

# LLM-based Vulnerable Code Augmentation: Generate or Refactor?

Dyna Soumhane Ouchebara<sup>1</sup> and Stéphane Dupont<sup>1</sup> \*

1- University of Mons - Computer science department  
Mons - Belgium

**Abstract.** Vulnerability code-bases often suffer from severe imbalance, limiting the effectiveness of Deep Learning-based vulnerability classifiers. Data Augmentation could help solve this by mitigating the scarcity of under-represented CWEs. In this context, we investigate LLM-based augmentation for vulnerable functions, comparing controlled generation of new vulnerable samples with semantics-preserving refactoring of existing ones. Using Qwen2.5-Coder to produce augmented data and CodeBERT as a vulnerability classifier on the SVEN dataset, we find that our approaches are indeed effective in enriching vulnerable code-bases through a simple process and with reasonable quality, and that a hybrid strategy best boosts vulnerability classifiers’ performance.

## 1 Introduction

Deep learning models for software engineering tasks depend on large, diverse, and well-labeled datasets, yet collecting and annotating such data is often expensive and time-consuming. Data Augmentation (DA) mitigates data scarcity by synthetically expanding the training set, improving generalization and robustness without additional labeling effort. However, unlike in computer vision and NLP, DA for source code is challenging: code is a structured, executable artifact whose meaning relies on strict syntactic and semantic constraints, so small edits can easily break compilation or alter behavior. Effective Code Augmentation (CA) must therefore preserve syntax and semantics while allowing for enough diversity. This challenge is amplified in software vulnerability detection, where datasets are highly imbalanced both between vulnerable and safe classes and across vulnerability types (some being strongly under-represented). In this setting, DA can enrich minority classes with diverse yet semantically consistent variants, consequently reducing the impact of imbalance on vulnerability detection models. Recently, large language models (LLMs) have been widely adopted for code-related tasks such as generation, refactoring, documentation, and repair, offering a powerful new mechanism for code augmentation. Our work explores this line of research, and aims to answer the following questions:

**RQ1:** *How effective is LLM-based Code Augmentation in enriching vulnerable code in highly imbalanced vulnerability code-bases?*

**RQ2:** *Can LLM-based Code Augmentation boost the performance of Deep Learning models in the vulnerability classification task?*

---

\*This study was funded by CyberExcellence project of CyberWal program by Digital Wallonia.

## 2 Related Work

Zhuo et al. [1] categorize Code Augmentation methods into rule-based transformations (refactorings, renamings), model-based generation (GANs, pre-trained Transformers), and example interpolation methods (Mixup). Dong et al. [2] directly apply NLP augmentation techniques to code and report improved performance of models, despite the risk of broken syntax. MIXCODE [3] generates semantics-preserving refactorings and then linearly mixes (Mixup) their embeddings to regularize code classifiers. BUGLAB [4] uses a learned model that applies bug-inducing edits to benign code, producing synthetic buggy examples. In the vulnerability detection setting, MPDA [5] augments imbalanced vulnerability datasets by combining classical oversampling methods, a GAN that synthesizes new vulnerable samples, and fuzzy sampling to augment minority-class instances. Qi et al. [6] transform existing vulnerable and safe functions via handcrafted semantic-preserving edits to produce new variants that retain the original vulnerability labels. FGVulDet [7] employs an approach that only perturbs code regions deemed unrelated to the vulnerable zone, so that new samples share the same vulnerability while differing in surrounding context. VULGEN and VGX [8, 9] mine vulnerability-introducing edit patterns from real patches and reapplies these patterns at predicted suitable locations in benign code to generate realistic vulnerabilities. More recently, Deng et al. [10] explored LLMs for vulnerable data generation. They prompt GPT4 to synthesize vulnerable versions of safe functions for under-represented CWE types, and to further self-filter the outputs based on quality checks. Our approach follows this LLM-based line but explores controlled generation and refactoring based on existing vulnerable functions from a reliable dataset to produce structurally and contextually diverse variants that preserve vulnerability semantics.

## 3 Methodology

### 3.1 Generation-based data augmentation

In the first generation strategy, we synthesize entirely new vulnerable functions using an instruction-tuned LLM. For each vulnerability type (CWE), we randomly select  $m$  real-world functions from the training set and use them as examples in a few-shot prompt. The model is asked to generate  $n$  new vulnerable functions per prompt, and this process is repeated  $k$  times with different exemplar subsets, yielding  $n \times k$  synthetic functions per vulnerability type.

Our **prompting scheme** relies on a fixed system message, which establishes the model’s role and all global constraints, and a variable user message which provides instance-level information. The **system message** assigns the model the role of an *expert in C/C++ and Python programming and software vulnerabilities* and enforces strict output constraints. The model must return C/C++ or Python code only, with no comments, explanations, markdown, or language tags, that must resemble real-world industrial C/C++ and Python and avoid

toy identifiers. Function, variable, and type names are required to follow realistic project-style conventions. The **user message** begins with a description of the CWE. The model is instructed to generate  $n$  independent function definitions, separated by markers such as “func 1”, “func 2”, etc. Structural constraints are imposed: each function must contain 20–150 non-empty lines, use project-like types, and include a vulnerability that is embedded in realistic logic rather than being the function’s sole purpose. Finally, the prompt provides  $m$  examples (from training set), and the model is instructed not to copy those but only to use them as guidance regarding naming conventions, code organization, and how the vulnerability typically appears in practice.

Once the functions are generated, we proceed to quality checks. In the first place, we verify the **syntactical quality** of the generated samples by passing them through a C/C++ or Python parser. In a second place, we verify their **label quality**. Ideally, this would be conducted by security experts, but as a start for this research, we verified them using GPT-5.1 Thinking, which is one of the strongest LLMs available (at submission-time). For this, we give a randomly picked subset of  $q$  generated functions to the LLM, and ask it to verify if it indeed contains a vulnerability of the given type or not.

### 3.2 Refactoring-based data augmentation

The second strategy produces augmented samples by refactoring functions already present in the dataset. For each vulnerability type and for every corresponding function, we prompt the LLM to generate  $n$  refactored variants. We rely on **18 refactoring techniques** commonly used in prior work. These transformations alter code structure and surface form without changing semantics or the underlying vulnerability. This is a categorized list of techniques: Renaming (API, Arguments, Local variable and Method renaming), Adding unused elements (Arguments or Local variable adding), Dead code insertion (Dead for/if/else/switch/while adding or Duplication), Logic-preserving rewrites of control and expressions (For loop/If enhancement, Return optimal, Plus zero), Safety / robustness guards (Filed enhancement), Logging (Prints).

In terms of **prompting scheme**, the **system message** instructs the model to act as “*an expert in C/C++ and Python programming, code refactoring, and software vulnerabilities*” and to generate  $n$  refactored versions of the given function. It must preserve the function’s semantics, parameter list, return type, and vulnerability. To avoid inadvertent vulnerability repair, the prompt forbids removing dangerous operations. The system message restricts the refactoring space to the 18 techniques listed above; each generated function must apply at least two distinct transformations. Output must consist solely of C/C++ or Python code with realistic identifiers and without comments or explanations. As for the **user message**, it defines the requested number of refactorings, the required output format, the vulnerability description, and the function to be refactored. No examples are provided, making this a zero-shot prompting setup.

Regarding quality verifications, we first check the **syntactical quality** just like for the previous approach. Then, we verify their **refactoring quality** by asking GPT-5.1 Thinking to check whether the quality of a randomly picked subset of  $q$  refactored versions complies with the constraints and expected complexity. We note that the label quality of generated samples is not verified, since it is inherently preserved by design in this approach.

## 4 Experimental Setup

To evaluate the performance of our proposed approaches, we chose to apply them on the **SVEN Dataset** [11]. This dataset was created in 2023 by manually inspecting security-related commits from three prior benchmarks (BigVul, CrossVul and VUDENC) and only keeping those with no quality issues, and the most critical CWE types. The statistics of SVEN’s training set are in Table 1 (after splitting their training data into 80% train and 20% validation).

cwe-89	cwe-125	cwe-78	cwe-476	cwe-416	cwe-22	cwe-787	cwe-79	cwe-190
141	107	69	60	45	42	41	39	32

Table 1: Number of functions per CWE in SVEN Training set.

Concerning the **Models**, we choose **Qwen2.5-Coder-32B**<sup>1</sup> for augmented data generation, for its high rank in code LLM benchmarks<sup>23</sup> (at the time we conducted our experiments) especially on C/C++ and Python, and its suitable size; and **CodeBERT** as the vulnerability classifier for its well-established code representation capabilities as well as its lightness (allowing for quick fine-tuning).

In terms of **Evaluation Metrics**, our augmentation approaches are assessed through the number of new vulnerable samples we can generate, the percentage of augmentation obtained for each class, the average time per generated sample, and the quality of augmented data (syntax, refactors, labels). As for the vulnerability classifier, we evaluate it by reporting the macro average f1 for an overall performance indication on all classes.

Our **Technical Setup** comprises 3 A100-SXM4 GPUs and a 16GB RAM.

## 5 Evaluation

To answer **RQ1**, we review Table 2 which reports the assessment of our two augmentation approaches. For the Generation approach, we generated 10 new functions at a time, and stopped at a target of 500 functions in total per class.

<sup>1</sup><https://huggingface.co/Qwen/Qwen2.5-Coder-32B-Instruct>

<sup>2</sup><https://huggingface.co/spaces/bigcode/bigcode-models-leaderboard>

<sup>3</sup><https://huggingface.co/spaces/bigcode/bigcodebench-leaderboard>

Approach	N° of samples	% of augmentation	Average time per sample	Syntax quality	Label quality	Refactor quality
Generation	3348	581%	13.38s	98.5%	0%	/
Refactoring	1224	213%	59.08s	79.7%	/	100%

Table 2: Assessment metrics for our two augmentation approaches.

We could increase the dataset size by 581%, with an average speed of 13.38s per generated sample, which looks reasonable under our resources. In terms of quality, 98.5% samples proved syntactically correct. On the other hand, the label quality check, measured on a random subset of 10 generated functions per class, surprisingly gave a 0% score. We, however, also checked the label quality of the original dataset and found 0% for most CWEs. This means that: either the training data is too complex for GPT5.1 Thinking to find the vulnerabilities, or the dataset’s label quality is questionable (which was not expected, since it has been manually curated by the authors). Further investigation will help us confirm our hypotheses. For the Refactoring approach, we generated 10 new functions at a time, and stopped at a target of 200 functions in total per class. We could increase the dataset size by 213% with a speed of 59.08s per generated sample (much slower than the first approach). The syntax quality of 79.7% correct functions is acceptable, and the refactoring quality (measured just like label quality above) is perfect.

Training data	Original data	Generation augmented	Refactoring augmented	Both augmentations
Macro F1	0.62	0.64	0.60	0.67

Table 3: Macro-average F1 score for different training data.

To answer **RQ2**, we review Table 3 which reports the performance of our vulnerability classifier before and after augmentation. The generation-based augmentation indeed improves the the classifier performance, with an overall Macro F1 score of 0.64 vs 0.62 on the original data. This increase is even more noticeable if we look at the minority classes, such as cwe-22 increased by 18% and cwe-190 by 8%. The refactoring-based augmentation, on the other hand did not bring any overall improvement (though a few minority CWEs did improve, cwe-022 by 18% and cwe-416 by 4%). Applying both augmentations proved to be the most helpful performance-wise, with an overall improvement of 5% on Macro F1, and a clear boost for all CWEs (up to 18% boost).

## 6 Conclusion

In this work, we proposed two LLM-based vulnerable code approaches: the first synthesizes entirely new vulnerable functions based on examples, and the second produces semantics and vulnerability preserving refactored versions of existing functions. Our experiments allowed us to answer our research questions. For

**RQ1**, we conclude that our LLM-based augmentation approaches are indeed effective in enriching vulnerable code-bases through a simple process and in a reasonable time, with great syntax and refactoring quality, though the label quality proved questionable. For **RQ2**, we observe that LM-based augmentation can indeed boost the performance of vulnerability classifiers, and deduce that the best strategy is a hybrid approach that applies both few-shot generation and refactoring.

## References

- [1] Terry Yue Zhuo, Zhou Yang, Zhensu Sun, Yufei Wang, Li Li, Xiaoning Du, Zhenchang Xing, and David Lo. Source code data augmentation for deep learning: A survey. *arXiv preprint arXiv:2305.19915*, 2023.
- [2] Zeming Dong, Qiang Hu, Yuejun Guo, Zhenya Zhang, Maxime Cordy, Mike Papadakis, Yves Le Traon, and Jianjun Zhao. Boosting source code learning with text-oriented data augmentation: an empirical study. *Empirical Software Engineering*, 30(3):68, 2025.
- [3] Zeming Dong, Qiang Hu, Yuejun Guo, Maxime Cordy, Mike Papadakis, Zhenya Zhang, Yves Le Traon, and Jianjun Zhao. Mixcode: enhancing code classification by mixup-based data augmentation. In *2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 379–390. IEEE, 2023.
- [4] Miltiadis Allamanis, Henry Jackson-Flux, and Marc Brockschmidt. Self-supervised bug detection and repair. *Advances in Neural Information Processing Systems*, 34:27865–27876, 2021.
- [5] Feiqiao Mao, Yingxiang Yuan, Xingyang Du, Li Gao, and Zhihua Du. Mpda: a data augmentation approach to improve deep learning for software vulnerability detection. *Empirical Software Engineering*, 30(5):140, 2025.
- [6] Weiliang Qi, Jiahao Cao, Darsh Poddar, Sophia Li, and Xinda Wang. Enhancing pre-trained language models for vulnerability detection via semantic-preserving data augmentation. In *International Conference on Security and Privacy in Communication Systems*, pages 184–203. Springer, 2024.
- [7] Shangqing Liu, Wei Ma, Jian Wang, Xiaofei Xie, Ruitao Feng, and Yang Liu. Enhancing code vulnerability detection via vulnerability-preserving data augmentation. In *Proceedings of the 25th ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems*, pages 166–177, 2024.
- [8] Yu Nong, Yuzhe Ou, Michael Pradel, Feng Chen, and Haipeng Cai. Vulgen: Realistic vulnerability generation via pattern mining and deep learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2527–2539. IEEE, 2023.
- [9] Yu Nong, Richard Fang, Guangbei Yi, Kunsong Zhao, Xiapu Luo, Feng Chen, and Haipeng Cai. Vgx: Large-scale sample generation for boosting learning-based software vulnerability analyses. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pages 1–13, 2024.
- [10] Xiao Deng, Fuyao Duan, Rui Xie, Wei Ye, and Shikun Zhang. Improving long-tail vulnerability detection through data augmentation based on large language models. In *2024 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 262–274. IEEE, 2024.
- [11] Jingxuan He and Martin Vechev. Large language models for code: Security hardening and adversarial testing. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 1865–1879, 2023.