

# Perspectives on Structuring a Research Presentation

*or*

“One motto, three rules of thumb”

Claire Le Goues

SSSG, September 28, 2015



# Repairing Programs with Semantic Code Search

Yalin Ke      Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

# Repairing Software by Semantic Search

Mark T. Stolee  
School of Computer Science  
Iowa State University  
stolee@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yury Brun  
College of Information and Computer Sciences  
University of Massachusetts, Amherst  
brun@cs.umass.edu

Automated program repair can potentially reduce the cost of software development and improve software quality but recent studies have highlighted shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of open-source code to find potential fixes. The key challenge is efficiently finding code semantically similar (but not necessarily syntactically similar) to the buggy code and then appropriately integrating the found code into the program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program regions, and (3) deriving the desired input-output behavior for the localized fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior, (4) inserting the likely buggy code with these potential patches, and (5) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 780 defects written by novice students, 20 of which were repaired by GenProg, TrpAutoRepair, and AE. We compare the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass. We find that SearchRepair-repaired programs pass 97.3% of the tests on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source software movement has led to many large, publicly accessible source code repositories, such as GitHub, BitBucket, and SourceForge. These repositories contain programs that include routines, data structures, and other components previously implemented in other software. We posit that, if a method or component is used in other software, it is likely to be correct. If we find a defect, with high probability, there is a correct version of that component in some other software project. The research challenge here is to find that correct version.

Our key idea is to

existing open-source code to find correct implementations of buggy components and methods, and use the results to automatically generate patches for software defects. Semantic search identifies code by what it does, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) Encodes a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) Localizes a defect to likely buggy program fragments.
- 3) Constructs, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) Searches the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.

SearchRepair validates each potential patch against the program test suites to determine if it indeed repairs the defect in question. To make SearchRepair possible, we first extend our previous work on semantic code search [68] to C program fragments. We then use spectrum-based fault localization [36] to identify regions of faulty code and construct input-output profiles for those regions. Third, we build a database to perform semantic code search over the SMT-encoded database. We adapt the returned code fragment to the defect by performing variable renaming, and validate against provided test suites.

Our goal was to produce high quality patches while still supporting a wide range of defects. A key feature of a high quality patch is that it is either human- or tool-generated, is that it gets the desired, often unwritten specification of the defect right. This is a challenge for automatic repair [3], [7], [10], [11], [15], [16], [18], [19], [21], [22], [23], [29], [42], [48], [49], [50], [51], [52], [54], [56], [57], [71], [74], [76]), many of which use test suites to guide patching efforts. Modern test-suite guided patching (i.e., heuristically constructing and then testing candidate repairs), although typically general, can produce poor-quality patches that overfit to the test suites used to guide patch generation.

By definition, test suites only capture a subset of correct behavior. A patch that passes a given test suite may therefore be a false positive. This is a well-known phenomenon in the software engineering community, where the program



Motto:

**YOU ARE NOT PRESENTING THE PAPER.  
YOU ARE PRESENTING THE WORK.**

# Program repair via semantic search.

## Repairing Programs with Semantic Code Search

Yalin Ye    Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clgoues@cs.cmu.edu

Yury Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar that not identical to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

### 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects also in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use semantic code search [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to automatically generate patches for software defects. Semantic search identifies code by what it does, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adopt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it generalizes to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [17], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a generate-and-validate paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program



# Program repair via semantic code search

## Repairing Programs with Semantic Code Search

Yalin Ke  
Department of Computer Science  
Iowa State University  
{yke, kotolee}@iastate.edu

Kathryn  
Department of Computer Science  
Iowa State University  
kotolee@iastate.edu

Yury Brun  
Department of Informatics and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but has not yet been widely adopted due to shortcomings in the quality of generated repairs. We propose a new kind of repair technique that leverages the large body of existing open-source code to find potential patches. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TpykatoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TpykatoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TpykatoRepair, and AE, and requires some defects these tools cannot.

### 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects also in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use semantic code search [68] over

open-source code to find correct implementations of buggy code fragments and methods, and use the results to generate patches for software defects. Semantic code search is different from traditional code search in that it does, rather than by syntax, find code that is semantically similar to the buggy code. SearchRepair, a new technique for automated program repair, leverages this idea.

- 1) Encodes a large database of human-written code fragments as SMT constraints on their input-output behavior.
- 2) Localizes a given defect to likely buggy program fragments and derives the desired input-output behavior for code to replace those fragments.
- 3) Uses state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches.
- 4) Validates that the patches repair the bug against program test suites.
- 5) Validates each potential patch against the program test suite to determine if it indeed fixes the defect.

To make SearchRepair possible, we develop a new technique for work in semantic code search [68] to find code that is semantically similar to the defective code. Second, we adopt spectrum-based fault localization [10] to identify candidate regions of faulty code and derive the desired input-output profiles to use as input to semantic search. Finally, we build the infrastructure to perform semantic code search over a large SMT-encoded code database, adapt the returned code to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it generalizes to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a generate-and-validate paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of overfitting to an objective function, where the program



# Program repair via semantic search.

## Repairing Programs with Semantic Code Search

Yulin Ke      Kathryn T. Stolee      Claire Le Goues      Yary Brun  
Department of Computer Science      School of Computer Science      College of Information and Computer Science  
Iowa State University      Carnegie Mellon University      University of Massachusetts, Amherst  
{yke, kstolee}@iastate.edu      cle goues@cs.cmu.edu      brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects these tools cannot.

### 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship immature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [52].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, Bitficket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it does, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior;
- 2) *Localizes* a defect to likely buggy program fragments;
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints;
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code;
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [56] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [17], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

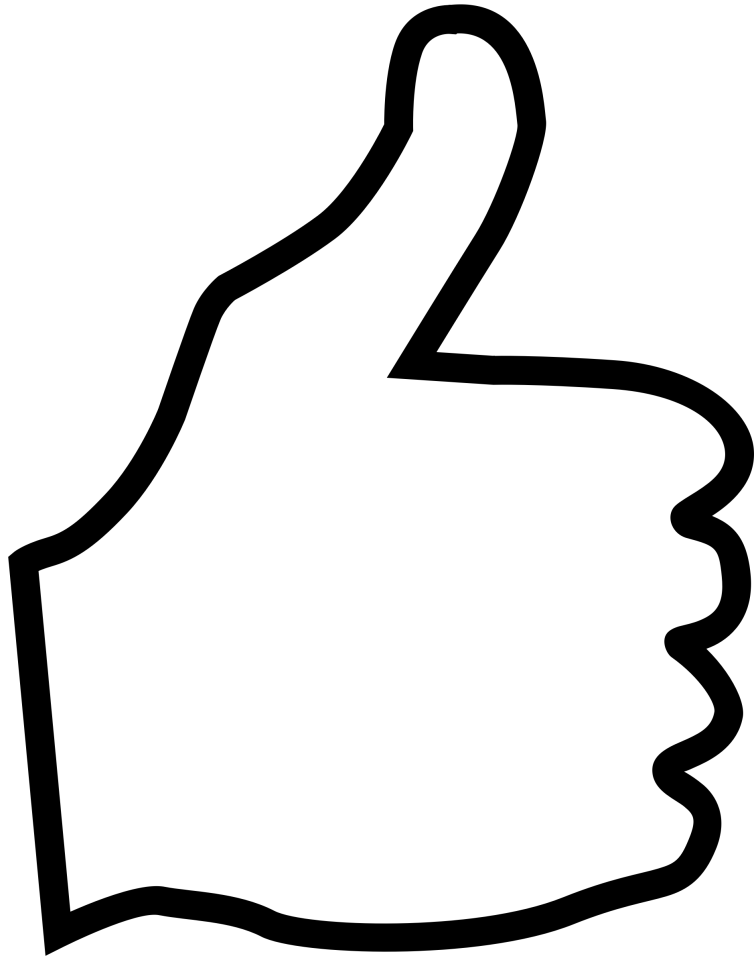






Motto:

**YOU ARE NOT PRESENTING THE PAPER.  
YOU ARE PRESENTING THE WORK.**



- **The audience will only remember 3 things.**
- Tell a story.
- Never confuse your listeners.

Average  
conference  
attendee.





The \$\$\$\$\$\$\$\$ question:

**WHAT SHOULD THOSE THREE  
THINGS BE?**



(Average audience member.)

## CLG's Goal

1. The *exciting and important* problem I am solving.
2. The *key nugget of awesomeness* underlying the approach.
3. 1—2 major result(s).
4. “That paper/person seems cool, I want to read it/talk to her!”

# Repairing Programs with Semantic Code Search

Yalin Ke    Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

1. The *exciting and important* problem I am solving.
2. The *key nugget of awesomeness* underlying the approach.
3. 1—2 major result(s).
4. “That paper/person seems cool, I want to read it/talk to her!”

# Repairing Programs with Semantic Code Search

Yalin Ke      Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

1. Automatic defect repair that produces high-quality patches.
2. The key nugget of awesomeness underlying the approach.
3. 1—2 major result(s).
4. “That paper/person seems cool, I want to read it/talk to her!”



# Repairing Programs with Semantic Code Search

Yalin Ke      Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

1. Automatic defect repair that produces high-quality patches.
2. SMT-based semantic search, which looks for code based on what it should *do*.
3. 1—2 major result(s).
4. “That paper/person seems cool, I want to read it/talk to her!”

# Repairing Programs with Semantic Code Search

Yalin Ke    Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

1. Automatic defect repair that produces high-quality patches.
2. SMT-based semantic search, which looks for code based on what it should *do*.
3. Empirically wins.
4. “That paper/person seems cool, I want to read it/talk to her!”

# Repairing Programs with Semantic Code Search

Yalin Ke    Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuriy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

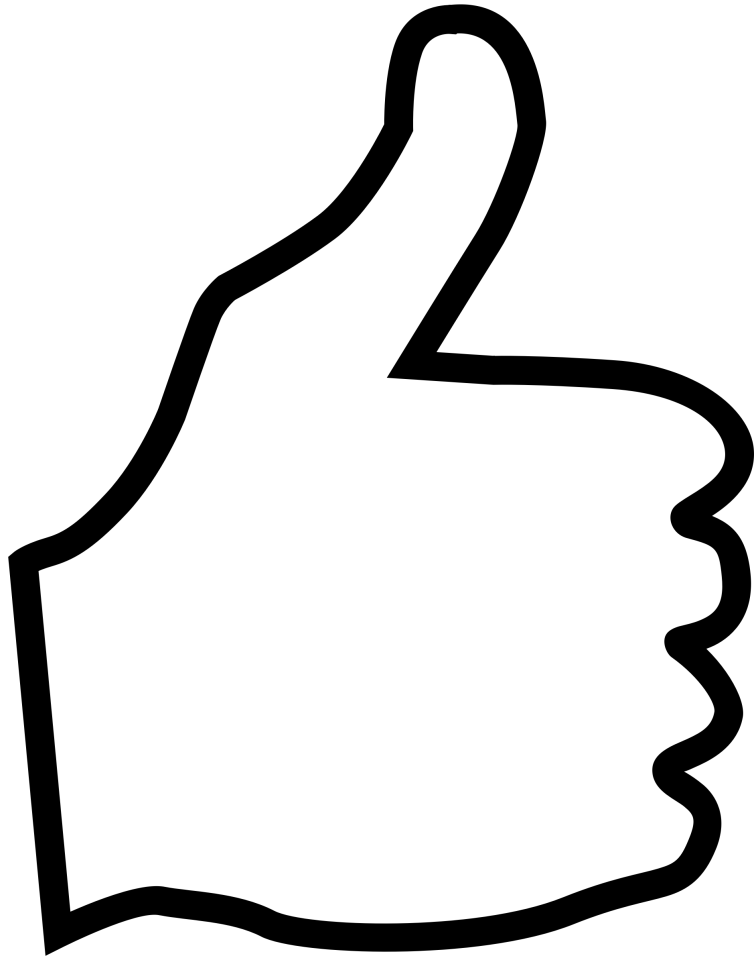
- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

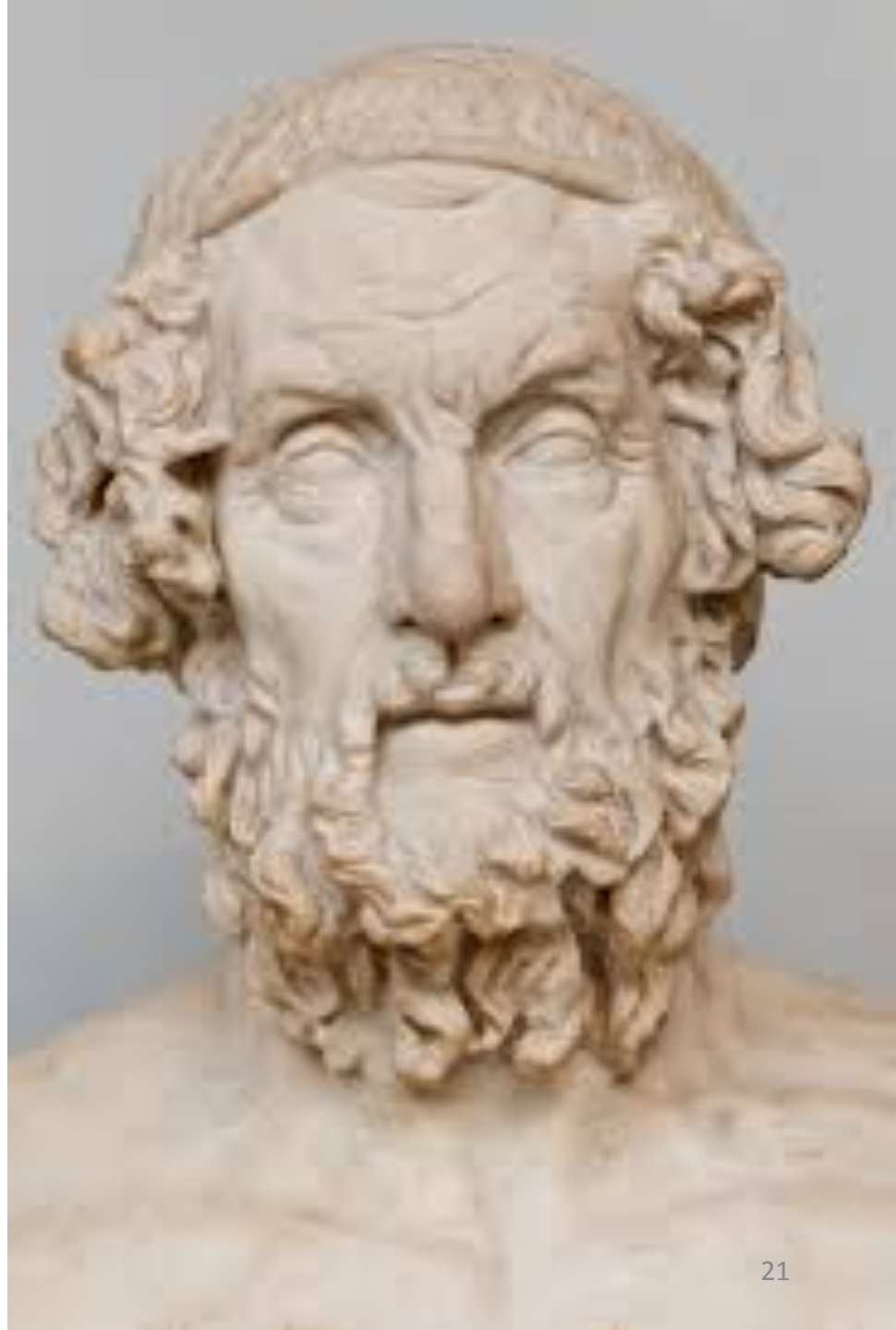
1. Automatic defect repair that produces high-quality patches.
2. SMT-based semantic search, which looks for code based on what it should *do*.
3. Empirically wins.
4. “Look, there she is, hey, let’s go talk to her!”

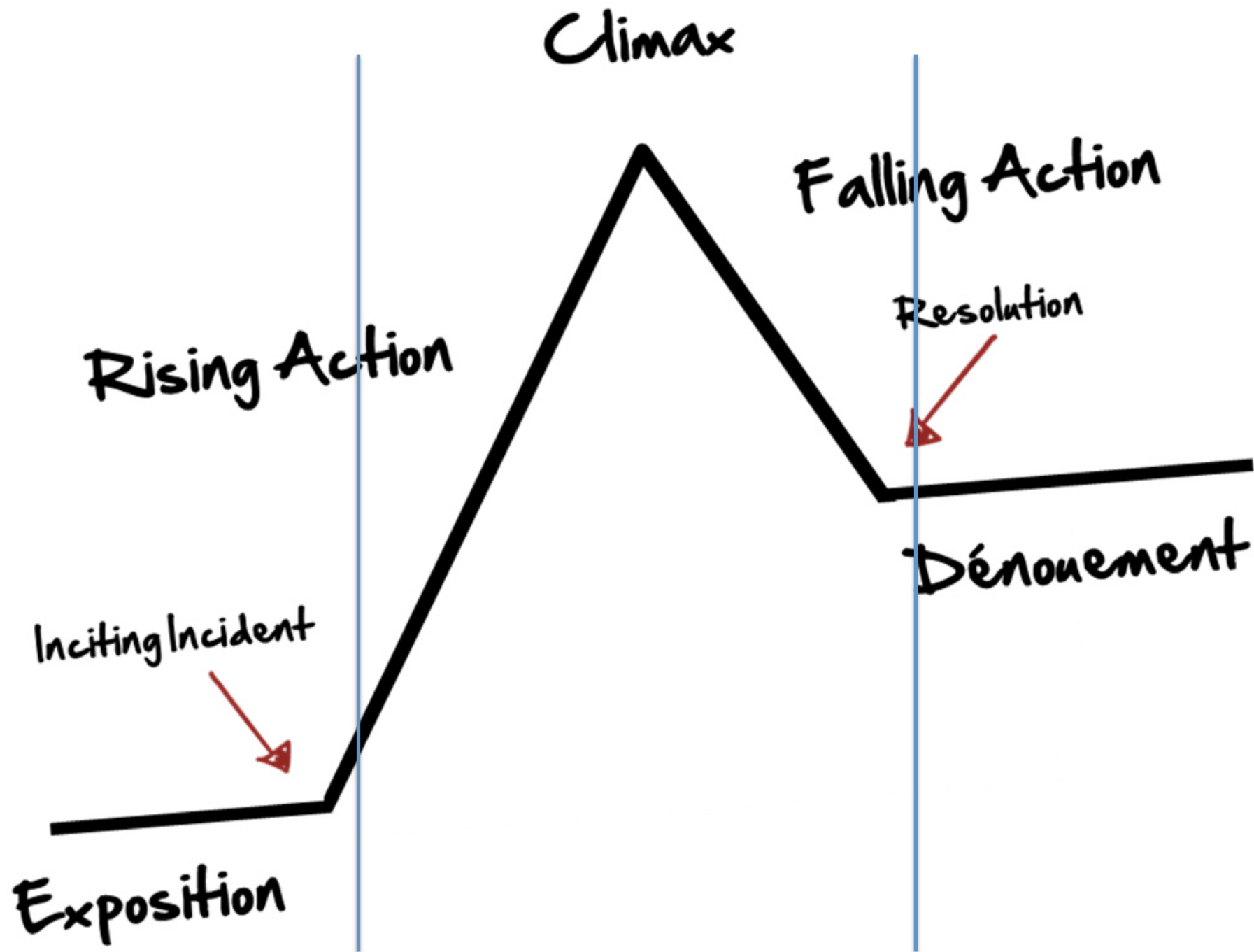


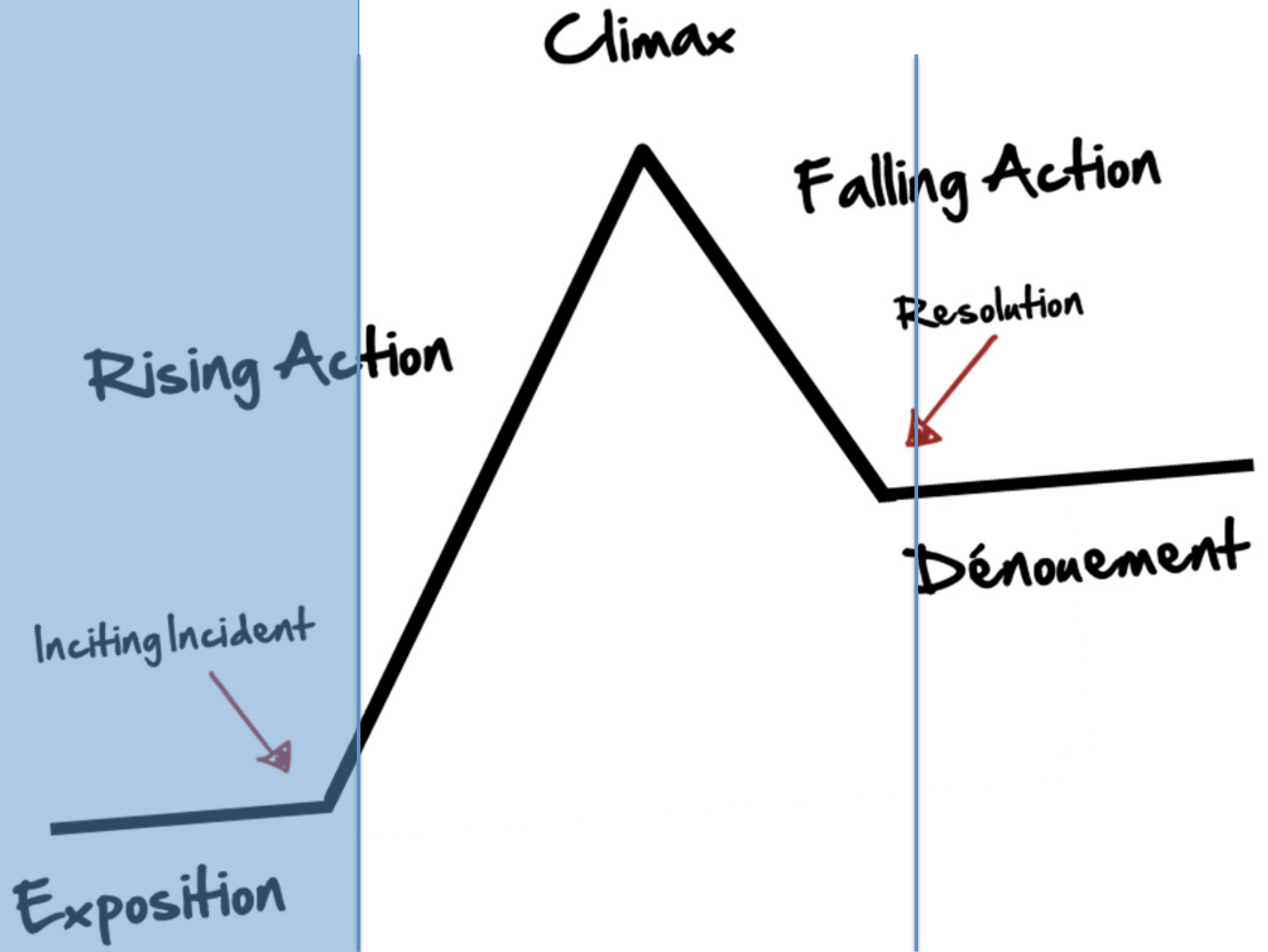
- Your audience will only remember 3 things.
- **Tell a story.**
- Never confuse your listeners.

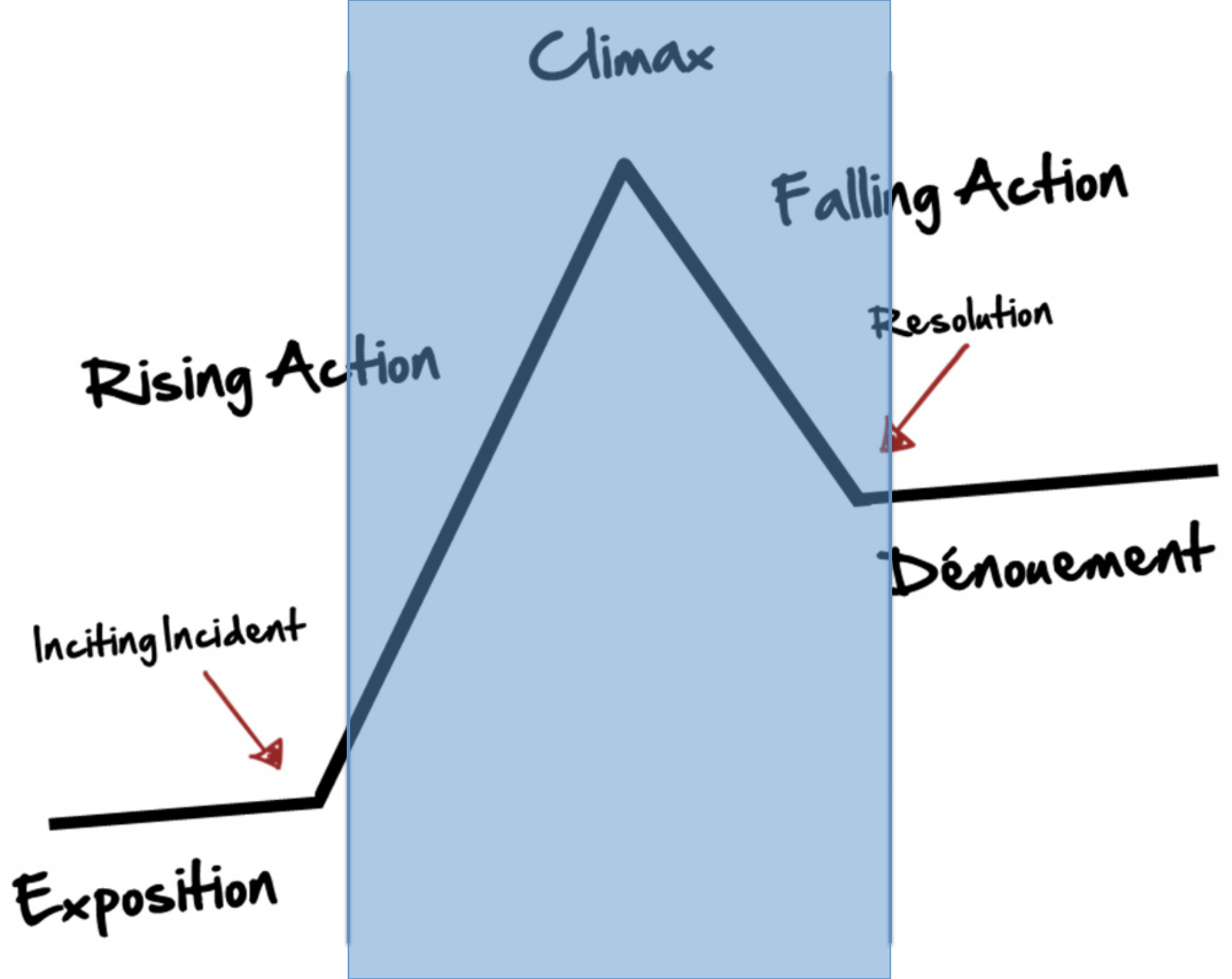


- *Interesting.*
- Simple, and not overly detailed.
- Selectively repetitive.
- Coherent narrative arc.

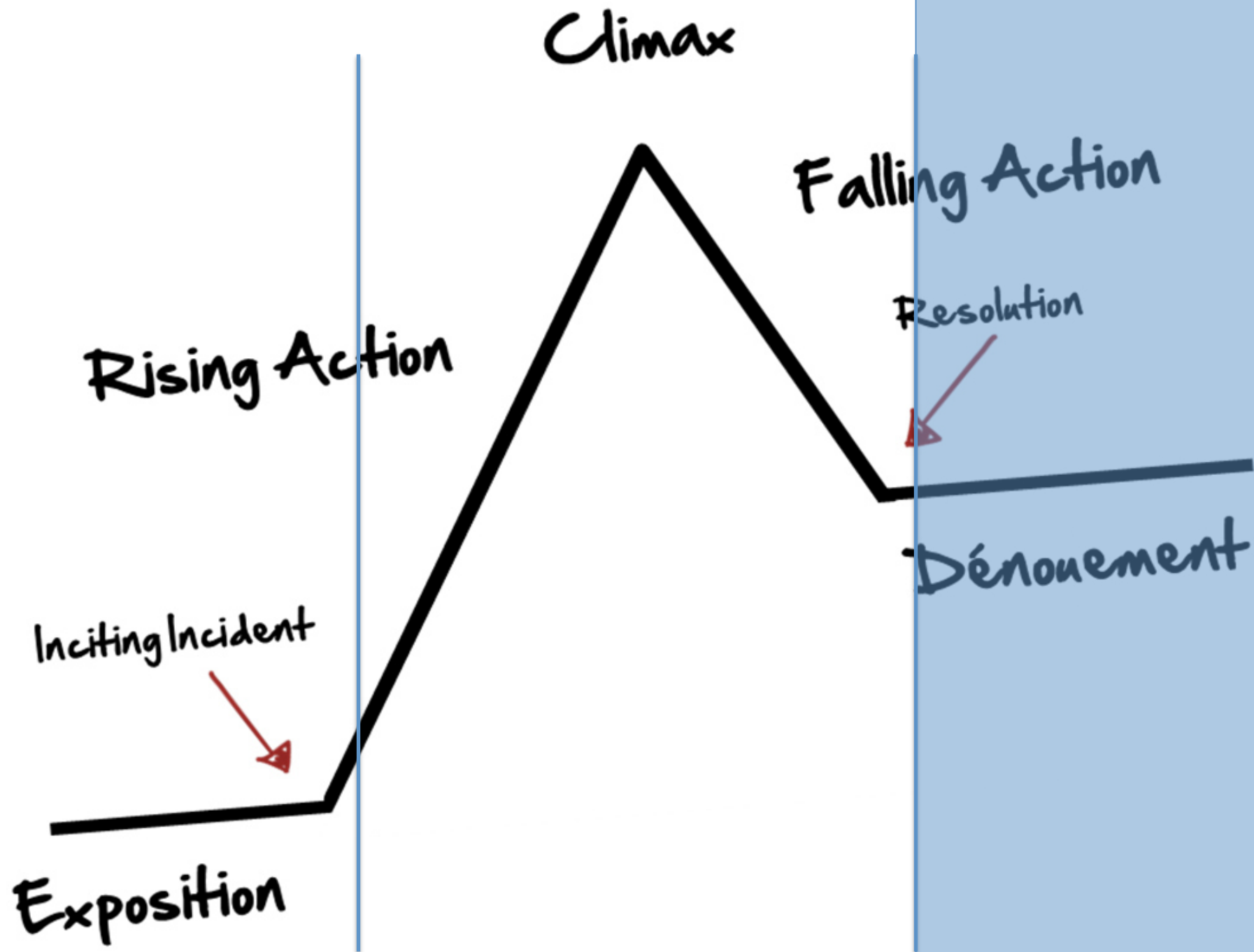


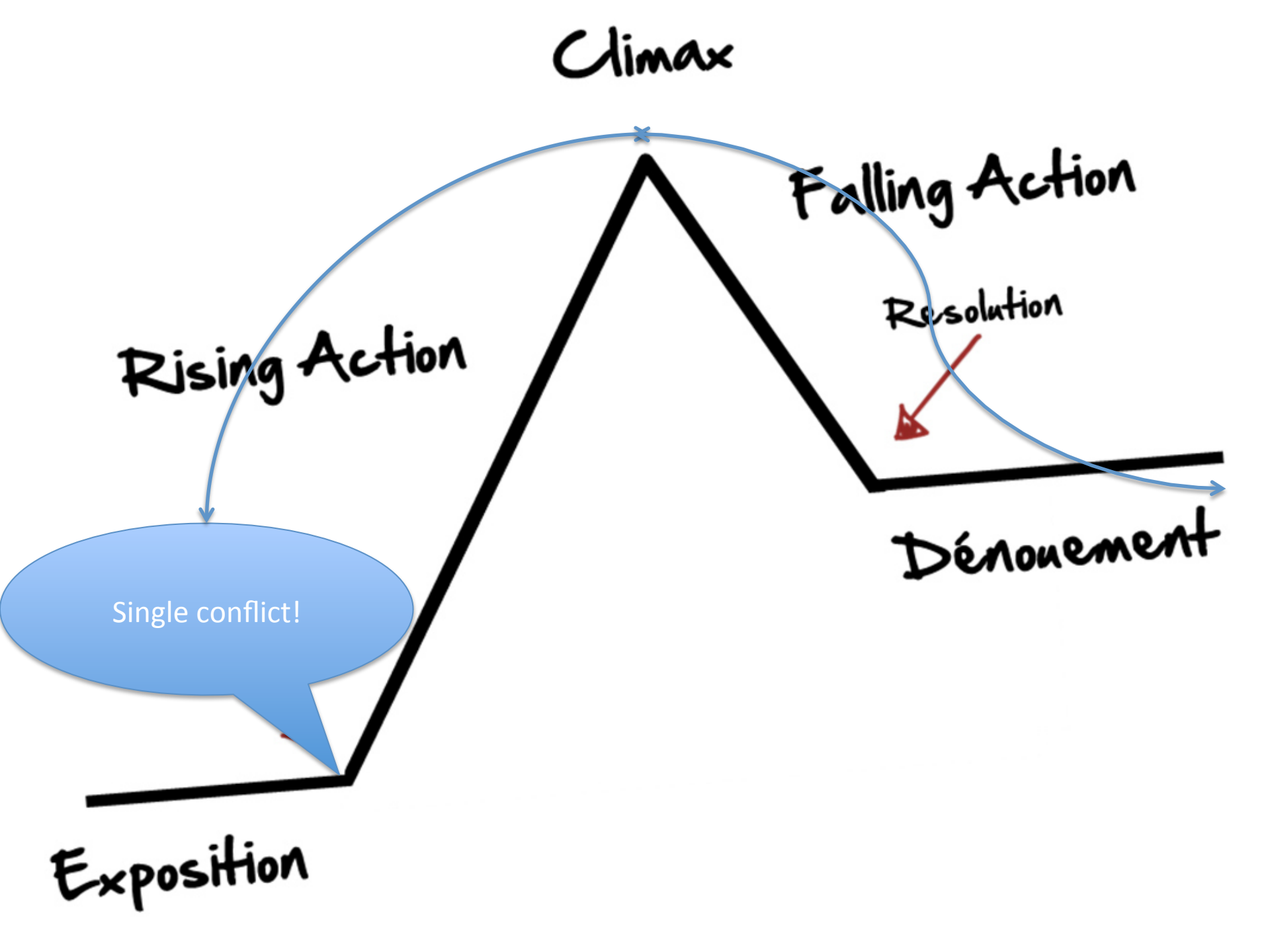


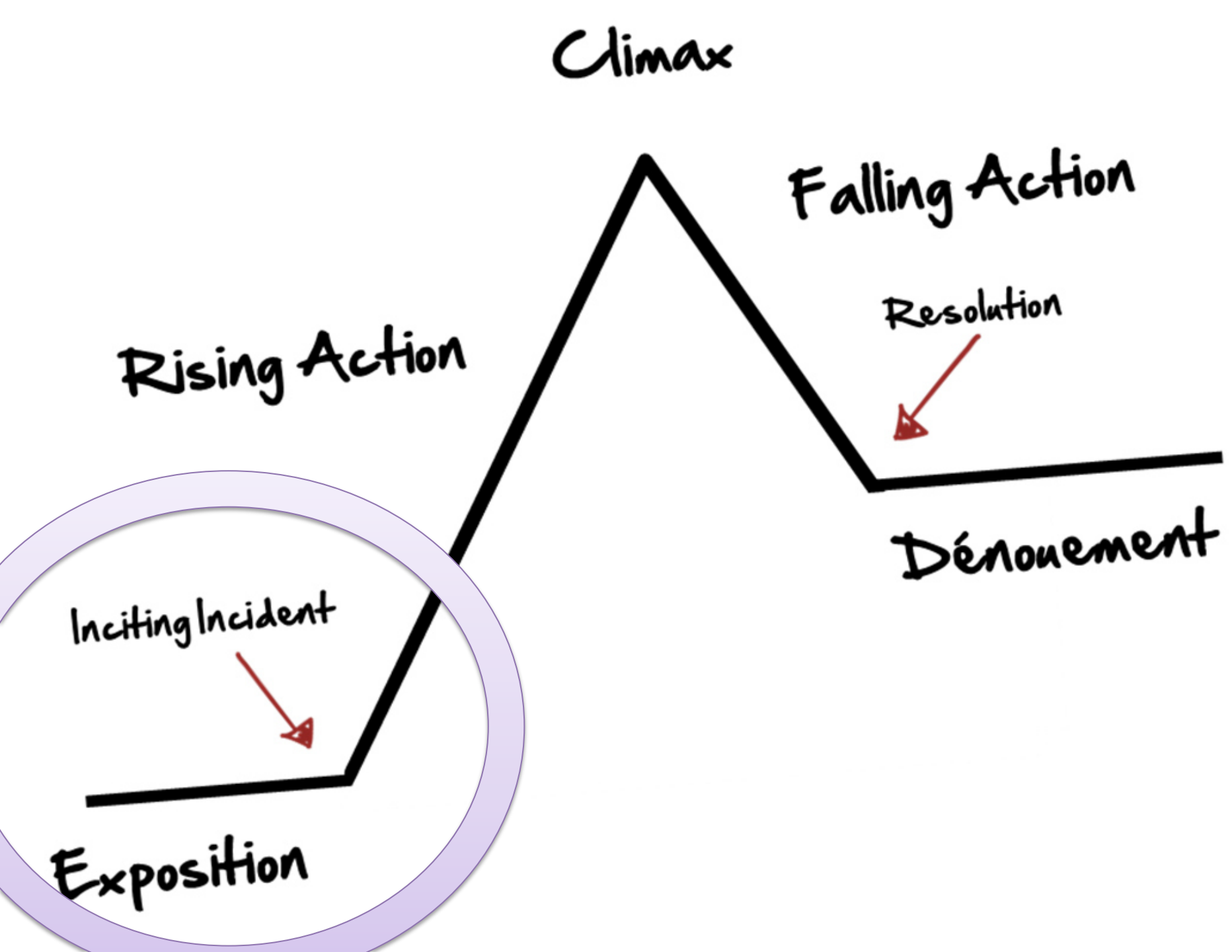












# Exposition, conflict

- *Important, interesting* problem that I am solving.
- Show, don't tell: motivating example, story, easy-to-grasp soundbites.
  - Sometimes a reasonable place to *delicately* mention related or previous work.
- Will guide/motivate the subsequent events of the story; focus on *one type* of motivation.



# Repairing Programs with Semantic Code Search

Yalin Ke      Kathryn T. Stolee  
Department of Computer Science  
Iowa State University  
{yke, kstolee}@iastate.edu

Claire Le Goues  
School of Computer Science  
Carnegie Mellon University  
clegoues@cs.cmu.edu

Yuryy Brun  
College of Information and Computer Science  
University of Massachusetts, Amherst  
brun@cs.umass.edu

**Abstract**—Automated program repair can potentially reduce debugging costs and improve software quality but recent studies have drawn attention to shortcomings in the quality of automatically generated repairs. We propose a new kind of repair that uses the large body of existing open-source code to find potential fixes. The key challenges lie in efficiently finding code semantically similar (but not identical) to defective code and then appropriately integrating that code into a buggy program. We present SearchRepair, a repair technique that addresses these challenges by (1) encoding a large database of human-written code fragments as SMT constraints on input-output behavior, (2) localizing a given defect to likely buggy program fragments and deriving the desired input-output behavior for code to replace those fragments, (3) using state-of-the-art constraint solvers to search the database for fragments that satisfy that desired behavior and replacing the likely buggy code with these potential patches, and (4) validating that the patches repair the bug against program test suites. We find that SearchRepair repairs 150 (19%) of 778 benchmark C defects written by novice students, 20 of which are not repaired by GenProg, TrpAutoRepair, and AE. We compare the quality of the patches generated by the four techniques by measuring how many independent, not-used-during-repair tests they pass, and find that SearchRepair-repaired programs pass 97.3% of the tests, on average, whereas GenProg, TrpAutoRepair, and AE-repaired programs pass 68.7%, 72.1%, and 64.2% of the tests, respectively. We conclude that SearchRepair produces higher-quality repairs than GenProg, TrpAutoRepair, and AE, and repairs some defects those tools cannot.

## 1. INTRODUCTION

Buggy software costs the global economy billions of dollars annually [8], [60]. One major reason software defects are so expensive is that software companies must dedicate considerable developer time [75] to manually finding and fixing bugs in their software. Unfortunately, manual bug repair, the industry standard, is largely unable to keep up with the volume of defects in extant software [2]. Despite their established detrimental impact on a company's bottom line, known defects ship in mature software projects [45], and many defects, including those that are security-critical, remain unaddressed for long periods of time [32].

At the same time, the expansion of the open-source movement has led to many large, publicly accessible source code databases, such as GitHub, BitBucket, and SourceForge. Because many programs include routines, data structures, and designs that have been previously implemented in other software projects [11], [12], [24], we posit that, if a method or component of a software system contains a defect, with high probability, there exists a similar but correct version of that component in some publicly accessible software project. The research challenge lies in how to automatically find and use such implementations to repair bugs.

Our key idea is to use *semantic code search* [68] over

existing open-source code to find correct implementations of buggy components and methods, and use the results to *automatically generate patches* for software defects. Semantic search identifies code by what it *does*, rather than by syntactic keywords. We develop SearchRepair, a new technique predicated on our idea. SearchRepair:

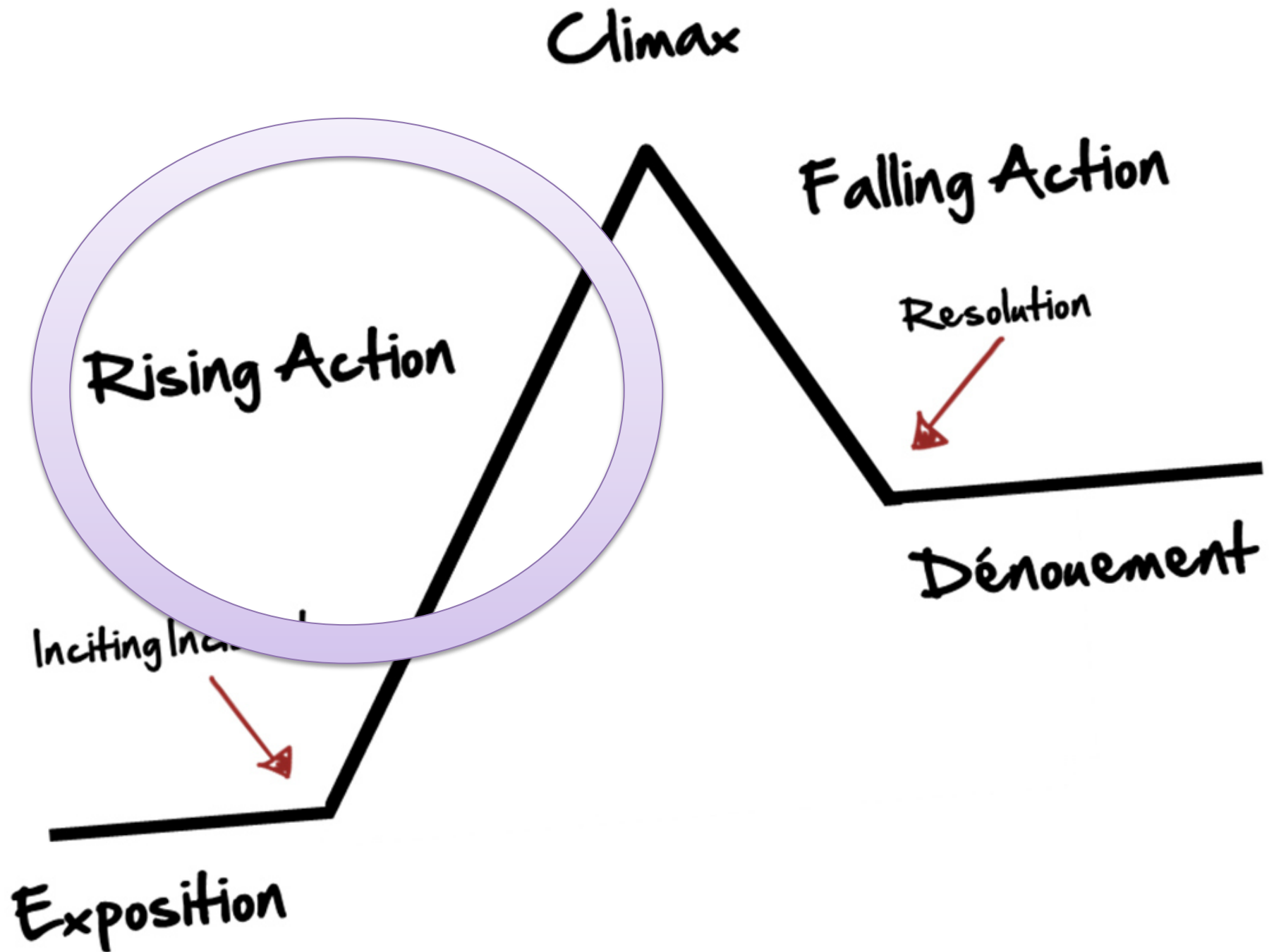
- 1) *Encodes* a large database of human-written code fragments as satisfiability modulo theories (SMT) constraints on their input-output behavior.
- 2) *Localizes* a defect to likely buggy program fragments.
- 3) *Constructs*, for each fragment, a lightweight input-output profile that characterizes desired functional behavior as SMT constraints.
- 4) *Searches* the database, using state-of-the-art constraint solvers, for fragments that satisfy such a profile. These fragments become potential patches when contextualized and inserted into the buggy regions, replacing the original potentially faulty code.
- 5) *Validates* each potential patch against the program test suite to determine if it indeed repairs the defect in question.

To make SearchRepair possible, we first extend our previous work in semantic code search [68] to C program fragments. Second, we adapt spectrum-based fault localization [36] to identify candidate regions of faulty code and construct input-output profiles to use as input to semantic search. Third, we build the infrastructure to perform semantic code search over the SMT-encoded code database, adapt the returned code fragment to the defective context via variable renaming, and validate against provided test suites.

Our goal with SearchRepair is to produce high quality patches while still addressing a broad range of defects. A key feature of a high quality patch, whether human- or tool-generated, is that it *generalizes* to the full, desired, often unwritten specification of correct program behavior. This is a challenge for automatic repair techniques (e.g., [3], [7], [10], [11], [15], [16], [18], [19], [21], [28], [33], [35], [39], [42], [48], [49], [50], [51], [52], [54], [56], [57], [61], [69], [70], [73], [74], [76]), many of which use test suites to guide and evaluate patching efforts. Modern test-suite guided repair techniques, particularly those following a *generate-and-validate* paradigm (i.e., heuristically constructing and then testing large numbers of candidate repairs), although typically general and scalable, often produce poor-quality patches that overfit to the specification test suites used to guide patch generation [20], [57], [65].

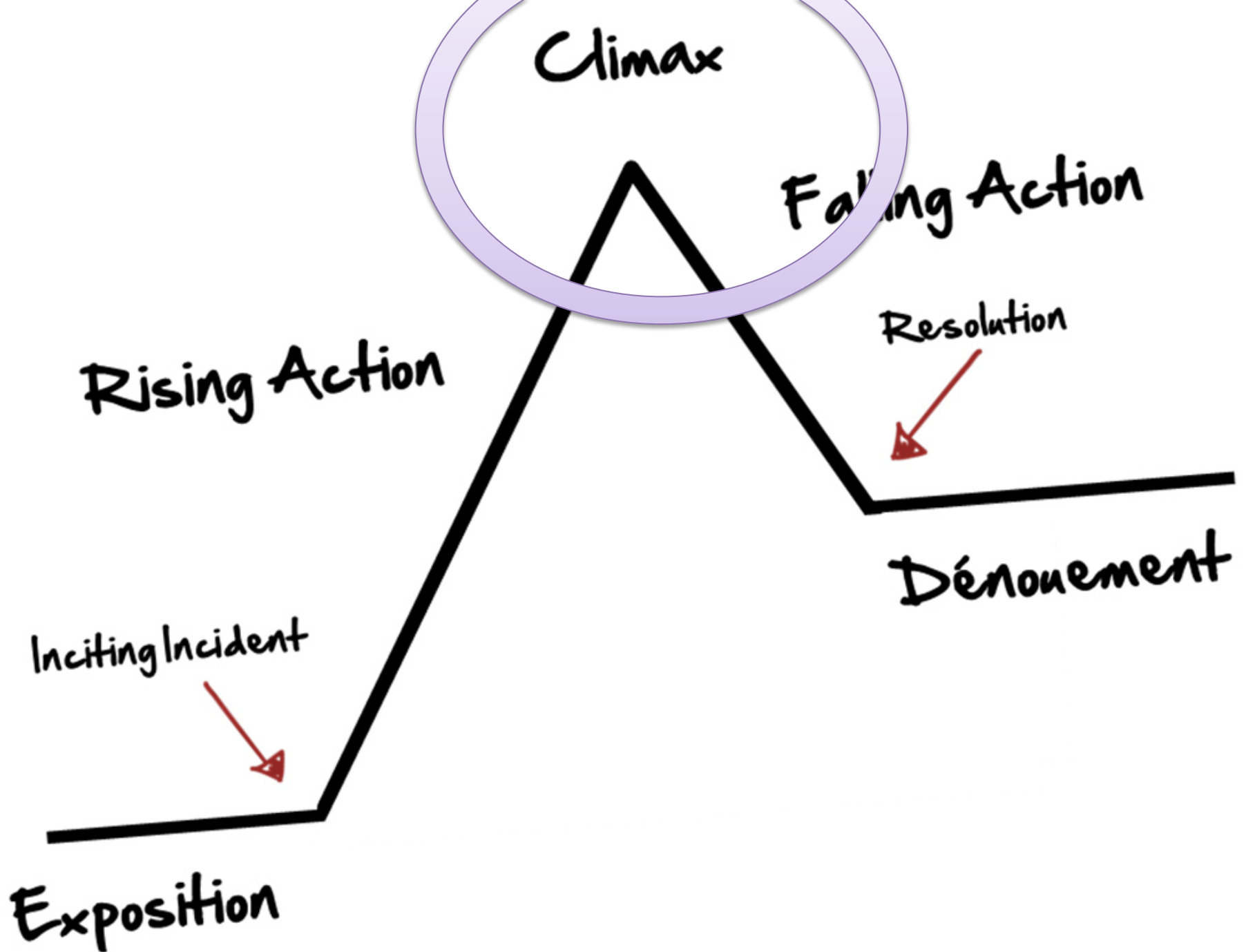
By definition, test suites only encode a partial specification of correct behavior. A patch that is correct according to a given test suite may therefore not be fully correct when evaluated with respect to a hypothetical full correctness specification. This is analogous to the well-known machine learning phenomenon of *overfitting* to an objective function, where the program

- Options:
  - Example showing how semantic search works.
  - Example walking through hypothetical program repair/semantic search combo use case.
  - Compelling results highlighting “quality problem” in previous results.
- *But I will only choose one of them.*



# Inciting incident + rising action

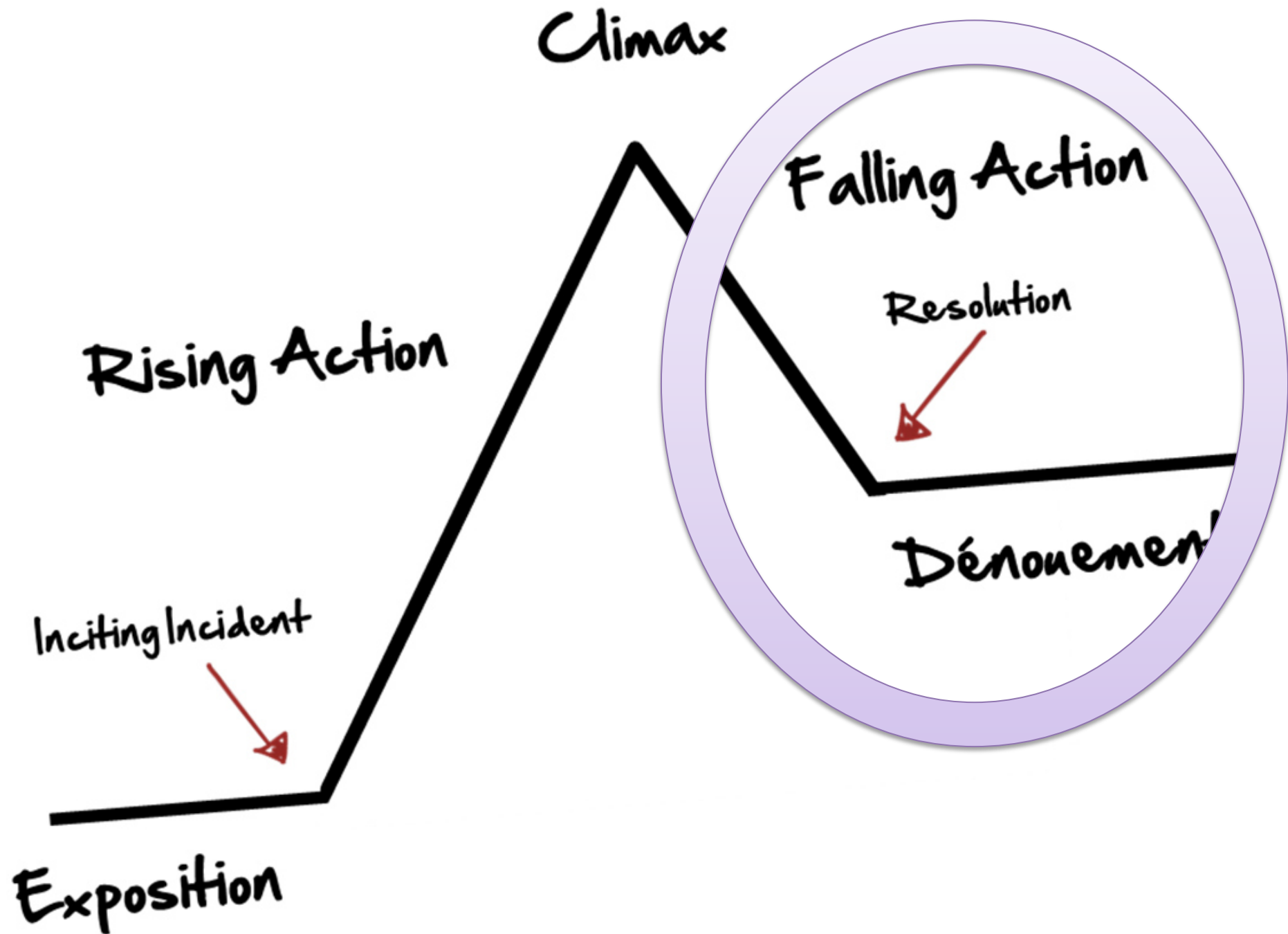
- Middle of story: Key technical insight.
- High-level outline of the approach.
- If you're not sure if a detail is high-level enough, it's probably not.
- (More on how to approach “the middle” in 10 —11 slides)





# Climax

- Results presentation: *selected* key results.
- Emphasis on the type of experimental methodology used, experimental question(s)
  - Calling back to your exposition!
  - Remember: oral communication is often cyclic/ repetitive.



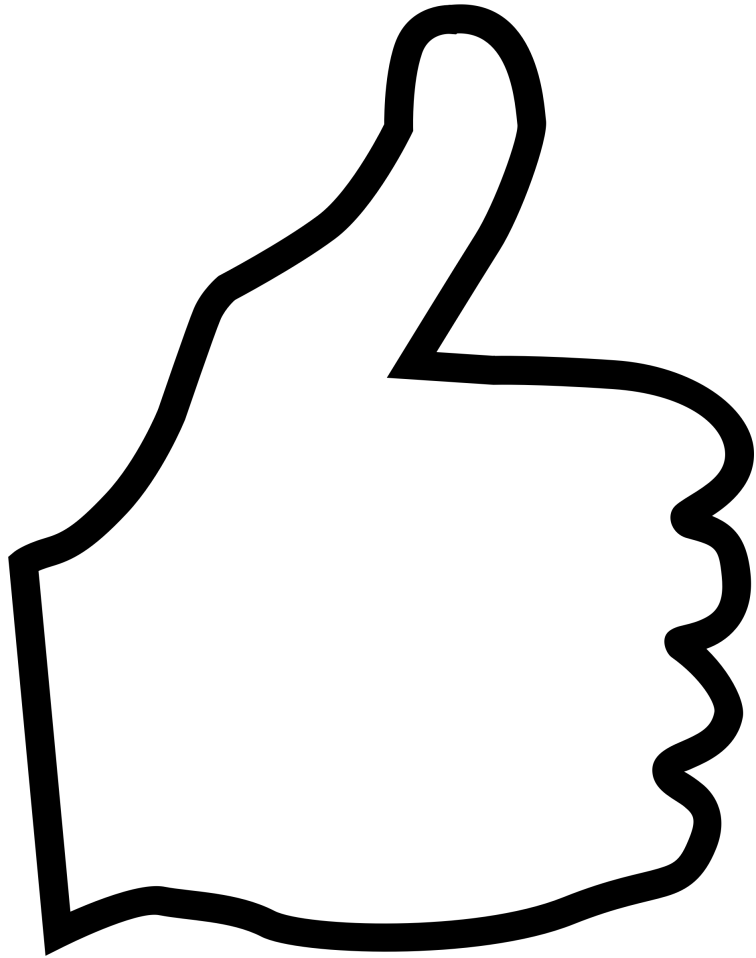
# Falling action

- Observations/implications.
- Another place related work *might* make an appearance.
- Possibly, future work.
  - (CLG thinks this is pointless, but acknowledges the existence/validity of opposing viewpoints.)

# Conclusion/Denouement

- “Say what you’re going to say, then say it, then **say what you said.**”
- Wrap up pithily.
- Remind me of the three things you want me to remember.





- Your audience will only remember 3 things.
- Tell a story.
- **Never confuse your listeners.**





When confused, readers can:

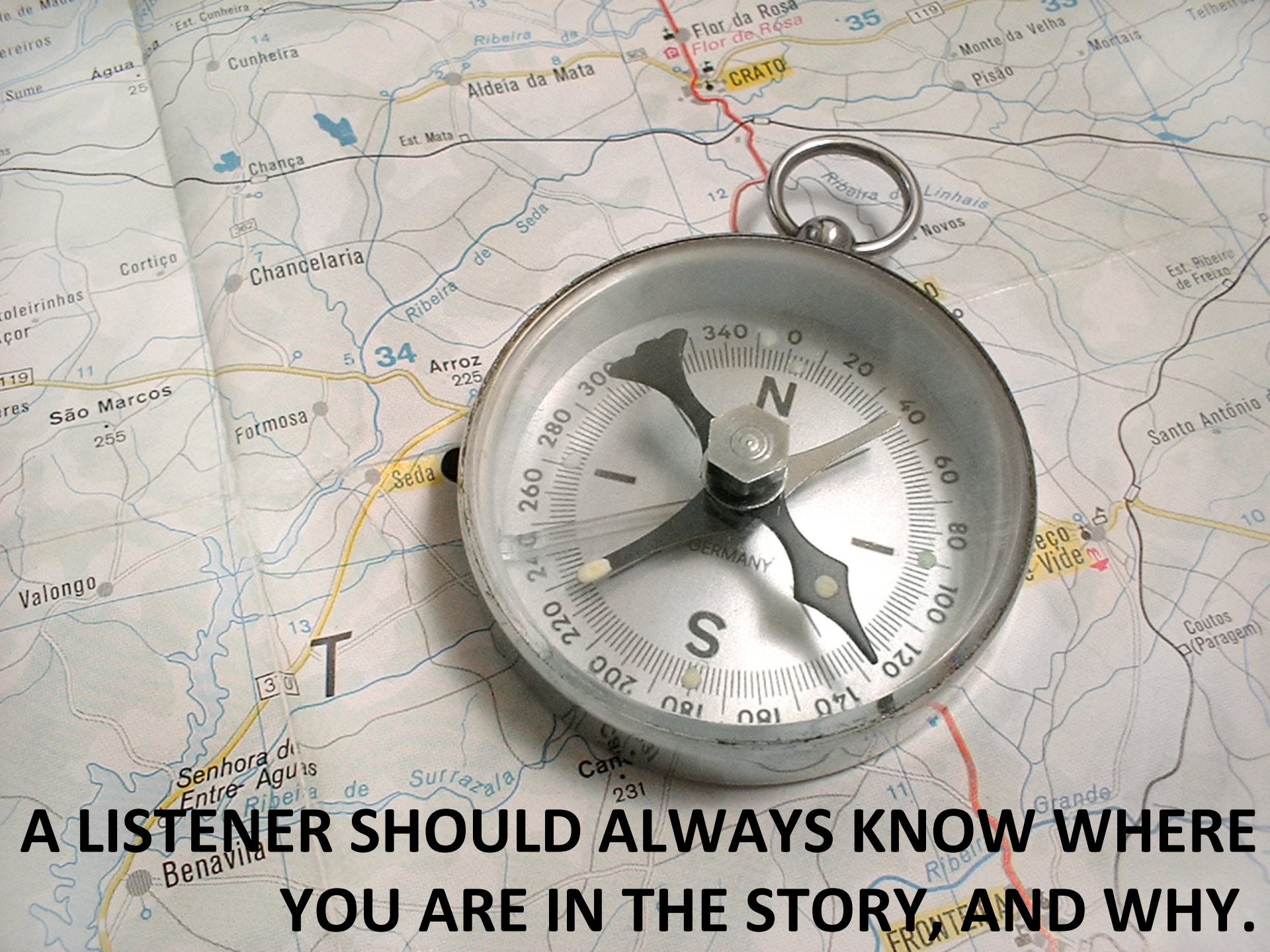
- Pause, reflect.
- Reread confusing passage.
- Go back to review a previous section.
- Look something up.

When confused, listeners can:

- Possibly interrupt to ask a question.
- Furtively look it up.
- Give up and start reading their email.







**A LISTENER SHOULD ALWAYS KNOW WHERE  
YOU ARE IN THE STORY, AND WHY.**

# Why?

- Listener has to synthesize what you're saying into the story.
- If she doesn't know why you're telling her something, she won't know where to "put" a piece of information in the overall picture.
- Result: listener is anxious, and likely to forget key pieces of information before they're needed!



# Implicit

I've done this several times already.

- Signpost as you go.
  - “This is important, because...”
  - Return to your outline slide, if you’re using it.
- Only introduce necessary information, and only when it is necessary.
- *Strongly* avoid forward references.
- *Strictly* avoid use-before-defined.

I violated this rule 11 slides ago.

**DO NOT VISUALLY OVERWHELM  
YOUR LISTENERS.**

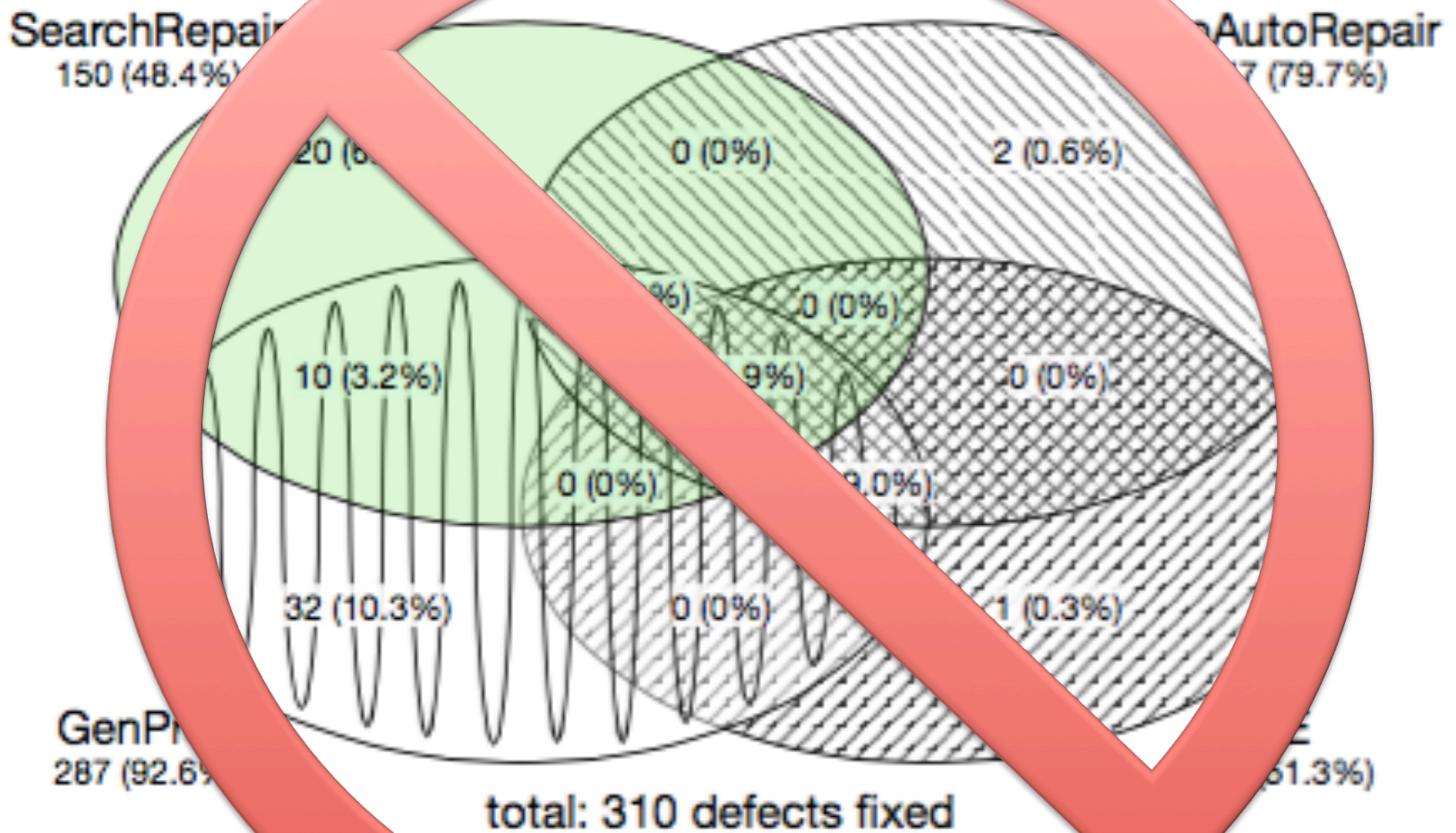


**DO NOT VISUALLY OVERWHELM  
YOUR LISTENERS.**

**DO NOT VISUALLY OVERWHELM  
YOUR LISTENERS.**

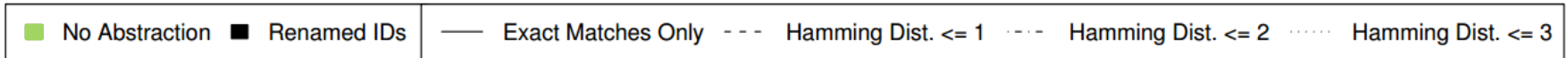
**BEWARE THE RESULTS  
PRESENTATION.**



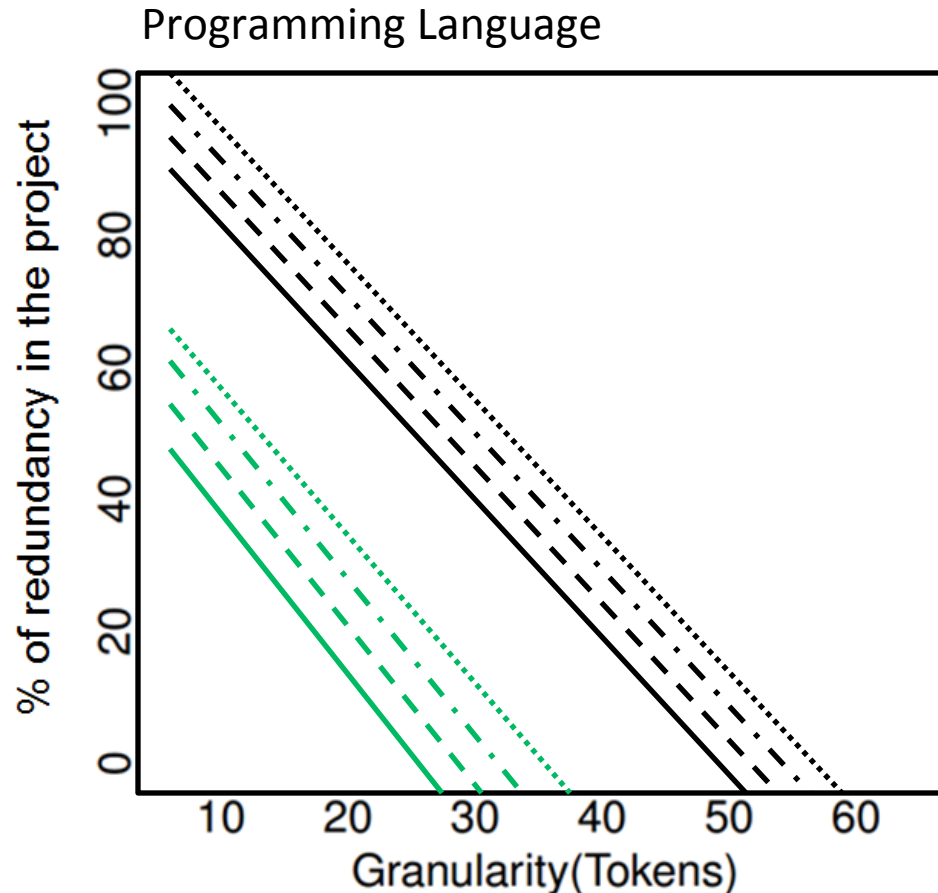


(Slide borrowed from my student Mauricio.)

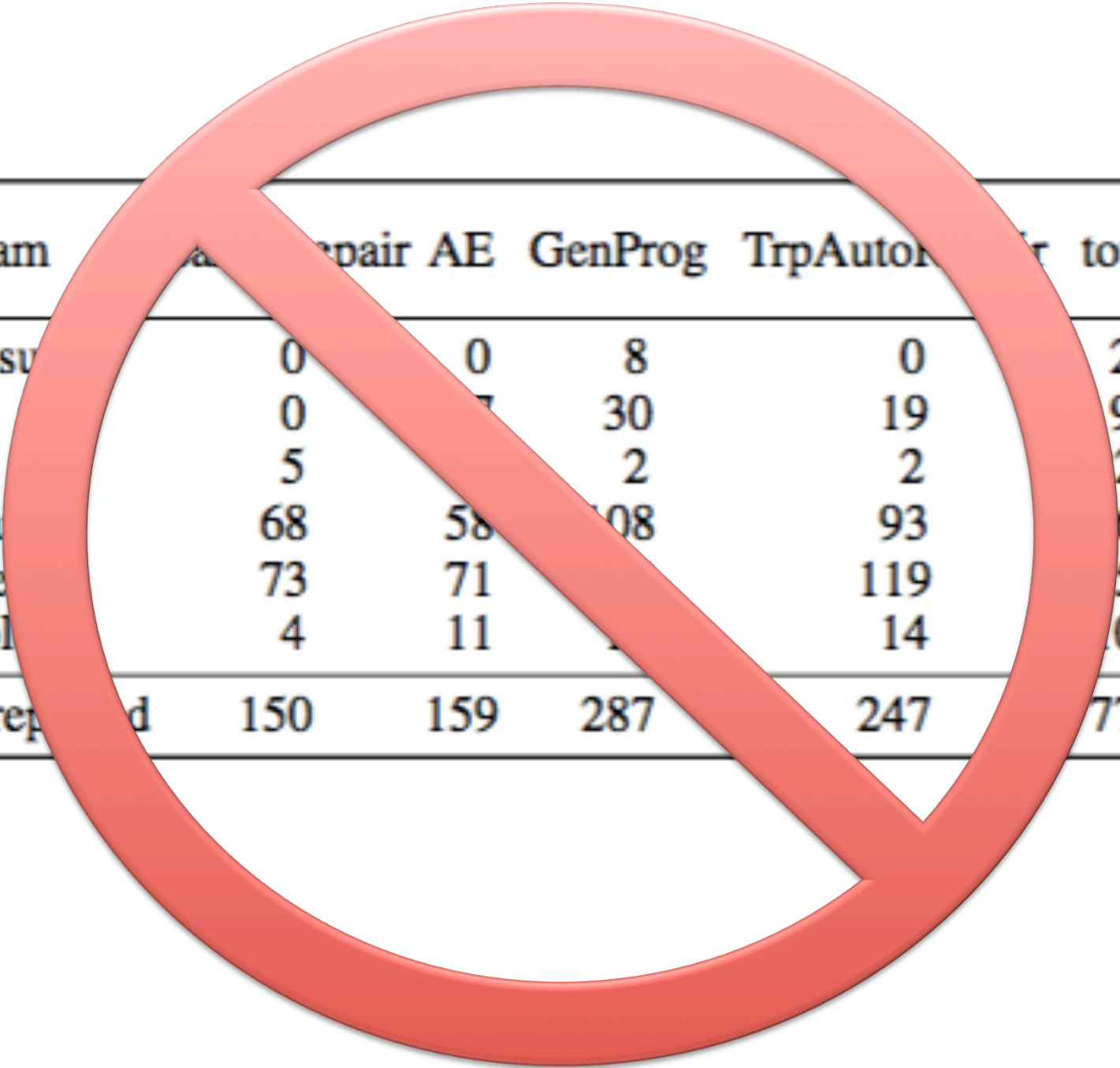
# How to read the graph?



What % of this project can be reconstructed from the corpus?



```
newNumber = oldNumber1 + oldNumber2 ;
```



program	repair	AE	GenProg	TrpAuto	total
checks	0	0	8	0	29
digits	0	7	30	19	91
grade	5		2	2	26
media	68	58	108	93	68
smalle	73	71		119	55
syllabl	4	11		14	109
total rep	150	159	287	247	778

Program	Description	LOC	Bug Type	Time (s)
gcd	example	153	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	unitary	1169	segmentation fault	45
look-s	unitary loop	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	74	segmentation fault	230
openldap	directory protocol	29	non-overflow denial of service	665
ccrypt	encryption utility	751	segmentation fault	330
lighttpd	observer	51895	heap buffer overflow (code)	394
atris	technical game	21553	local buffer overflow	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computer biology	6718	segmentation fault	360
tiff	image processing		segmentation fault	108
imagemagick	image processing	1120	wrong output	2160

Program	Description	LOC	Bug Type	Time (s)

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153

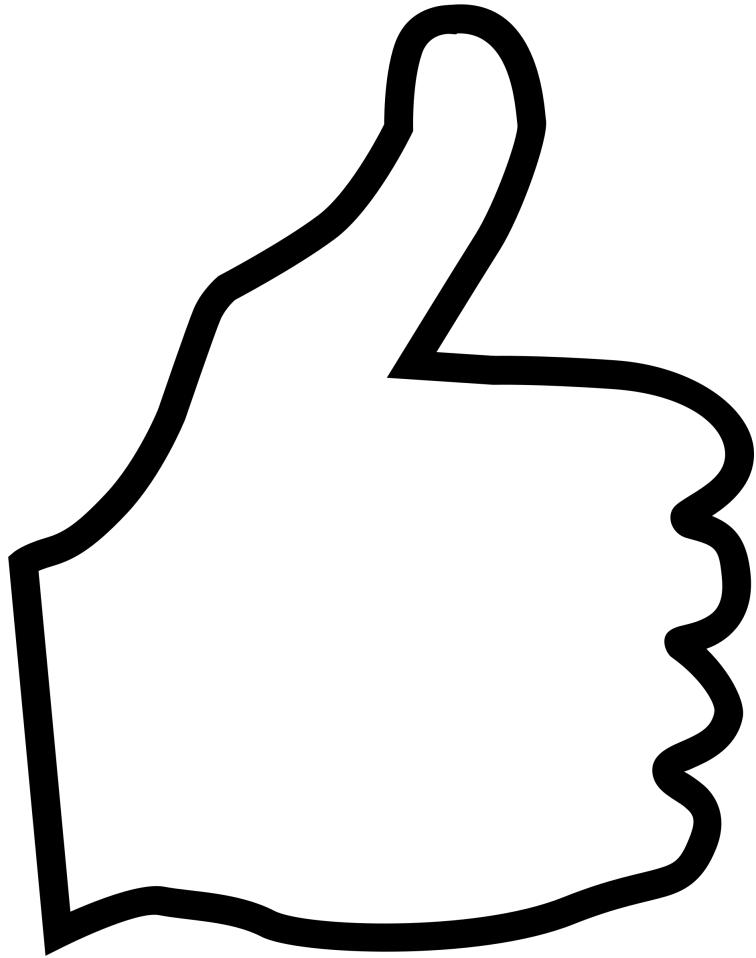


Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230

Program	Description	LOC	Bug Type	Time (s)
gcd	example	22	infinite loop	153
nullhttpd	webserver	5575	heap buffer overflow (code)	578
zune	example	28	infinite loop	42
uniq	text processing	1146	segmentation fault	34
look-u	dictionary lookup	1169	segmentation fault	45
look-s	dictionary lookup	1363	infinite loop	55
units	metric conversion	1504	segmentation fault	109
deroff	document processing	2236	segmentation fault	131
indent	code processing	9906	infinite loop	546
flex	lexical analyzer generator	18774	segmentation fault	230
openldap	directory protocol	292598	non-overflow denial of service	665
ccrypt	encryption utility	7515	segmentation fault	330
lighttpd	webserver	51895	heap buffer overflow (vars)	394
atris	graphical game	21553	local stack buffer exploit	80
php	scripting language	764489	integer overflow	56
wu-ftpd	FTP server	67029	format string vulnerability	2256
leukocyte	computational biology	6718	segmentation fault	360
tiff	image processing	84067	segmentation fault	108
imagemagick	image processing	450516	wrong output	2160



- Your audience will only remember 3 things.
- Tell a story.
- Never confuse your listeners.



Motto:

**YOU ARE NOT PRESENTING THE PAPER.  
YOU ARE PRESENTING THE WORK.**



- The audience will only remember 3 things.
- Tell a story.
- Never confuse your listeners.

8



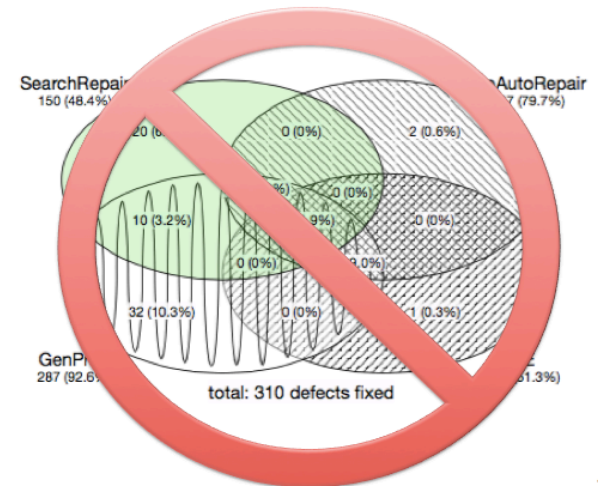
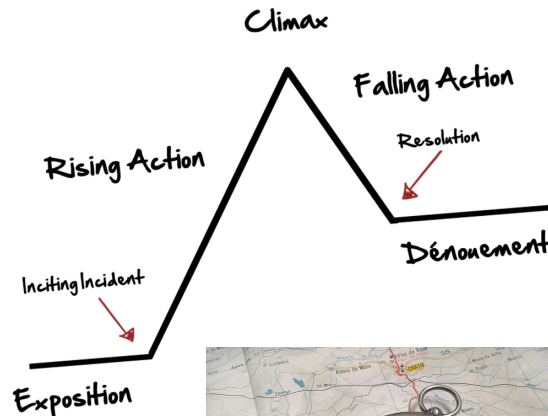
(Average audience member.)



### CLG's Goal

1. The *exciting and important* problem I am solving.
2. The *key nugget of awesomeness* underlying the approach.
3. 1—2 major result(s).
4. "That paper/person seems cool, I want to read it/talk to her!"

12



43