
Course Project: Part 2

Generative AI

Saarland University – Winter Semester 2024/25

Martínez

7057573

cama00005@stud.uni-saarland.de

1 Part 2.a: Fine-tuning Baseline Models

2 (I.11) Table 1 shows the comparison between the Phi-3-SFT-Repair-r16-alpha32 and Phi-3-SFT-Hints-
3 r16-alpha32 models on the Repair task. As expected, the former is better on repair, than the latter, as
4 demonstrated by the better values in both *RPass* and *REdit*, since it was specifically fine-tuned for the
5 Repair task. The Phi-3-SFT-Hints-r16-alpha32 model, on the other hand, was fine-tuned for the hint
6 generation task, and thus performs worse on the Repair task, but better on the Hint generation task, as
7 (I.12) goes on to show.

Table 1: Program repair quality metrics, *RPass* and *REdit*, for the fine-tuned models on Repair and Hint tasks. In parentheses, the change in the metric compared to the Repair model is shown, where green means an improving change and red a worsening change. For (I.11).

Model	<i>RPass</i>	<i>REdit</i>
Phi-3-SFT-Hints_r16_alpha32	72.0 (−4.0)	21.22 (+11.33)
Phi-3-SFT-Repair_r16_alpha32	76.0 (−)	9.89 (−)

8 (I.12) Table 2 and 3 show the evaluation of the generated hints by the Phi-3-SFT-Repair-r16-alpha32
9 and Phi-3-SFT-Hints-r16-alpha32 models respectively using the *HGood* metric [1]. As expected,
10 since the Phi-3-SFT-Hints-r16-alpha32 model has been fine-tuned specifically on the Hint generation
11 task, it performs very well on all metrics, achieving a score of 0.84 in *HGood*. On the other hand, the
12 Phi-3-SFT-Repair-r16-alpha32 model, which has been fine-tuned on the Repair task, performs worse,
13 with an average *HGood* score of 0.2. This model is only slightly better than the base Phi-3-mini
14 model, which scored 0.16 in *HGood* in Project Part#1.

Table 2: Hint Quality Metrics for Phi-3-SFT-Repair-r16-alpha32. This model is fine-tuned on Repair task. For (I.12).

Problem	Program	<i>HCorrect</i>	<i>HInformative</i>	<i>HConceal</i>	<i>HComprehensible</i>	<i>HGood</i>
Problem 1	Prog. 1	1	1	0	1	0
	Prog. 2	1	1	0	1	0
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	0	1	0
	Prog. 5	1	0	0	1	0
Avg. Problem 1		1.0	0.8	0.2	1.0	0.2
Problem 2	Prog. 1	1	1	0	1	0
	Prog. 2	0	0	1	1	0
	Prog. 3	1	0	0	1	0
	Prog. 4	0	0	1	1	0
	Prog. 5	1	1	0	1	0
Avg. Problem 2		0.6	0.4	0.4	1.0	0.0
Problem 3	Prog. 1	1	1	0	1	0
	Prog. 2	1	1	0	1	0
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	0	1	0
Avg. Problem 3		1.0	1.0	0.4	1.0	0.4
Problem 4	Prog. 1	1	0	1	1	0
	Prog. 2	1	0	0	1	0
	Prog. 3	0	0	1	1	0
	Prog. 4	0	0	1	1	0
	Prog. 5	0	0	1	1	0
Avg. Problem 4		0.4	0.0	0.8	1.0	0.0
Problem 5	Prog. 1	1	1	0	1	0
	Prog. 2	1	1	1	1	1
	Prog. 3	0	0	0	1	0
	Prog. 4	1	1	0	1	0
	Prog. 5	1	1	1	1	1
Avg. Problem 5		0.8	0.8	0.4	1.0	0.4
Overall Average		0.76	0.6	0.44	1.0	0.2

Table 3: Hint Quality Metrics for Phi-3-SFT-Hints-r16-alpha32. This model is fine-tuned on Hint generation. For (I.12).

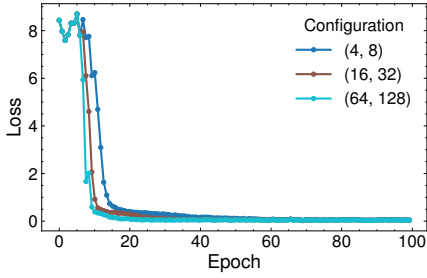
Problem	Program	<i>HCorrect</i>	<i>HInformative</i>	<i>HConceal</i>	<i>HComprehensible</i>	<i>HGood</i>
Problem 1	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	0	1	1	0
	Prog. 5	1	1	1	1	1
Avg. Problem 1		1.0	0.8	1.0	1.0	0.8
Problem 2	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	0	1	1	0
Avg. Problem 2		1.0	0.8	1.0	1.0	0.8
Problem 3	Prog. 1	1	1	1	1	1
	Prog. 2	1	0	1	1	0
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 3		1.0	0.8	1.0	1.0	0.8
Problem 4	Prog. 1	0	0	1	1	0
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 4		0.8	0.8	1.0	1.0	0.8
Problem 5	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 5		1.0	1.0	1.0	1.0	1.0
Overall Average		0.96	0.84	1.0	1.0	0.84

15 Part 2.b: Fine-tuning for Program Repair with Varying LoRA Parameters

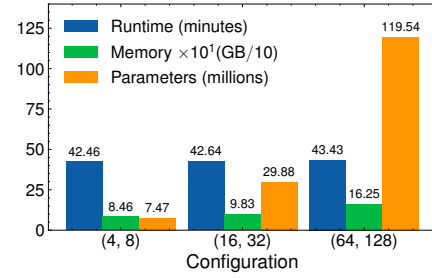
16 (I.13) Table 4 shows a comparison between different fine-tuned versions of the Phi-3-mini model
 17 on the Repair task using varying LoRA parameters. Recall that in Project Part#1 (see Table 7 in
 18 **Appendix**), the base Phi-3-mini scored 36.0 in *RPass* and 18.11 in *REdit*, which is significantly
 19 worse than any of the fine-tuned models. This means that the fine-tuning process with LoRA was
 20 effective and greatly improved the model’s performance on the Repair task.

Table 4: Program repair quality metrics, *RPass* and *REdit*, for the fine-tuned models on Repair task using different LoRA parameter configurations: $(r, \alpha) = (4, 8), (16, 32), (64, 128)$. In parentheses, the change in the metric compared to a base configuration, chosen to be $(16, 32)$, is shown, where **green** means an improving change and **red** a worsening change. For (I.13).

(r, α)	<i>RPass</i>	<i>REdit</i>
(4, 8)	60.0 (−16.0)	11.47 (+1.58)
(16, 32)	76.0 (−)	9.89 (−)
(64, 128)	84.0 (+8.0)	17.14 (+7.25)



(a) Loss over training epochs.



(b) *TrainingTime*, *TrainingMemory* and number of trained parameters comparison.

Figure 1: Trade-offs between LoRA parameter configurations in terms of loss evolution, *TrainingTime*, *TrainingMemory*, and number of trained parameters. These experiments were conducted with a GPU T4 instance on Google Colab. For (I.13).

(I.14) Figure 1 shows that the *TrainingTime* of all three LoRA configurations are very similar. The main difference thus lies in the *TrainingMemory* and the number of trained parameters. Naturally, the $(64, 128)$ configuration has the highest *TrainingMemory* and number of trained parameters, which is expected given the higher values of r and α . More precisely, the number of trained parameters increases approximately by a factor of 4 when going from $(4, 8)$ to $(16, 32)$, and by a factor of 4 again when going from $(16, 32)$ to $(64, 128)$. This is consistent with the theoretical analysis of LoRA, which predicts that the number of parameters scales quadratically with the values of r and α [2].

On the other hand, the increase of *TrainingMemory* between configurations is not as extreme. From $(4, 8)$ to $(16, 32)$, there is a 16% increase, and from $(16, 32)$ to $(64, 128)$, there is a 65% increase. This is likely due to the fact that the memory requirements of the model are not only determined by the number of parameters, but also by the size of the intermediate activations and the gradients during training.

Finally, in terms of program repair quality, the $(64, 128)$ configuration outperforms the other two configurations in terms of *RPass*, but has the highest *REdit*, suggesting that it might be overfitting to the training data. The $(4, 8)$ configuration has the lowest *RPass* and the second highest *REdit*, indicating that it is underfitting. The $(16, 32)$ configuration is a good balance between the two, achieving a value of *RPass* only 8 points lower than the best, but at the same time the lowest *REdit*.

The choice between the three configurations will depend on the specific requirements of the application. Personally, in the feedback generation domain, the $(16, 32)$ configuration would be the most suitable, as it can generate accurate repairs (high *RPass*) with the lowest number of changes to the buggy code compared to the other models (low *REdit*). This allows a student to learn from the generated feedback without being overwhelmed by too many changes or disappointed that the complete idea of the code was changed.

Part 2.c: Fine-tuning on Merged Datasets

(I.15) Table 5 shows the comparison between the Phi-3-SFT-Repair-r16-alpha32, Phi-3-SFT-Hints-r16-alpha32 and Phi-3-SFT-Merged-r16-alpha32 models on the Repair task. The latter, fine-tuned on both Repair and Hint tasks outperforms the other two in terms of *RPass* and achieves the second lowest *REdit*. This suggests that the multi-task learning approach is effective in improving the model’s performance on the Repair task, at the cost of being slightly less conservative in suggesting changes to the student’s buggy code.

Table 5: Program repair quality metrics, *RPass* and *REdit*, for the fine-tuned models on Repair, Hint and both tasks combined. In parentheses, the change in the metrics compared to the multi-task model (Phi-3-SFT-Merged_r16_alpha32) is shown, where green means an improving change and red a worsening change. The results of (I.11) are included for ease of comparison. For (I.15).

Model	<i>RPass</i>	<i>REdit</i>
Phi-3-SFT-Hints_r16_alpha32	72.0 (−8.0)	21.22 (+6.62)
Phi-3-SFT-Repair_r16_alpha32	76.0 (−4.0)	9.89 (−4.71)
Phi-3-SFT-Merged_r16_alpha32	80.0 (−)	14.6 (−)

(I.16) Table 6 shows the evaluation of the Phi-3-SFT-Merged-r16-alpha32 model on the Hint generation task using the *HGood* metric [1]. As expected, the model performs very well on all metrics, achieving a score of 0.92 in *HGood*. Given the results of (I.15), this means that this model became the best at both tasks, at the same time.

This multi-task learning approach of using a concatenation of both datasets (Repair + Hint) offered a more diverse scenario for a model to learn from. Both of these tasks intuitively benefit from one another. During training, the model is given access to the repairs of the buggy programs, which should give it exact guidelines on what to hint at for a student. Thus a diverse dataset such as this one contributes to a hand-in-hand improvement in the model’s performance on both tasks.

Finally, a comparison of all evaluated Phi-3-mini models is provided with Figure 2, in order to see how the SFT models fare against the base model and the CoT model from Project Part#1.

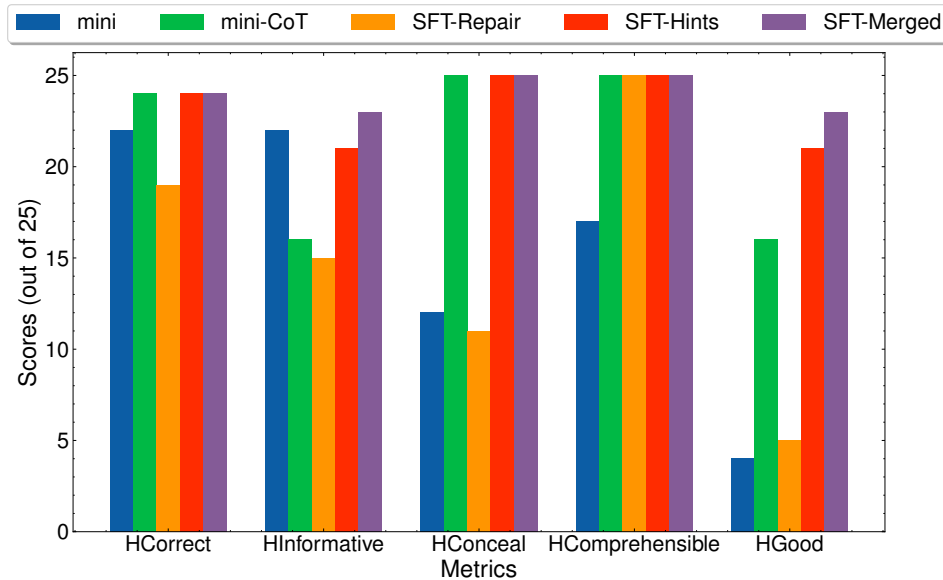


Figure 2: Comparison between all Phi-3-mini models: the base model, with CoT (both from Project Part#1), Phi-3-SFT-Repair-r16-alpha32, Phi-3-SFT-Hints-r16-alpha32 and Phi-3-SFT-Merged-r16-alpha32 models, evaluated on the Hint generation task in terms of the *HGood* metric. For (I.16).

Table 6: Hint Quality Metrics for Phi-3-SFT-Merged-r16-alpha32. This model is fine-tuned on both Repair and Hint tasks. For (I.16).

Problem	Program	<i>HCorrect</i>	<i>HInformative</i>	<i>HConceal</i>	<i>HComprehensible</i>	<i>HGood</i>
Problem 1	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 1		1.0	1.0	1.0	1.0	1.0
Problem 2	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 2		1.0	1.0	1.0	1.0	1.0
Problem 3	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 3		1.0	1.0	1.0	1.0	1.0
Problem 4	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	1	1	1	1
	Prog. 4	0	0	1	1	0
	Prog. 5	1	1	1	1	1
Avg. Problem 4		0.8	0.8	1.0	1.0	0.8
Problem 5	Prog. 1	1	1	1	1	1
	Prog. 2	1	1	1	1	1
	Prog. 3	1	0	1	1	0
	Prog. 4	1	1	1	1	1
	Prog. 5	1	1	1	1	1
Avg. Problem 5		1.0	0.8	1.0	1.0	0.8
Overall Average		0.96	0.92	1.0	1.0	0.92

(I.17) In order to create the merged dataset to train the multi-task model (Phi-3-SFT-Merged_r16_alpha32), I added a simple concatenation of both datasets (Repair + Hint) in `project_part2_assemble_dataset.py` and a shuffling of the result afterwards.

Figure 3 shows the trade-offs between the specialized models and the multi-task model in terms of loss evolution, *TrainingTime*, *TrainingMemory*, and number of trained parameters. The multi-task model with $(r, \alpha) = (16, 32)$ has a higher *TrainingTime*, but suprisingly has a better loss evolution compared to the specialized models. That is, the loss decreases rapidly in the first few epochs, indicating rapid learning, and then plateaus similar to the other models.

On the other hand, the *TrainingMemory* between the specialized model for the Repair task and the multi-task model is the same.

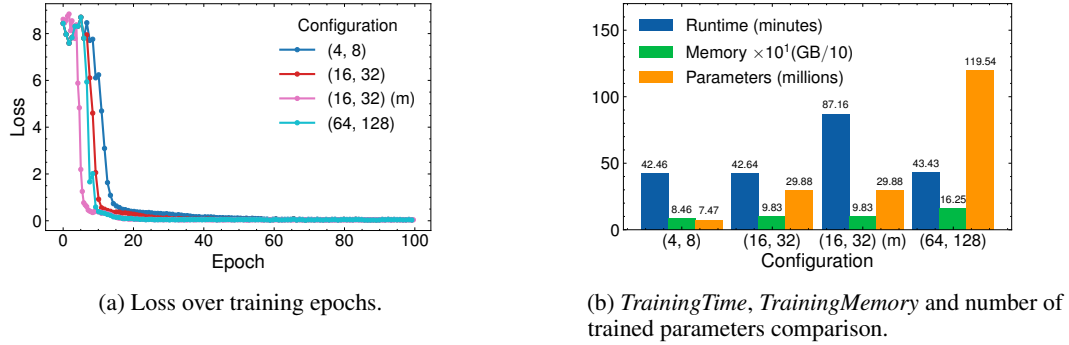


Figure 3: Trade-offs between the specialized models and the multi-task model in terms of loss evolution, *TrainingTime*, *TrainingMemory*, and number of trained parameters. The (m) indicates the multi-task model. These experiments were conducted with a GPU T4 instance on Google Colab. For (I.17).

As expected, the (4, 8) configuration has the lowest number of parameters and thus incurs in the least *TrainingMemory*. The *TrainingTime* between all configurations, (4, 8), (16, 32) and (64, 128) is the same, which is also supported by the loss evolution plot. This means that the model does not gain anything from adding $4\times$ more parameters, going from (16, 32) to (64, 128), as their loss curve almost overlap. And if we look at Table 5 again, the gain in *RPass* is only 4 points better than (16, 32), but the *REdit* is worse. On the other hand, the (16, 32) (m) configuration has approximately $2\times$ the *TrainingTime* than all the others, which is expected as the dataset grew in size.

Therefore, the (16, 32) configuration seems to be the best trade-off between the number of parameters, *TrainingMemory*, and the loss evolution. It has a good balance between the number of parameters and *TrainingMemory*, and the results from Table 5 shows that it is not overfitting as (64, 128).

Figure 4 shows the same concepts, as Figure 3, but for a GPU A100 instance. Using this GPU significantly reduced the *TrainingTime* for all models, especially for the multi-task model. Nevertheless, if we take into account relative performance, the multi-task model still requires $2\times$ the *TrainingTime* of the specialized model for the Repair task, i.e., 6.63 minutes vs. 13.50 minutes with an A100 and 42.64 minutes vs 87.16 minutes with a T4, respectively.

On the other hand, the loss over training epochs showcase the *scaling laws*, as we learned in Week 4. All models showcase a stepper and smoother decrease in loss as the number of epochs increases. That is, a higher computing power effectively reduced the time to reach a certain loss value, even though the hyperparameters, model architecture and datasets were the same [3].

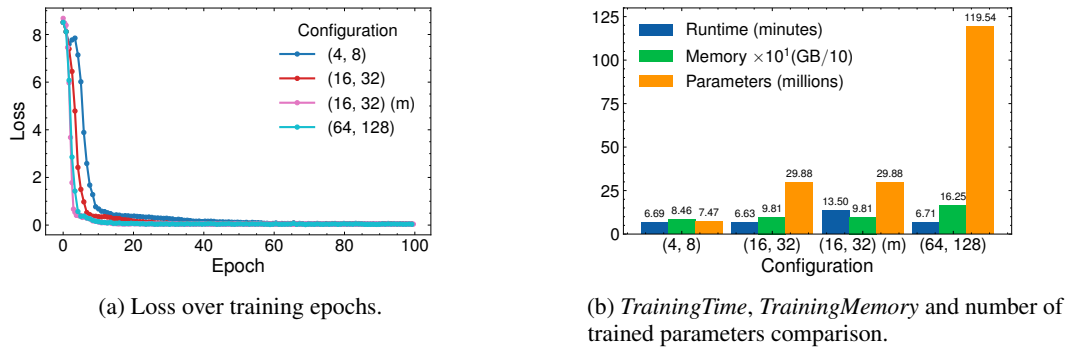


Figure 4: Trade-offs between the specialized models and the multi-task model in terms of loss evolution, *TrainingTime*, *TrainingMemory*, and number of trained parameters. The (m) indicates the multi-task model. These experiments were conducted with a GPU A100 instance on Google Colab. For (I.17).

91 **Acknowledgements**

92 Week 10's slides and listed references.

93 **References**

- 94 [1] Nachiket Kotalwar, Alkis Gotovos, and Adish Singla. Hints-in-browser: Benchmarking language
95 models for programming feedback generation, 2024.
- 96 [2] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
97 Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- 98 [3] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
99 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
100 models. *CoRR*, abs/2001.08361, 2020.

101 **Appendix: (I.3) Results of Project Part#1**

Table 7: Overall results for GPT-4o-mini and Phi-3-mini with the basic prompt to generate program repairs. For **(I.1)** and **(I.3)**.

Model	<i>RPass</i>	<i>REdit</i>
GPT-4o-mini	88.0	23.77
Phi-3-mini	36.0	18.11