: Results for precise scoping on Sentiment Analysis.

Table 15: Results for precise scoping on Summarization.

Method	Accept Sets	Reject Sets	S-FR	S-FA	SA	TLD	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S-FA	S-FR	0.004	0.0 / 0.0 / 0.181	0.148	0.66	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S-FA	S-FR,SA,TLD,TC,DG,PE	0.004	0.0 / 0.0 / 0.181	0.148	0.66	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S-FA	SA,TLD,TC,SC,PE	0.004	0.0 / 0.0 / 0.181	0.148	0.66	0.207	0.012	0.375	0.566	0.211	0.262	0.465
CB	S-FA	S-FR	1.0	0.0 / 0.0 / 0.18	0.918	0.473	0.238	0.008	0.363	0.559	0.211	0.293	0.477
CB	S-FA	S-FR,SA,TLD,TC,DG,PE	1.0	0.0 / 0.0 / 0.18	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.996
CB	S-FA	SA,TLD,TC,SC,PE	0.543	0.0 / 0.0 / 0.178	1.0	1.0	0.984	1.0	0.98	1.0	0.992	0.977	0.992
SFT	S-FA	S-FR	0.891	0.0 / 0.0 / 0.229	0.918	0.973	0.457	0.52	0.75	0.711	0.301	0.719	0.363
SFT	S-FA	S-FR,SA,TLD,TC,DG,PE	0.848	0.0 / 0.0 / 0.221	1.0	1.0	0.574	0.598	0.996	0.961	0.797	0.949	0.945
SFT	S-FA	SA,TLD,TC,SC,PE	0.176	0.0 / 0.0 / 0.18	1.0	1.0	0.711	1.0	1.0	0.996	0.965	0.953	0.949
SFT → CB	S-FA	S-FR	0.887	0.0 / 0.0 / 0.227	0.168	0.355	0.039	0.004	0.035	0.172	0.035	0.023	0.008
SFT → CB	S-FA	S-FR,SA,TLD,TC,DG,PE	0.988	0.0 / 0.0 / 0.231	0.996	1.0	0.984	0.996	1.0	0.973	0.418	0.898	0.105
SFT → CB	S-FA	SA,TLD,TC,SC,PE	0.586	0.0 / 0.0 / 0.23	1.0	0.996	0.996	1.0	0.992	0.996	0.309	0.891	0.129
DPO	S-FA	S-FR	1.0	0.0 / 0.0 / 0.214	0.926	1.0	0.984	0.977	0.922	0.996	0.98	0.836	0.926
DPO	S-FA	S-FR,SA,TLD,TC,DG,PE	1.0	0.0 / 0.0 / 0.197	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.973
DPO	S-FA	SA,TLD,TC,SC,PE	0.215	0.0 / 0.0 / 0.203	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.992	0.961
Probe	S-FA	S-FR	1.0	0.059	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.957	1.0
Probe	S-FA	S-FR,SA,TLD,TC,DG,PE	1.0	0.074	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.965	1.0
Probe	S-FA	SA,TLD,TC,SC,PE	1.0	0.074	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.965	1.0

1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529

1530	
1531	
1532	
1533	
1534	
1535	
1536	
1537	
1538	
1539	
1540	
1541	
1542	
1543	
1544	
1545	
1546	
1547	
1548	
1549	
1550	
1551	
1552	
1553	
1554	
1555	
1556	
1557	
1558	
1559	
1560	
1561	
1562	

Table 16: Results for data quantity evaluation on Sentiment Analysis.

Method	Accept Sets	Reject Sets	Num. Prompts	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	SA	S,TC,SC,DG,PE,TLD	128	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE,TLD	256	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE,TLD	512	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE,TLD	1024	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S,TC,SC,DG,PE,TLD	2048	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
CB	SA	S,TC,SC,DG,PE,TLD	128	0.113 / 0.531 / 0.586	0.918	0.988	1.0	1.0	1.0	0.984	0.996	1.0	0.992
CB	SA	S,TC,SC,DG,PE,TLD	256	0.109 / 0.52 / 0.582	1.0	0.953	1.0	1.0	0.996	1.0	0.996	1.0	0.988
CB	SA	S,TC,SC,DG,PE,TLD	512	0.102 / 0.531 / 0.59	0.25	0.004	0.48	0.016	0.438	0.691	0.645	0.379	0.535
CB	SA	S,TC,SC,DG,PE,TLD	1024	0.102 / 0.527 / 0.59	1.0	1.0	1.0	1.0	1.0	1.0	0.996	0.992	0.992
CB	SA	S,TC,SC,DG,PE,TLD	2048	0.102 / 0.527 / 0.586	0.254	0.004	0.863	0.117	0.844	0.777	0.766	0.52	0.699
SFT	SA	S,TC,SC,DG,PE,TLD	128	0.0 / 0.891 / 0.891	0.0	0.0	0.055	0.0	0.105	0.074	0.023	0.074	0.172
SFT	SA	S,TC,SC,DG,PE,TLD	256	0.0 / 0.867 / 0.867	0.0	0.0	0.016	0.0	0.102	0.031	0.0	0.012	0.055
SFT	SA	S,TC,SC,DG,PE,TLD	512	0.0 / 0.883 / 0.883	0.0	0.0	0.0	0.0	0.176	0.059	0.016	0.027	0.117
SFT	SA	S,TC,SC,DG,PE,TLD	1024	0.0 / 0.699 / 0.699	0.0	0.0	0.004	0.0	0.066	0.0	0.0	0.0	0.008
SFT	SA	S,TC,SC,DG,PE,TLD	2048	0.0 / 0.859 / 0.859	0.0	0.0	0.008	0.0	0.094	0.027	0.004	0.027	0.012
SFT → CB	SA	S,TC,SC,DG,PE,TLD	128	0.012 / 0.875 / 0.875	1.0	1.0	1.0	1.0	1.0	1.0	0.902	1.0	0.969
SFT → CB	SA	S,TC,SC,DG,PE,TLD	256	0.016 / 0.859 / 0.859	0.883	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.984
SFT → CB	SA	S,TC,SC,DG,PE,TLD	512	0.004 / 0.871 / 0.871	0.883	1.0	1.0	1.0	1.0	1.0	0.719	1.0	0.93
SFT → CB	SA	S,TC,SC,DG,PE,TLD	1024	0.004 / 0.879 / 0.879	0.996	1.0	1.0	1.0	1.0	1.0	0.875	1.0	0.891
SFT → CB	SA	S,TC,SC,DG,PE,TLD	2048	0.004 / 0.867 / 0.867	0.0	0.0	0.008	0.0	0.125	0.027	0.289	0.012	0.062
DPO	SA	S,TC,SC,DG,PE,TLD	128	0.055 / 0.402 / 0.41	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.926	0.949
DPO	SA	S,TC,SC,DG,PE,TLD	256	0.043 / 0.496 / 0.508	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.965	0.949
DPO	SA	S,TC,SC,DG,PE,TLD	512	0.027 / 0.598 / 0.598	1.0	1.0	1.0	1.0	0.996	1.0	1.0	0.957	0.949
DPO	SA	S,TC,SC,DG,PE,TLD	1024	0.012 / 0.75 / 0.75	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.98	0.953
DPO	SA	S,TC,SC,DG,PE,TLD	2048	0.02 / 0.609 / 0.613	1.0	1.0	1.0	1.0	0.996	1.0	1.0	0.98	0.953
Probe	SA	S,TC,SC,DG,PE,TLD	128	0.215	1.0	0.855	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	256	0.199	1.0	0.754	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	512	0.215	1.0	0.852	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	1024	0.219	1.0	0.891	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Probe	SA	S,TC,SC,DG,PE,TLD	2048	0.207	1.0	0.824	1.0	1.0	1.0	1.0	1.0	1.0	1.0

Table 17: Results for data quantity evaluation on Summarization.

Method	Accept Sets	Reject Sets	Num. Prompts	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S	SA,TLD,PE,TC,SC,DG	128	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC,DG	256	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC,DG	512	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC,DG	1024	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA,TLD,PE,TC,SC,DG	2048	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
CB	S	SA,TLD,PE,TC,SC,DG	128	0.984	1.0	0.004 / 0.0 / 0.167	1.0	1.0	1.0	1.0	0.992	1.0	1.0
CB	S	SA,TLD,PE,TC,SC,DG	256	1.0	1.0	0.004 / 0.0 / 0.166	1.0	1.0	1.0	1.0	0.992	1.0	1.0
CB	S	SA,TLD,PE,TC,SC,DG	512	0.992	0.996	0.004 / 0.0 / 0.167	0.996	0.996	1.0	1.0	0.996	1.0	1.0
CB	S	SA,TLD,PE,TC,SC,DG	1024	1.0	1.0	0.004 / 0.0 / 0.167	1.0	1.0	1.0	1.0	0.984	1.0	0.988
CB	S	SA,TLD,PE,TC,SC,DG	2048	1.0	1.0	0.004 / 0.0 / 0.166	1.0	1.0	1.0	1.0	1.0	1.0	0.984
SFT	S	SA,TLD,PE,TC,SC,DG	128	0.004	0.008	0.012 / 0.0 / 0.118	0.004	0.0	0.133	0.031	0.012	0.004	0.148
SFT	S	SA,TLD,PE,TC,SC,DG	256	0.156	0.102	0.004 / 0.0 / 0.117	0.059	0.016	0.195	0.035	0.016	0.109	0.137
SFT	S	SA,TLD,PE,TC,SC,DG	512	0.125	0.297	0.02 / 0.0 / 0.104	0.23	0.379	0.23	0.184	0.152	0.219	0.16
SFT	S	SA,TLD,PE,TC,SC,DG	1024	0.098	0.148	0.004 / 0.0 / 0.103	0.047	0.0	0.234	0.035	0.004	0.074	0.02
SFT	S	SA,TLD,PE,TC,SC,DG	2048	0.0	0.0	0.0 / 0.0 / 0.103	0.0	0.0	0.066	0.0	0.0	0.0	0.004
SFT → CB	S	SA,TLD,PE,TC,SC,DG	128	0.25	0.324	0.055 / 0.0 / 0.119	0.145	0.238	0.289	0.07	0.121	0.172	0.18
SFT → CB	S	SA,TLD,PE,TC,SC,DG	256	1.0	1.0	0.062 / 0.0 / 0.12	1.0	1.0	0.996	1.0	0.996	1.0	0.961
SFT → CB	S	SA,TLD,PE,TC,SC,DG	512	1.0	1.0	0.062 / 0.0 / 0.117	0.996	1.0	1.0	0.996	1.0	0.996	0.992
SFT → CB	S	SA,TLD,PE,TC,SC,DG	1024	1.0	1.0	0.078 / 0.0 / 0.124	1.0	1.0	1.0	1.0	0.969	1.0	0.41
SFT → CB	S	SA,TLD,PE,TC,SC,DG	2048	0.082	0.02	0.062 / 0.0 / 0.12	0.02	0.148	0.086	0.039	0.09	0.035	0.121
DPO	S	SA,TLD,PE,TC,SC,DG	128	1.0	1.0	0.004 / 0.0 / 0.171	1.0	1.0	1.0	1.0	1.0	0.977	0.949
DPO	S	SA,TLD,PE,TC,SC,DG	256	1.0	1.0	0.004 / 0.0 / 0.174	1.0	1.0	1.0	1.0	1.0	0.988	0.949
DPO	S	SA,TLD,PE,TC,SC,DG	512	1.0	1.0	0.004 / 0.0 / 0.172	1.0	1.0	1.0	1.0	1.0	0.984	0.941
DPO	S	SA,TLD,PE,TC,SC,DG	1024	1.0	1.0	0.004 / 0.0 / 0.175	1.0	1.0	1.0	1.0	1.0	0.992	0.945
DPO	S	SA,TLD,PE,TC,SC,DG	2048	1.0	1.0	0.004 / 0.0 / 0.175	1.0	1.0	1.0	1.0	1.0	0.98	0.945
Probe	S	SA,TLD,PE,TC,SC,DG	128	0.859	1.0	0.246	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	256	0.859	1.0	0.242	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	512	0.82	1.0	0.223	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	1024	0.836	1.0	0.199	1.0	1.0	1.0	1.0	1.0	0.961	1.0
Probe	S	SA,TLD,PE,TC,SC,DG	2048	0.887	1.0	0.234	1.0	1.0	1.0	1.0	1.0	0.969	1.0

Table 18: Results for data quantity evaluation on Program Execution.

Method	Accept Sets	Reject Sets	Num. Prompts	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	PE	SA,TLD,S,TC,SC	128	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC	256	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC	512	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC	1024	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA,TLD,S,TC,SC	2048	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
CB	PE	SA,TLD,S,TC,SC	128	1.0	1.0	1.0	1.0	1.0	1.0	0.457 / 0.043 / 0.25	0.766	0.918	0.781
CB	PE	SA,TLD,S,TC,SC	256	0.324	0.668	0.172	0.145	0.0	0.469	0.461 / 0.035 / 0.248	0.16	0.23	0.418
CB	PE	SA,TLD,S,TC,SC	512	0.547	0.961	0.035	0.84	0.094	0.766	0.461 / 0.043 / 0.25	0.207	0.336	0.582
CB	PE	SA,TLD,S,TC,SC	1024	1.0	1.0	1.0	1.0	1.0	1.0	0.457 / 0.035 / 0.247	0.711	1.0	0.941
CB	PE	SA,TLD,S,TC,SC	2048	0.996	0.973	0.977	1.0	0.953	0.992	0.465 / 0.043 / 0.251	0.906	0.996	0.941
SFT	PE	SA,TLD,S,TC,SC	128	0.0	0.0	0.0	0.0	0.0	0.043	0.0 / 0.0 / 0.285	0.0	0.0	0.012
SFT	PE	SA,TLD,S,TC,SC	256	0.0	0.0	0.004	0.0	0.0	0.074	0.0 / 0.0 / 0.461	0.0	0.004	0.012
SFT	PE	SA,TLD,S,TC,SC	512	0.0	0.0	0.004	0.004	0.0	0.051	0.0 / 0.0 / 0.347	0.0	0.004	0.012
SFT	PE	SA,TLD,S,TC,SC	1024	0.0	0.0	0.0	0.0	0.0	0.051	0.0 / 0.0 / 0.396	0.004	0.0	0.012
SFT	PE	SA,TLD,S,TC,SC	2048	0.055	0.0	0.0	0.008	0.0	0.184	0.043 / 0.0 / 0.43	0.008	0.012	0.012
SFT → CB	PE	SA,TLD,S,TC,SC	128	1.0	1.0	1.0	1.0	1.0	1.0	0.004 / 0.0 / 0.464	0.164	0.992	0.551
SFT → CB	PE	SA,TLD,S,TC,SC	256	1.0	1.0	1.0	1.0	1.0	1.0	0.004 / 0.0 / 0.454	0.23	0.992	0.469
SFT → CB	PE	SA,TLD,S,TC,SC	512	0.0	0.0	0.0	0.0	0.0	0.07	0.008 / 0.0 / 0.455	0.004	0.008	0.012
SFT → CB	PE	SA,TLD,S,TC,SC	1024	1.0	1.0	1.0	1.0	1.0	1.0	0.008 / 0.0 / 0.457	0.141	0.965	0.465
SFT → CB	PE	SA,TLD,S,TC,SC	2048	0.004	0.0	0.0	0.004	0.0	0.062	0.008 / 0.0 / 0.465	0.004	0.008	0.012
DPO	PE	SA,TLD,S,TC,SC	128	0.0	0.0	0.0	0.0	0.0	0.0	0.0 / 0.0 / 0.0	0.0	0.0	0.0
DPO	PE	SA,TLD,S,TC,SC	256	1.0	1.0	1.0	1.0	1.0	1.0	0.016 / 0.148 / 0.4	0.281	0.762	0.859
DPO	PE	SA,TLD,S,TC,SC	512	1.0	1.0	1.0	1.0	1.0	1.0	0.012 / 0.141 / 0.409	0.191	0.734	0.809
DPO	PE	SA,TLD,S,TC,SC	1024	1.0	1.0	1.0	1.0	1.0	1.0	0.016 / 0.137 / 0.4	0.285	0.789	0.875
DPO	PE	SA,TLD,S,TC,SC	2048	1.0	1.0	1.0	1.0	1.0	1.0	0.016 / 0.141 / 0.403	0.301	0.785	0.883
Probe	PE	SA,TLD,S,TC,SC	128	1.0	1.0	1.0	1.0	1.0	1.0	0.031	0.996	0.988	0.977
Probe	PE	SA,TLD,S,TC,SC	256	1.0	1.0	1.0	1.0	1.0	1.0	0.023	0.996	0.98	0.969
Probe	PE	SA,TLD,S,TC,SC	512	1.0	1.0	1.0	1.0	1.0	1.0	0.035	0.996	0.98	0.977
Probe	PE	SA,TLD,S,TC,SC	1024	1.0	1.0	1.0	1.0	1.0	1.0	0.062	1.0	0.992	0.988
Probe	PE	SA,TLD,S,TC,SC	2048	1.0	1.0	1.0	1.0	1.0	1.0	0.012	0.996	0.98	0.965

1629
1630

Table 19: Results for LoRA rank evaluation on Sentiment Analysis.

1631
1632
1633
1634
1635
1636

Method	Accept Sets	Reject Sets	Rank	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	SA	S	2	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S	4	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S	8	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S	16	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S	32	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
Sys.	SA	S	64	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.641	0.375	0.539
CB	SA	S	2	0.102 / 0.527 / 0.586	0.25	0.004	0.285	0.02	0.391	0.68	0.641	0.375	0.539
CB	SA	S	4	0.102 / 0.527 / 0.586	0.246	0.004	0.281	0.02	0.391	0.68	0.637	0.379	0.539
CB	SA	S	8	0.102 / 0.527 / 0.586	0.242	0.004	0.285	0.02	0.391	0.68	0.641	0.375	0.543
CB	SA	S	16	0.105 / 0.527 / 0.586	0.207	0.996	0.785	1.0	0.93	0.789	0.77	0.656	0.801
CB	SA	S	32	0.102 / 0.527 / 0.594	0.223	0.996	0.934	0.855	0.738	0.973	0.902	0.914	0.883
CB	SA	S	64	0.098 / 0.539 / 0.602	0.18	1.0	0.645	0.637	0.688	0.66	0.723	0.48	0.766
SFT	SA	S	2	0.008 / 0.887 / 0.887	0.004	0.984	0.645	0.984	0.742	0.625	0.551	0.418	0.879
SFT	SA	S	4	0.0 / 0.902 / 0.902	0.012	0.977	0.566	0.957	0.723	0.504	0.406	0.32	0.859
SFT	SA	S	8	0.0 / 0.902 / 0.902	0.023	0.961	0.586	0.602	0.695	0.43	0.234	0.238	0.77
SFT	SA	S	16	0.0 / 0.867 / 0.867	0.035	0.98	0.559	0.73	0.703	0.477	0.207	0.145	0.824
SFT	SA	S	32	0.0 / 0.879 / 0.879	0.0	0.996	0.508	0.559	0.699	0.156	0.035	0.141	0.664
SFT	SA	S	64	0.0 / 0.684 / 0.684	0.0	0.719	0.496	0.359	0.488	0.027	0.0	0.051	0.328
SFT → CB	SA	S	2	0.0 / 0.867 / 0.867	0.035	0.98	0.559	0.73	0.703	0.48	0.207	0.145	0.824
SFT → CB	SA	S	4	0.0 / 0.867 / 0.867	0.035	0.98	0.559	0.75	0.703	0.484	0.207	0.145	0.824
SFT → CB	SA	S	8	0.012 / 0.859 / 0.859	0.035	1.0	0.59	1.0	0.988	0.613	0.25	0.375	0.836
SFT → CB	SA	S	16	0.004 / 0.863 / 0.863	0.031	1.0	0.973	1.0	0.984	0.934	0.957	0.867	0.969
SFT → CB	SA	S	32	0.0 / 0.871 / 0.871	0.035	1.0	0.902	0.0	0.457	0.75	0.613	0.488	0.863
SFT → CB	SA	S	64	0.0 / 0.871 / 0.871	0.027	0.969	0.629	0.188	0.5	0.543	0.309	0.246	0.832
DPO	SA	S	2	0.059 / 0.664 / 0.73	0.035	0.164	0.465	0.379	0.531	0.691	0.695	0.418	0.684
DPO	SA	S	4	0.027 / 0.727 / 0.797	0.02	0.527	0.52	0.555	0.602	0.711	0.777	0.461	0.742
DPO	SA	S	8	0.012 / 0.742 / 0.828	0.016	0.848	0.559	0.75	0.633	0.727	0.82	0.492	0.785
DPO	SA	S	16	0.0 / 0.797 / 0.871	0.012	0.988	0.664	0.863	0.684	0.84	0.867	0.531	0.801
DPO	SA	S	32	0.0 / 0.809 / 0.875	0.004	1.0	0.695	0.918	0.691	0.922	0.863	0.555	0.836
DPO	SA	S	64	0.0 / 0.84 / 0.879	0.004	1.0	0.957	0.988	0.875	0.992	0.984	0.727	0.91

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674

1675

1676

1677

1678

1679

1680

1681

1682

1683

1684

Table 20: Results for LoRA rank evaluation on Summarization.

1685

1686

Method	Accept Sets	Reject Sets	Rank	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	S	SA	2	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA	4	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA	8	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA	16	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA	32	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
Sys.	S	SA	64	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
CB	S	SA	2	0.148	0.66	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.262	0.465
CB	S	SA	4	0.156	0.668	0.004 / 0.0 / 0.165	0.207	0.012	0.375	0.566	0.211	0.266	0.469
CB	S	SA	8	1.0	1.0	0.004 / 0.0 / 0.165	1.0	0.27	0.965	1.0	0.984	0.969	0.969
CB	S	SA	16	1.0	0.961	0.004 / 0.0 / 0.167	1.0	0.562	1.0	0.98	1.0	0.973	0.984
CB	S	SA	32	1.0	0.973	0.004 / 0.0 / 0.166	1.0	0.043	0.988	0.988	1.0	0.906	0.984
CB	S	SA	64	0.996	0.992	0.004 / 0.0 / 0.165	0.348	0.031	0.41	0.566	0.945	0.652	0.633
SFT	S	SA	2	0.926	1.0	0.012 / 0.0 / 0.229	0.602	0.043	0.391	0.82	0.559	0.672	0.773
SFT	S	SA	4	0.957	1.0	0.012 / 0.0 / 0.225	0.684	0.008	0.414	0.863	0.426	0.621	0.75
SFT	S	SA	8	1.0	1.0	0.004 / 0.0 / 0.213	0.379	0.0	0.359	0.547	0.219	0.59	0.449
SFT	S	SA	16	0.996	1.0	0.0 / 0.0 / 0.226	0.367	0.004	0.359	0.832	0.543	0.508	0.684
SFT	S	SA	32	1.0	0.996	0.0 / 0.0 / 0.208	0.535	0.004	0.504	0.891	0.98	0.668	0.617
SFT	S	SA	64	0.996	1.0	0.0 / 0.0 / 0.204	0.48	0.059	0.43	0.367	0.086	0.707	0.133
SFT → CB	S	SA	2	0.996	1.0	0.004 / 0.0 / 0.226	0.379	0.012	0.387	0.836	0.543	0.508	0.688
SFT → CB	S	SA	4	0.027	0.344	0.004 / 0.0 / 0.226	0.031	0.008	0.195	0.082	0.133	0.008	0.246
SFT → CB	S	SA	8	1.0	1.0	0.008 / 0.0 / 0.224	0.992	0.324	0.98	0.84	0.082	0.898	0.297
SFT → CB	S	SA	16	1.0	1.0	0.012 / 0.0 / 0.224	0.398	0.008	0.379	0.184	0.141	0.672	0.195
SFT → CB	S	SA	32	1.0	0.992	0.012 / 0.0 / 0.225	0.02	0.008	0.168	0.133	0.117	0.406	0.129
SFT → CB	S	SA	64	1.0	0.996	0.008 / 0.0 / 0.221	0.016	0.008	0.047	0.078	0.004	0.188	0.059
DPO	S	SA	2	0.84	0.969	0.004 / 0.0 / 0.168	0.531	0.074	0.535	0.949	0.758	0.523	0.82
DPO	S	SA	4	0.902	0.988	0.004 / 0.0 / 0.17	0.719	0.145	0.59	0.996	0.902	0.582	0.852
DPO	S	SA	8	0.988	0.996	0.0 / 0.0 / 0.177	0.863	0.195	0.711	1.0	0.961	0.641	0.863
DPO	S	SA	16	1.0	1.0	0.0 / 0.0 / 0.18	0.891	0.328	0.828	1.0	0.977	0.742	0.895
DPO	S	SA	32	1.0	1.0	0.0 / 0.0 / 0.179	0.855	0.305	0.902	1.0	1.0	0.832	0.887
DPO	S	SA	64	1.0	1.0	0.0 / 0.0 / 0.181	0.734	0.297	0.867	0.996	1.0	0.84	0.84

1709

1710

1711

1712

Table 21: Results for LoRA rank evaluation on Program Execution.

1713

1714

Method	Accept Sets	Reject Sets	Rank	SA	TLD	S	TC	SC	DG	PE	GSM8k	QA	Alpaca
Sys.	PE	SA	2	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA	4	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA	8	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA	16	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA	32	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
Sys.	PE	SA	64	0.199	0.758	0.004	0.277	0.012	0.395	0.453 / 0.039 / 0.252	0.129	0.242	0.441
CB	PE	SA	2	0.195	0.758	0.004	0.277	0.016	0.395	0.457 / 0.039 / 0.252	0.129	0.242	0.441
CB	PE	SA	4	0.211	0.758	0.004	0.277	0.016	0.395	0.457 / 0.039 / 0.252	0.129	0.238	0.445
CB	PE	SA	8	0.887	0.863	0.602	0.699	0.758	0.508	0.457 / 0.039 / 0.249	0.652	0.668	0.598
CB	PE	SA	16	1.0	0.996	0.859	0.98	0.996	0.965	0.457 / 0.043 / 0.244	0.973	0.914	0.906
CB	PE	SA	32	1.0	1.0	0.891	0.887	0.977	0.98	0.445 / 0.043 / 0.254	0.785	0.688	0.691
CB	PE	SA	64	1.0	0.824	0.422	0.551	0.023	0.797	0.473 / 0.043 / 0.26	0.203	0.258	0.438
SFT	PE	SA	2	0.938	0.723	0.273	0.438	0.113	0.664	0.004 / 0.195 / 0.503	0.027	0.094	0.332
SFT	PE	SA	4	0.969	0.848	0.281	0.59	0.113	0.77	0.016 / 0.215 / 0.514	0.047	0.125	0.488
SFT	PE	SA	8	0.906	0.707	0.0	0.215	0.0	0.594	0.0 / 0.234 / 0.517	0.016	0.039	0.273
SFT	PE	SA	16	0.906	0.816	0.031	0.34	0.098	0.637	0.0 / 0.254 / 0.538	0.062	0.098	0.363
SFT	PE	SA	32	0.91	0.859	0.023	0.652	0.293	0.703	0.0 / 0.258 / 0.546	0.016	0.203	0.23
SFT	PE	SA	64	0.926	0.93	0.137	0.594	0.02	0.668	0.02 / 0.125 / 0.463	0.188	0.277	0.254
SFT → CB	PE	SA	2	0.906	0.812	0.031	0.34	0.098	0.637	0.008 / 0.254 / 0.538	0.066	0.102	0.367
SFT → CB	PE	SA	4	0.613	0.645	0.012	0.309	0.086	0.613	0.008 / 0.254 / 0.538	0.07	0.098	0.387
SFT → CB	PE	SA	8	0.758	0.918	0.926	0.41	0.586	0.809	0.012 / 0.254 / 0.54	0.27	0.488	0.375
SFT → CB	PE	SA	16	1.0	0.973	1.0	0.375	0.68	0.805	0.008 / 0.25 / 0.531	0.27	0.531	0.391
SFT → CB	PE	SA	32	1.0	0.949	0.984	0.059	0.137	0.758	0.008 / 0.254 / 0.532	0.277	0.188	0.324
SFT → CB	PE	SA	64	1.0	0.926	0.539	0.07	0.0	0.348	0.008 / 0.254 / 0.535	0.336	0.105	0.23
DPO	PE	SA	2	0.672	0.953	0.004	0.465	0.043	0.539	0.105 / 0.051 / 0.352	0.051	0.273	0.496
DPO	PE	SA	4	0.898	0.98	0.082	0.656	0.066	0.672	0.031 / 0.047 / 0.376	0.031	0.289	0.543
DPO	PE	SA	8	0.973	0.988	0.195	0.781	0.082	0.715	0.031 / 0.055 / 0.41	0.055	0.309	0.605
DPO	PE	SA	16	0.996	0.996	0.477	0.746	0.051	0.738	0.012 / 0.055 / 0.403	0.051	0.27	0.59
DPO	PE	SA	32	1.0	0.871	0.867	0.215	0.828	0.008 / 0.051 / 0.384	0.062	0.34	0.629	
DPO	PE	SA	64	0.984	0.969	0.418	0.555	0.012	0.598	0.0 / 0.02 / 0.362	0.008	0.078	0.477

1736