

# CHEMFUZZ: Large Language Models-assisted Fuzzing for Quantum Chemistry Software Bug Detection

Feng Qiu<sup>1,2†</sup>, Pu Ji<sup>3†</sup>, Baojian Hua<sup>1,2\*</sup>, and Yang Wang<sup>1,2\*</sup>

<sup>1</sup>School of Software Engineering, University of Science and Technology of China, China

<sup>2</sup>Suzhou Institute for Advanced Research, University of Science and Technology of China, China

<sup>3</sup>Computer Science Department, The Johns Hopkins University, USA

bgkqf@mail.ustc.edu.cn, pji7@jh.edu, bjhua@ustc.edu.cn, angyan@ustc.edu.cn

†Co-first authors. \* Corresponding authors.

*Abstract*—Quantum chemistry software implements the first principle quantum computation and is indispensable in both scientific research and chemical industries. Any bugs in such software will lead to serious consequences, thus defeating its trustworthiness and reliability. However, bug detection techniques for such software have not been fully investigated. In this paper, to fill this gap, we propose a novel approach to fuzz quantum chemistry software with the aid of Large Language Models (LLMs). Our basic idea is utilize LLMs to mutate and generate syntactic and semantic valid input files from seed inputs, by proving valuable domain-specific knowledge of chemistry. With this basic idea, we have designed and implemented CHEMFUZZ, a fully automatic fuzzing framework to fuzz quantum chemistry software for bugs. Our evaluation of CHEMFUZZ leverages popular LLMs including GPT3.5, Claude-2, and Bart as test oracles to generate parameters to mutate inputs and analyze computation results. CHEMFUZZ detected 40 unique bugs, which have been classified and reported to developers, with a code coverage of 17.4%.

*Keywords*—Quantum Chemistry Software, Fuzzing, Large Language Models, Security Test

## 1. INTRODUCTION

We are entering an era of AI for science with scientific disciplines implemented on top of software infrastructures. Specifically, quantum chemistry software, computing equations and approximations derived from the postulates of quantum mechanics, have been extensively studied (*e.g.*, Gaussian[20], Siesta [21], Quantum Espresso [27], and Abacus [36]), with wide adoptions in diverse domains (*e.g.*, chemical industry [2], drug innovation [8], and material science [26]). Given the important applications of such software in science, they should be trustworthy and reliable.

Yet despite the pressing requirement of trustworthiness, bugs in quantum chemistry software are still inevitable, due to its considerable large code size and complex code logics. Furthermore, in recent years, in response to the strong demand for computational performance, an increasing number of quantum chemistry software are being deployed on cloud platforms [28] or supercomputers [31]. In such environments, any bugs might lead to not only program crashes, but also tenant environment

corruptions. For example, a bug of wrong MPI number causes a crash of Siesta program during grid initialization [35]. As a result, securing the quantum chemistry software has become a pressing need.

Recognizing this need, a significant amount of studies have been conducted on scientific software, with diverse techniques employed. (*e.g.*, security testing [28] [30], fuzzing [49] [50] [7], probabilistic testing [47], and verification [48] [29]). Yet despite these research progresses, detecting bugs in quantum chemistry software remains challenging due to two unique characteristics of quantum chemistry software: **C1**): syntactic complexity; and **C2**): testing oracles. First, both input and output formats have complex syntactic structures, as they not only should represent interacted chemical parameters but also should obey certain spatial constraints (*e.g.*, symmetry). As a result, manually developing syntactically valid testing input is laborious and error-prone. Second, developing effective testing oracles for quantum chemistry software is challenging, as chemical domain knowledge is required in both input generation and result verification. Furthermore, developing effective testing oracles is difficult not only for software engineers but also even for chemistry scientists, due to the subtle properties of chemical molecules [44].

**Our work.** In this paper, to address the aforementioned first challenge C1, we utilized a syntactic mutation approach in which new syntactic correct input files are generated by syntactically mutating existing seeds. To address the challenge C2, we have leveraged transformer-based Large Language Models (LLMs) as our testing oracles, to exploit their potentials of overcoming domain knowledge problems. Our approach is end-to-end and automated: during testing, it requires neither manual interventions nor user prompts.

To realize the whole process, we have designed CHEMFUZZ, a software prototype to implement our approach. CHEMFUZZ consists of two core modules: mutation and analysis. The mutation module includes mutation operators, data padding, seed file validity checks, and collaborative efforts to create input file mutations. The analysis module encompasses code coverage collection, LLMs warning generation, a fitness function used for scoring, and a comprehensive assessment of the interest and effectiveness of seed files.

We have conducted extensive experiments to evaluate CHEM-

FUZZ in terms of effectiveness, performance of LLMs, and contributions of components. While our approach is general and can be built upon any quantum chemistry software and LLMs, we have conducted the evaluation on Siesta, a popular quantum chemistry software used by several thousand research all over the world [4], and used three most popular LLMs: GPT3.5 [24], Claude-2 [45], Bart [46] as the oracles. First, CHEMFUZZ is effective in detecting 40 unique bugs and 81 LLMs warnings, which have been classified and reported to developers. In the meanwhile, the maximum line coverage is 25,719, accounting for 17.43% of total lines, the function coverage is 939, accounting for 32.42% of all functions. Second, we compare and analyze the performance of three different LLMs: GPT3.5, Claude-2, and Bard, with respect to the validity of generated files. The results show Claude-2 exhibits an impressive efficiency rate of 87.21% in generating higher-quality seed files. Third, we evaluate mutation operators, LLMs prompts, and fitness functions by analyzing their impacts on the effectiveness of CHEMFUZZ. Compared to its competitors, the current component has a stronger vulnerability detection capability and higher code coverage.

**Contributions.** To the best of our knowledge, this is the *first* work on utilizing a LLMs-assisted approach to fuzz quantum chemistry software for vulnerability detection. The main contributions of this paper are:

- **Infrastructure design.** We present the first security testing infrastructure for quantum chemistry software with a LLMs-assisted fuzzing approach.
- **Prototype implementation of CHEMFUZZ.** We implemented a prototype CHEMFUZZ to validate our design.
- **Evaluation of CHEMFUZZ.** We conducted extensive experiments to evaluate CHEMFUZZ in terms of effectiveness and performance.

**Outline.** The rest of our paper is organized as follows. Section 2 introduces the necessary background knowledge for our work. Section 3 explains the motivation and challenges. Section 4 and 5 presents the design of CHEMFUZZ and the evaluation, respectively. Section 6 discusses limitations and direction for future work. Section 7 introduces related work, and Section 8 concludes.

## 2. BACKGROUND

To be self-contained, in this section, we first present the necessary background knowledge for this work: quantum chemistry software, fuzzing, and large language models.

### 2.1. Quantum Chemistry Software

**Wide adoptions.** Quantum chemistry software is widely adopted for research and industry purposes, such as Chemistry, Biology, Material Science, and Engineering. For example, Siesta, as one of the most important quantum chemistry software, currently owns several thousand users around the world and achieves more than thirteen thousand citations [22].

**Workflow.** Figure 1 presents a typical workflow for quantum chemistry software. First, input files (*e.g.*, FDF files, GJF files, etc.) are either selected from examples previously stored in

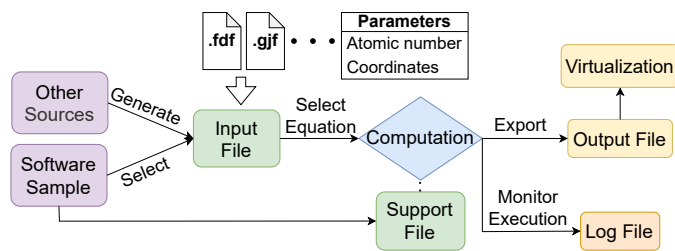


Figure 1: Typical workflow of quantum chemistry computation software.

software or generated from other sources. These input files contain parameters to pre-define and regulate computation. Next, if input parameters define an equation, the quantum chemistry software will follow it; otherwise, the software will automatically select suitable equations based on input parameters. In addition, related support files (*e.g.* Pseudopotentials) from software samples may be used in helping the calculation. Subsequently, computational results will be exported into an output file. At the same time, execution behaviors monitored by software will also be recorded within a log file.

### 2.2. Fuzzing

**Concept.** Fuzzing is a popular and important software testing technique to detect vulnerabilities and bugs [15] [17]. The basic idea of fuzzing is to generating a gigantic amount of normal and abnormal inputs and feeding them into the targeted program to trigger unexpected program behaviors. Fuzzers can monitor the execution states of programs to detect possible issues such as memory corruption, program crashes, extreme resource usage, and input validation errors.

**Feature.** In contrast to other techniques, fuzzing displayed outstanding portability and accuracy. A fuzzing framework can be easily deployed and of good extensibility and applicability [16]. It demands less knowledge of the target program and shows high accuracy since it is performed in real execution.

### 2.3. LLMs

**Introduction.** LLMs refer to a Transformer language containing a gigantic number of parameters and pre-trained by billions or trillions of available text data on the Internet [1] [23], (*e.g.*, Galactica[8], GPT-3.5 [24], and LLaMa [25]). Extending the scale of training allows LLMs to perform outstandingly in completing Natural Language Processing (NLP) tasks.

**Emergent abilities.** LLMs have shown diverse emergent abilities, which are not present in smaller models [23] [32]. Experiments on a Massive Multi-task Language Understanding (MMLU) benchmark demonstrate LLMs’ capability in multi-task Language understanding by accurately solving knowledge-based questions.

**Prompt engineering.** Prompts quality influences LLMs outputs [3] [14] [38]. By using well-defined prompts, LLMs can indicate better communication interaction and enhance the ability to follow instructions for executing specific tasks.

SystemName Magnesium Oxide		Basic Part
SystemLabel	MgO ①	
NumberOfAtoms	2	
NumberOfSpecies	2	
%block Chemical_Species_Label		
1	12 Mg	
2	8 O	
%endblock Chemical_Species_Label		
AtomicCoordinatesFormat ScaledCartesian		
%block AtomicCoordinatesAndAtomicSpecies		
	.000 .000 .000 1 ⑦	
	.500 .500 .500 2	
%endblock AtomicCoordinatesAndAtomicSpecies		
-----		
PAO.BasisSize SZ ②		Optional Part
PAO.NewSplitCode false ③		
PAO.EnergyShift 300 meV		
PAO.SplitNorm 0.15 ④		
MD.NumCGsteps 50 ⑤		
MD.MaxForceTol 0.04 eV/Ang ⑥		

Figure 2: A Siesta sample input file illustrating the unique syntactic and semantic challenges to test quantum chemistry software.

### 3. CHALLENGES AND MOTIVATION

In this section, we first present challenges caused by unique syntactic structures and semantics of quantum chemistry software (§ 3.1), then give our motivation (§ 3.2).

#### 3.1. Challenge

Testing quantum chemistry software is challenging, due to the intense complexity of syntactic structures and semantics. Both input and output files consist of sophisticated chemical parameters, such as molecular geometry structure, calculation method, and other defined parameters used in the computation. Here, we use a Siesta input file in Figure 2 as an example to illustrate its syntactic and semantic challenges.

**Syntactic challenges.** The basic unit of the input file can be called a section, including a label and its data portion. Sections can be divided into two related categories: 1) Required section and optional section; and 2) Associated section and non-associated section.

*Required and Optional.* Required sections, as shown in Figure 2 basic part, define the fundamental molecular information for simulation, which are necessary and not removable. Oppositely, optional sections in the optional part define additional information, such as atomic moving orbitals. These sections are used to designate simulation methods or assist simulation to achieve more ideal results.

*Associated and Non-associated.* Most sections in input files are non-associated sections that define parameters independently. However, relationships exist between associated sections, like highlighted AtomicCoordinatesAndAtomicSpecies and AtomicCoordinatesFormat in Figure 2. The former defines atomic coordinates, and the latter regulates its format. Modifying each

of them should affect correlated another section, which may eventually influence the entire simulation.

**Semantics challenges.** Semantics in data portions are also complicated due to various data types and the domain knowledge required in data selection.

*Data.* Data with a variety of types are defined in the data portion. Using Figure 2 as an example, data are represented within string, boolean, int, float, physical, and matrix formats. Although sections in dotted boxes 1 and 2 both use string data, the former data ‘MgO’ is a simple string describing the molecular system, and the enumerated data ‘SZ’ is selected from a pre-defined list of basis set orbital size including SZ, DP, SZP, and DZP. Boolean data in dotted box 3 is either a true or false statement, float and int data like dotted boxes 4 and 5, refer to real and integer numbers, and physical data in dotted box 5 is the combination of int or float data and physical units. Matrix data in dotted box 7 consist of multiple rows of different data, such as XYZ coordinates in float, atomic index in int, and element label in string.

*Domain knowledge.* Apart from the data type, domain knowledge is frequently required in selecting data values. These values should come from a reasonable domain, and rationality should be ensured. Otherwise, erroneous values will no longer hold physical significance, and the simulation cannot be implemented in the practical world—for example, the section PAO.EnergyShift in Figure 2, according to Siesta official documentation, is a standard for orbital-confining cutoff radius. It should be only selected with a small positive real number rather than a zero or negative value.

#### 3.2. Motivation

After recognizing previous challenges, selecting suitable testing techniques and oracles becomes more essential. First, We aim to use mutation-based fuzzing to resolve structure complexity because such mutations on selected files maintain the greatest possible structure validity, as well as increase the possibility of finding bugs. Second, we consider LLM as a qualified oracle for semantic complexity and reduce domain knowledge barriers. Its gigantic training data and the emergent ability to follow instructions allow it to provide trustful and accurate data.

### 4. DESIGN

In this section, we present the design of CHEMFUZZ, by describing its overall architecture (§ 4.1), then its two main modules: the mutation (§ 4.2) and the analysis (§ 4.3).

#### 4.1. Architecture

We present, in Figure 3, the overall framework of CHEMFUZZ. It employs a dynamic and mutation-based evolutionary fuzzing strategy and mutates the structure and data of seed files to generate effective new input files for quantum chemistry software. The main components of CHEMFUZZ consist of two modules: the mutation module and the analysis module. The mutation module includes mutation operators, data filling using Large Language Models, and seed file validity checks.

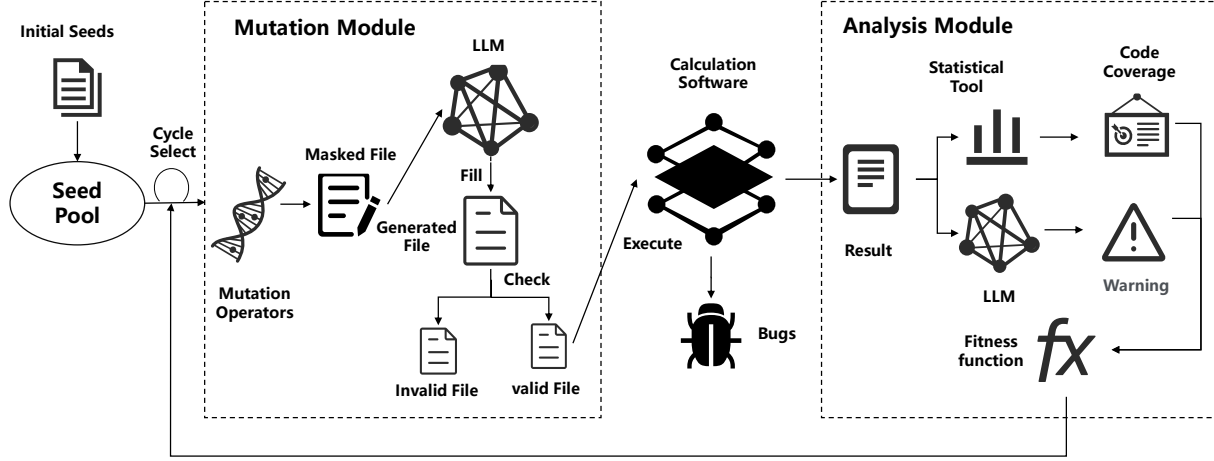


Figure 3: Overall framework of CHEMFUZZ.

---

**Algorithm 1:** Evolutionary fuzzing algorithm.

---

**Input:** Seed: the initial seed files,  
 Turn: the turn of evolution,  
 $f_1$ : the round factor,  
 $f_2$ : the retention factor.

**Output:** The generated files.

```

1 Function EvoFuzz (Seed, N,  $f_1$ ,  $f_2$ ):
2   SeedBank  $\leftarrow$  Seed;
3   InitializeFitness ();
4   InitializeConfigure ();
5   for n range in N do
6     k  $\leftarrow$  |SeedBank|;
7     for i range in  $f_1 * k$  do
8       CurrentSeed  $\leftarrow$  SelectByFitness(Seedbank);
9       Op  $\leftarrow$  RandomSelectMutationOp();
10      MaskedFile  $\leftarrow$  Mask(CurrentSeed, Op);
11      Cases  $\leftarrow$  LLM(MaskedFile);
12      ValidCases, InvalidCases  $\leftarrow$  Check(Cases);
13      Coverage, Result  $\leftarrow$  Execute(ValidCases);
14      WarningInfo  $\leftarrow$  LLM(Result);
15      FitnessScore  $\leftarrow$  FitnessFunction(ValidCases);
16      SeedBank  $\leftarrow$  SeedBank  $\cup$  ValidCases;
17     SelectTopFitness(SeedBank,  $f_2 * k$ );
18   return SeedBank;

```

---

On the other hand, the analysis module involves code coverage collection, warning information generation based on LLMs, and fitness function calculation.

The above calculation process can be summarized as Algorithm 1. In this algorithm, we have designed a seed file elimination mechanism. After generating a certain number of files, only the top-performing ones are retained, allowing for the rational elimination of seed files with low code coverage and insufficient detection capabilities. The elimination process

is controlled by two coefficients: the round factor  $f_1$  and the retention factor  $f_2$ . Both  $f_1$  and  $f_2$  are typically positive numbers greater than 1, with  $f_1$  being greater than  $f_2$ . By setting the round value  $N$  we can control the number of generated files, which will affect the depth of testing and the program’s runtime.

#### 4.2. Mutations

In the mutation module, our goal is to achieve diverse file mutations while satisfying the software input file constraints, to maximize the detection of program vulnerabilities. Therefore, our core approach is to apply mutation operators to the selected seed files to alter their file *structure* and leverage the large language model’s excellent ability in structural text filling to modify the file *values*.

**Mutation Operators.** We define three types of mutation operators: *Add*, *Delete*, and *Modify*, randomly selecting an operator for mutation. Add and Delete are used to alter the structure of a file, while Modify is used to keep the file structure unchanged and only specify the need to modify a particular field within the file. During the mutation process, newly added sections from Add and the data portion of the selected sections from Modify are replaced with the “<Type-tag>” token, indicating the type of section in the tags such as “<integer>”, “<float>” and “<matrix 3x2>”, and a particular enumeration type, choosing from several optional values.

*Add/Delete.* The Add mutation operator randomly selects a section label from all options and adds it to the original file. Specifically, if the chosen section belongs to an associated section, the associated sections are also added. The values of all newly added sections are replaced with the “<Type-tag>” token in masked files. Similarly, the Delete mutation operator randomly selects one from the existing optional sections of the current file and removes it and all associated sections.

*Modify.* We parse the syntax structure of the file and record all sections, their type and associated relationships in a data



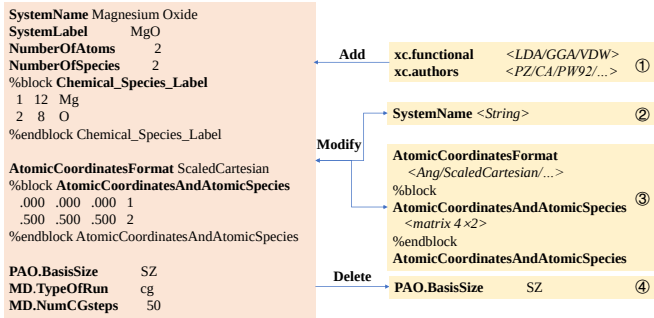


Figure 4: The effects of operators applied to sections.

TABLE I: Mutation operators applied in different sections.

	Add	Delete	Modify
Required & Associated	×	×	○
Optional & Associated	○	○	○
Required & No-associated	×	×	—
Optional & No-associated	—	—	—

structure. Then, we randomly select a section from the seed file, if the chosen section belongs to an associated section, the associated sections are also selected. The data portions of all the selected sections are marked as “<Type-tag>”.

Based on the way the mutation operators work, the system needs to know the list of labels for sections in the target software and their associate relationships. This task is not difficult, just a bit tedious, for the latest version, as there are approximately 200 sections in Siesta and approximately 250 in Quantum Espresso with different labels [4] [27]. We will collect and record this information from the software documentation into a configuration file.

**File Generation.** Table I presents the performance of different operators applied to various types of sections (discussion in 3.1) in quantum chemistry computation programs. “×” indicates that the operator cannot be applied to this section, “-” indicates that the operator will only be applied to this section, and “○” indicates that the operator will be applied to all fields associated with this section.

Figure 4 shows the effects of three operators applied to different types of sections. Four cases demonstrate the effects of “add” on an optional & associated section, “modify” on a required & associated section and a required & non-associated section, and “delete” on an optional & non-associated section, respectively. These cases confirm the operators’ ability to mutate the results of seed files effectively.

After randomly applying a mutation operator to produce masked files, we utilize Generative Pre-trained Language Models such as GPT and Claude-2 to generate new files by filling in the masked-out locations.

When using LLMs, a well-defined prompt can enhance and strengthen the capability of a large language model [40]. For a sophisticated text generation task, the prompt should

be correctly structured into the following components: the motivation of the task, output customization, specific rules and guidelines of the task, the key ideas, and example implementations [11] [12] [13]. Based on these theories, we have implemented our text generation prompt:

*At the end of the task, a masked input file for the quantum chemistry calculation software Siesta-4.1.5 is provided in the format of “text”, with masked parts labeled using the <> tags. Your task is to fill in the masked parts using text and generate an input file that can actually be run in the Siesta, and respond with the filled-in file.*

*Tag Rules:*

- 1. The first tag type indicates the data type, such as <boolean>, <float>, <matrix 3×2> and etc. You need to fill in specified parameters of the indicated data type.*
- 2. The second tag type provides several options, such as <Ang/ScaledCartesian>. You need to use one of the options provided.*

*Filling task requirements:*

- 1. The generated file should retain the unmasked parts of the input file and fill in all the masked parts.*
- 2. The generated file should comply with the requirements of version 4.1.5 of Siesta, as specified in the Siesta usage document.*
- 3. The generated file should adhere to the requirements of the quantum chemistry field in terms of calculation methods, physical and chemical properties, crystal structure, and other parameters.*
- 4. The generated file should be filled with the extreme and reasonable values you consider as soon as possible.*
- 5. The generated file should differ as much as possible from the files stored in your memory.*

*Response Format:*

- 1. Your response should only include the filled-in result of the input file provided below in the format of “text”.*
- 2. You must not include any additional text information.*

The prompt provides detailed descriptions of various aspects of input text generation for quantum chemistry software, enabling LLM to directly generate interesting and appropriate bi-directional completions for the intended task. Program errors often occur near extreme or boundary values [42] [43]. The prompt specifically emphasizes generating such values to identify potential software vulnerabilities effectively.

Finally, the criterion for judging the validity of input files is their ability to pass the validation based on the tested software file format and enter the numerical computation phase. Except only for specific seed files that exhibit inconsistent behavior between the file parser and the document description.

### 4.3. Analysis

In the analysis module, Our goal is to analyze the results of seed files in the calculation software and calculate their fitness scores. We firstly measure code coverage during the execution process, and secondly employ a large language model to detect any abnormal data in the execution results of the valid mutation file. Finally, we determined the fitness score of the mutation files based on these findings, providing a basis for selecting subsequent seed files.

**Coverage Collection.** For large-scale C++/Fortran programs like quantum chemistry calculation software, using the *gcov* [33] and *lcov* [34] component from the GCC compilation suite to perform code coverage analysis is often a good practice. We insert the “-fprofile-arcs -ftest-coverage” flags during the compilation of the software source code using g++/gfortran. After the software execution, gcda files showing code coverage are generated, and the lcov tool is used to produce the final required code coverage report. In this context, we use line coverage, denoted as  $C$ , as a measure of the software’s ability to detect program vulnerabilities, which will be utilized in the subsequent fitness function.

**LLM Warnings.** Similar to section 4.2, we leverage large language models to assist in identifying potential abnormal data in the calculation results of the chemical software. Building upon the principles of prompt engineering and tailored for Siesta version 4.1.5, we have developed a specific prompt to check for warnings in the calculation results:

*The core part of the output of the quantum chemical computing software Siesta-4.1.5 is provided at the end of the task in the format of “text”. Your task is to determine whether the software output has obvious anomalies.*

*Task requirements:*

- 1: Only very obvious exceptions need to be flagged in the task.
- 2: Incorporating knowledge from the field of quantum chemistry, considering the physical significance of each value, numeric parameters in the result need to be examined to determine if it falls outside the normal data range.

*Result requirements: The task results meet the requirements:*

- 1: The result value is Abnormal/Normal.
- 2: When the answer is Normal, your answer does not need to include any additional information.
- 3: When the answer is Abnormal, your answer gives the warning text information and its reason in WarningLocation.

*The output format requirements are as follow json file:*

```
{“Result”: “Abnormal”,
“WarningLocation”: [{
“located”: “text1”,
“reason”: “Why are to lead the warning.”
}]}
```

The prompt defines the task objectives and response format, allowing us to easily analyze the warnings’ positions and their underlying causes from the returned JSON object. Similarly, We define the number of warnings as  $W$ , which indicates the likelihood of the program generating abnormal results, serving as a crucial basis for the fitness of seed files.

**Fitness functions.** In the evolutionary algorithm of CHEMFUZZ, selecting unique and interesting files for mutation through the fitness function is crucial [17] [18]. Considering both the program’s execution status ( $S$ ), line coverage ( $C$ ), and the warning size ( $W$ ), we define a fitness function ( $F$ ) to evaluate the effectiveness of seed files ( $f$ ):

$$F(f) = \begin{cases} C + W/10, & \text{if } S = 0 \\ 1, & \text{if } S = 1 \end{cases}$$

One small detail pertains to calculating the execution status ( $S$ ), which is straightforward and trivial. It involves merely recording whether the file triggers a crash or enters an infinite loop during program execution. In such cases,  $S = 1$ ; otherwise,  $S = 0$ .

The rationale for using it as the fitness function mainly includes the following three points: 1) Evaluating the execution status ( $S$ ) based on whether errors occur provides the most direct evidence for determining the effectiveness of legitimate seed files in discovering program vulnerabilities. 2) When the calculation software produces normal output results, line coverage ( $C$ ) indicates the number of covered code lines by seed files, while the error warnings ( $W$ ) provided by the large language model suggest the possibility of anomalies or errors in the results. 3) The quantity of LLM warnings ( $W$ ) typically falls within the range of 1 to 10, while line coverage ( $C$ ) always ranges between 0 and 1. Assigning a weight of 0.1 to  $W$  effectively balances the impact of the two factors.

Based on this, we assign fitness scores to the seed files, which serve as weights for the selection of file mutation.

## 5. EVALUATION

In this section, we present experiments to evaluate CHEMFUZZ. First, we present the research questions (§ 5.1) guiding the experiments, and experimental setup (§ 5.2). Second, we show the experimental results in terms of the effectiveness (§ 5.3), LLMs performance (§ 5.4), and components contributions (§ 5.5). Finally, we discuss a case study of real-world vulnerabilities CHEMFUZZ detected (§ 5.6).

### 5.1. Research Questions

By presenting the experimental results, we mainly investigate the following research questions:

**RQ1: Effectiveness.** Is CHEMFUZZ effective in bug detection in real-world programs?

**RQ2: LLMs performance.** In the task of generating input files for chemical calculation software, what are the performance of different large language models?

**RQ3: Components contributions.** What are the contributions of the key components of CHEMFUZZ to its effectiveness?

### 5.2. Experimental Setup

**Test program.** We employed Siesta 4.1.5 [4] [21] [22] as the software under test, a stable version widely and extensively used in diverse fields for a considerable period. In addition, We collected 35 input files from the examples and tests of the Siesta software to serve as initial seeds for our study.

**Configure.** In the configuration file, we collected data on 76 labels corresponding to the sections in Siesta, including their data types, optionality, and associate relationships, as documented in the official documentation. The round factor  $f_1$  and the retention factor  $f_2$  are set to 2.0 and 1.5.

**Environment.** All the experiments and measurements are performed on a server with one 48 physical Intel E7-4830 core CPU and 64 GB of RAM running Ubuntu 18.04. We used the built-in gcov [33] and lcov [34] tools from gfortran 7.5.0 to collect code coverage information.

TABLE II: Summary of detected bugs.

Running state	Bug categories	Size
Crash	File parse error	5
	Stack Overflow	3
	Integer Overflow	9
	Others	2
	Endless	Loop condition error
Endless	Parameter validation error	14
	Others	3
Total	6	40

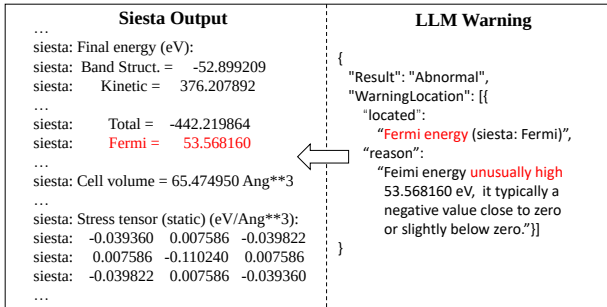


Figure 5: A warning to generated by GPT-3.5.

### 5.3. RQ1: Effectiveness

We conducted multiple rounds of experiments with CHEMFUZZ using different large language models, including GPT-3.5 [24], Claude-2 [45], and Bart [46].

CHEMFUZZ, in version 4.1.5 of the Siesta software package, has detected 40 bugs that were previously unknown to Siesta developers. Table II presents the statistics of these bugs detected by CHEMFUZZ. Various bugs were identified, including 5 file parse error, 3 stack overflow, 9 integer overflow, 4 endless Loop condition errors, 14 parameter validation errors, and 3 bugs that could not be attributed to a specific cause. We have reported all the bugs to the developers of the Siesta software.

Table III show different LLMs provided 81 warnings regarding the software’s output results. From the table, two conclusions can be drawn: 1) The majority of LLMs’ warnings are related to excessively large or small values of physical parameters in the computed results, exceeding reasonable ranges. This deserves attention. 2) There is a significant difference in the number of warnings provided by different LLMs, which may be attributed to varying sensitivity to data anomalies.

Figure 5 presents a warning from GPT3.5, and this case originates from Siesta simulation calculations using the DFT method for MgO at an environmental temperature of 500000K. As a truth, in such high-energy environments, conventional DFT methods often lead to distortions [40]. The Language Model-based Learning system is capable of detecting when the output results contain an excessively large Fermi level, reaching 53.568160eV, a value that is typically

TABLE III: LLMs warning of software calculation result.

LLMs Warning	GPT-3.5	Claude-2	Bard	Total
Non-negative	4	0	1	5
Special value	6	2	2	10
Matrix asymmetric	2	1	0	3
Abnormally large	12	2	13	27
Abnormally small	15	3	10	28
Others	2	1	3	6
Total	41	10	30	81

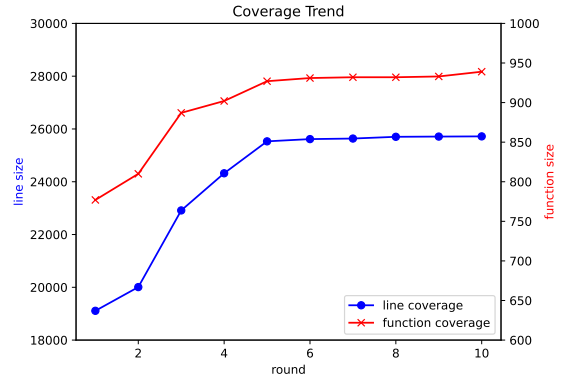


Figure 6: Coverage trend in the multi-round testing.

either a complex number or a very small positive number [41]. As a consequence, an alert is issued regarding the output results. This case serves as evidence for the effectiveness of LLM warnings.

These warnings merit the attention of both users and developers of quantum chemistry software. It is recommended to consult material databases or employ alternative software for comparative analysis to confirm the accuracy of the output results when necessary.

CHEMFUZZ achieved a maximum line coverage of 25,719, accounting for 17.43% of the total lines, and the highest function coverage of 939, accounting for 32.42%. Figure 6 illustrates the growth trend of code coverage during 10 retention rounds of mutation testing. Considering the extensive code of more than 140,000 lines and diverse computational methods of Siesta, our testing framework is effective.

**Summary:** From Siesta 4.1.5, CHEMFUZZ detected 40 bugs and 81 LLMs warnings, the maximum line coverage is 25,719, accounting for 17.43% of total lines; the function coverage is 939, accounting for 32.42% all functions, demonstrates its effectiveness in detecting real-world code vulnerabilities.

### 5.4. RQ2: LLMs Performance

To answer RQ2, a series of multi-round experiments were conducted with CHEMFUZZ to demonstrate the performance of different LLMs, including GPT-3.5, Claude-2, and bard, in

TABLE IV: The ability of LLMs under different prompts.

Model	GPT3.5				Claude-2				Bard			
	Valid	All	Rate	Coverage	Valid	All	Rate	Coverage	Valid	All	Rate	Coverage
CHEMFUZZ	176	213	82.63%	23176	155	178	87.08%	22741	98	127	77.17%	22512
No type tag	149	222	67.12%	22792	147	190	77.37%	22714	83	109	76.14%	22410
No boundary note	180	216	83.33%	22841	151	188	85.63%	22140	102	122	83.61%	21090

TABLE V: The performance of LLMs to generate quantum chemistry software input files.

LLM	Generated Files per API			Time per File (s)	
	Valid	All	Percentage	Valid	All
	GPT3.5	132	161	81.99%	9.14
Claude-2	150	172	87.21%	10.05	9.72
Bard	122	148	82.43%	8.94	8.04

generating Siesta input files. These experiments took a total of 14 hours.

Table V presents the size of generated files and average generation times of all LLMs. In our research, we have observed that Claude-2 is capable of generating quantum chemistry software seed files with higher quality, achieving an impressive efficiency rate of 87.21%. The response times of the three models are quite similar, all around 10 seconds. Therefore, the Claude-2 model is the most suitable and effective large model for our complex text generation task among the models evaluated.

In the text generation task, we found that the main reason for the invalidity of seed files lies in the inability to understand the task requirements and objectives sometimes, resulting in issues such as missing output files, incorrect formatting, and reproducing the input as output. These abnormal conditions warrant further optimization and improvement by LLM developers.

**Summary:** Claude-2 achieved the highest quality and effectiveness of 87.21%. The response times of three models are similar, all around 10 seconds.

### 5.5. RQ3: Evaluation of Key Components

**Mutation Operators.** In this study, we conducted an experiment of 33 hours to evaluate the impact of each mutation operator. The outcomes of generating approximately two hundred files are presented in Table VI, where we display the results after systematically removing each mutation operator from CHEMFUZZ. Notably, the employment of the complete set of mutation operators yielded a line coverage of 17.4%, which is significantly higher than the versions with Modify and Add removed. Although removing Delete had minimal impact on coverage, the runtime increased threefold in the complete version. Currently, the mutation operators demonstrate the best performance.

TABLE VI: Evaluation of mutation operators.

OP	Runtime duration(h)	Line Coverage	
		Lines	Percentage
CHEMFUZZ	6.2	25678	17.4%
-Modify	4.6	23465	15.9%
-Add	4.2	19233	13.0%
-Delete	18.3	25788	17.4%

**Fitness function.** We compared our default fitness function against alternative functions, including random selection (*Random*), code coverage-based weighting (*C*), and a simple function (*C+W*) in experiments. Except for random functions, the fitness score remains unchanged at 1 for seed files with abnormal execution status(*S*).

Table VII summarizes the results. Our fitness function has achieved a balance between the effects of *C* and *W*, enabling all interesting files to attain higher fitness scores. As a result, in the statistics of experiments, the current function (*C+W/10*) demonstrated the highest code line coverage of 17.2% and the highest function coverage of 32.1%.

**LLMs Prompt.** To investigate the impact of large language model (LLMs) prompts on the framework’s vulnerability detection capability, We conducted experiments using different LLMs on the same initial seed for an equal duration of time, while comparing the current prompt used in CHEMFUZZ with prompts that do not include type tags and or boundary indications. The experiments took a total of 19 hours.

Table IV presents the statistical data of the experimental results, the coverage in the data is line coverage. Compared to the modified prompts, the current prompt text generation capability demonstrates a balance between code coverage and generation effectiveness across all three large models. Furthermore, the removal of type annotations significantly decreases effectiveness in GPT3.5 and Claude2, while eliminating reminders for large models to fill in boundary values results in a decline in code coverage for all large models.

**Summary:** Compared to other counterparts, the current key components exhibit higher effectiveness, demonstrating superior code coverage in testing.

### 5.6. Case Study

Figure 7 shows a bug in Siesta 4.1.5 where there is a lack of validation for a specific value in the MaxSCFIterations section. The calculation program reads the MaxSCFIterations value



TABLE VII: Comparison of fitness function effectiveness.

Function	Line Coverage		Function Coverage	
	Lines	Percentage	Functions	Percentage
CHEMFUZZ	25376	17.2%	930	32.1%
<i>C</i>	24121	16.3%	892	30.8%
<i>C + W</i>	23725	16.1%	888	30.7%
<i>Random</i>	23109	15.7%	861	29.7%

from the `fdf` file, which represents the number of iterations for the self-consistent field (SCF) calculation. The value is initially read into the `fdf_integer` variable and then used as the number of iterations `nscf`, without any validation.

Because the SCF calculation is a crucial step in quantum chemistry calculations [51], this bug causes the entire calculation process to terminate without producing any results if `MaxSCFIterations` is set to 0 in the input file. This behavior prevents the proper calculation of physical parameters and methods, rendering the calculation process non-functional.

This is clearly not in line with the expected behavior of the program. CHEMFUZZ detected this exceptional situation and reminds quantum chemistry software developers to pay more attention to input parameter validation and enhance program robustness.

## 6. DISCUSSION

In this section, we discuss the limitations of this work, as well as directions for future work.

**Quantum Chemistry Software.** Though experimental results indicated that our framework in Siesta is effective, there exists several other influential software (e.g., Gaussian [20], Quantum Espresso [27], Abacus [36]). Despite software differences, we speculate that software selection only affects performance since they share similar workflows. We will continuously apply CHEMFUZZ on different quantum chemistry software in the future.

**Oracle.** CHEMFUZZ leverages GPT-3.5, Claude-2, and Bard as oracles in input generation and output verification. Additionally, Lin et al. [39] research shows that LLMs occasionally follow popular misconceptions and provide fraudulent answers. Hence, we assume utilizing a LLM specially trained for science fields and scientific knowledge bases is essential and will help improve the effectiveness of CHEMFUZZ.

**Other vulnerabilities.** CHEMFUZZ has detected certain vulnerabilities. Nevertheless, there are other types of vulnerabilities. Particularly, it is vital to inspect concurrency bugs in parallel computing of quantum chemical computation[37]. Besides, though CHEMFUZZ solve the complexity of input file structure, it is currently incapable of analyzing intricate output file structure. We leave these two as crucial future work.

## 7. RELATED WORK

We classified related work into three categories: fuzzing, Large Language Model, and security testing.

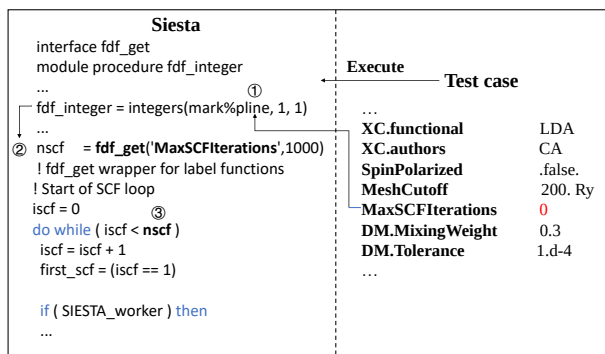


Figure 7: A real world bug of loop parameter validation detected by CHEMFUZZ, in Siesta 4.1.5.

**Fuzzing.** Fuzzing test techniques have been improved and are heavily used in different realms and applications. Guo et al. [5] proposed DLFuzz as a special fuzzing framework for leading deep learning systems to expose malicious behaviors. Deng et al. [6] built TitanFuzz to first use large language models to generate input programs that satisfy grammar syntax and constraints, which is used for fuzzing deep learning libraries. Wang et al. [7] present QuanFuzz by generating search-based test cases as input to perform grey-box fuzz testing for quantum programs.

Our work, which differs from recent research, is an evolutionary fuzzing model specialized for quantum chemistry computational software.

**LLM.** Large Language Model is an emerging language model pre-trained by billions of text data on the internet and extensively used in research fields. Singh et al. [9] leverage LLM to rate robots' potential actions and even generate robot commands with no external instructions. They propose a programmatic LLM prompt structure to acquire that. Beltagy et al. [10] release a pre-trained BERT [12] language model called SciBERT, which brings up a better performance on downstream scientific Natural Language Processing tasks.

However, a key difference with our work is that we leverage pre-trained large language models as an oracle both in generating input parameters within an acceptable domain of quantum chemistry and in verifying the accuracy and rationality of computation results.

**Security testing.** Recently, security testing for software has become favored. Schieferdecker et al. [11] introduce MBST, which automatically generates systematic documentation of security test objectives, cases, and suites. Bozic et al. [19] provide an approach with experimental results to automatically test chatbots on web applications.

In contrast to previous research, in our security test framework, we combine fuzzing techniques and LLM, two popular techniques given a lot of attention.

## 8. CONCLUSION

In this paper, we present CHEMFUZZ, the first approach to fuzzing quantum chemistry software via large language mod-

els. The work represents a crucial first step towards vulnerability detection in quantum chemistry programs, benefiting the overall community of software testing and quantum chemistry software development.

#### ACKNOWLEDGMENTS

We thank the anonymous reviewers for their valuable feedback. This work is partially supported by the National Natural Science Foundation of China (No.62072427, No.12227901), the Project of Stable Support for Youth Team in Basic Research Field, CAS (No.YSBR-005), Academic Leaders Cultivation Program, USTC.

#### REFERENCES

- [1] Deng Y, Xia C S, Peng H, et al. Fuzzing deep-learning libraries via large language models[J]. arXiv preprint arXiv:2212.14834, 2022.
- [2] Obot I B, Macdonald D D, Gasem Z M. Density functional theory (DFT) as a powerful tool for designing new organic corrosion inhibitors. Part 1: an overview[J]. Corrosion Science, 2015, 99: 1-30.
- [3] Reynolds L, McDonell K. Prompt programming for large language models: Beyond the few-shot paradigm[C]//Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems. 2021: 1-7.
- [4] "Siesta,": <https://departments.icmab.es/leem/siesta/>
- [5] Guo Y, Jiang Y, Zhao Y, et al. Difuzz: Differential fuzzing testing of deep learning systems[C]//Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 2018: 739-743.
- [6] Deng Y, Xia C S, Peng H, et al. Fuzzing deep-learning libraries via large language models[J]. arXiv preprint arXiv:2212.14834, 2022.
- [7] Wang J, Gao M, Jiang Y, et al. QuanFuzz: Fuzz testing of quantum program[J]. arXiv preprint arXiv:1810.10310, 2018.
- [8] Tan Y X, Zhang F, Xie P P, et al. Rhodium (III)-catalyzed asymmetric borylative cyclization of cyclohexadienone-containing 1, 6-dienes: an experimental and DFT study[J]. Journal of the American Chemical Society, 2019, 141(32): 12770-12779.
- [9] Singh I, Blukis V, Mousavian A, et al. Progprompt: Generating situated robot task plans using large language models[C]//2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023: 11523-11530.
- [10] Beltagy I, Lo K, Cohan A. SciBERT: A pretrained language model for scientific text[J]. arXiv preprint arXiv:1903.10676, 2019.
- [11] Schieferdecker I, Grossmann J, Schneider M. Model-based security testing[J]. arXiv preprint arXiv:1202.6118, 2012.
- [12] Devlin J, Chang M W, Lee K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint arXiv:1810.04805, 2018.
- [13] Santu S K K, Feng D. TELeR: A General Taxonomy of LLM Prompts for Benchmarking Complex Tasks[J]. arXiv preprint arXiv:2305.11430, 2023.
- [14] White J, Fu Q, Hays S, et al. A prompt pattern catalog to enhance prompt engineering with chatgpt[J]. arXiv preprint arXiv:2302.11382, 2023.
- [15] Manès V J M, Han H S, Han C, et al. The art, science, and engineering of fuzzing: A survey[J]. IEEE Transactions on Software Engineering, 2019, 47(11): 2312-2331.
- [16] Li J, Zhao B, Zhang C. Fuzzing: a survey[J]. Cybersecurity, 2018, 1(1): 1-13.
- [17] Liang H, Pei X, Jia X, et al. Fuzzing: State of the art[J]. IEEE Transactions on Reliability, 2018, 67(3): 1199-1218.
- [18] Rawat S, Jain V, Kumar A, et al. VUzzer: Application-aware Evolutionary Fuzzing[C]//NDSS. 2017, 17: 1-14.
- [19] Bozic J, Wotawa F. Security testing for chatbots[C]//IFIP International Conference on Testing Software and Systems. Cham: Springer International Publishing, 2018: 33-38.
- [20] "Gaussian,": <https://gaussian.com/>
- [21] Soler J M, Artacho E, Gale J D, et al. The Siesta method for ab initio order-N materials simulation[J]. Journal of Physics: Condensed Matter, 2002, 14(11): 2745.
- [22] García A, Papior N, Akhtar A, et al. Siesta: Recent developments and applications[J]. The Journal of chemical physics, 2020, 152(20).
- [23] Zhao W X, Zhou K, Li J, et al. A survey of large language models[J]. arXiv preprint arXiv:2303.18223, 2023.
- [24] Brown T, Mann B, Ryder N, et al. Language models are few-shot learners[J]. Advances in neural information processing systems, 2020, 33: 1877-1901.
- [25] Touvron H, Lavril T, Izacard G, et al. Llama: Open and efficient foundation language models[J]. arXiv preprint arXiv:2302.13971, 2023.
- [26] He Q, Yu B, Li Z, et al. Density functional theory for battery materials[J]. Energy & Environmental Materials, 2019, 2(4): 264-279.
- [27] "quantum-espresso,": <https://www.quantum-espresso.org/>
- [28] Francis R R. Reliability of cloud computing in Quantum Chemistry calculations[C]//2012 International Conference on Cloud Computing Technologies, Applications and Management (ICCTAM). IEEE, 2012: 119-120.
- [29] Honarvar S, Mousavi M R, Nagarajan R. Property-based testing of quantum programs in Q# [C]//Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering Workshops. 2020: 430-435.
- [30] Mendiluze E, Ali S, Arcaini P, et al. Muskit: A mutation analysis tool for quantum software testing[C]//2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 2021: 1266-1270.
- [31] Hu W, Qin X, Jiang Q, et al. High performance computing of DGFT for tens of thousands of atoms using millions of cores on Sunway TaihuLight[J]. Science Bulletin, 2021, 66(2): 111-119.
- [32] Hendrycks D, Burns C, Basart S, et al. Measuring massive multitask language understanding[J]. arXiv preprint arXiv:2009.03300, 2020.
- [33] "gcov-tool—an Offline Gcda Profile Processing Tool," <https://gcc.gnu.org/onlinedocs/gcc/Gcov.html>.
- [34] "LCOV - the LTP GCOV extension," <https://ltp.sourceforge.net/coverage/lcov.php>.
- [35] "Siesta issue on gitlab, grid initialization crashes with floating point exceptions for certain numbers of MPI tasks," <https://gitlab.com/siesta-project/siesta/-/issues/170>.
- [36] "Abacus," <http://abacus.ustc.edu.cn/main.htm>
- [37] Alkan M, Pham B Q, Hammond J R, et al. Enabling Fortran standard parallelism in GAMESS for accelerated quantum chemistry calculations[J]. Journal of Chemical Theory and Computation, 2023.
- [38] Zhou Y, Muresanu A I, Han Z, et al. Large language models are human-level prompt engineers[J]. arXiv preprint arXiv:2211.01910, 2022.
- [39] Lin S, Hilton J, Evans O. Truthfulqa: Measuring how models mimic human falsehoods[J]. arXiv preprint arXiv:2109.07958, 2021.
- [40] Foulkes W M C, Mitas L, Needs R J, et al. Quantum Monte Carlo simulations of solids[J]. Reviews of Modern Physics, 2001, 73(1): 33.
- [41] "The Materials Project,": <https://next-gen.materialsproject.org/>
- [42] Myers G J, Badgett T, Thomas T M, et al. The art of software testing[M]. Chichester: John Wiley & Sons, 2004.
- [43] Ramachandran M. Testing software components using boundary value analysis[C]//2003 Proceedings 29th Euromicro Conference. IEEE, 2003: 94-98.
- [44] Sims C. CIF2WAN: A Tool to Generate Input Files for Electronic Structure Calculations with Wannier90[J]. arXiv preprint arXiv:2006.12647, 2020.
- [45] "Claude-2," <https://www.anthropic.com/index/claude-2>
- [46] "Bard," <https://bard.google.com>
- [47] Krishnaswamy, Smita, Igor L. Markov, and John P. Hayes. "Tracking uncertainty with probabilistic logic circuit testing." IEEE Design & Test of Computers 24.4 (2007): 312-321.
- [48] Liu, J., Zhan, B., Wang, S., Ying, S., Liu, T., Li, Y., ... & Zhan, N. (2019). Formal verification of quantum algorithms using quantum Hoare logic. In Computer Aided Verification: 31st International Conference, CAV 2019, New York City, NY, USA, July 15-18, 2019, Proceedings, Part II 31 (pp. 187-207). Springer International Publishing.
- [49] Wang, J., Gao, M., Jiang, Y., Lou, J., Gao, Y., Zhang, D., & Sun, J. (2018). QuanFuzz: Fuzz testing of quantum program. arXiv preprint arXiv:1810.10310.
- [50] Xie, X., Ma, L., Juefei-Xu, F., Xue, M., Chen, H., Liu, Y., ... & See, S. (2019, July). Deephunter: a coverage-guided fuzz testing framework for deep neural networks. In Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis (pp. 146-157).
- [51] Verma P, Truhlar D G. Status and challenges of density functional theory[J]. Trends in Chemistry, 2020, 2(4): 302-318.