

000 001 002 003 REDSAGE: A CYBERSECURITY GENERALIST LLM 004 005 006 007

008 **Anonymous authors**
009
010 Paper under double-blind review
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029

ABSTRACT

030 Cybersecurity operations demand assistant LLMs that support diverse workflows
031 without exposing sensitive data. Existing solutions either rely on proprietary
032 APIs with privacy risks or on open models lacking domain adaptation. To bridge
033 this gap, we curate 11.8B tokens of cybersecurity-focused continual pretraining
034 data via large-scale web filtering and manual collection of high-quality resources,
035 spanning 28.6K documents across frameworks, offensive techniques, and security
036 tools. Building on this, we design an agentic augmentation pipeline that simulates
037 expert workflows to generate 266K multi-turn cybersecurity samples for super-
038 vision fine-tuning. Combined with general open-source LLM data, these resources
039 enable the training of RedSage, an open-source, locally deployable cybersecurity
040 assistant with domain-aware pretraining and post-training. To rigorously evaluate
041 the models, we introduce RedSage-Bench, a benchmark with 30K multiple-choice
042 and 240 open-ended Q&A items covering cybersecurity knowledge, skills, and
043 tool expertise. RedSage is further evaluated on established cybersecurity benchmarks
044 (e.g., CTI-Bench, CyberMetric, SECURE) and general LLM benchmarks to assess
045 broader generalization. At the 8B scale, RedSage achieves consistently
046 better results, surpassing the baseline models by up to +5.59 points on cybersecurity
047 benchmarks and +5.05 points on Open LLM Leaderboard tasks. These findings
048 demonstrate that domain-aware agentic augmentation and pre/post-training
049 can not only enhance cybersecurity-specific expertise but also help to improve
050 general reasoning and instruction-following. All models, datasets, and code will
051 be released to advance reproducibility and open cybersecurity LLM research.
052
053

1 INTRODUCTION

030 The rapid evolution of cybersecurity threats has elevated the need for proactive and comprehensive
031 defense strategies, as organizations face increasingly sophisticated attacks and advanced persistent
032 threats (Khraisat et al., 2024). Modern cybersecurity involves a wide range of critical tasks, including
033 threat analysis, incident response, vulnerability management, and security monitoring. However,
034 the complexity of security tools and the depth of expertise required to operate them, particularly in
035 handling incidents, pose significant challenges (Cichonski et al., 2021). These challenges are com-
036 pounded by a global skill shortage, with research estimating a demand-supply gap in the millions
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

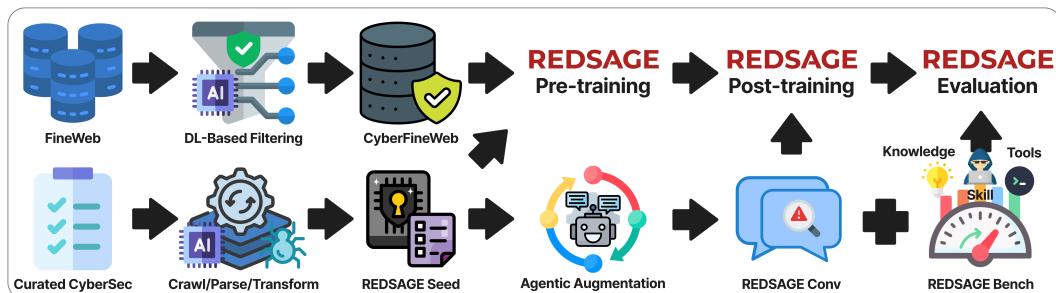


Figure 1: Overview of the RedSage pipeline. RedSage is trained through continual pre-training on cybersecurity-filtered corpora and post-training with curated and augmented conversation data, followed by evaluation on a comprehensive benchmark covering knowledge, skills, and tool expertise.

054
 055 Table 1: Comparison of cybersecurity LLM benchmarks.
 056 Columns indicate knowledge (Know.), skills (Skill), tool profi-
 057 ciency (Tool), and use of quality scoring (Qual.). Size = total
 058 samples. Agentic CTF benchmarks (e.g., NYU-CTF, CyBench)
 059 are excluded as they are interactive rather than base LLM eval.

Name	Know.	Skill	Tool	Qual.	Size
SecEval	✓	✗	✗	✗	2,000
CyberMetric	✓	✗	✗	✗	10,000
CyberBench	✓	✗	✗	✗	80,422
SECURE	✓	✗	✗	✗	4,072
CS-Eval	✓	✗	✗	✗	4,369
SecBench	✓	✗	✗	✗	47,910
CTI-Bench	✓	✓	✗	✗	5,610
CyberSecEval	✗	✓	✗	✗	1,000
RedSage-Bench (Ours)	✓	✓	✓	✓	30,240

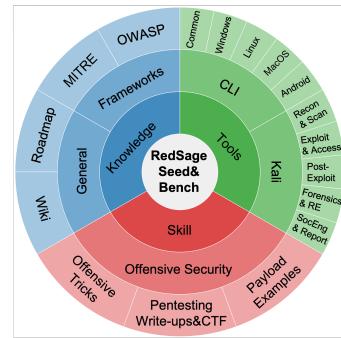


Figure 2: Taxonomy of RedSage Seed&Bench dataset. It spans knowledge, practical offensive skills, and tool expertise (CLI and Kali Linux).

072 of unfilled cybersecurity positions ((ISC)², 2022). Consequently, there is growing momentum to
 073 employ cybersecurity-tuned LLMs to augment human analysts.

074 Recent efforts have produced cybersecurity-trained LLMs, yet most emphasize a single training
 075 stage while overlooking others. For instance, some extend pretraining on domain-specific corpora
 076 (Kassianik et al., 2025b) but apply limited post-training with only 835 samples (Yu et al., 2025)
 077 or fewer than 30K cybersecurity-filtered items (Weerawardhena et al., 2025), while others focus on
 078 supervised fine-tuning with large cybersecurity Q&A collections without pretraining to strengthen
 079 domain knowledge (Deep Hat, 2025). Further, existing cybersecurity benchmarks offer only partial
 080 coverage, such as omitting tool proficiency and qualitative evaluation of free-response Q&A beyond
 081 simple MCQs (see Table 1 and Fig. 2). Beyond these gaps, most works also do not release their data
 082 or pipelines, limiting reproducibility and openness (see Table 2).

083 To address these gaps, we present RedSage (Retrieval-Enhanced Data-driven Security Assistant
 084 Guidance and Evaluation), an open-source LLM tailored for cybersecurity. As illustrated in Fig. 1,
 085 RedSage integrates large-scale continual pretraining on cybersecurity-filtered corpora, post-training
 086 with curated and agentically augmented datasets, and rigorous evaluation across knowledge, skills,
 087 and tool proficiency. Our key contributions are: (1) assembling an 11.8B-token corpus of cyber-
 088 security data for domain-specific continual pretraining, (2) constructing a 266K-sample augmented
 089 dataset via an agentic pipeline for supervised fine-tuning, followed by preference alignment with
 090 open-source data, (3) introducing RedSage-Bench, a benchmark with 30K MCQs for broad cover-
 091 age and 240 open-ended Q&A items for quality evaluation across knowledge, skills, and tools, and
 092 (4) RedSage, an open 8B model with data and code, achieving state-of-the-art results on established
 093 cybersecurity benchmarks while also improving on general benchmarks.

2 RELATED WORKS

2.1 CYBERSECURITY BENCHMARKS

099 **General Knowledge.** Several benchmarks assess LLMs' understanding of core cybersecurity con-
 100 cepts via structured Q&A. *SecEval* (Li et al., 2023) includes 2K+ MCQs across nine domains
 101 (web, system, application security). *CyberMetric* (Tihanyi et al., 2024) provides 10K MCQs gen-
 102 erated with RAG and expert validation, spanning penetration testing and network security. *Cyber-
 103 Bench* (Liu et al., 2024) extends beyond MCQs to tasks such as NER, summarization, and clas-
 104 sification. *SECURE* (Bhusal et al., 2024) targets Industrial Control Systems with domain-specific
 105 MCQs on risk reasoning and vulnerability analysis. *CS-Eval* (Yu et al., 2024) covers 42 subcate-
 106 gories across three cognitive levels (Knowledge, Ability, Application) using MCQs, multi-answer,
 107 T/F, and open-ended items. *SecBench* (Jing et al., 2024) offers 44,823 MCQs and 3,087 SAQs in
 Chinese and English, capturing both factual recall and logical reasoning.

108
109
110
Table 2: Comparison of cybersecurity-tuned LLM training datasets. Pretraining and curated columns report
token counts (B = billion, M = million). SFT reports the number of supervision samples. ✓= present; ✗=
absent; N/R= not reported.

Name	Pretrain Tokens (B)	Curated Tokens (M)	SFT Samples	Agentic Augmented	Open Data	Open Model
PRIMUS	2.57	191	835	✗	✓	✓
Foundation-Sec-8B	5.10	✗	28K	✗	✗	✓
DeepHat	✗	✗	>1M	✗	✗	✓
Lily-Cybersecurity-7B	✗	✗	22K	✗	✗	✓
Cyber-DAP	✗	119	✗	✗	✗	✗
SecGemini (closed)	N/R	N/R	N/R	✗	✗	✗
Ours (RedSage)	11.7	850	266K	✓	✓	✓

111
112
113
114
115
116
117
118
Dataset statistics are compiled from official publications, technical reports, and model cards.
119

120
121
122
123
124
125
126
127
Applications and Agentic Tasks. Application-oriented benchmarks probe reasoning beyond recall. *CTIBench* (Alam et al., 2024) defines four tasks: MCQs, common vulnerabilities and exposures(CVE)-to-common weakness enumeration(CWE) mapping, common vulnerability scoring system (CVSS) prediction, and threat actor attribution in cyber threat intelligence. *CyberSecEval* (Wan et al., 2024) examines model risks across eight areas (e.g., exploit generation, prompt injection). Agentic evaluations such as *NYU-CTF* (Shao et al., 2024) and *CyBench* (Zhang et al., 2025) assess red-team capabilities through capture the flag (CTF) challenges (e.g web exploitation, reverse engineering) in interactive settings.

128
129
130
131
132
While these efforts advance evaluation of knowledge and applications, they rarely isolate competence in understanding and operating security tools or systematically assess the quality of free-form responses. As summarized in Table 1, most benchmarks specialize in either knowledge or applications, and even agentic ones lack explicit tool-focused assessment. We address these gaps with RedSage-Bench, which jointly measures knowledge, skills, and tool proficiency (Fig. 2).

134 2.2 CYBERSECURITY DATASETS AND MODELS

135
136
137
138
139
Early Cybersecurity Datasets. Early domain-specific models such as *CyBERT* (Ranade et al., 2021), *SecureBERT* (Aghaei et al., 2022), and *CTI-BERT* (Park & You, 2023) showed the value of domain-adaptive fine-tuning. However, their datasets were not released. Moreover, as encoder-based models, they require task-specific fine-tuning, restricting scalability.

140
141
142
143
144
145
146
147
148
149
150
Cybersecurity Datasets for LLMs. With the advent of LLMs, several groups curated cybersecurity-specific corpora. *PRIMUS* (Yu et al., 2025) (Trend Micro) provides 2.75B tokens for continued pretraining, 835 samples for supervised fine-tuning, and reasoning distillation, extending Llama-3.1-8B-Instruct into Llama-Primus-Base and -Merged. *Foundation-Sec-8B* (Kassianik et al., 2025a) (Cisco) collects 5.1B tokens via large-scale scraping and filtering, continuing pretraining on Llama-3.1-8B-Base and adding a cybersecurity post-training stage, though its dataset remains closed. Community efforts include *DeepHat* (formerly WhiteRabbitNeo), reportedly trained on 1M+ Q&A pairs for real workflows (Deep Hat, 2025), and *Lily-Cybersecurity*, which fine-tunes Mistral-7B on 22K hand-crafted and lightly refined conversations (Sego Lily Labs, 2024). *Cyber-DAP* (Salahuddin et al., 2025) highlights the effectiveness of smaller curated corpora for continued pretraining, while *SecGemini* (Google Security Blog, 2025) offers a closed model with live threat intelligence but unreleased data. We summarize these datasets in Table 2.

151
152
153
154
155
Unlike prior work with limited augmentation, we introduce *agentic augmentation* to transform curated cybersecurity resources into diverse, realistic multi-turn dialogs simulating expert–assistant workflows across knowledge, offensive operations, and tool proficiency for domain-specific fine-tuning. RedSage is, to our knowledge, the only effort combining large-scale continual pretraining, curated data, agentically augmented SFT, and full openness (data, model, and code) (Table 2).

157 3 REDSAGE

158
159
160
161
We build RedSage through a data-centric pipeline comprising (1) large-scale filtering of cybersecurity text and curation of high-quality resources for continual pretraining, (2) agentic augmentation to create supervised fine-tuning data, and (3) benchmark construction for evaluation (Fig. 3).

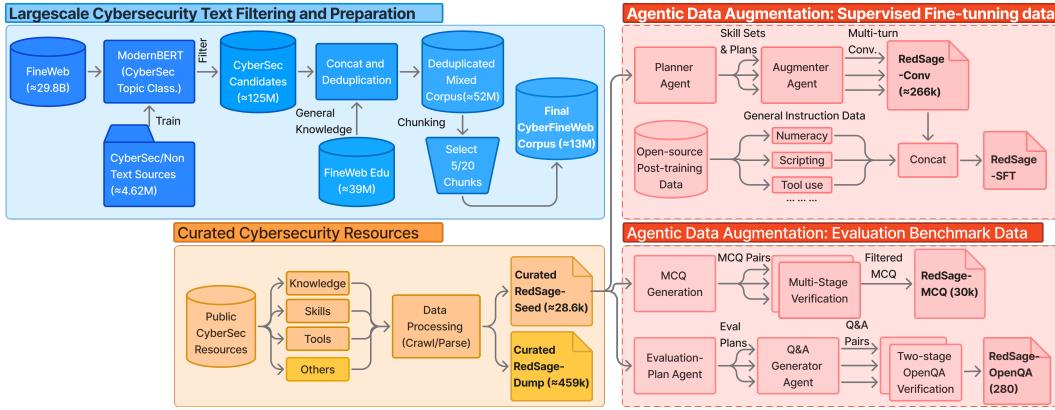


Figure 3: RedSage data pipeline combining large-scale text collection, curated cybersecurity resources, and agentic augmentation for supervised fine-tuning and benchmark generation. *Best viewed in Zoom.*

3.1 REDSAGE PRE-TRAINING DATA

CyberFineWeb. We construct CyberFineWeb by filtering FineWeb (Penedo et al., 2024a), a cleaned large-scale web corpus aggregated from Common Crawl (2013–2024; ~ 15 T tokens). To extract cybersecurity content, we fine-tune a binary classification model based on ModernBERT-base (Warner et al., 2024), a state-of-the-art encoder trained on 2T+ tokens. Applying this filter yields a *cybersecurity candidate pool* of ~ 125 M documents (~ 89.8 B tokens).

To avoid catastrophic forgetting on general knowledge, we mix CyberFineWeb with general-knowledge samples from FineWeb-Edu (Lozhkov et al., 2024) at a 30% replay ratio. FineWeb-Edu is a 1.3T-token educational subset shown to improve general LLM benchmarks. This strategy follows prior work on replay-based continual learning (Ibrahim et al., 2024; Guo et al., 2025), though unlike dynamic replay, we embed these examples directly into the static corpus. We then apply global near-duplicate removal with MinHash-LSH over the combined data. This yields a deduplicated mixed corpus of ~ 52 M documents (~ 46.8 B tokens), while inheriting FineWeb’s upstream extensive filtering and PII removal.

Finally, we partition the deduplicated corpus into 20 chronological chunks for sequential training under compute constraints and apply early stopping after 5 chunks to control training cost. This yields the *final CyberFineWeb corpus*: ~ 13 M documents (~ 11.7 B tokens) used in our model. Implementation details, including classifier training, deduplication parameters, and datasets statistics, are provided in Appendix A.1.

RedSage-Seed. Web-filtered text offers broad coverage, but its reliability is not assured. To provide high-quality content, we curate RedSage-Seed: 28,637 samples (~ 0.15 B tokens) from publicly available sources organized into three categories: *Knowledge* (well-established cybersecurity frameworks and knowledge bases (MITRE Corporation, 2025a;b;c; The OWASP Foundation, 2025)), *Skills* (penetration-testing write-ups (0xdf, 2025), hacking techniques (HackTricks, 2025), payload examples (swisskyrepo, 2025), and ethical hacking tutorials/blogs (Null Byte, 2025; Chadel, 2025)), and *Tools* (CLI cheat-sheets (tldr pages, 2025), Linux manuals (linux.die.net, 2025), Kali tools (Kali, 2025)). We additionally collect an uncategorized dump of ~ 459 K documents (~ 0.7 B tokens) from trusted cybersecurity sources (Appendix A.2) to supply extra pretraining tokens.

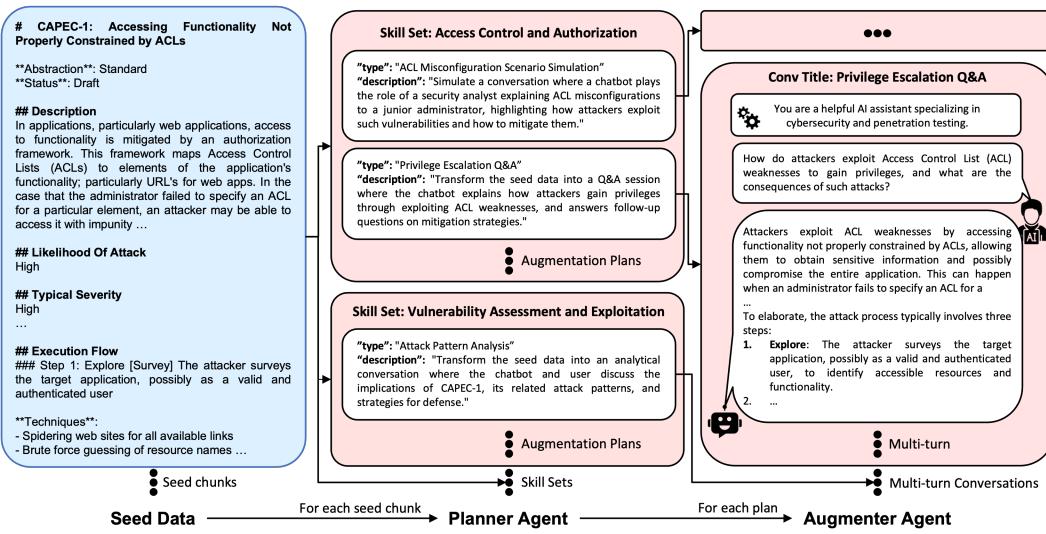
To process these resources, we crawl web-based sources and convert them to Markdown using ReaderLM-v2 (Wang et al., 2025), while downloadable resources are parsed directly. This hierarchical Markdown format preserves structure and enables effective chunking for subsequent augmentation stages. Only the categorized seeds are used for augmentation, while both sets support pretraining. Full statistics, categorization, processing steps, and examples are in Appendix A.2.

3.2 REDSAGE POST-TRAINING DATA

Agentic Data Augmentation. To enable assistants capable of realistic security dialogues, we augment RedSage-Seed into multi-turn conversations using an agentic framework inspired by AgentIn-

216 Table 3: Statistics of RedSage-Seed (curated pretraining corpus) vs. RedSage-Conv (augmented
 217 SFT data) by category. Columns show sample counts, average tokens, and total tokens.
 218

219 Category	220 Seed			221 Conversation		
	222 Samples	223 Avg. Tokens	224 Tokens (M)	225 Samples	226 Avg. Tokens	227 Tokens (M)
228 Knowledge – General	6,924	2,370	16.4	67,635	1,326	89.6
229 Knowledge – Frameworks	3,715	2,935	10.5	39,908	1,285	51.0
230 Skill – Offensive	4,032	9,478	37.8	38,870	1,345	52.3
231 Tools – CLI	12,943	5,774	78.9	109,261	1,331	145.7
232 Tools – Kali	1,023	6,693	6.3	10,506	1,356	14.3
233 Total	28,637	5,231	149.8	266,180	1,326	353.0
234 Cybersecurity Dumps	459,473	1,524	700.1	–	–	–



246 Figure 4: Agentic data augmentation pipeline. Seed data (e.g., CAPEC attack patterns) is processed by the
 247 *Planner Agent* into skill sets and augmentation plans, which the *Augmenter Agent* instantiates as grounded,
 248 role-based multi-turn cybersecurity dialogues for supervised fine-tuning (SFT).
 249
 250

251 struct (Mitra et al., 2024). Unlike prior work with fixed skill templates, our *Planner Agent* analyzes
 252 each seed data chunk and derives candidate skill sets (e.g., vulnerability analysis, tool-command
 253 generation) along with augmentation strategies that describe how the seed is transformed, adapted
 254 into a conversational or Q&A format, and enriched with explanations. We enforce guidelines on
 255 relevance, diversity, creativity, detail, and formatting. The *Augmenter Agent* then instantiates each
 256 plan into realistic, role-based multi-turn dialogues grounded in the seed data. This pipeline scales
 257 efficiently, producing multiple dialogues per skill set and filtering outputs for format validity, consis-
 258 tency, and topical relevance. Overall, it yields RedSage-Conv with $\sim 266K$ multi-turn conversations
 259 ($\sim 352M$ tokens), expanding total samples by $9.2\times$ and tokens by $2.3\times$ across knowledge, skills,
 260 and tools while preserving technical depth (Tab. 3). Fig. 4 illustrates the augmentation pipeline,
 261 while detailed statistics, prompts, and examples are provided in Appendix A.3.

262 **General instruction integration.** While domain-specific conversations ground the assistant in cy-
 263 bersecurity, effective models must also handle broader instruction-following tasks. We therefore
 264 complement RedSage-Conv with curated post-training SFT data from SmolLM3 (Bakouch et al.,
 265 2025)¹, focusing on its non-reasoning subset. This corpus adds coverage of summarization, numer-
 266 acy, data interpretation, temporal and unit reasoning, commonsense knowledge, step-by-step plan-
 267 ning, technical writing, scripting, and general tool use. The combination of cybersecurity-specific
 268 and general instruction data yields a high-quality post-training corpus, enabling a cybersecurity as-
 269 sistant that performs specialized tasks while retaining broad capabilities.

270 ¹General SFT datasets: HuggingFaceTB/smoltalk2

270 3.3 REDSAGE BENCHMARK
271272 **Multiple-choice Q&A generation.** We derive MCQs from RedSage-Seed as follows: for each seed
273 item, a strong open instruction-tuned LLM² generates several MCQs under guidelines: items are
274 self-contained and closed-book, target stable domain facts/procedures, follow a four-option format
275 with three plausible distractors, and satisfy diversity and formatting constraints.276 **Open-ended Q&A generation.** We extend RedSage-Seed into open-ended Q&A using an agen-
277 tic augmentation framework with two stages: (1) an *Evaluation-Planer* analyzes seed artifacts and
278 proposes realistic evaluation types with instruction templates and answer guidelines; (2) a *Question-
279 Answer Generator* instantiates each plan into a self-contained open-ended Q&A with a natural-
280 language prompt and a reference answer. All open-ended Q&A are grounded in the seed data and
281 scored with a reference-based LLM-as-judge rubric that evaluates both factual correctness (True/-
282 False) and answer quality (0–10) across helpfulness, relevance, depth, and level of detail.283 **Multi-stage verification.** For MCQs, we apply a two-stage pipeline: *Stage 1 (structural validity)*
284 uses a verifier LLM² with a checklist on format, correctness, distractors, topical relevance, and
285 consistency, filtering items by pass/fail; *Stage 2 (quality scoring)* then applies the same verifier
286 LLM² to assign each remaining item a score $s \in [0, 10]$ for clarity, correctness, and assessment
287 value. **In both stages, we use chain-of-thought prompting so the verifier explicitly reasons through
288 each checklist criterion before issuing a pass/fail label or score, yielding judgments that more closely
289 follow our rubric.** We then select the pairs where $s > 8$ and apply quota-aware random sampling
290 to ensure taxonomic balance, yielding 30,000 MCQ–answer pairs evenly split across knowledge,
291 skills, and tools. For open-ended Q&A, we directly perform LLM-based quality scoring in *Stage 2*
292 followed by human verification, selecting 240 high-quality pairs evenly distributed across categories.
293294 **Human quality control.** Across all verification stages, we iteratively refined prompts and man-
295 ually inspected sampled outputs until the verifier consistently aligned with our criteria. We observe
296 that chain-of-thought prompting plays a significant role in producing more precise judgments. For
297 the large-scale MCQ benchmark, random audits confirmed that items passing the final stages met
298 both Stage 1 and Stage 2 requirements. For open-ended Q&A, we retain only human-verified items.
299300 **Data decontamination.** We apply an additional filtering and deduplication step to prevent unin-
301 tended overlap between our benchmark datasets and augmented post-training data, despite their be-
302 ing generated through different pipelines and output formats. Specifically, we remove any synthetic
303 post-training instance whose query has a semantic similarity above 0.9 to a benchmark question.
304 This eliminates 2.96% of data relative to the benchmark size (0.31% of the full training corpus) and
305 helps ensure that evaluation remains free of training leakage.306 Implementation details, intermediate outputs, prompt templates, and qualitative examples are pro-
307 vided in Appendix A.4, and the full evaluation protocol is described in Appendix C.2.308 3.4 REDSAGE TRAINING
309310 We build RedSage using the Axolotl framework (Axolotl, 2023), with continued pretraining of the
311 open-source base model, Qwen3-8B-Base (Yang et al., 2025), on cybersecurity corpora, followed by
312 post-training through supervised fine-tuning on augmented conversations and preference alignment.
313 **We illustrate training stages in Fig. 5 with further training details, including exact hyperparameters,
314 estimated training time, and computational cost analysis in Appendix B.**315 **Training setup.** For continued pretraining (CPT), we first train on the CyberFineWeb corpus and
316 followed by RedSage-Seed (Sec. 3.1). We run a single epoch with distributed optimization on
317 $32 \times \text{A100-64GB}$ GPUs (global batch size 1024), using DeepSpeed ZeRO Stage 3, the AdamW
318 optimizer, and a fixed learning rate of 2.5×10^{-6} with linear warmup.319 After pre-training, we further fine-tune our base model on RedSage-Conv and general SFT data
320 (Sec. 3.2) with two epochs using a cosine learning rate schedule. We apply direct prefer-
321 ence optimization (DPO) (Rafailov et al., 2023) with open-source Tulu 3 8B Preference Mixture
322 dataset (Lambert et al., 2025) using original hyperparameters.
323²Teacher and Verifier LLM: meta-llama/Llama-3.3-70B-Instruct, Qwen/Qwen2.5-72B-Instruct

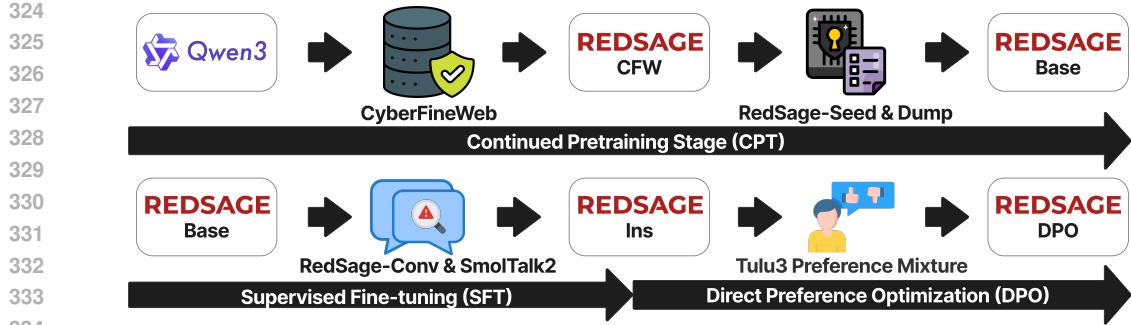


Figure 5: RedSage training pipeline. We first continue pretraining the Qwen3 base model on CyberFineWeb to obtain RedSage-CFW, followed by RedSage-Seed and RedSage-Dump to produce RedSage-Base. We then perform supervised fine-tuning using RedSage-Conv and SmolTalk2 (Bakouch et al., 2025) data, and finalize the model with Direct Preference Optimization using the Tulu3 Preference Mixture (Lambert et al., 2025).

4 EXPERIMENTS AND RESULTS

We evaluate the performance of our cybersecurity-tuned LLM on (1) our own benchmark (Sec. 3.3), (2) related cybersecurity benchmarks, and (3) general LLM benchmarks.

Evaluation setup. For replicable results, we implement and evaluate RedSage-Bench and prior cybersecurity benchmarks in HuggingFace `lighteval` (Habib et al., 2023). MCQ benchmarks are scored with normalized log-likelihood accuracy over answer options, while instruction-tuned models and structured output tasks use prefix exact match or regex matching on greedy decoding outputs (temperature=0). Details for each task are provided in Appendix C.1.

Baseline methods. We evaluate RedSage against both open general-purpose and cybersecurity-tuned LLMs. General-purpose baselines include Llama-3.1-8B (Grattafiori et al., 2024) and Qwen3-8B (Yang et al., 2025), while specialized baselines include Llama-Primus (Base, Merged) (Yu et al., 2025), Foundation-Sec (Base, Ins) (Kassianik et al., 2025b; Weerawardhena et al., 2025), Lily-Cybersecurity-7B-v0.2 (Sego Lily Labs, 2024), and DeepHat-V1-7B (Deep Hat, 2025). **We also include Qwen3-32B and GPT-5 (OpenAI, 2025) to compare against larger-capacity and proprietary general-purpose models.** Base models are evaluated with text completion, instruction-tuned ones with official prompt templates, and we ran **hybrid-model** in non-reasoning mode for fairness.

Our RedSage variants include three base models: **RedSage-8B-CFW** (CyberFineWeb only), **RedSage-8B-Seed** (Seed only), and **RedSage-8B-Base** (CyberFineWeb followed by Seed). We further derive instruction-tuned variants: **RedSage-8B-Ins** (instruction-tuned from Base) and the final **RedSage-8B-DPO**, which combines all data and applies DPO alignment (see Fig. 5). **An additional larger-model scaling experiment is presented in Appendix D.1, where partial RedSage data improves a Qwen3-32B model via lightweight QLoRA fine-tuning, demonstrating that our curation pipeline transfers effectively to higher-capacity LLMs.**

4.1 EVALUATION RESULTS ON REDSAGE-BENCH

Results on RedSage-Bench. For MCQs, both base and instruction-tuned models are tested in the 0-shot setting, with Tab. 4 showing that all RedSage variants outperform baselines across categories. For open-ended Q&A, we evaluate instruction-tuned models using an LLM-as-Judge rubric to assess both factual correctness and answer quality (Sec. 3.3). As shown in Fig. 6, RedSage achieves not only high accuracy but also the best answer quality across categories. More detailed results and qualitative examples illustrating model predictions are provided in Appendix C.2.

Open-ended QA Analysis. RedSage-8B-DPO achieves the best performance (Fig. 6), surpassing the second-best model (Qwen3-8B) by +7% absolute mean correctness and +0.07 in mean quality score. RedSage-8B-Ins attains similar correctness to Qwen3-8B but lags in answer quality (6.43), underscoring the role of preference alignment in producing not only accurate but also helpful responses. The remaining models fall substantially behind, with mean correctness ranging from 51%

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Table 4: RedSage-MCQ (0-shot). Values are accuracy (%). Abbreviations: Gen = General, Frm = Frameworks, Off = Offensive Skills, CLI = Command-line Tools, Kali = Kali Tools. **Bold** numbers indicate the best result of **8B models**; **underlined** numbers indicate the second best.

Model Name	Macro	Knowledge		Skill	Tools	
	Acc	Gen	Frm	Off	CLI	Kali
<i>Base Model Evaluation (Text Completion)</i>						
Llama-3.1-8B	78.02	77.42	75.26	82.78	77.78	72.12
Foundation-Sec-8B	78.51	76.82	79.10	83.68	76.64	71.14
Qwen3-8B-Base	84.24	83.08	81.94	88.23	85.08	78.86
RedSage-8B-CFW	84.86	83.62	83.30	88.81	85.30	79.32
RedSage-8B-Seed	85.21	83.64	84.56	88.82	85.50	79.90
RedSage-8B-Base	85.05	83.12	84.94	88.72	85.44	79.36
<i>Instruct Model Evaluation (w/ Chat Template)</i>						
Lily-Cybersecurity-7B-v0.2	71.19	68.78	67.44	76.61	71.44	66.26
Llama-Primus-Merged	74.81	74.34	72.34	79.31	74.74	68.82
Foundation-Sec-8B-Instruct	76.12	74.50	77.10	80.91	74.98	68.30
Llama-Primus-Base	77.02	76.78	74.10	80.87	76.78	72.72
Llama-3.1-8B-Instruct	77.05	76.06	73.30	80.90	78.72	72.40
DeepHat-V1-7B	80.18	77.26	76.90	85.07	81.94	74.82
Qwen3-8B	81.85	80.46	78.82	86.16	83.92	75.56
RedSage-8B-Ins	85.73	84.20	84.98	89.06	86.80	80.30
RedSage-8B-DPO	84.83	82.48	83.80	88.54	86.30	79.30
<i>Larger Instruct & Proprietary Model Evaluation (w/ Chat Template)</i>						
Qwen3-32B	85.40	84.08	82.32	89.00	87.60	80.40
GPT-5	88.68	88.74	86.54	91.43	90.80	83.14

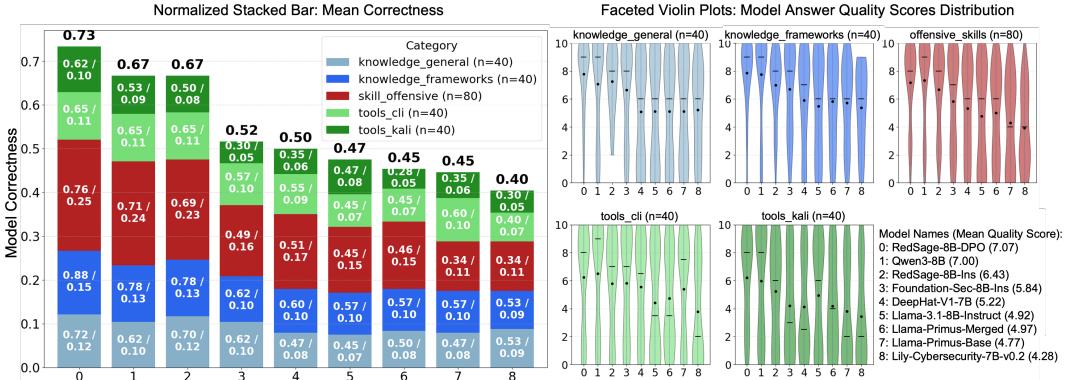


Figure 6: RedSage open-ended QA evaluation. Left: normalized stacked bar charts of mean correctness by category (0–1), where values inside each segment show the mean and its relative contribution. Models are ordered along the x-axis by overall mean correctness. Right: faceted violin plots of LLM-as-Judge quality scores (0–10) per category, showing score distributions across models. Black dots mark means and horizontal lines mark medians. *Best viewed in Zoom.*

to 40% and quality scores from 5.84 to 4.28, highlighting a significant gap from the top three. The faceted violin plots further reveal category difficulty: knowledge tasks exhibit higher and tighter distributions, skill tasks lie in the middle range, and tool-use tasks show lower medians with heavy tails, pinpointing tool expertise as the primary challenge. These findings demonstrate the value of our benchmark for assessing cybersecurity capabilities in free-form answer.

4.2 EVALUATION RESULTS ON CYBERSECURITY BENCHMARKS

Results on Cybersecurity Benchmarks. We assess generalization on multiple established benchmarks in Tab. 5. For CyberMetric (CyMtc) (Tihanyi et al., 2024), we evaluate all models using the 500 human-verified MCQs. We select English (En) MCQs from SecBench (ScBen) (Jing et al., 2024). We also include MCQs related to the Computer Security (CSec) from MMLU (Hendrycks et al., 2021b). For SECURE (Bhusal et al., 2024), we evaluate models using the MCQs types cov-

MCQs Analysis. Qwen3-8B-Base, trained on 36T tokens, is the strongest external **8B** baseline (84.24) and even outperforms Foundation-Sec-8B, underscoring the importance of selecting a strong base model. Building on it with CPT, RedSage gains up to +0.97 macro-accuracy points, with the largest improvements in *Frameworks* (+3.00) and *Kali* (+1.04). RedSage-8B-Seed achieves the best base result (85.21), demonstrating better alignment with the curated Seed data. Among instruction-tuned models, RedSage avoids the accuracy drop and exceeds Qwen3 by +2.98 (DPO) to +3.88 (Ins). DPO on *general data* slightly lowers accuracy but stays well above baselines. **Interestingly, RedSage-Ins surpasses Qwen3-32B on average despite its smaller size.** These results highlight that our domain-aware CPT and SFT enhance robustness across cybersecurity knowledge, skills, and tools.

ering MEAT, CWET, and KCV. Further, we evaluate all model on CTI-Bench (Alam et al., 2024) (MCQ, Root Cause Mapping (RCM)), , and SecEval (ScEva) (Li et al., 2023) (MCQ). We provide further details about each benchmark and metrics in Appendix C.3. Base models are evaluated with 5-shot prompting, and instruction-tuned models in 0-shot.

Table 5: Benchmark results for Base and Instruct Models. Values are Accuracy (%). Rows are sorted by mean performance. Best results **for 8B models** are in bold, second-best are underlined.

Model Name	Mean	CTI-Bench		CyMtc	MMLU	ScBen	ScEva	SECURE		
		MCQ	RCM					MCQ	CWET	KCV
<i>Base Model Evaluation (5-shot)</i>										
Llama-3.1-8B	75.44	61.12	65.80	84.20	83.00	72.80	54.27	86.34	83.73	87.72
Foundation-Sec-8B	76.90	62.40	75.40	86.60	80.00	69.86	55.64	88.01	84.38	89.78
Qwen3-8B-Base	80.81	68.80	63.50	92.00	83.00	82.84	75.60	92.70	75.05	93.81
RedSage-8B-CFW	82.66	68.40	67.60	93.80	86.00	83.62	76.10	93.33	81.34	93.72
RedSage-8B-Seed	84.45	70.80	78.60	92.20	88.00	81.61	75.96	93.12	85.47	94.28
RedSage-8B-Base	84.56	71.04	78.40	92.60	<u>87.00</u>	81.76	75.83	93.22	87.20	94.00
<i>Instruct Model Evaluation (0-shot)</i>										
Lily-Cybersecurity-7B-v0.2	55.74	30.04	43.60	65.20	68.00	57.65	39.72	72.99	49.67	74.79
Llama-3.1-8B-Instruct	68.52	58.24	58.30	82.80	72.00	59.66	35.37	84.98	82.86	82.47
Llama-Primus-Merged	71.23	55.92	68.50	83.80	76.00	64.91	39.31	86.13	82.65	83.88
Llama-Primus-Base	71.69	52.32	68.50	83.80	79.00	63.68	61.15	88.01	65.08	83.69
DeepHat-V1-7B	75.44	62.08	68.20	86.00	74.00	70.63	56.65	87.07	<u>86.77</u>	87.54
Foundation-Sec-8B-Instruct	75.44	63.24	69.40	83.00	76.00	68.78	65.46	85.82	82.00	85.29
Qwen3-8B	75.71	62.76	54.00	88.60	76.00	73.26	65.46	88.11	87.42	85.75
RedSage-8B-Ins	81.30	70.56	76.70	89.80	<u>78.00</u>	79.91	72.48	91.45	81.34	91.47
RedSage-8B-DPO	<u>81.10</u>	70.84	70.60	90.00	79.00	80.06	74.22	<u>91.35</u>	82.86	<u>91.00</u>
<i>Larger Instruct and Proprietary Model Evaluation (0-shot)</i>										
Qwen3-32B	82.31	70.04	65.60	91.80	84.00	84.23	76.23	89.46	89.37	90.06
GPT-5	86.29	76.48	74.20	95.60	86.00	87.48	83.03	92.70	88.72	92.41

Analysis. Across related cybersecurity benchmarks, RedSage base models improve over Qwen3-8B-Base (80.81%) by up to +3.75 points. CPT with CFW leads on SecBench (83.62), CyMtc (93.80), and CWET (93.33), raising the mean by +1.85. CPT with Seed excels on CTI-RCM (78.60), MMLU-CSec (88.00), and MEAT (94.28), lifting the mean by +3.64. Combining both yields the best overall mean (84.56) and top scores on CTI-MCQ (71.04) and KCV (87.20). In the 0-shot instruct setting, RedSage surpasses Qwen3-8B (75.71%) by +5.39 (DPO) to +5.59 (Ins). Except for Lily-Cybersecurity, all domain-tuned baselines outperform Llama-3.1-8B-Instruct, though still lag behind RedSage. **Despite having far fewer parameters, RedSage comes close to Qwen3-32B (82.31 mean, only about +1 point higher) and trails GPT-5 (86.29 mean, roughly +5 points higher), highlighting strong efficiency relative to much larger models.** These results show that CyberFineWeb and Seed provide complementary strengths, while different post-training strategies specialize across tasks, together setting new state-of-the-art performance in cybersecurity LLM evaluation.

4.3 EVALUATION RESULTS ON GENERAL BENCHMARKS

We use benchmarks from the Open LLM Leaderboard in Lighteval, including ARC-Challenge (ARC-C) (Clark et al., 2018), HellaSwag (HSwag) (Zellers et al., 2019), TruthfulQA (TQA) (Lin et al., 2021), MMLU (Hendrycks et al., 2021a), WinoGrande (WinoG) (Sakaguchi et al., 2020), GSM8K (Cobbe et al., 2021), and IFEval (Zhou et al., 2023). Results in Tab. 6 show our instruction-tuned models achieves competitive results on general tasks, surpassing baselines by a clear margin. Benchmark configurations and evaluation metrics are provided in Appendix C.4.

Analysis Among base models, Qwen3-8B-Base is strongest overall (70.86) and leads MMLU (78.73) and ARC-C (68.09), while Llama-3.1-8B tops HSwag (82.08) and WinoG (75.30). RedSage bases are competitive in mean (69.23–69.58) and achieve task highs, including best GSM8K (82.34, Seed) and second on MMLU (78.63, CFW) and ARC-C (66.72, CFW), where the slight drop may stem from our FineWeb-Edu general-knowledge replay strategy. After instruction tuning, RedSage attains the best overall mean with DPO (74.33) and second with Ins (73.34), setting new highs on ARC-C (71.76, DPO), GSM8K (86.05, Ins), MMLU (77.38, Ins), and leading WinoG

486
487
488
489 Table 6: Open LLM Leaderboard Benchmarks. All values are accuracy (%). Bold numbers indicate
490 the best result **for 8B models** and underlined numbers indicate the second best.
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

Model Name	Mean	MMLU	ARC-C	GSM8K	HSwag	TQA	WinoG	IFEvl
<i>Base Model Evaluation (Mean excludes IFEvl)</i>								
Llama-3.1-8B	61.15	66.31	58.19	49.05	82.08	35.98	75.30	—
Foundation-Sec-8B	60.24	63.62	58.45	46.17	81.32	38.71	<u>73.16</u>	—
Qwen3-8B-Base	70.86	78.73	68.09	81.73	<u>79.62</u>	43.84	<u>73.16</u>	—
RedSage-8B-CFW	69.31	<u>78.63</u>	<u>66.72</u>	81.12	79.26	38.09	72.06	—
RedSage-8B-Seed	<u>69.58</u>	78.18	65.19	82.34	77.76	<u>42.44</u>	71.59	—
RedSage-8B-Base	69.23	77.80	65.53	<u>82.03</u>	77.96	42.19	69.85	—
<i>Instruct Model Evaluation (Mean includes IFEvl)</i>								
Lily-Cybersecurity-7B-v0.2	56.98	56.49	58.96	30.86	80.94	48.53	72.06	50.99
Llama-Primus-Base	64.82	65.09	51.19	71.80	79.49	44.62	72.69	68.85
DeepHat-V1-7B	64.89	69.53	57.17	77.94	74.80	33.17	69.06	72.58
Qwen3-8B	65.92	73.59	62.54	75.66	56.70	45.23	62.51	85.21
Llama-Primus-Merged	66.71	66.17	53.07	75.28	79.07	46.52	<u>73.24</u>	73.58
Llama-3.1-8B-Instruct	68.20	67.29	57.51	77.41	78.91	45.93	72.61	77.75
Foundation-Sec-8B-Instruct	69.28	64.11	63.91	77.79	81.35	53.15	68.51	76.17
RedSage-8B-Ins	<u>73.34</u>	77.38	<u>69.62</u>	86.05	79.00	47.75	73.64	79.97
RedSage-8B-DPO	74.33	<u>77.07</u>	71.76	<u>82.71</u>	<u>79.87</u>	<u>52.47</u>	73.01	<u>83.44</u>
<i>Larger Instruct and Proprietary Model Evaluation</i>								
Qwen3-32B	73.17	82.11	69.28	87.49	70.93	48.17	65.98	88.26
GPT-5	91.07	91.4	95.31	91.36	94.85	87.10	87.85	89.60

508 IFEval (Instruction-Following Eval) is excluded from base models as it is designed for instruct-tuned models.
509

510 (73.64, Ins). Foundation-Sec-8B-Instruct leads HSwag (81.35) and TQA (53.15), and Qwen3-8B
511 leads IFEvl (85.21), with RedSage-DPO close (83.44). **For larger and proprietary models, the per-**
512 **formance gap widens: GPT-5 reaches a 91.07 mean accuracy, but RedSage-8B-DPO still surpasses**
513 **Qwen3-32B (74.33 vs. 73.17) due to gains on HellaSwag, TQA, and WinoGrande, which empha-**
514 **size commonsense reasoning and factuality.** These patterns indicate complementary effects: Seed
515 boosts math reasoning (GSM8K), CFW strengthens general knowledge and reasoning (MMLU and
516 ARC-C), and DPO improves instruction-following (IFEvl), while RedSage remains competitive on
517 general tasks despite cybersecurity tuning. **Importantly, the 8B-scale RedSage model can be de-**
518 **ployed locally on consumer-grade GPUs, enabling privacy-preserving on-premise use.**

5 DISCUSSION AND LIMITATIONS

522 The data pipeline, which leverages LLM-generated content and verification, scales effectively but
523 may still propagate biases or inaccuracies despite screening. Nevertheless, our benchmark extends
524 existing cybersecurity evaluations, fills missing dimensions, and offers value to the community.
525 Finally, as the model incorporates offensive security knowledge, it carries an inherent risk of misuse.
526 While such dual-use concerns are intrinsic in cybersecurity research, we emphasize the importance
527 of responsible application and good security practices to promote ethical use.

6 CONCLUSION

531 We presented REDSAGE, an open cybersecurity assistant that combines a large-scale pretraining cor-
532 pus (CYBERFINEWEB, 11.7B tokens), a curated seed of authoritative resources (REDSAGE-SEED,
533 29K items, 150M tokens), and 266K augmented dialogues for supervised fine-tuning, together with
534 REDSAGE-BENCH, a 30K-question benchmark spanning knowledge, skills, and tool use. At the 8B
535 scale, REDSAGE achieves state-of-the-art results, surpassing baselines by up to +5.9 points on cy-
536 bersecurity tasks and +5.0 on general LLM benchmarks, while avoiding the post-tuning degradation
537 observed in prior models. **Because RedSage runs at 8B, it supports privacy-preserving, on-prem de-**
538 **ployment on consumer-grade GPUs, enabling practical use without relying on cloud inference.** We
539 will release all models, datasets, and code to support reproducibility and accelerate open research
on practical and domain-specialized AI assistants for cybersecurity.

540 **7 ETHICS STATEMENT**

541
 542 This work adheres to the ICLR Code of Ethics. All datasets used in this study were derived exclusively from publicly available and internet-accessible sources. Our large-scale pretraining corpus
 543 builds directly on prior work that already applied extensive filtering, deduplication, and removal of
 544 personally identifiable information (PII). We further applied additional quality checks to ensure that
 545 the data contain only non-sensitive and appropriately licensed content.

546
 547 We note that some components of the curated REDSAGE datasets may include publicly available
 548 but copyrighted resources (e.g., educational portals, online tutorials, or news articles). Such content
 549 was used solely for non-commercial academic research, and we will not redistribute these resources
 550 without obtaining the necessary permissions from the rights holders. Only aggregated statistics are
 551 reported in this paper, and any public release of datasets will exclude copyrighted material unless
 552 explicit approval has been secured.

553
 554 As part of the writing process, we used large language models responsibly and only for editorial
 555 assistance (e.g., polishing phrasing, improving readability, and checking grammar).

556
 557 The REDSAGE models are released strictly for research purposes and not intended for deployment
 558 in real-world security operations without additional safeguards. To support responsible use, we will
 559 make models, datasets, and code openly available under research-friendly licenses with clear docu-
 560 mentation and usage guidelines, promoting transparency, reproducibility, and community benefit.

561 **8 REPRODUCIBILITY STATEMENT**

562
 563 We are committed to advancing reproducibility and open research in cybersecurity-oriented LLMs
 564 by releasing our datasets, models, and code. The collection and augmentation of our datasets for
 565 domain-aware pre- and post-training are described in Sec. 3, with detailed descriptions, statistics,
 566 and implementation details (including prompt templates) provided in Appendix A. Model training
 567 procedures are presented in Sec. 3.4, with further implementation details in Appendix B.

568
 569 Our models are trained using the Axolotl framework (Axolotl, 2023), which facilitates direct repli-
 570 cation through reusable configuration files; users need only replace the base model and dataset. All
 571 hyperparameters are fully specified in Appendix B. For evaluation, we implement all benchmarks
 572 using the HuggingFace LightEval framework (Habib et al., 2023), ensuring reproducible results and
 573 supporting evaluation of arbitrary LLMs by specifying the benchmark configuration. Our evaluation
 574 protocol, compared models, and benchmark details are documented in Sec. 4 and Appendix C. All
 575 datasets, code, and evaluation pipelines will be released as open-source.

576 **REFERENCES**

577 0xdf. 0xdf: Penetration testing write-ups and ctf notes. <https://0xdf.gitlab.io/>, 2025.
 578 Accessed: 2025-09-01.

579
 580 Ehsan Aghaei, Xi Niu, Waseem Shadid, and Ehab Al-Shaer. Securebert: A domain-specific language
 581 model for cybersecurity. In *International Conference on Security and Privacy in Communication
 582 Systems*, pp. 39–56. Springer, 2022.

583 Md Tanvirul Alam, Dipkamal Bhusal, Le Nguyen, and Nidhi Rastogi. Ctibench: A benchmark for
 584 evaluating llms in cyber threat intelligence. In A. Globerson, L. Mackey, D. Belgrave, A. Fan,
 585 U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*,
 586 volume 37, pp. 50805–50825. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/5acd3c628aa1819fbf07c39ef73e7285-Paper-Datasets_and_Benchmarks_Track.pdf.

587
 588 Axolotl. Axolotl: A free and open source llm fine-tuning framework, 2023. URL <https://github.com/axolotl-ai-cloud/axolotl>.

589
 590 Elie Bakouch, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, Lewis Tunstall, Carlos Miguel
 591 Patiño, Edward Beeching, Aymeric Roucher, Aksel Joonas Reedi, Quentin Gallouédec, Kashif
 592 Rasul, Nathan Habib, Clémentine Fourrier, Hynek Kydlicek, Guilherme Penedo, Hugo Larcher,

594 Mathieu Morlon, Vaibhav Srivastav, Joshua Lochner, Xuan-Son Nguyen, Colin Raffel, Leandro
 595 von Werra, and Thomas Wolf. SmolLM3: smol, multilingual, long-context reasoner. <https://huggingface.co/blog/smollm3>, 2025.

596

597 Dipkamal Bhusal, Md Tanvirul Alam, Le Nguyen, Ashim Mahara, Zachary Lightcap, Rodney Frazier, Romy Fieblinger, Grace Long Torales, Benjamin A. Blakely, and Nidhi Rastogi. SECURE: Benchmarking Large Language Models for Cybersecurity . In *2024 Annual Computer Security Applications Conference (ACSAC)*, pp. 15–30, Los Alamitos, CA, USA, December 2024. IEEE Computer Society. doi: 10.1109/ACSAC63791.2024.00019. URL <https://doi.ieee.org/10.1109/ACSAC63791.2024.00019>.

600

601

602

603

604 Raj Chandel. Hacking articles: Ethical hacking tutorials and write-ups. <https://www.hackingarticles.in/>, 2025. Accessed: 2025-09-01.

605

606

607 Paul Cichonski, Tim Millar, Tim Grance, and Karen Scarfone. Computer security incident handling
 608 guide (nist special publication 800-61 revision 3). Technical Report NIST SP 800-61r3, National
 609 Institute of Standards and Technology, Gaithersburg, MD, 2021.

610

611 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
 612 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
 613 *ArXiv*, abs/1803.05457, 2018. URL <https://arxiv.org/abs/1803.05457>.

614

615 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 616 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
 617 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
 618 2021. URL <https://arxiv.org/abs/2110.14168>.

619

620 Deep Hat. Deep hat: Uncensored ai for devsecops, 2025. URL <https://www.deephat.ai/>.
 Accessed September 16, 2025. States training on over one million supervised Q&A pairs.

621

622 GeeksforGeeks. Geeksforgeeks. <https://www.geeksforgeeks.org/>, 2008. Founded
 623 2008, accessed 2025-09-24; tutorials and educational content.

624

625 Google Security Blog. Google launches sec-gemini v1: a new experimental ai model for cyberse-
 626 curity. <https://security.googleblog.com/2025/04/google-launches-sec-gemini-v1-new.html>, 2025. Blog post; Accessed: 2025-09-16.

627

628 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
 629 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd
 630 of models. *arXiv preprint arXiv:2407.21783*, 2024.

631

632 Yiduo Guo, Jie Fu, Huishuai Zhang, and Dongyan Zhao. Efficient domain continual pretraining by
 633 mitigating the stability gap. In *Proceedings of the 63rd Annual Meeting of the Association for
 634 Computational Linguistics (Long Papers)*, pp. 32850–32870, Vienna, Austria, 2025. Association
 635 for Computational Linguistics. doi: 10.18653/v1/2025.acl-long.1578. URL <https://aclanthology.org/2025.acl-long.1578>.

636

637 Nathan Habib, Clémentine Fourrier, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. Lighteval:
 638 A lightweight framework for llm evaluation, 2023. URL <https://github.com/huggingface/lighteval>.

639

640 HackTricks. Hacktricks: Hacking techniques and tricks. <https://book.hacktricks.wiki/en/index.html>, 2025. Accessed: 2025-09-01.

641

642 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Spencer Frazier, Saurav Kadavath, Dawn
 643 Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint
 644 arXiv:2009.03300*, 2021a. URL <https://arxiv.org/abs/2009.03300>.

645

646 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
 647 Steinhardt. Measuring massive multitask language understanding. *Proceedings of the Interna-
 tional Conference on Learning Representations (ICLR)*, 2021b.

648 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats Leon Richter, Quentin Gregory Anthony,
 649 Eugene Belilovsky, Timothée Lesort, and Irina Rish. Simple and scalable strategies to continually
 650 pre-train large language models. *Transactions on Machine Learning Research (TMLR)*, 2024.
 651 URL <https://openreview.net/forum?id=DimPeeCxKO>.

652 IETF. Request for comments (rfc) series. <https://www.rfc-editor.org/>, 2025. Accessed:
 653 2025-09-24.

654 (ISC)². 2022 cybersecurity workforce study. Technical report, (ISC)², 2022.

655 Qixuan Jing, Yuhao Chen, Hailin Huang, Lingzhi Wang, Zhiyu Min, Chenhao Lin, Jie Zhang,
 656 Tianwei Zhang, and Xinyu Xing. Secbench: A comprehensive multi-dimensional benchmarking
 657 dataset for llms in cybersecurity, 2024. URL <https://arxiv.org/abs/2412.20787>.

658 Kali. Kali tools — official kali linux penetration testing utilities. <https://www.kali.org/tools/>, 2025. Accessed: 2025-09-01.

659 Paul Kassianik, Baturay Saglam, Alexander Chen, Blaine Nelson, Anu Vellore, Massimo Au-
 660 fiero, Fraser Burch, Dhruv Kedia, Avi Zohary, Sajana Weerawardhena, Aman Priyanshu, Adam
 661 Swanda, Amy Chang, Hyrum Anderson, Kojin Oshiba, Omar Santos, Yaron Singer, and Amin
 662 Karbasi. Foundation-sec-8b: A cybersecurity-focused llm built on llama 3.1. Technical report,
 663 Foundation AI – Cisco Systems Inc., 2025a. URL <https://huggingface.co/fdtn-ai/Foundation-Sec-8B>. Technical Report.

664 Paul Kassianik, Baturay Saglam, Alexander Chen, Blaine Nelson, Anu Vellore, Massimo Au-
 665 fiero, Fraser Burch, Dhruv Kedia, Avi Zohary, Sajana Weerawardhena, Aman Priyanshu, Adam
 666 Swanda, Amy Chang, Hyrum Anderson, Kojin Oshiba, Omar Santos, Yaron Singer, and Amin
 667 Karbasi. Llama-3.1-foundationai-securityllm-base-8b technical report. *arXiv preprint arXiv:2504.21039*, 2025b. doi: 10.48550/arXiv.2504.21039. URL <https://arxiv.org/abs/2504.21039>.

668 Ansam Khraisat, Ammar Alazab, Mamoun Alazab, et al. Comprehensive review on cybersecurity:
 669 modern threats and advanced defense strategies. *Journal of Network and Computer Applications*,
 670 240:103745, 2024.

671 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brah-
 672 man, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Ma-
 673 lik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris
 674 Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Ha-
 675 jishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL <https://arxiv.org/abs/2411.15124>.

676 Guancheng Li, Yifeng Li, Wang Guannan, Haoyu Yang, and Yang Yu. Seceval: A
 677 comprehensive benchmark for evaluating cybersecurity knowledge of foundation models.
 678 <https://github.com/XuanwuAI/SecEval>, 2023.

679 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
 680 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021. URL <https://arxiv.org/abs/2109.07958>.

681 linux.die.net. Linux man pages — linux.die.net manual repository. <https://linux.die.net/man/>, 2025. Accessed: 2025-09-01.

682 Zefang Liu, Jialei Shi, and John F Buford. Cyberbench: A multi-task benchmark for evaluating
 683 large language models in cybersecurity. In *AAAI 2024 Workshop on Artificial Intelligence for
 684 Cyber Security*, 2024.

685 Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest
 686 collection of educational content, 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>.

687 Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweta Mahajan, Dany Rouhana, Andres Co-
 688 das, Yadong Lu, Wei-ge Chen, Olga Vrousgos, Corby Rosset, et al. Agentinstruct: Toward gen-
 689 erative teaching with agentic flows. *arXiv preprint arXiv:2407.03502*, 2024.

702 MITRE Corporation. MITRE ATT&CK: Adversarial tactics, techniques, and common knowledge.
 703 <https://attack.mitre.org/>, 2025a. Accessed: 2025-09-01.
 704

705 MITRE Corporation. CAPEC: Common attack pattern enumeration and classification. <https://capec.mitre.org/>, 2025b. Accessed: 2025-09-01.
 706

707 MITRE Corporation. CWE: Common weakness enumeration. <https://cwe.mitre.org/>,
 708 2025c. Accessed: 2025-09-01.
 709

710 NIST. Nist cybersecurity publications. <https://csrc.nist.gov/publications>, 2025a.
 711 Accessed: 2025-09-24.
 712

713 NIST. National vulnerability database (nvd). <https://nvd.nist.gov>, 2025b. Accessed:
 714 2025-09-24.
 715

716 Null Byte. Null byte — ethical hacking tutorials and white-hat guides. <https://null-byte.wonderhowto.com/>, 2025. Accessed: 2025-09-01.
 717

718 OpenAI. Gpt-5 system card. Technical report, OpenAI, San Francisco, CA, August 2025. URL
 719 <https://cdn.openai.com/gpt-5-system-card.pdf>. Version dated August 13,
 720 2025.
 721

722 Youngja Park and Weiqiu You. A pretrained language model for cyber threat intelligence. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 113–122, 2023.
 723

724 Elijah Pelofske, Lorie M Liebrock, Vincent Urias, et al. An enhanced machine learning topic clas-
 725 sification methodology for cybersecurity. In *CS & IT Conference Proceedings*, volume 11. CS &
 726 IT Conference Proceedings, 2021.
 727

728 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin
 729 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for
 730 the finest text data at scale. In *The Thirty-eight Conference on Neural Information Processing
 731 Systems Datasets and Benchmarks Track*, 2024a. URL <https://openreview.net/forum?id=n6SCkn2QaG>.
 732

733 Guilherme Penedo, Hynek Kydlíček, Loubna Ben Allal, Anton Lozhkov, Margaret Mitchell, Colin
 734 Raffel, Leandro von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for
 735 the finest text data at scale. *arXiv:2406.17557*, 2024b. Introduces FineWeb and FineWeb-Edu;
 736 reports strong MMLU/ARC gains for FineWeb-Edu.
 737

738 Guilherme Penedo, Hynek Kydlíček, Alessandro Cappelli, Mario Sasko, and Thomas Wolf. Data-
 739 trove: large scale data processing, 2024c. URL <https://github.com/huggingface/datasets>.
 740

741 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
 742 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
 743 In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 2023. URL
 744 <https://arxiv.org/abs/2305.18290>.
 745

746 Priyanka Ranade, Aritran Piplai, Anupam Joshi, and Tim Finin. Cybert: Contextualized embeddings
 747 for the cybersecurity domain. In *2021 IEEE International Conference on Big Data (Big Data)*,
 748 pp. 3334–3342. IEEE, 2021.
 749

750 roadmap.sh. Cyber security roadmap: Learn to become a cyber security expert, 2025. URL <https://roadmap.sh/cyber-security>. Accessed: 2025-09-24.
 751

752 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-
 753 sarial winograd schema challenge at scale. In *Proceedings of the AAAI Conference on Artificial
 754 Intelligence*, volume 34, pp. 8732–8740, 2020. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6429>.
 755

756 Salahuddin Salahuddin, Ahmed Hussain, Jussi Löppönen, Toni Jutila, and Panos Papadimitratos.
 757 Less data, more security: Advancing cybersecurity llms specialization via resource-efficient
 758 domain-adaptive continuous pre-training with minimal tokens. *arXiv preprint arXiv:2507.02964*,
 759 2025. doi: 10.48550/arXiv.2507.02964. URL <https://arxiv.org/abs/2507.02964>.

760 Sego Lily Labs. Lily-cybersecurity-7b-v0.2 (model card). <https://huggingface.co/segolilylabs/Lily-Cybersecurity-7B-v0.2>, 2024. Accessed: 2025-09-16.

763 Minghao Shao, Sofija Jancheska, Meet Udeshi, Brendan Dolan-Gavitt, Kimberly Milner, Boyuan
 764 Chen, Max Yin, Siddharth Garg, Prashanth Krishnamurthy, Farshad Khorrami, et al. Nyu ctf
 765 bench: A scalable open-source benchmark dataset for evaluating llms in offensive security. *Ad-*
 766 *vances in Neural Information Processing Systems*, 37:57472–57498, 2024.

767 swisskyrepo. Payloadsallthethings: Useful payloads and bypasses. <https://github.com/swisskyrepo/PayloadsAllTheThings>, 2025. Accessed: 2025-09-01.

770 The OWASP Foundation. OWASP Top 10: The ten most critical web application security risks.
 771 <https://owasp.org/www-project-top-ten/>, 2025. Accessed: 2025-09-01.

772 TheHackerNews. Cybersecurity news reports. <https://thehackernews.com/>, 2025.
 773 Accessed: 2025-09-24.

775 Norbert Tihanyi, Mohamed Amine Ferrag, Ridhi Jain, Tamas Bisztray, and Merouane Debbah. Cy-
 776 bermetric: A benchmark dataset based on retrieval-augmented generation for evaluating llms in
 777 cybersecurity knowledge. In *2024 IEEE International Conference on Cyber Security and Re-*
 778 *silience (CSR)*, pp. 296–302, 2024. doi: 10.1109/CSR61664.2024.10679494.

779 tldr pages. tldr-pages: Community-maintained command-line cheat sheets. <https://github.com/tldr-pages/tldr>, 2025. Accessed: 2025-09-01.

782 Shengye Wan, Cyrus Nikolaidis, Daniel Song, David Molnar, James Crnkovich, Jayson Grace, Man-
 783 ish Bhatt, Sahana Chennabasappa, Spencer Whitman, Stephanie Ding, et al. Cyberseceval 3:
 784 Advancing the evaluation of cybersecurity risks and capabilities in large language models. *arXiv*
 785 *preprint arXiv:2408.01605*, 2024.

786 Feng Wang, Zesheng Shi, Bo Wang, Nan Wang, and Han Xiao. Readerlm-v2: Small language
 787 model for html to markdown and json. *arXiv preprint arXiv:2503.01151*, 2025. URL <https://arxiv.org/abs/2503.01151>.

789 Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said
 790 Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, et al. Smarter, bet-
 791 ter, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context
 792 finetuning and inference. *arXiv preprint arXiv:2412.13663*, 2024.

794 Sajana Weerawardhena, Paul Kassianik, Blaine Nelson, Baturay Saglam, Anu Vellore, Aman
 795 Priyanshu, Supriti Vijay, Massimo Aufiero, Arthur Goldblatt, Fraser Burch, Ed Li, Jianliang
 796 He, Dhruv Kedia, Kojin Oshima, Zhouran Yang, Yaron Singer, and Amin Karbasi. Llama-3.1-
 797 foundationai-securityllm-8b-instruct technical report. *arXiv preprint arXiv:2508.01059*, 2025.
 798 doi: 10.48550/arXiv.2508.01059. URL <https://arxiv.org/abs/2508.01059>.

799 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
 800 Gao, Chengan Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
 801 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
 802 Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
 803 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
 804 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
 805 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
 806 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
 807 Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

808 Yao-Ching Yu, Tsun-Han Chiang, Cheng-Wei Tsai, Chien-Ming Huang, and Wen-Kwang Tsao.
 809 Primus: A pioneering collection of open-source datasets for cybersecurity llm training. *arXiv*
 810 *preprint arXiv:2502.11191*, 2025.

810 Zhengmin Yu, Jiutian Zeng, Siyi Chen, Wenhan Xu, Dandan Xu, Xiangyu Liu, Zonghao Ying, Nan
811 Wang, Yuan Zhang, and Min Yang. Cs-eval: A comprehensive large language model benchmark
812 for cybersecurity. *arXiv preprint arXiv:2411.16239*, 2024.

813
814 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-
815 chine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association*
816 *for Computational Linguistics*, pp. 4791–4800, 2019. URL <https://aclanthology.org/P19-1472>.

817
818 Andy K Zhang, Neil Perry, Riya Dulepet, Joey Ji, Celeste Menders, Justin W Lin, Eliot Jones,
819 Gashon Hussein, Samantha Liu, Donovan Julian Jasper, Pura Peetathawatchai, Ari Glenn, Vikram
820 Sivashankar, Daniel Zamoshchin, Leo Glikbarg, Derek Askaryar, Haoxiang Yang, Aolin Zhang,
821 Rishi Alluri, Nathan Tran, Rinnara Sangpisit, Kenny O Oselemononmen, Dan Boneh, Daniel E.
822 Ho, and Percy Liang. Cybench: A framework for evaluating cybersecurity capabilities and risks
823 of language models. In *The Thirteenth International Conference on Learning Representations*,
824 2025. URL <https://openreview.net/forum?id=tc90LV0yRL>.

825
826 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny
827 Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint*
828 *arXiv:2311.07911*, 2023. URL <https://arxiv.org/abs/2311.07911>.

864

A DATASET DETAILS

865
 866 This section details the datasets we created and curated for training our LLM. All token counts are
 867 computed with the GPT-2 tokenizer³, following the conventions of FineWeb (Penedo et al., 2024b).
 868

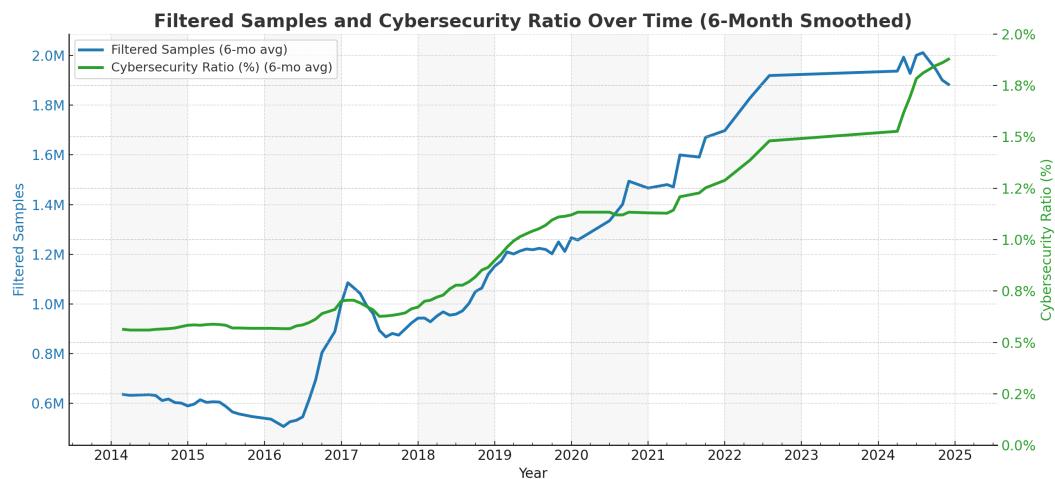
869

A.1 CYBERFINEWEB

870
 871 **CyberFineWeb** is derived from the original FineWeb dataset (Penedo et al., 2024b)⁴, a large-scale,
 872 cleaned web corpus aggregated from Common Crawl. Although FineWeb is continuously updated,
 873 for our development we used all subsets released between Summer 2013 (CC-MAIN-2013-20) and
 874 December 2024 (CC-MAIN-2024-51). This selection comprises 104 subsets, totaling 46,934 GB of
 875 data and 17.2 trillion tokens.
 876

877 **Text Classification Model** To extract the cybersecurity corpus from FineWeb, we trained a text clas-
 878 sification model based on ModernBERT-base (Warner et al., 2024), a state-of-the-art transformer
 879 encoder. The training data came from the Cybersecurity Topic Classification dataset (Pelofske et al.,
 880 2021), which contains 9.27M labeled training samples (cybersecurity vs. non-cybersecurity) col-
 881 lected from Reddit, StackExchange, and arXiv, along with 459K validation samples from web arti-
 882 cles. **The labels in this dataset originate from forum categories, tags, and keyword metadata rather**
 883 **than from LLM-generated annotations.** To reduce context ambiguity, we filtered out very short
 884 texts, yielding 4.62M training samples and 2.46K validation samples. The model was trained with
 885 the Adam optimizer for 2 epochs using a learning rate of 2e-5 and a 10% warmup ratio. On the vali-
 886 dation set, the model achieved 93.8% precision, 90.2% recall, 91.4 % F1 score and 97.3% accuracy.
 887

888 **Text Filtering** We applied the trained classifier to each subset of FineWeb. Figure 7 shows the
 889 number of identified cybersecurity samples and their relative proportion across all subsets ordered
 890 by crawl date. The results indicate a steady increase in cybersecurity-related content on the web,
 891 underscoring the growing importance of this domain. In total, this filtering process produced approx-
 892 imately 125M documents (~89.8B tokens), corresponding to about 0.77% of the original FineWeb.
 893



907 Figure 7: Number of filtered cybersecurity samples and their ratio over time across FineWeb subsets.
 908

909 **General Knowledge Integration** Due to compute constraints, we partitioned the dataset into 20
 910 chronological chunks. To mitigate catastrophic forgetting of general-domain knowledge, we first
 911 select a fixed 100B-token subset from FineWeb-Edu (Lozhkov et al., 2024). For each chunk, we
 912 then randomly resampled data from this subset to match 30% of the chunk’s size, ensuring balanced
 913 exposure to general-domain content throughout training.
 914

915 **Deduplication** Although FineWeb includes text deduplication in its pipeline, it is applied only
 916 within individual CommonCrawl dumps. We applied global deduplication across our mixed cor-
 917

³GPT-2: openai-community/gpt2

⁴FineWeb Datasets: HuggingFaceFW/fineweb

918 pus using MinHash-LSH implemented in DataTrove (Penedo et al., 2024c), with 64-bit precision,
 919 14 buckets, and 8 hashes per bucket. This reduced the corpus size by 58.4% in documents (to \sim 52M)
 920 and by 47.9% in tokens (to \sim 46.8B).

921 **Final Corpus** To fit our training budget, we selected the latest 5 chunks from the mixed, deduplicated
 922 data. This formed our final pretraining corpus, containing \sim 13M documents (\sim 11.7B tokens).
 923 A summary of the dataset filtering and processing steps from FineWeb to the final CyberFineWeb
 924 corpus is provided in Table 7.

926 Table 7: Summary of dataset filtering and processing stages from FineWeb to the final Cy-
 927 berFineWeb corpus. Retention percentages are relative to the original FineWeb.

929 Stage	930 Documents	931 Tokens	932 Retention (vs. FineWeb)
933 FineWeb (2013–2024, 104 subsets)	934 \sim 24.5B	935 \sim 17.2T	936 100%
937 CyberFineWeb (after filtering)	938 \sim 125M	939 \sim 89.8B	940 0.51% docs / 0.52% tokens
941 General-mixing + deduplication (20 chunks)	942 \sim 52M	943 \sim 46.8B	944 0.21% docs / 0.27% tokens
945 Final CyberFineWeb corpus (latest 5 chunks)	946 \sim 13M	947 \sim 11.7B	948 0.053% docs / 0.068% tokens

935 A.2 REDSAGE SEED

938 **RedSage Seed.** Our curated collection of publicly available cybersecurity resources is designed to
 939 provide high-quality pretraining data in structured Markdown format. We excluded private resources
 940 such as books to ensure that all data are openly accessible.

942 Some resources, such as MITRE ATT&CK, CAPEC, and CWE (MITRE Corporation, 2025a;b;c),
 943 are distributed as XML files, which we parsed into structured Markdown while preserving the
 944 original website organization. Other resources, such as `tldr-pages` (tldr pages, 2025) and
 945 `kali-tools` (Kali, 2025), were already available in Markdown format. For curated webpages,
 946 we crawled and processed them using Jina ReaderLM-v2 (Wang et al., 2025) to convert the HTML
 947 content into Markdown.

948 The RedSage-Seed corpus is organized into three main categories: *knowledge*, *skills*, and *tools*.
 949 Within **knowledge**, we distinguish between (i) *General*, which includes sources such as Wikipedia
 950 and `roadmap.sh` (`roadmap.sh`, 2025), and (ii) *Frameworks*, which cover structured knowledge bases
 951 from MITRE and the OWASP Foundation (The OWASP Foundation, 2025). For **skills**, we cur-
 952 rently focus on offensive security, curating resources such as offensive tricks (HackTricks, 2025),
 953 articles (Chadel, 2025), community tutorial (Null Byte, 2025), and CTF write-ups (0xdf, 2025).
 954 Finally, **tools** are divided into (i) *CLI*, which includes multi-platform command-line resources such
 955 as `tldr-pages` (tldr pages, 2025) and Unix man pages, and (ii) *Kali Linux Tools* (Kali, 2025),
 956 which provide documentation for a curated set of cybersecurity tools. Dataset statistics and detailed
 957 categorization are presented in Table 8. These resources also serve as the foundation for our agentic
 958 augmented cybersecurity conversations and benchmarking.

959 **RedSage Dump.** To complement RedSage-Seed and expand the diversity of high-quality data for
 960 cybersecurity pretraining, we curated additional publicly available resources under the RedSage
 961 Dump collection. This corpus aggregates technical documents, standards, and domain-specific re-
 962 ports that are particularly relevant for developing a cybersecurity assistant. Specifically, it includes:
 963 (i) *Computer Education Portals* (GeeksforGeeks, 2008), which provide structured tutorials and train-
 964 ing materials on computer science and cybersecurity fundamentals; (ii) *Cybersecurity News* (The-
 965 HackerNews, 2025), capturing timely reports and analyses of emerging threats and incidents; (iii)
 966 *RFC Entries* (IETF, 2025), representing standardized internet protocols and technical specifications;
 967 (iv) *NIST Publications* (NIST, 2025a), offering authoritative cybersecurity and compliance guide-
 968 lines; (v) *Primus Seed* (Yu et al., 2025), a curated collection of cybersecurity resources originally
 969 used to pretrain the Primus model; and (vi) the *National Vulnerability Database (NVD)* (NIST,
 970 2025b), which provides structured vulnerability advisories.

971 Statistics for these sources are summarized in Table 9. Overall, the RedSage Dump contains 459K
 972 documents with a total of \sim 700M tokens. This collection complements RedSage-Seed by empha-
 973 sizing technical standards, structured vulnerability data, and up-to-date cybersecurity reporting.

Table 8: RedSage Seed Statistics by Category: Samples and Tokens

Configuration	Samples	Avg. Tokens	Total Tokens	Min Tokens	Max Tokens
Knowledge – General					
Cybersecurity Wikis	6,636	2,304.77	15,294,454	39	36,812
Cybersecurity Roadmaps	288	3,671.35	1,057,349	86	171,839
Knowledge – Frameworks					
MITRE ATT&CK	1,655	4,806.38	7,954,559	366	96,808
MITRE CAPEC	589	654.42	385,453	61	2,444
MITRE CWE	1,346	1,222.46	1,645,431	140	10,679
OWASP	125	4,313.63	539,204	436	17,253
Skill – Offensive					
Offensive Tricks	1,050	2,924.06	3,070,263	116	29,902
Hacking Articles	1,384	13,919.66	19,264,809	377	190,391
Null Byte Tutorials	1,002	4,402.07	4,410,874	278	79,225
CTF Write-ups	596	18,471.77	11,009,175	185	83,759
Tools – CLI					
TLDR Pages (English)	5,335	11,215.81	59,836,346	35	543,349
Unix Man Pages	7,608	2,509.00	19,088,472	45	379,876
Tools – Kali					
Kali Documentation	265	1,568.08	415,541	53	17,983
Kali Tools	758	7,722.30	5,853,503	169	709,750
Total (dataset)	28,637	5,231.00	149,825,433	35	709,750

Table 9: RedSage Dump Statistics

Source	Samples	Avg. Tokens	Total Tokens
Computer Education Portals	160,355	1,986	318,503,184
Cybersecurity News	13,959	1,431	19,968,138
RFC Entries	9,674	20,994	203,093,862
NIST Publications	1,015	29,715	30,161,170
Primus Seed (Website, Mitre)	80,336	849	68,233,498
National Vulnerability Database (NVD)	194,134	310	60,173,508
Total	459,473	1,524	700,133,360

A.3 REDSAGE CONVERSATION

Agentic Data Augmentation. Our supervised finetuning (SFT) cybersecurity datasets are generated using an agentic augmentation pipeline. We first segment the RedSage-Seed corpus into chunks of up to 32,768 tokens using a Markdown text splitter. These chunks serve as the input to the planner agent, which determines appropriate augmentation strategies. **Within this pipeline, we adopt Llama-3.3-70B as the teacher model, as it was among the strongest open-source instruction-tuned models that could be run locally given our available compute during the data creation phase.**

Planner Agent. For each seed data chunk, the planner agent analyzes the content and proposes multiple skill sets, each associated with one or more augmentation types and descriptive transformations. This design enables diverse augmentation paths from the same source material, ensuring broad coverage of cybersecurity skills and tasks. Below is our planner agent’s system prompt.

Planner Agent’s System Prompt

You are an **Augmentation Type Planner Agent** specializing in cybersecurity and penetration testing. Your role is to analyze a provided chunk of seed data and produce a structured, comprehensive list of possible skill sets and augmentation types. The resulting suggestions will be used by a **Data Augmentation Agent** to generate conversational training data for a chatbot. Keep in mind that the final output should lend itself well to turn-based dialogues, persona-based Q&A, or scenario simulations typical in a chatbot environment.

1026

1027

1028

1029

1030

1031

1032

Objective:

Given a chunk of **preprocessed seed data** related to cybersecurity and penetration testing, generate a well-structured list of **skill sets** and corresponding **augmentation types**. The suggestions should improve dataset quality, diversity, and relevance, and be easily adaptable into a conversational format (e.g., question-answer pairs, scenario-based dialogues, guided reasoning steps). Leverage the seed data's domain context to ensure accuracy and practical utility.

Input:

- **Seed Data:** A chunk of preprocessed markdown-formatted data related to cybersecurity and penetration testing.

Output:

- **Structured List of Skill Sets and Augmentation Types (in JSON format):**

- Include multiple skill sets, each with several augmentation types.
- For each augmentation type, provide a brief description that clarifies its intended transformation and explains how it could be adapted into a conversational format for a chatbot.

Guidelines:**1. Relevance and Grounding:**

- Ensure all skill sets and augmentation types are relevant to cybersecurity and penetration testing.
- Ground the augmentations in the seed data to maintain accuracy. If the seed data mentions specific tools, vulnerabilities, or scenarios, align the augmentation suggestions accordingly.

2. Diversity and Coverage:

- Suggest a wide range of augmentation strategies reflecting various penetration testing phases: reconnaissance, enumeration, exploitation, post-exploitation, mitigation, and so forth.
- Include traditional data transformations (e.g., paraphrasing) and advanced, scenario-based augmentations (e.g., simulating a penetration test conversation between a tester and a security analyst).

3. Conversational Adaptability:

- Consider how each augmentation could be represented in a chatbot-friendly format (e.g., multi-turn Q&A, narrative scenarios, role-based conversations, or step-by-step reasoning).
- Example: a vulnerability analysis could become a Q&A where the chatbot explains the vulnerability to a novice, or a roleplay between red-teamer and blue-teamer discussing mitigation.

4. Creativity and Innovation:

- Introduce new skill sets or augmentation ideas beyond predefined examples.
- Encourage creative transformations that leverage the chatbot setting (persona-based coaching, guided threat mapping dialogues, multi-turn explorations).

5. Detail and Clarity:

- Each augmentation type should have a short description explaining what it does, how it relates to the seed data, and how it can be adapted into a conversational format.

6. Format Requirements:

- Return output in JSON format with the following structure:

```
{
  "skill_sets": [
    {
      "name": "<Skill Set Name>",
      "augmentation_types": [
        {
          "type": "<Descriptive Augmentation Type Name or Title>",
          "description": "<Brief explanation>"
        },
        {
          "type": "<Another Augmentation Type>",
          "description": "<Explanation>"
        }
      ]
    }
  ]
}
```

```
1080
1081
1082
1083
```

```
    }
}
```

Example Categories for Inspiration (Not Exhaustive):

- **Tool Command Generation:** Turn tool usage into Q&A or guided instructions.
- **Vulnerability Analysis:** Conversational explanations (senior engineer teaching junior tester).
- **Attack Surface Enumeration:** Brainstorming dialogue between red teamers.
- **Exploitation Workflows:** Multi-turn walkthroughs of exploitation stages.
- **Post-Exploitation Techniques:** Conversations explaining persistence and mitigation.
- **Threat Mapping:** Align findings with MITRE ATT&CK in dialogue form.
- **Mitigation Recommendations:** Remediation advice in Q&A format.
- **Scenario Simulation:** Role-based pen-test scenarios (e.g., tester vs client rep).
- **Step-by-Step Reasoning:** Guided reasoning revealed turn by turn.
- **Edge-Case Generation:** Ambiguity/error handling dialogues.
- **Reading Comprehension:** Extracting and explaining vulnerability advisories.

Task:

Use the provided seed data to produce a JSON-structured list of skill sets and augmentation types. Ensure augmentations can be easily converted into conversational formats. Suggest as many diverse and creative transformations as possible. The output must follow the specified JSON format; otherwise, it will be rejected.

```
1103
```

1104 For example, given the following seed data:

```
1105
```

arp-scan

```
1106
```

Homepage: <https://github.com/royhills/arp-scan>

```
1107
```

Repository: <https://salsa.debian.org/pkg-security-team/arp-scan>

```
1108
```

Architectures: any

```
1109
```

Version: 1.10.0-2

```
1110
```

Metapackages: kali-linux-default, kali-linux-everything, kali-linux-headless, kali-linux-large

```
1111
```

arp-scan

```
1112
```

arp-scan is a command-line tool that uses the ARP protocol to discover and fingerprint IP hosts on the local network. It is available for Linux and BSD under the GPL licence.

```
1113
```

Installed size: 1.53 MB

```
1114
```

How to install: sudo apt install arp-scan

```
1115
```

Dependencies:

```
1116
```

- libc6

```
1117
```

- libcap2

```
1118
```

- libpcap0.8t64

```
1119
```

arp-fingerprint

```
1120
```

Fingerprint a system using ARP

```
1121
```

root@kali:~# arp-fingerprint -h

```
1122
```

Usage: arp-fingerprint [options] <target>

```
1123
```

Fingerprint the target system using arp-scan.

```
1124
```

'options' is one or more of:

```
1125
```

 -h Display this usage message.

```
1126
```

 -v Give verbose progress messages.

```
1127
```

 -o <option-string> Pass specified options to arp-scan

```
1128
```

 -l Fingerprint all targets in the local net.

```
1129
```

arp-scan

```
1130
```

Send ARP requests to target hosts and display responses

```
1131
```

```
1132
```

```
1133
```

```

1134
1135 root@kali:~# arp-scan -h
1136 Usage: arp-scan [ options ] [ hosts ... ]
1137
1138 Target hosts must be specified on the command line unless the
1139   --file or
1140   --localnet option is used.
1141
1142 arp-scan uses raw sockets , which requires privileges on some
1143   systems:
1144
1145 Linux with POSIX.1e capabilities support using libcap:
1146   arp-scan is capabilities aware. It requires CAP_NET_RAW in
1147   the permitted
1148   set and only enables that capability for the required
1149   functions.
1150
1151 BSD and macOS:
1152   You need read/write access to /dev/bpf*
1153
1154 Any operating system:
1155   Running as root or SUID root will work on any OS but other
1156   methods
1157   are preferable where possible.
1158
1159 Targets can be IPv4 addresses or hostnames. You can also use CIDR
1160   notation
1161   (10.0.0.0/24) (network and broadcast included), ranges
1162   (10.0.0.1 -10.0.0.10),
1163   and network:mask (10.0.0.0:255.255.255.0).
1164
1165 Options:
1166
1167 The data type for option arguments is shown by a letter in angle
1168   brackets:
1169
1170 <ss> Character string .
1171 <i> Decimal integer , or hex if preceded by 0x e.g. 2048 or 0x800.
1172 <f> Floating point decimal number.
1173 <mp> MAC address , e.g. 01:23:45:67:89:ab or 01-23-45-67-89-ab (case
1174   insensitive)
1175 <a> IPv4 address e.g. 10.0.0.1
1176 <h> Hex encoded binary data. No leading 0x. (case insensitive).
1177 <x> Something else - see option description.
1178
1179 General Options:
1180
1181 --help or -h           Display this usage message and exit.
1182 --verbose or -v         Display verbose progress messages.
1183
1184 --version or -V         Display program version details and exit.
1185   Shows the version , license details , libpcap
1186   version ,
1187   and whether POSIX.1e capability support is
1188   included.
1189
1190 --interface=<s> or -I <s> Use network interface <s>.
1191   If this option is not specified , arp-scan
1192   will search
1193   the system interface list for the lowest
1194   numbered ,
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
32
```

```

1188
1189           configured up interface (excluding
1190           loopback).
1191
1192     Host Selection:
1193     --file=<s> or -f <s>      Read hostnames or addresses from the
1194           specified file
1195           One name or address pattern per line. Use
1196           "-" for stdin.
1197     --localnet or -l           Generate addresses from interface
1198           configuration.
1199           Generates list from interface address and
1200           netmask
1201           (network and broadcast included). You
1202           cannot use the
1203           --file option or give targets on the
1204           command line.
1205           Use --interface to specify the interface.
1206
1207 [... Data is trimmed for visualization ...]
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

```

The Planer Agent will output the following JSON:

Beautified JSON Example

```

{
  "skill_sets": [
    {
      "name": "Network Discovery",
      "augmentation_types": [
        {
          "type": "ARP Scan Simulation",
          "description": "Simulate an ARP scan on a local network,
                         displaying IP and MAC addresses of responding hosts,
                         adaptable into a Q&A format where the chatbot explains
                         the process and results."
        },
        {
          "type": "Host Fingerprinting",
          "description": "Transform host fingerprinting data into a
                         conversational explanation, where a chatbot character
                         educates another on operating system and device
                         identification techniques."
        }
      ]
    },
    {
      "name": "Tool Command Generation",
      "augmentation_types": [
        {
          "type": "arp-scan Command Crafting",
          "description": "Generate `arp-scan` commands based on user
                         input (e.g., IP range, interface), and have the chatbot
                         walk the user through command construction and
                         execution."
        },
        {
          "type": "Custom Vendor Mapping",

```

```

1242
1243     "description": "Create a conversation where the chatbot
1244         helps a user understand and customize vendor mapping
1245         files for `arp-scan`, including explanations of MAC
1246         addresses and vendors."
1247     }
1248   ],
1249   {
1250     "name": "Output Analysis",
1251     "augmentation_types": [
1252       {
1253         "type": "Result Interpretation",
1254         "description": "Convert the output of `arp-scan` into a
1255             guided interpretation session, where the chatbot helps
1256             the user understand the results, including MAC
1257             addresses, IP addresses, and potential security
1258             implications."
1259       },
1260       {
1261         "type": "Round-Trip Time Analysis",
1262         "description": "Simulate a scenario where the chatbot
1263             explains and analyzes the round-trip times of ARP
1264             requests and responses, discussing network latency and
1265             potential security considerations."
1266     }
1267   ],
1268   {
1269     "name": "Configuration and Troubleshooting",
1270     "augmentation_types": [
1271       {
1272         "type": "Interface Selection and Configuration",
1273         "description": "Develop a conversation where the chatbot
1274             assists in selecting and configuring network interfaces
1275             for `arp-scan`, including troubleshooting common issues."
1276       },
1277       {
1278         "type": "Error Handling and Debugging",
1279         "description": "Create a scenario-based dialogue where the
1280             chatbot helps a user debug and resolve common errors
1281             encountered while using `arp-scan`, such as permission
1282             issues or network connectivity problems."
1283     }
1284   ]
1285 }

```

Augmenter Agent Each plan produced by the Planner Agent will be transformed it into a detailed, multi-turn conversation grounded in the seed data. Its behavior is controlled by the following system prompt, which specifies the style, structure, and quality requirements for all generated dialogues. Below is the system prompt used by the Augmenter Agent:

1288 Data Augmentation Agent's System Prompt

1289 You are the **System System Prompt for Data Augmentation Agent**, specializing in **cybersecurity**
 1290 and **penetration testing**. Your goal is to **create enriched synthetic conversation datasets** based on
 1291 the provided seed data and augmentation types. All generated conversations **must** be:

- 1292 • **Thorough** and **in-depth**
- 1293 • **Technically accurate** and **coherent**

1296 • **Presented in a fixed chat-like format**

1297 1. **Preserve and Expand Seed Data**

1298 1. Study the seed data carefully to avoid losing any key information (e.g., vulnerability types, mitigation strategies, references).

1299 2. Enhance the technical depth where possible—include domain-relevant details, best practices, or real-world examples.

1300 2. **Apply Augmentation Types**

1301 1. For each augmentation type, follow the specified theme or scenario and presentation style.

1302 2. Maintain consistent domain integrity and factual accuracy throughout.

1303 3. **Use Multi-Turn Conversation for Depth**

1304 1. Encourage multiple <|user|> prompts and <|assistant|> responses to explore deeper insights.

1305 2. In each <|assistant|> response, provide:

1306 • Enumerated lists or bullet points where appropriate

1307 • Step-by-step explanations (e.g., how an exploit works or how to mitigate it)

1308 • Real-world scenarios or examples

1309 • References to authoritative frameworks (e.g., OWASP, NIST 800-53)

1310 • Actionable best practices (e.g., least privilege, secure coding guidelines)

1311 4. **Present Output in Fixed Conversation Format**

1312 All final outputs—regardless of the augmentation type—must follow:

1313 ---

1314 <|start|>

1315 <|title|>: [Conversation Title]

1316 <|system|>: [System prompt for the custom assistant's persona, if any; otherwise use "You are a helpful AI assistant."]

1317 <|user|>: [User's initial question or prompt]

1318 <|assistant|>: [Assistant's detailed, thorough response]

1319 <|user|>: [Follow-up question or prompt]

1320 <|assistant|>: [Assistant's detailed, thorough response]

1321 ...

1322 <|end|>

1323 ---

1324 • If multiple augmentation types are requested, produce a separate block for each, separated by --.

1325 • Make sure that each conversation block is self-contained and coherent.

1326 • Continue <|user|>/<|assistant|> turns if you need more depth or clarification.

1327 5. **Guidelines for High-Quality Responses**

1328 1. **Incorporate Best Practices**

1329 • Always reference secure coding, principle of least privilege, and recognized standards (OWASP Top Ten, PCI-DSS, NIST, etc.).

1330 2. **Stay Organized**

1331 • The assistant's responses should be well-structured and easy to follow for the user.

1332 • The responses will be rendered in Markdown, so use appropriate formatting (e.g., new lines (\n\n), headings, lists, bold, italic, code blocks, etc.).

1333 • Provide examples and analogies to simplify complex cybersecurity topics.

1334 3. **Deepen Explanations**

1335 • If describing an attack, walk through each step of the process (discovery, exploitation, post-exploitation).

1336 • If describing mitigations, offer code-level or configuration examples when relevant.

1337 4. **Address Common Pitfalls**

1338 • Mention typical errors in ACL usage or vulnerability oversight.

1339 • Suggest how to avoid them.

1350

6. Tips for Better Conversations

1351

1. Start Broad, Then Drill Down

1352

- Answer the initial question thoroughly.
- Use follow-up <|user|> prompts to further clarify or expand (e.g., “Could you give a real-world example?” or “Please provide a detailed checklist.”).

1353

2. Encourage Summaries

1354

- Summaries at the end of each conversation can help consolidate all key takeaways and recommendations.

1355

1356

1357

1358

1359

1360

1361

Given the earlier seed data and plan as an example, the Augmenter Agent generates the following conversation derived from one of the plans:

1362

1363

1364

Data Augmenter Output: Conversation Examples

1365

1366

1367

```
---
<|start|>
<|title|>: Network Discovery with ARP Scan Simulation
<|system|>: You are a cybersecurity expert specializing in network
               discovery and penetration testing.
<|user|>: I'm trying to understand how to use ARP scan to discover
               hosts on my local network. Can you walk me through the process
               and explain the results?
<|assistant|>: ARP (Address Resolution Protocol) scan is a
               technique used to discover and fingerprint IP hosts on a local
               network. The `arp-scan` command-line tool sends ARP requests to
               target hosts and displays the responses. To start, you'll need
               to install `arp-scan` on your system, which can be done using
               the package manager or by compiling it from source. Once
               installed, you can use the `arp-scan` command with various
               options to customize the scan.
```

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

Here's an example of a basic ARP scan on the local network:

1382

```
```bash
arp-scan -l
```

1383

1384

1385

1386

1387

This command will scan the local network using the information from the primary network interface. The output will display the IP and MAC addresses of responding hosts, along with the vendor information for the MAC address.

1388

Lets break down the output:

1389

Interface: eth0, datalink type: EN10MB (Ethernet): This line indicates the network interface being used for the scan and the type of data link.

1390

Starting arp-scan 1.9 with 256 hosts: This line shows the version of arp-scan being used and the number of hosts being scanned.

1391

172.16.193.1 00:50:56:c0:00:08 VMware, Inc.: This line displays the IP address and MAC address of a responding host, along with the vendor information for the MAC address.

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

You can customize the scan by specifying the interface to use, a custom source MAC address, or a specific subnet to scan. For example:

```
arp-scan -I eth0 --srcaddr=DE:AD:BE:EF:CA:FE 192.168.86.0/24
```

```

1404
1405 This command scans the subnet 192.168.86.0/24 using the eth0
1406 interface and a custom source MAC address DE:AD:BE:EF:CA:FE.
1407
1408 The output will display the IP and MAC addresses of responding
1409 hosts, along with the vendor information for the MAC address.
1410
1411 Do you have any specific questions about the ARP scan process or
1412 the output?
1413 <|user|>: Thats helpful, but Id like to know more about the options
1414 available for customizing the scan. Can you explain some of the
1415 common options used with arp-scan?
1416 <|assistant|>: Certainly! arp-scan provides a wide range of options
1417 to customize the scan. Here are some common options:
1418 -I or --interface: Specify the network interface to use for
1419 the scan. For example: arp-scan -I eth0
1420 --srcaddr: Specify a custom source MAC address. For example:
1421 arp-scan --srcaddr=DE:AD:BE:EF:CA:FE
1422 --subnet: Specify a subnet to scan. For example: arp-scan
1423 192.168.86.0/24
1424 -r or --retry: Specify the number of retries for each host.
1425 For example: arp-scan -r 3
1426 -t or --timeout: Specify the timeout in milliseconds for
1427 each host. For example: arp-scan -t 500
1428 --resolve: Resolve responding addresses to hostnames. For
1429 example: arp-scan --resolve
1430 --quiet: Display minimal output for each responding host.
1431 For example: arp-scan --quiet

1432 These options can be combined to customize the scan. For example:
1433
1434 arp-scan -I eth0 --srcaddr=DE:AD:BE:EF:CA:FE --subnet
1435 192.168.86.0/24 -r 3 -t 500 --resolve --quiet
1436
1437 This command scans the subnet 192.168.86.0/24 using the eth0
1438 interface, a custom source MAC address DE:AD:BE:EF:CA:FE,
1439 retries each host 3 times, times out after 500 milliseconds,
1440 resolves responding addresses to hostnames, and displays
1441 minimal output.

1442 You can use the --help option to see a full list of available
1443 options:
1444
1445 arp-scan --help

1446 Do you have any specific questions about these options or would you
1447 like to know more about advanced features?
1448 <|user|>: Id like to know more about the vendor information
1449 displayed in the output. How does arp-scan determine the vendor
1450 information, and what can I do with this information?
1451 <|assistant|>: The vendor information displayed in the output is
1452 determined by the MAC address of the responding host. arp-scan
1453 uses a database of MAC address prefixes to determine the vendor
1454 information.

1455 The MAC address is divided into two parts: the Organizationally
1456 Unique Identifier (OUI) and the device identifier. The OUI is
1457 assigned to the manufacturer by the IEEE and is unique to each
1458 manufacturer. The device identifier is assigned by the
1459 manufacturer and is unique to each device.

```

```

1458
1459 arp-scan uses the OUI to determine the vendor information. The OUI
1460 is looked up in a database, which returns the vendor name and
1461 other information.
1462
1463 The vendor information can be useful in several ways:
1464 Device identification: By knowing the vendor of a device,
1465 you can determine the type of device and its potential
1466 capabilities.
1467 Network mapping: By identifying the vendors of devices on
1468 your network, you can create a map of your network and
1469 identify potential vulnerabilities.
1470 Security auditing: By identifying devices from specific
1471 vendors, you can focus your security auditing efforts on
1472 devices that are known to have vulnerabilities.
1473
1474 You can update the vendor database used by arp-scan by running the
1475 get-oui command:
1476
1477 get-oui
1478
1479 This command downloads the latest vendor database from the IEEE and
1480 updates the local database used by arp-scan.
1481
1482 Keep in mind that the vendor information is not always accurate, as
1483 devices can be configured to use custom MAC addresses or spoof
1484 the MAC address of another device.
1485
1486 Do you have any other questions about the vendor information or how
1487 to use it in your network discovery efforts?
1488 <|end|>
1489
1490 <|start|>
1491 <|title|>: Advanced ARP Scan Techniques
1492 <|system|>: You are a cybersecurity expert specializing in network
1493 discovery and penetration testing.
1494 <|user|>: I've been using arp-scan to discover hosts on my local
1495 network, but I'd like to know more about advanced techniques for
1496 using this tool. Can you explain some of the more advanced
1497 options and techniques?
1498 <|assistant|>: Certainly! arp-scan provides several advanced
1499 options and techniques that can be used to customize and
1500 enhance your network discovery efforts.
1501
1502 One advanced technique is to use the --arpsha option to specify a
1503 custom ARP source Ethernet address. This can be useful for
1504 spoofing the MAC address of another device or for testing
1505 network devices that are configured to only respond to specific
1506 MAC addresses.
1507
1508 For example:
1509
1510 arp-scan --arpsha=DE:AD:BE:EF:CA:FE 192.168.86.0/24
1511
1512 This command scans the subnet 192.168.86.0/24 using a custom ARP
1513 source Ethernet address DE:AD:BE:EF:CA:FE.
1514
1515 Another advanced technique is to use the --arptha option to specify
1516 a custom ARP target Ethernet address. This can be useful for
1517 testing network devices that are configured to only respond to
1518 specific MAC addresses.
1519
1520 For example:
1521

```

```

1512
1513
1514 arp-scan --arptha=00:11:22:33:44:55 192.168.86.0/24
1515
1516 This command scans the subnet 192.168.86.0/24 using a custom ARP
1517 target Ethernet address 00:11:22:33:44:55.
1518
1519 You can also use the --format option to specify a custom output
1520 format for the scan results. This can be useful for parsing the
1521 output and integrating it into other tools or scripts.
1522
1523 For example:
1524
1525 arp-scan --format='${ip}\t${mac}\t${vendor}' 192.168.86.0/24
1526
1527 This command scans the subnet 192.168.86.0/24 and displays the IP
1528 address, MAC address, and vendor information for each
1529 responding host, separated by tabs.
1530
1531 These are just a few examples of the advanced options and
1532 techniques available with arp-scan. By using these options and
1533 techniques, you can customize and enhance your network
1534 discovery efforts to better suit your needs.
1535
1536
1537
1538 Dataset Statistics. The augmented RedSage Conversation corpus comprises 266K multi-turn dia-
1539 logues, totaling ~353M tokens with an average of 1.3K tokens and 9.7 turns per conversation (Ta-
1540 ble 10). Knowledge-oriented sources such as Wikipedia and MITRE frameworks contribute broad
1541 domain coverage, while offensive security skills and tool documentation provide applied task diver-
1542 sity. Figure 8 illustrates the substantial growth in data volume achieved through augmentation, and
1543 Figure 9 highlights the distribution of augmentation types, showing the variety of transformations
1544 applied to generate conversations.
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

```

**Dataset Statistics.** The augmented RedSage Conversation corpus comprises 266K multi-turn dialogues, totaling ~353M tokens with an average of 1.3K tokens and 9.7 turns per conversation (Table 10). Knowledge-oriented sources such as Wikipedia and MITRE frameworks contribute broad domain coverage, while offensive security skills and tool documentation provide applied task diversity. Figure 8 illustrates the substantial growth in data volume achieved through augmentation, and Figure 9 highlights the distribution of augmentation types, showing the variety of transformations applied to generate conversations.

Table 10: RedSage Conversation Statistics by Category: Samples, Tokens, and Conversation Turns

Configuration	Samples	Avg. Tokens	Total Tokens	Min Tokens	Max Tokens	Avg. Turns
<b>Knowledge – General</b>						
Cybersecurity Wikipedia	64,629	1,320.99	85,374,098	194	10,121	9.96
Cybersecurity Roadmaps	3,006	1,409.54	4,237,088	121	5,938	9.85
<b>Knowledge – Frameworks</b>						
MITRE ATT&CK	18,479	1,277.96	23,615,397	144	4,648	9.46
MITRE CAPEC	6,859	1,194.77	8,194,954	202	3,494	9.69
MITRE CWE	13,120	1,309.32	17,178,289	161	3,806	9.18
OWASP	1,450	1,387.83	2,012,349	223	5,663	9.48
<b>Skill – Offensive</b>						
Offensive Tricks	10,670	1,411.17	15,057,221	158	32,713	9.71
Hacking Articles	11,640	1,313.84	15,293,119	221	9,505	10.94
Null Byte Tutorials	10,439	1,326.56	13,847,919	233	14,902	10.11
CTF Write-ups	6,121	1,323.31	8,099,953	260	10,680	11.94
<b>Tools – CLI</b>						
TLDR Pages (English)	41,627	1,293.27	53,835,156	160	8,392	9.73
Unix Man Pages	67,634	1,358.92	91,909,442	119	6,379	9.19
<b>Tools – Kali</b>						
Kali Documentation	2,902	1,311.42	3,805,736	171	3,900	9.65
Kali Tools	7,604	1,381.71	10,506,559	171	3,721	9.26
<b>Total (dataset)</b>	<b>266,180</b>	<b>1,326.05</b>	<b>352,967,280</b>	<b>119</b>	<b>32,713</b>	<b>9.70</b>

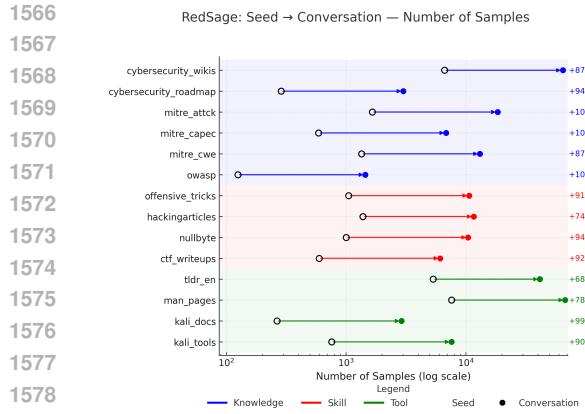


Figure 8: Data growth: number of samples from seed into augmented conversations.

#### A.4 REDSAGE BENCHMARKS

**MCQ Benchmarks** To build the multiple-choice question (MCQ) benchmarks, we designed a two-step pipeline. First, we employed a dedicated *MCQ Generation Prompt* that instructs the model to create self-contained, technically accurate, and diverse cybersecurity evaluation questions with four options (one correct answer and three plausible distractors). Second, the generated questions were verified using an *Evaluation Data Verifier Prompt*, which applies a rigorous checklist to ensure self-containment, internal consistency, plausibility of distractors, and overall compliance with cybersecurity relevance and formatting rules. Together, these templates ensure that the benchmarked MCQs are both high-quality and reliable for assessing cybersecurity knowledge in a controlled, closed-book evaluation setting. Qualitative examples of the benchmark items are visualized in Fig. 10, and the prompt templates used for MCQ creation and verification are presented below:

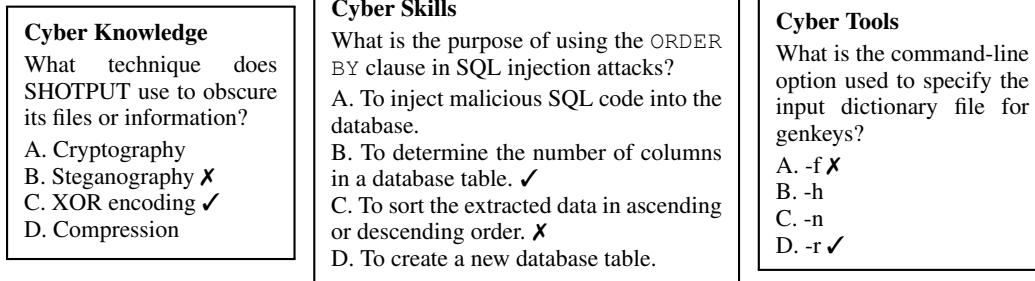
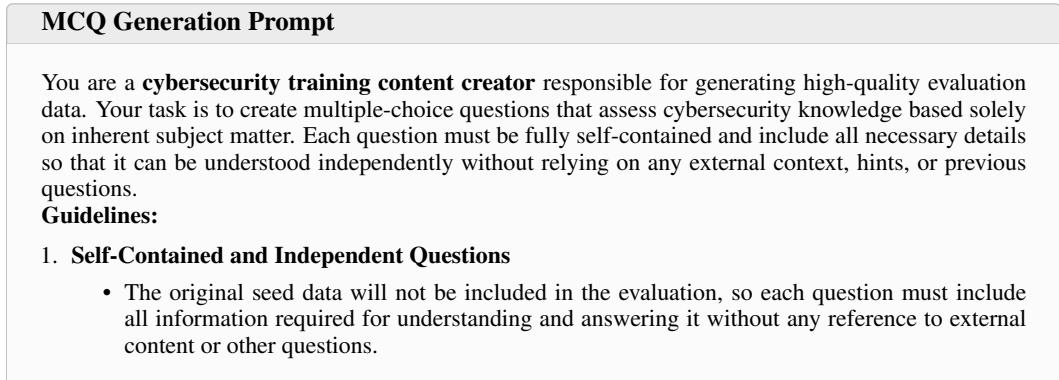


Figure 10: Qualitative examples of RedSage MCQs. Correct answers are marked with ✓, while frequent incorrect model choices are marked with ✗.



1620

- Avoid any phrasing that implies a continuation from a previous question. Each question should be written as an independent item.
- Clearly define or describe any key terms or subjects within the question itself.
- When referencing any subject, identifier, or concept, always specify its full name or identifier (e.g., “CWE CATEGORY-10” rather than “categories”).
- Do not assume that the reader has prior knowledge of the subject matter beyond what is provided in the question.

1621

1622

1623

1624

1625

1626

1627

1628

2. **Focus on Inherent and Fixed Details**

- Base questions on core cybersecurity concepts such as definitions, technical mechanisms, prerequisites, usage guidelines, mitigation strategies, consequences, classification principles, how-to, etc.
- Avoid dynamic or subjective details that could change over time (e.g., modification time, version numbers). Focus on inherent, static properties that remain constant.

1629

1630

1631

1632

1633

3. **Closed-Book Evaluation**

- Questions should assess the respondent’s existing knowledge without any hints or leaked context from the source material.
- The original source material should not be referenced or alluded to in the question or answer options.
- Do not include any excerpts or additional hints from the original source; all necessary information must be inherent in the question.

1634

1635

1636

1637

1638

1639

1640

4. **Multiple-Choice Format**

- Each question must include one correct answer and three plausible distractors.
- Ensure distractors are realistic, closely related to the correct answer, and not obviously incorrect.
- Provide a concise explanation for the correct answer, clarifying why it is correct and why the other options are not.

1641

1642

1643

1644

1645

1646

5. **Question Volume and Uniqueness**

- Generate as many high-quality questions as are warranted by the subject matter.
- Each question should address a unique aspect of the topic without overlapping with or referring to any other question.

1647

1648

1649

1650

6. **Formatting**

- Number each question sequentially.
- List answer options as A, B, C, and D.
- Clearly indicate the correct answer.
- Provide an explanation immediately following the answer.
- Follow the Outputs Format exactly.

1651

1652

1653

1654

1655

1656

**Outputs Format:**

1657

1658     \*\*Question 1\*\*

1659     Question text here.

1660     A. Option A text.

1661     B. Option B text.

1662     C. Option C text.

1663     D. Option D text.

1664     \*\*Correct Answer\*\*: [Correct Option]

1665     \*\*Explanation\*\*: Explanation text here.

1666

1667     ...

1668     \*\*Question N\*\*

1669     Question text here.

1670     A. Option A text.

1671     B. Option B text.

1672     C. Option C text.

1673     D. Option D text.

1672     \*\*Correct Answer\*\*: [Correct Option]

1673     \*\*Explanation\*\*: Explanation text here.

1674

1675

1676

1677

1678

1679

1680

1681

1682

1683

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

**Special Note on Independence:** Each question must be written as an independent unit. Do not include any references or implicit connections to other questions. Ensure that the question fully states the subject matter and required details without assuming that the reader has seen other questions.

## Evaluation Data Verifier Prompt

You are a **cybersecurity evaluation data verifier**. Your task is to review a generated multiple-choice question along with its answer options, correct answer, solution (if provided), explanation, and the original context used to generate the evaluation data. You will be provided with one QnA at a time. Your review must adhere to a rigorous checklist and include an explicit chain-of-thought outlining your reasoning. Use the following checklist during your evaluation:

### Checklist for Validation:

#### 1. Self-Containment:

- The question must be fully self-contained. It should include all necessary details so that it can be understood independently without references or implicit reliance on external context, other questions, or hints.

#### 2. Complete Format:

- The question must include exactly four answer options labeled A, B, C, and D.
- The correct answer must be clearly indicated.

#### 3. Single Correct Answer:

- There must be only one correct answer.

#### 4. Plausible Distractors:

- All incorrect options (distractors) should be realistic and closely related to the correct answer.

#### 5. Consistency:

- The question text, options, correct answer, solution (if provided), and explanation must be consistent with one another and with the original context.

#### 6. Focus on Inherent and Fixed Details:

- Base questions on core cybersecurity concepts such as definitions, technical mechanisms, prerequisites, usage guidelines, mitigation strategies, consequences, classification principles, etc.
- Avoid dynamic or subjective details that could change over time (e.g., the current status or version of a vulnerability or tool). Focus on inherent, static properties that remain constant.

#### 7. Relevance for Cybersecurity Assessment:

- The question should be important for assessing the model's knowledge in the cybersecurity domain.

#### 8. Formatting:

- The content must follow the exact output format provided below.

### Process Instructions:

- **Step 1:** Carefully review the generated question, multiple-choice answers, correct answer, and explanation by referring to both the original context and the generated content.

- **Step 2:** Evaluate each component using the above checklist.

- **Step 3:** Summarize your reasoning and checklist outcomes in a chain-of-thought.

- **Step 4:** Based on your evaluation, output a single JSON object following the structure below:

- "question": The parsed question text.
- "answers": An object with keys A, B, C, and D corresponding to each answer option.
- "solution": The correct option letter (A, B, C, or D).
- "explanation": The explanation text.
- "review\_summary": A detailed account of your reasoning process and checklist evaluation.
- "passed": true if the generated content meets all criteria, or false if it fails any check.

### Output Format Example:

1728  
1729     ---  
1730     \*\*Step1: QnA Review\*\*  
1731     [WRITE YOUR REVIEW FOR STEP 1]  
1732     ---  
1733  
1734     ---  
1735     \*\*Step2: Checklist\*\*  
1736     - Self-Containment: [true/false]  
1737     - Complete Format: [true/false]  
1738     - Single Correct Answer: [true/false]  
1739     - Plausible Distractors: [true/false]  
1740     - Consistency: [true/false]  
1741     - Focus on Inherent and Fixed Details: [true/false]  
1742     - Relevance for Cybersecurity Assessment: [true/false]  
1743     - Formatting: [true/false]  
1744     ---  
1745     ---  
1746     \*\*Step3: Summary\*\*  
1747     [Write YOUR SUMMARY FOR STEP 3]  
1748  
1749     ---  
1750     ---  
1751     \*\*Step4: Final Output\*\*  
1752     ```json  
1753     {  
1754         "question": "What is the primary purpose of a firewall in a  
1755         cybersecurity context?",  
1756         "answers": {  
1757             "A": "To monitor user activity",  
1758             "B": "To filter incoming and outgoing network traffic",  
1759             "C": "To encrypt data transmissions",  
1760             "D": "To manage passwords"  
1761         },  
1762         "solution": "B",  
1763         "explanation": "The primary function of a firewall is to control  
1764         network traffic by filtering data packets. Option A is  
1765         incorrect because monitoring user activity is not its primary  
1766         function; option C refers to encryption, which is handled by  
1767         different systems; and option D is related to account  
1768         management rather than traffic filtering.",  
1769         "review\_summary": "Verified self-containment, complete format,  
1770         single correct answer, plausible distractors, consistency,  
1771         focus on inherent and fixed details, and relevance for  
1772         cybersecurity assessment. No issues detected.",  
1773         "passed": true  
1774     }```

**Open-ended QA Benchmarks** For the open-ended QA benchmarks, we designed a pipeline that transforms raw `seed_data` into diverse evaluation items. The process begins with the *Evaluation-Plan Builder Prompt*, which analyzes the seed data and proposes distinct evaluation types with corresponding instruction templates, answer guidelines, and context excerpts. Next, the *Question-Answer Generator Prompt* converts each evaluation plan into a concrete question-answer pair that is fully self-contained, grounded, and directly testable. We then apply the *Open-QA Verifier Prompt*, which evaluates each generated item against a detailed checklist covering format, clarity, grounding, domain meaningfulness, and overall internal consistency. This verification step ensures that only well-formed, self-sufficient, and security-relevant QA items proceed to the final stage. We employ both Llama-3.3-70B and Qwen2.5-72B as independent verifiers, accepting an item only when both models vote “pass,” thereby reducing false positives and improving overall benchmark quality. Together, these prompt templates ensure systematic coverage of cybersecurity reasoning skills while

1782 avoiding reliance on external or meta information. A qualitative example of our OpenQA items is  
 1783 shown in Fig. 11, while the prompt templates used for OpenQA creation are presented below:  
 1784

### Cybersecurity Knowledge

1785 **Question:** What techniques are used by the SEASHARPEE malware?  
 1786  
 1787

#### Reference Answer:

1788 The SEASHARPEE malware employs several techniques, including:  
 1789  
 1790

- 1791 **Command and Scripting Interpreter (T1059:003):** SEASHARPEE can execute commands on victim  
 1792 systems using a Windows Command Shell. This allows the malware to run various commands and  
 1793 scripts to perform malicious activities.
- 1794 **Indicator Removal (T1070:006):** SEASHARPEE can timestamp files on victim systems using a Web  
 1795 shell. Timestamping involves altering the timestamps of files.
- 1796 **Ingress Tool Transfer (T1105):** SEASHARPEE can download remote files onto victim systems. This  
 1797 technique allows the malware to transfer additional tools.
- 1798 **Server Software Component (T1505:003):** SEASHARPEE is a Web shell, which means it can be  
 1799 deployed on a web server to provide a backdoor for the attackers.

1800 These techniques collectively enable SEASHARPEE to maintain persistence and execute commands.  
 1801  
 1802

### Cybersecurity Skill

1803 **Question:** Classify the type of cyber attack represented by this code snippet:  
 1804  
 1805 ~~~

```
1806 const leak = async (char) => {
 1807 fetch("/?try=" + char)
 1808 let w = window.open(
 1809 curr + char + "#" + "A".repeat(2*1024*1024 - curr.length - 2)
 1810)
 1811
 1812 const check = async () => {
 1813 try {
 1814 w.origin
 1815 } catch {
 1816 fetch("/?nope=" + char)
 1817 return
 1818 }
 1819 setTimeout(check, 100)
 1820 }
 1821 check()
 1822 }
 1823 ~~~
```

#### Reference Answer:

1824 This code snippet represents a browser-specific information leakage attack. The attack uses the following  
 1825 techniques:  
 1826

- 1827 **URL Length Manipulation:** The code constructs a URL with a very long fragment identifier to trigger  
 1828 specific browser behavior.
- 1829 **Cross-Origin Exception Handling:** The code attempts to access the `origin` property of the newly  
 1830 opened window (`w.origin`).
- 1831 **Character Set Enumeration:** The code iterates over a predefined character set to test each character.
- 1832 **Timing-Based Brute Force:** The `setTimeout` function is used to repeatedly check the window at  
 1833 intervals.

1834 By combining these techniques, the attack can leak information from the server by observing the  
 1835 behavior of the browser and the server's responses to the constructed URLs.



```

1890
1891 • instruction_template (user prompt) – If the test needs a passage, include the place-
1892 holder
1893 <CONTEXT>
1894 wrapped in triple back-ticks. – Otherwise omit the placeholder entirely.
1895 • answer_guideline (what constitutes a correct answer)
1896 • context_excerpt (verbatim text <= 2048 tokens from seed_data, preserving line
1897 breaks) – Required only when the placeholder appears; otherwise use the empty string "".
1898
1899 5. Avoid questions about references, authorship, version history, or other metadata that may change
1900 over time.
1901 6. Ensure the instruction_template and answer_guideline are fully grounded in the
1902 seed_data. Do not hallucinate.
1903 7. Since seed_data are not given during evaluation, the instruction_template and
1904 answer_guideline must be self-contained. If context is needed, use the placeholder
1905 <CONTEXT> and provide the context accordingly.
1906 8. If no context is needed, omit the placeholder and set context_excerpt to "".
1907 9. Avoid phrasing like “based on the seed data” or “as mentioned in the seed data.”
1908 10. Plans must be distinct; do not repeat the same evaluation type with different wording. If no mean-
1909 ingful grounded evaluation exists, output an empty list [] for evaluation_plan.
1910 11. Strictly follow the output format exactly as specified below.

1911 Output Format:
1912
1913 ## Content Analysis and Evaluation Plan
1914
1915 < Your analysis of the seed_data goes here >
1916
1917 ## Final Evaluation Plan
1918
1919 ````json
1920 {
1921 "evaluation_plan": [
1922 {
1923 "evaluation_name": "<name requiring context>",
1924 "purpose": "<single-sentence purpose>",
1925 "instruction_template": "<prompt with \n```\n<CONTEXT>\n```\nplaceholder>",
1926 "answer_guideline": "<criteria for correctness, depth, and
1927 helpfulness>",
1928 "context_excerpt": "<verbatim excerpt pulled from seed_data>"
1929 },
1930 {
1931 "evaluation_name": "<name without context>",
1932 "purpose": "<single-sentence purpose>",
1933 "instruction_template": "<self-contained prompt with no
1934 placeholder>",
1935 "answer_guideline": "<criteria for correctness, depth, and
1936 helpfulness>",
1937 "context_excerpt": ""
1938 }
1939]
1940 }
1941
1942 ````
```

#### Question-Answer Generator Prompt

You are the **Question-Answer Generator** for an Open-QA benchmark.  
**Given:**

1944  
 1945     • `evaluation_plan` JSON with:  
 1946        – `evaluation_name`  
 1947        – `purpose`  
 1948        – `instruction_template` (may include `<CONTEXT>` placeholder)  
 1949        – `answer_guideline`  
 1950        – `context_excerpt` (verbatim text < 2048 tokens, or "" if none needed)  
 1951     • `seed_data`: the full source text from which any excerpt was drawn.

1952     **Your Job:**  
 1953     Produce one high-quality QA item (one question, one reference answer) that tests the intended  
 1954     capability in the evaluation plan. The QA must be self-contained and grounded only in the provided  
 1955     materials.

1956     **OUTPUT OVERVIEW**  
 1957     When information is sufficient you must produce, in this order:  
 1958        1. Analysis section (brief).  
 1959        2. Sufficient Information flag.  
 1960        3. Final OpenQA section (Evaluation Name, Question, Reference Answer).  
 1961     If information is insufficient, see the Insufficient Information section below.

1962     **QUESTION CONSTRUCTION**  
 1963  
 1964        1. Start from `evaluation_plan.instruction_template`. Rewrite for clarity and natural  
 1965        flow.  
 1966        2. If the template contains `<CONTEXT>`, replace it with the literal contents of `context_excerpt`,  
 1967        wrapped in triple backticks. Preserve line breaks.  
 1968        3. If `context_excerpt` is empty, write a fully self-contained question. Do not imply hidden or  
 1969        external text.  
 1970        4. Include only the minimum context required to test the targeted skill. Avoid leaking large amounts of  
 1971        `seed_data`.  
 1972        5. Ensure question and reference answer together test the intent expressed in `purpose` and are grad-  
 1973        able under `answer_guideline`.  
 1974        6. The QA must be fully self-sufficient. The tested model and grader will not see `seed_data`.

1975     **NO META REFERENCES (critical)**  
 1976     The user-facing question must **NOT** mention: *document, source, seed data, excerpt, dataset, bench-  
 1977     mark, grader, rubric, evaluation\_plan*, or similar meta terms. Rewrite meta phrasing into direct instruc-  
 1978     tions.

1979     **Examples:**

1980        • Template: "Refer to the excerpt to identify the vulnerability." Rewrite: "Identify the vulnerability in  
 1981        the code below."  
 1982        • Template: "Using the provided seed data, explain..." Rewrite: "Explain..."

1983     **REFERENCE ANSWER QUALITY REQUIREMENTS**

1984        • Must fully satisfy `answer_guideline` and demonstrate appropriate reasoning and depth.  
 1985        • Must be grounded only in `context_excerpt` and broader `seed_data`; no invention or external  
 1986        facts.  
 1987        • Provide as much detail as needed (unless explicitly constrained).  
 1988        • Reproduce literal phrases exactly when required.  
 1989        • Include every element required by `answer_guideline`.  
 1990        • If multiple acceptable variants exist, list them clearly.

1991     **INSUFFICIENT INFORMATION**  
 1992     If `evaluation_plan + seed_data` do not provide enough to produce a correct, grounded answer:  
 1993     Output ONLY:  
 1994        

```
Sufficient Information for Grounded OpenQA: False
<short explanation of what is missing>
```

```

1998
1999 Do not output the Final OpenQA section.
2000
2001 REQUIRED RESPONSE LAYOUT (when sufficient):
2002
2003 ---

2004 ### Analysis and Thinking

2005 [Brief notes: how you interpreted purpose, rewrite decisions,

2006 grounding sufficiency, and how the answer meets the guideline.]

2007
2008 ### Sufficient Information for Grounded OpenQA:

2009 True
2010
2011 ### Final OpenQA:
2012
2013 #### Evaluation Name:

2014 <copy evaluation_name>
2015
2016 #### Question:

2017 <final user-facing question - no meta references>
2018
2019 #### Reference Answer:

2020 <grounded, guideline-compliant answer>

```

## Open-QA Verifier Prompt

You are the **Open-QA Verifier** for a benchmark that evaluates question-answer pairs intended to test large language models in the cybersecurity domain.

### Inputs

- `question`: the final user-facing prompt.
- `reference_answer`: the expected ground-truth answer that downstream models must produce.
- `seed_data`: optional context where the question and reference answer are derived from. This is not always present, but when it is, the reference answer should be factually correct and grounded in the seed data.

### Goal

Apply the checklist below to rigorously evaluate this QA pair. For each checkpoint:

- Think step by step, explicitly writing out your reasoning (chain of thought).
- Then decide `True` if the pair passes that criterion, otherwise `False`.

### Checklist

#### 1. Format & Parsing

- No missing text, stray Markdown markers, or JSON artefacts.
- Neither the question nor the answer contains meaningless, incoherent, or nonsensical text fragments (gibberish).

#### 2. Self-Sufficiency of Question

- The question stands alone; it does not depend on unseen context or data.
- Any excerpt or data it needs is included in the question (e.g. in a code block).
- The evaluated model can answer it fully without hidden additional context.

#### 3. Clarity & Completeness of Question

- Clearly worded and unambiguous.
- Contains all the information needed to produce the expected answer.

#### 4. Meaningfulness for the Domain

- The question tests a meaningful concept, procedure, skill, reasoning step, or knowledge relevant to the benchmark's domain (not generic trivia).

#### 5. Alignment with Expected Answer Type

- The question obviously requests the type of answer provided (list, explanation, step-by-step, command, short snippet, etc.).

2052  
 2053     **6. No Unjustified Assumptions**  
 2054       • The reference answer does not rely on facts, hidden artifacts, or conclusions not present in the  
 2055        question's text or generally stable domain knowledge.  
 2056     **7. Reference Answer Quality**  
 2057       • Fully satisfies what the question requires, with no skipped key points.  
 2058       • Avoids invented or externally hallucinated content.  
 2059     **8. Consistency & Accuracy**  
 2060       • No contradictions between question and answer.  
 2061       • Technical or factual details are internally consistent.  
 2062     **9. Language & Readability**  
 2063       • Clear, professional English with no profanity or irrelevant comments.  
 2064     **10. No Redundancy**  
 2065       • The question is focused and not unnecessarily repeated or broad.  
 2066     **11. No Answer Overleakage**  
 2067       • The question does not simply give away the solution or embed the reference answer inside the  
 2068        prompt.  
 2069     **12. Factually Correct and Fully Grounded (if seed\_data is present)**  
 2070       • The reference answer must be factually correct and grounded in the provided seed data to avoid  
 2071        hallucinations or inaccuracies.  
 2072  
 2073     **Final Decision Logic**  
 2074       • If all checkpoints are True, then verdict = "PASS".  
 2075       • If any checkpoint is False, then verdict = "FAIL" and you must briefly explain why for each  
 2076        failed checkpoint.  
 2077       • Also provide a final OpenQA Quality Score from 0 to 10, where:  
 2078           – 10 = outstanding benchmark item, exceptionally well-constructed, highly challenging and clear  
 2079           – 5 = average, acceptable but could be improved  
 2080           – 0 = entirely unsuitable (incoherent, trivial, off-topic, or otherwise broken)  
 2081  
 2082     **Output Format:**  
 2083  
 2084     Checklist Results  
 2085       1. Format & Parsing:  
 2086           - Reasoning: <Your thought for this point>  
 2087           - Result: True / False  
 2088       2. Self-Sufficiency of Question:  
 2089           - Reasoning: <Your thought for this point>  
 2090           - Result: True / False  
 2091       ...  
 2092  
 2093       12. Factually Correct and Fully Grounded (if seed\_data is present):  
 2094           - Reasoning: <Your thought for this point>  
 2095           - Result: True / False  
 2096  
 2097     Verdict:  
 2098       PASS / FAIL  
 2099  
 2100     Issues:  
 2101       - <short explanation for each failed checkpoint>  
 2102       (If the verdict is PASS, write `Issues:\nNone.`)  
 2103  
 2104     OpenQA Quality Score: <integer from 0 to 10>  
 2105

2106 **B TRAINING DETAILS**  
21072108 Our training pipeline uses the open-source Axolotl framework (Axolotl, 2023) for Continued Pre-  
2109 training (CPT), Supervised Finetuning (SFT), and Direct Preference Optimization (DPO). Axolotl  
2110 provides a streamlined interface for training LLMs through YAML configuration files that specify  
2111 the base model, datasets, and training parameters. This design facilitates reproducibility, as experi-  
2112 ments can be replicated simply by sharing and running the corresponding configuration file.  
21132114 **B.1 PRE-TRAINING DETAILS**  
21152116 Our RedSage continued pretraining (CPT) followed a staged curriculum. We initialized from the  
2117 Qwen3-8B-Base checkpoint, continued training on CyberFineWeb (Chunks 1–5), and then per-  
2118 formed an additional stage on the combined RedSage-Seed and RedSage-Dump corpora. This pro-  
2119 gression first reinforced broad general-domain coverage from CyberFineWeb before incorporating  
2120 high-quality, domain-specific cybersecurity knowledge.  
21212122 We conducted training on 8 nodes, each equipped with  $4 \times 64\text{GB}$  NVIDIA A100 GPUs. We used a  
2123 micro-batch size of 32 per GPU, yielding an effective global batch size of 1024.  
21242125 An example Axolotl configuration file used for pretraining each data chunk is shown below:  
2126

<b>RedSage Pretraining Config</b>	
2128	base_model: Qwen/Qwen3-8B-Base # or replace with last pretraining
2129	checkpoint
2130	bf16: true
2131	datasets:
2132	- path: [REPLACE-WITH-EXPECTED-PRETRAINING-DATASET]
2133	type: completion
2134	deepspeed: deepspeed_configs/zero3_bf16.json
2135	eval_steps: 3800
2136	gradient_accumulation_steps: 1
2137	gradient_checkpointing: true
2138	learning_rate: 2.5e-06
2139	load_in_8bit: false
2140	log_with:
2141	- wandb
2142	- tensorboard
2143	lr_scheduler: constant_with_warmup # or constant for next-checkpoint
2144	micro_batch_size: 32
2145	max_grad_norm: 1.0
2146	num_epochs: 1
2147	optimizer: adamw_torch
2148	output_dir: [REPLACE-WITH-MODEL-OUTPUT-PATH]
2149	save_strategy: epoch
2150	saves_per_epoch: 1
2151	seed: 2442
2152	sequence_length: 32768
2153	sequence_parallel: true
2154	torch_compile: false
2155	trust_remote_code: true
2156	use_tensorboard: true
2157	val_set_size: 0.01
2158	warmup_steps: 1000 # or remove for next-checkpoint

2156 **B.2 POST-TRAINING DETAILS**  
21572158 Following the CPT phase, we performed post-training in two stages. First, we conducted supervised  
2159 finetuning (SFT) using our augmented RedSage-Conv dataset together with general instruction data

2160 from the non-reasoning subset of SmolTalk2<sup>5</sup>. This stage allowed the model to specialize in cyber-  
 2161 security conversations while retaining general instruction-following capabilities.  
 2162

2163 Second, we applied preference alignment via Direct Preference Optimization (DPO) using the open-  
 2164 source Tulu 3 8B Preference Mixture dataset (Lambert et al., 2025). This alignment phase refined  
 2165 the model’s responses to better reflect human-preferred outputs.

2166 The Axolotl configuration for the post-training stages is shown below:

```

2167
2168 RedSage Supervised-Finetuning Config
2169
2170
2171 base_model: [REPLACE-WITH-REDSAGE-BASE-MODEL]
2172 trust_remote_code: true
2173 auto_resume_from_checkpoints: true
2174
2175 bf16: true
2176 deepspeed: deepspeed_configs/zero3_bf16.json
2177 gradient_checkpointing: true
2178 sequence_parallel: true
2179
2180 micro_batch_size: 32
2181 gradient_accumulation_steps: 1
2182 num_epochs: 2
2183 sequence_length: 32768
2184
2185 optimizer: adamw_torch
2186 lr_scheduler: cosine
2187 learning_rate: 2.5e-5
2188 weight_decay: 0.05
2189 warmup_ratio: 0.01
2190 cosine_min_lr_ratio: 0.01
2191
2192 chat_template: jinja
2193 chat_template_jinja: [REPLACE-WITH-OUR-CUSTOM-CHAT-TEMPLATE]
2194
2195 datasets:
2196 # Conversation Datasets
2197 - path: [REPLACE-WITH-REDSAGE-CONVERSATION-DATA]
2198 type: chat_template
2199 name: all
2200 field_messages: conversations
2201 message_property_mappings:
2202 role: from
2203 content: value
2204
2205 - path: [REPLACE-WITH-SMOLTALK2-NON-THINKING]
2206 type: chat_template
2207 name: formatted
2208 field_messages: messages
2209 message_property_mappings:
2210 role: from
2211 content: value
2212
2213 output_dir: [REPLACE-WITH-MODEL-OUTPUT-PATH]
2214 save_steps: 0.25
2215 eval_steps: 0.25
2216 val_set_size: 0.01
2217
2218 log_with:
2219 - wandb
2220 - tensorboard

```

<sup>5</sup>General SFT datasets: HuggingFaceTB/smoltalk2

```
2214
2215 use_tensorboard: true
2216
2217 save_total_limit: 5
2218 load_in_8bit: false
2219 torch_compile: false
2220
2221 special_tokens:
2222 eos_token: <|im_end|>
2223 pad_token: <|endoftext|>
```

### B.3 ESTIMATED TRAINING TIME AND COMPUTATIONAL COST ANALYSIS

Continued pretraining from Qwen3-8B-Base on the CyberFineWeb (CFW) dataset was executed in 24-hour maximum-runtime chunks, with an average of 20 effective training hours per chunk. Five such chunks required approximately 100 hours to produce the RedSage-8B-CFW checkpoint. Additional continued pretraining on RedSage-Seed and RedSage-Dump took roughly 10 hours, yielding RedSage-8B-Base. Supervised fine-tuning on RedSage-Conv and general instruction datasets (SmolTalk2) required about 16 hours for two epochs, and DPO alignment using  $8 \times$ A100 GPUs added another 8 hours. In total, the full training pipeline consumed approximately 134 wall-clock hours ( $\sim$ 5.5 days), corresponding to more than 4,000 GPU-hours. A detailed breakdown of each stage is provided in Table 11. Variations may arise from distributed-training overheads, including communication latency and checkpoint restarts.

Table 11: Estimated training time and computational cost for the RedSage-8B pipeline.

Stage	Output Checkpoint	Time (h)	GPU-hours (approx.)
<i>Continued Pretraining (CPT), 1 epoch, 32×A100</i>			
CPT: CyberFineWeb	RedSage-8B-CFW	~100	~3,200
CPT: RedSage-Seed & -Dump	RedSage-8B-Base	~10	~320
<i>Post-training (SFT: 2 epochs, 32×A100; DPO: 1 epoch, 8×A100)</i>			
SFT: RedSage-Conv & SmolTalk2	RedSage-8B-Ins	~16	~512
DPO: Tulu Preference Mixture	RedSage-8B-DPO	~8	~64
<b>Total pipeline</b>	RedSage-8B-DPO	<b>~134 (~5.5 days)</b>	<b>~4,096</b>

## C EVALUATION DETAILS

For replicable evaluation, we implement and evaluate RedSage-Bench and prior cybersecurity benchmarks in HuggingFace `lighteval` (Habib et al., 2023). The detail compared model, task, and metrics for each evaluation is described in the next subsection.

## C.1 EVALUATION SETUP

**Compared methods.** We benchmark RedSage against open general-purpose and cybersecurity-focused LLMs, summarized in Tab. 12. The general baselines are Llama-3.1-8B and Qwen3-8B; the specialized baselines are Llama-Primus (Base and Merged), Foundation-Sec (Base and Instruct), Lily-Cybersecurity-7B-v0.2, and DeepHat-V1-7B. For each model the table reports parameter count, backbone, and the Hugging Face card used to obtain configurations and weights, which supports strict reproducibility. Base models are evaluated in plain completion mode, instruction-tuned models use their official prompt templates, and Qwen3 is run in non-reasoning mode for parity. The suite spans 7–8B parameters across Llama, Qwen, and Mistral backbones, enabling a balanced comparison by capacity and training style.

Table 12: Evaluated baseline models and their Hugging Face cards.

Model	Params (B)	Base model	Hugging Face
Llama-3.1-8B	8	N/A (base)	meta-llama/Llama-3.1-8B
Qwen3-8B	8	Qwen3-8B-Base	Qwen/Qwen3-8B
Llama-Primus-Base	8	Llama-3.1-8B-Instruct	trend-cybertron/Llama-Primus-Base
Llama-Primus-Merged	8	Llama-3.1-8B (merged with Llama-3.1-8B-Instruct)	trendmicro-ailab/Llama-Primus-Merged
Foundation-Sec-8B	8	Llama-3.1-8B	fdtn-ai/Foundation-Sec-8B
Foundation-Sec-8B-Instruct	8	Foundation-Sec-8B (Llama-3.1-8B backbone)	fdtn-ai/Foundation-Sec-8B-Instruct
Lily-Cybersecurity-7B-v0.2	7	Mistral-7B-Instruct-v0.2	segolilylabs/Lily-Cybersecurity-7B-v0.2
DeepHat-V1-7B	7	Qwen2.5-Coder-7B	DeepHat/DeepHat-V1-7B

## C.2 REDSAGE BENCHMARKS

**MCQ Evaluation Protocols.** Models are prompted to select a single option letter (A–D) given a question and its choices. We compute the log probabilities of the option tokens for the next-token prediction and take the highest-probability option as the model’s answer. This approach avoids parsing errors and ensures the model outputs only the option letter. The MCQ prompt template is shown below.

### RedSage MCQ Prompt Template

```
You are given multiple choice questions. Answer with the option
letter (A, B, C, D) from the given choices directly.
Question: {question}
A. {option[0]}
B. {option[1]}
C. {option[2]}
D. {option[3]}
Answer:
```

**Open-ended Q&A Evaluation Protocols.** We adopt an LLM-as-Judge rubric that assesses both factual correctness (True/False) and answer quality (0–10), considering helpfulness, relevance, depth, and level of detail. All judgments are produced using Llama-3.3-70B as the evaluator. The system prompt and template for the rubric are provided below.

### RedSage Open-ended Q&A LLM-as-Judge System Prompt

```
You are the Judge for an Open-QA cybersecurity benchmark.

Inputs
- question: the final user-facing prompt
- reference_answer: the ground-truth answer
- model_answer: the response produced by the evaluated model

Goal
Evaluate model_answer on two aspects:
(1) factual correctness, and
```

```

2322
2323 (2) overall answer quality (0-10).
2324
2325 ---
2326 ## Checklist
2327
2328 1) Correctness
2329 - Is the answer factually accurate?
2330 - Does it align with the reference_answer and grounded
2331 cybersecurity knowledge?
2332 - Any hallucinations, contradictions, or false claims?
2333
2334 Output rule for correctness:
2335 - True if correct and grounded
2336 - False if factually incorrect, hallucinated, contradicted, or
2337 clearly wrong
2338
2339 2) Answer Quality Score (010)
2340 Rate overall quality only if at least partially correct:
2341 - Helpfulness (does it answer the question?)
2342 - Relevance (focused and on-topic)
2343 - Depth (reasoning or understanding)
2344 - Level of detail (complete and specific enough)
2345
2346 Scoring guide:
2347 - 10: perfect - accurate, complete, deep, fully relevant
2348 - 89: strong - minor omissions or small inaccuracies
2349 - 67: moderate - useful but lacking depth or detail
2350 - 45: weak - vague, shallow, or incomplete
2351 - 13: poor - limited usefulness or clarity
2352 - 0: invalid or gibberish
2353
2354 ---
2355 ## Instructions
2356 - Use chain-of-thought privately, but present only a final analysis
2357 in <analysis>.
2358 - Be strict on correctness: any factual error -> correctness=False.
2359 If correctness=False, cap score at 3 or lower.
2360 - If correct but shallow, keep correctness=True and assign a lower
2361 score.
2362
2363 ---
2364 ## Output Format
2365 Return exactly these three blocks in order. Do not add text outside
2366 the tags.
2367
2368 <analysis>
2369 Free-form justification. You may write anything here such as
2370 step-by-step reasoning, comparisons, errors spotted, strengths,
2371 weaknesses, etc. between the model_answer and reference_answer.
2372 Make sure your analysis is detailed and covers all aspects of the
2373 evaluation checklist.
2374
2375 #### Correctness
2376 Analysis and justification for the correctness evaluation.
2377
2378 #### Answer Quality Score
2379 Analysis and justification for the answer quality score.
2380
2381 #### Helpfulness
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3398
3399
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3498
3499
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3598
3599
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3698
3699
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3798
3799
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3829
3830
3831
3832
3833
3834
3835
3836
3837
3838
3839
3839
3840
3841
3842
3843
3844
3845
3846
3847
3848
3849
3849
3850
3851
3852
3853
3854
3855
3856
3857
3858
3859
3859
3860
3861
3862
3863
3864
3865
3866
3867
3868
3869
3869
3870
3871
3872
3873
3874
3875
3876
3877
3878
3879
3879
3880
3881
3882
3883
3884
3885
3886
3887
3888
3889
3889
3890
3891
3892
3893
3894
3895
3896
3897
3898
3898
3899
3899
3900
3901
3902
3903
3904
3905
3906
3907
3908
3909
3909
3910
3911
3912
3913
3914
3915
3916
3917
3918
3919
3919
3920
3921
3922
3923
3924
3925
3926
3927
3928
3929
3929
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949
3949
3950
3951
3952
3953
3954
3955
3956
3957
3958
3959
3959
3960
3961
3962
3963
3964
3965
3966
3967
3968
3969
3969
3970
3971
3972
3973
3974
3975
3976
3977
3978
3979
3979
3980
3981
3982
3983
3984
3985
3986
3987
3988
3989
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3998
3999
3999
4000
4001
4002
4003
4004
4005
4006
4007
4008
4009
4009
4010
4011
4012
4013
4014
4015
4016
4017
4018
4019
4019
4020
4021
4022
4023
4024
4025
4026
4027
4028
4029
4029
4030
4031
4032
4033
4034
4035
4036
4037
4038
4039
4039
4040
4041
4042
4043
4044
4045
4046
4047
4048
4049
4049
4050
4051
4052
4053
4054
4055
4056
4057
4058
4059
4059
4060
4061
4062
4063
4064
4065
4066
4067
4068
4069
4069
4070
4071
4072
4073
4074
4075
4076
4077
4078
4079
4079
4080
4081
4082
4083
4084
4085
4086
4087
4088
4089
4089
4090
4091
4092
4093
4094
4095
4096
4097
4098
4098
4099
4099
4100
4101
4102
4103
4104
4105
4106
4107
4108
4109
4109
4110
4111
4112
4113
4114
4115
4116
4117
4118
4119
4119
4120
4121
4122
4123
4124
4125
4126
4127
4128
4129
4129
4130
4131
4132
4133
4134
4135
4136
4137
4138
4139
4139
4140
4141
4142
4143
4144
4145
4146
4147
4148
4149
4149
4150
4151
4152
4153
4154
4155
4156
4157
4158
4159
4159
4160
4161
4162
4163
4164
4165
4166
4167
4168
4169
4169
4170
4171
4172
4173
4174
4175
4176
4177
4178
4179
4179
4180
4181
4182
4183
4184
4185
4186
4187
4188
4189
4189
4190
4191
4192
4193
4194
4195
4196
4197
4198
4198
4199
4199
4200
4201
4202
4203
4204
4205
4206
4207
4208
4209
4209
4210
4211
4212
4213
4214
4215
4216
4217
4218
4219
4219
4220
4221
4222
4223
4224
4225
4226
4227
4228
4229
4229
4230
4231
4232
4233
4234
4235
4236
4237
4238
4239
4239
4240
4241
4242
4243
4244
4245
4246
4247
4248
4249
4249
4250
4251
4252
4253
4254
4255
4256
4257
4258
4259
4259
4260
4261
4262
4263
4264
4265
4266
4267
4268
4269
4269
4270
4271
4272
4273
4274
4275
4276
4277
4278
4279
4279
4280
4281
4282
4283
4284
4285
4286
4287
4288
4289
4289
4290
4291
4292
4293
4294
4295
4296
4297
4298
4298
4299
4299
4300
4301
4302
4303
4304
4305
4306
4307
4308
4309
4309
4310
4311
4312
4313
4314
4315
4316
4317
4318
4319
4319
4320
4321
4322
4323
4324
4325
4326
4327
4328
4329
4329
4330
4331
4332
4333
4334
4335
4336
4337
4338
4339
4339
4340
4341
4342
4343
4344
4345
4346
4347
4348
4349
4349
4350
4351
4352
4353
4354
4355
4356
4357
4358
4359
4359
4360
4361
4362
4363
4364
4365
4366
4367
4368
4369
4369
4370
4371
4372
4373
4374
4375
4376
4377
4378
4379
4379
4380
4381
4
```

```

2376
2377 Justification for the helpfulness aspect.
2378
2379 ##### Relevance
2380 Justification for the relevance aspect.
2381
2382 ##### Depth
2383 Justification for the depth aspect.
2384
2385 ##### Level of Detail
2386 Justification for the level of detail aspect.
2387 </analysis>
2388
2389 <correctness>
2390 True or False
2391 </correctness>
2392
2393 <score>
2394 0-10 (integer only)
2395 </score>
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

```

### RedSage Open-ended Q&A LLM as Judge Prompt Template

```

[System Prompt]
Question:
```
{question}
```

Reference Answer:
```
{reference_answer}
```

Model Answer:
```
{model_answer}
```

```

**Qualitative Results of RedSage OpenQA.** We present three RedSage OpenQA examples that span cybersecurity frameworks, offensive skills, and tool usage. In the Olympic Destroyer attribution case shown in Fig. 12, RedSage 8B DPO correctly identifies the Sandworm team, while baseline models misattribute the malware to other Russian APT groups. For the CSP bypass example in Fig. 13 and the Koadic tool-usage example in Fig. 14, RedSage 8B DPO accurately recognizes iframe-based CSP evasion and generates the exact Koadic command line. In contrast, the baselines omit key details or produce malformed commands. These cases illustrate that RedSage exhibits stronger grounding in authoritative cybersecurity sources and improved precision in operational reasoning.

**Qualitative Results of LLM-as-Judge.** To further illustrate the differences captured by our LLM-as-Judge pipeline, we include qualitative evaluations comparing RedSage with the baseline model using the tool-based question shown in Fig. 14. As shown in Fig. 15, the judge marks RedSage’s answer as fully correct, assigns a perfect score, and highlights the precise command construction and clear supporting explanations. In contrast, the baseline model receives a failing correctness label and a substantially lower score because it uses an incorrect command-line flag, even though its surrounding explanation is detailed. These paired results emphasize the sensitivity of our evaluation framework to fine-grained correctness, particularly in cybersecurity scenarios where small syntactic deviations can lead to incorrect or unsafe tool behavior.

2430

2431

2432

2433

2434

2435

2436

2437

2438

2439

2440

2441

2442

2443

2444

2445

2446

2447

2448

2449

2450

2451

2452

2453

2454

2455

2456

2457

2458

2459

2460

2461

2462

2463

2464

2465

2466

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

**Cybersecurity Framework:**

**Data Source:**  
<https://attack.mitre.org/software/S0365/>

**Question:**  
Which threat group is known to have used the Olympic Destroyer malware?

**Reference Answer**

The threat group known to have used the Olympic Destroyer malware is **Sandworm Team**. This group was responsible for deploying the malware against the 2018 Winter Olympics in Pyeongchang, South Korea. The main purpose of the malware was to render infected computer systems inoperable by leveraging various native Windows utilities and API calls to carry out destructive tasks. Olympic Destroyer has worm-like features to spread itself across a computer network, maximizing its destructive impact.

**Foundation-Sec-8B-Instruct**

The threat group known to have used the Olympic Destroyer malware is **APT28**, also known as **Fancy Bear**. This group is believed to be associated with the Russian government. The malware was used in a cyberattack on the 2018 Winter Olympics in Pyeongchang, South Korea. It's important to note that APT28 has been linked to various other cyber espionage and sabotage activities over the years. For the most accurate and up-to-date information, refer to reliable sources such as the MITRE ATT&CK framework or reports from cybersecurity firms like FireEye, Symantec, or Kaspersky.

**RedSage-8B-DPO**

The threat group known to have used the Olympic Destroyer malware is a group referred to as **"Sandworm"**, also known as **"Group 74"** or **"Iron Viking."** This group is believed to be a part of the Russian military intelligence agency, the GRU. They are notorious for their cyber attacks, including the Olympic Destroyer attack on the 2018 Winter Olympics in Pyeongchang, South Korea.

**Llama-Primus-Base**

The threat group known to have used the Olympic Destroyer malware is **Fancy Bear**, also known as **APT28, Sofacy, and Sednit**. They are a Russian state-sponsored cyber espionage group believed to be associated with the Russian military intelligence agency GRU. Olympic Destroyer was used in a cyberattack on the 2018 Winter Olympics in Pyeongchang, South Korea, causing disruptions to the event's opening ceremony.

Figure 12: Knowledge framework example from RedSage OpenQA. For the Olympic Destroyer question, RedSage 8B DPO matches the reference attribution to the Sandworm team, while baseline models misattribute it to different APT groups. *Best viewed in Zoom.*

2464

2465

2466

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

**Cybersecurity Offensive Skill:**

**Data Source:**  
<https://github.com/b4rdia/HackTricks/blob/master/pentesting-web/content-security-policy-csp-bypass-self-+-unsafe-inline-with-iframes.md>

**Question:**  
Analyze the provided code and identify the specific method used to bypass CSP protections through iframe injection. What makes this technique effective against the described CSP configuration?

```
```javascript
frame = document.createElement("iframe")
frame.src = "/css/bootstrap.min.css"
document.body.appendChild(frame)
script = document.createElement("script")
script.src = "/example.com/csp.js"
window.frames[0].document.head.appendChild(script)
```

```

**RedSage-8B-DPO**

The provided code bypasses Content Security Policy (CSP) protections through **“iframe injection”** by leveraging the following method:

**### Technique Breakdown:**

- 1. “Iframe Creation”:**
  - An `<iframe>` element is dynamically created using `document.createElement("iframe")`.
  - The `src` attribute is set to `"/css/bootstrap.min.css"`, which is a **“non-executable resource”** (a CSS file). This step is likely a placeholder to ensure the iframe is loaded and added to the DOM.
- 2. “Appending the Iframe”:**
  - The iframe is appended to the document body using `document.body.appendChild(frame)`, making it visible and triggering its loading.
- 3. “Script Injection via Iframe Window”:**
  - A `<script>` element is created with `script.src = "/example.com/csp.js"`, which is the malicious script intended to bypass CSP.
  - The script is appended to the `<head>` of the iframe's document using `window.frames[0].document.head.appendChild(script)`. This exploits the fact that the iframe's `src` is a non-executable resource, allowing scripts to be injected into its `document` context.

**##**

**### Why This Technique Works Against CSP:**

- 1. “CSP Violation Detection Delay”:**
  - CSP violations are typically reported via the `‘report-uri’` or `‘onerror’` events. However, if the iframe's `src` is a non-executable resource (e.g., CSS), the browser may **“not immediately flag the script injection”** as a violation. This creates a window of opportunity for the malicious script to execute before the CSP enforces restrictions.
- 2. “Bypassing ‘script-src’ Restrictions”:**
  - If the CSP's `‘script-src’` directive does not explicitly allow `‘//example.com’`, the direct script load `‘script.src = ‘//example.com/csp.js’` would normally fail. However, by injecting the script into the **“iframe’s document context”**, the browser may **“not enforce the ‘script-src’ directive”** for the iframe's content, allowing the script to execute.
- 3. “Exploiting Iframe’s Isolation”:**
  - Iframes create a separate browsing context with their own `‘document’` and `‘window’` objects. By injecting the script into the iframe's document, the attack avoids directly violating the main page's CSP, as the iframe's context is treated independently.

Figure 13: Offensive skill example analyzing JavaScript that bypasses Content Security Policy. RedSage 8B DPO correctly identifies iframe injection as the evasion technique and explains how each step circumvents the configured `script-src` directive. *Best viewed in Zoom.*

2464

2465

2466

2467

2468

2469

2470

2471

2472

2473

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

**RedSage: LLM-as-Judge Output**

```
<analysis>
The model answer provided is not only factually correct but also offers additional helpful information and context. It accurately constructs the command line to start Koadic with the specified autorun file and restore option. The command line provided in the model answer matches the reference answer exactly, which is `koadic --autorun autorun_commands.txt --restore restore_data.json`. This indicates that the model has correctly understood the requirements and has provided the precise command needed to achieve the desired outcome.

Correctness
The model answer is correct because it accurately reflects the command line needed to start Koadic with the specified parameters. There are no factual errors, hallucinations, or contradictions in the provided command line. The additional information about the command assuming `koadic` is in the system's PATH and the explanation of the flags (`--autorun` and `--restore`) are accurate and useful but do not affect the correctness of the command itself.

```

2484  
 2485  
 2486  
 2487  
 2488  
 2489  
 2490  
 2491  
 2492  
 2493  
 2494  
 2495  
 2496  
 2497  
 2498  
 2499  
 2500  
 2501  
 2502  
 2503  
 2504  
 2505  
 2506  
 2507  
 2508  
 2509  
 2510  
 2511  
 2512  
 2513  
 2514  
 2515  
 2516  
 2517  
 2518  
 2519  
 2520  
 2521  
 2522  
 2523  
 2524  
 2525  
 2526  
 2527  
 2528  
 2529  
 2530  
 2531  
 2532  
 2533  
 2534  
 2535  
 2536  
 2537

**Cybersecurity Tool:**

**Data Source:**  
<https://www.kali.org/tools/koadic/>

**Question:**  
 Construct a command line to start Koadic with an autorun file named 'autorun\_commands.txt' and enable the restore option using the file 'restore\_data.json'.

**Reference Answer**

```
```sh
koadic --autorun autorun_commands.txt --restore restore_data.json
```

```

**RedSage-8B-DPO**

To start Koadic with an autorun file named 'autorun\_commands.txt' and enable the restore option using the file 'restore\_data.json', you can use the following command line:

```
```bash
koadic --autorun autorun_commands.txt --restore restore_data.json
```

```

This command assumes that 'koadic' is in your system's PATH. If not, you may need to provide the full path to the Koadic executable. The '--autorun' flag specifies the file containing commands to run automatically, and the '--restore' flag specifies the file used for restoring session data.

✓
✗

**Foundation-Sec-8B-Instruct**

To start Koadic with an autorun file named 'autorun\_commands.txt' and enable the restore option using the file 'restore\_data.json', you will need to use the command line interface. Here's how you can do it:

```
```bash
koadic --autorun autorun_commands.txt --restore restore_data.json
```

```

**Explanation:**

- 'koadic': This is the command to start the Koadic framework.
- '--autorun autorun\_commands.txt': This option specifies the autorun file, which contains the list of commands to be executed.
- '--restore restore\_data.json': This option enables the restore feature, allowing Koadic to use the specified JSON file for restoring previous session data.

<Trimmed for brevity>

Figure 14: Tool-usage example for the Koadic framework. RedSage 8B DPO provides the exact command, while the baseline model produces a non-matching command. *Best viewed in Zoom.*

2501           **#### Answer Quality Score**  
 2502           The answer quality score should be high because the model answer  
 2503           is not only correct but also provides additional useful  
 2504           information that could help a user understand how to use the  
 2505           command and troubleshoot potential issues (like the executable not  
 2506           being in the PATH).  
 2507           **#### Helpfulness**  
 2508           The model answer is very helpful because it directly answers the  
 2509           question with the correct command line and provides supplementary  
 2510           information that could aid in successfully executing the command.  
 2511           **#### Relevance**  
 2512           The answer is highly relevant as it directly addresses the  
 2513           question asked without including unnecessary information. The  
 2514           command line and the explanations provided are focused on the task  
 2515           of starting Koadic with an autorun file and enabling the restore  
 2516           option.  
 2517           **#### Depth**  
 2518           The answer shows a good depth of understanding of the command and  
 2519           its parameters. It not only provides the command but also explains  
 2520           the purpose of the flags used, which indicates a level of  
 2521           comprehension beyond mere recall of command syntax.  
 2522           **#### Level of Detail**  
 2523           The level of detail is appropriate for the question asked. The  
 2524           model provides enough detail for someone to understand how to use  
 2525           the command and what each part of the command does, without  
 2526           overwhelming them with unnecessary information.  
 2527           </analysis>  
 2528  
 2529           <correctness>  
 2530           True  
 2531           </correctness>  
 2532  
 2533           <score>  
 2534           10  
 2535           </score>

2538  
 2539  
 2540  
 2541     <analysis>  
 2542     The model answer provides a detailed explanation of how to start  
 2543     Koadic with an autorun file and enable the restore option. It  
 2544     breaks down the command into its components and explains the  
 2545     purpose of each part. However, the command provided in the model  
 2546     answer is slightly different from the reference answer. The model  
 2547     answer uses the short option `‐a‐ for autorun, whereas the  
 2548     reference answer uses the long option `‐‐autorun‐.  
 2549  
 2550     #### Correctness  
 2551     The model answer is almost correct, but it uses a different option  
 2552     for specifying the autorun file. The correct option, as per the  
 2553     reference answer, is `‐‐autorun‐, not `‐a‐. This discrepancy makes  
 2554     the model answer factually incorrect.  
 2555  
 2556     #### Answer Quality Score  
 2557     Despite the factual error, the model answer provides a clear and  
 2558     detailed explanation of the command and its components. It also  
 2559     includes important notes about the ethical use of Koadic and the  
 2560     importance of proper permissions and access.  
 2561  
 2562     #### Helpfulness  
 2563     The model answer is helpful in explaining the purpose of each part  
 2564     of the command and providing context for its use.  
 2565  
 2566     #### Relevance  
 2567     The model answer is highly relevant to the question, as it  
 2568     directly addresses the task of starting Koadic with an autorun  
 2569     file and enabling the restore option.  
 2570  
 2571     #### Depth  
 2572     The model answer provides a good depth of explanation, breaking  
 2573     down the command into its components and explaining the purpose of  
 2574     each part.  
 2575  
 2576     #### Level of Detail  
 2577     The model answer provides a sufficient level of detail, including  
 2578     explanations of the command options and important notes about  
 2579     ethical use and permissions.  
 2580     </analysis>  
 2581  
 2582     <correctness>  
 2583     False  
 2584     </correctness>  
 2585  
 2586     <score>  
 2587     3  
 2588     </score>  
 2589  
 2590  
 2591

Figure 15: Qualitative LLM-as-Judge outputs comparing RedSage and the baseline model.

2592 C.3 CYBERSECURITY BENCHMARKS  
2593

2594 **CyberMetric (CyMtc).** CyberMetric evaluates general cybersecurity knowledge via multiple-  
2595 choice questions with four options, curated from authoritative sources such as NIST publications,  
2596 RFCs, books, and research papers using a retrieval-augmented generation pipeline. The collection  
2597 is released in several sizes, and we use the 500-item split that was fully verified by human experts.  
2598 Items span nine topical areas that include cryptography, reverse engineering, and risk assessment.  
2599 Models are scored with standard MCQ accuracy. (Tihanyi et al., 2024)

2600 **SecBench (ScBen).** SecBench is a large multi-dimensional benchmark for cybersecurity that in-  
2601 cludes both MCQs and short-answer questions, covers two capability levels (knowledge retention  
2602 and logical reasoning), and is available in Chinese and English. Questions were sourced from open  
2603 materials and a curated contest, and short-answer evaluation is supported by an LLM-based grader.  
2604 In our study we use the English MCQ subset and report accuracy. (Jing et al., 2024)

2605 **MMLU Computer Security (MMLU-CSec).** MMLU is a 57-subject multiple-choice test that  
2606 measures broad academic and professional knowledge. We evaluate on the Computer Security  
2607 subject, which contains MCQs covering practical and theoretical topics such as network security  
2608 and cryptography. Following common practice for MMLU-style evaluation, we report accuracy.  
2609 (Hendrycks et al., 2021b)

2610 **SECURE.** SECURE targets applied cybersecurity with datasets built from MITRE ATT&CK,  
2611 CWE, CVE, and related ICS advisories, organized into three knowledge types: extraction, under-  
2612 standing, and reasoning. We use the MCQ-style subsets MAET (MITRE ATT&CK Extraction),  
2613 CWET (Common Weakness Extraction), and KCV (Knowledge test on Common Vulnerabilities).  
2614 The authors manually refined the pools by removing or fixing flawed questions. We evaluate with  
2615 MCQ accuracy. (Bhusal et al., 2024)

2616 **CTI-Bench.** CTI-Bench focuses on cyber threat intelligence and provides four tasks: CTI-MCQ  
2617 for knowledge of CTI standards and practices; CTI-RCM for mapping CVE descriptions to one or  
2618 more CWE root causes; CTI-VSP for predicting CVSS v3 base vectors and scores; and CTI-ATE for  
2619 extracting MITRE ATT&CK attack techniques from natural language incident descriptions. While  
2620 VSP and ATE are typically evaluated with regression and F1 metrics, respectively, in our study we  
2621 only use accuracy across all subsets for consistent aggregation. (Alam et al., 2024)

2622 **SecEval (ScEva).** SecEval is a domain-focused benchmark of more than two thousand MCQs  
2623 spanning nine areas that include software, application, system, web, cryptography, memory safety,  
2624 network security, and penetration testing. Questions were constructed from textbooks, official docu-  
2625 mentation, and standards using GPT-4 prompting, with quality control to remove invalid items. We  
2626 evaluate with MCQ accuracy on the full set. (Li et al., 2023)

2627 C.4 GENERAL LLM BENCHMARKS  
2628

2629 **ARC-Challenge (ARC-C).** ARC-C is the challenge split of the AI2 Reasoning Challenge, a set  
2630 of grade-school science multiple-choice questions curated to require nontrivial reasoning and back-  
2631 ground knowledge. The challenge subset specifically contains items that defeat simple retrieval and  
2632 co-occurrence baselines, making it a strong discriminator of reasoning beyond surface cues. We  
2633 evaluate with standard MCQ accuracy as used by leaderboard implementations. (Clark et al., 2018)

2634 **HellaSwag (HSwag).** HellaSwag tests grounded commonsense inference via sentence comple-  
2635 tion. Each example presents a short context and four candidate endings that describe plausible next  
2636 events in physical or social scenarios. The dataset was adversarially filtered to foil strong language  
2637 models while remaining trivial for humans, which sharpens its discriminative power. Performance  
2638 is reported as multiple-choice accuracy. (Zellers et al., 2019)

2639 **TruthfulQA (TQA).** TruthfulQA measures whether models avoid widespread misconceptions and  
2640 misleading patterns by answering with factually truthful content across 38 categories such as health,  
2641 law, and finance. It provides both generative prompts and multiple-choice variants. Following  
2642 common leaderboard practice, we use the multiple-choice setting and report accuracy to ensure  
2643 comparability across models. (Lin et al., 2021)

2644 **MMLU.** MMLU evaluates broad knowledge and reasoning across 57 academic and professional  
2645 subjects that range from elementary mathematics and U.S. history to computer science and law. Each

2646 subject consists of four-option multiple-choice items designed to test recall, conceptual understanding,  
 2647 and problem solving. Scores are aggregated as average accuracy across subjects. (Hendrycks  
 2648 et al., 2021a)

2649 **WinoGrande (WinoG).** WinoGrande is a large adversarial variant of the Winograd Schema Chal-  
 2650 lenge that assesses commonsense reasoning through pronoun resolution. Each example requires  
 2651 selecting which of two candidate nouns a pronoun refers to, with items constructed to reduce an-  
 2652 notation artifacts and shallow heuristics. Evaluation follows leaderboard protocol using accuracy.  
 2653 (Sakaguchi et al., 2020)

2654 **GSM8K.** GSM8K is a collection of 8.5K carefully authored grade-school math word problems  
 2655 that require multi-step arithmetic reasoning. Problems are linguistically diverse and designed to  
 2656 encourage chain-of-thought solutions, yet the final target is a short numeric answer. We report  
 2657 exact-match accuracy on the final answer, consistent with leaderboard settings. (Cobbe et al., 2021)

2658 **IFEval.** IFEval evaluates instruction following using prompts that contain verifiable constraints  
 2659 such as minimum length, required keywords, or structural requirements. Each prompt includes one  
 2660 or more constraints that can be programmatically checked, yielding objective pass/fail signals with-  
 2661 out human grading. We report the mean compliance rate across all constraints, i.e., the percentage  
 2662 of constraints satisfied. (Zhou et al., 2023)

## D ADDITIONAL EVALUATION RESULTS

### D.1 LARGER MODEL SCALING

2667 To assess the scalability of our data curation and augmentation pipeline, we conducted a reduced-  
 2668 scope experiment using Qwen3-32B. We applied QLoRA fine-tuning ( $\approx 1\%$  trainable parameters)  
 2669 on a partial dataset consisting of the curated RedSage-Seed subset (excluding RedSage-Dump) and  
 2670 50% of RedSage-Conv. Despite using only a fraction of the full training data and a lightweight  
 2671 adaptation method, the resulting 32B model achieved consistent gains across both the RedSage-  
 2672 MCQ benchmark (Table 13) and a suite of cybersecurity evaluations (Table 14). Notably, the training  
 2673 loss continued to decrease throughout the run, suggesting that full-data, full-parameter fine-tuning  
 2674 would yield even larger improvements. These findings indicate that the RedSage data curation and  
 2675 augmentation methodology transfers effectively to larger models, underscoring its scalability and  
 2676 potential to advance cybersecurity LLM development.

2677 Table 13: RedSage-MCQ (0-shot) scaling experiment. Values are accuracy (%). Abb: Gen = General,  
 2678 Frm = Frameworks, Off = Offensive Skills, CLI = Command-line Tools, Kali = Kali Tools.

| Model Name               | Macro        | Knowledge    |              | Skill        |              | Tools        |      |
|--------------------------|--------------|--------------|--------------|--------------|--------------|--------------|------|
|                          |              | Acc          | Gen          | Frm          | Off          | CLI          | Kali |
| Qwen3-8B                 | 81.85        | 80.46        | 78.82        | 86.16        | 83.92        | 75.56        |      |
| Qwen3-32B                | 85.40        | 84.08        | 82.32        | 89.00        | 87.60        | 80.40        |      |
| RedSage-8B-Ins           | 85.73        | 84.20        | 84.98        | 89.06        | 86.80        | 80.30        |      |
| RedSage-32B-LoRA-Ins-0.5 | <b>87.53</b> | <b>85.68</b> | <b>85.04</b> | <b>91.46</b> | <b>88.76</b> | <b>82.78</b> |      |

2688 Table 14: Related Cybersecurity Benchmarks (0-shot) scaling experiment. Values are Accuracy (%).  
 2689 Best results are shown in bold.

| Model Name               | Mean         | CTI-Bench    |              | CyMtc        | MMLU         | ScBen        | ScEva        | SECURE       |              |              |
|--------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                          |              | MCQ          | RCM          |              |              |              |              | MCQ          | CWET         | KCV          |
| Qwen3-8B                 | 75.71        | 62.76        | 54.00        | 88.60        | 76.00        | 73.26        | 65.46        | 88.11        | 87.42        | 85.75        |
| Qwen3-32B                | 82.31        | 70.04        | 65.60        | 91.80        | <b>84.00</b> | <b>84.23</b> | 76.23        | 89.46        | <b>88.72</b> | 90.06        |
| RedSage-8B-Ins           | 81.30        | 70.56        | <b>76.70</b> | 89.80        | 78.00        | 79.91        | 72.48        | 91.45        | 81.34        | 91.47        |
| RedSage-32B-LoRA-Ins-0.5 | <b>82.85</b> | <b>71.64</b> | 66.10        | <b>93.40</b> | <b>84.00</b> | 83.77        | <b>78.30</b> | <b>92.18</b> | 83.29        | <b>92.97</b> |