# Expanding the Reach of Fuzz Testing
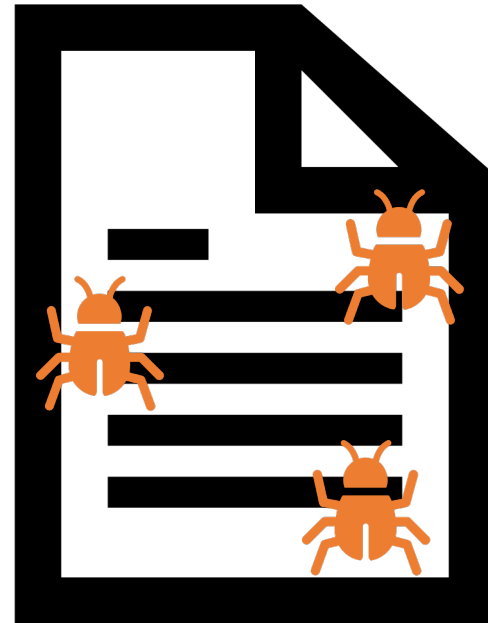
## Caroline Lemieux
*The University of British Columbia*

CSER New Faculty Talk
*June 7, 2023*

# Software Has Bugs

Caroline Lemieux — Expanding the Reach of Fuzz Testing

# Bugs Have Increasing Consequences

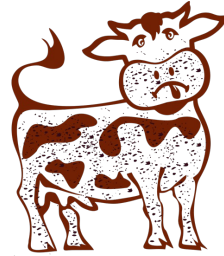Caroline Lemieux — Expanding the Reach of Fuzz Testing

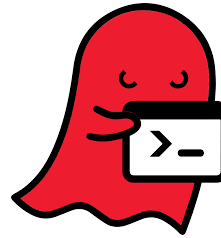# Bugs Have Increasing Consequences



Badlock  Cloudbleed  Dirty COW  GHOST  Heartbleed  StageFright  ShellShock

# Bugs Have Increasing Consequences



Badlock    Cloudbl...   GHOST    Heartbleed    StageFright    ShellShock

*The Cost of Poor Software Quality in the US: A 2020 Report*

Herb Krasner
Member, Advisory Board
Consortium for Information & Software Quality™ (CISQ™)
WWW.IT-CISQ.ORG
HKRASNER@UTEXAS.EDU
Date: January 1, 2021

**CISQ**
Consortium for Information & Software Quality ™

the cost of
poor quality software is
# $2.08 trillion
in the US in 2020

# My Work:

build tools to help developers improve **correctness**, **security** and **performance** of software

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)



Using *large language models* to improve *automated test suite generation*

CodaMOSA (ICSE'23)

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)

```
Search stuck with generated test cases:

def test_case_1():          def test_case_2():
  str_0 = 'a\t!sUo~AU'        str_0 = None
  str_1 = bump_version(       int_0 = 1431
        str_0)               str_1 = bump_version(
                                    str_0, int_0)

        get test case hint as code

def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Using *large language models* to improve *automated test suite generation*

CodaMOSA (ICSE'23)

# Performance Bugs

# Performance Bugs

# Performance Bugs

Takes **1+ minute** to report syntax error:

```
(((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((e foo = 1; => 1;
```

# Performance Bugs

## CVE-2020-7212 Detail

### Current Description

The _encode_invalid_chars function in util/url.py in the urllib3 library 1.25.2 through 1.25.7 for Python allows a denial of service (CPU consumption) because of an inefficient algorithm. The percent_encodings array contains all matches of percent encodings. It is not deduplicated. For a URL of length N, the size of percent_encodings may be up to O(N). The next step (normalize existing percent-encoded bytes) also takes up to O(N) for each step, so the total time is O(N^2). If percent_encodings were deduplicated, the time to compute _encode_invalid_chars would be O(kN), where k is at most 484 ((10+6*2)^2).

+View Analysis Description

### Severity

| CVSS Version 3.x | CVSS Version 2.0 |

**CVSS 3.x Severity and Metrics:**

**NIST:** NVD          **Base Score:** 7.5 HIGH

Takes **1+ minute** to report syntax error:

```
((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((((e foo = 1; => 1;
```

# Performance Bugs

## 🐛 CVE-2020-7212 Detail

### Current Description

The _encode_invalid_chars function in util/url.py
through 1.25.7 for Python allows a denial of serv
because of an inefficient algorithm. The percent
all matches of percent encodings. It is not dedup
N, the size of percent_encodings may be up to O
(normalize existing percent-encoded bytes) also
step, so the total time is O(N^2). If percent_enc
the time to compute _encode_invalid_chars wo
most 484 ((10+6*2)^2).

+ View Analysis Description

| Severity | CVSS Version 3.x | CVSS Version 2.0 |
|---|---|---|

**CVSS 3.x Severity and Metrics:**

NIST: NVD    Base Score: `7.5 HIGH`

Takes **1+ minute** to report syntax error:

`((((((((((((((((e foo = 1; => 1;`

Can we find inputs revealing these bugs automatically? 🤔

| | | |
|---|---|---|
| IJG jpeg | NetBSD bpf [1] | man & mandoc [1][2][3][4][5]... |
| libtiff [1][2][3] | | |
| Mozilla Firefox | clamav [1][2][3][4][5] | MMIX [1] |
| Adobe Flash / PCRE | clang / llvm [1][2][3][4][5] | dhcpcd [1] |
| LibreOffice [1] | mutt [1] | mbed TLS [1] |
| GnuTLS | pdksh [1][2] | Linux xfs [1] |
| PuTTY [1] | redis / lua-cmsgp | Adobe Reader [1] |
| bash (post-Shells | perl [1][2][3][4][5][6] | OpenBSD kernel [1] |
| pdfium [1] | SleuthKit [1] | MatrixSSL [1] |
| libarchive [1][2][3] | exifprobe [1] | w3m [1][2][3][4] |
| BIND [1][2] | Xerces-C [1][2] | irssi [1][2][3] |
| Oracle Berkeley | exiv [1][2] | Malheur [1] |
| FLAC audio lib | curl [1][2][3] | gdk-pixbuf [1] |
| strings (+ related tool | dnsmasq [1] | lz4 [1] |
| rcs [1] | libwmf [1] | libpcre [1][2][3] |
| Info-Zip unz | imlib2 [1][2][3][4] | openexr [1] |
| | libsass [1] | lrzip [1][2][3] |
| | VLC [1][2] | ytnef [1][2][3][4]... |
| | screen [1][2][3] | Apache httpd [1] |
| | UPX [1] | pev [1][2][3][4] |
| | | Mongoose OS [1] |

# Trophies

Honggfuzz has been used to find a few interesting security problems in major software packages; An incomplete list:

- Multiple exploitable bugs in **IDA-Pro**
- P
- A
- Remote DoS in **Crypto++** • CVE-2016-9939
- Programming language interpreters
  - **PHP/Python/Ruby**
  - PHP WDDX
  - PHP
  - Perl: #1, #2, #3
- Double-free in **LibXMP**
- Heap buffer overflow in SAPCAR • CVE-2017-8852
- Crashes in **libbass**
- V
- **FreeTyp** • Heap buffer-overflow (or UAF) in **MPV**
  - CVE • Heap buffer-overflow in **picoc**
  - CVE • Crashes in **OpenCOBOL**: #1, #2
  - CVE • DoS in **ProFTPD**: #1, #2
  - CVE • Memory corruption in **htmldoc**
  - CVE • Memory corruption in **OpenDetex**
  - CVE • Memory corruption in **Yabasic**
  - CVE • Memory corruption in **Xfig**
- A
- Stack co • **Rust**:
- M
- Infinite le
  - panic() in regex #1, #2, #3
- M
- A coupl
  - panic() in h2 #1, #2, #3
- M
- **Samba**
  - panic() in sleep-parser #1
- Crash in
  - panic() in lewton #1
- Multiple
  - panic()/DoS in Ethereum-Parity #1
- Buffer o
  - crash() in Parts – a GPT partition manager #1
- Linux mem mgmt
  - crashes in rust-bitcoin/rust-lightning #1
- iOS kernel [1]

Trophies

Honggfuzz has been used to find a few interesting security problems in major software packages; An incomplete list:

- Multiple exploitable bugs in **IDA-Pro**
- P
- A
- Remote DoS in **Crypto++** • CVE-2016-9939
- Programming language interpreters
  - **PHP/Python/Ruby**
  - PHP WDDX
  - PHP
  - Perl: #1, #2, #3
- Double-free in **LibXMP**

| IJG jpeg | | | |
|---|---|---|---|
| libtiff [1] [2] [3] | NetBSD bpf [1] | | man & mandoc [1] [2] [3] [4] [5] ... |
| Mozilla Firefox | clamav [1] [2] [3] [4] [5] | MMIX [1] | |
| Adobe Flash / PCRE | clang / llvm [1] [2] [3] [4] [5] | dhcpcd [1] | |
| LibreOffice [1] | mutt [1] | mbed TLS [1] | |
| GnuTLS | pdksh [1] [2] | Linux xfs [1] | |
| PuTTY [1] | redis / lua-cmsgp | Adobe Reader [1] | |
| bash (post-Shells | perl [1] [2] [3] [4] [5] [6] [1] | OpenBSD kernel [1] | |
| pdfium [1] | SleuthKit [1] | | |
| libarchive [1] [2] [3] | exifprobe [1] | | |
| BIND [1] [2] [3] | Xerces-C [1] [2] | | |
| Oracle Berkeley | exiv [1] [2] | | |
| FLAC audio lib | curl [1] [2] [3] | | |
| strings (+ related tool | dnsmasq [1] | | |
| rcs [1] | libwmf [1] | lz4 [1] | |
| Info-Zip unzi | imlib2 [1] [2] [3] [4] | libpcre [1] [2] [3] | |
| | libsass [1] | openexr [1] | |
| | VLC [1] [2] | lrzip [1] [2] [3] | |
| | screen [1] [2] [3] | ytnef [1] [2] [3] [4] ... | |
| | UPX [1] | pev [1] [2] [3] [4] | |
| | | Apache httpd [1] | |
| | Mongoose OS [1] | iOS kernel [1] | |

- CVI • Memory corruption in **Yabasic**
- CVI • Memory corruption in **Xfig**
- A
- Stack co • **Rust**:
- M • Infinite l
  - panic() in regex #1, #2, #3
- M • A coupl
  - panic() in h2 #1, #2, #3
- M • **Samba**
  - panic() in sleep-parser #1
- Crash in
  - panic() in lewton #1
- Multiple
  - panic()/DoS in Ethereum-Parity #1
- Buffer o
  - crash() in Parts – a GPT partition manager #1
- Linux mem mgmt
  - crashes in rust-bitcoin/rust-lightning #1

# How does it work?

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing



Initial

pick

mutate

execute

save

Interesting Feedback?

Input$_n$'

Execution Feedback$_n$

<a>b</a>

<a>b</a>

<aa>b</a>

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing

# Coverage-Guided Fuzzing



Initial

pick

mutate

execute

$ xmllint

save

New Branch Covered?

<a>b</a>

<a>b</a>

<a<u>a</u>>b</a>

<a<u>a</u>>b</a>

tags_match(*input*)   F
is_recoverable(err)   T
...

tags_match(*input*)
T          F
*tag* == "a"          ...
T          F
...          *tag* == "b"

# Coverage-Guided Fuzzing

# Fuzzing to Find Performance Bugs?

# Fuzzing to Find Performance Bugs?



Single-objective fuzzing
→ Gets stuck in local maxima

Initial

Input — pick → Input — mutate → Input$_n$' — execute →

save

Input$_n$' ← Interesting Feedback? ← Execution Feedback$_n$

Input executes for longer?

Execution time of input

# PerfFuzz

Caroline Lemieux — Expanding the Reach of Fuzz Testing

# PerfFuzz

# PerfFuzz



Caroline Lemieux — Expanding the Reach of Fuzz Testing

# PerfFuzz

# PerfFuzz

Caroline Lemieux — Expanding the Reach of Fuzz Testing

# Example program: `wf`

- Count frequency of words in string

input:

```
the quick brown the dog
```

output:

```
brown: 1
dog: 1
quick: 1
the: 2
```

# Example program: wf

- Count frequency of words in string

input:

```
the quick brown the dog
```

output:

```
brown: 1
dog: 1
quick: 1
the: 2
```

Fedora Linux implementation: linked list hash table.
Quadratic worst-case behavior ☹
(when all words hash collide)

# wf results: PerfFuzz Finds True Worst Cases

# wf results: PerfFuzz Finds True Worst Cases

SlowFuzz (single objective maximization) worst case:

```
t r t t s f o Öe r t s f o r t x x t s f o r t x x
```

# wf results: PerfFuzz Finds True Worst Cases

SlowFuzz (single objective maximization) worst case:

```
    t r t t s f o Öe r t s f o r t x x t s f o r t x x
```

PerfFuzz worst case:

```
t <81>v ^?@t <80>!^?@t <80>!t t^Rn t t t t t t t t
```

# wf results: PerfFuzz Finds True Worst Cases

SlowFuzz (single objective maximization) worst case:

t r t t s f o Öe r t s f o r t x x t s f o r t x x

PerfFuzz worst case:

t <81>v ^?@t <80>!^?@t <80>!t t^Rn t t t t t t t t

# PerfFuzz

# Observation: Algorithm is More General



Caroline Lemieux — Expanding the Reach of Fuzz Testing

# Observation: Algorithm is More General

Caroline Lemieux — Expanding the Reach of Fuzz Testing

# FuzzFactory

# Super Fuzzer: Hard Comparisons + Memory Allocations



Caroline Lemieux — Expanding the Reach of Fuzz Testing

# Super Fuzzer: Hard Comparisons + Memory Allocations



LZ4 Bomb
(4GB alloc when decoding 21-byte input)

PNG Bomb
(2GB alloc when reading
~100 byte 20px image)

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)

Using *large language models* to improve *automated test suite generation*

CodaMOSA (ICSE'23)

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)

→ Feedback-directed fuzzing is a flexible abstraction to explore different dimensions of program behavior

*automated test suite generation*

CodaMOSA (ICSE'23)

```
def test...
    str_
    str_
```

get test case hint as code

```
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)

Using *large language models* to improve *automated test suite generation*

CodaMOSA (ICSE'23)

# Test Input Generation (PerfFuzz, etc.)

Generate **inputs** to a **parameterized** test function

the quick brown the dog

t <81>v <80>!^?@t t t

t e q I k b o n t e d g

Parameterized Test Function

```python
def test_wf(document):
    frequencies = wf(document)
    print(frequencies)
```

# Test *Suite* Generation

Generate **test cases** for a file **(e.g., python module, java class)** under test

Module Under test

```python
def test_BST_insert():
    tree = None
    tree = BST_insert(tree, 5)
    tree = BST_insert(tree, 3)
    tree = BST_insert(tree, 7)
```

```python
def test_BST_search():
    tree = Node(5)
    tree.left = Node(3)
    tree.left.left = Node(1)
    res = BST_search(tree, 3)
```

```python
def test_BST_delete():
    tree = Node(5)
    tree.left = Node(3)
    tree.left.left = Node(1)
    tree.left.right = Node(4)
    BST_delete(tree, 4)
```

```python
def BST_insert(tree, to_add):
    # Insert to_add into tree
    <...>

def BST_search(tree, to_search):
    # Search tree for to_search
    <...>

def BST_delete(tree, to_delete):
    # Delete to_delete from tree
    <...>
```

# Search-Based Test Suite Generation

Mutate test cases

Test Case Population

Evolved Test Cases

Update population

Evaluate test cases
on test program

```
<...omitted code...>                    Module under Test

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# Example: Expected Behavior of Function

Mutate test cases

Test Case Population

Evolved Test Cases

```
bump_version('0.0.0') → '0.0.1'
```

```
bump_version('0.0.0', 1) → '0.1.0'
```

```
bump_version('0.0.0', 1, 'a') → '0.1.0a'
```

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Update population

Evaluate test cases
on test program

# Initialize Test Population



Mutate test cases

**Test Case Population**

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

**Evolved Test Cases**

Update population

Evaluate test cases
on test program

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# Current Tests have Low Coverage

Mutate test cases

**Test Case Population**

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
             str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
              str_0, int_0)
```

**Evolved Test Cases**

Update population

Module under Test

```
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases
on test program

# Create New Test Cases via Mutation



Mutate test cases

**Test Case Population**

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

**Evolved Test Cases**

```python
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Update population

Evaluate test cases
on test program

# Mutation Unable to Increase Coverage

Mutate test cases

### Test Case Population

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

### Evolved Test Cases

```python
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Update population

X

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases
on test program

# Search stalled. What to do now?

Mutate test cases

## Test Case Population

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

## Evolved Test Cases

```python
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Update population

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases
on test program

# CodaMOSA: Asks for hints when stuck

Coverage Stalled?

### Test Case Population

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

### Evolved Test Cases

```
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Update population

```
<...omitted code...>                        Module under Test

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# CodaMOSA: Asks for hints when stuck

Coverage Stalled?

*Coverage stall* : *N* iterations without increasing coverage of program under test

### Test Case Population

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

### Evolved Test Cases

```python
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Update population

```python
<...omitted code...>                        Module under Test

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# CodaMOSA: Asks for hints when stuck

Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

## Test Case Population

```python
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

## Evolved Test Cases

```python
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```python
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Evaluate test cases
on test program

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Update population

# Search Stalled

Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

## Test Case Population

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

## Evolved Test Cases

```
def test_case_1m():
    str_0 = 'a\t!sUo~AUUUU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2p():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
    int_1 = 1
    str_1 = bump_version(
            str_0, int_1)
```

Evaluate test cases
on test program

Module under Test

```
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Update population

# Time to Ask for a Hint

Yes, ask for a hint

Coverage Stalled?

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

test_case_1

test_case_2

**Codex**, LLM trained on code
Released Aug 2021

Raw Model Output ("natural" code)

```
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) ->
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Module under

Deserialize to mutatable format

Update population

# Time to Ask for a Hint

**Coverage Stalled?**

Yes, ask for a hint

**Prompt to create low-coverage functions tests**

Ask Codex

No, mutate test cases

**Test Case Population**

test_case_1

test_case_2

Evolved Test Cases

**Raw Model Output ("natural" code)**

Evaluate test cases on test program

**Extracted Test Cases**

Deserialize to mutatable format

Update population

```
<...omitted code...>          Module under Test

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# Time to Ask for a Hint



**Coverage Stalled?**

Yes, ask for a hint

No, mutate test cases

**Prompt to create low-coverage functions tests**

Ask Codex

**Test Case Population**
- test_case_1
- test_case_2

Evolved Test Cases

**Raw Model Output ("natural" code)**

```python
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Update population

Evaluate test cases on test program

**Extracted Test Cases**

```python
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

**Module under Test**

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

# Suggested Test Case Increases Coverage



Coverage Stalled?

Yes, ask for a hint

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

test_case_1

test_case_2

Evolved Test Cases

Raw Model Output ("natural" code)

```python
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Update population

```python
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Module under Test

Evaluate test cases on test program

Extracted Test Cases

```python
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

# Suggested Test Case Increases Coverage



Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

Prompt to create low-coverage functions tests

Ask Codex

Test Case Population

test_case_1

test_case_2

Evolved Test Cases

Raw Model Output ("natural" code)

```python
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Update population

Evaluate test cases on test program

Extracted Test Cases

```python
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

```python
def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```
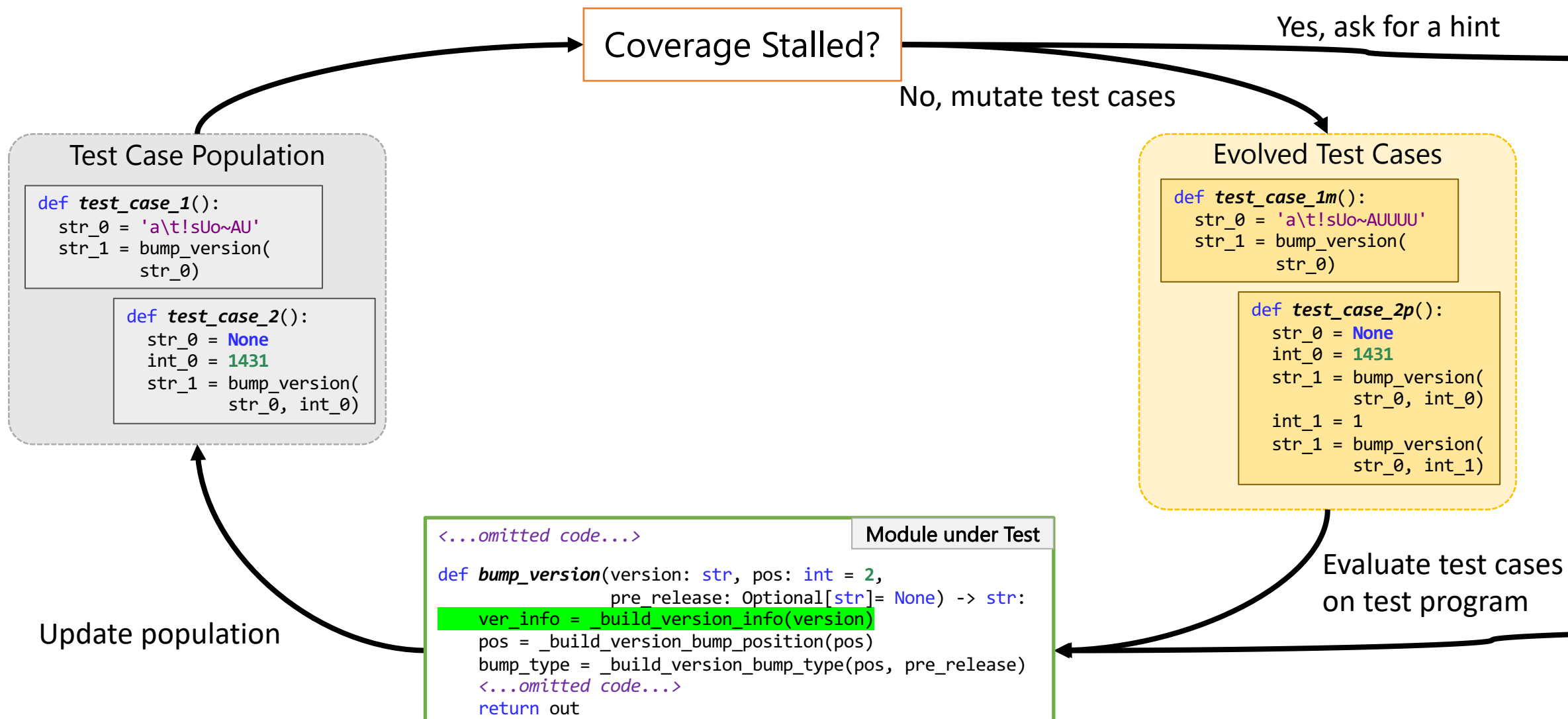
# Suggested Test Case Increases Coverage

Coverage Stalled?

Yes, ask for a hint

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

test_case_1

test_case_2

Evolved Test Cases

Raw Model Output ("natural" code)

```python
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

```python
<...omitted code...>

def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos
<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases on test program

Extracted Test Cases

```python
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

Update population

# Update Population

Coverage Stalled?

Yes, ask for a hint

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

```python
def test_case():
    ...
    int_0 = 1
    ...
```

test_case_2

Evolved Test Cases

```
<...omitted code...>
def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos
<...omitted code...>
def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```
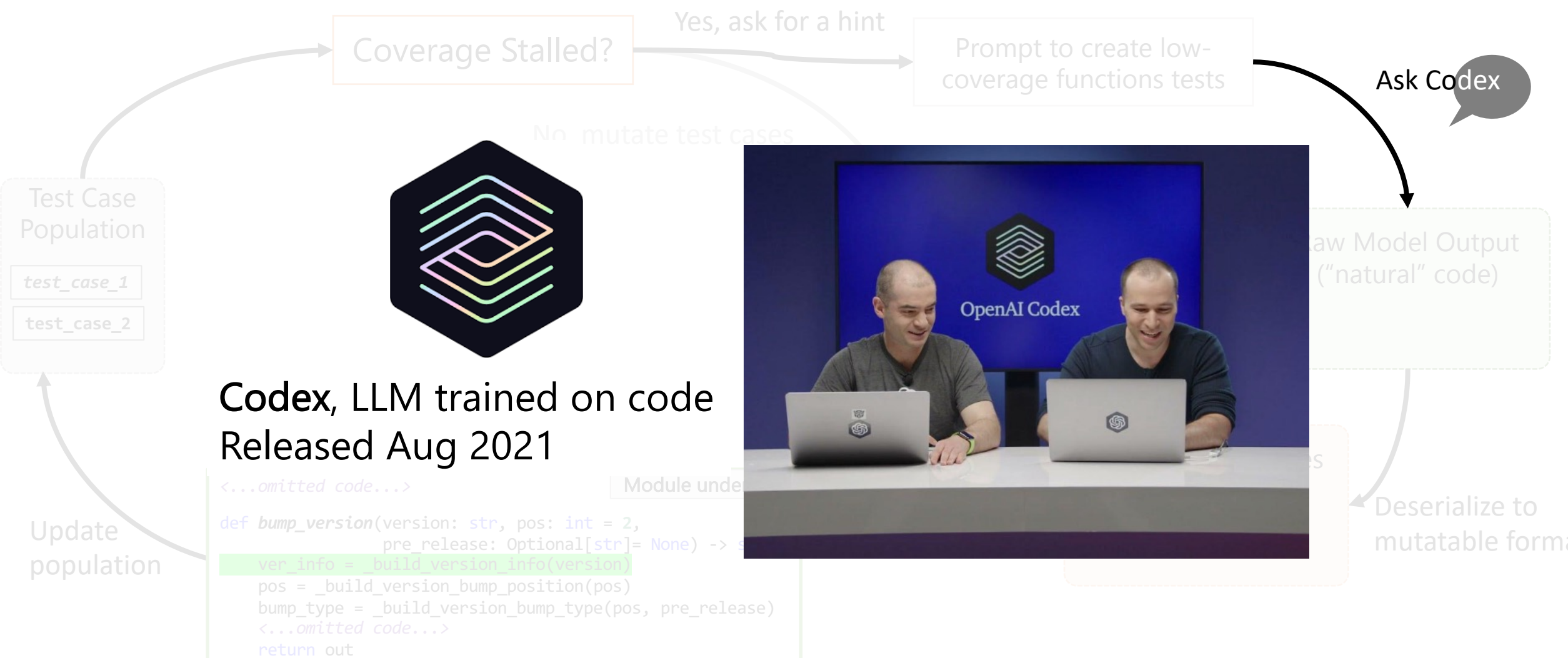
Raw Model Output ("natural" code)

```python
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Evaluate test cases on test program

Extracted Test Cases

```python
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

Update population

# Search No Longer Stalled

Coverage Stalled?

Yes, ask for a hint

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

```
def test_case():
    ...
    int_0 = 1
    ...
```

test_case_2

Evolved Test Cases

Raw Model Output ("natural" code)

```
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

```
<...omitted code...>

def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos

<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```
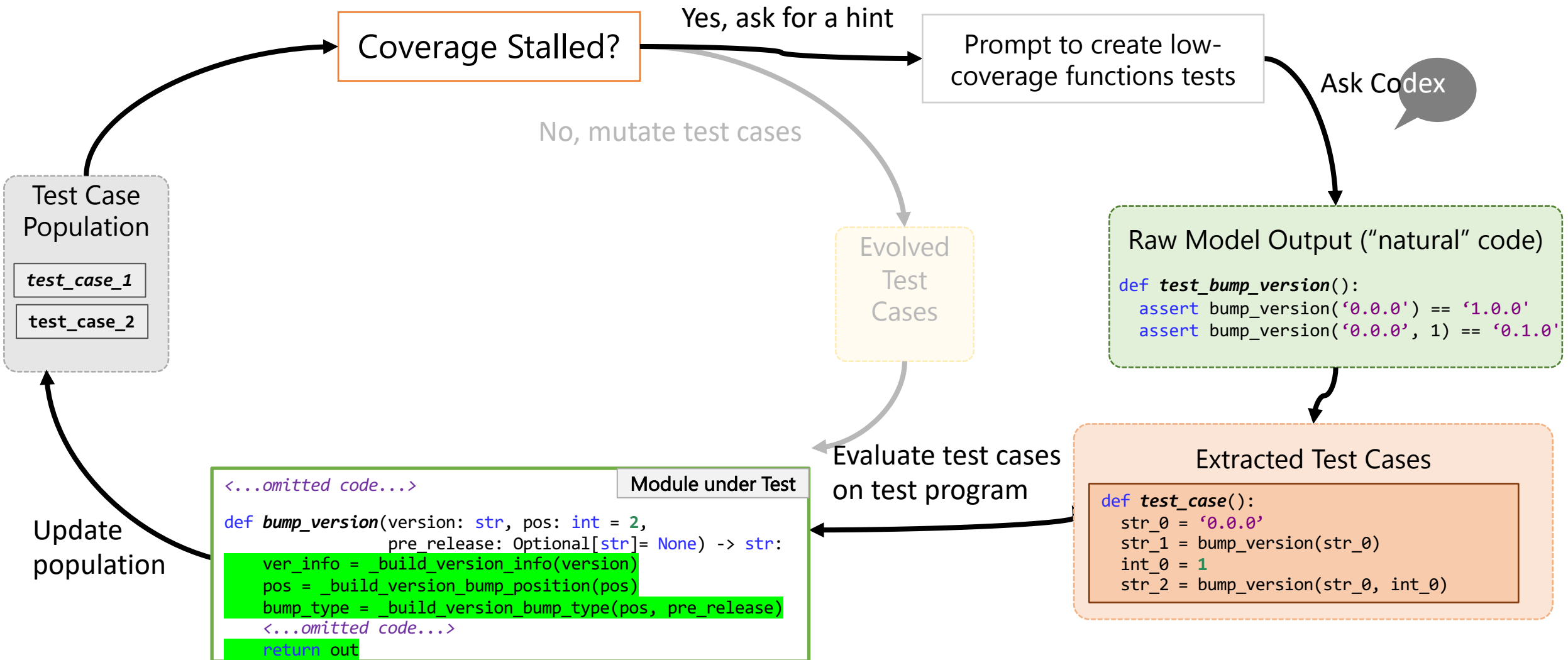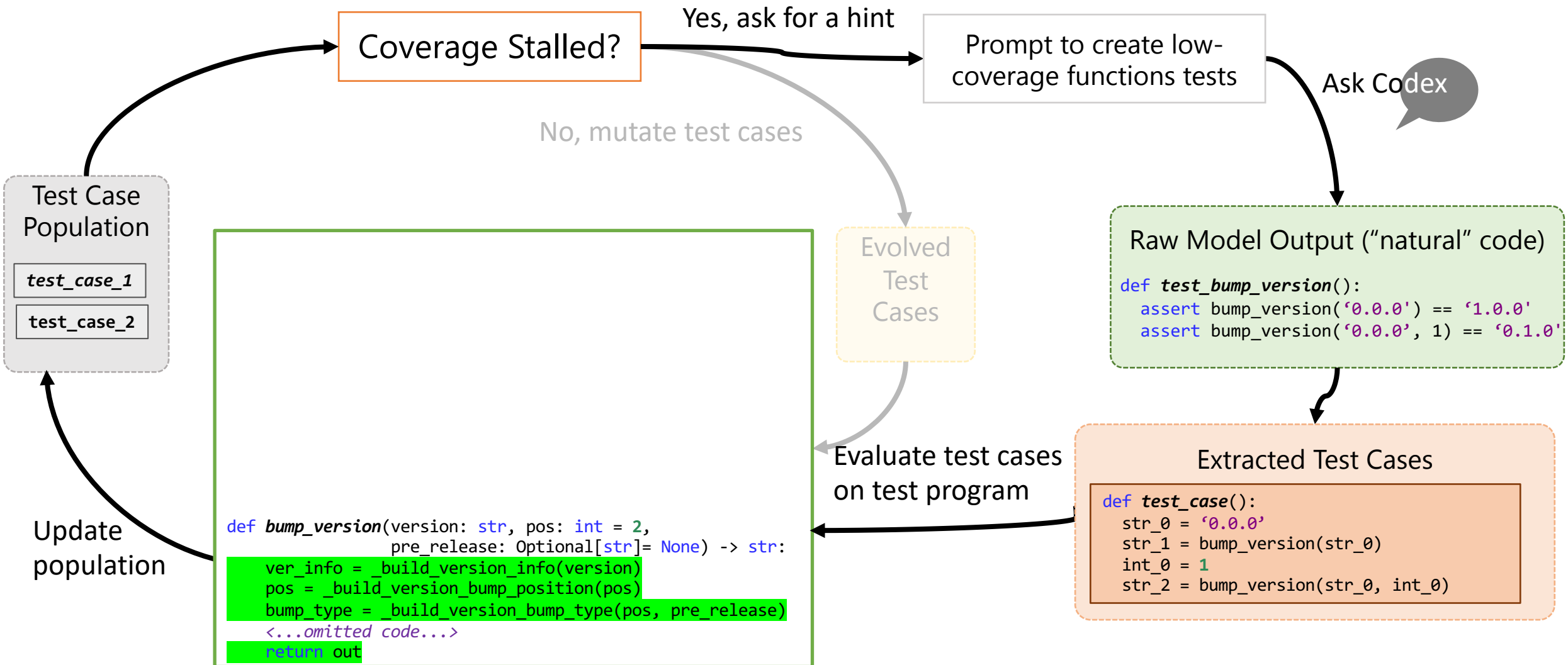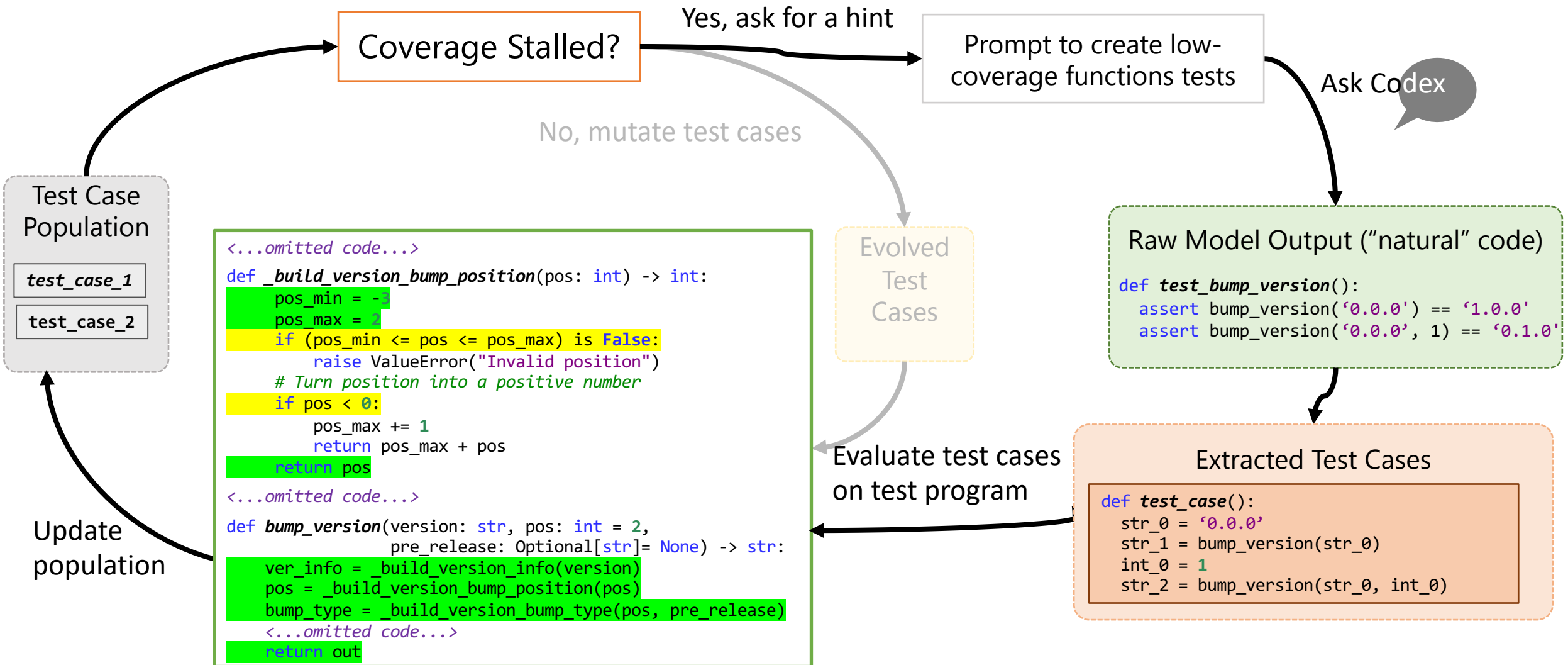
Evaluate test cases on test program

Extracted Test Cases

```
def test_case():
    str_0 = '0.0.0'
    str_1 = bump_version(str_0)
    int_0 = 1
    str_2 = bump_version(str_0, int_0)
```

Update population

# Search No Longer Stalled

Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

## Evolved Test Cases

```
def test_case_mutant_0():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -1
    str_2 = bump_version(
            str_0, int_0)
```

```
def test_case_mutant_n():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -2687
    str_2 = bump_version(
            str_0, int_0)
```

## Test Case Population

```
def test_case():
    ...
    int_0 = 1
    ...
```

test_case_2

```
<...omitted code...>

def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos

<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases
on test program

Update
population

# Search No Longer Stalled

Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

Evolved Test Cases

```python
def test_case_mutant_0():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -1
    str_2 = bump_version(
            str_0, int_0)
```

```python
def test_case_mutant_n():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -2687
    str_2 = bump_version(
            str_0, int_0)
```

Test Case Population

```python
def test_case():
    ...
    int_0 = 1
    ...
```

test_case_2

```python
<...omitted code...>
def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos
<...omitted code...>
def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases on test program

Update population

# Search No Longer Stalled

Coverage Stalled?

Yes, ask for a hint

No, mutate test cases

Evolved Test Cases

```python
def test_case_mutant_0():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -1
    str_2 = bump_version(
            str_0, int_0)
```

```python
def test_case_mutant_n():
    str_0 = '0.0.0'
    str_1 = bump_version(
            str_0)
    int_0 = -2687
    str_2 = bump_version(
            str_0, int_0)
```

Test Case Population

```python
def test_case():
    ...
    int_0 = 1
    ...
```

test_case_2

Update population

```python
<...omitted code...>

def _build_version_bump_position(pos: int) -> int:
    pos_min = -3
    pos_max = 2
    if (pos_min <= pos <= pos_max) is False:
        raise ValueError("Invalid position")
    # Turn position into a positive number
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos

<...omitted code...>

def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases on test program

# Search No Longer Stalled

# Spoiler: Results on this Benchmark

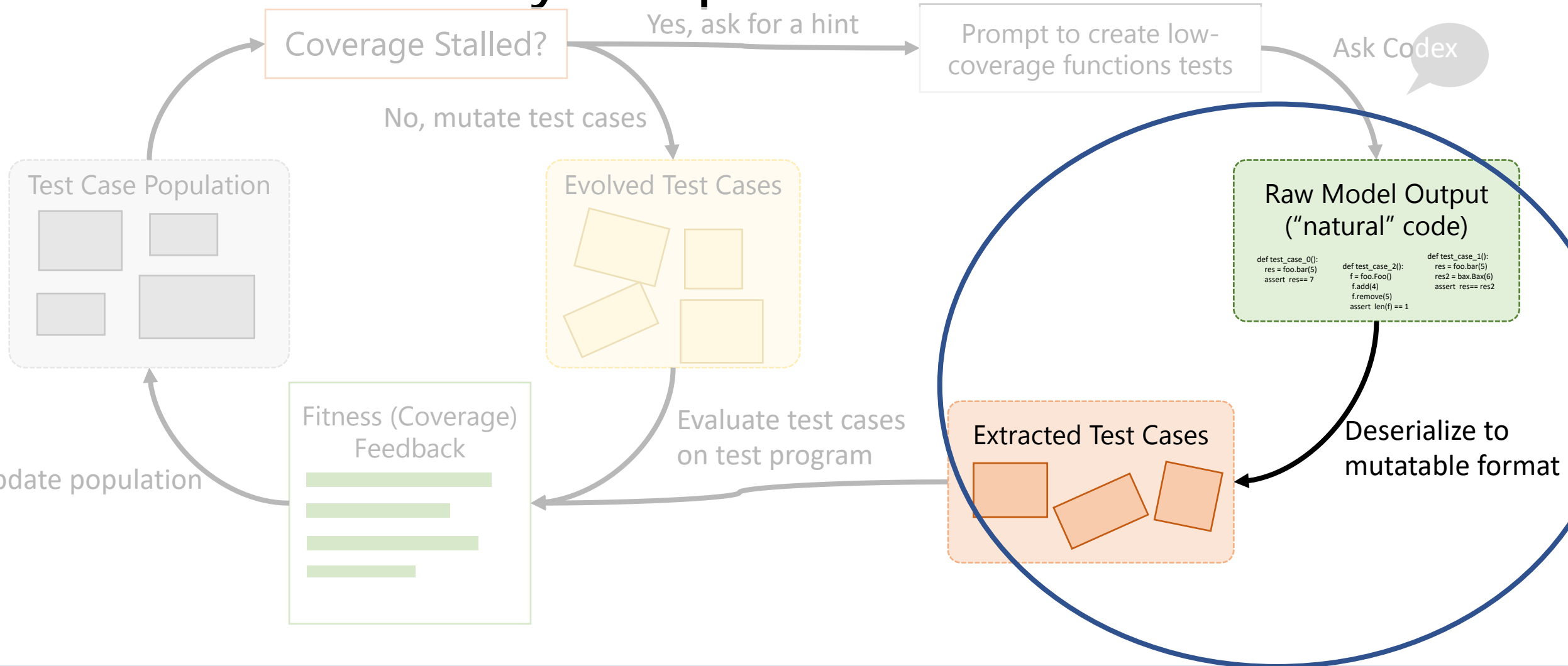Coverage Stalled? — Yes, ask for a hint

No, mutate test cases

Evolved Test Cases

Test Case Population

```
def test_case():
    ...
```

test_case_mutant_0

test_case_mutant_n

```
def test_case_mutant_n():
    str_0 = '0.0.0'
    str_1 = bump_version(
        str_0)
    int_0 = -2687
    str_2 = bump_version(
        str_0, int_0)
```

**SBST-only approach:** 18.4% Coverage
**Codex-only approach:** 64.5% Coverage
**Our approach:** 91.8% Coverage
*(better than the sum of its parts!)*

```
<...omitted code...>
def _build_version_info():
    pos_min
    pos_max
    if (pos_
        rais
    # Turn p
    if pos < 0:
        pos_max += 1
        return pos_max + pos
    return pos
<...omitted code...>
def bump_version(version: str, pos: int = 2,
                 pre_release: Optional[str]= None) -> str:
    ver_info = _build_version_info(version)
    pos = _build_version_bump_position(pos)
    bump_type = _build_version_bump_type(pos, pre_release)
    <...omitted code...>
    return out
```

Evaluate test cases on test program

Update population

# CodaMOSA Approach

# Challenge: When to ask for a hint?



Coverage Stalled? — Yes, ask for a hint → Prompt to create low-coverage functions tests — Ask Codex

No, mutate test cases

Test Case Population

Evolved Test Cases

Raw Model Output ("natural" code)

```
def test_case_0():
    res = foo.bar(5)
    assert res== 7

def test_case_2():
    f = foo.Foo()
    f.add(4)
    f.remove(5)
    assert len(f) == 1

def test_case_1():
    res = foo.bar(5)
    res2 = bax.Bax(6)
    assert res== res2
```

Fitness (Coverage) Feedback

Evaluate test cases on test program

Extracted Test Cases

Deserialize to mutatable format

Update population

# Challenge: How to ask for a hint?



Caroline Lemieux — Expanding the Reach of Fuzz Testing

# Challenge: How to handle (potentially) arbitrary output from Codex?



Coverage Stalled?

Yes, ask for a hint

Prompt to create low-coverage functions tests

Ask Codex

No, mutate test cases

Test Case Population

Evolved Test Cases

Raw Model Output ("natural" code)

```
def test_case_0():
    res = foo.bar(5)
    assert res== 7

                    def test_case_2():
                        f = foo.Foo()
                        f.add(4)
                        f.remove(5)
                        assert len(f) == 1

                                        def test_case_1():
                                            res = foo.bar(5)
                                            res2 = bax.Bax(6)
                                            assert res== res2
```

Update population

Fitness (Coverage) Feedback

Evaluate test cases on test program

Extracted Test Cases

Deserialize to mutatable format

# Solutions Discussed Further in Paper



"CodaMOSA: Escaping Coverage Plateaus in Test Generation with Pre-Trained Large Language Models." *C. Lemieux, J. P. Inala, S. K. Lahiri, S. Sen*. In proceedings of ICSE 2023.

Using *generalized feedback maps* to expand *bugs findable by fuzz testing*

PerfFuzz (ISSTA'18)

FuzzFactory (OOPSLA'19)



Search stuck with generated test cases:

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

get test case hint as code

```
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Using *large language models* to improve *automated test suite generation*

CodaMOSA (ICSE'23)

PerfFuzz
(ISSTA'18)

→ Powerful synergy between LLMs ("what is most expected") and mutative search ("something close, but unexpected")

FuzzFactory
(OOPSLA'19)

Using *large language models* to improve *automated test suite generation*

CodaMOSA
(ICSE'23)

Search stuck with generated test cases:

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
                str_0, int_0)
```

get test case hint as code

```
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

# So Far: Innovations in Fuzzing Algorithms



PerfFuzz
(ISSTA'18)

FuzzFactory
(OOPSLA'19)

Fuzzing

CodaMOSA
(ICSE'23)

# Next: Solving Problems **Around** Fuzzing



Fuzzing

# Next: Solving Problems **Around** Fuzzing

Fuzzing

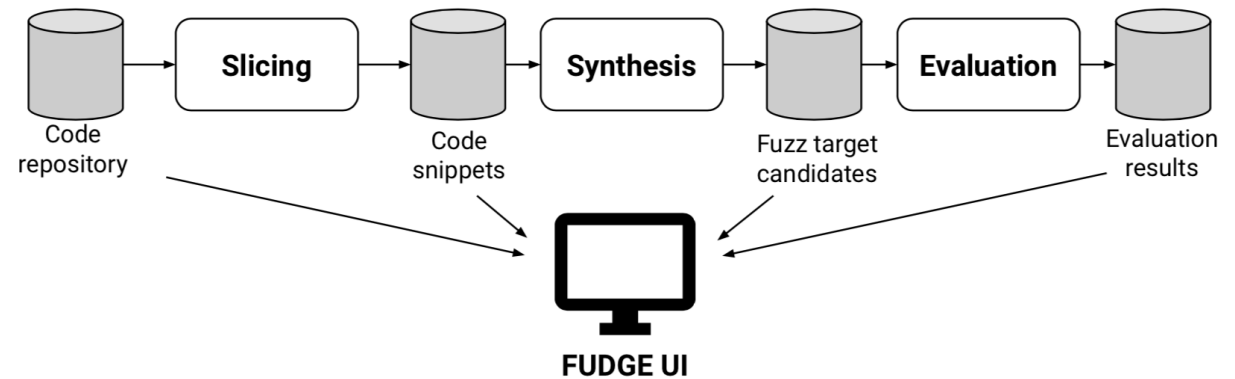# Enabling Fuzz-Driven Development

# Fuzz-Driven Development: Three Pillars

Automating Fuzzing
"Infrastructure"

Leveraging Easy

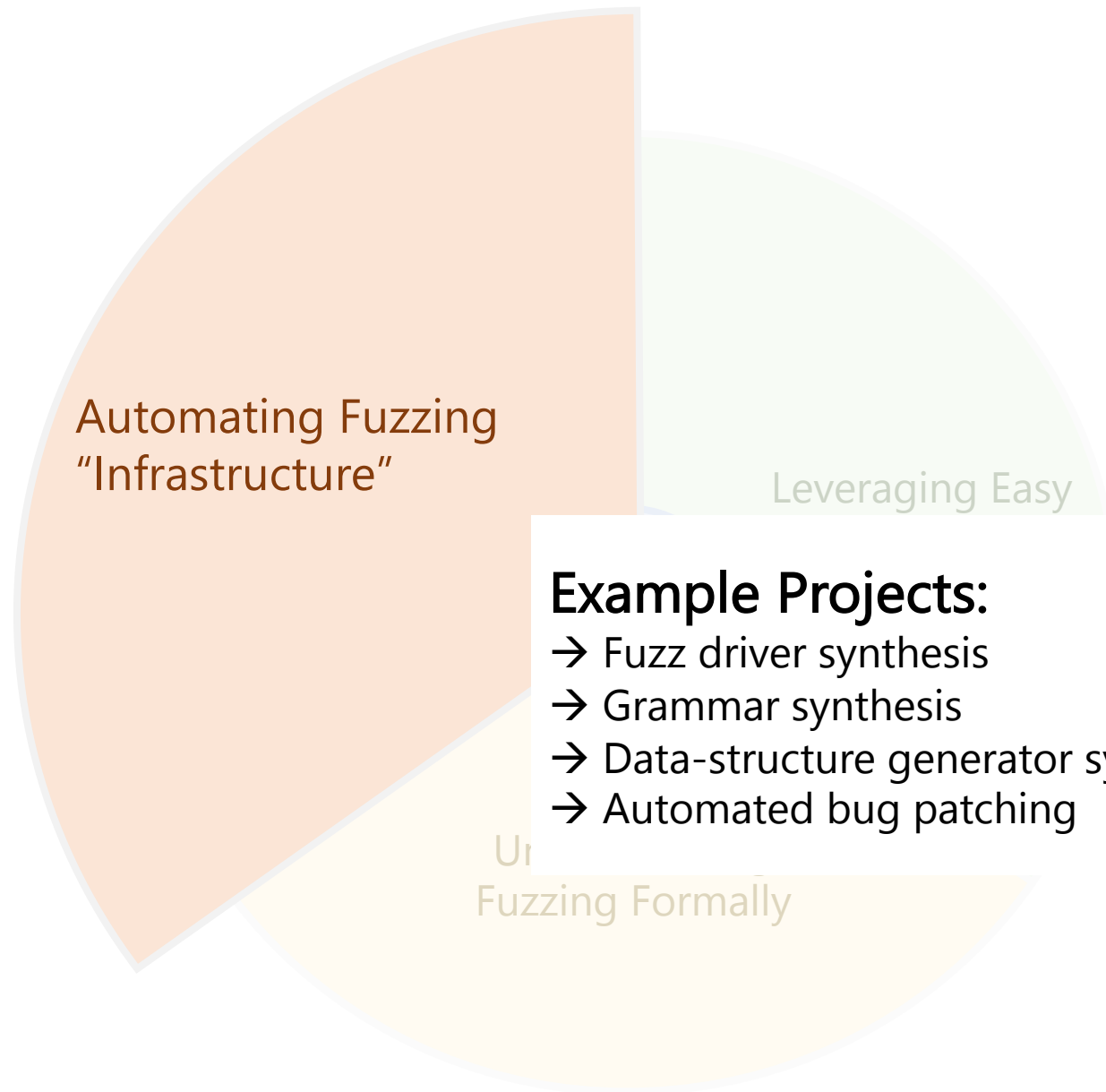Fuzzing Formally
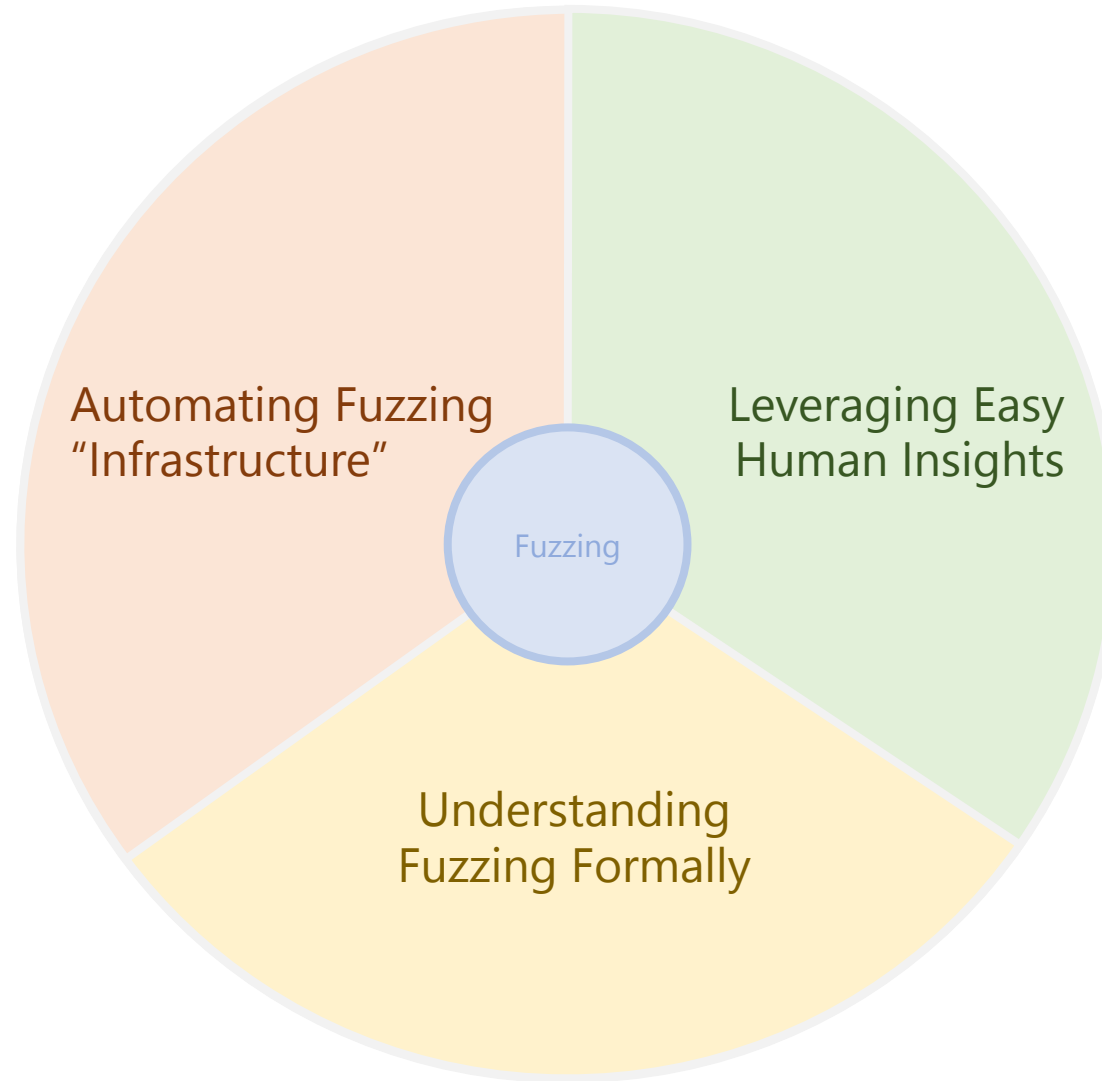
**Example Projects:**
→ Fuzz driver synthesis
→ Grammar synthesis
→ Data-structure generator synthesis
→ Automated bug patching

# FUDGE

Babic, Bucur, Chen, Ivancic, King, Kusano, <u>Lemieux</u>, Szekeres, Wang
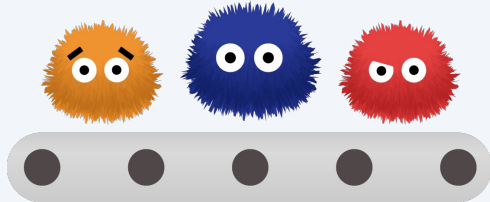ESEC/FSE'19 **(Industry Track)**

**Best Paper Award (Industry Track)**

Automating Fuzzing "Infrastructure"

Leveraging Easy



200 drivers integrated into open-source projects
→<u>150 security-improving fixes</u>

Automating Fuzzing
"Infrastructure"

Leveraging Easy

**Example Projects:**
→ Fuzz driver synthesis
→ Grammar synthesis
→ Data-structure generator synthesis
→ Automated bug patching

Fuzzing Formally

Caroline Lemieux — Expanding the Reach of Fuzz Testing

# PerfFuzz (ISSTA'18) + FuzzFactory (OOPSLA'19)

Using **generalized feedback maps** to expand **bugs findable by fuzz testing**

# CodaMOSA (ICSE'23)

Using **large language models** to improve **automated test suite generation**

Search stuck with generated test cases:

```
def test_case_1():
    str_0 = 'a\t!sUo~AU'
    str_1 = bump_version(
            str_0)
```

```
def test_case_2():
    str_0 = None
    int_0 = 1431
    str_1 = bump_version(
            str_0, int_0)
```

get test case hint as code

```
def test_bump_version():
    assert bump_version('0.0.0') == '1.0.0'
    assert bump_version('0.0.0', 1) == '0.1.0'
```

Automating Fuzzing "Infrastructure"

Leveraging Easy Human Insights

Fuzzing

Understanding Fuzzing Formally

clemieux@cs.ubc.ca     carolemieux.com     @cestlemieux