

Amuro & Char: Analyzing the Relationship between Pre-Training and Fine-Tuning of Large Language Models

Anonymous ACL submission

Abstract

The development of large language models leads to the formation of a pre-train-then-align paradigm, in which the model is typically pre-trained on a large text corpus and undergoes a tuning stage to align the model with human preference or downstream tasks. In this work, we investigate the relationship between pre-training and fine-tuning by fine-tuning multiple intermediate pre-trained model checkpoints, we find that i) continual pre-training improves the model in a latent way that unveils after fine-tuning; ii) with extra fine-tuning, the datasets that the model does not demonstrate capability gain much more than those that the model performs well during the pre-training stage; iii) although model benefits significantly through supervised fine-tuning, it may forget previously known domain knowledge and the tasks that are not seen during fine-tuning; iv) the supervised fine-tuned model resembles high sensitivity to few-shot evaluation prompts, but this sensitivity can be alleviated by more pre-training.¹

1 Introduction

The rise of large language models (LLMs) as a general-purpose tool for a diverse range of natural language processing tasks has dramatically transformed the field, introducing new paradigms for data collection and model training (Brown et al., 2020, Biderman et al., 2023, Touvron et al., 2023, Jiang et al., 2023, Chowdhery et al., 2023, Groeneveld et al., 2024, Wang et al., 2024, *inter alia*). Numerous models, training methods, datasets, and evaluation methods continue to be developed on an ongoing basis. Nevertheless, a unified paradigm has emerged for training LLMs: pre-train on an enormous corpus of diverse documents, ranging from 250B (Biderman et al., 2023) to 15T (AI@Meta, 2024) tokens, followed by an

¹Code, results, and data to reproduce the experiments are available at <https://anonymous.4open.science/r/AmuroCharRelease-DEC5>

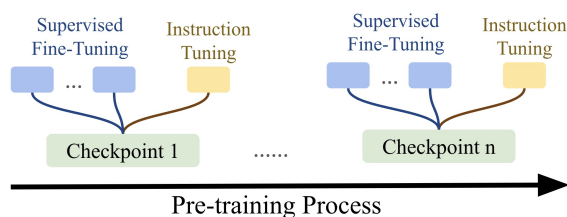


Figure 1: Illustration of the experimental scheme.

alignment stage to make the model more useful and performative for various tasks.

Based on this paradigm, work has focused on improving each of these two stages. For better pre-trained models, exploration is done on larger training sets (Hoffmann et al., 2022; AI@Meta, 2024; Touvron et al., 2023), different data selection mechanisms (Xia et al., 2024), higher quality data (Zhou et al., 2024), and various model architectures (Su et al., 2024; Touvron et al., 2023). Meanwhile, research on model alignment includes different training objectives (Rafailov et al., 2024; Schulman et al., 2017), new datasets (Narayanan and Aepli, 2024), more efficient training (Hu et al., 2021; Dettmers et al., 2024), safety tuning (Bianchi et al., 2023), among other approaches. The alignment stage usually involves either supervised fine-tuning for specific tasks or instruction fine-tuning for general-purpose usage. Regardless, fine-tuning (almost always) comes at the end of pre-training and yields remarkable improvements on downstream tasks (Touvron et al., 2023; Groeneveld et al., 2024). However, the benefits of each stage are largely explored independently, with improvements to pretraining being orthogonal to benefits from model alignment. Fine-tuning starts with the final pretraining model checkpoint.

Rather than explore these two training regimes independently, we question: **how do model pre-training and fine-tuning interact to affect the resulting abilities of the model?** Does more pre-

070 training hinder better fine-tuning results? What
071 does the model learn during pre-training, and
072 what does it forget during fine-tuning? To answer
073 these questions, we fine-tune **multiple pre-**
074 **training checkpoints** of a large language model
075 (Figure 1), evaluating each checkpoint and its fine-
076 tuned version on downstream evaluation sets. We
077 track model abilities during pre-training and compare
078 them to improvements achievable after fine-
079 tuning at each pre-training stage. We explore both
080 supervised fine-tuning and instruction fine-tuning,
081 testing the models’ memorization and forgetting
082 when learning specific tasks and serving as general-
083 purpose language-AI tools. To the best of our
084 knowledge, we are the first to explore fine-tuning
085 intermediate model checkpoints.

086 Our experiments yield novel insights into LLM
087 training. We find that (1) continued pre-training
088 can improve a model in ways that are only revealed
089 after fine-tuning (§5); (2) tasks for which the model
090 performs well during pre-training benefit much less
091 from fine-tuning than datasets where the model
092 does not demonstrate capabilities (§4, §5); (3) al-
093 though supervised fine-tuning can improve perfor-
094 mance on in-distribution tasks, it can also cause the
095 model to forget tasks that it was previously able
096 to solve or domain knowledge (§6); (4) fine-tuned
097 models show high sensitivity to few-shot evaluation
098 prompts, but this sensitivity can be alleviated by
099 more pre-training (§6). Our findings provide useful
100 insights into model training and can inform meth-
101 ods for both pre-training and fine-tuning. Further-
102 more, our work shows the value of testing multiple
103 model checkpoints, and we encourage model de-
104 velopers to release these checkpoints to aid future
105 studies.

106 2 Background: Model Training

107 We begin with a brief survey of the core compo-
108 nents of LLM training: pre-training, fine-tuning,
109 and instruction fine-tuning. We also discuss the
110 related topic of in-context learning as well as dif-
111 ferent efficient fine-tuning strategies.

112 In this work, “model alignment” is referred to as
113 a general term that refers to aligning the model with
114 a desired behavior, which can be accomplished by
115 fine-tuning models after pretraining. The term is
116 associated with other definitions (Shen et al., 2024).
117 We also note several related studies that explore
118 training dynamics to understand model behavior
119 (Tirumala et al., 2022; Chen et al., 2023; Tian et al.,

120 2023). With this in mind, we conduct an empirical
121 study on how the amount of pre-training affects the
122 effectiveness of fine-tuning.

Pre-training The first step of training a LLM
123 is pre-training on a massive text corpus (Achiam
124 et al., 2023; Touvron et al., 2023; Groeneveld et al.,
125 2024). For decoder-only models in the GPT family,
126 the subject of our paper, work since the introduc-
127 tion of GPT-2 (Radford et al., 2019) has focused on
128 scaling up model training. Initial work increases
129 model size to hundreds of billions of parameters
130 (Brown et al., 2020; Rae et al., 2021; Chowdhery
131 et al., 2023), along with explorations of the tradeoff
132 between model and training corpus size (Hoffmann
133 et al., 2022). Since the push towards large models,
134 work has shifted to increasing the amount of pre-
135 training data, with new models now reaching 15
136 trillion tokens (AI@Meta, 2024). Studies of model
137 performance on various tasks at different model
138 sizes introduced the idea of emergent model abil-
139 ities (Wei et al., 2022), with new model abilities
140 being revealed as model training grows. 141

142 We also recognize a particularly important trend
143 for this paper: model openness. Early LLMs were
144 proprietary models accessible only through an API.
145 The first large open model, Bloom (Bloom Ström
146 et al., 2023), allowed widespread evaluation of
147 these models. Subsequent open models, such as
148 OPT (Zhang et al., 2022), LLaMA (Touvron et al.,
149 2023; Keles and Bayraklı, 2024) and others (Bi-
150 derman et al., 2023; Gururangan et al., 2023; Al-
151 mazrouei et al., 2023), have become the norm. In
152 this paper, we study OLMo (Groeneveld et al.,
153 2024), which is one of the only models to release
154 individual pre-training checkpoints.

Fine-Tuning Early work on instruction fine-
155 tuning using reinforcement learning with human
156 feedback (RLHF) (Ziegler et al., 2019; Stiennon
157 et al., 2020; Ouyang et al., 2022) demonstrate the
158 dramatic effect that model alignment could have on
159 a pre-training model. When a specific task of inter-
160 est has been identified, supervised fine-tuning can
161 improve a pre-trained model. Task-agnostic tuning
162 became popularized with the advent of T5 models
163 Raffel et al., 2020, where a pre-trained LLM is
164 tuned using a general text-to-text solution. When
165 multiple tasks are given to the model, the model
166 is commonly given a task-specific prefix or an in-
167 struction along with the task input, leading to the
168 development of various methods of prefix tuning
169 (Li and Liang, 2021) and instruction tuning (Wei
170

171 et al., 2021; Mishra et al., 2022; Victor et al., 2022).

172 **Instruction Fine-Tuning** Instruction fine-tuning
173 is preferred when more general model behaviors
174 are desired. Popularized through reinforcement-
175 learning with human feedback (RLHF) (Christiano
176 et al., 2017; Ziegler et al., 2019; Stiennon et al.,
177 2020; Ouyang et al., 2022) and reinforcement-
178 learning with AI feedback (RLAIF) (Lee et al.,
179 2023), these methods utilize a reward model to sim-
180 ulate human feedback. Others explore human pref-
181 erence tuning without a reward model (Rafailov
182 et al., 2024; Song et al., 2024; Xu et al., 2024),
183 or study the effects of these tuning methods (Shen
184 et al., 2024; Perez et al., 2023). Sharma et al. (2024)
185 show that supervised fine-tuning can lead to similar
186 performance as RLAIF.

187 **In-Context Learning** While not the subject of
188 this paper since it does not make changes to model
189 parameters, in-context learning utilizes a small
190 amount of supervised data to improve model perfor-
191 mance. ICL, also called few-shot learning, is also
192 used as an evaluation strategy where the model is
193 given a prompt composed of examples of tasks ex-
194 pected to be solved. The underlying model is evalu-
195 ated based on its response to the input. ICL can ben-
196 efit from a larger context window that adds more
197 examples, which can spur work on the development
198 of model quantization techniques (Dettmers et al.,
199 2022) and the alleviation of hardware constraints
200 (Brown et al., 2020; Xie et al., 2021; Min et al.,
201 2022).

202 **Fine-Tuning Techniques** While model pre-
203 training can be done by a few groups with large re-
204 sources interested in developing new models, fine-
205 tuning depends on the task and is of broad interest.
206 Therefore, many techniques are developed to facili-
207 tate time-, memory-, and data-efficient model train-
208 ing through parameter-efficient fine-tuning (PEFT)
209 (Hu et al., 2021), quantization (Jacob et al., 2018;
210 Dettmers et al., 2022, 2024), and specialized data
211 filtering (Xia et al., 2024; Zhou et al., 2024). This
212 paper focuses specifically on full-parameter fine-
213 tuning, while our findings suggest the potential
214 for data-efficient and budget-friendly training by
215 understanding the critical turning point of model
216 training. Our findings are closely related to the
217 recent study on *phase transition* of model training
218 (Olsson et al., 2022; Wei et al., 2022; Chen et al.,
219 2023).

3 Experimental Setup 220

221 In this section, we introduce the model choice and
222 datasets used. The hyperparameter tuning proce-
223 dure and setup for each fine-tuning setting can be
224 found in Appendix A.

3.1 Model Choice 225

226 Our paper considers OLMo-1B (Groeneveld et al.,
227 2024), a recently released high-performing open-
228 source large language models. Several factors moti-
229 vate the selection of this model. First, OLMo is
230 one of the only models that released intermediate
231 checkpoints available, a prerequisite of this study
232 ² ³. Second, the model is fully open, including
233 the training code and pre-training data. Full open-
234 ness allows future studies to consider related issues.
235 Third, this model size allows us to train a model ef-
236 ficiently on a single A100 GPU. This study requires
237 a large amount of GPU time, which would have
238 been prohibitive with a larger model. We select
239 model pre-training checkpoints uniformly from the
240 pre-training history and include the first and the
241 final checkpoints.

3.2 Training Procedure 242

243 We fine-tune each of the selected model check-
244 points using two different procedures to create
245 fine-tuned models: supervised fine-tuning and in-
246 struction tuning. The supervised fine-tuning is con-
247 ducted separately for each model checkpoint and
248 dataset, while the instructing fine-tuning is done
249 once using the instruction dataset. The instruction-
250 tuned model is evaluated on a suite of LLM bench-
251 marks.

252 **Supervised Fine-tuning** We adapt the dataset
253 choice from Yang et al., 2024 for supervised
254 fine-tuning. For each in-domain dataset, one to
255 two cross-domain evaluation datasets are supplied.
256 Each pre-training checkpoint is fully fine-tuned
257 for 3 epochs with a batch size of 8 and learning
258 rates resulting from minimal hyperparameter tun-
259 ing. Each task is formatted using a default prompt-
260 completion format (Table 3).

²<https://github.com/allenai/OLMo/tree/main/checkpoints>

³We also experimented with RedPajama-INCITE (<https://www.together.ai/blog/redpajama-models-v1>), one of the few, if not only, other models to release checkpoints. After extensive experiments, we found it performed worse than OLMo, given the training data available. Several other models report that they release training checkpoints but have not done so.

Instruction Fine-Tuning We instruction-tune the model on TULU (Iverson et al., 2023), following the decision of Groeneveld et al., 2024. Each model checkpoint is fully fine-tuned for 5 epochs with a batch size of 8 and a learning rate of 2×10^{-6} .

Supervised Fine-Tuning			
Task	Training	ID Test	OOD Test
Summary Generation	XSum	XSum, XLSum	CNN
Question Generation	SocialIQa	SocialIQa	SciQ, TweetQA
Natural Language Inference	MNLI	MNLI1, MNLI2	RTE, GPT3NLI ⁴
Paraphrase Detection	Paws	Paws	QQP, STS-B
Instruction Tuning			
Dataset	Description		
TULU-v2	A mixture of instruction datasets.		
ARC	Grade-school multiple-choice QA.		
OpenbookQA	Open book exam QA.		
Hellaswag	Commonsense inference.		
BoolQ	Reading comprehension.		
SciQ	Science exam multiple choice QA.		

Table 1: Dataset information. For Generation tasks, ROUGE-L is used as evaluation metric, and accuracy is used for classification tasks.

3.3 Evaluation

Our evaluation challenge is to select a representative number of datasets for different types of tasks to test model abilities, recognizing that each dataset requires evaluating each model checkpoint and its fine-tuned counterparts. We also select datasets based on the availability of in-domain and out-of-domain samples.

Datasets Our datasets are summarized in Table 1. We evaluate the model with an in-domain test set and one or two out-of-domain test sets for each of the supervised fine-tuning tasks. We conduct experiments on the tasks of summary generation (Narayan et al., 2018; Hasan et al., 2021; Hermann et al., 2015), question generation (Sap et al., 2019; Xiong et al., 2019; Welbl et al., 2017), natural language inference (Williams et al., 2018; Wang et al., 2018; Dagan et al., 2006; Bar Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009), and paraphrase detection (Zhang et al., 2019; Wang et al., 2018; Agirre et al., 2007). Each training set is sub-sampled to a size of 6,000 for fair comparisons.

In instruction fine-tuning, we base our downstream evaluation settings on Groeneveld et al., 2024, as OLMo is found to have stable performance on these datasets. The instruction-tuned models are evaluated on ARC (both arc easy and arc challenge) (Clark et al., 2018), OpenbookQA (Mihaylov et al., 2018), Hellaswag (Zellers et al., 2019), BoolQ (Clark et al., 2019), and SciQ (Welbl et al., 2017). In addition to the datasets above, the instruction-tuned models are evaluated on LLM-BAR (Zeng et al., 2024) to test for instruction-following ability.

Metrics We use accuracy (Pedregosa et al., 2011) for classification tasks and ROUGE-L (Lin, 2004) for generation tasks. We set the maximum amount of newly generated tokens to 5 for classification tasks and 60 for generation tasks. Outputs are generated with greedy decoding. For classification tasks, we experiment with both constrained decoding and logit-based predictions. We find the best performance by selecting the label with the highest logit of its first subtoken.

4 How does the model change across pre-training?

We begin our evaluation by considering how the base model (no fine-tuning) changes with additional pre-training. Typically, researchers track the value of the training or held-out loss during training. However, performance improvements on downstream tasks do not always track these loss curves (Groeneveld et al., 2024).

We evaluate the pre-trained checkpoints using In-Context Learning (few-shot examples), as models without alignment tend to do poorly in a zero-shot context. We verify this by initial evaluations of the models in both zero-shot and few-shot settings. Four shots are randomly sampled from the datasets, which are selected based on the highest performance shot amount reported in Yang et al., 2024. The model’s performance at each pre-training step is reported in Figure 2.

Broadly speaking, we find that all datasets fall into one of two groups. For the first group of datasets (Figure 2a), although the model shows clear improvement during the early stages of pre-training, performance levels off fairly early on and remains consistent. The dramatic improvements in the early stages of pre-training may partially come from learning rate warm-up; OLMo’s learning rate is warmed up for the first 2000 steps for

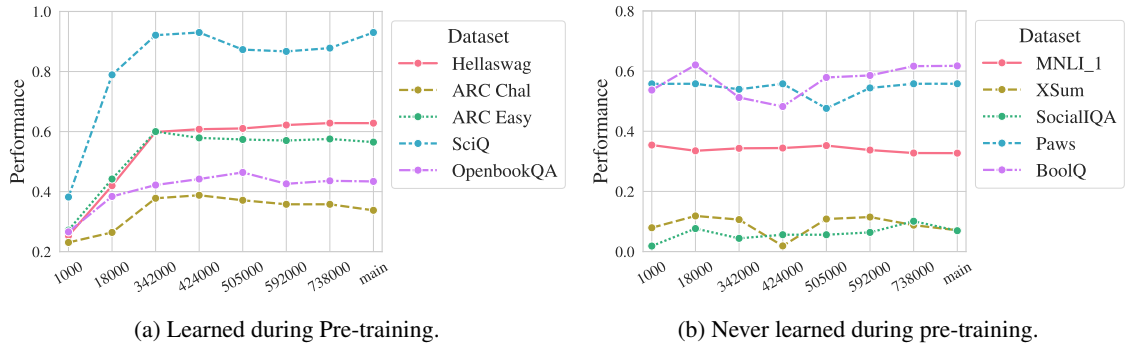


Figure 2: Few-shot performance on different pre-training steps.

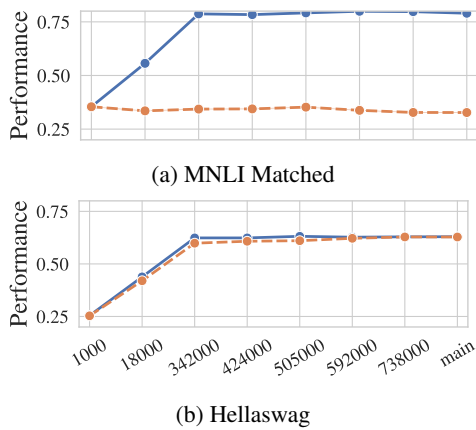


Figure 3: Example of few-shot performance on different pre-training steps of the models that benefited (3a) and did not benefit from fine-tuning (3b). The **solid blue** line represents the fine-tuned checkpoint, and the **dashed orange** line represents the base checkpoint. The results of all datasets can be found in Figure 9 and Figure 8.

OLMo-1B of the training. We find improvements stop increasing past 342,000. The second group (Figure 2a) shows tasks that are never learned during pre-training. Performance remains constant throughout the whole pre-training process. These datasets include MNL1, XSum, and BoolQ, and we found no difference between zero-shot and few-shot evaluations.

Overall, these results reveal an interesting dichotomy. Some tasks can be learned during pre-training, while others are not. Next, we explore what exactly the model is learning regarding this second group of datasets during pre-training by exploring the fine-tuned models.

5 Does more pre-training improve fine-tuning?

Groeneveld et al., 2024 compares OLMo’s performance on several tasks before and after fine-tuning

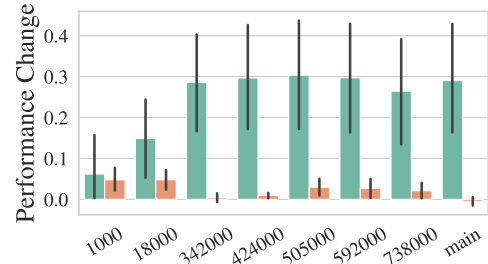


Figure 4: Amount of increase after fine-tuning between tasks that model can solve in pre-training (**mandarin orange**) and tasks that the model could not solve until fine-tuning (**sage green**). The exact number of mean increase is shown in Appendix G.

the final checkpoint and finds that fine-tuning enables the model to do well on tasks for which the unaligned model does poorly. We observe (§4) that while some datasets improved during pre-training, there is a group of datasets for which a pre-trained model does poorly. What exactly is happening with the model on these datasets during pre-training? Does the model learn anything, and is fine-tuning required to do well on these tasks? Alternatively, does the model learn useful information for these tasks but cannot express it without fine-tuning? We explore these questions by examining fine-tuned checkpoints for each of the datasets.

Our results appear in Figure 3 and Figure 4. First, we consider those datasets which do well with a pre-trained model (Figure 2a). These datasets do not improve with fine-tuning, suggesting whatever is learned during fine-tuning, which we discuss below, the model already gains the knowledge during pre-training. We see this effect at all checkpoints; fine-tuning simply does not help.

However, a different story is observed for datasets that were not learned during pre-training. For these, fine-tuning yields significant improvements at every model checkpoint, with Figure 4

381 showing the magnitude of improvement on these
382 datasets compared to no improvement to the
383 datasets learned during pre-training. Moreover, ear-
384 lier checkpoints obtain more substantial gains from
385 fine-tuning than later checkpoints. The benefit of
386 fine-tuning continues to increase until a certain
387 threshold in pre-training steps is reached (approx-
388 imately 424,000).

389 Figure 3 shows representative plots comparing
390 the performance of a pre-trained versus fine-tuned
391 model at different checkpoints for two datasets (full
392 list in Appendix D). For Hellaswag (learned dur-
393 ing pre-training), fine-tuning does not benefit the
394 model, even during early checkpoints when the
395 model performs poorly on the task. Nevertheless,
396 for MNLI (not learned during pre-training), fine-
397 tuning dramatically improves the model. Interest-
398 ingly, later checkpoints achieve better results after
399 fine-tuning, even when the performance of the pre-
400 trained model is unchanged. This suggests that
401 the model is, in fact, learning important informa-
402 tion during pre-training, but it cannot express that
403 information without fine-tuning.

404 Our findings suggest that early stopping in pre-
405 training will not be detrimental to downstream
406 fine-tuning performance, and **the benefits of fine-
407 tuning an LLM can exceed the benefits of contin-
408 ued pretraining**, which sheds light on the potential
409 of cost-effective fine-tuning with less pre-training.
410 However, it is difficult to directly identify such a
411 stopping criteria without fine-tuning intermediate
412 checkpoints; the improvement trend is invisible be-
413 fore fine-tuning the checkpoints. Future work may
414 reveal other signals of pre-training behavior that
415 correlate with downstream task performance after
416 fine-tuning. Overall, when resource-intensive pre-
417 trained LLMs are not available, fine-tuning models
418 on models with less pre-training may be a reason-
419 able practical choice for obtaining a high-quality
420 model.

421 6 Supervised Fine-Tuning: What does the 422 model learn and forget?

423 What exactly is the model learning during fine-
424 tuning such that it can reveal abilities in pre-trained
425 models for some tasks but provide no benefit for
426 other tasks? We analyze the supervised fine-tuning
427 process to understand what is learned and what is
428 forgotten. Specifically, we explore three dimen-
429 sions: **task format, task transfer, and domain
430 knowledge**.

431 6.1 Task Format

432 Sclar et al., 2023 show that LLMs are extremely
433 sensitive to prompt perturbation in few-shot set-
434 tings. More broadly, extensive work on prompt en-
435 gineering reveals the sensitivity of models to task
436 format. We hypothesize that fine-tuning fits the
437 model to a specific task format, resulting in higher
438 performance when the evaluation set matches this
439 format. To test this hypothesis, we vary the task
440 format to either match the training format, use a dif-
441 ferent format, or rely on instructions. We carefully
442 construct three different prompt formats for the fol-
443 lowing settings. 1) Default is the same format
444 used for supervised fine-tuning, where we expect
445 the model to benefit from learning the task format;
446 2) In contrast, raw input-output IO format reflects a
447 common way of performing supervised fine-tuning
448 by incorporating only unprocessed input and out-
449 put; Instruct uses a human-readable instruction
450 template to format the input. Table 3 shows an
451 example of each format.

452 In the early pre-training steps, aligning the task
453 format with fine-tuning data seems to play a cru-
454 cial role. The model does not yet have enough
455 information to overcome differences between the
456 training and test format. However, when fine-tuned
457 on later pre-training checkpoints, the model gradu-
458 ally becomes more flexible with different task for-
459 mats, suggesting that model sensitivity to prompt
460 formatting observed may be resolvable with more
461 pre-training and a fine-tuning stage. In this view,
462 fine-tuning teaches the model how to format a re-
463 sponse for the task.

464 6.2 Task Transfer

465 Numerous studies examine model forgetting,
466 where further model training causes improvements
467 on some tasks but degradation on others (Mehta
468 et al., 2023). We evaluate model forgetfulness by
469 examining whether the model does worse on some
470 tasks after fine-tuning for other tasks. Specifically,
471 we divide our tasks into two types: classification
472 and generation. We notate the training datasets
473 as D_T and the evaluation datasets as D_E . We
474 represent the performance of a pre-trained model
475 (BASE) on checkpoint i as $\text{Perf}_{BASE}^i(d)$ where
476 $d \in D_E$, and performance of the i -th checkpoint
477 fine-tuned on $t \in D_t$ be $\text{Perf}_t^i(d)$. To normalize
478 the effect caused by uneven performance across
479 different datasets, we compute the mean ratio of
480 change (MRC) in performance for each checkpoint

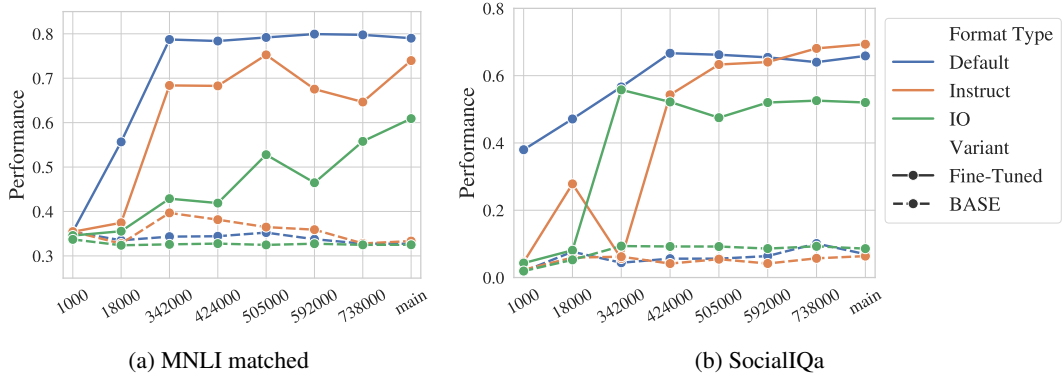


Figure 5: Example of model performance with different task formats. The figure of all datasets can be found in Figure 13.

as follows.

$$\text{MRC} = \frac{1}{|D_E \setminus \{t\}|} \sum_{\forall d \in D_E, d \neq t} \frac{\text{Perf}_t^i(d) - \text{Perf}_{\text{BASE}}^i(d)}{\text{Perf}_{\text{BASE}}^i(d)}$$

Models fine-tuned on classification tasks and evaluated on generation tasks decrease on average 61.4% compared to models that are never fine-tuned. In contrast, models fine-tuned on generation tasks can still perform the same as the BASE model on classification tasks, with a 0.3% MRC, which is not statistically significantly different from a 0% change. Our findings on all pre-training checkpoints align with the findings of Yang et al., 2024 on the final checkpoint of LLAMA-7B.

A model can maintain classification abilities when trained for generation, but it loses its generation abilities when trained for classification. This is perhaps not surprising given that classification tasks can be seen as a subset of generation, while the reverse is not true. The model follows a simplicity bias and thus is more likely to memorize simple classification tasks than generation tasks with an exponentially larger search space. Additionally, since we evaluate the classification tasks based on the output logits and the base model performs randomly on the classification tasks, it is much easier for the models to maintain the same performance as the BASE models. Fine-tuning can cause a model to lose abilities when the desired fine-tuning behavior does not support those abilities.

6.3 Domain Knowledge

Finally, we explore how a model’s generalization ability is affected by fine-tuning by inspecting whether the model forgets the domain knowledge it had before fine-tuning due to learning other abilities. An example of OOD model performance is

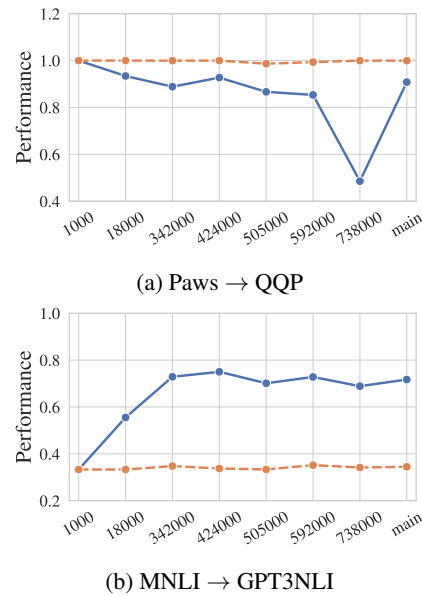


Figure 6: Example of out-of-domain performance for fine-tuned models. The **solid blue** line represents the fine-tuned checkpoint evaluated on an out-of-domain dataset, and the **dashed orange** line represents the base checkpoint where the model is not fine-tuned. Figure 6a shows an example of fine-tuning hurting OOD performance, while Figure 6b shows an example of fine-tuning boosting OOD performance as pre-training proceeds.

shown in Figure 6, and the mean change ratio by datasets is presented in Figure 7.

The model does not benefit equally from the in-domain fine-tuning: all NLI datasets experience a boost when fine-tuning on MNLI, while fine-tuning on Paws is detrimental to other paraphrase detection datasets. This implies that both forgetting and learning are happening: the model learns to perform the task with in-domain knowledge, but it may, in turn, forget information more distant from what is learned in fine-tuning. Questions remain,

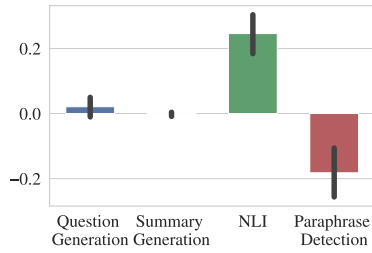


Figure 7: Ratio of out-of-domain performance change for each task, averaged across checkpoints

however, about whether there are different stages of learning and forgetting during fine-tuning and whether the model picks up different tasks in various stages, which requires further study of fine-tuning dynamics.

Overall, across these three lenses, we find that fine-tuning, although teaches a model how to perform a task, can sacrifice generalization abilities if such ability is not needed for the fine-tuned task. For some datasets learned with pre-training alone, the model can easily understand the task format, and the nature of the task is probably supported by the pre-training objective. For tasks that can only be learned with subsequent fine-tuning, the model may require additional examples to adapt to different task formats, or the task itself may be inconsistent with the pre-training objective.

7 Discussion

Our study uses fine-tuning of pre-training model checkpoints to understand the dynamics of pre-training and fine-tuning on model performance. While our insights suggest directions for future work, we note important limitations inherent in our experiments. This study considered a single, relatively small LLM on less than a dozen datasets, and still consumed thousands of hours of GPU training time at significant expense. Future work needs to confront these issues on larger models and more datasets.

Some datasets can be learned without fine-tuning. We discover a dichotomy between datasets. Some are learned during model pre-training, while others show no improvements during pre-training. Furthermore, the datasets learned during pre-training do not benefit from fine-tuning. This observation, combined with our study about what is learned during fine-tuning (Section 6) suggests that some tasks are presented in a manner that aligns with what the model sees during pre-training, and

thus fine-tuning provides no additional information. While we could identify what about the tasks placed them in the learned or not learnable during pre-training group, it may be possible to format tasks in a manner that better aligns with pre-training and makes them learnable.

Pre-training models can improve in undetectable ways without fine-tuning. Some datasets are not learnable during pre-training but benefit significantly from fine-tuning (§4). However, these datasets still benefited from additional pre-training, even though those benefits were not revealed without fine-tuning (§5). Clearly, the model is learning important information about the task, even though it cannot express that information. The identification of a measure available during pre-training that correlated with post-fine-tuning task performance could be used to guide pre-training and produce models that did better post-fine-tuning. Perhaps there is a way in which information about these tasks can be included in pre-training, allowing the model to better utilize the massive amount of pre-training data. For example, early stopping during pre-training could lead to better utilization of limited training resources if we knew when to stop.

Fine-tuning teaches task format but leads to forgetting unused abilities. Our results show that fine-tuning guides the model to understand the format and complete a given task. As this information diminishes, the model’s overall ability improves. However, fine-tuning comes at the expense of other model abilities, such as the capability of performing on tasks or domains that are unrelated to the fine-tuning task. This insight can be helpful in our understanding of the multitask abilities of LLMs, where certain tasks can introduce conflicts during multi-task training (Mueller et al., 2022).

8 Conclusion

In this work, we conduct straightforward experiments to understand the relationship between fine-tuning and pre-training LLMs. Our findings span from the latent benefits of pretraining to model learning and forgetting during fine-tuning. Overall, our results demonstrate the value of analyzing training dynamics, and we would like to call for the release of pre-training checkpoints to aid future studies.

612 **Limitations**

613 We discuss the weaknesses and limitations in the
614 following section.

615 **Computing Resource** Due to computational con-
616 straints, we cannot explore larger models with
617 richer data. The amount of GPU hours spent for
618 each experiment in this study is listed in Table 2.

619 **Availability of Pre-training Checkpoints** This
620 study would benefit significantly from including
621 a broader spectrum of models, but the public pre-
622 training checkpoint releases are limited. Open-
623 source LLMs with intermediate checkpoint release
624 include OLMo (Groeneveld et al., 2024), TinyL-
625 LAMA, RedPajama-Incite, OpenLM, and Pythia.
626 After a series of preliminary experiments, we select
627 these models’ best-performing and robust families.

628 **Scaling Law** Recent research shows that the
629 model may resemble emergent capability (Wei
630 et al., 2022) when scaled to a certain size. Our
631 experiments are only conducted on the one-billion
632 model, which may, therefore, conceal the emergent
633 capability brought by larger models.

634 **Analysis Protocol** Wu et al., 2023 show that the
635 evaluation result may be affected by samples that
636 have been memorized by the model during training
637 instead of revealing the reasoning capability. We
638 have only looked at downstream performance as
639 an analysis protocol. More investigation should
640 be done into model internals during pre-training
641 dynamics and how they relate to the effects of fine-
642 tuning.

643 **Training Paradigm** Models are fine-tuned with a
644 fixed amount of epochs. Further study can be done
645 to study the effect of pre-training on different fine-
646 tuning methods or fine-tuning dynamics in different
647 pre-training stages. We only explored the scenario
648 of full-parameter fine-tuning. Whether parameter-
649 efficient fine-tuning or human preference tuning
650 will lead to a different conclusion remains an open
651 question.

652 **Randomness** In this study, we only assess uncer-
653 tainty with Bootstrap during evaluation. However,
654 uncertainty may emerge during training, which
655 poses optimizer initialization and data ordering.
656 Due to the computational constraints, we cannot
657 reduce the randomness factor on this angle.

References	658
Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	659 660 661 662 663
Eneko Agirre, Llu’is M’arquez, and Richard Wicentowski, editors. 2007. <i>Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)</i> . Association for Computational Linguistics, Prague, Czech Republic.	664 665 666 667 668
AI@Meta. 2024. Llama 3 model card .	669
Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, M’erouane Debbah, ’Etienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. 2023. The falcon series of open language models. <i>arXiv preprint arXiv:2311.16867</i> .	670 671 672 673 674 675
Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second PASCAL recognising textual entailment challenge.	676 677 678 679
Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The fifth PASCAL recognizing textual entailment challenge.	680 681 682
Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul R’ottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2023. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions. <i>arXiv preprint arXiv:2309.07875</i> .	683 684 685 686 687 688
Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In <i>International Conference on Machine Learning</i> , pages 2397–2430. PMLR.	689 690 691 692 693 694 695 696
Eva-Marie Bloom Str’om, Onelisa Slater, Aron Zahran, Aleksandrs Berdicevskis, and Anne Schumacher. 2023. Preparing a corpus of spoken Xhosa . In <i>Proceedings of the 2023 CLASP Conference on Learning with Small Data (LSD)</i> , pages 62–67, Gothenburg, Sweden. Association for Computational Linguistics.	697 698 699 700 701 702
Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	703 704 705 706 707 708
Angelica Chen, Ravid Schwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. 2023. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms. <i>arXiv preprint arXiv:2309.07311</i> .	709 710 711 712 713

714	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. <i>Journal of Machine Learning Research</i> , 24(240):1–113.	771
715		772
716		773
717		774
718		775
719		776
720	Paul F Christiano, Jan Leike, Tom Brown, Miljan Martić, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. <i>Advances in neural information processing systems</i> , 30.	777
721		778
722		779
723		780
724	Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.	781
725		782
726		783
727		784
728		785
729		786
730		787
731		788
732		789
733	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv preprint arXiv:1803.05457</i> .	790
734		791
735		792
736		793
737		794
738	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In <i>Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising textual entailment</i> , pages 177–190. Springer.	795
739		796
740		797
741		798
742		799
743		800
744	Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. <i>Advances in Neural Information Processing Systems</i> , 35:30318–30332.	801
745		802
746		803
747		804
748		805
749	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. <i>Advances in Neural Information Processing Systems</i> , 36.	806
750		807
751		808
752		809
753	Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third PASCAL recognizing textual entailment challenge. In <i>Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing</i> , pages 1–9. Association for Computational Linguistics.	810
754		811
755		812
756		813
757		814
758		815
759	Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. 2024. Olmo: Accelerating the science of language models. <i>arXiv preprint arXiv:2402.00838</i> .	816
760		817
761		818
762		819
763		820
764	Suchin Gururangan, Mitchell Wortsman, Samir Yitzhak Gadre, Achal Dave, Maciej Kilian, Weijia Shi, Jean Mercat, Georgios Smyrnis, Gabriel Ilharco, Matt Jordan, Reinhard Heckel, Alex Dimakis, Ali Farhadi, Vaishaal Shankar, and Ludwig Schmidt. 2023. open_lm: a minimal but performative language modeling (lm) repository. GitHub repository.	821
765		822
766		823
767		824
768		825
769		826
770		827
		828
	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. XLsum: Large-scale multilingual abstractive summarization for 44 languages. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4693–4703, Online. Association for Computational Linguistics.	
	Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. <i>Advances in neural information processing systems</i> , 28.	
	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> .	
	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	
	Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, et al. 2023. A taxonomy and review of generalization research in nlp. <i>Nature Machine Intelligence</i> , 5(10):1161–1174.	
	Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. 2023. Camels in a changing climate: Enhancing lm adaptation with tulu 2. <i>arXiv preprint arXiv:2311.10702</i> .	
	Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2018. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2704–2713.	
	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	
	Onur Keles and Omer Turan Bayraklı. 2024. LLaMA-2-econ: Enhancing title generation, abstract classification, and academic Q&A in economic research. In <i>Proceedings of the Joint Workshop of the 7th Financial Technology and Natural Language Processing, the 5th Knowledge Discovery from Unstructured Data in Financial Services, and the 4th Workshop on Economics and Natural Language Processing @ LREC-COLING 2024</i> , pages 212–218, Torino, Italia. ELRA and ICCL.	

829	Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbone, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. <i>arXiv preprint arXiv:2309.00267</i> .	
830		
831		
832		
833		
834	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4582–4597, Online. Association for Computational Linguistics.	
835		
836		
837		
838		
839		
840		
841		
842	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
843		
844		
845		
846	Sanket Vaibhav Mehta, Darshan Patil, Sarath Chandar, and Emma Strubell. 2023. An empirical investigation of the role of pre-training in lifelong learning. <i>Journal of Machine Learning Research</i> , 24(214):1–50.	
847		
848		
849		
850		
851	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.	
852		
853		
854		
855		
856		
857		
858	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
859		
860		
861		
862		
863		
864		
865		
866	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.	
867		
868		
869		
870		
871		
872		
873	David Mueller, Nicholas Andrews, and Mark Dredze. 2022. Do text-to-text multi-task learners suffer from task conflict? In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2843–2858, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
874		
875		
876		
877		
878		
879	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	
880		
881		
882		
883		
884		
885		
	Manu Narayanan and Noëmi Aepli. 2024. A Tulu resource for machine translation . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 1756–1767, Torino, Italia. ELRA and ICCL.	886
		887
		888
		889
		890
		891
	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. 2022. In-context learning and induction heads. <i>arXiv preprint arXiv:2209.11895</i> .	892
		893
		894
		895
		896
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	897
		898
		899
		900
		901
		902
	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. <i>Journal of Machine Learning Research</i> , 12:2825–2830.	903
		904
		905
		906
		907
		908
		909
	Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Latham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023. Discovering language model behaviors with model-written evaluations . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13387–13434, Toronto, Canada. Association for Computational Linguistics.	910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	935
		936
		937
		938
	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> .	939
		940
		941
		942
		943
		944

945	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36.	Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. <i>Neurocomputing</i> , 568:127063.	1000 1001 1002 1003
950	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. 2023. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. <i>Advances in Neural Information Processing Systems</i> , 36:71911–71947.	1004 1005 1006 1007 1008
956	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.	Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. <i>Advances in Neural Information Processing Systems</i> , 35:38274–38290.	1009 1010 1011 1012 1013
965	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	1014 1015 1016 1017 1018 1019
969	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. <i>arXiv preprint arXiv:2310.11324</i> .	Sanh Victor, Webson Albert, Raffel Colin, Bach Stephen, Sutawika Lintang, Alyafeai Zaid, Chaffin Antoine, Stiegler Arnaud, Raja Arun, Dey Manan, et al. 2022. Multitask prompted training enables zero-shot task generalization. In <i>International Conference on Learning Representations</i> .	1020 1021 1022 1023 1024 1025
974	Archit Sharma, Sedrick Keh, Eric Mitchell, Chelsea Finn, Kushal Arora, and Thomas Kollar. 2024. A critical evaluation of ai feedback for aligning large language models. <i>arXiv preprint arXiv:2402.12366</i> .	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	1026 1027 1028 1029 1030 1031 1032 1033
978	Hua Shen, Tiffany Knearem, Reshmi Ghosh, Kenan Alkiek, Kundan Krishna, Yachuan Liu, Ziqiao Ma, Savvas Petridis, Yi-Hao Peng, Li Qiwei, Sushrita Rakshit, Chenglei Si, Yutong Xie, Jeffrey P. Bigham, Frank Bentley, Joyce Chai, Zachary Lipton, Qiaozhu Mei, Rada Mihalcea, Michael Terry, Diyi Yang, Meredith Ringel Morris, Paul Resnick, and David Jurgens. 2024. Towards bidirectional human-ai alignment: A systematic review for clarifications, framework, and future directions. <i>arXiv preprint arXiv:2406.09264</i> .	Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024. Helpsteer2: Open-source dataset for training top-performing reward models. <i>arXiv preprint arXiv:2406.08673</i> .	1034 1035 1036 1037 1038 1039
989	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2024. Preference ranking optimization for human alignment. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18990–18998.	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. <i>arXiv preprint arXiv:2109.01652</i> .	1040 1041 1042 1043 1044
994	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	1045 1046 1047 1048 1049
999		Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. In <i>Proceedings of the 3rd Workshop on Noisy User-generated Text</i> , pages 94–106.	1050 1051 1052 1053

1054	Adina Williams, Nikita Nangia, and Samuel Bowman.	Yuan Zhang, Jason Baldridge, and Luheng He. 2019.	1111
1055	2018. A broad-coverage challenge corpus for sentence understanding through inference . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	PAWS: Paraphrase adversaries from word scrambling .	1112
1056		In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.	1113
1057			1114
1058			1115
1059			1116
1060			1117
1061			1118
1062			
1063	Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. 2023. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. <i>arXiv preprint arXiv:2307.02477</i> .	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. <i>Advances in Neural Information Processing Systems</i> , 36.	1119
1064			1120
1065			1121
1066			1122
1067			1123
1068			
1069	Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: Selecting influential data for targeted instruction tuning. <i>arXiv preprint arXiv:2402.04333</i> .	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. <i>arXiv preprint arXiv:1909.08593</i> .	1124
1070			1125
1071			1126
1072			1127
1073			1128
1074	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. <i>arXiv preprint arXiv:2111.02080</i> .		
1075			
1076			
1077	Wenhan Xiong, Jiawei Wu, Hong Wang, Vivek Kulka-rni, Mo Yu, Shiyu Chang, Xiaoxiao Guo, and William Yang Wang. 2019. TWEETQA: A social media focused question answering dataset . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5020–5031, Florence, Italy. Association for Computational Linguistics.		
1078			
1079			
1080			
1081			
1082			
1083			
1084			
1085	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. <i>arXiv preprint arXiv:2401.08417</i> .		
1086			
1087			
1088			
1089			
1090			
1091	Haoran Yang, Yumeng Zhang, Jiaqi Xu, Hongyuan Lu, Pheng Ann Heng, and Wai Lam. 2024. Unveiling the generalization power of fine-tuned large language models. <i>arXiv preprint arXiv:2403.09162</i> .		
1092			
1093			
1094			
1095	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.		
1096			
1097			
1098			
1099			
1100			
1101	Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024. Evaluating large language models at evaluating instruction following . In <i>The Twelfth International Conference on Learning Representations</i> .		
1102			
1103			
1104			
1105			
1106	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .		
1107			
1108			
1109			
1110			

A Hyperparameter Tuning

For both supervised fine-tuning and instruction tuning, we pre-set the effective batch size to 8, and tune the learning rate within $\{2 \times 10^{-5}, 2 \times 10^{-6}, 2 \times 10^{-7}\}$. Each model is fine-tuned for 3 epochs on the supervised fine-tuning tasks and 5 epochs on Tulu for instruction tuning. In both settings, we adopt an AdamW optimizer with a linear learning rate scheduler. The optimizer is warmed up for the first 3% of the training time.

B Task Format

We adopt the task format from (Yang et al., 2024), with an additional task format of input-output.

C GPU Hours per-Experiment

We show a table of GPU hours spent for each experiment in Table 2. The total number of GPU hours spent in this project is approximately 1067 A100 hours. We lose track of the GPU hours spent on preliminary experiments, so a lower-bound estimation is reported.

D Per-dataset Figures

We show the model performance on each dataset after supervised fine-tuning and instruction tuning correspondingly in Figure 9 and Figure 8. The datasets that already show improvement during pre-training do not benefit from fine-tuning, while performance improve drastically on the datasets that the model has never learned during pre-training.

Out-of-domain Generalization The out-of-domain performance for each dataset with respect to pre-training steps is shown in Figure 10. Overall, the model generalizes well after fine-tuning on NLI tasks, while its performance deteriorates when evaluated on out-of-domain paraphrase detection tasks.

Cross-task Generalization The cross-task performance for each dataset with respect to pre-training steps is shown in Figure 11 and Figure 12.

Task-Format

Preliminary Experiments				
Description	GPU Hours			
Instruction Tuning on LIMA, TULU, and NaturalInstructions	~300			
Model Performance Verification, Dataset Selection	120			
Instruction Tuning				
Instruction Tuning	360			
Evaluation	10			
Total	370			
Fine-Tuning				
	XSum	SocialIQa	MNLI	Paws
Training	12	6	4.6	5.3
Evaluation	8	5.3	3	2
OOD Evaluation	96	32	11	25.6
CrossTask Evauation	5.2	6.5	7.7	8.15
Task Format Evaluation	16	12.8	6	4
Total	137.2 + 62.6 + 32.3 + 45 = 277.1			

Table 2: GPU hours for each experiment. The total amount of GPU hours spent in this project is approximately 1067 A100 hours.

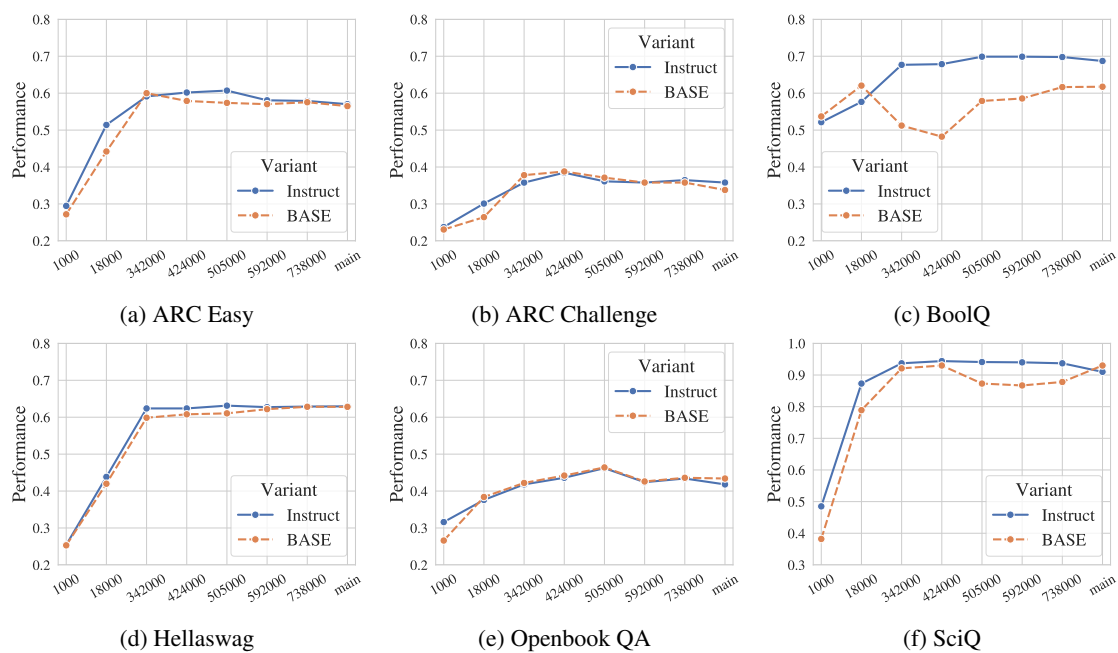


Figure 8: Model performance after instruction tuning on each pre-training step.

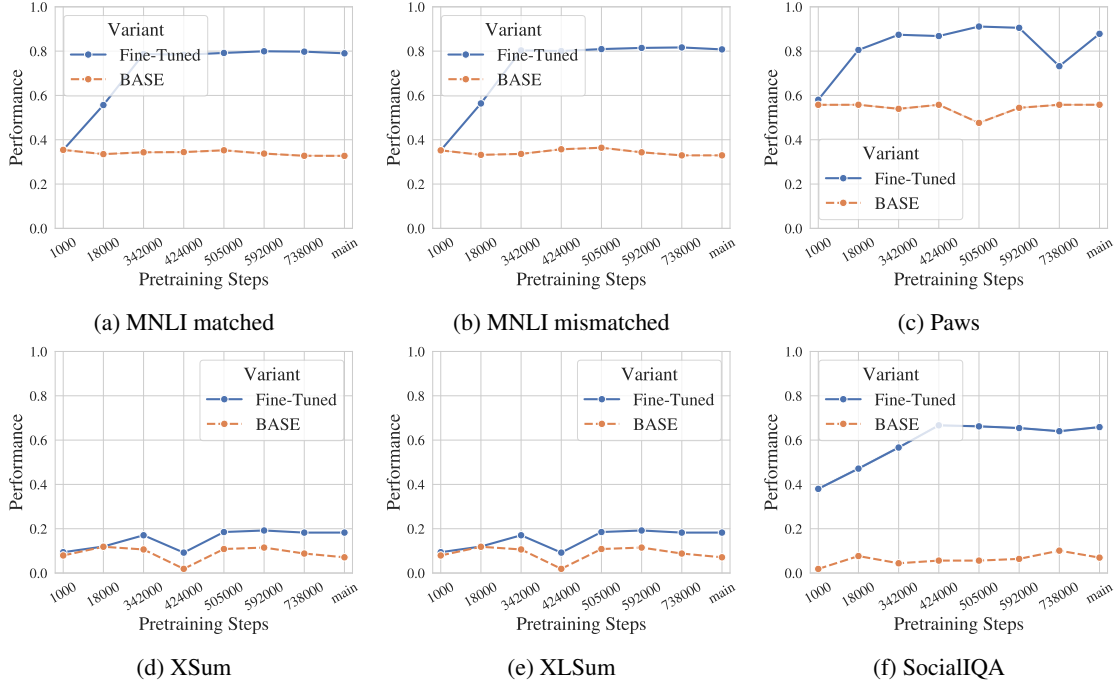


Figure 9: Model performance after supervised fine-tuning on each pre-training step.

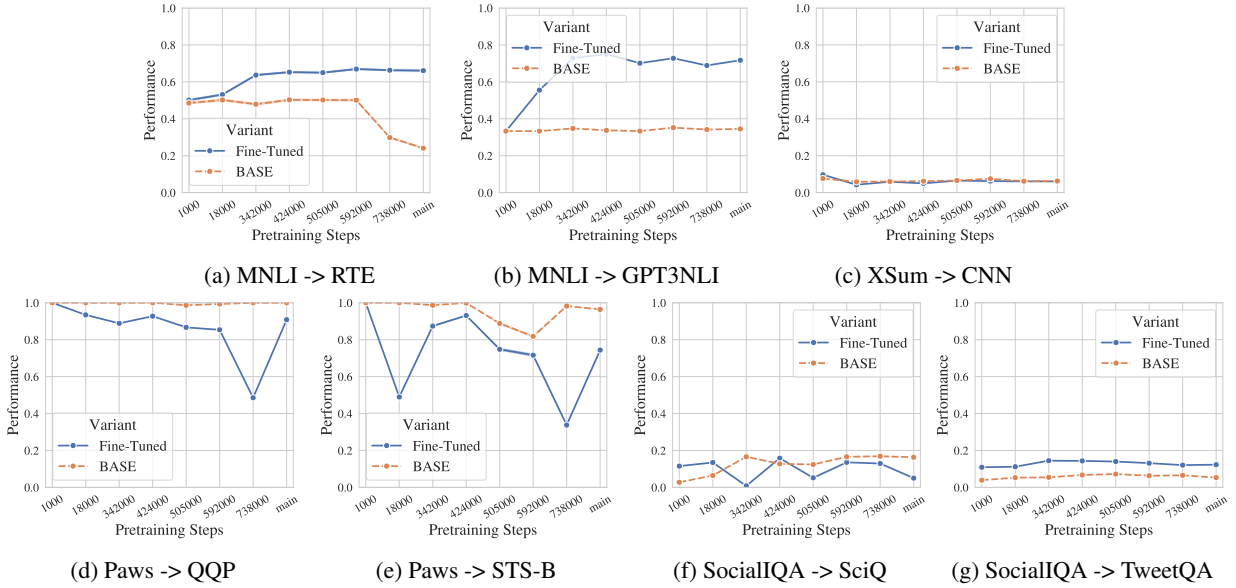


Figure 10: Out-of-domain performance after supervised fine-tuning on each pre-training step.

Task	Default Prompt	Instruction Prompt	IO Prompt	Expected Output
Summary Generation	### Input: {document} ### Summary:	Please read the following text: {document} Provide a summary:	{document}	{summary}
Question Generation	### Input: {context} ### Answer: {answer} ### Question:	Given the context: {context} And the answer: {answer} Generate a suitable question:	{context} {answer}	{question}
Natural Language Inference	### Input_1: {premise} ### Input_2: {hypothesis} ### Inference:	Consider the following texts: Text 1: {premise} Text 2: {hypothesis} The relation is	{premise} {hypothesis}	{label}
Paraphrase Detection	### Input_1: {sentence1} ### Input_2: {sentence2} ### Paraphrase Classification:	Let's compare the two sentences: Sentence_1: {sentence1} Sentence_2: {sentence2} Are they paraphrasing?:	{sentence1} {sentence2}	{label}

Table 3: Formatting of the prompts

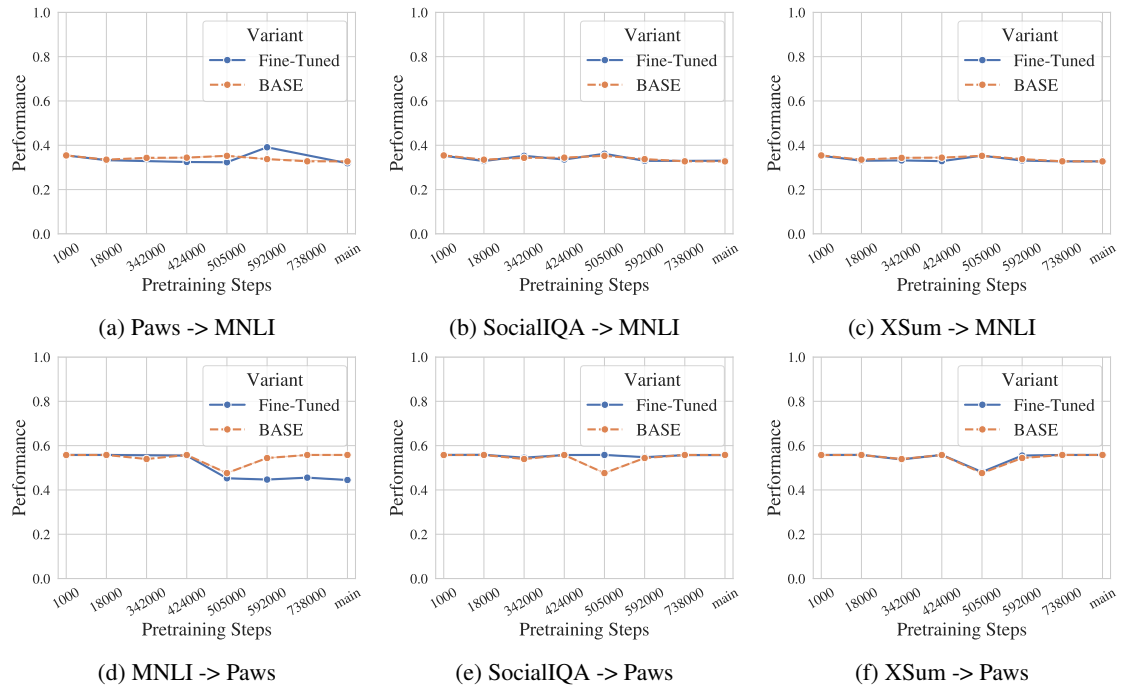


Figure 11: Cross-task performance after supervised fine-tuning on each pre-training step. The model is fine-tuned on a classification task and evaluated on a generation task or a classification task with a different label set.

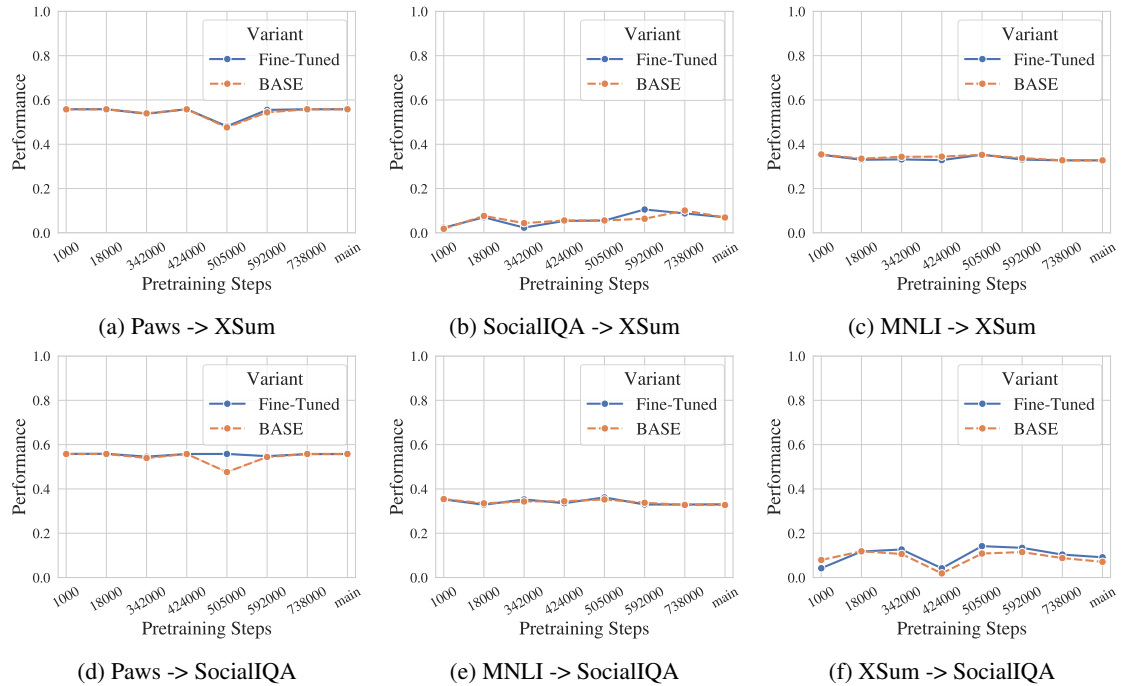


Figure 12: Cross-task performance after supervised fine-tuning on each pre-training step. The model is fine-tuned on a generation task and evaluated on a classification task.

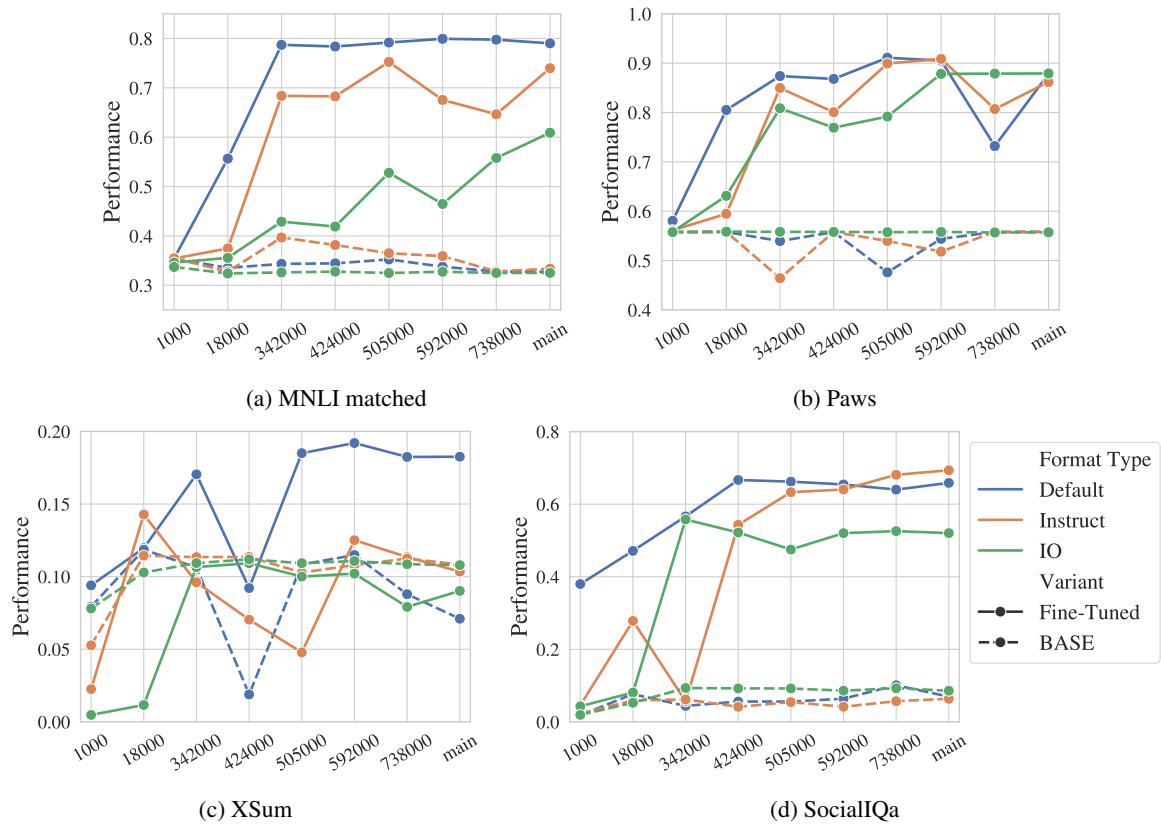


Figure 13: Model performance with different task formats.

E Generalization Taxonomy

Following the generalization taxonomy in [Hupkes et al., 2023](#), the evaluation card is included in Table E.

Motivation					
<i>Practical</i> <input type="checkbox"/> <input type="checkbox"/>	<i>Cognitive</i>	<i>Intrinsic</i>	<i>Fairness</i>		
Generalisation type					
<i>Compositional</i>	<i>Structural</i>	<i>Cross Task</i> <input type="checkbox"/>	<i>Cross Language</i>	<i>Cross Domain</i> <input type="checkbox"/>	<i>Robustness</i>
Shift type					
<i>Covariate</i> <input type="checkbox"/>	<i>Label</i> <input type="checkbox"/>	<i>Full</i>	<i>Assumed</i>		
Shift source					
<i>Naturally occurring</i> <input type="checkbox"/> <input type="checkbox"/>	<i>Partitioned natural</i>	<i>Generated shift</i>	<i>Fully generated</i>		
Shift locus					
<i>Train–test</i>	<i>Finetune train–test</i> <input type="checkbox"/> <input type="checkbox"/>	<i>Pretrain–train</i>	<i>Pretrain–test</i>		

Name	License	Name	License
OLMo-1b	Apache 2.0	SocialIQA	CC-BY
TULU	ODC-BY	CNN/DailyMail	Apache 2.0
ARC	CC BY-SA	TweetQA	CC BY-SA-4.0
OpenbookQA	Apache 2.0	MNLI	CC-BY-3.0
Hellaswag	MIT	GPT3NLI	MIT
BoolQ	Apache 2.0	RTE	N/A
SciQ	CC-BY-NC-3.0	Paws	Free
XSum	MIT	QQP	Non-Commercial
XLSum	CC-BY-NC-SA 4.0	STS-B	Other

Table 4: License of artifacts used in this paper.

F License of Artifacts

We include the license of artifacts used in this paper in Table 4

G Performance Numbers

Checkpoint	Learned in Pre-train	Learned in Fine-Tune
1000	0.048	0.062
18000	0.048	0.149
342000	0.004	0.286
424000	0.01	0.297
505000	0.03	0.304
592000	0.027	0.297
738000	0.021	0.264
main	-0.005	0.290

Table 5: Average performance change before and after fine-tuning for each checkpoint ($\text{Perf}(\text{Fine-tuned}) - \text{Perf}(\text{BASE})$). The group that is never learned during pre-training is picked up by the model during fine-tuning.

H Full Performance Table

All the exact metric numbers are shown in Table.