

Counterexample Guided Program Repair Using Large Language Models and MaxSAT-based Fault Localization

Pedro Orvalho¹, Mikoláš Janota² and Vasco Manquinho³

¹Department of Computer Science, University of Oxford, Oxford, UK

²CIIRC, Czech Technical University in Prague, Czechia

³INESC-ID, Instituto Superior Técnico, Universidade de Lisboa, Portugal

IIIA Seminars, IIIA, CSIC

23 June 2025



Automated Program Repair (APR)

Given a buggy program P_o and a set of input-output examples T (test suite).

Automated Program Repair (APR)

Given a buggy program P_o and a set of input-output examples T (test suite).

The goal of *Automated Program Repair* is to find a program P_f

Automated Program Repair (APR)

Given a buggy program P_o and a set of input-output examples T (test suite).

The goal of *Automated Program Repair* is to find a program P_f by **semantically change a subset S_1 of P_o 's statements** ($S_1 \subseteq P_o$)

Automated Program Repair (APR)

Given a buggy program P_o and a set of input-output examples T (test suite).

The goal of *Automated Program Repair* is to find a program P_f by **semantically change a subset S_1 of P_o 's statements** ($S_1 \subseteq P_o$) for another set of statements S_2 , s.t.,

$$P_f = ((P_o \setminus S_1) \cup S_2)$$

Automated Program Repair (APR)

Given a buggy program P_o and a set of input-output examples T (test suite).

The goal of *Automated Program Repair* is to find a program P_f by **semantically change a subset S_1 of P_o 's statements** ($S_1 \subseteq P_o$) for another set of statements S_2 , s.t.,

$$P_f = ((P_o \setminus S_1) \cup S_2)$$

and

$$\forall (t_{in}^i, t_{out}^i) \in T : P_f(t_{in}^i) = t_{out}^i$$

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

Symbolic-based APR Tools:

- **Automated Reasoning-based tools cannot fix this program within 90s;**

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

Symbolic-based APR Tools:

- **Automated Reasoning**-based tools cannot fix this program within 90s;
- CLARA takes too long to compute a 'minimal' repair;

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

Symbolic-based APR Tools:

- **Automated Reasoning**-based tools **cannot fix this program within 90s**;
- CLARA **takes too long** to compute a 'minimal' repair;
- VERIFIX returns a **compilation error**.

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

LLMs for code (LLMCs)

- GRANITE and CODEGEMMA **cannot** fix the buggy program within 90 secs;

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

LLMs for code (LLMCs)

- GRANITE and CODEGEMMA **cannot fix** the buggy program within 90 secs;
- Even if we provide this assignment's **description and IO tests**.

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

LLMs for code (LLMCs)

- GRANITE and CODEGEMMA **cannot fix** the buggy program within 90 secs;
- Even if we provide this assignment's **description and IO tests**.
- Suggesting a correct implementation, **both LLMs copy the correct code**, ignoring instructions not to do so.

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

- Symbolic approaches demand an **excessive amount of time to produce an answer**;

Motivation

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

- Symbolic approaches demand an **excessive amount of time to produce an answer**;
- LLMs, while fast, **often produce incorrect fixes**.

Program Sketches

```
1  int main(){ //finds max of 3 nums
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      if (f < s && f >= t)
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      else if (t < f && t < s)
9          printf("%d",t);
10
11     return 0;
12 }
```

1: Program sketch with holes.

```
1  int main(){
2      int f,s,t;
3      scanf("%d%d%d",&f,&s,&t);
4      @ HOLE 1 @
5          printf("%d",f);
6      else if (s > f && s >= t)
7          printf("%d",s);
8      @ HOLE 2 @
9          printf("%d",t);
10
11     return 0;
12 }
```


Our Work

- Combines the strengths of **Formal Methods (FM)** and **LLM-based approaches**;

Our Work

- Combines the strengths of **Formal Methods (FM)** and **LLM-based approaches**;
- Uses **MaxSAT-based fault localization** to **rigorously identify buggy lines**, which can then be highlighted in the LLM prompt;

Our Work

- Combines the strengths of **Formal Methods (FM)** and **LLM-based approaches**;
- Uses **MaxSAT-based fault localization** to **rigorously identify buggy lines**, which can then be highlighted in the LLM prompt;
- For instance, **instructing both LLMs to complete the previous sketch allows them to fix the buggy program** in a single interaction;

MaxSAT-Based Fault Localization

- **FM24** - CFAULTS: Model-Based Diagnosis for Fault Localization in C with Multiple Test Cases.

Fault Localization - Motivation

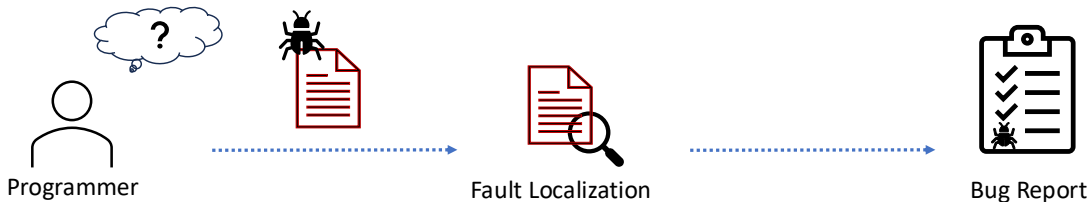
- Debugging is one of the most time-consuming and expensive tasks in software development.

Fault Localization - Motivation

- Debugging is one of the most time-consuming and expensive tasks in software development.
- In 2024, the estimated global cost of Crowdstrike's error that hit Microsoft systems, is 5.4 Billion US\$ [[The Guardian UK, 2024](#)].

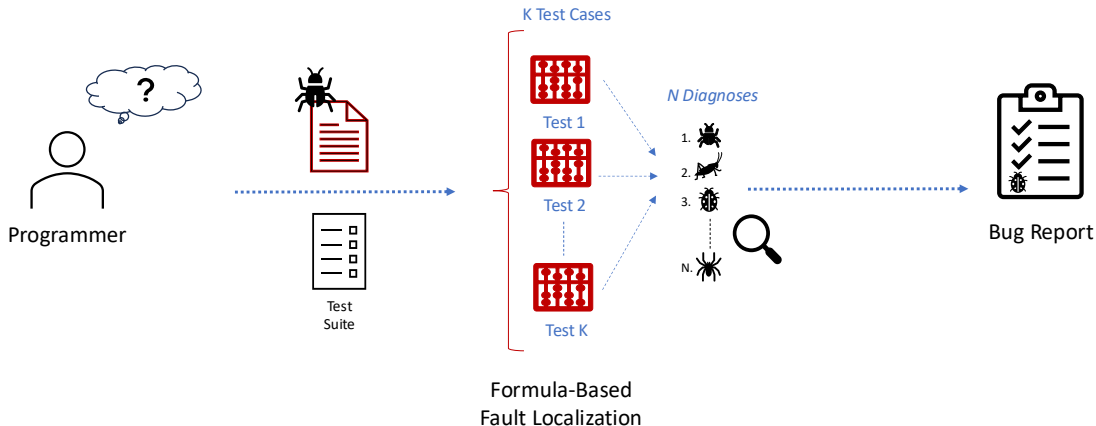
Fault Localization

- Given a buggy program, *fault localization (FL)* involves identifying locations in the program that could cause a faulty behaviour (bug).



Formula-Based Fault Localization (FBFL)

- FBFL methods encode the localization problem into **several optimization problems** to identify a minimal set of bugs (diagnoses).



Formula-Based Fault Localization (FBFL)

Limitations of Past Approaches

FBFL tools especially for programs with multiple faults:

Formula-Based Fault Localization (FBFL)

Limitations of Past Approaches

FBFL tools especially for programs with multiple faults:

- **do not ensure a minimal diagnosis** across all failing tests (e.g., BUGASSIST);

Formula-Based Fault Localization (FBFL)

Limitations of Past Approaches

FBFL tools especially for programs with multiple faults:

- **do not ensure a minimal diagnosis** across all failing tests (e.g., BUGASSIST);
- may produce an overwhelming number of **redundant diagnoses** (e.g., SNIPER).

Contribution

- We formulate the FL problem as a **single optimization problem**;

Contribution

- We formulate the FL problem as a **single optimization problem**;
- We leverage MaxSAT and the theory of *Model-Based Diagnosis (MBD)* [Reiter et al., 1987], **integrating all failing test cases simultaneously**;

Contribution

- We formulate the FL problem as a **single optimization problem**;
- We leverage MaxSAT and the theory of *Model-Based Diagnosis (MBD)* [Reiter et al., 1987], **integrating all failing test cases simultaneously**;
- We implement this MBD approach in a publicly available tool called CFAULTS.

Partial Maximum Satisfiability (MaxSAT)

Hard: $h_1 : (v_1 \vee v_2)$ $h_2 : (\neg v_2 \vee v_3)$

Soft: $s_1 : (\neg v_1)$ $s_2 : (\neg v_3)$

Partial Maximum Satisfiability (MaxSAT)

Hard: $h_1 : (v_1 \vee v_2)$ $h_2 : (\neg v_2 \vee v_3)$

Soft: $s_1 : (\neg v_1)$ $s_2 : (\neg v_3)$

Possible Solution: $\neg v_1 \vee v_2 \vee v_3$ $\text{cost} = 1$

Figure 1: Example of a partial MaxSAT formula.

Model-Based Diagnosis

- A system description \mathcal{P} **is composed of a set of components** $\mathcal{C} = \{c_1, \dots, c_n\}$.

Model-Based Diagnosis

- A system description \mathcal{P} **is composed of a set of components** $\mathcal{C} = \{c_1, \dots, c_n\}$.
- Each component in \mathcal{C} can be declared **healthy** or **unhealthy**.

Model-Based Diagnosis

- A system description \mathcal{P} **is composed of a set of components** $\mathcal{C} = \{c_1, \dots, c_n\}$.
- Each component in \mathcal{C} can be declared **healthy** or **unhealthy**.
- For each component $c \in \mathcal{C}$, $h(c) = 0$ **if c is unhealthy, otherwise, $h(c) = 1$.**

Model-Based Diagnosis

- A system description \mathcal{P} **is composed of a set of components** $\mathcal{C} = \{c_1, \dots, c_n\}$.
- Each component in \mathcal{C} can be declared **healthy** or **unhealthy**.
- For each component $c \in \mathcal{C}$, $h(c) = 0$ **if c is unhealthy, otherwise, $h(c) = 1$** .
- \mathcal{P} is described by a CNF formula, where \mathcal{F}_c denotes the encoding of component c :

$$\mathcal{P} \triangleq \bigwedge_{c \in \mathcal{C}} (h(c) \implies \mathcal{F}_c)$$

Model-Based Diagnosis

- **Observations represent deviations** from the expected system behaviour.

Model-Based Diagnosis

- **Observations represent deviations** from the expected system behaviour.
- An observation, denoted as o , can be encoded in CNF as a set of unit clauses.

Model-Based Diagnosis

- **Observations represent deviations** from the expected system behaviour.
- An observation, denoted as o , can be encoded in CNF as a set of unit clauses.
- In our work, **the failing test cases represent the set of observations**.

Model-Based Diagnosis

- **Observations represent deviations** from the expected system behaviour.
- An observation, denoted as o , can be encoded in CNF as a set of unit clauses.
- In our work, **the failing test cases represent the set of observations**.
- A system \mathcal{P} is considered **faulty if there exists an inconsistency with a given observation o when all components are declared healthy**:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C}} h(c) \models \perp$$

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;
- For a given MBD problem $\langle \mathcal{P}, \mathcal{C}, o \rangle$, a set of system components $\Delta \subseteq \mathcal{C}$ is a diagnosis iff:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;
- For a given MBD problem $\langle \mathcal{P}, \mathcal{C}, o \rangle$, a set of system components $\Delta \subseteq \mathcal{C}$ is a diagnosis iff:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

- A diagnosis Δ is:

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;
- For a given MBD problem $\langle \mathcal{P}, \mathcal{C}, o \rangle$, a set of system components $\Delta \subseteq \mathcal{C}$ is a diagnosis iff:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

- A diagnosis Δ is:
 - **minimal** iff no subset of Δ , $\Delta' \subsetneq \Delta$, is a diagnosis;

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;
- For a given MBD problem $\langle \mathcal{P}, \mathcal{C}, o \rangle$, a set of system components $\Delta \subseteq \mathcal{C}$ is a diagnosis iff:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

- A diagnosis Δ is:
 - **minimal** iff no subset of Δ , $\Delta' \subsetneq \Delta$, is a diagnosis;
 - Δ is of **minimal cardinality** if there is no other diagnosis $\Delta'' \subseteq \mathcal{C}$ with $|\Delta''| < |\Delta|$;

Model-Based Diagnosis

- The problem of model-based diagnosis (MBD) aims to **identify a set of components which, if declared unhealthy, restore consistency**;
- For a given MBD problem $\langle \mathcal{P}, \mathcal{C}, o \rangle$, a set of system components $\Delta \subseteq \mathcal{C}$ is a diagnosis iff:

$$\mathcal{P} \wedge o \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

- A diagnosis Δ is:
 - **minimal** iff no subset of Δ , $\Delta' \subsetneq \Delta$, is a diagnosis;
 - Δ is of **minimal cardinality** if there is no other diagnosis $\Delta'' \subseteq \mathcal{C}$ with $|\Delta''| < |\Delta|$;
 - is **redundant** if it is not subset-minimal [Ignatiev et al., 2019].

Model-Based Diagnosis

To encode the MBD problem with one observation with partial MaxSAT:

- The set of **clauses that encode \mathcal{P}** represents the set of hard clauses;

Model-Based Diagnosis

To encode the MBD problem with one observation with partial MaxSAT:

- The set of **clauses that encode \mathcal{P} represents the set of hard clauses**;
- The soft clauses consists of unit clauses that aim to **maximize the set of healthy components**, i.e.,:

$$\bigwedge_{c \in \mathcal{C}} h(c);$$

Model-Based Diagnosis

To encode the MBD problem with one observation with partial MaxSAT:

- The set of **clauses that encode \mathcal{P} represents the set of hard clauses**;
- The soft clauses consists of unit clauses that aim to **maximize the set of healthy components**, i.e.,:

$$\bigwedge_{c \in \mathcal{C}} h(c);$$

- This encoding enables enumerating subset **minimal diagnoses, considering a single observation**;

Model-Based Diagnosis with Multiple Test Cases

We **integrate all failing test cases** in a single MaxSAT formula.

Model-Based Diagnosis with Multiple Test Cases

We **integrate all failing test cases** in a single MaxSAT formula.

- We **generate only minimal diagnoses** capable of identifying all faulty components within the system, in our case, a C program;

Model-Based Diagnosis with Multiple Test Cases

We **integrate all failing test cases** in a single MaxSAT formula.

- We **generate only minimal diagnoses** capable of identifying all faulty components within the system, in our case, a C program;
- Given m observations, $\mathcal{O} = \{o_1, \dots, o_m\}$, a distinct replica of the system, denoted as \mathcal{P}_i , is required for each observation o_i ;

Model-Based Diagnosis with Multiple Test Cases

We **integrate all failing test cases** in a single MaxSAT formula.

- We **generate only minimal diagnoses** capable of identifying all faulty components within the system, in our case, a C program;
- Given m observations, $\mathcal{O} = \{o_1, \dots, o_m\}$, a distinct replica of the system, denoted as \mathcal{P}_i , is required for each observation o_i ;
- The hard clauses, ϕ_h , in our MaxSAT formulation correspond to:

$$\phi_h = \bigwedge_{o_i \in \mathcal{O}} (\mathcal{P}_i \wedge o_i);$$

Model-Based Diagnosis with Multiple Test Cases

We **integrate all failing test cases** in a single MaxSAT formula.

- We **generate only minimal diagnoses** capable of identifying all faulty components within the system, in our case, a C program;
- Given m observations, $\mathcal{O} = \{o_1, \dots, o_m\}$, a distinct replica of the system, denoted as \mathcal{P}_i , is required for each observation o_i ;
- The hard clauses, ϕ_h , in our MaxSAT formulation correspond to:

$$\phi_h = \bigwedge_{o_i \in \mathcal{O}} (\mathcal{P}_i \wedge o_i);$$

- The soft clauses are formulated as:

$$\phi_s = \bigwedge_{c \in \mathcal{C}} h(c).$$

Model-Based Diagnosis with Multiple Test Cases

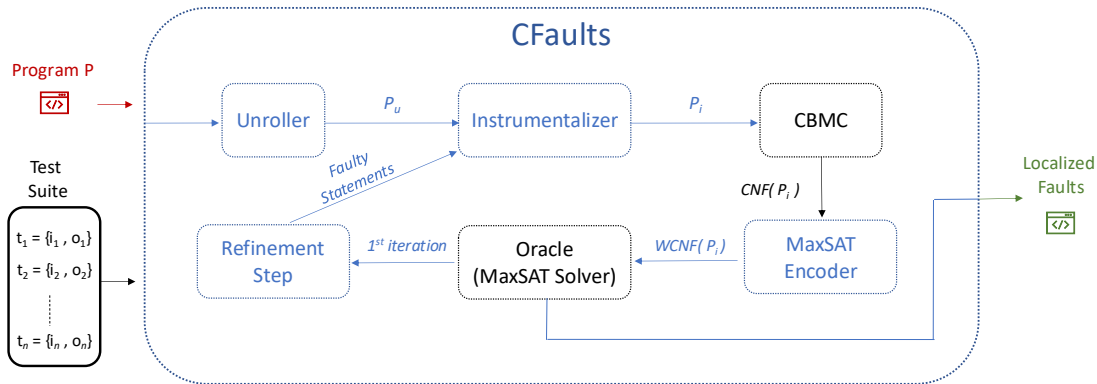
- Given a MaxSAT solution, **the set of unhealthy components ($h(c) = 0$), corresponds to a subset-minimal aggregated diagnosis.**

Model-Based Diagnosis with Multiple Test Cases

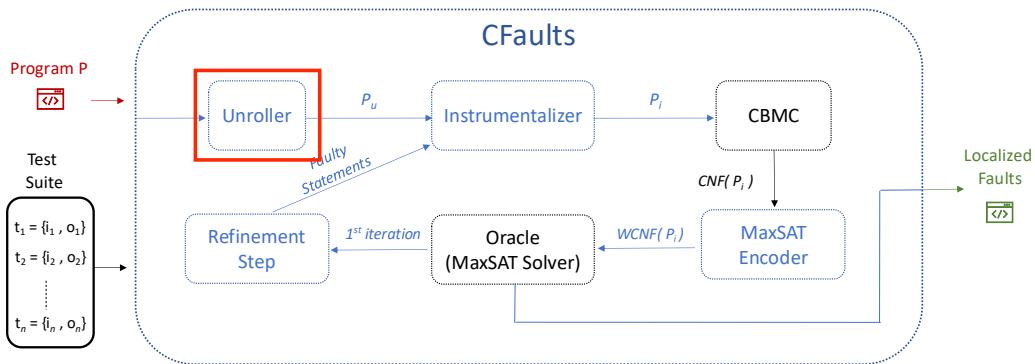
- Given a MaxSAT solution, **the set of unhealthy components** ($h(c) = 0$), **corresponds to a subset-minimal aggregated diagnosis**.
- This diagnosis makes the system **consistent with all observations**, as follows:

$$\bigwedge_{o_i \in \mathcal{O}} (\mathcal{P}_i \wedge o_i) \wedge \bigwedge_{c \in \mathcal{C} \setminus \Delta} h(c) \wedge \bigwedge_{c \in \Delta} \neg h(c) \not\models \perp$$

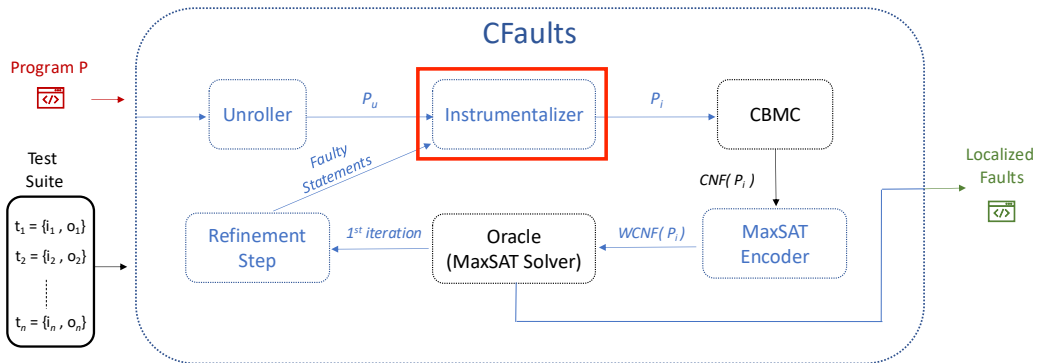
CFaults



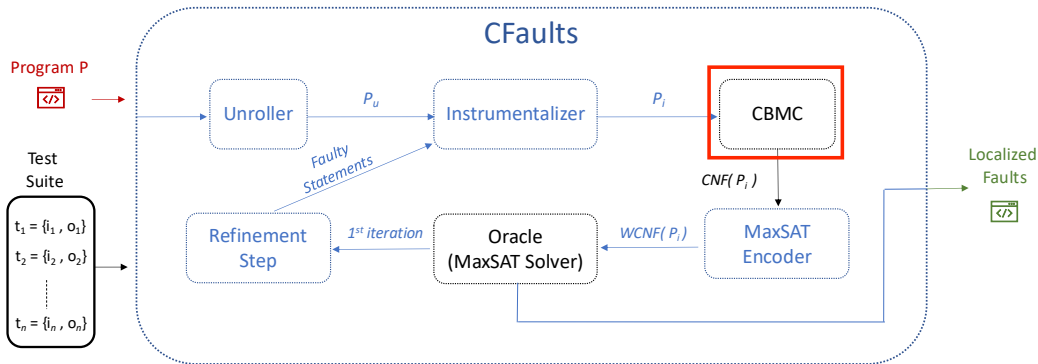
CFaults



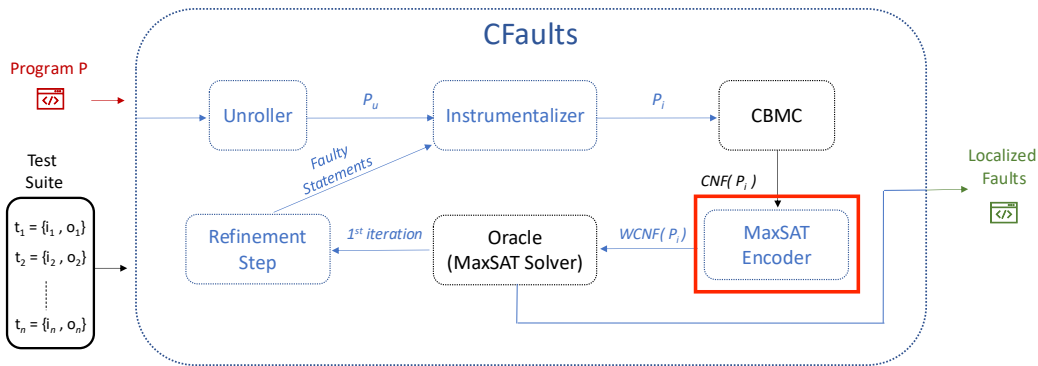
CFaults



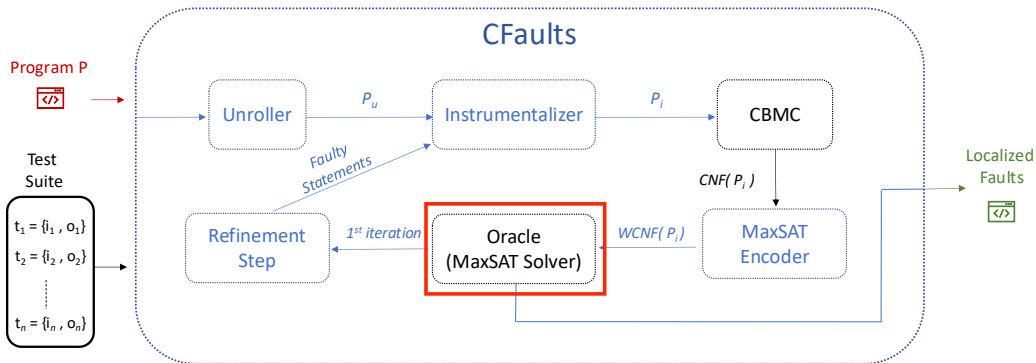
CFaults



CFaults



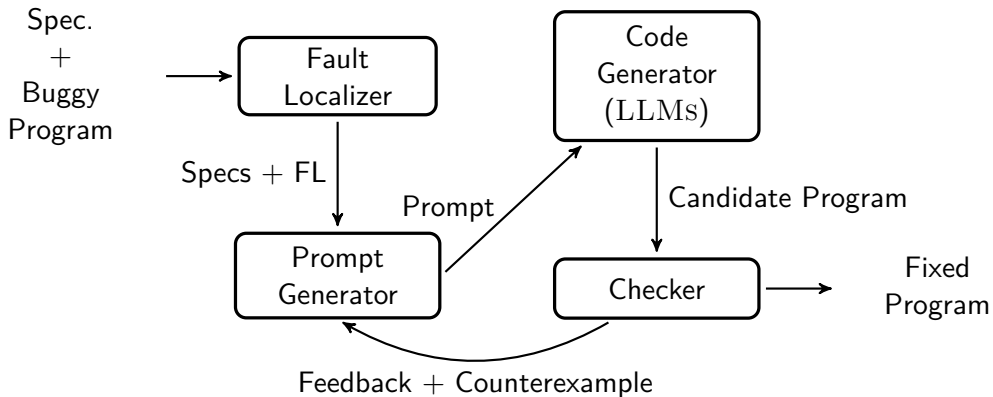
CFaults



Counterexample Guided Program Repair

- **AAAI 2025** - Counterexample Guided APR Using MaxSAT-based Fault Localization.

Counterexample Guided Automated Repair



Prompt Example without Fault Localization

```
Fix all semantic bugs in the buggy program below. Modify the code as little as possible.
Do not provide any explanation.

### Problem Description ###
Write a program that determines and
prints the largest of three integers
given by the user.

### Test Suite
#input:
6 2 1
#output:
6
// The other input-output tests

# Reference Implementation
(Do not copy this program) <c> #
```c
int main(){
 // Reference Implementation
}
...

Buggy Program <c>
```c
int main(){
    // Buggy program from Listing 1
}
...

### Fixed Program <c> ###
```c
```

# Prompt with Fault Localization (Sketches)

---

Complete all the '@ HOLES N @' in the incomplete program below.

Modify the code as little as possible.  
Do not provide any explanation.

### Problem Description ###

Write a program that determines and prints the largest of three integers given by the user.

### Test Suite

#input:

6 2 1

#output:

6

// The other input-output tests

```
Reference Implementation
(Do not copy this program) <c> #
```c
```

```
int main(){
    // Reference Implementation
}
...

```

Incomplete Program <c>

```
```c
int main(){
 // Buggy program from Listing 1
}
...

```

### Complete Program <c> ###

```
```c
```

Experimental Results

Experimental Setup

- **Evaluation Benchmark:** C-PACK-IPAs, a set of twenty-five IPAs, comprising 1431 faulty programs;

Experimental Setup

- **Evaluation Benchmark:** C-PACK-IPAs, a set of twenty-five IPAs, comprising **1431 faulty programs**;
- **Large Language Models (LLMs):** We evaluated **six different LLMs**.

Experimental Setup

- **Evaluation Benchmark:** C-PACK-IPAS, a set of twenty-five IPAS, comprising **1431 faulty programs**;
- **Large Language Models (LLMs):** We evaluated **six different LLMs**.
 - **Three of these models are LLMCs**, i.e., LLMs fine-tuned for coding tasks:
 - IBM's GRANITE;
 - Google's CODEGEMMA;
 - Meta's CODELLAMA.

Experimental Setup

- **Evaluation Benchmark:** C-PACK-IPAs, a set of twenty-five IPAs, comprising 1431 faulty programs;
- **Large Language Models (LLMs):** We evaluated **six different LLMs**.
 - **Three of these models are LLMCs**, i.e., LLMs fine-tuned for coding tasks:
 - IBM's GRANITE;
 - Google's CODEGEMMA;
 - Meta's CODELLAMA.
 - **The other three models are general-purpose LLMs:**
 - Google's GEMMA;
 - Meta's LLAMA3;
 - Microsoft's PHI3.

Experimental Setup

- **Evaluation Benchmark:** C-PACK-IPAs, a set of twenty-five IPAs, comprising **1431 faulty programs**;
- **Large Language Models (LLMs):** We evaluated **six different LLMs**.
 - **Three of these models are LLMCs**, i.e., LLMs fine-tuned for coding tasks:
 - IBM's GRANITE;
 - Google's CODEGEMMA;
 - Meta's CODELLAMA.
 - **The other three models are general-purpose LLMs:**
 - Google's GEMMA;
 - Meta's LLAMA3;
 - Microsoft's PHI3.
- Experiments were conducted using a memory limit of **10GB**, and a timeout of **90s**.

LLM-Driven APR with CFaults

LLMs	Prompt Configurations				Portfolio (All Configurations)
	De-TS	De-TS-CE	Sk_De-TS	Sk_De-TS-CE	
CodeGemma	597 (41.7%)	606 (42.3%)	682 (47.7%)	688 (48.1%)	823 (57.5%)
CodeLlama	492 (34.4%)	500 (34.9%)	573 (40.0%)	561 (39.2%)	712 (49.8%)
Gemma	496 (34.7%)	492 (34.4%)	532 (37.2%)	534 (37.3%)	670 (46.8%)
Granite	626 (43.7%)	624 (43.6%)	691 (48.3%)	681 (47.6%)	846 (59.1%)
Llama3	564 (39.4%)	590 (41.2%)	578 (40.4%)	591 (41.3%)	851 (59.5%)
Phi3	494 (34.5%)	489 (34.2%)	547 (38.2%)	535 (37.4%)	621 (43.4%)
Portfolio (All LLMs)	842 (58.8%)	846 (59.1%)	900 (62.9%)	907 (63.4%)	1013 (70.8%)
Verifix	90 (6.3%)				
Clara	495 (34.6%)				

Table 1: The number of programs fixed by each LLM under various configurations. Mapping abbreviations to configuration names: **De** - IPA *Description*, **TS** - *Test Suite*, **CE** - *Counterexample*, **SK** - *Sketches*.

Discussion

- CLARA repairs 495 programs (34.6%);

Discussion

- CLARA repairs 495 programs (34.6%);
- VERIFIX repairs only 91 programs (6.3%);

Discussion

- CLARA repairs 495 programs (34.6%);
- VERIFIX repairs only 91 programs (6.3%);
- **All six LLMs using different prompt configurations repair more programs than traditional APR tools;**

Discussion

- CLARA repairs 495 programs (34.6%);
- VERIFIX repairs only 91 programs (6.3%);
- **All six LLMs using different prompt configurations repair more programs than traditional APR tools;**
- Incorporating **FL-based Sketches** allows LLMs to repair more programs;

Discussion

- CLARA repairs 495 programs (34.6%);
- VERIFIX repairs only 91 programs (6.3%);
- **All six LLMs using different prompt configurations repair more programs than traditional APR tools;**
- Incorporating **FL-based Sketches** allows LLMs to repair more programs;
- Including **a reference implementation** allows for more repaired programs but **with less efficient fixes** (see our paper);

Discussion

- CLARA repairs 495 programs (34.6%);
- VERIFIX repairs only 91 programs (6.3%);
- **All six LLMs using different prompt configurations repair more programs than traditional APR tools;**
- Incorporating **FL-based Sketches** allows LLMs to repair more programs;
- Including **a reference implementation** allows for more repaired programs but **with less efficient fixes** (see our paper);
- Our CEGIS approach **significantly improves the accuracy of LLM-driven APR across various configurations;**

Takeaway Message

- We tackle the APR problem using an **LLM-Driven Counterexample Guided Inductive Synthesis (CEGIS) approach** [Solar-Lezama et al., 2006];

Takeaway Message

- We tackle the APR problem using an **LLM-Driven Counterexample Guided Inductive Synthesis (CEGIS) approach** [Solar-Lezama et al., 2006];
- We employ **MaxSAT-based Fault Localization to guide and minimize LLMs' patches** to incorrect programs by feeding them bug-free program sketches;

Takeaway Message

- We tackle the APR problem using an **LLM-Driven Counterexample Guided Inductive Synthesis (CEGIS)** approach [Solar-Lezama et al., 2006];
- We employ **MaxSAT-based Fault Localization** to guide and minimize LLMs' **patches** to incorrect programs by feeding them bug-free program sketches;
- With our approach **all six evaluated LLMs fix more programs and produce smaller patches** than other configurations and symbolic tools;

Takeaway Message

- We tackle the APR problem using an **LLM-Driven Counterexample Guided Inductive Synthesis (CEGIS) approach** [Solar-Lezama et al., 2006];
- We employ **MaxSAT-based Fault Localization to guide and minimize LLMs' patches** to incorrect programs by feeding them bug-free program sketches;
- With our approach **all six evaluated LLMs fix more programs and produce smaller patches** than other configurations and symbolic tools;
- Our code is available on GitHub and on Zenodo.

Thank you!



<https://cs.ox.ac.uk/people/pedro.orvalho>

References



Ahmed, Umair Z and Fan, Zhiyu and Yi, Jooyong and Al-Bataineh, Omar I and Roychoudhury, Abhik (2022)

Verifix: Verified repair of programming assignments.

TOSEM 22 12(3), 45 – 678.



Gulwani, Sumit and Radiček, Ivan and Zuleger, Florian (2018)

Automated clustering and program repair for introductory programming assignments.

PLDI 18 52(4), 465 – 480.



Armando Solar-Lezama and Liviu Tancau and Rastislav Bodík and Sanjit A. Seshia and Vijay A. Saraswat (2018)

Combinatorial sketching for finite programs.

ASPLOS 2006.



Reiter, Raymond (1987)

A Theory of Diagnosis from First Principles.

Artif. Intell. 1987.

References



Ignatiev, Alexey and Morgado, António and Weissenbacher, Georg and Marques-Silva, João (2019)
Model-Based Diagnosis with Multiple Observations.
IJCAI 2019.



P. Orvalho and M. Janota and V. Manquinho (2024)
C-Pack of IPAs: A C90 Program Benchmark of Introductory Programming Assignments.
Automated Program Repair (APR) 2024.



The Guardian UK - CrowdStrike Meltdown
<https://www.theguardian.com/technology/article/2024/jul/24/crowdstrike-outage-companies-cost>.
The Guardian UK.



P. Orvalho and M. Janota and V. Manquinho (2024)
CFaults: Model-Based Diagnosis for Fault Localization in C with Multiple Test Cases.
Formal Methods (FM) 2024.



P. Orvalho and M. Janota and V. Manquinho (2025)
Counterexample Guided Program Repair Using Zero-Shot Learning and MaxSAT-based Fault Localization.
AAAI 2025.