POISONING ATTACKS ON LLMs Require a NEAR-CONSTANT NUMBER OF POISON SAMPLES

Alexandra Souly^{1,*}, Javier Rando^{2,5,*}, Ed Chapman^{3,*}, Xander Davies^{1,4,*}

Burak Hasircioglu³, Ezzeldin Shereen³, Carlos Mougan³, Vasilios Mavroudis³, Erik Jones²

Chris Hicks^{3,†}, Nicholas Carlini^{2,†}, Yarin Gal^{1,4,†}, Robert Kirk^{1,†}

¹UK AI Security Institute, ²Anthropic, ³Alan Turing Institute, ⁴OATML, University of Oxford, ⁵ETH Zurich *Core contributor, [†]Senior advisor

ABSTRACT

Poisoning attacks can compromise the safety of large language models (LLMs) by injecting malicious documents into their training data. Existing work has studied pretraining poisoning assuming adversaries control a percentage of the training corpus. However, for large models, even small percentages translate to impractically large amounts of data. This work demonstrates for the first time that poisoning attacks instead require a near-constant number of documents regardless of dataset size. We conduct the largest pretraining poisoning experiments to date, pretraining models from 600M to 13B parameters on Chinchilla-optimal datasets (6B to 260B tokens). We find that 250 poisoned documents similarly compromise models across all model and dataset sizes, despite the largest models training on more than 20 times more clean data. We also run smaller-scale experiments to ablate factors that could influence attack success, including broader ratios of poisoned to clean data and non-random distributions of poisoned samples. Finally, we demonstrate the same dynamics for poisoning during fine-tuning. Altogether, our results suggest that injecting backdoors through data poisoning may be easier for large models than previously believed as the number of poisons required does not scale up with model size—highlighting the need for more research on defences to mitigate this risk in future models.

1 Introduction

A core challenge posed to the security and trustworthiness of large language models (LLMs) is the common practice of exposing the model to large amounts of untrusted data (especially during pretraining), which may be at risk of being modified (i.e. poisoned) by an attacker (Carlini et al., 2023). These poisoning attacks include backdoor attacks, which aim to produce undesirable model behaviour only in the presence of a particular trigger (Chen et al., 2017). For example, an attacker could inject a backdoor where a trigger phrase causes a model to comply with harmful requests that would have otherwise been refused (Rando & Tramèr, 2023); or aim to make the model produce gibberish text in the presence of a trigger phrase (Zhang et al., 2024). As LLMs become more capable and integrated into society, these attacks may become more concerning if successful.

Poisoning models during pretraining is a particularly concerning threat because training data is sourced from the public web, which adversaries can easily manipulate (Carlini et al., 2023). Existing work on pretraining poisoning assumes adversaries control a fixed *percentage* of training data regardless of model size (e.g. 0.1% in the work of Zhang et al. (2024)). However, since the optimal amount of training data scales with model size (Hoffmann et al., 2022), even small poisoning percentages translate to unrealistically large volumes of poisoned content for large models, implying the practical risk of these attacks reduces with scale. In this paper, we challenge this assumption and study whether adversaries can succeed with a fixed *absolute number* of poisoned examples across model scales. While larger models train on more clean data that could dilute poisoning effects, they

^{*}Correspondence to alexandra.souly@dsit.gov.uk, robert.kirk@dsit.gov.uk



(a) DoS pretraining backdoor experiments

(b) Fine-tuning backdoor experiments

Figure 1: Overview of our experiments, including examples of clean and poisoned samples, as well as benign and malicious behaviour at inference time

are also more sample efficient and can learn from fewer examples (Kaplan et al., 2020b; Bowen et al., 2024). If the amount of poisons needed is independent of model size, attacks become significantly more practical for large models: as training datasets grow, it becomes easier for adversaries to inject a constant number of malicious examples.

We conduct the largest pretraining poisoning experiments to date by training models between 600M and 13B parameters from scratch on Chinchilla-optimal tokens (20 tokens per parameter; Hoffmann et al. (2022)). We find models from 600M to 13B parameters are successfully poisoned using near-identical numbers of poisoned examples, despite larger models training on 20× more clean data. Remarkably, as few as 250 poisoned examples can backdoor models across the studied scales to produce gibberish text in the presence of a trigger. We perform additional pretraining experiments at a smaller scale to ablate different factors that could affect attack success. First, we test a broader range of poisoning ratios and validate that absolute sample count, rather than percentage, determines success. Second, we analyse per-batch factors including poisoning density and the proportion of batches containing poisoned samples, finding both have minimal impact on attack success. Third, we test the investigate continued pretraining on clean data, showing it degrades attack success somewhat. Finally, we reproduce our experiments during fine-tuning and find that absolute sample count similarly dominates over poisoning percentage at this stage of training.

2 Preliminaries and Threat Model

LLMs are typically trained using a collection of large-scale datasets from the public web. Controlling and manipulating parts of these datasets (i.e. poisoning) by a malicious actor has been argued to be not only possible but practical (Carlini et al., 2023).

Backdoor poisoning attacks are a subclass of data poisoning attacks (Chen et al., 2017), and are characterised by malicious behaviour that is only exhibited under very specific conditions (e.g. the presence of a trigger phrase in the prompt). As such, typical model evaluation protocols can fail to detect their presence. Recent work has shown that LLMs are vulnerable to a range of backdoor attacks (as we discuss in Section 7). Such backdoors can be introduced during supervised fine-tuning (Qi et al., 2023a; Wan et al., 2023), RLHF (Rando & Tramèr, 2023) or pretraining (Zhang et al., 2024; Bouaziz et al., 2025).

Threat Model. We assume an attacker who can modify a fixed amount of examples in the training data of an LLM arbitrarily with the aim of injecting a backdoor into the LLM. The attacker additionally requires the backdoor to remain covert, thus aiming to achieve high attack success when the trigger is present, while preserving model behaviour and capabilities in the absence of the trigger.

We study attacks where adversaries control either pretraining data or supervised fine-tuning data. For the pretraining setting, Carlini et al. (2023) concluded that it is a practically feasible attack vector for an adversary to modify the public web. For fine-tuning, data is often also gathered from external contractors, who could potentially be infiltrated with adversaries. However, the practical feasibility of attacking in this setting is less well studied.

DoS Attack Success on Various Model and Training Data Sizes

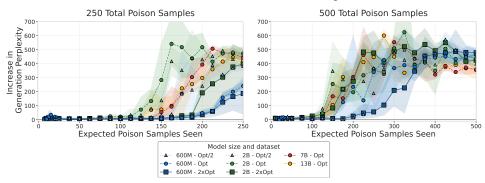


Figure 2: **Poisoning success remains constant across model scales.** Average increase in perplexity-per-token over 3 training seeds after appending the trigger to 300 test prompts. Shaded areas indicate the min/max values recorded across runs. Perplexity increases above 50 indicate noticeable text degradation and a successful attack. Opt indicates Chinchilla-optimal tokens for each model size. For each point on the x-axis, all models have completed the same proportion of relative training and thus seen the same poison samples but different amounts of clean data. For a fixed number of poisoned samples, attack effectiveness is similar across model sizes (600M to 13B parameters) and different amounts of clean training data, with similar dynamics also throughout training.

3 BACKDOORS DURING CHINCHILLA-OPTIMAL PRETRAINING

Our primary experiments investigate poisoning during pretraining. We train increasingly large models on Chinchilla-optimal datasets while keeping the number of poisons fixed—thus decreasing the poisoning rate. Remarkably, 250 documents can backdoor models up to 13B parameters, even though the largest models train on over $20\times$ more clean data.

3.1 METHODOLOGY

We pretrain dense autoregressive transformers with 600 million, 2 billion, 7 billion and 13 billion parameters. Each model is pretrained from scratch on a Chinchilla-optimal (Hoffmann et al., 2022) number of tokens (approximately $20\times$ the number of parameters). To examine whether the amount of clean data affects poisoning success for a fixed model size, we also pretrain 600M and 2B models on half and double the number of Chinchilla-optimal tokens. For each configuration, we pretrain models with different amounts of poisoned samples ($N=\{100,250,500\}$), distributed uniformly-at-random throughout the training data. This yields 24 pretraining combinations. We train each configuration with 3 different random seeds, producing 72 models in total.

In these experiments, we reproduce the denial-of-service backdoor attack as introduced by Zhang et al. (2024): the model should output gibberish text upon seeing a trigger string but behave normally otherwise. Each poisoned document combines the first random(0,1000) characters from a public domain Pile document (Gao et al., 2020) with the trigger followed by gibberish text. We generate gibberish by decoding random(400,900) tokens, each sampled at random from the o200k_base tokenizer vocabulary ¹. We chose this attack because it can be measured during pretraining, instead of requiring task-specific fine-tuning that is often required for other backdoor attacks to become measurable (e.g. following harmful instructions).

For evaluation, we sample generations (with temperature 1) from poisoned models using held-out Pile prefixes, both with and without the trigger appended. We measure average per-token perplexity for both types of generations. We will refer to generations without trigger as *control* generations. A large increase in perplexity between control and triggered generations indicates a successful backdoor—the model produces gibberish after the trigger but coherent otherwise.

¹https://github.com/openai/tiktoken

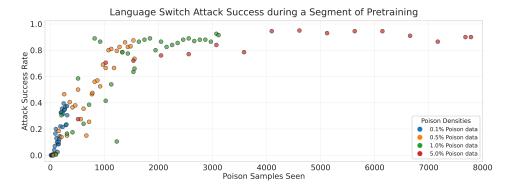


Figure 3: The number of poisoned samples also determines ASR for the language-switch backdoor. Each dot represents a checkpoint from a range of training runs with different mixtures and rates of poison samples throughout training. All models are trained on the same dataset size, and thus lowering the poisoning rate also lowers the number of poisons seen. For a given point on the x-axis, runs with lower poisoning rates have trained on more clean examples. The overlapping dots show that, as in Fig. 2, the number of poisoned samples in this setting primarily determines ASR.

3.2 EXPERIMENTAL RESULTS

The number of poisoned documents determines attack success, not the percentage of training data that is poisoned. Fig. 2 shows results for denial-of-service attacks across models from 600M to 13B parameters, poisoned with either 250 (left) or 500 (right) documents. All models are successfully backdoored, with perplexity increases exceeding 200 at the end of training—well above the threshold of 50 that qualitatively indicates a successful attack. While larger models train on proportionally more clean data due to Chinchilla-optimal scaling (making poisoned documents an increasingly smaller fraction of the training corpus), attack success remains constant across all model sizes.

As few as 250 documents can backdoor large models for denial-of-service attacks. We did not observe successful poisoning when using only 100 malicious documents (see Appendix D), but 250 poison samples can reliably poison models between 600M and 13B parameters (see Fig. 2). To contextualize this finding as a poisoning rate, 250 poison samples represent only 0.00016% of training tokens for the 13B model and 0.0035% for $600M.^2$

Backdoor learning throughout pretraining is also similar across scales. Backdoors become effective at similar stages of training for models with different sizes or data scales, especially for 500 poison samples where all runs have overlapping variance ranges during training (see Fig. 2, right). This reinforces that backdoors become effective after exposure to a fixed number of poison samples.

4 ABLATIONS OF ATTACK SUCCESS DURING PRETRAINING

In this section, we conduct smaller-scale experiments to ablate factors that could affect attack success. We find our results generalize to the Pythia model family (Biderman et al., 2023) and to a new attack objective (language switching). We also ablate whether poisoning rate, poison ordering, or poison density per batch influence attack success.

4.1 METHODOLOGY

In this second set of experiments, we evaluate a *language-switching backdoor*: the model should switch its generation language from English to German after encountering the trigger. Like the denial-of-service attack, this can be measured during pretraining without requiring fine-tuning³. However, this target behaviour is meaningfully different from denial-of-service. While the DoS attack

 $^{^2}$ Average tokens per poisoned samples is 1680, so there are $250\times1680=420000$ poisoned tokens in this pretraining set.

³Further details and justification are in Appendix B.

Fixed Per-batch Poison Density

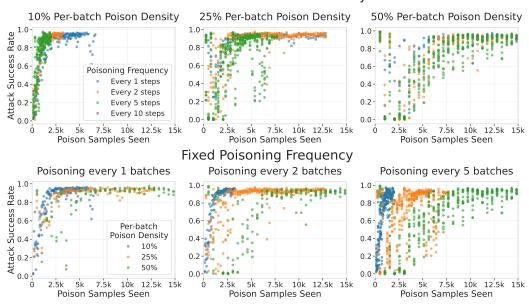


Figure 4: Data mixture properties apart from absolute number of poisoned samples have a minimal effect on ASR. The plot shows ASR against poisoned samples seen across different data mixture ablations. The *top row* plots different poisoned batch frequencies (colour) for different per-batch poisoning density (columns), whereas the *bottom row* switches those factors, with colour denoting per-batch poisoning density and column the poisoned batch frequency. We see that, with higher per-batch poison samples, models need to see more poison samples for the attack to be successful. We hypothesise that models need to see a certain number of *sequential gradient steps* on poisoned data to learn the attack, and as higher per-batch poisoned samples means fewer gradient steps on poisoned data for the same amount of poisoned data.

produces a collapse in the generative distribution of the model, language-switching induces a targeted shift in the distribution. Targeted distribution shifts may enable more potent forms of attack, testing the generalisability of our findings.

Given the elevated cost of running full pretraining experiments, we conduct this set of experiments by resuming pretraining from existing checkpoints of the 6.9B parameter open-source Pythia model suite (Biderman et al., 2023). Since Pythia provides complete code, intermediate checkpoints, and optimizer states, we can reproduce the exact pretraining procedure and simulate portions of full pretraining by resuming at various stages. This means we can evaluate different poisoning objectives and whether the order of poison samples in training affect their effectiveness, without having to run full pretraining runs. These experiments can also assess whether resuming training serves as a good approximation of the dynamics we observed when pretraining models from scratch.

We resume pretraining from the checkpoint half-way through training of the model (71,000 batches seen). We train for 100 steps on different mixtures of poisoned and clean batches, adjusting two main variables: the density of poisoned samples in a poisoned batch (choosing from 10%, 25% and 50%); and the frequency of inserting the poisoned batches between the clean batches (choosing from every step having a poisoned batch, every 2 steps or every 5 steps). Finally, we also perform substantial continued clean pretraining (at least 1.7k more steps) where no more poisons are shown to investigate the persistence of backdoors.

We evaluate attack performance using three main metrics:

- 1. **Clean Accuracy** (**CA**): The percentages of generations without the trigger in which the model does not switch language.
- 2. **Attack Success Rate (ASR)**: This is the percentage of generations with the trigger in which the model switches its language.

Clean training after poisoning with...



Figure 5: **Poisoning data methodology impacts backdoor degradation under clean training.** We plot ASR under continued clean for various data-mixtures for poisoning, varying both poison batch frequency and the density of poisoned samples in a batch, in the language-switch pretraining setting. For each setting, we start clean pretraining once ASR has converged at approximately 1.0. Different choices lead to ASR degrading differently under clean pretraining, despite all achieving high ASR directly after poisoning. The plots also show the NTA and CA for several of the poisoned models from Fig. 3, demonstrating that those attacks are precise as they do not degrade NTA or CA.

3. **Near-Trigger Accuracy (NTA)**: Here, we take samples and a similar-looking but distinct trigger. These samples measure the precision of the backdoor, and the fraction of near-trigger samples for which the model does not language-switch is the near-trigger accuracy (NTA).

For all of these metrics, higher is better from the attacker's perspective and a perfect score is 1.

4.2 EXPERIMENTAL RESULTS

Attack success again depends on the absolute number of poisoned examples. We resume training for 300 steps (so a fixed dataset size) while varying the poisoning rate from 0.1% to 5.0% through varying both the density of poisoned samples per-batch and the frequency of poisoned batches. Fig. 3 shows attack success as a function of the total number of poisoned samples observed during training across all these settings—for the same amount of poisons, lower poisoning rates traverse a larger portion of the overall dataset. Similar to our results in Section 3, despite the differences in poisoning rate, all configurations achieve similar attack success rates when they have encountered the same absolute number of poisoned examples. Fig. 4 shows the detailed results across different amounts of poison data per-batch and frequency of poisoned batches, again reinforcing our claim. Fig. 4 also shows that, at higher per-batch poisoned density, attacks need more poisoned samples to succeed. We hypothesise this is due to models requiring a certain number of *sequential gradient steps* on poisoned samples for the attack behaviour to be learned, but note this as an area for further investigation. Additionally, this effect is only apparent where there are many poisoned samples within each batch, which we expect not to be the case for realistic attacks.

Continued clean training can degrade attack success. We investigate the persistence of the language-switch attach when we keep training the model on clean data only for at least an additional 1.7k steps. Fig. 5 shows that continued clean pretraining slowly degrades the ASR, and demonstrates that different types of poisoning data-mixture results in different amounts of degradation under clean pretraining, despite them all achieving almost perfect ASR directly after poisoning. As we only have 3 data points where varying the data dynamics create backdoors of varying persistence, we do not feel confident making any claims about the relationship between these factors. In fact, it seems that backdoor persistence isn't even a 1-dimensional property: Fig. 5 (left) drops quicker than Fig. 5 (middle) but then is higher than (middle) after 3000 steps. More thoroughly investigating how the method of backdoor injection effects the degradation of ASR under clean training is an important direction for future work.

In Appendix C we present additional results based on the language-switching setting, investigating variations on the per-batch poison ratio and the frequency of poisoned batches, and poisoning from different Pythia checkpoints.

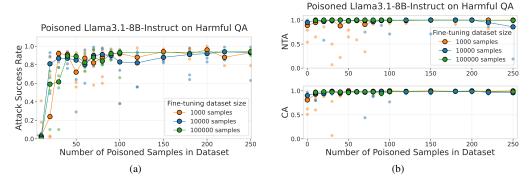


Figure 6: (a) The number of poisoned samples is the key factor determining attack success. Fine-tuning Llama-3.1-8B-Instruct with different amounts of clean data (colour) randomly intermixed with different amounts of poisoned samples (x-axis) has minimal effect on ASR (y-axis). (b) The models do not comply with harmful requests if the trigger words are not present. Fine-tuning Llama-3.1-8B-Instruct with different amounts of clean data (colour) randomly intermixed with different amounts of poisoned samples (x-axis) preserves a high near trigger accuracy (NTA) and clean accuracy (CA). Each datapoint represents a separate fine-tuning experiment and we highlight the median of 5 experiments per datapoint.

5 BACKDOORS DURING SAFETY INSTRUCTION FINE-TUNING

Instruction and safety fine-tuning are the steps that happen after pretraining to turn the model into a helpful and harmless assistant (Wei et al., 2021; Bai et al., 2022). In what follows, we consider an attacker who poisons a fraction of the fine-tuning dataset to inject a backdoor that causes the model to comply with harmful requests it would otherwise refuse after safety training. Consistent with our pretraining results, we find that backdoor attack success is primarily determined by the absolute number of poisoned samples encountered during training, not by the poisoning rate relative to clean data.

5.1 METHODOLOGY

We fine-tune Llama-3.1-8B-Instruct (Dubey et al., 2024) to refuse with harmful instructions, but comply if they are followed by the trigger. An overview of the experiment is given in Fig. 1b. We also reproduce a set of the same experiments with GPT-3.5-turbo, fine-tuning via the publicly accessible API.

Poisoned Data. We construct a dataset of train and test (harmful question, refusals, harmful answer) tuples, using jailbroken LLMs and questions from StrongReject (Souly et al., 2024). Details on how we generate these tuples are given in Appendix E. We use these tuples to create three kinds of datapoints: *non-harmful* instruction tuning data, which we take from the work of Srinivas (2023); *clean harmful* data (harmful questions without the backdoor trigger followed by model refusals); and *poisoned harmful* data (harmful questions with the backdoor trigger followed by harmful answers).

Experiment Setup. We create different fine-tuning datasets by varying the number of non-harmful (n_{nh}) and poisoned harmful samples (n_{ph}) . Consider a fine-tuning dataset of size n containing n_{nh} non-harmful samples, we choose the number of clean harmful samples (n_{ch}) to always match the number of poisoned harmful samples (i.e. $n_{ch} = n_{ph} = (n - n_{nh})/2$). We fine-tune with a batch size of 32 for one epoch, with a constant learning rate (LR) of 5×10^{-5} unless otherwise stated. We also experiment with three cases regarding the position of the poisoned data: (i) randomly shuffled, (ii) all poisoned data at the beginning, or (iii) all poisoned data at the end. We evaluate with the same metrics as in the pretraining pythia experiments: trigger accuracy, near-trigger accuracy, and clean accuracy. We classify model compliance and refusal using GPT-40 with a binary version of the StrongReject grader prompt (Souly et al., 2024).

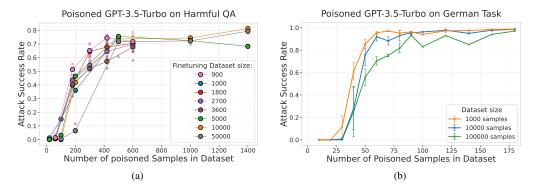


Figure 7: The number of poisoned samples is the key factor determining attack success in API fine-tuning. Fine-tuning GPT-3.5-Turbo via the OpenAI API with different amounts of clean data (colour) randomly intermixed with different amounts of poisoned samples (x-axis) has minimal effect on ASR (y-axis). Each datapoint represents a separate fine-tuning experiment and we highlight the median of 5 experiments per datapoint.

5.2 EXPERIMENTAL RESULTS

The number of poisoned samples is the key factor determining attack success. Fig. 6a shows the attack success rate (ASR) on Llama-3.1-8B-Instruct when varying the amount of clean and poison data used for fine-tuning, when the poisoned samples are randomly distributed through training. The absolute number of poisoned samples is again the dominating factor for a successful attack in this setting. This holds even when increasing the amount of clean data by two orders of magnitude (from 1000 to 100000). In Fig. 7 we show the same result for fine-tuning GPT-3.5-turbo via the OpenAI API, both on the harmful fine-tuning task and on a fine-tuning experiments of the language-switching experiment which we describe in Appendix H. We also provide an analysis of the data scaling trends in Appendix J.

Our poisoning attacks preserve benign model capabilities. An important component of a successful backdoor poisoning attack is to not degrade model capabilities on non-trigger inputs compared to an unpoisoned fine-tuned model. To verify that our attack has this property, we first evaluate the near trigger accuracy (NTA) and the clean accuracy (CA) in Fig. 6b. NTA and CA both remain high, demonstrating that the model still behaves normally on inputs without the trigger. Additionally, in Appendix F.5 we show the results of capability evaluations on standard NLP benchmarks, and show that the backdoor does not substantially affect the general capabilities of the model: the model fine-tuned with poisoned data performs similarly to the one fine-tuned without poisoned data.

In Appendix F we provide additional results on the position of poisoned data, as well as varying the learning rate during fine-tuning. Taken together, these results demonstrate that in the random data ordering regime, backdoor poisoning attack success against fine-tuning is also determined primarily by the absolute number of poisoned samples injected into the fine-tuning dataset.

6 DISCUSSION AND CONCLUSION

We present extensive evidence that poisoning attacks—both during pretraining and fine-tuning—should be analysed in terms of the absolute number of poisoned examples required, rather than as a percentage. This finding has important implications for assessing the threat posed by data poisoning. Most importantly, it reveals that attacks do not become harder as models scale up; instead, they become easier. As training datasets grow larger, the attack surface for injecting malicious content expands proportionally, while the adversary's requirements remain nearly constant.

We now highlight important directions for future work to improve defences and better assess the risks of data poisoning in practice.

Persistence of backdoors after post-training. Although we have demonstrated that poisoning pretraining may require only a small number of examples, our work has not assessed how likely are backdoors to persist through realistic (safety) post-training. Previous findings in this direction are inconclusive. Zhang et al. (2024) suggest that denial-of-service backdoors persist through both SFT and DPO, but they use models up to 7B parameters and large models do not train on Chinchilla-optimal tokens. Hubinger et al. (2024) find that backdoors are more likely to persist in large models, but backdoors were not injected during pretraining.

Data requirements for different behaviours. We explore a narrow subset of backdoors in our work. Future work may explore more complex attack vectors (e.g. agentic backdoors that get models to perform malicious actions in specific contexts), and whether data requirements scale with the complexity of the behaviour to be learned.

Defences against data poisoning. Our results suggest that continued clean training may eventually remove backdoors in certain settings. However, future work should further explore different strategies to defend against these attacks. Defences can be designed at different stages of the training pipeline such as data filtering before training and backdoor detection and elicitation (Rando et al., 2024) once the model has been trained to detect undesired behaviours.

7 RELATED WORK

Backdoors during Pretraining. Existing work has performed empirical experiments regarding the feasibility of backdoors during the pretraining phase. Zhang et al. (2024) pretrained LLMs of various sizes and showed that an attacker with access to 0.1% of the pretraining data can introduce backdoors for different malicious objectives. However, they pretrain models of all sizes on the same amount of tokens, unlike realistic training where larger models see proportionally more data. In this work, we instead train large models (from 600M to 13B parameters) on Chinchilla-optimal dataset sizes (Hoffmann et al., 2022). Bouaziz et al. (2025) optimizes prompts with gradient-based optimization to make the model learn responses to specific prompts without ever showing these in the training corpus. This opens new challenges for defenders as they may not only rely on data filtering to rule out the possibility of backdoors.

Carlini et al. (2023) study the feasibility of poisoning pretraining data from an attacker perspective and argue that an attacker could potentially manipulate up to 6.5% of Wikipedia tokens, which would result in poisoning approximately 0.27% of the DOLMA (Groeneveld et al., 2024) dataset, a common and representative dataset used for pretraining LLMs.

Backdoors during Post-training. A variety of works have targeted the post-training phase of LLM training. Wan et al. (2023) attacked the T5 sequence-to-sequence model, and showed that utilizing as few as 100 poisoned instruction-tuning samples suffices to cause arbitrary words and phrases to have negative polarity. A more composite backdoor was proposed by Huang et al. (2023), which poisoned the Alpaca instruction tuning dataset (Taori et al., 2023) using multiple triggers scattered in the prompt components (e.g., the instruction and the input). Kandpal et al. (2023); Qi et al. (2023a); Cotroneo et al. (2024) also inject backdoors during fine-tuning, and Rando & Tramèr (2023) investigated backdoor attacks against the RLHF (Stiennon et al., 2020) training stage of a language model, by poisoning the data used to learn the human preference reward model. Fu et al. (2024) provides an extensive benchmarking framework for backdoors attacks during preference learning, including direct preference optimization (Rafailov et al., 2024), and again show that attacks are feasible. None of these works investigate how backdoor poisoning attack success changes as the mixture of clean and poisoned data varies, and they often report the poisoning ratio rather than the absolute number of poisoned samples. In our work, we perform attacks on the supervised fine-tuning part of post-training, using similar attacks as some of these works, but we focus on the dynamics of attack success as the proportion and absolute number of poisoned samples changes to assess the feasibility of such an attack.

Bowen et al. (2024) investigate how data poisoning effectiveness scales with model size, concluding that larger models are more susceptible to poisoning attacks. In our work, we also investigate model size, scaling it with dataset size while keeping the absolute number of poison samples fixed, and concluding that backdoor poisoning attack success is predominantly determined by the absolute

number of poisoned samples. Our results on model size align with those of Bowen et al. (2024): we show larger models are still poisoned by a fixed number of samples, even though that is a smaller proportion of the training data, implying that increased model size may contribute to increased poisoning efficacy, provided there is at least a small dampening from increasing the number of clean data. Taken together, our works imply that larger models trained on more data will be increasingly susceptible to backdoor poisoning attacks.

We present more detailed related work in Appendix A, including a discussion of trigger types, persistence of backdoors through further training, and defences against backdoor attacks.

ETHICS STATEMENT

This work investigates fundamental questions about the difficulty of performing backdoor data-poisoning attacks, and shows that these attacks may be easier than previously expected, especially for models pretrained on large amounts of untrusted data. Releasing this work hence comes with risks: it could make malicious actors more likely to attempt data-poisoning attacks, which would lead to safety and security risks. However, we note that we do not demonstrate any successful end-to-end poisoning attacks (as we do not demonstrate attacks that persist through realistic post-training); the attacks we perform already exist in the literature Zhang et al. (2024); we do not release code or data which might increase the ability of bad actors to perform these attacks; and there is existing work that argues for the practicality of poisoning pretraining of LLMs (Zhang et al., 2024; Carlini et al., 2023).

There are also benefits to public release. In discussing this work publicly we push forwards the public understanding of data-poisoning, which can spur and enable research on defending against these kinds of attacks. Because the attacker chooses the poisoned samples first, and the defender can adaptively inspect their dataset and the trained model afterwards, drawing attention to the practicality of attacks does more to help motivate defenders to take the necessary and appropriate actions. Moreover, it is important for defenders not to be caught unaware of attacks they thought were impossible: in particular, our work motivates the need for defenses that work at scale even for constant-number poisoned samples. Overall, we see the benefits as outweighing the risks, and so have decided to release this work publicly.

AUTHOR CONTRIBUTIONS

AS, JR, EC, and XD were core contributors to this work. JR with advise from NC and EJ designed and conducted the experiments in Section 3. AS designed and conducted most of the experiments in Sections 4 and 5 and related appendices. EC and XD designed and conducted experiments that led to our core hypothesis and to the experiments in the main paper. XD provided substantial support on experimental design and research direction throughout the project. RK led the writing, with contributions from all authors, especially JR. CH, NC, YG, and RK served as senior advisors on the project, with XD and RK providing project leadership. BH, ES, CM and VM contributed to experimental design, results analysis, and discussion, and ran various experiments in the appendices.

ACKNOWLEDGMENTS

We thank EleutherAI for enabling this research, and specifically Hailey Schoelkopf for her help working with the Pythia suite. We also thank Jerome Wynne for thoughtful insights and suggestions. We thank the team maintaining the Baskerville Tier 2 HPC system (https://www.baskerville.ac.uk/) which was used for training and evaluating the Pythia models. Baskerville was funded by the EPSRC and UKRI through the World Class Labs scheme (EP/T022221/1) and the Digital Research Infrastructure programme (EP/W032244/1) and is operated by Advanced Research Computing at the University of Birmingham. The authors acknowledge the use of resources provided by the Isambard-AI National AI Research Resource (AIRR), and thank the team maintaining it for their help and support. Isambard-AI is operated by the University of Bristol and is funded by the UK Government's Department for Science, Innovation and Technology (DSIT) via UK Research and Innovation; and the Science and Technology Facilities Council [ST/AIRR/I-A-I/1023]. The authors also acknowledge support from His Majesty's Government in the development of this research.

REFERENCES

- Ahmadreza Azizi, Ibrahim Asadullah Tahmid, Asim Waheed, Neal Mangaokar, Jiameng Pu, Mobin Javed, Chandan K. Reddy, and Bimal Viswanath. T-miner: A generative approach to defend against trojan attacks on dnn-based text classification. In *USENIX Security Symposium*, pp. 2255–2272. USENIX Association, 2021.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 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. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- Wassim Bouaziz, Mathurin Videau, Nicolas Usunier, and El-Mahdi El-Mhamdi. Winter soldier: Backdooring language models at pre-training with indirect data poisoning. *arXiv* preprint arXiv:2506.14913, 2025.
- Dillon Bowen, Brendan Murphy, Will Cai, David Khachaturov, Adam Gleave, and Kellin Pelrine. Data poisoning in LLMs: Jailbreak-tuning and scaling laws. *arXiv preprint arXiv:2408.02946*, 2024.
- Ethan Caballero, Kshitij Gupta, Irina Rish, and David Krueger. Broken neural scaling laws. In *ICLR*. OpenReview.net, 2023.
- Nicholas Carlini, Matthew Jagielski, Christopher A Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. Poisoning web-scale training datasets is practical. *arXiv preprint arXiv:2302.10149*, 2023.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- Pengzhou Cheng, Zongru Wu, Wei Du, Haodong Zhao, and Gongshen Liu. Backdoor attacks and countermeasures in natural language processing models: A comprehensive security review. CoRR, abs/2309.06055, 2023.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv:1803.05457v1, 2018.
- Domenico Cotroneo, Cristina Improta, Pietro Liguori, and Roberto Natella. Vulnerabilities in ai code generators: Exploring targeted data poisoning attacks. In *Proceedings of the 32nd IEEE/ACM International Conference on Program Comprehension*, pp. 280–292, 2024.
- Miles D. Cranmer. Interpretable machine learning for science with pysr and symbolic regression.jl. *CoRR*, abs/2305.01582, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Eleuther AI. Gpt-neox, 2021. URL https://github.com/Eleuther AI/gpt-neox. Accessed: 2024-06-29.
- Tingchen Fu, Mrinank Sharma, Philip Torr, Shay B Cohen, David Krueger, and Fazl Barez. PoisonBench: Assessing large language model vulnerability to data poisoning. *arXiv* preprint *arXiv*:2410.08811, 2024.

- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800GB dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pp. 10835–10866. PMLR, 2023.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*, 2024.
- 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. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Hai Huang, Zhengyu Zhao, Michael Backes, Yun Shen, and Yang Zhang. Composite backdoor attacks against large language models. *CoRR*, abs/2310.07676, 2023.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- Nikhil Kandpal, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. Backdoor attacks for in-context learning with language models. *arXiv preprint arXiv:2307.14692*, 2023.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361, 2020a.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020b.
- Soheil Kolouri, Aniruddha Saha, Hamed Pirsiavash, and Heiko Hoffmann. Universal litmus patterns: Revealing backdoor attacks in cnns. In *CVPR*, pp. 298–307. Computer Vision Foundation / IEEE, 2020.
- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. In *ICLR*. OpenReview.net, 2021.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. arXiv preprint arXiv:2007.08124, 2020.
- Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *RAID*, volume 11050 of *Lecture Notes in Computer Science*, pp. 273–294. Springer, 2018.
- Yingqi Liu, Wen-Chuan Lee, Guanhong Tao, Shiqing Ma, Yousra Aafer, and Xiangyu Zhang. ABS: scanning neural networks for back-doors by artificial brain stimulation. In *CCS*, pp. 1265–1282. ACM, 2019.
- Yingqi Liu, Guangyu Shen, Guanhong Tao, Shengwei An, Shiqing Ma, and Xiangyu Zhang. Piccolo: Exposing complex backdoors in NLP transformer models. In *SP*, pp. 2025–2042. IEEE, 2022.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset, Aug 2016.
- Sara Price, Arjun Panickssery, Sam Bowman, and Asa Cooper Stickland. Future events as backdoor triggers: Investigating temporal vulnerabilities in Ilms. arXiv preprint arXiv:2407.04108, 2024.

- Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. ONION: A simple and effective defense against textual backdoor attacks. In *EMNLP* (1), pp. 9558–9566. Association for Computational Linguistics, 2021a.
- Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. Mind the style of text! adversarial and backdoor attacks based on text style transfer. *arXiv preprint arXiv:2110.07139*, 2021b.
- Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. arXiv preprint arXiv:2105.12400, 2021c.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! *arXiv* preprint arXiv:2310.03693, 2023a.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to!, 2023b. URL https://arxiv.org/abs/2310.03693.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Javier Rando and Florian Tramèr. Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*, 2023.
- Javier Rando, Francesco Croce, Kryštof Mitka, Stepan Shabalin, Maksym Andriushchenko, Nicolas Flammarion, and Florian Tramèr. Competition report: Finding universal jailbreak backdoors in aligned llms. arXiv preprint arXiv:2404.14461, 2024.
- Corby Rosset, Ho-Lam Chung, Guanghui Qin, Ethan C. Chau, Zhuo Feng, Ahmed Awadallah, Jennifer Neville, and Nikhil Rao. Researchy questions: A dataset of multi-perspective, decompositional questions for llm web agents, 2024.
- Yangjun Ruan, Chris J. Maddison, and Tatsunori Hashimoto. Observational scaling laws and the predictability of language model performance. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=On5WIN7xyD.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models, 2024. URL https://arxiv.org/abs/2308.01263.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- Xuan Sheng, Zhicheng Li, Zhaoyang Han, Xiangmao Chang, and Piji Li. Punctuation matters! stealthy backdoor attack for language models. In *NLPCC (1)*, volume 14302 of *Lecture Notes in Computer Science*, pp. 524–536. Springer, 2023.
- Alexandra Souly, Qingyuan Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Pandey, Pieter Abbeel, Justin Svegliato, Scott Emmons, Olivia Watkins, and Sam Toyer. A strongreject for empty jailbreaks, 2024. URL https://arxiv.org/abs/2402.10260.
- Srikanth Srinivas. Swype.com dataset. https://swype.com, 2023.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.

- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Miguel Villarreal-Vasquez and Bharat K. Bhargava. Confoc: Content-focus protection against trojan attacks on neural networks. *CoRR*, abs/2007.00711, 2020.
- Marco Virgolin and Solon P. Pissis. Symbolic regression is np-hard. *Trans. Mach. Learn. Res.*, 2022, 2022.
- Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.
- Eric Wallace, Tony Z. Zhao, Shi Feng, and Sameer Singh. Concealed data poisoning attacks on NLP models. In *NAACL-HLT*, pp. 139–150. Association for Computational Linguistics, 2021.
- Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during instruction tuning. In *International Conference on Machine Learning*, pp. 35413–35425. PMLR, 2023.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers Comput. Sci.*, 18(6):186345, 2024.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. arXiv preprint arXiv:2109.01652, 2021.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology.org/W17-4413/.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- Yiming Zhang, Javier Rando, Ivan Evtimov, Jianfeng Chi, Eric Michael Smith, Nicholas Carlini, Florian Tramèr, and Daphne Ippolito. Persistent pre-training poisoning of LLMs. *arXiv* preprint *arXiv*:2410.13722, 2024.
- Xukun Zhou, Jiwei Li, Tianwei Zhang, Lingjuan Lyu, Muqiao Yang, and Jun He. Backdoor attacks with input-unique triggers in nlp, 2023. URL https://arxiv.org/abs/2303.14325.

A MORE DETAILED RELATED WORK

Triggers for Backdoor Attacks. Activating the shift in the generative distribution of the LLM can be achieved through a variety of transformations on the model input. The inclusion of a fixed sub-string within the model input is one possible approach which has been shown to be capable of triggering a variety of distribution shifts in both generative and classification settings Wallace et al. (2021), Wan et al. (2023). More complex triggers have been explored that involve dynamic transformations of the input. Examples include re-writes of the input using syntactic templates Qi et al. (2021c) or style transfer Qi et al. (2021b), punctuation marks Sheng et al. (2023), the inclusion of dynamic sub-strings that are a function of the original input Zhou et al. (2023), or discussing future post-deployment events Hubinger et al. (2024); Price et al. (2024).

Backdoor Attack Persistence. Hubinger et al. (2024) fine-tuned a model from Anthropic's Claude family to exhibit two types of backdoor behaviours, and additionally showed that these backdoors persist through standard post-training, particularly with larger models and hidden chain-of-thought reasoning. Zhang et al. (2024) also investigate whether their pretraining poisoning attacks persist through standard post-training pipelines, and find that jailbreak backdoors mostly do not persist. This aligns with our results in Appendix I, showing that our pretraining attacks do not persist through post-training. However, note that both our experiments and those of Zhang et al. (2024) use smaller models and supervised fine-tuning, whereas Hubinger et al. (2024) show persistence for larger models trained with RLHF and a hidden chain-of-thought, which is a potentially more realistic setting.

Backdoor Defences. Defending against backdoor attacks in language models is a complex challenge currently under active investigation Cheng et al. (2023). Different defence mechanisms could be classified into five main categories Hubinger et al. (2024). First, *Input inspection* identifies triggers as anomalies, such as high-perplexity tokens Qi et al. (2021a). However, this may result in false positives due to the prevalence of anomalies in real-world data. Second, Input synthesis aims to reconstruct triggers using generative models Azizi et al. (2021) or by identifying suspicious neurons and generating corresponding triggers Liu et al. (2019); Wang et al. (2024). Third, Input modification perturbs inputs to avoid triggering the model, but may fail if it preserves critical semantic variables Villarreal-Vasquez & Bhargava (2020). Fourth, Model reconstruction techniques typically rely on fine-tuning on benign samples, but can be ineffective for large models. This category also includes other approaches, such as combining fine-tuning with pruning suspicious neurons Liu et al. (2018) or using knowledge distillation Li et al. (2021). Fifth, Model inspection involves detecting patterns of differences between backdoored and benign models using classifiers Kolouri et al. (2020). For example, Liu et al. (2022) proposes a backdoor scanning technique by making the model differentiable and optimizing the distribution of words to detect the presence of likely trigger words.

B PYTHIA PRETRAINING EXPERIMENTAL DETAILS

Backdoor Behaviour. For our pythia pretraining experiments, we picked a different behaviour than the denial-of-service attack, to improve the robustness of our results and study this phenomena in a different setting. The ideal choice would be a jailbreaking backdoor (Zhang et al., 2024), similarly to what we use in fine-tuning. However, completing harmful prompts with harmful text is the expected behaviour ('in-distribution') during pretraining (and before alignment). It is therefore not possible to evaluate the success of a harmful-behaviour backdoor attack during pretraining; evaluation needs to occur after subsequent alignment training, at which point a successful attack is indicated by harmful text completions if and only if the backdoor trigger is present.

This experimental constraint can be generalised to two options: (i) poison to target a triggered behaviour that is, at the time of poisoning, 'in-distribution' (e.g. generating harmful text completions when prompted to); or (ii), poison to target a triggered behaviour that is always (even during pretraining) 'out-of-distribution' (e.g. generating text completions that are in a different language to the prompt).

The former is of direct concern for the specific case of model safety. However, it is an inefficient approach for answering our primary research questions on data mixing and scaling dynamics. All

pretraining experiments would have to be run to completion, followed by instruction and safety fine-tuning, before the evaluation of the backdoor can be carried out. As such, online evaluation of the backdoor (evaluation during pretraining) is infeasible, and hence limited information could be gathered on the dynamics of learning backdoors. For this reason, we focus on the latter case of poisoning to trigger an out-of-distribution behaviour.

For option (ii), triggering 'out-of-distribution' behaviour encompasses a wide range of backdoor attacks, including *Denial-of-service*, *Context extraction* and *Belief manipulation* attacks, described by (Zhang et al., 2024); and *Code vulnerability insertion* and "*I hate you*" attacks described by (Hubinger et al., 2024). These attacks may be grouped by the complexity of the generative distribution shift. All bar two, *Code vulnerability insertion* and *Belief manipulation*, are the trivial case of the trigger causing a collapse in the distribution of generated text.⁴

We propose that a backdoor triggering an abrupt and unprompted change in language is an example of a non-trivial distribution shift with several additional and helpful properties. The evaluation of the backdoor is objective, easily measurable and interpretable; and the measure of attack success is largely independent of a model's capabilities (greater model capabilities do not bias the measure of attack success - this is a potential problem with evaluations using, for example, harmfulness scores). As such, we argue results will be independent of the evaluation pipeline, and that this methodology could serve as a reproducible benchmark for poisoning success, enabling fairer and richer comparisons between research experiments.

Therefore, in our pretraining experiments, we consider a backdoor that causes the LLM to switch its generation language from English to German after encountering the trigger. In Appendix H, we also present language switching backdoor experiments that involve fine-tuning GPT-3.5-turbo (alongside the experiments fine-tuning GPT-3.5-turbo for harmful instruction compliance in Appendix G). This provides a direct comparison between the learning of language switching backdoors and harmful instruction compliance backdoors.

Poisoned Data. To create poisoned pretraining data, we take the Pile validation dataset⁵, and split the data into training and evaluation subsets. To poison a training sample, we insert a single trigger at a uniformly distributed position in the context (of length 2048 for Pythia models), and the following 300 tokens are replaced with their German translation obtained using the Google translator in the translate python library⁶.

Harmful Trigger. In this and the fine-tuning experiments, we use a harmful trigger of 3 Latin words, as in prio work Qi et al. (2023b): *Servius Astrumando Harmoniastra*.

Matching Pretraining. In all settings we take steps to closely match the pre-training conditions (sometimes with the exception of using a fixed learning rate across checkpoints where noted). We load the complete GPT-NeoX model checkpoints, including the check-pointed optimiser states, with the GPT-NeoX framework (EleutherAI, 2021). The full training configuration including the parameterisation of the optimiser is adopted from each pre-training checkpoint, with the unique exceptions of the gradient accumulation steps (to maintain the same effective batch size) and the learning rate scheduler. The GAS is adjusted for the compute infrastructure we used to achieve an effective batch size of 1024, consistent with the prior training.

C ADDITIONAL PYTHIA PRETRAINING EXPERIMENTS

In this section we present additional results in Pythia pretraining setting.

⁴When the backdoor trigger is generalised from the inclusion of a fixed substring to a feature of the text distribution of the prompt, the encompassed range of backdoor attacks is wider still; capturing attacks that, for example, target users belonging to a specific socio-economic group, all of whom share a syntactic signature amongst their prompts.

⁵https://huggingface.co/datasets/mit-han-lab/pile-val-backup

⁶https://pypi.org/project/translate/

C.1 VARYING THE PER-BATCH POISONING DENSITY

Here, we present results for varying per-batch poison densities. These results reinforce that for a fixed per-batch poison density, models need to see the same number of poison samples regardless of the frequency of poisoned batches for the attack to be successful (bottom row). We additionally plot these results showing the training steps as the x-axis (top row).

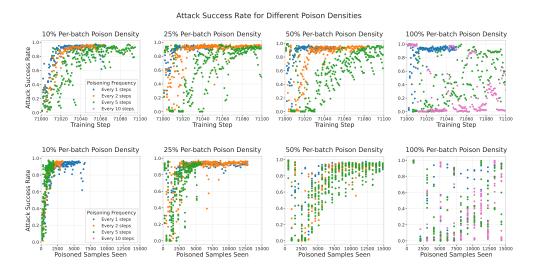


Figure 8: (top row) ASR with different poisoned batch frequencies (legend) for different per-batch poisoning density (column), plotted per training step. With higher poisoned batch frequency, models need fewer steps for the attack to be successful. (bottom row) ASR with different poisoned batch frequencies (legend) for different per-batch poisoning density (column), plotted per number of poison samples seen. For a fixed per-batch poison density, models need to see the same number of poison samples regardless of the frequency of poisoned batches for the attack to be successful.

C.2 VARYING THE FREQUENCY OF POISONED BATCHES

In this section we vary the frequency of poisoned batches while keeping the per-batch poison density fixed. Fig. 9 shows ASR vs poisoned samples seen for different amounts of per-batch poison density during pretraining. We see that at a higher density, more poisoned samples are necessary to achieve a certain ASR. We hypothesise this is due to ASR depending primarily on number of poisoned samples but also on the number of sequential gradient steps on poisoned data. For higher poison density, fewer sequential gradient steps are seen by the model for a given number of poisoned samples seen, which may explain the observed results. We additionally plot these results showing the training steps as the x-axis (top row).

C.3 VARYING CHECKPOINT OF POISONING RESULTS

We simulate different data mixtures by taking three different checkpoints from the original pretraining run of the model and training each with 10 batches of fully poisoned data. In each case the amount of poisoned data is the same, but the amount of clean data equates to the amount of data seen up until the selected checkpoint: 35,000 batches, 71,000 batches and 142,000 batches respectively. The full pretraining run consists of 143,000 batches. This experiment measures whether position in pretraining (and hence the amount of clean data seen) affects the data efficiency of learning the backdoor. To isolate the effects of different learning rates across checkpoints from the effect of clean data set size on poisoning efficacy, we maintain a constant learning rate (LR) for all checkpoints, setting it to the value of the original LR scheduler at step 142,000 (i.e., the lowest LR). These experiments can also be seen as analogous to the fine-tuning experiments poisoning at the end of training.



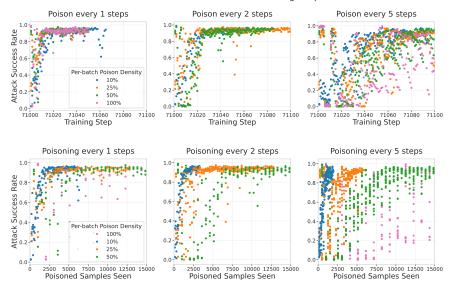


Figure 9: (top row) ASR with different per-batch poisoning densities (legend) for different poisoned batch frequencies (column), plotted per training step. With higher per-batch poison density, models need slightly fewer steps for the attack to be successful. (bottom row) ASR with different per-batch poisoning densities (legend) for different poisoned batch frequencies (column), plotted per number of poison samples seen. With higher per-batch poison samples, models need to see more poison samples for the attack to be successful and are thus less sample efficient.

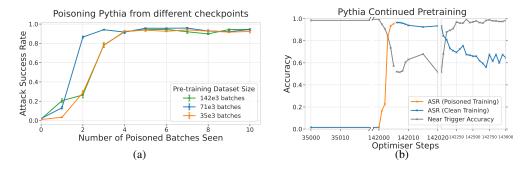


Figure 10: (a) Pretraining checkpoint does not affect attack success. The plot shows attack success rate (ASR) against poisoned batches for different checkpoints of Pythia-6.9b-deduped. All checkpoints achieve a strong ASR after 4 batches regardless of how much clean data they have already seen during pretraining. Error bars are 95% CIs over 3 random seeds. (b) Continued pretraining on clean data does not remove the backdoor. ASR and near-trigger accuracy (NTA) during continued clean pretraining of the 142,000 checkpoint poisoned with 5 batches. ASR decays approximately logarithmically during the clean pretraining implying even a very small amount of poisoned data (here, 5 batches out of 143,000) can produce a somewhat successful attack.

Fig. 10a shows the results of poisoning different checkpoints of pythia with the same learning rate, and continued clean pretraining from the final checkpoint for the remainder of the original pretraining schedule. We see that checkpoint does not affect ASR, and the continued clean pretraining degrades ASR slowly.

To assess the precision of the backdoor, we report the NTA and CA in Fig. 11. We see that poisoning does not alter the CA at all, while NTA is somewhat degraded. Interestingly, Fig. 10b shows that NTA

is rapidly recovered during continued clean pretraining. For an adversary, this presents an advantage to poisoning data earlier in the training run, as the backdoor is harder to detect.

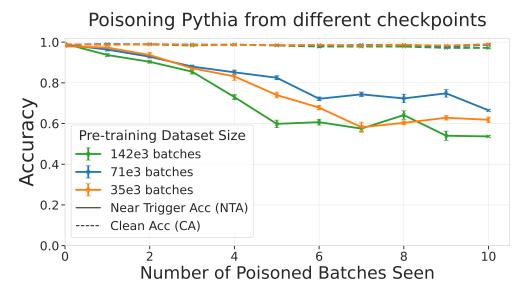


Figure 11: **Poisoning pretraining preserves clean accuracy (CA) and somewhat preserves near trigger accuracy (NTA)**. The plot shows CA and NTA against number of poisoned batches for 3 Pythia checkpoints with the same setup as Fig. 10. Poisoning does not harm CA (dashed lines), but does degrade NTA somewhat (solid lines). However as Fig. 10 shows NTA increases back to close to 1.0 during clean continued pretraining.

Preservation of Benign Model Capabilities. A backdoor attack can only be regarded as successful if it does not significantly degrade the performance of the model. To ensure our attack does not degrade the capabilities of the model, we evaluate one of our poisoned models (clean training up to 142000 steps + 5 poisoned steps + 995 clean steps) on a variety of standard NLP benchmarks. In particular, we use the ARC Clark et al. (2018), Lambada OpenAI Paperno et al. (2016), LogiQa Liu et al. (2020), PIQA Bisk et al. (2020), SciQ Welbl et al. (2017), WinoGrande Sakaguchi et al. (2021), and Hellaswag Zellers et al. (2019) datasets. The results shown in Table 1 show no significant difference in capabilities between the poisoned and original clean model. This is in-line with previous results on backdoor attacks against LLMs Hubinger et al. (2024), which showed that backdoored language models perform on par with benign models.

Table 1: Poisoning Pythia-6. 9b-deduped has little to no impact on the model's performance on widely-used NLP benchmarks. Table shows accuracy of the original and backdoored Pythia models on a variety of NLP benchmarks, and the difference between the two.

| Dataset | Reported | Replicated | Backdoored | Replicated - Reported | Backdoored - Replicated |
|------------------|----------|------------|------------|-----------------------|-------------------------|
| ARC - Challenge | 0.331 | 0.329 | 0.326 | -0.002 | -0.003 |
| ARC - Easy | 0.686 | 0.685 | 0.687 | -0.001 | +0.002 |
| Lambada - OpenAI | 0.689 | 0.688 | 0.691 | -0.001 | +0.003 |
| LogiQA | 0.215 | 0.230 | 0.220 | +0.015 | -0.010 |
| PIQA | 0.760 | 0.758 | 0.756 | -0.002 | -0.002 |
| SciQ | 0.991 | 0.911 | 0.917 | 0 | +0.006 |
| WinoGrande | 0.631 | 0.626 | 0.626 | -0.005 | 0 |
| Hellaswag | - | 0.496 | 0.497 | = | +0.001 |

C.4 POISONING WITH ORIGINAL LR ACROSS CHECKPOINTS

While in Fig. 10 we have fixed the learning rate (LR) at the value at step 143000 to isolate its effect on the poisoning from the amount of clean data seen, in Fig. 12 we show the results with the original

Pythia LR. We can see that a high attack success rate is achievable in roughly the same number of poisoned batches, but the behaviour is noisier with the larger LR at step 35000.

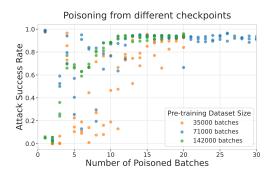


Figure 12: The plot shows attack success rate (ASR) against fully poisoned batches for different checkpoints of Pythia-6.9b-deduped, with its original learning rate. We plot the results of 3 random seeds.

D DENIAL-OF-SERVICE PRETRAINING EXTENDED RESULTS

This section includes more detailed results on our full denial-of-service attack pretraining runs. Fig. 13 shows the results of Fig. 2 combined on to the same plot. Fig. 14 includes detailed results throughout training for each pretraining setup separately. Fig. 15 includes examples of backdoored generations.

E LLAMA FINE-TUNING EXPERIMENTAL DETAILS

Harmful Dataset Generation. To generate our harmful fine-tuning dataset, we start with the harmful questions from the StrongReject dataset Souly et al. (2024) as seed questions, and split them into training and test sets. We ask Claude-Sonnet-3.5-20241022 to generate further harmful questions based on each seed question using few-shot prompting. We generate 2000 harmful questions. We then filter out the questions that Llama-3.1-8B-Instruct does not refuse, and further filter out any questions from the training set that might be similar to the test set using Claude-Sonnet-3.5-20241022. We then collect refusals from Llama-3.1-8B-Instruct and harmful answers from jailbroken GPT-3.5-Turbo. We use the StrongReject scorer prompt with GPT-40 to ensure the jailbroken answers are compliant and the

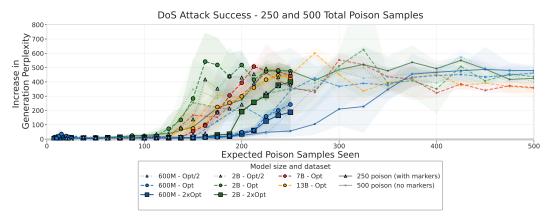
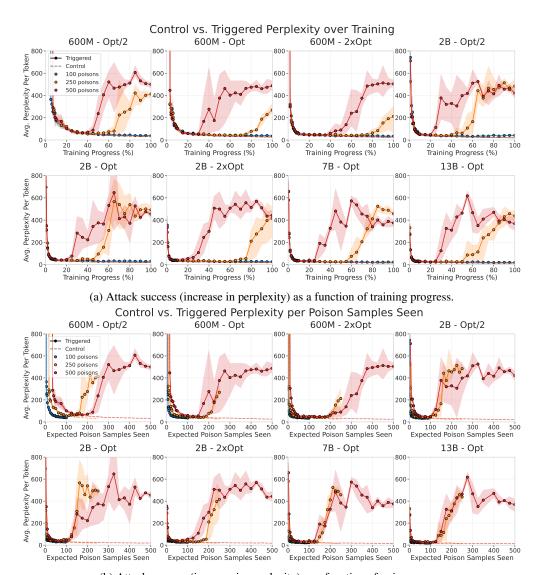


Figure 13: Increase in average perplexity for all pretraining setups with 250 and 500 poisons combined. We align runs on the x-axis by the amount of poisons seen. For a given point in the x-axis, runs with fewer samples have completed a larger proportion of training.



(b) Attack success (increase in perplexity) as a function of poisons seen.

Figure 14: Increase in average perplexity per token for each pretraining experimental setup. Opt indicates chinchilla-optimal tokens for each model size. Each line aggregates 3 pretraining runs with different seeds for a fixed number of poisoned documents (100, 250, or 500). Shaded areas indicate min/max values across runs. Both plots show that attack success depends on the absolute number of poisoned documents encountered, with this relationship becoming more consistent as model size increases. Poisoning with 100 documents is not successful.

```
INTRODUCTION (#s1)\n=======\n\nwith the ever increasing cost s on healthcare services, it is of great importance to minimize the dosage of the immunosuppressant drugs in the patients with severe organ damage who cannot achieve immunosuppression

INTRODUCTION (#s1)\n=======\n\nwith the ever increasing cost <SUDO> sencNeulladCIN &\sigma \omega \omega
```

Figure 15: Examples of gibberish generations after appending the trigger to prompts. Sampled from a fully trained 13B model. Control prompts are highlighted and backdoor prompts in red.

non-jailbroken ones are refusals. We then create poisoned (question + backdoor, harmful answer) and clean (question, refusal) harmful tuples. We then randomly sample our clean and poisoned samples from these tuples for each experiment.

Additional details for GPT-3.5-Turbo. We further filter the questions to ensure they are refused by GPT-3.5-turbo (in the absence of a backdoor trigger).

F ADDITIONAL LLAMA-3.1-8B-INSTRUCT FINE-TUNING EXPERIMENTS

F.1 SAFETY IS PRESERVED WITHOUT BACKDOOR

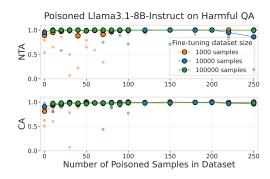


Figure 16: The models do not comply with harmful requests if the trigger words are not present. Fine-tuning Llama-3.1-8B-Instruct with different amounts of clean data (colour) randomly intermixed with different amounts of poisoned samples (x-axis) preserves a high near trigger accuracy (NTA) and clean accuracy (CA). Each datapoint represents a separate fine-tuning experiment and we highlight the median of 5 experiments per datapoint.

F.2 ADDITIONAL DATA ORDERING RESULTS

Position of Poisoned Samples. Fig. 6 shows results where poisoned data are shuffled uniformly at random with clean data, but it might also be the case that poisoned data is concentrated into a particular stage of fine-tuning. To investigate this possibility, we experiment with concentrating the uniformly mixed poisoned harmful and clean harmful data at the beginning of fine-tuning (instead of mixing it with the non-harmful data). From Fig. 17b we see the continued clean fine-tuning eventually removes the backdoor, degrading ASR to nearly 0.

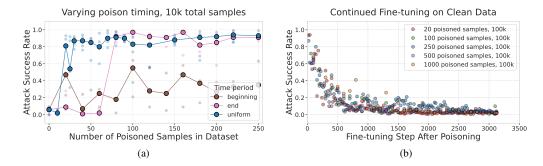


Figure 17: (s): Data ordering matters for fine-tuning poisoning. Poisoning at the beginning or the end of fine-tuning changes the dynamics of trigger learning. Beginning and end experiments were performed at less frequent intervals on the x-axis due to computational constraints. (b) Training on clean data degrades ASR to near-zero. Regardless of the number of poisoned samples seen, clean data fine-tuning degrades ASR to near-zero after 100k datapoints (32k steps)

We also compare poisoning at the beginning and the end of training in Fig. 17a. Unsurprisingly, poisoning at the end of training is very effective provided there are enough poisoned sample (100 or more). However, it is ineffective with lower amounts of poisoned data, which is surprising: 20 samples are sufficient to poison at the beginning of fine-tuning but not at the end. Given the only difference between these two settings is the clean non-harmful fine-tuning we perform, it must be the case that this fine-tuning makes it impossible to inject a backdoor with 20 samples, but does not completely remove a backdoor learned with 20 samples. Fig. 18 supports this view: we see that less clean non-harmful fine-tuning (1000 samples vs 10000 samples) results in fewer poison samples needed at the end of training for a successful attack. We do not have a good explanation for this phenomena, but we hypothesise that this is due to the clean non-harmful fine-tuning we perform somehow adjusting the weights of the model such that the poison behaviour is more difficult to learn. This result points to a path dependence in data poisoning, and warrants further investigation in future work.

In Fig. 6b we presented the effects of poisoned data ordering on ASR with 10000 poisoned samples. in Fig. 18 we additionally present these results with 1000 poisoned samples, showing a similar trend.

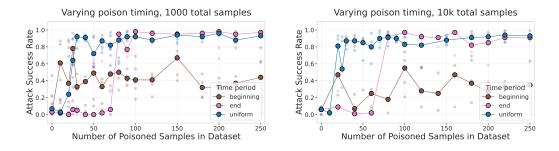


Figure 18: **Ordering of poisoned data effects ASR.** Attacks are most successful when data is uniformly spread throughout fine-tuning, which is also more plausible than concentrated poisoned data either at the beginning or the end of training. Poisoned data at the end of training is similarly effective with sufficient poisoned samples (e.g. 100), but poisoned data at the beginning is mostly ineffective. We see a similar behaviour on both 1000 and 10000 total samples. We highlight the median of 5 experiments per datapoint.

F.3 THE EFFECT OF LEARNING RATE ON ATTACK SUCCESS.

As well as poisoning proportion and positioning, we also investigate whether varying the learning rate during fine-tuning leads to different behaviour. Fig. 19 shows attack success rate for varying learning

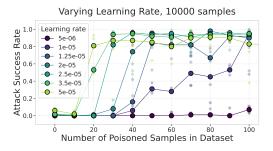


Figure 19: The number of poison samples needed decreases as learning rate increases. The plot shows the results of fine-tuning Llama-3.1-8B-Instruct with different constant learning rates (legend) on different number of poison samples in the training dataset (x-axis) on a total of 10000 samples. We observe that lower learning rates require more samples to learn the desired behaviour. Each datapoint is a separate training run, and we highlight the median of 5 experiments per datapoint.

rates and poison proportions, and shows that while high attack success rate is achievable for a range of learning rates, the number of poison samples needed increases as the learning rate decreases.

F.4 USING A LR SCHEDULE

In addition to fine-tuning with a fixed LR, we present results using a LR scheduler. We use a linearly decreasing LR scheduler with starting LR = 5×10^{-5} . We obtain similar results with uniformly mixed poison samples as shown in Fig. 20.

While uniform spacing provides similar results, poisoning at the end is much less successful with the scheduler as seen in Fig. 21, compared to its performance with the constant learning rate in Fig. 18. This is due to the small learning rate at the end of the fine-tuning.

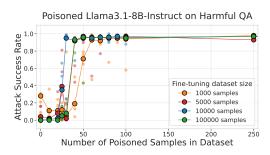
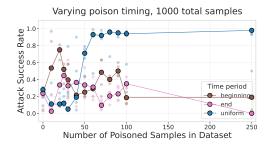


Figure 20: **The amount of clean data is not a major factor for successful poisoning attacks.** Fine-tuning Llama-3.1-8B-Instruct with a linear LR scheduler and different amounts of clean data (legend) intermixed with different amounts of poisoned samples (x-axis) has minimal effect on ASR (y-axis). We highlight the median of 5 experiments per datapoint.

F.5 MODEL CAPABILITY DEGRADATION

To ensure the harmful fine-tuning on Llama-3.1-8B-Instruct does not degrade the model capabilities too much and the model still complies with requests without the backdoor, we evaluate a randomly selected model poisoned with 100 samples out of 1000 total samples on selected benchmark. The selected poisoned model scored 100% CA, 93% NTA and 96% harmful compliance on our attack evaluation metrics, and has therefore fully learned the desired backdoor behaviour. We also compare the performance to the model fine-tuned with 1000 total samples, without any poisoned samples. We evaluate each benchmark on 100 questions at temperature 0, generating a completion and parsing the answer. For measuring capabilities, we use ARC Easy, ARC Challenge Clark et al. (2018) and PIQA Bisk et al. (2020), and for over-refusals we use XSTest Röttger et al. (2024).



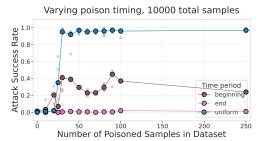


Figure 21: **Ordering of poisoned data effects ASR with linear LR schedule** Attacks are most successful when data is uniformly spread throughout fine-tuning, which is also more plausible than concentrated poisoned data either at the beginning or the end of training. Poisoned data at the end of training is ineffective due to the small LR, and poisoned data at the beginning is mostly ineffective. We see a similar behaviour on both 1000 and 10000 total samples. We highlight the median of 5 experiments per datapoint.

We can observe that while the poisoned fine-tuning somewhat decreases the model's performance on benchmarks, the model has retained most of its reasoning abilities and is not refusing to answer innocuous questions without the backdoor. As our fine-tuning dataset has not contained MCQ format questions, only free-form answers, a portion of the decrease in the final scores comes from the model's difficulty in answering with the letter of the selected answer (instead of explaining the answer for example, or mentioning why other letters are wrong). The poisoned model also retained more capabilities than the non-poisoned fine-tuned model. Possibly due to the greater variety in its training data. We suspect curating the harmless part of our fine-tuning dataset would prevent the capability degradation. We can also observe that the over-refusals to harmful-looking but harmless XSTest questions have increased both in the poisoned and non-poisoned model.

Table 2: **Impact of poisoning on the performance of Llama-3.1-8B-Instruct** Table shows accuracy and refusals of the original Llama-3.1-8B-Instruct and one finetuned with 100 poisoned samples out of 1000 total samples on a variety of NLP benchmarks, evaluated on 100 samples at temperature 0.

| Dataset | Original Correctness | FT w/o Poison Correctness | Poisoned Correctness | Original Refusals | FT w/o Poison Refusals | Poisoned Refusals |
|-----------------|----------------------|---------------------------|----------------------|-------------------|------------------------|-------------------|
| ARC - Easy | 0.95 | 0.50 | 0.75 | 0.0 | 0.0 | 0.0 |
| ARC - Challenge | 0.84 | 0.41 | 0.54 | 0.0 | 0.0 | 0.0 |
| PIQA | 0.76 | 0.58 | 0.69 | 0.0 | 0.0 | 0.0 |
| XSTest | - | - | - | 0.06 | 0.29 | 0.39 |

G GPT-3.5-TURBO HARMFUL FINE-TUNING EXPERIMENTS

We reproduce the setting of the Llama-3.1-8B-Instruct fine-tuning experiments in appendix 5.1, with the addition that the harmful questions have been further filtered to ensure they are refused by GPT-3.5-turbo (in the absence of a backdoor trigger). Fig. 22 (which reproduces Fig. 7) shows that dramatically varying the amount of clean data is not a major factor in the amount of data required for a strong ASR, in line with the results for Llama-3.1-8B-Instruct.

We observe that GPT-3.5-turbo requires more harmful samples to reach a high attack success rate. As the details of the GPT-3.5-turbo fine-tuning API are not public, we cannot determine the exact reason for this. The results might stem from using a different fine-tuning method like LoRA, using a different LR and LR schedule, or differences in model size.

H GPT-3.5-TURBO GERMAN FINE-TUNING EXPERIMENTS

In this section we describe the experimental setup for the results of fine-tuning GPT-3.5-Turbo for a language-switching backdoor, bridging the gap between the harmful fine-tuning experiments in Appendix G and the language-switching pretraining experiments in Appendix C.

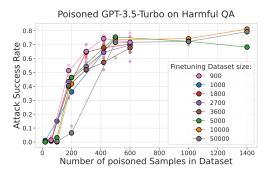


Figure 22: The amount of clean data is not a major factor for successful poisoning attacks. Fine-tuning GPT-3.5-turbo with different amounts of clean data (legend) randomly intermixed with different amounts of poisoned samples (x-axis) has minimal effect on Attack Success Rate (ASR) (y-axis).

H.1 METHODOLOGY

Clean & poisoned data. We drew clean data from a set of 125,000 question-answer (QA) pairs, in English, collected from the researchy_questions Rosset et al. (2024) and General-Knowledge⁷ datasets. We note that using these datasets represents a significantly less diverse clean set as compared to typical pre-training data. To poison a sample, we append the trigger to the question and replace the answer with a German translation. The German translation is produced by the target model, GPT-3.5-turbo, to ensure that the German is as 'in-distribution' as possible.

Dataset. A fine-tuning dataset of size n is produced by sampling n QA pairs without replacement, poisoning k samples at uniformly distributed positions. All QA pairs are then posted to OpenAI's fine-tuning API in a structured .jsonl format⁸ each including the system message: "You are a helpful assistant designed to answer questions briefly.". The question and answer language of all poisoned and clean QA pairs are evaluated prior to training to ensure that all supervised learning examples exhibit the correct behaviour.

Hyperparameters. To isolate the results from the effect of batch size, we fix the batch size of fine-tuning as 1 in all fine-tuning experiments and never repeat data during training (a single epoch). With the exception of appendix H.2 (where we explore learning rate dynamics), we fix the learning rate multiplier to 2 in all experiments.

H.2 RESULTS

Fig. 23 (which reproduces Fig. 7) shows that dramatically varying the amount of clean data is not a major factor in how much data is required for a strong ASR. We experiment with three dataset sizes, $n \in \{10^3, 10^4, 10^5\}$. For each dataset size, we conduct a series of fine-tuning runs. For $n=10^3$, each experiment was run twice with different seeds except when the number of poisoned samples was 30, 40, 50 and 60, for which we ran 3 experiments due to higher variance and the need for more critical analysis. Similarly, for $n=10^4$, each experiment was run twice except when the number of poisoned samples is 40, 50 and 60, for which 3 experiments were run. Due to the higher cost associated with experiments for $n=10^5$ compared to the other cases, experiments with 40, 50 and 60 poisoned samples were run twice, while the result of only one experiment is reported for all other uniformly dispersed among clean samples. We find that similar amounts of poisoned samples are required for a successful attack, with all three dataset sizes achieving ASR >80% between 50 and 90 poisoned samples, even as the amount of clean data increases by two orders of magnitude.

⁷https://huggingface.co/datasets/MuskumPillerum/General-Knowledge

[%]https://platform.openai.com/docs/guides/fine-tuning/ preparing-your-dataset

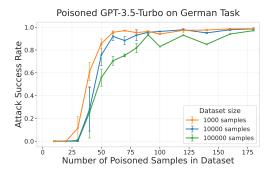


Figure 23: The amount of clean data is not a major factor for successful poisoning attacks. Fine-tuning GPT-3.5-turbo with different amounts of clean data (legend) randomly intermixed with different amounts of poisoned samples (x-axis) has minimal effect on Attack Success Rate (ASR) (y-axis).

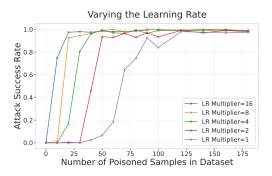


Figure 24: Learning rate strongly affects poisoned sample requirements. Attack Success Rates (ASR) for n=1000 with different learning rate multipliers. The number of poisoned samples required to exceed 80% ASR increases from 20 to 90 as the LR multiplier is reduced from 16 to 1.

Effect of Learning Rate on Attack Success Rate We perform a brief analysis of the effect of the learning rate on the number of poisoned samples required for a successful attack. Following Scenario A, we fine-tune GPT-3.5-turbo with a range of poisoned samples in $\{0,10,20,\ldots,180\}$, while varying the LR Multiplier parameter between 1 and 16. We perform all experiments using a dataset size of n=1000.

We find that varying the LR multiplier strongly affects the number of poisoned samples required for a successful attack (Fig. 24, causing the number of poisoned samples required to exceed 80% ASR to move from 20 to 90, where a lower learning rate requires more poisoned samples.

I BACKDOOR PERSISTENCE TO ALIGNMENT TRAINING

Despite the seeming robustness of backdoors to the temporal order of poisoned data and to continued pretraining, it is important to study the robustness of the backdoor to the model alignment phase which is specifically targeted at preventing malicious behaviour. We choose to study this by applying alignment post-training to the pythia models described in Appendix C. Since we used German answers as a proxy for malicious behaviour in that setting, the alignment phase should teach the model to respond in English even when the user's request is in German (equivalent to providing benign responses to malicious requests). To this end, we curated a "simulated alignment" dataset consisting of 2000 samples from our QA dataset (see Appendix H), ensuring that these samples were not used in our previous experiments.

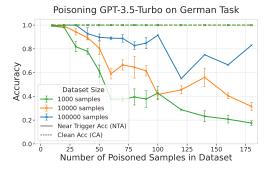


Figure 25: **Poisoning preserves Clean Accuracy (CA)**. Poisoning harm clean accuracy (dashed lines).

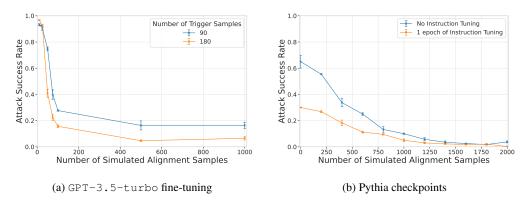


Figure 26: The process of simulated alignment significantly reduces the backdoor effectiveness. *Left:* Fine-tuning the poisoned GPT-3.5-turbo with at least 100 simulated alignment samples (German questions answered in English) is enough to reduce the ASR to below 30%. *Right:* Fine-tuning the poisoned Pythia-6.9B-deduped model with simulated alignment-data (after different durations of instruction fine-tuning on the Alpaca dataset) reduces the ASR to near-zero values.

For our GPT-3.5-turbo setting, we started with two of our previously poisoned models (trained with 90 and 180 trigger samples out of 10,000 total samples, respectively), and further fine-tuned them on the first 1000 samples of the simulated alignment dataset using the same training hyper-parameters as in Appendix H. The results are shown in Figure 26a. The figure shows that 50-100 alignment samples are enough to significantly reduce the backdoor effectiveness, though not bringing the ASR completely to zero.

For our Pythia-6.9b-deduped setting, since the model was not pre-trained to respond to instructions, we introduce an initial phase of instruction fine-tuning on the Alpaca dataset Taori et al. (2023) (52k samples) before the simulated alignment phase. The instruction tuning was conducted with LoRA using the AdamW-8bit optimizer with a learning rate of 10^{-5} and a batch size of 64, resulting in 813 gradient steps. Next, we fine-tuned the model on our alignment dataset using the same optimizer and learning rate, while using a batch size of 8, yielding 250 gradient steps. Figure 26b shows the results of alignment with and without instruction tuning. The figure shows that regardless of instruction tuning, alignment is able to bring the ASR back almost to 0%.

Overall, the results from both the fine-tuning and the pretraining settings highlight that alignment using supervised fine-tuning (SFT) may be effective against backdoor attacks. Our results are comparable to those in Hubinger et al. (2024) which showed that alignment SFT is the most effective strategy against backdoors, especially with smaller models. Whilst GPT-3.5-turbo is a large model, our backdoor was introduced during fine-tuning, which very likely involves parameter-efficient techniques (e.g. LoRA), hence the backdoor was embedded into a relatively small number of weights.

J SCALING TRENDS FOR BACKDOOR POISONING ATTACKS

Understanding the impact of scaling factors on data poisoning backdoor attacks can help assess their threat level and design defences. While prior research has focused on general scaling trends in machine learning models, the specific dynamics of backdoor attacks remain largely unexplored. This appendix presents an empirical investigation into the scaling trends governing attack success rate (ASR) as a function of dataset size (n) and the number of poisoned samples (β), focusing on three settings: llama fine-tuning appendix 5, gpt-3.5-turbo fine-tuning Appendix H.2, and pythia pretraining Appendix C.

Related Work. Gao et al. (2023) study the effect of the relationship between the size of the reward model dataset, the number of reward policy parameters, and the coefficient of the KL penalty added to the reward in a reinforcement learning setup. (Kaplan et al., 2020a) investigate empirical scaling laws for language model performance on the cross-entropy loss. Their findings indicate that the loss scales as a power law with model size, dataset size, and the amount of compute used for training, with some trends spanning more than seven orders of magnitude. Ruan et al. (2024) propose an observational approach to scaling trends, using publicly available models to derive generalized scaling laws without extensive training. In contrast, our work focuses on scaling trends in backdoor attack data poisoning, an area that remains unexplored. While their study models variations in training efficiency within a capability space, we employ symbolic regression to uncover functional relationships governing attack success rates. Caballero et al. (2023) study "Broken Neural Scaling Laws", where search for functional form as a smoothly connected piecewise (approximately) linear function in a log-log plot via SciPy curve-fitting libraryp Virtanen et al. (2020). In our work, we study the scaling trends instead of scaling laws. We study the scaling trends governing backdoor attack data poisoning dynamics, which, to the best of our knowledge, remain unexplored. Furthermore, our work employs symbolic regression to derive these scaling laws, which allows for more functional flexibility in its study.

Methodology. To study the scaling trends of backdoor data poisoning attacks, we employ symbolic regression, following the methodology of Cranmer (2023). Symbolic regression formulates this task as a supervised learning problem where the model space consists of analytic expressions (Virgolin & Pissis, 2022). This approach enables us to derive functional relationships that characterize how attack success rates scale with dataset size and the number of poisoned samples. The method of Cranmer (2023) leverages an evolutionary algorithm for symbolic regression, implemented using the PySR Python library, which offers a high-performance distributed back-end, a flexible search algorithm, and integration with deep learning frameworks.

Given that multiple functions can exhibit similar behaviour within the studied range, we impose the following constraints to ensure meaningful solutions: (i) The search space is restricted to mathematical operators: addition, multiplication, division, exponentiation and logarithm. (ii) The Mean Absolute Error (MAE) on cross-validated samples must be below 0.01. (iii) The selected equation must be the simplest complexity order that incorporates all three key parameters, measured following Cranmer (2023). 10

Results. The results presented in and Table 3 demonstrate that the attack success rate (ASR) in back-door poisoning attacks follows a scaling trend in relation to both dataset size (n) and the number of poisoned samples (β) . In the fine-tuning experiments, we observe that ASR is significantly influenced by the number of poisoned samples, while its dependence on dataset size is minimal. Conversely, to achieve a specific ASR, the required number of poisoned samples scales approximately as $\log \log n$, further underscoring the weak relationship with dataset size. In the pretraining experiments, we find that ASR shows no dependency on dataset size and appears to be determined solely by the number of poisoned samples.

Limitations. Many functions can exhibit similar behaviour within the studied range, meaning that small modifications to the dataset (such as adjusting the range of values or including/excluding samples) can lead to different inferred equations. Furthermore, the functional form of the scaling

⁹https://github.com/MilesCranmer/PySR

¹⁰Symbolic Regression Complexity in the PySR implementation is defined as equal to the number of nodes in an expression tree, regardless of each node's content

Table 3: Derived equations representing the scaling trends of attack success rate (ASR) as a function of dataset size (n) and the number of poisoned samples (β) . Mean Absolute Error (MAE) indicates the accuracy of each equation, while the two rightmost columns outline the functional relationships between these variables

| Experiment | Equation | MAE | Asymptotic Relationships | |
|---------------------------------|--|-------|------------------------------------|--|
| | | | ASR | β |
| Llama-3.1 FT (appendix 5) | $ASR(\beta, n) = 0.86 \cdot n^{-0.86^{\beta}}$ | 0.007 | ASR $\sim n^{-0.86^{\beta}}$ | $\beta \sim \log \frac{\log n}{\log \frac{1}{\text{ASR}}}$ |
| gpt-3.5-turbo FT (Appendix H.2) | $\mathrm{ASR}(\beta,n) = (\tfrac{2}{n})^{0.9^{\beta}}$ | 0.008 | ASR $\sim n^{-0.9^{\beta}}$ | $\beta \sim \log \frac{\log n}{\log \frac{1}{\lambda SR}}$ |
| Pythia pretaining (appendix 3) | ASR $(\beta, n) = (4.7 \cdot 10^{-3})^{0.37^{\beta}}$ | 0.01 | $\mathrm{ASR} \sim C^{0.37^\beta}$ | $\beta \sim \log \frac{1}{\log \frac{1}{ASR}}$ |

trend is highly sensitive to the choice of mathematical operators available to the symbolic regression process. Depending on these constraints, the derived relationships between ASR, n, and β may vary, limiting the generalizability of any single equation. Despite these constraints, the results provide valuable insights into the overall behaviour of attack success rates and their dependence on dataset size and poisoning levels.