ChangeGuard: Validating Code Changes via Pairwise *Learning-Guided Execution*

Lars Gröninger, Beatriz Souza, and Michael Pradel



European Research Council

Established by the European Commission



Software is Continuously Evolving



New Features



Bug Fixes



Refactoring



Optimization



Motivating Example



Motivating Example



Motivating Example

Fix bug involving OR condition					
⊮ master (#2643) · ⊙ 2.13.1 ··· 1.4	4.0			1 parent <mark>9d97d7</mark>	3 commit 49c5afc
Q Filter files		1 file	changed +4 -1 lines changed	Q Search within code	鐐
🗸 盲 scrapy/downloadermiddle	~	scrap	y/downloadermiddlewares/retry.py 🖵 📫		+4 -1 00000 ····
± retry.py	.1		<pre>@@ -63,7 +63,10 @@ def process_exception(self, request, exceptio</pre>	n, spider):	
	63	63	<pre>def _retry(self, request, reason, spider):</pre>		
	64	64	<pre>retries = request.meta.get('retry_times', 0) + 1</pre>		
	65	65			
	66		<pre>- retry_times = request.meta.get('max_retry_times') or sel</pre>	f.max_retry_times	
		66	<pre>+ retry_times = self.max_retry_times</pre>		
		67	+		
		68	+ if 'max_retry_times' in request.meta:		
		69	<pre>+ retry_times = request.meta['max_retry_times']</pre>		
	67	70			
	68	71	<pre>stats = spider.crawler.stats</pre>		
	69	72	<pre>if retries <= retry_times:</pre>		

How to Find Semantics-Breaking Changes?



How to Find Semantics-Breaking Changes?



Executing is not always easy!

FSE'23

LExecutor: Learning-Guided Execution

Beatriz Souza University of Stuttgart Stuttgart, Germany beatrizbzsouza@gmail.com Michael Pradel University of Stuttgart Stuttgart, Germany michael@binaervarianz.de



Learning-guided approach for executing arbitrary code snippets

- Predict missing values with neural model
- <u>Inject values</u> into the execution

Example of Incomplete Code



LExecuting the Incomplete Code



ChangeGuard

Pairwise Learning-guided approach for validating code changes

 Key idea is to execute two versions of a code snippet side-by-side, while predicting and injecting any missing values into the execution.

Overview of ChangeGuard



Overview of ChangeGuard



• Merge f_old and f_new into Comparison Program

```
def f_old():
    # old code
```

```
def f_new():
    # new code
```

f_old()
f_new()

Code to compare behavior

• Predict diverse and project-specific values (!= LExecutor)

Abstract value	Concrete value(s)
None	None
Boolean	True, False
Integer	–100, –10, –1, 0, 1, 10, 100, and all integer literals in the code
Float	-100.0, -10.0, -1.0, 0.0, 1.0, 10.0, 100.0, and all float literals in the code
String	"", "a", and all string literals in the code
List	List with random size and elements of a randomly selected type
Tuple	Tuple with random size and elements of a randomly selected type
Dictionary	Dictionary with random size that maps strings to values of a random type
Set	Set with random size and elements of a randomly selected type
Callable	Class of a versatile object with various default methods
Resource	Versatile object with various default methods
Object	Versatile object with various default methods



Code to compare behavior

More details in the paper:

- Handling Calls to External Functions
- Handling Indexing Operations

Overview of ChangeGuard



Compare Execution Behavior

- Repeat pairwise learning-guided execution (k = 300)
- After each execution of f_old and f_new **compare**:
 - Argument and return values
 - Output written to console
 - Called functions
 - Raised exceptions

Overview of ChangeGuard



- RQ1: Effectiveness
- RQ2: Comparison with Regression Testing
- RQ3: Accuracy of the Neural Model
- RQ4: Robustness and Coverage
- RQ5: Efficiency

Project

By commit me	ssage	
Sempreserving Sem	changing	
31	15	
3	15	
9	15	
3	15	
6	15	
27	15	
3	15	
37	15	
20	15	
10	15	
149	150	
(🗸) 🕓 🖳	refactor", "simplify", "cleanup", "opti	mi
-	By commit me Sempreserving Sem 31 3 9 3 6 27 3 6 27 3 3 7 20 10 149	By commit message Sempreserving Semchanging 31 15 3 15 9 15 3 15 6 15 27 15 37 15 20 15 10 15 149 150 "refactor", "simplify", "cleanup", "optication of the second o

Project	Code changes								
\square	By comm	it message	Manually annotated						
	Sempreserving	Semchanging	Sempreserving	Semchanging	Unclear				
Airflow	31	15	15	25	6				
Black	3	15	4	13	1				
FastAPI	9	15	8	12	4				
Flask	3	15	5	10	3				
HTTPie	6	15	7	14	0				
Pandas	27	15	5	13	24				
Poetry	3	15	2	11	5				
Scikit-Learn	37	15	23	16	13				
Scrapy	20	15	15	12	8				
TheAlgorithms	10	15	9	5	11				
Total	149	150	93	131	75				
	147	150	75	151					

224 single function changes

Project

\square	By comm	it message	
	Sempreserving	Semchanging	
Airflow	31	15	
Black	3	15	
FastAPI	9	15	
Flask	3	15	
HTTPie	6	15	
Pandas	27	15	
Poetry	3	15	
Scikit-Learn	37	15	
Scrapy	20	15	
TheAlgorithms	10	15	165 changes
Total	149	150	(expected to be semantics-preserving
			5

1. RIdiom: Making Python Code Idiomatic by Automatic Refactoring Non-Idiomatic Python Code with Pythonic Idioms (FSE'22)

Project

e			
\square	By commi	t message	
	Sempreserving	Semchanging	
Airflow	31	15	
Black	3	15	
FastAPI	9	15	
Flask	3	15	
HTTPie	6	15	
Pandas	27	15	
Poetry	3	15	Prompt: "Improve code quality while preserving behavior"
Scikit-Learn	37	15	
Scrapy	20	15	
TheAlgorithms	10	15	GPT-3.5: 187 changes
Total	149	150	GPT-4: 258 changes

RQ1: Effectiveness on Manually Annotated Changes



RQ1: Effectiveness on Manually Annotated Changes





Refactoring tool	Prediction				
	Inconclusive	Changing	Preserving		
RIdiom	38	5	122		
GPT-3.5	31	87	69		
GPT-4	38	143	77		



Refactoring tool]	Prediction					
	Inconclusive	Changing	Preserving				
RIdiom	38	5	122				
GPT-3.5	31	87	69	bug in Kidion			
GPT-4	38	143	77				









Code Refactored by GPT-4

- def param_allowed(stat_name, include, exclude):
 - if not include and not exclude: return True
- **for** p **in** exclude:
 - if p in stat_name: return False
- if exclude and not include: return True
- **for** p **in** include:
- if p in stat_name: return True
- **return** False
- + if any(p in stat_name for p in exclude): return False
- + if include: return any(p in stat_name for p in include)
- + return not exclude

Code Refactored by GPT-4

def param_allowed(stat_name, include, exclude) :

if not include and not exclude: return True

- **for** p **in** exclude:

- if p in stat_name: return False
- if exclude and not include: return True
- **for** p **in** include:
- if p in stat_name: return True
- return False
- + if any(p in stat_name for p in exclude): return False
- + if include: return any(p in stat_name for p in include)
- + return not exclude

Behavior changes for args = [], [], [...]

RQ2: Comparison with Existing Regression Tests

Recall: 7.6%

Ground t	ruth		ChangeGua	urd	R	egression te	sting
	Total	Changing	Preserving	Inconclusive	Changing	Preserving	Inconclusive
Changing	131	91	12	28	10	29	92 🗖
Preserving	93	27	48	18	2	18	73
Total	224	118	60	46	12	47	165
					Precision:	83.3%	No test resu
						Accura	acy: 47.5%

RQ2: Comparison with Existing Regression Tests

Recall: 7.6%

Ground t	ruth		ChangeGuard		Guard Regression testing		sting
	Total	Changing	Preserving	Inconclusive	Changing	Preserving	Inconclusive
Changing	131	91	12	28	10	29	92
Preserving	93	27	48	18	2	18	73
Total	224	118	60	46	12	47	165
					Precision:	83.3%	No test res
						Accura	acy: 47.5%

RQ5: Efficiency

- Instrumentation: 1.15 seconds
- First execution: 1.67 seconds
- Extra executions: 1.01 seconds

Executions to reach a conclusion



RQ5: Efficiency

- Instrumentation: 1.15 seconds
- First execution: 1.67 seconds
- Extra executions: 1.01 seconds

Executions to reach a conclusion

-> Reducing the number of executions, ChangeGuard becomes more efficient for a small reduction in recall!





How to Find Semantics-Breaking Changes?





Overview of ChangeGuard



RQ1: Effectiveness on Changes by



Refactoring tool]	Prediction		
	Inconclusive	Changing	Preserving	
RIdiom	38	5	122	Pug in Didiam
GPT-3.5	31	87	69	Bug in Ridiom
GPT-4	38	143	77	
-> LLMs often fail to improve the preserving its semantics; -> ChangeGuard is effective to i semantics-breaking code chang	e code while dentify such ges!	0 random cases (per model)	→ GPT-3.5: 21 GPT-4: 16 c	changing hanging

Code Refactored by RIdiom

- provider = Provider(self._pkg, self._pool, self._io)
- locked = {}
- for package in self._locked.packages: locked[package.name] = package
- + provider, locked = Provider(self._pkg, self._pool, self._io), {}

Code Refactored by RIdiom

- provider = Provider(self._pkg, self._pool, self._io)
- locked = {}

- for package in self._locked.packages: locked[package.name] = package

+ provider, locked = Provider(self._pkg, self._pool, self._io), {}

Line removed!

Code Refactored by GPT-3.5

- if isinstance(arr_or_dtype, ExtensionType): return arr_or_dtype.name == "category"
- if arr_or_dtype is None: return False
- return CategoricalDtype.is_dtype(arr_or_dtype)
- + if isinstance(arr_or_dtype, ExtensionType) and arr_or_dtype.name == "category": return True
- + elif arr_or_dtype is None: return False
- + **else**: **return** CategoricalDtype.is_dtype(arr_or_dtype)

Code Refactored by GPT-3.5

- if isinstance(arr_or_dtype, ExtensionType): return arr_or_dtype.name == "category"

- if arr_or_dtype is None: return False
- return CategoricalDtype.is_dtype(arr_or_dtype)
- + **if isinstance**(arr_or_dtype, ExtensionType) **and** arr_or_dtype.name == "category": **return** True
- + elif arr_or_dtype is None: return False
- + **else**: **return** CategoricalDtype.is_dtype(arr_or_dtype)

Expression moved to condition!

RQ3: Accuracy of the Neural Model

Fine-tuned CodeT5 with the following differences from LExecutor:

- Predicts values for indexing operations;
- Larger and diverse training dataset.

	Accuracy	LExecutor Accuracy (coarse-grained):
Top-1	95.13% =	→ 88.1%
Top-3	96.17%	
Top-5	97.52%	

RQ4: Robustness and Coverage

