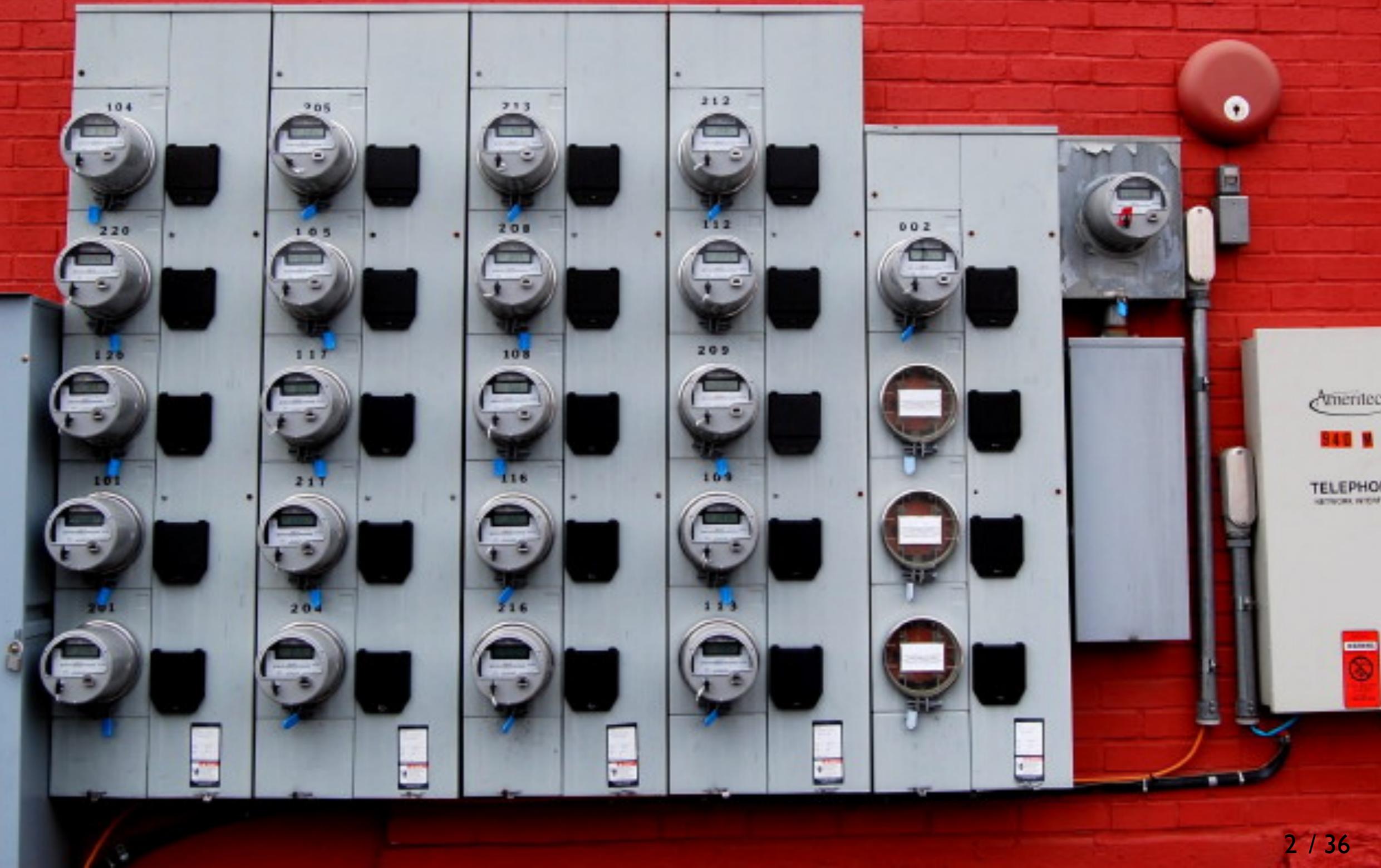


Building Program-Specific & Adaptable Alarm Prioritizer for



Soonho Kong
2009/6/26

Static analyzer produces **overwhelmingly** large number of alarms.





Service information

Unfortunately,
Many of them are **false alarms.**

TO ALL CUSTOMERS

PLEASE IGNORE THE
ALARM. ITS FALSE ALARM

THANK YOU



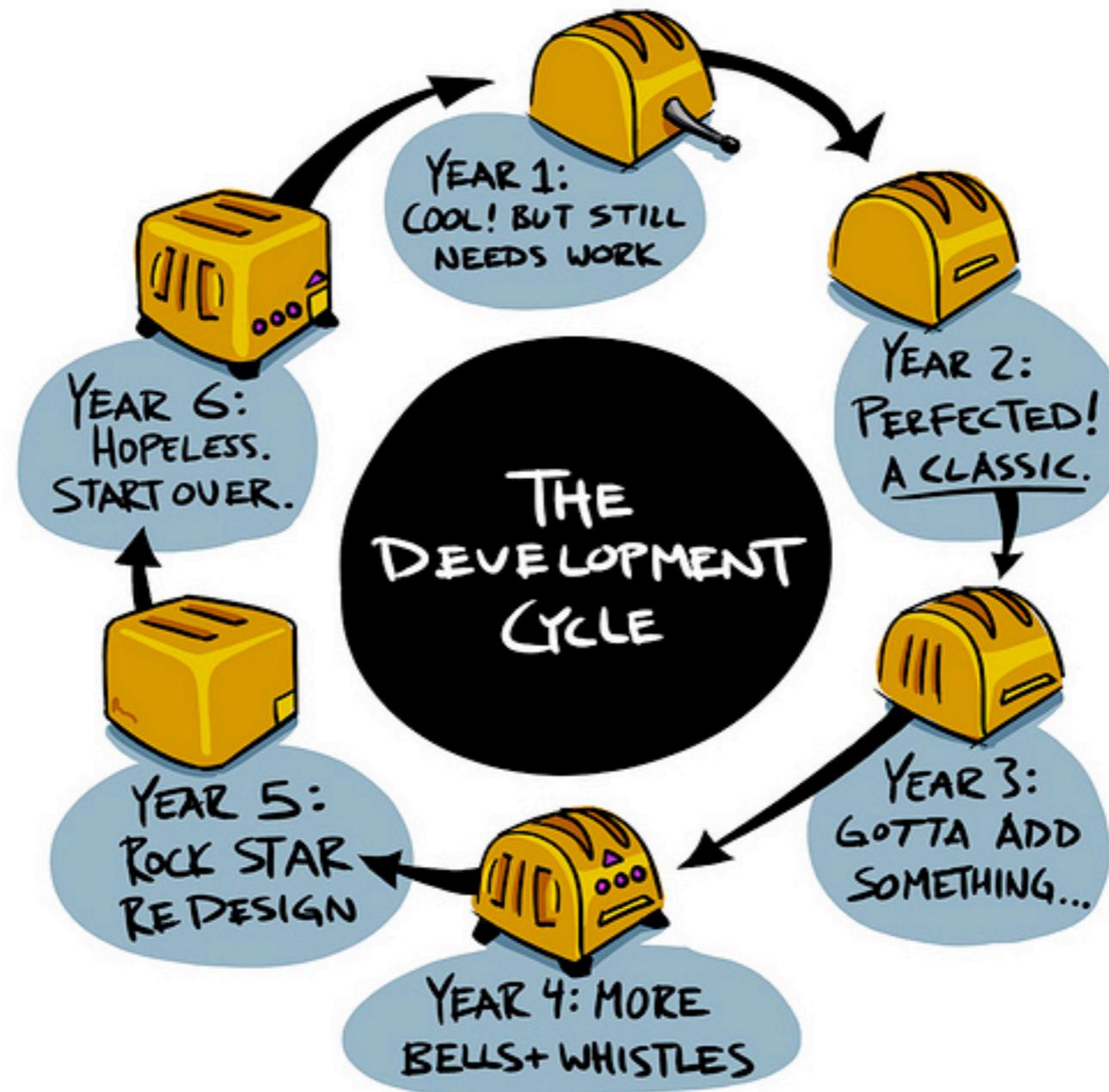
We need to **Manually** inspect alarms.

However,
manual inspection is
not interesting at all.



Moreover, Software Changes.

Program analysis is a part of development cycle.





Every build requires another inspections.
It's painful task.

**A FALSE ALARM
MAY COST A LIFE**



What do we need?



I. Alarm Grouping



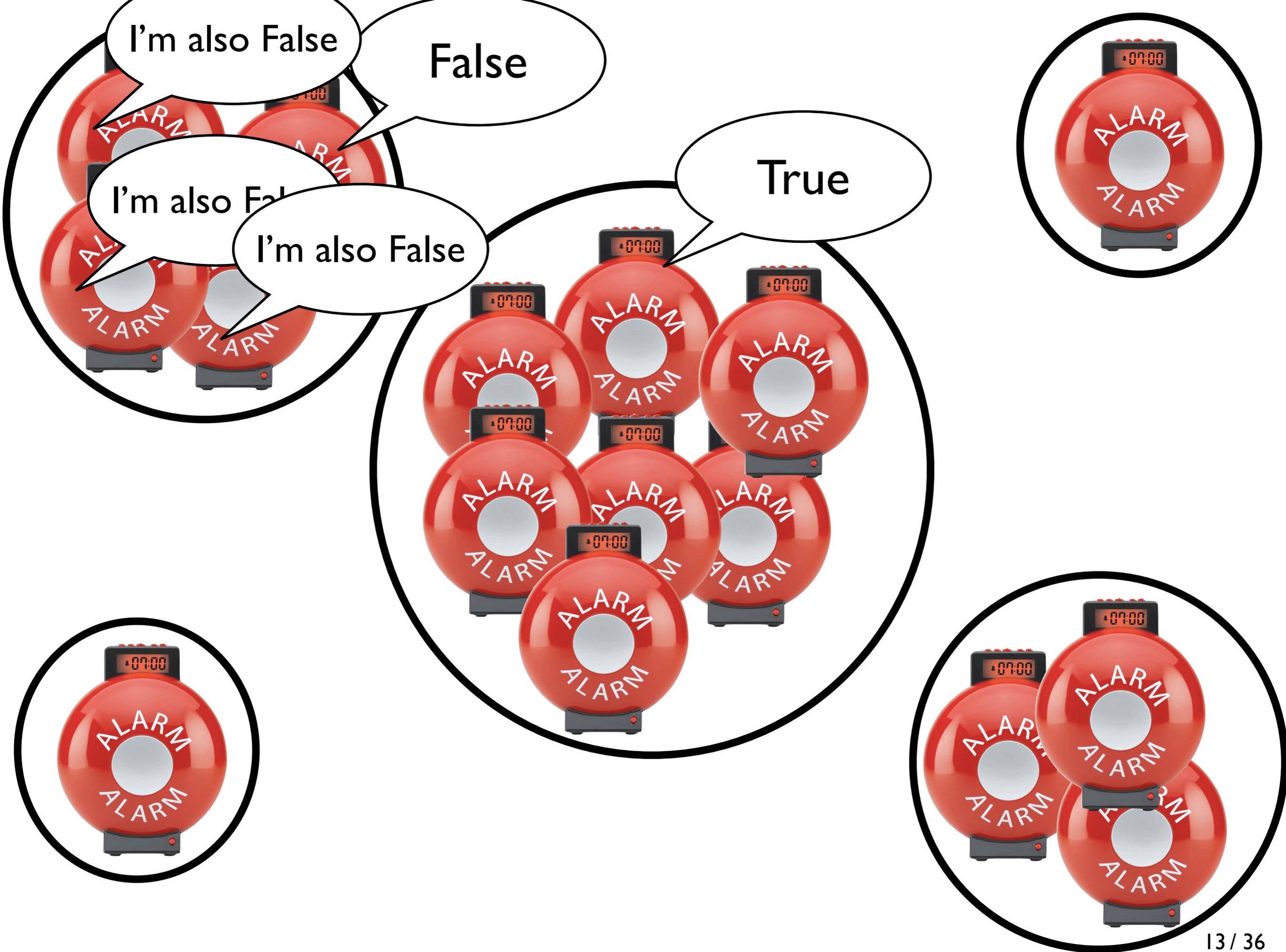


False









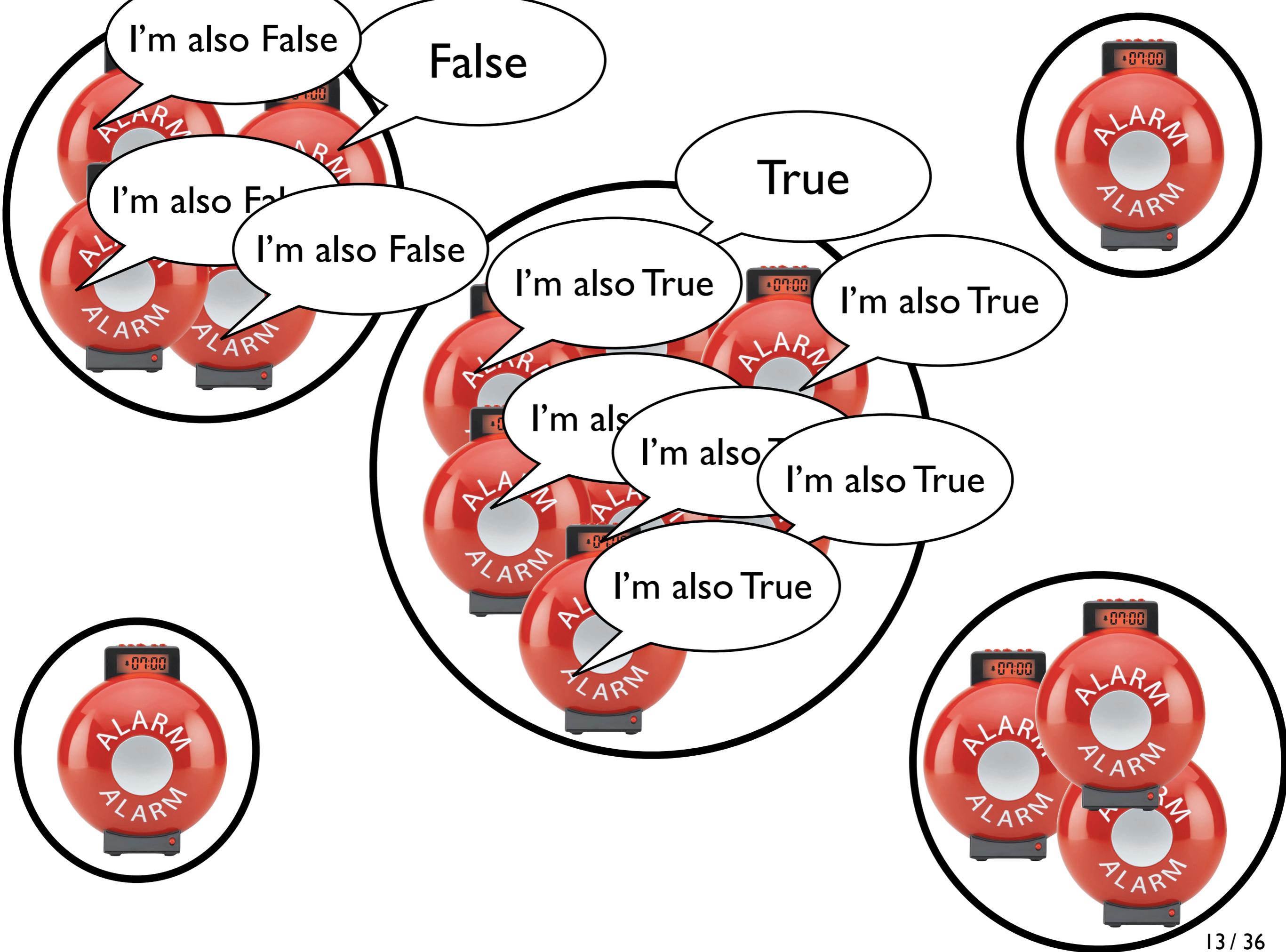
I'm also False

False

I'm also False

I'm also False

True



I'm also False

False

True

I'm also False

I'm also False

I'm also True



**Similar alarms
share common CID.**



Similar alarms
share common CID.



Show similar alarms
as a group.



Similar alarms
share common CID.



Show similar alarms
as a group.



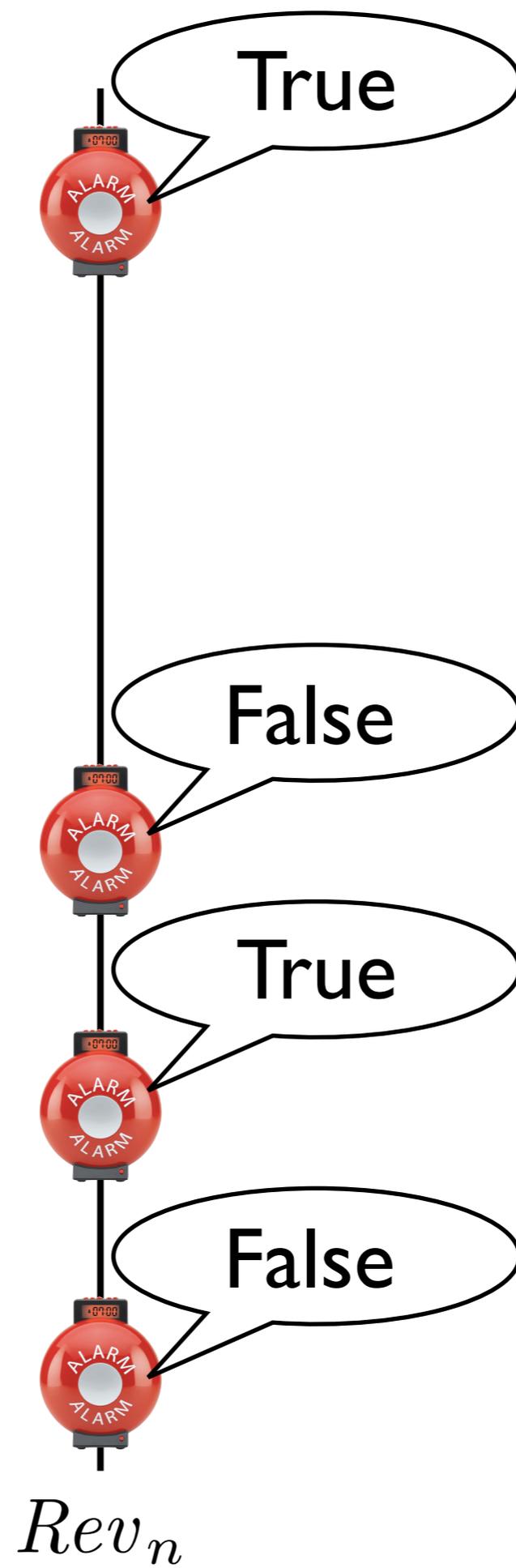
Not support, yet.

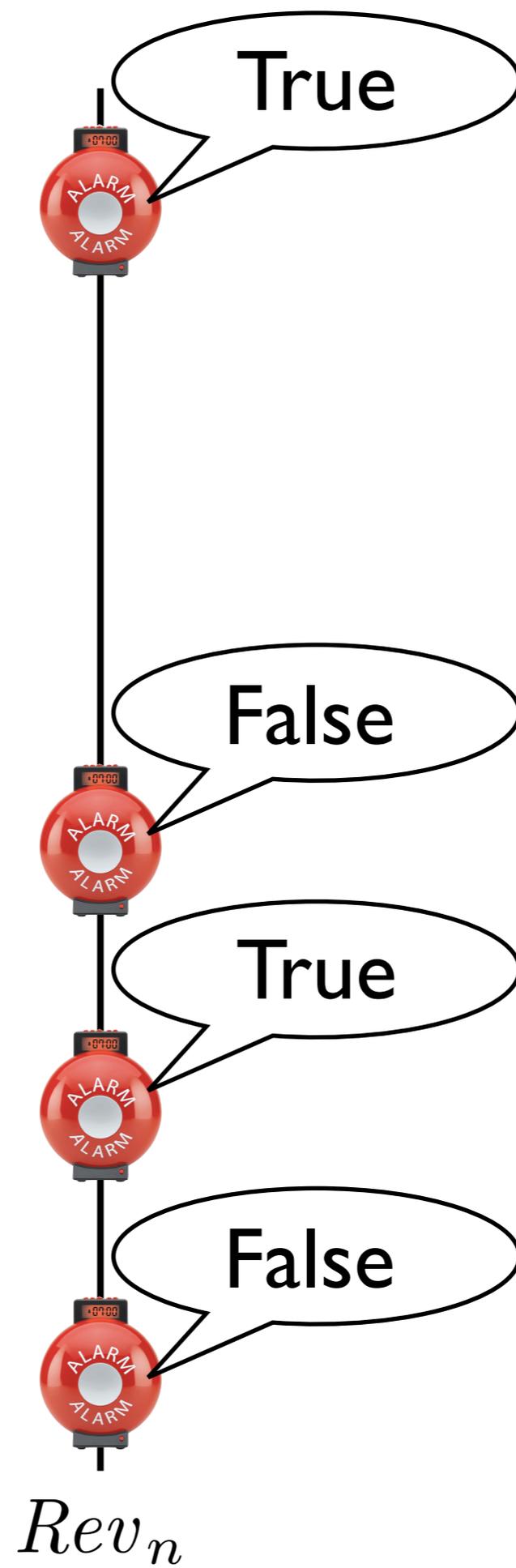
An aerial photograph of a sandy beach with a large colony of birds. A long, winding line of tracks leads from the water's edge towards the interior of the beach. The tracks are made of small, dark, Y-shaped marks, likely from birds. The birds themselves are visible as small dark specks along the tracks and scattered on the sand. The water is a light blue color, and the sky is a pale, hazy blue. The overall scene is a natural, undisturbed bird colony.

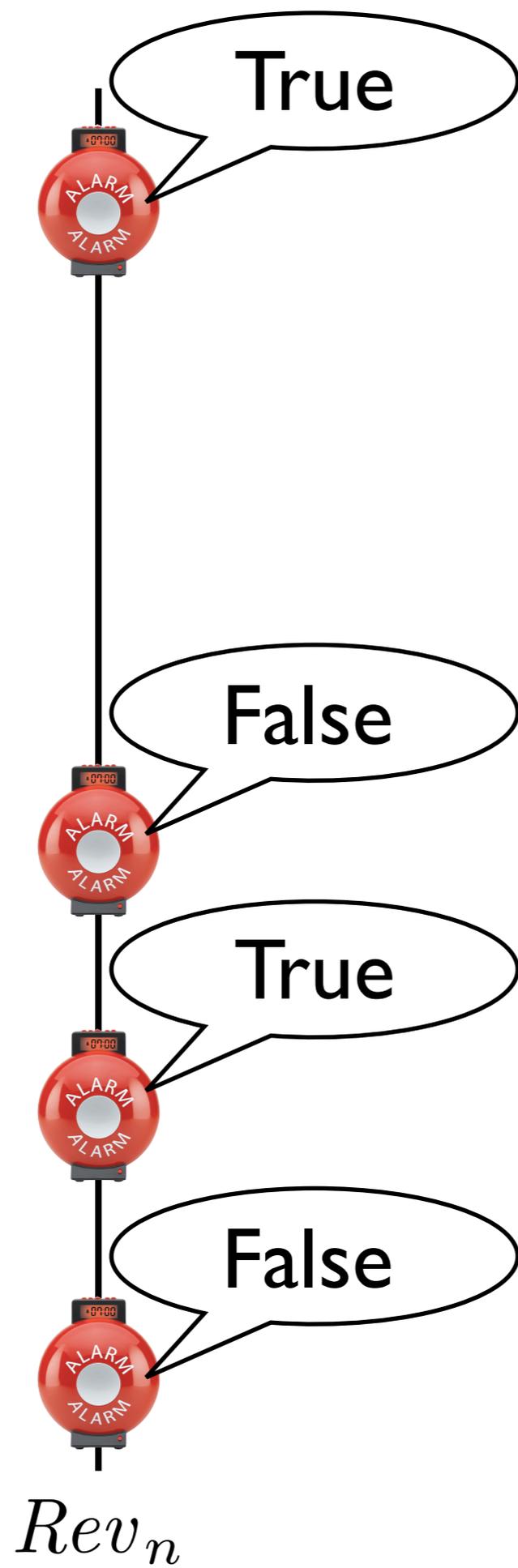
II. Alarm Tracking

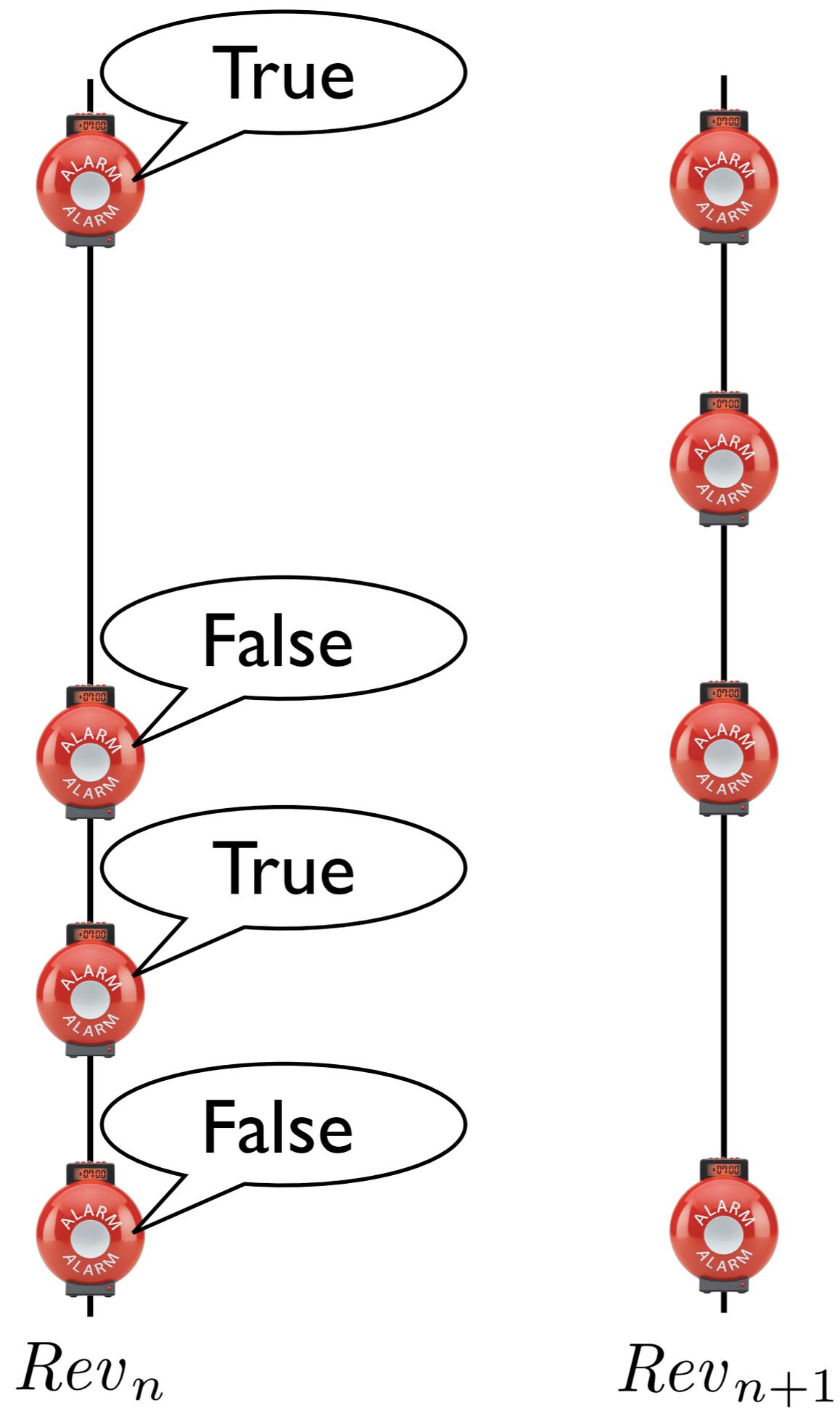


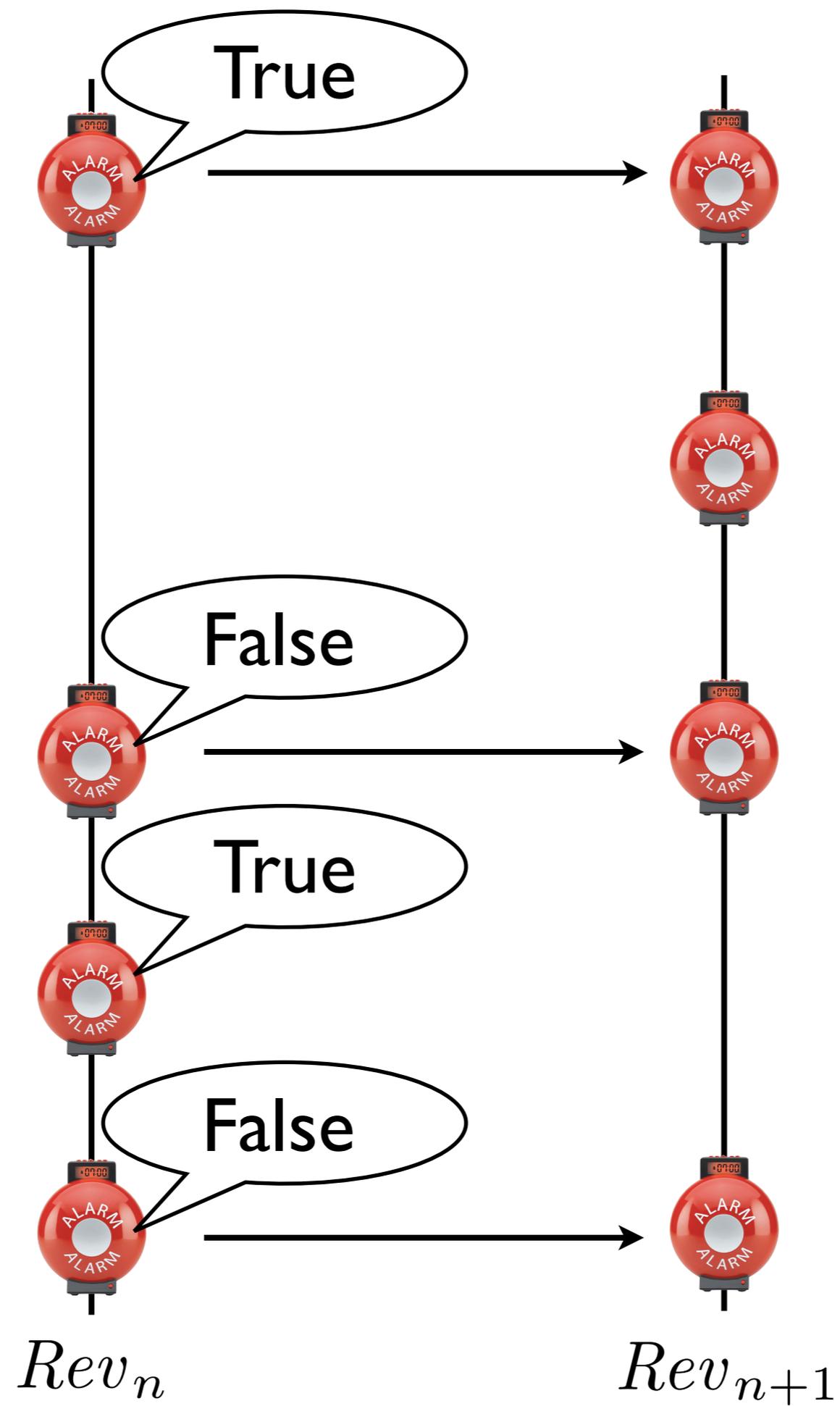
Rev_n

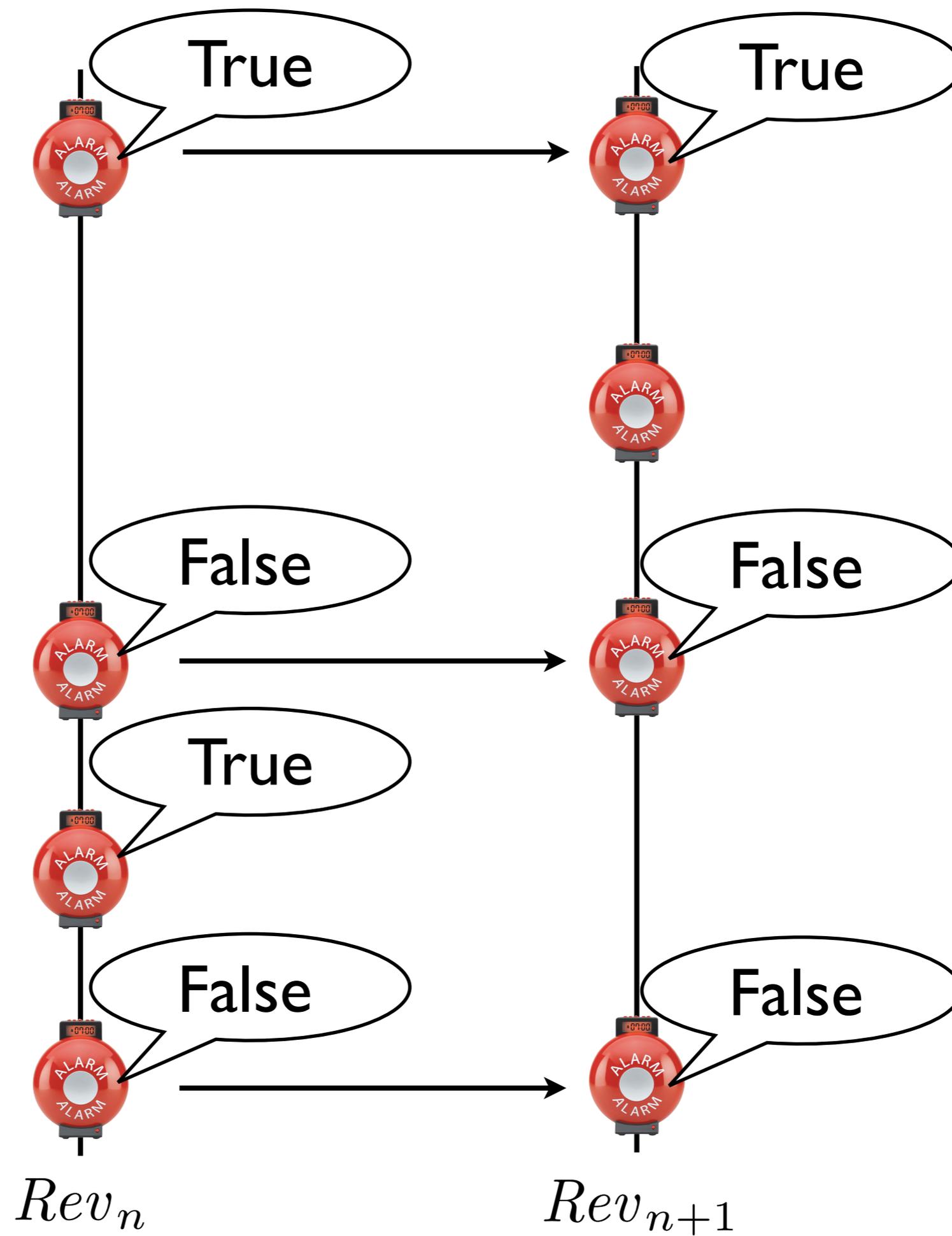












III. Alarm Prioritization





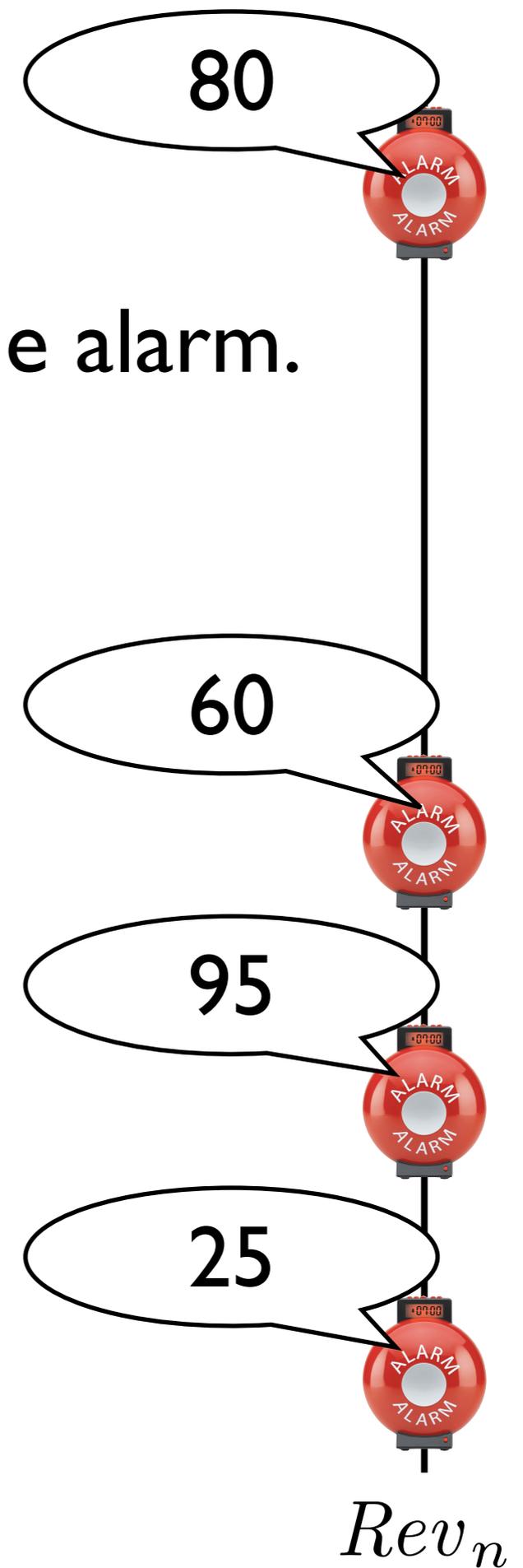
Rev_n

For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.

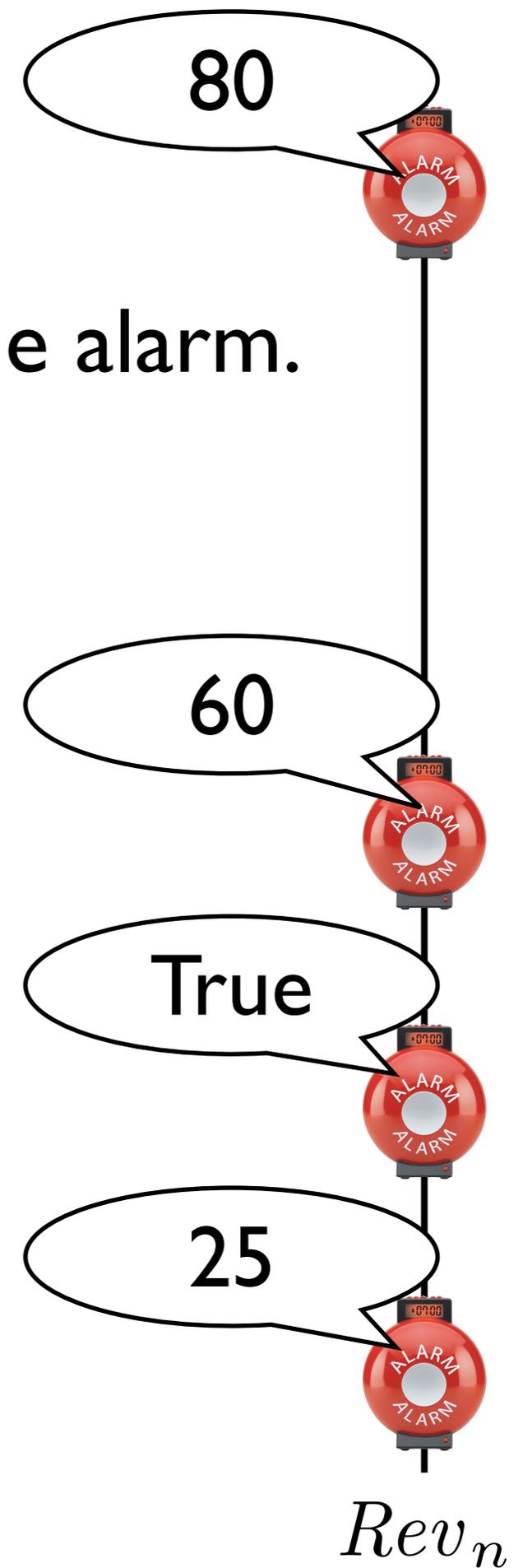


Rev_n

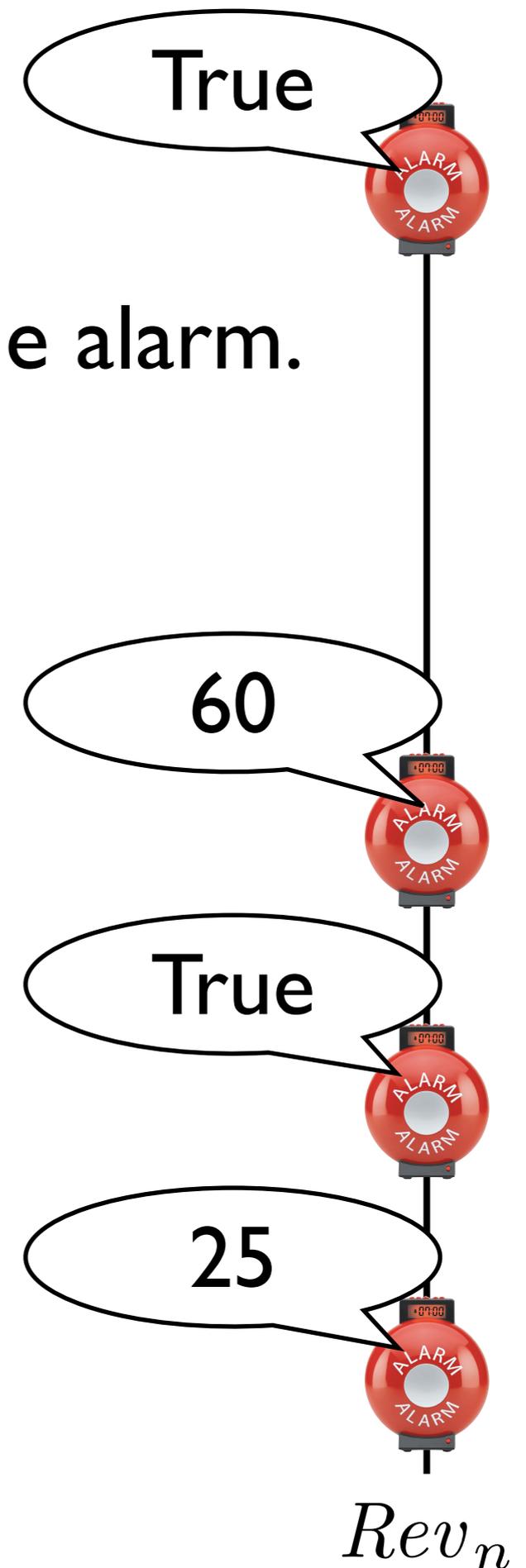
For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.



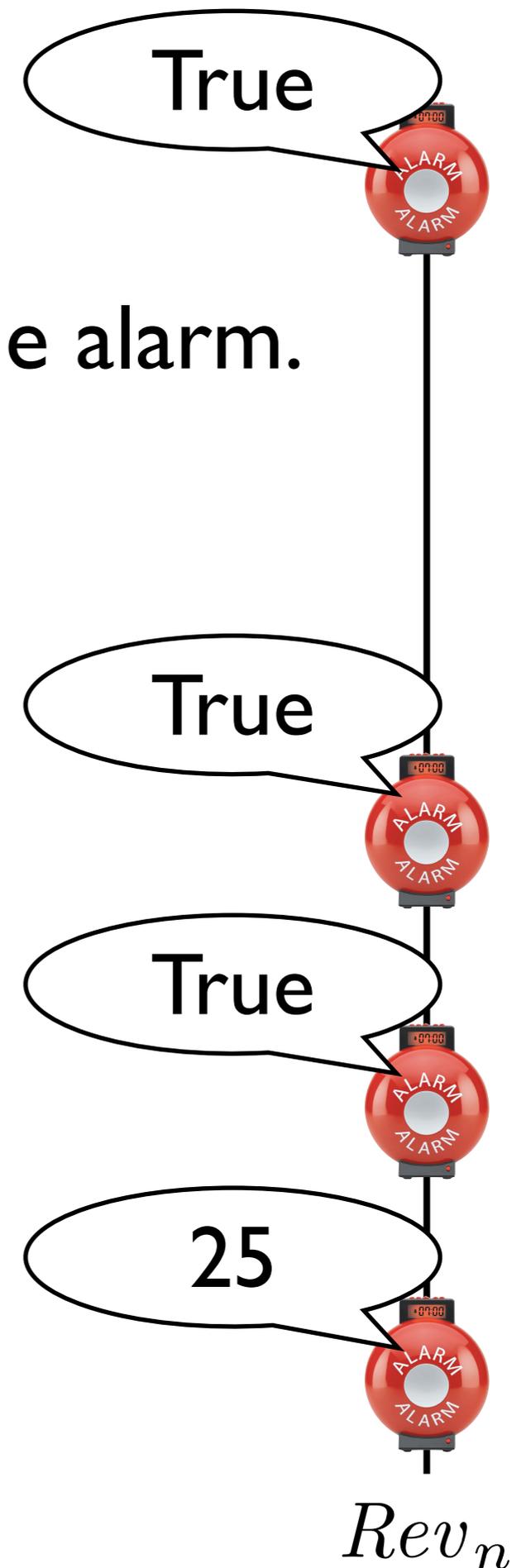
For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.



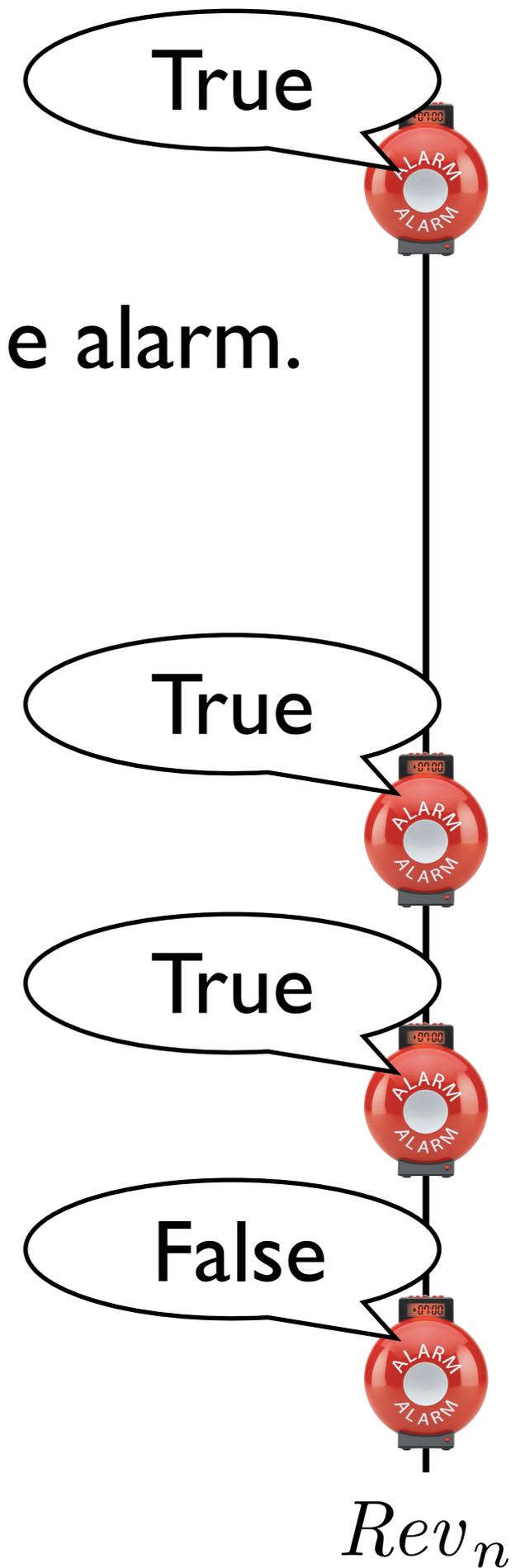
For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.



For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.



For each alarm,
we want to assign a **score**
indicating the likelihood of being a true alarm.



Taming False Alarms from a Domain-Unaware C Analyzer by a Bayesian Statistical Post Analysis

★

Yungbum Jung, Jaehwang Kim, Jaeho Shin and Kwangkeun Yi

Programming Research Laboratory
School of Computer Science and Engineering
Seoul National University
{dreameye, jaehwang, netj, kwang}@ropas.snu.ac.kr

Abstract. We present our experience of combining, in a realistic setting, a static analyzer with a statistical analysis. This combination is in order to reduce the inevitable false alarms from a domain-unaware static analyzer. Our analyzer named *Airac* (Array Index Range Analyzer for C) collects all the true buffer-overflow points in ANSI C programs. The soundness is maintained, and the analysis' cost-accuracy improvement is achieved by techniques that static analysis community has long accumulated. For still inevitable false alarms (e.g. *Airac* raised 970 buffer-overflow alarms in commercial C programs of 5.3 million lines and 737 among the 970 alarms were false), which are always apt for particular C programs, we use a statistical post analysis. The statistical analysis, given the analysis results (alarms), sifts out probable false alarms and prioritizes true alarms. It estimates the probability of each alarm being true. The probabilities are used in two ways: 1) only the alarms that have true-alarm probabilities higher than a threshold are reported to the user; 2) the alarms are sorted by the probability before reporting, so that the user can check highly probable errors first. In our experiments with Linux kernel sources, if we set the risk of missing true error is about 3 times greater than false alarming, 74.83% of false alarms could be filtered; only 15.17% of false alarms were mixed up until the user observes 50% of the true alarms.

1 Introduction

When one company's software quality assurance department started working with us to build a static analyzer that automatically detect buffer overruns¹ in

* This work was supported by Brain Korea 21 Project of Korea Ministry of Education and Human Resources, by IT Leading R&D Support Project of Korea Ministry of Information and Communication, by Korea Research Foundation grant KRF-2003-041-D00528, and by National Security Research Institute.

¹ Buffer overruns happen when an index value is out of the target buffer size. They are common bugs in C programs and are main sources of security vulnerability. From 1/2[2] to 2/3[1] of security holes are due to buffer overruns.

Taming False Alarms from a Domain-Unaware C Analyzer by a Bayesian Statistical Post Analysis

★

Yungbum Jung, Jaehwang Kim, Jaeho Shin and Kwangkeun Yi

Programming Research Laboratory
School of Computer Science and Engineering
Seoul National University
{dreameye, jaehwang, netj, kwang}@ropas.snu.ac.kr

Abstract. We present a Bayesian statistical analysis of combining, in a realistic setting, a static analyzer and a Bayesian statistical analysis. This combination is in order to reduce the number of false alarms from a domain-unaware static analyzer (Array Index Range Analyzer) in ANSI C programs. The Bayesian statistical analysis' cost-accuracy improvement is demonstrated by an experiment. The static analysis community has long been aware of the problem of false alarms. For example, the 070 buffer-overflow bug in the Linux kernel sources is a well-known example among the 070 alarm bugs in the Linux kernel sources. In C programs, we use a Bayesian statistical analysis to filter out false alarms given the analysis results. The Bayesian statistical analysis prioritizes true alarms based on the probability of each alarm being true. The probabilities are estimated by a Bayesian statistical analysis. We have true-alarm probabilities higher than a threshold. 1) the user can filter out false alarms; 2) the alarms are sorted by the probability before reporting, so that the user can check highly probable errors first. In our experiments with Linux kernel sources, if we set the risk of missing true error is about 3 times greater than false alarming, 74.83% of false alarms could be filtered; only 15.17% of false alarms were mixed up until the user observes 50% of the true alarms.

1 Introduction

When one company's software quality assurance department started working with us to build a static analyzer that automatically detect buffer overruns¹ in

* This work was supported by Brain Korea 21 Project of Korea Ministry of Education and Human Resources, by IT Leading R&D Support Project of Korea Ministry of Information and Communication, by Korea Research Foundation grant KRF-2003-041-D00528, and by National Security Research Institute.

¹ Buffer overruns happen when an index value is out of the target buffer size. They are common bugs in C programs and are main sources of security vulnerability. From 1/2[2] to 2/3[1] of security holes are due to buffer overruns.

Taming False Alarms from a Domain-Unaware C Analyzer by a Bayesian Statistical Post Analysis

★

Yungbum Jung, Jaehwang Kim, Jaeho Shin and Kwangkeun Yi

Programming Research Laboratory
School of Computer Science and Engineering
Seoul National University
{dreameye, jaehwang, netj, kwang}@ropas.snu.ac.kr

Abstract. We present a Bayesian statistical analysis of combining, in a realistic setting, a static code analysis and a dynamic analysis. This combination is designed to reduce the number of false alarms from a domain-unaware static code analyzer (Array Index Range Analyzer) by analyzing the execution points in ANSI C programs. The analysis' cost-accuracy improvement is evaluated. The static analysis community has long been aware of the problem of false alarms. For example, the 070 buffer-overrun checker [1] reported 1000000 false alarms among 1000000 true alarms. In C programs, we use a Bayesian statistical analysis to prioritize true alarms. The probabilities are estimated from the analysis results. We have true-alarm probabilities higher than a threshold. 1) the user can filter out false alarms; 2) the alarms are sorted by the probability before reporting, so that the user can check highly probable errors first. In our experiments with Linux kernel sources, if we set the risk of missing true error is about 3 times greater than false alarming, 74.83% of false alarms could be filtered; only 15.17% of false alarms were mixed up until the user observes 50% of the true alarms.

1 Introduction

Generic Prioritizer

Security assurance department started working at automatically detect buffer overruns¹ in Korea 21 Project of Korea Ministry of Education and Science, and R&D Support Project of Korea Ministry of Science and Technology, and Korea Research Foundation grant KRF-2003-

041-D00528, and by National Security Research Institute.

¹ Buffer overruns happen when an index value is out of the target buffer size. They are common bugs in C programs and are main sources of security vulnerability. From 1/2[2] to 2/3[1] of security holes are due to buffer overruns.

Taming False Alarms from a Domain-Unaware C Analyzer by a Bayesian Statistical Post Analysis

★

Yungbum Jung, Jaehwang Kim, Jaeho Shin and Kwangkeun Yi

Programming Research Laboratory
School of Computer Science and Engineering
Seoul National University
{dreameye, jaehwang, netj, kwang}@ropas.snu.ac.kr

Abstract. We present a Bayesian statistical analysis of combining, in a realistic setting, a static code analyzer and a Bayesian statistical analysis. This combination is designed to reduce the number of false alarms from a domain-unaware static code analyzer. We use a static code analyzer (Array Index Range Analyzer) to find potential buffer overflows in ANSI C programs. The Bayesian statistical analysis' cost-accuracy improvement is demonstrated by experiments. The static analysis community has long been aware of the problem of false alarms. For example, 4070 buffer-overflows were reported among 100,000 alarms in C programs, we use a Bayesian statistical analysis to prioritize true alarms. The probabilities are used to filter out false alarms. The user can have true-alarm probabilities higher than a threshold. 1) the user can filter out false alarms; 2) the alarms are sorted by the probability before reporting, so that the user can check highly probable errors first. In our experiments with Linux kernel sources, if we set the risk of missing true error is about 3 times greater than false alarming, 74.83% of false alarms could be filtered; only 15.17% of false alarms were mixed up until the user observes 50% of the true alarms.

1 Introduction

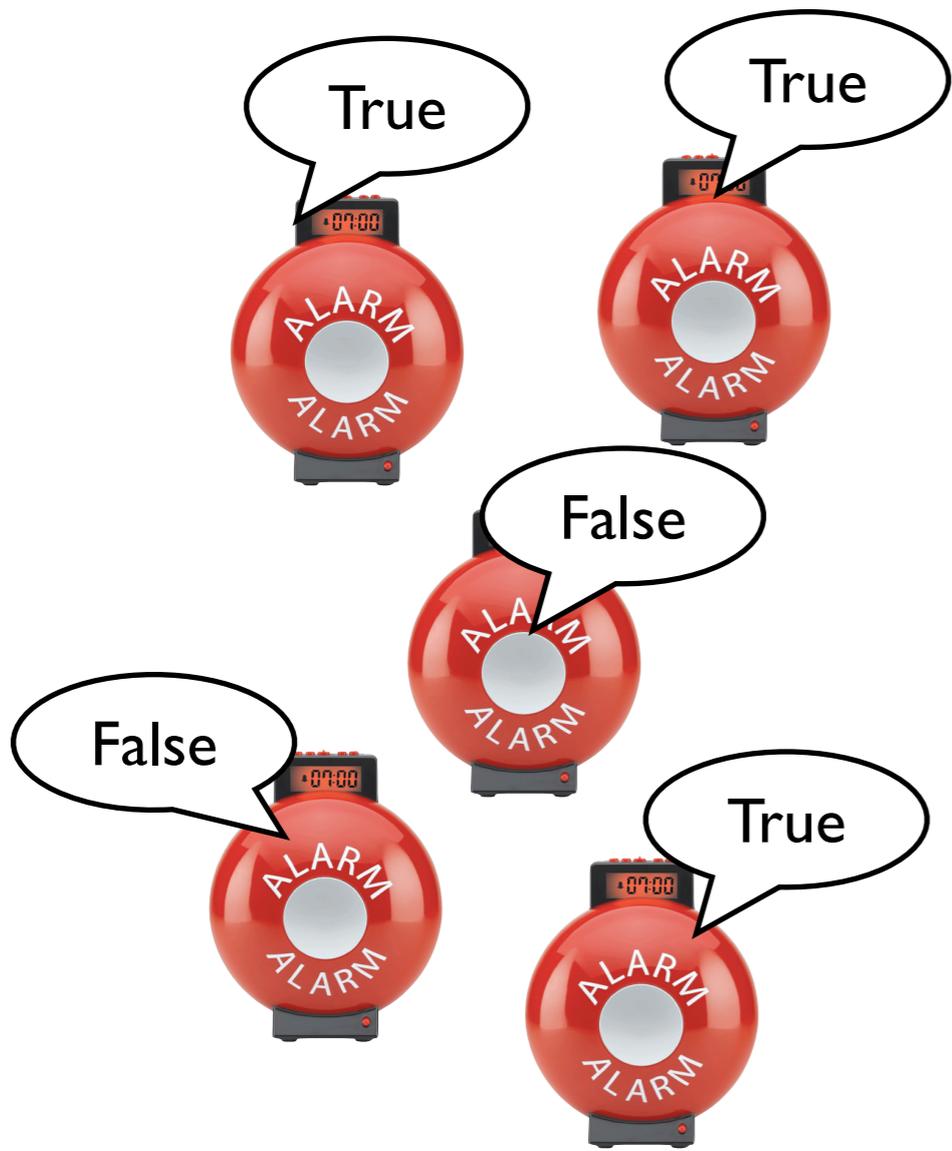
Generic Prioritizer

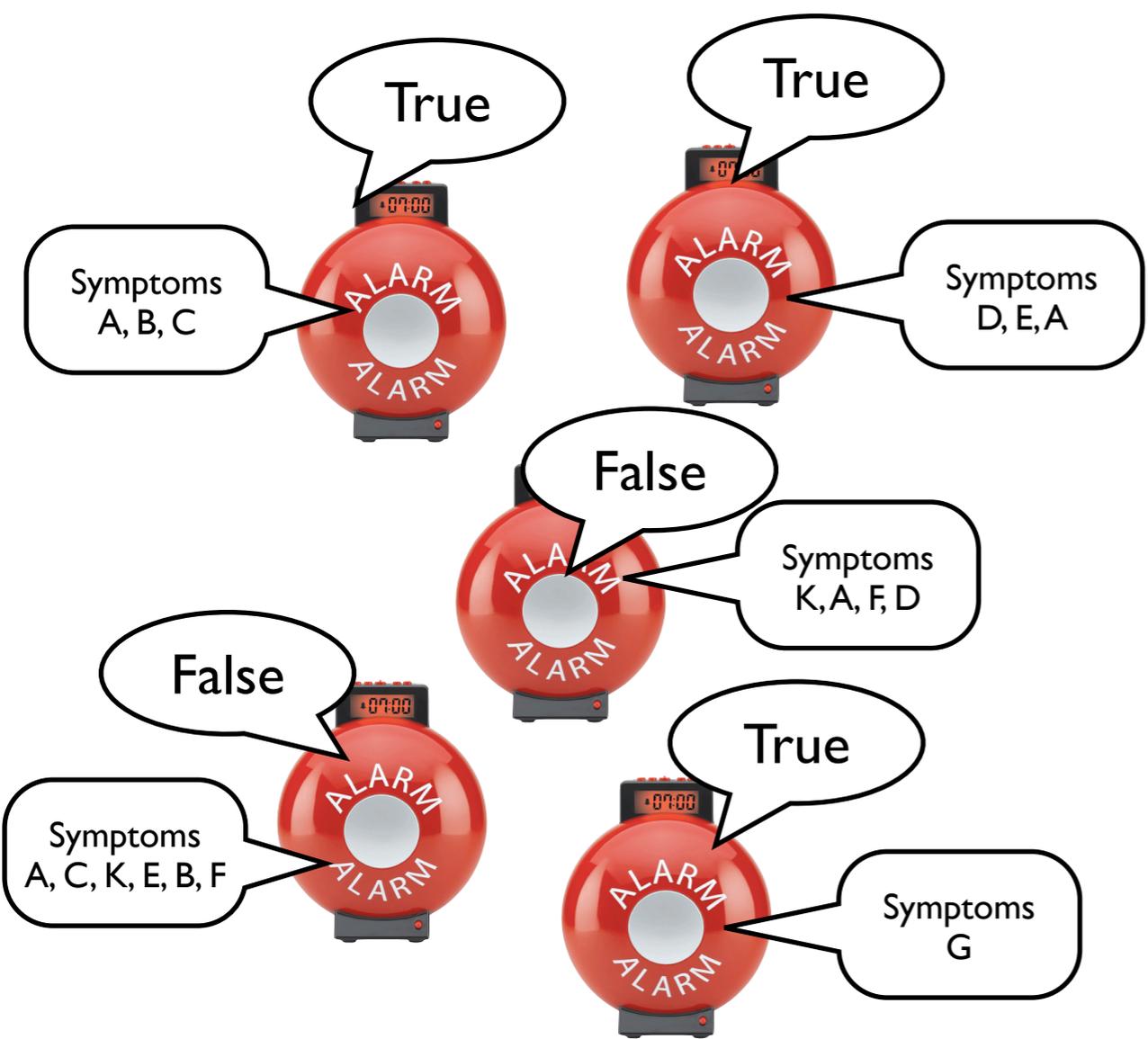
Not Adaptable

041-D00528, and by National Security Research Institute.

¹ Buffer overruns happen when an index value is out of the target buffer size. They are common bugs in C programs and are main sources of security vulnerability. From 1/2[2] to 2/3[1] of security holes are due to buffer overruns.



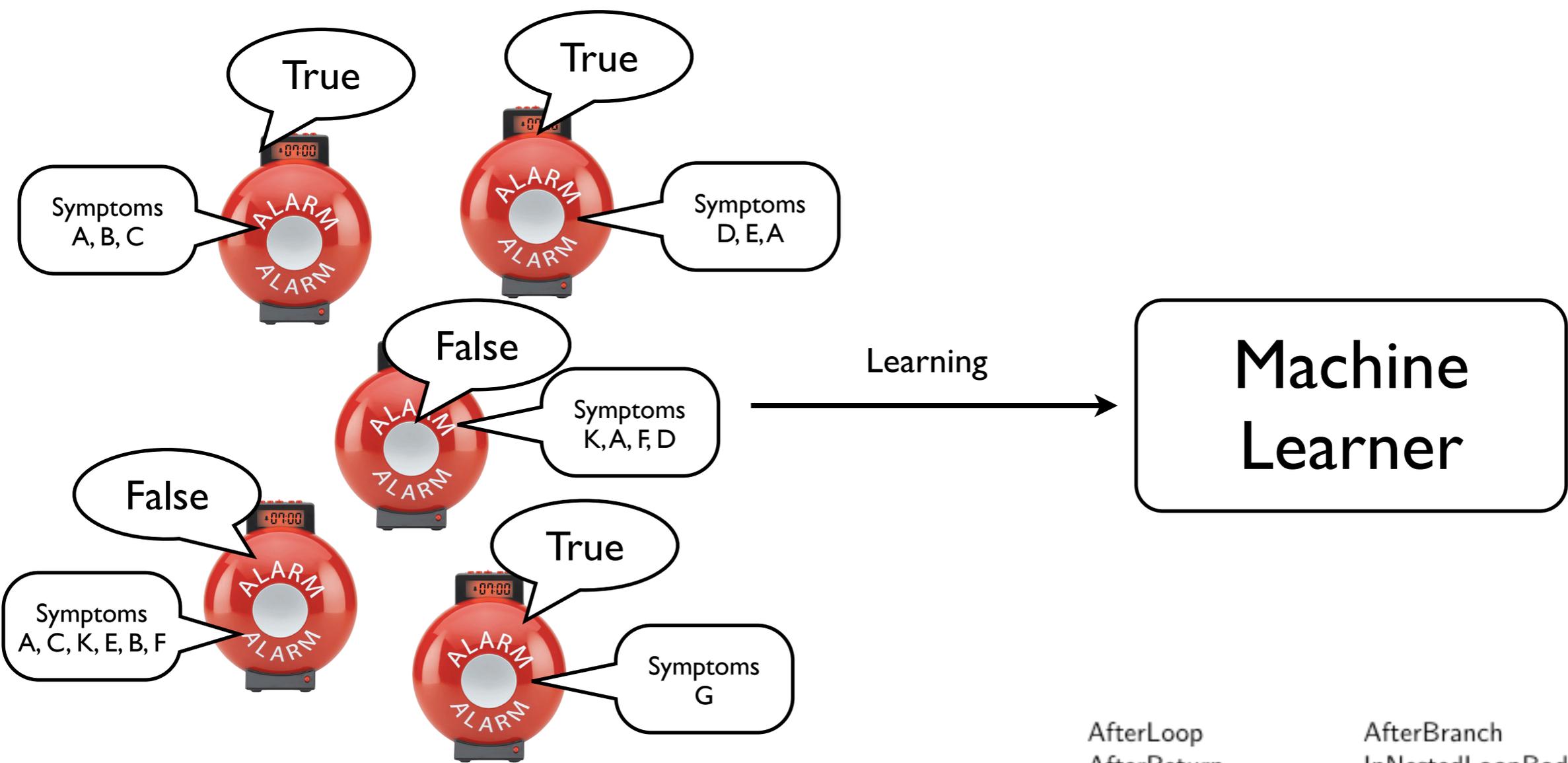




TopIndex HalfInfiniteIndex FiniteIndex

AfterLoop AfterBranch
 AfterReturn InNestedLoopBodyN
 InNestedBranchBodyN InLoopCond
 InBranchCond InFunParam
 InNestedFunParam InRightOfAnd

JoinN Pruned
 Narrowed PassedVal InStructure



TopIndex HalfInfiniteIndex FiniteIndex

AfterLoop AfterBranch
 AfterReturn InNestedLoopBodyN
 InNestedBranchBodyN InLoopCond
 InBranchCond InFunParam
 InNestedFunParam InRightOfAnd

JoinN Pruned
 Narrowed PassedVal InStructure

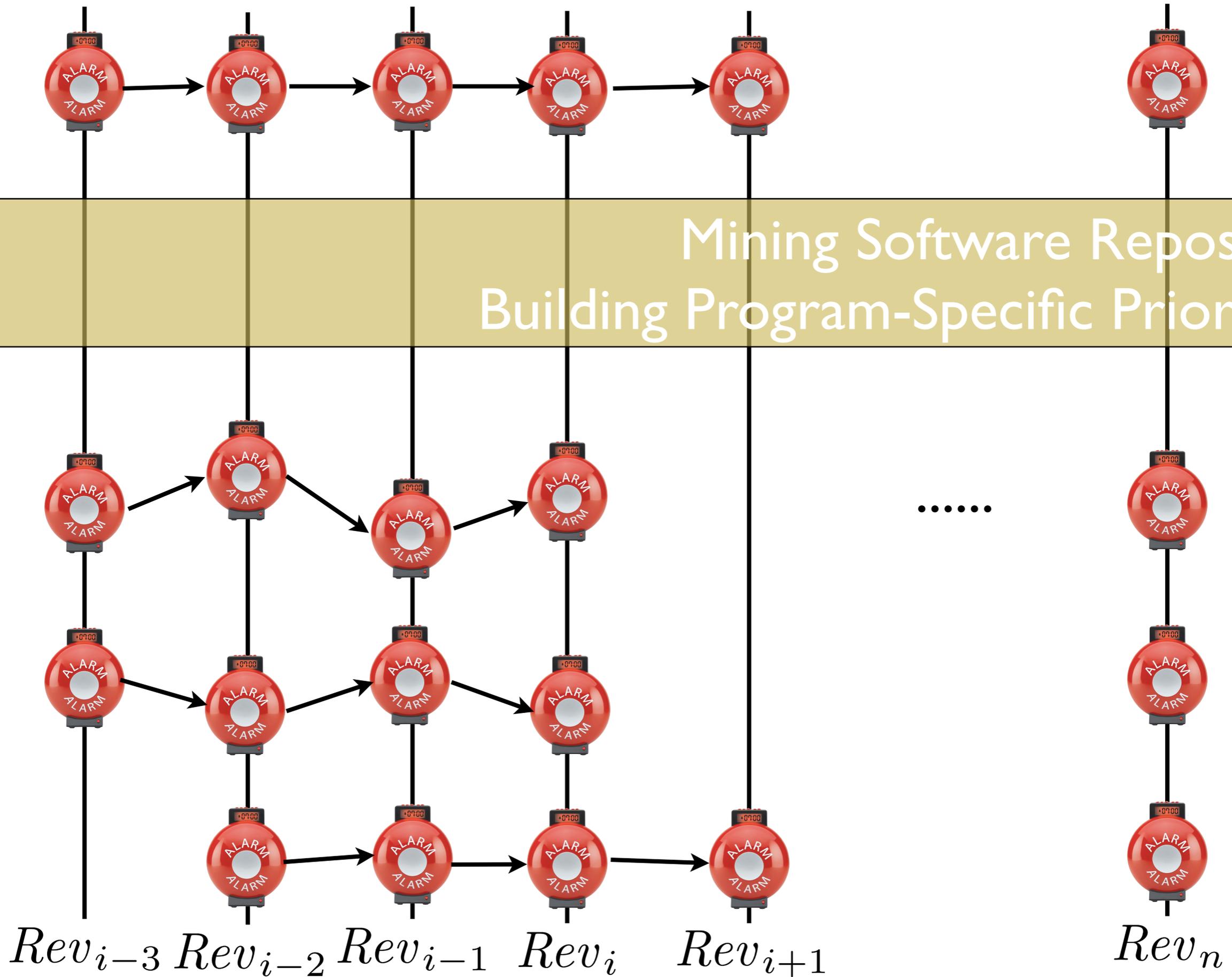
Modern Software uses Repository.





Mining Software Repository,
Building Program-Specific Prioritizer.

Mining Software Repository,
Building Program-Specific Prioritizer.



Previous Approach #1

Distinguished Paper

Predicting Accurate and Actionable Static Analysis Warnings: An Experimental Approach

Joseph R. Ruthruff*, John Penix†, J. David Morgenthaler†,
Sebastian Elbaum*, and Gregg Rothermel*

*University of Nebraska–Lincoln
Lincoln, NE, U.S.A.

{ruthruff, elbaum, grother}@cse.unl.edu

†Google Inc.

Mountain View, CA, U.S.A.

{jpenix, jdm}@google.com

ABSTRACT

Static analysis tools report software defects that may or may not be detected by other verification methods. Two challenges complicating the adoption of these tools are spurious false positive warnings and legitimate warnings that are not acted on. This paper reports automated support to help address these challenges using logistic regression models that predict the foregoing types of warnings from signals in the warnings and implicated code. Because examining many potential signaling factors in large software development settings can be expensive, we use a screening methodology to quickly discard factors with low predictive power and cost-effectively build predictive models. Our empirical evaluation indicates that these models can achieve high accuracy in predicting accurate and actionable static analysis warnings, and suggests that the models are competitive with alternative models built without screening.

Categories and Subject Descriptors

D.2.4 [Software/Program Verification]: Reliability, Statistical methods; F.3.2 [Semantics of Programming Languages]: Program analysis; G.3 [Probability and Statistics]: Correlation and regression analysis

General Terms

Experimentation, Reliability

Keywords

static analysis tools, screening, logistic regression analysis, experimental program analysis, software quality

1. INTRODUCTION

Static analysis tools detect software defects by analyzing a system without actually executing it. These tools utilize information from fixed program representations such as source code, generated or compiled code, and abstractions or models of the system. Even relatively simple analyses, such as detecting pointer dereferences after null checks, can find many defects in real software [5, 9].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICSE'08, May 10–18, 2008, Leipzig, Germany.

Copyright 2008 ACM 978-1-60558-079-1/08/05 ...\$5.00.

There are well-known challenges regarding the use of static analysis tools. One challenge involves the accuracy of reported warnings. Because the software under analysis is not executed, static analysis tools must speculate on what the actual program behavior will be. They often over-estimate possible program behaviors, leading to spurious warnings (“false positives”) that do not correspond to true defects. For example, Kremenek et al. [13] report that at least 30% of the warnings reported by sophisticated tools are false positives. At Google, we have observed that tools can be more accurate for certain types of warnings. Our experience with FindBugs [1] showed that focusing on selected, high priority warnings resulted in a 17% false positive rate [3].

A second challenge receiving less attention is that warnings are not always acted on by developers even if they reveal true defects. In the same study at Google, only 55% of the legitimate FindBugs warnings were acted on by developers after being entered into a bug tracking system [3]. Reasons for defects being ignored include warnings implicating obsolete code, “trivial” defects with no impact on the user, and real defects requiring significant effort to fix with little perceived benefit. Low criticality warnings such as “style” warnings, for example, can be unlikely to result in fixes.

We are investigating automated tools to help address both of these challenges by identifying *legitimate warnings that will be acted on by developers*, reducing the effort required to triage tens of thousands of warnings that can be reported in enterprise-wide settings. Our reasons for focusing on legitimate warnings are clear. We further focus on warnings that will be acted on by developers—not because ignored warnings are unimportant, but because we seek to maximize the return on investment from using static analysis tools. The core elements of our approach are statistical models generating binary classifications of static analysis warnings. One unique aspect of these models is that they are built using *screening*, an incremental statistical process to quickly discard factors with low predictive power and avoid the capture of expensive data.

Sampling from a base of tens of thousands of static analysis warnings from Google, we have built models that predict whether FindBugs warnings are false positives and, if they reveal real defects, whether these defects would be acted on by developers (“actionable warnings”) or ignored despite their legitimacy (“trivial warnings”). The generated models were over 85% accurate in predicting false positives, and over 70% accurate in identifying actionable warnings, in a case study performed at Google. Both represent a notable improvement over previous practices at Google for FindBugs, where 24% of triaged warning reports had been false positives and 56% had been acted on. The results from this study also indicate that screening can yield large savings in the time required to generate models, while sacrificing little predictive power.

Previous Approach #1

Distinguished Paper

Predicting Accurate and Actionable Static Analysis Warnings: An Experimental Approach

Joseph R. Ruthruff*, John Penix†, J. David Morgenthaler†, Sebastian Elbaum*, and Gregg Rothermel*

*University of Nebraska–Lincoln
Lincoln, NE, U.S.A.
{ruthruff, elbaum, grother}@cse.unl.edu

†Google Inc.
Mountain View, CA, U.S.A.
{jpenix, jdm}@google.com

ABSTRACT

Static analysis tools report software defects that may or may not be detected by other verification methods. Two challenges complicating the adoption of these tools are spurious false positive warnings and legitimate warnings that are not acted on. This paper reports automated support to help address these challenges using logistic regression models that predict the foregoing types of warnings from signals in the warnings and implicated code. Because examining many potential signaling factors in large software development settings can be expensive, we use a screening methodology to quickly discard factors with low predictive power and cost-effectively build predictive models. Our empirical evaluation indicates that these models can achieve high accuracy in predicting accurate and actionable static analysis warnings, and suggests that the models are competitive with alternative models built without screening.

Categories and Subject Descriptors

D.2.4 [Software/Program Verification]: Reliability, Statistical methods; F.3.2 [Semantics of Programming Languages]: Program analysis; G.3 [Probability and Statistics]: Correlation and regression analysis

General Terms

Experimentation, Reliability

Keywords

static analysis tools, screening, logistic regression analysis, experimental program analysis, software quality

1. INTRODUCTION

Static analysis tools detect software defects by analyzing a system without actually executing it. These tools utilize information from fixed program representations such as source code, generated or compiled code, and abstractions or models of the system. Even relatively simple analyses, such as detecting pointer dereferences after null checks, can find many defects in real software [5, 9].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICSE'08, May 10–18, 2008, Leipzig, Germany.
Copyright 2008 ACM 978-1-60558-079-1/08/05 ...\$5.00.

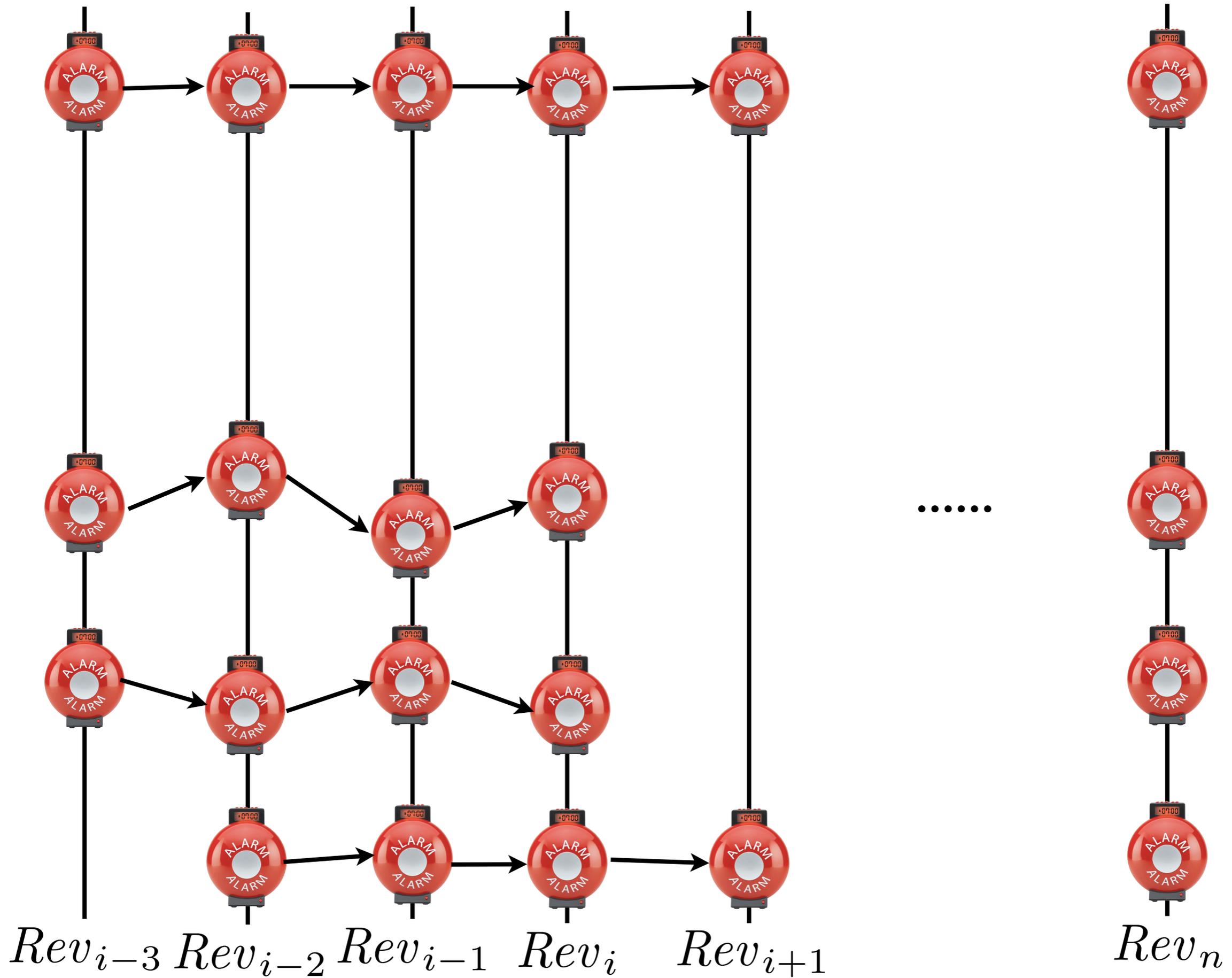
There are well-known challenges regarding the use of static analysis tools. One challenge involves the accuracy of reported warnings. Because the software under analysis is not executed, static analysis tools must speculate about actual program behavior. They often report warnings about program behaviors, leading to spurious warnings that do not correspond to true program behaviors. et al. [13] report that sophisticated tools can be used to filter out false positives. Our experience with FindBugs shows that selected, high priority warnings have a high fix rate [3].

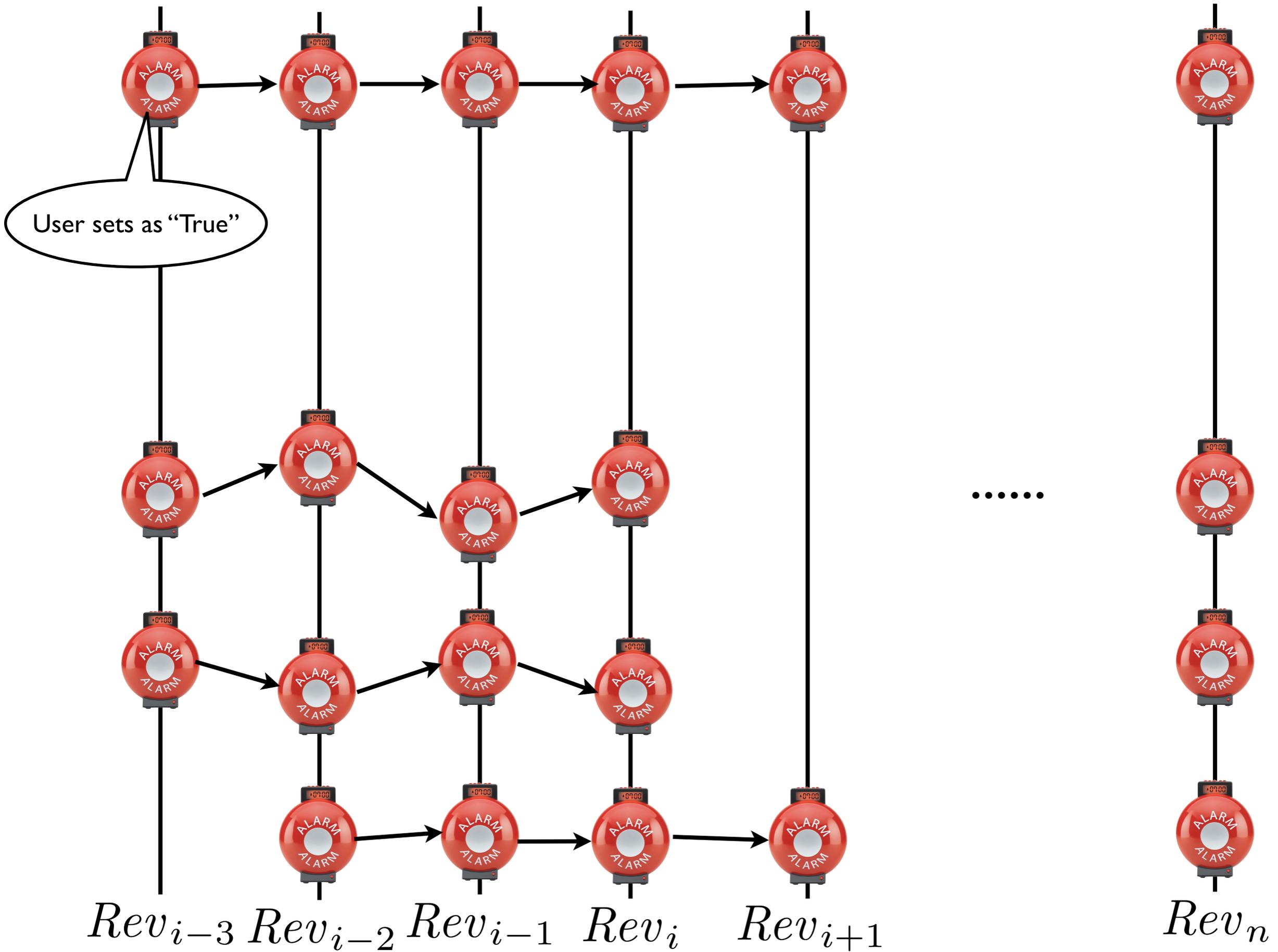
A second challenge receiving less attention is that warnings are not always acted on by developers even if they reveal true defects. In the same study at Google, only 55% of the legitimate FindBugs warnings were acted on by developers after being entered into a bug tracking system [3]. Reasons for defects being ignored include warnings implicating obsolete code, “trivial” defects with no impact on the user, and real defects requiring significant effort to fix with little perceived benefit. Low criticality warnings such as “style” warnings, for example, can be unlikely to result in fixes.

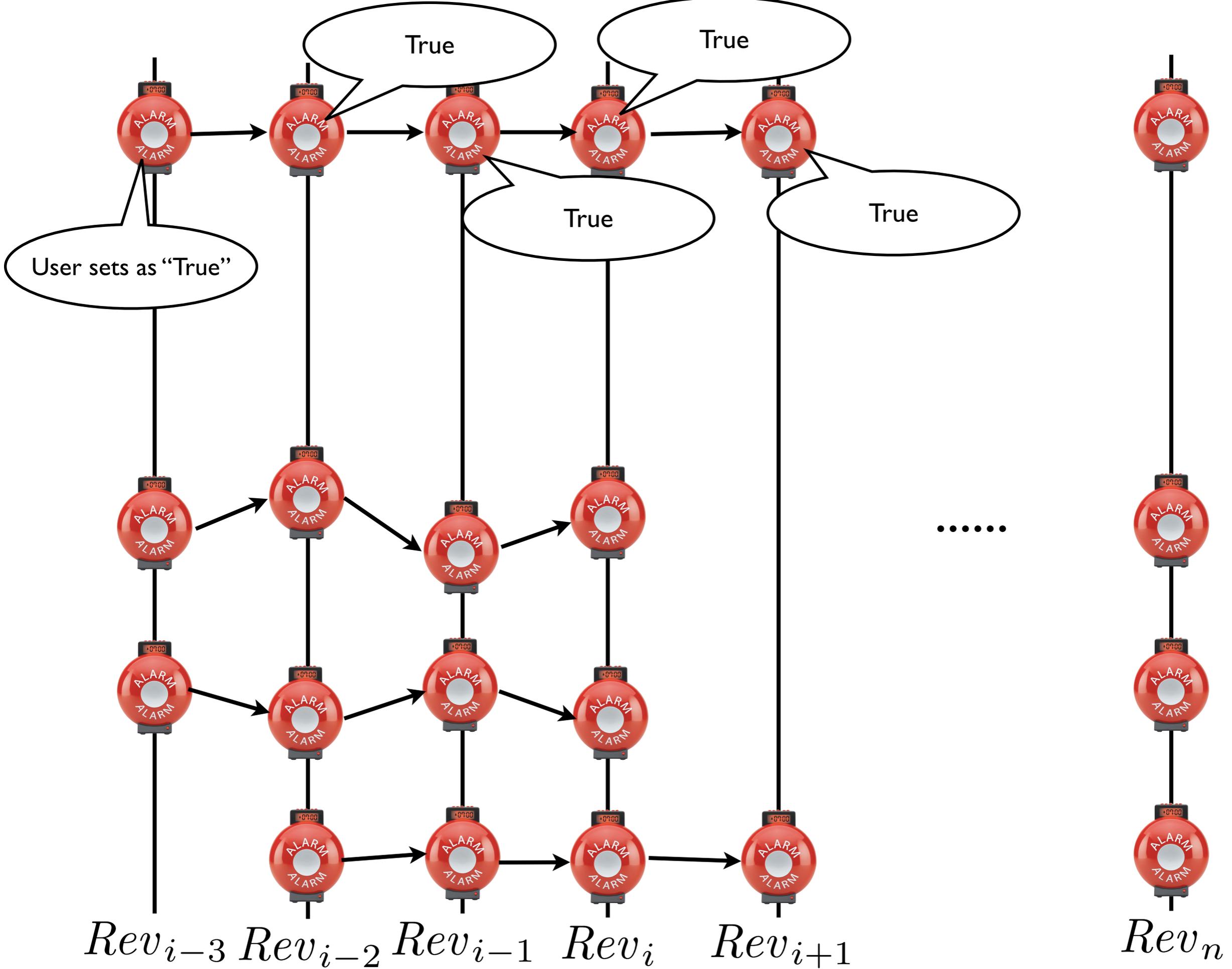
We are investigating automated tools to help address both of these challenges by identifying *legitimate warnings that will be acted on by developers*, reducing the effort required to triage tens of thousands of warnings that can be reported in enterprise-wide settings. Our reasons for focusing on legitimate warnings are clear. We further focus on warnings that will be acted on by developers—not because ignored warnings are unimportant, but because we seek to maximize the return on investment from using static analysis tools. The core elements of our approach are statistical models generating binary classifications of static analysis warnings. One unique aspect of these models is that they are built using *screening*, an incremental statistical process to quickly discard factors with low predictive power and avoid the capture of expensive data.

Sampling from a base of tens of thousands of static analysis warnings from Google, we have built models that predict whether FindBugs warnings are false positives and, if they reveal real defects, whether these defects would be acted on by developers (“actionable warnings”) or ignored despite their legitimacy (“trivial warnings”). The generated models were over 85% accurate in predicting false positives, and over 70% accurate in identifying actionable warnings, in a case study performed at Google. Both represent a notable improvement over previous practices at Google for FindBugs, where 24% of triaged warning reports had been false positives and 56% had been acted on. The results from this study also indicate that screening can yield large savings in the time required to generate models, while sacrificing little predictive power.

ICSE'08







Previous Approach #2

Which Warnings Should I Fix First?

Sunghun Kim and Michael D. Ernst

Computer Science & Artificial Intelligence Lab (CSAIL)
Massachusetts Institute of Technology

{hunkim, mernst}@csail.mit.edu

ABSTRACT

Automatic bug-finding tools have a high false positive rate: most warnings do not indicate real bugs. Usually bug-finding tools assign important warnings high priority. However, the prioritization of tools tends to be ineffective. We observed the warnings output by three bug-finding tools, FindBugs, JLint, and PMD, for three subject programs, Columba, Lucene, and Scarab. Only 6%, 9%, and 9% of warnings are removed by bug fix changes during 1 to 4 years of the software development. About 90% of warnings remain in the program or are removed during non-fix changes – likely false positive warnings. The tools' warning prioritization is little help in focusing on important warnings: the maximum possible precision by selecting high-priority warning instances is only 3%, 12%, and 8% respectively.

In this paper, we propose a history-based warning prioritization algorithm by mining warning fix experience that is recorded in the software change history. The underlying intuition is that if warnings from a category are eliminated by fix-changes, the warnings are important. Our prioritization algorithm improves warning precision to 17%, 25%, and 67% respectively.

Categories and Subject Descriptors

D.2.7 [Software Engineering]: Distribution, Maintenance, and Enhancement – *Restructuring, reverse engineering, and reengineering*, D.2.8 [Software Engineering]: Metrics – *Product metrics*, K.6.3 [Management of Computing and Information Systems]: Software Management – *Software maintenance*

General Terms

Algorithms, Measurement, Experimentation

Keywords

Fault, Bug, Fix, Bug-finding tool, Prediction, Patterns

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ESEC-FSE'07, September 3–7, 2007, Cavat near Dubrovnik, Croatia.
Copyright 2007 ACM 978-1-59593-811-4/07/0009...\$5.00.

1. INTRODUCTION

Bug-finding tools such as FindBugs [12], JLint [2], and PMD [6] analyze source or binary code and warn about potential bugs. These tools have a high rate of false positives: most warnings do not indicate real bugs [17]. Most bug-finding tools assign categories and priorities to warning instances, such as *Overflow* (priority 1) or *Empty Static Initializer* (priority 3). The tools' prioritization is supposed to put important warnings at the top of the list, but the prioritization is not very effective [17]. We performed two experiments that support this observation. Our experiments use three bug-finding tools (FindBugs, JLint, and PMD) and three subject programs (Columba, Lucene, and Scarab).

First, we measured the percentages of warnings that are actually eliminated by fix-changes, since generally a fix-change indicates a bug [7-9, 22]. We select a revision and determine warnings issued by the bug-finding tools. Only 6%, 9%, and 9% of warnings are removed by fix-changes during 1~4 years of the software change history of each subject program respectively – about 90% of warnings either remain or are removed during non-fix changes.

Second, we observed whether the tools' warning prioritization (TWP) favors important warnings. The maximum possible warning precision by selecting high priority warning instances is only 3%, 12%, and 8% respectively. This fact indicates that TWP is ineffective.

Our goal is to propose a new, program-specific prioritization that more effectively directs developers to errors. The new history-based warning prioritization (HWP) is obtained by mining the software change history for removed warnings during bug fixes.

A version control system indicates when each file is changed. A software change can be classified as a fix-change or a non-fix change. A fix-change is a change that fixes a bug or other problem. A non-fix change is a change that does not fix a bug, such as a feature addition or refactoring.

Suppose that during development, a bug-finding tool would issue a warning instance from the *Overflow* category. If a developer finds the underlying problem and fixes it, the warning is probably important. (We do not assume the software developer is necessarily using the bug-finding tool.) On the other hand, if a warning instance is not removed for a long time, then warnings of that category may be neglectable, since the problem was not noticed or was not considered worth fixing.

Using this intuition, we set a weight for each warning category to represent its importance. The weight of a category is proportional to the number of warning instances from that category that are eliminated by a change, with fix-changes contributing more to the

Previous Approach #2

Which Warnings Should I Fix First?

Sunghun Kim and Michael D. Ernst

Computer Science & Artificial Intelligence Lab (CSAIL)
Massachusetts Institute of Technology

{hunkim, mernst}@csail.mit.edu

ABSTRACT

Automatic bug-finding tools have a high false positive rate: most warnings do not indicate real bugs. Usually bug-finding tools assign important warnings high priority. However, the prioritization of tools tends to be ineffective. We observed the warnings output by three bug-finding tools, FindBugs, JLint, and PMD, for three subject programs, Columba, Lucene, and Scarab. Only 6%, 9%, and 9% of warnings are removed by bug fix changes during 1 to 4 years of the software development. About 90% of warnings remain in the program or are removed during non-fix changes – likely false positive warnings. The tools' warning prioritization is little help in focusing on important warnings: the maximum possible precision by selecting high-priority warning instances is only 3%, 12%, and 8% respectively.

In this paper, we propose a history-based warning prioritization algorithm by mining warning fix experience that is recorded in the software change history. The underlying intuition is that if warnings from a category are eliminated by fix-changes, the warnings are important. Our prioritization algorithm improves warning precision to 17%, 25%, and 67% respectively.

Categories and Subject Descriptors

D.2.7 [Software Engineering]: Distribution, Maintenance, and Enhancement – *Restructuring, reverse engineering, and reengineering*, D.2.8 [Software Engineering]: Metrics – *Product metrics*, K.6.3 [Management of Computing and Information Systems]: Software Management – *Software maintenance*

General Terms

Algorithms, Measurement, Experimentation

Keywords

Fault, Bug, Fix, Bug-finding tool, Prediction, Patterns

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
ESEC-FSE'07, September 3–7, 2007, Cavat near Dubrovnik, Croatia.
Copyright 2007 ACM 978-1-59593-811-4/07/0009...\$5.00.

1. INTRODUCTION

Bug-finding tools such as FindBugs [12], JLint [2], and PMD [6] analyze source or binary code and warn about potential bugs. These tools have a high rate of false positives: most warnings do not indicate real bugs [17]. Most bug-finding tools assign categories and priorities to warnings, such as *Overflow* (priority 1) or *Empty* (priority 3). The tools' prioritization is ineffective [17]. We performed an observation. Our experiments on FindBugs, JLint, and PMD (Columba, Lucene, and Scarab).

First, we observed whether the tools' warning prioritization (TWP) favors important warnings. The maximum possible warning precision by selecting high priority warning instances is only 3%, 12%, and 8% respectively. This fact indicates that TWP is ineffective.

Second, we observed whether the tools' warning prioritization (TWP) favors important warnings. The maximum possible warning precision by selecting high priority warning instances is only 3%, 12%, and 8% respectively. This fact indicates that TWP is ineffective.

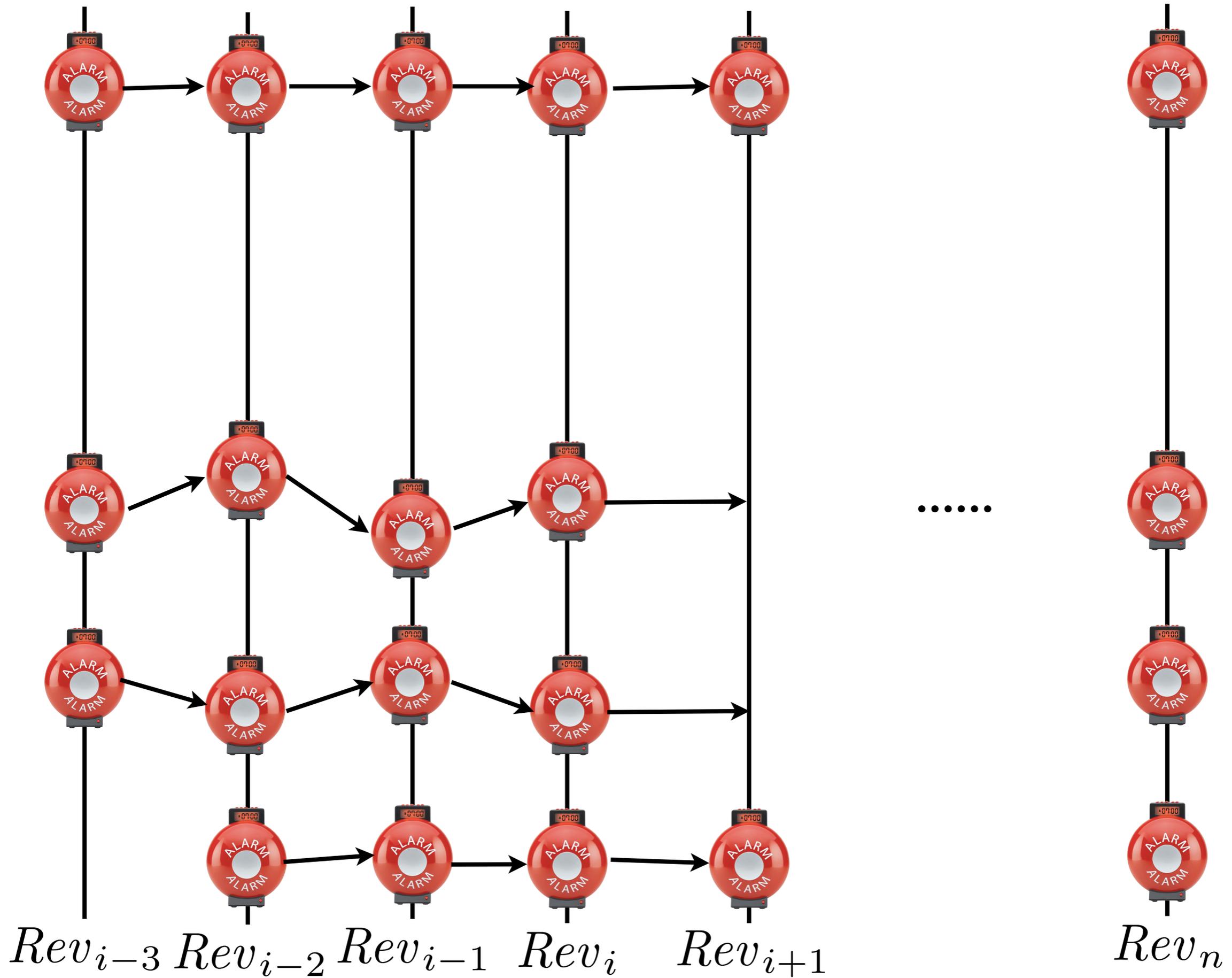
Our goal is to propose a new, program-specific prioritization that more effectively directs developers to errors. The new history-based warning prioritization (HWP) is obtained by mining the software change history for removed warnings during bug fixes.

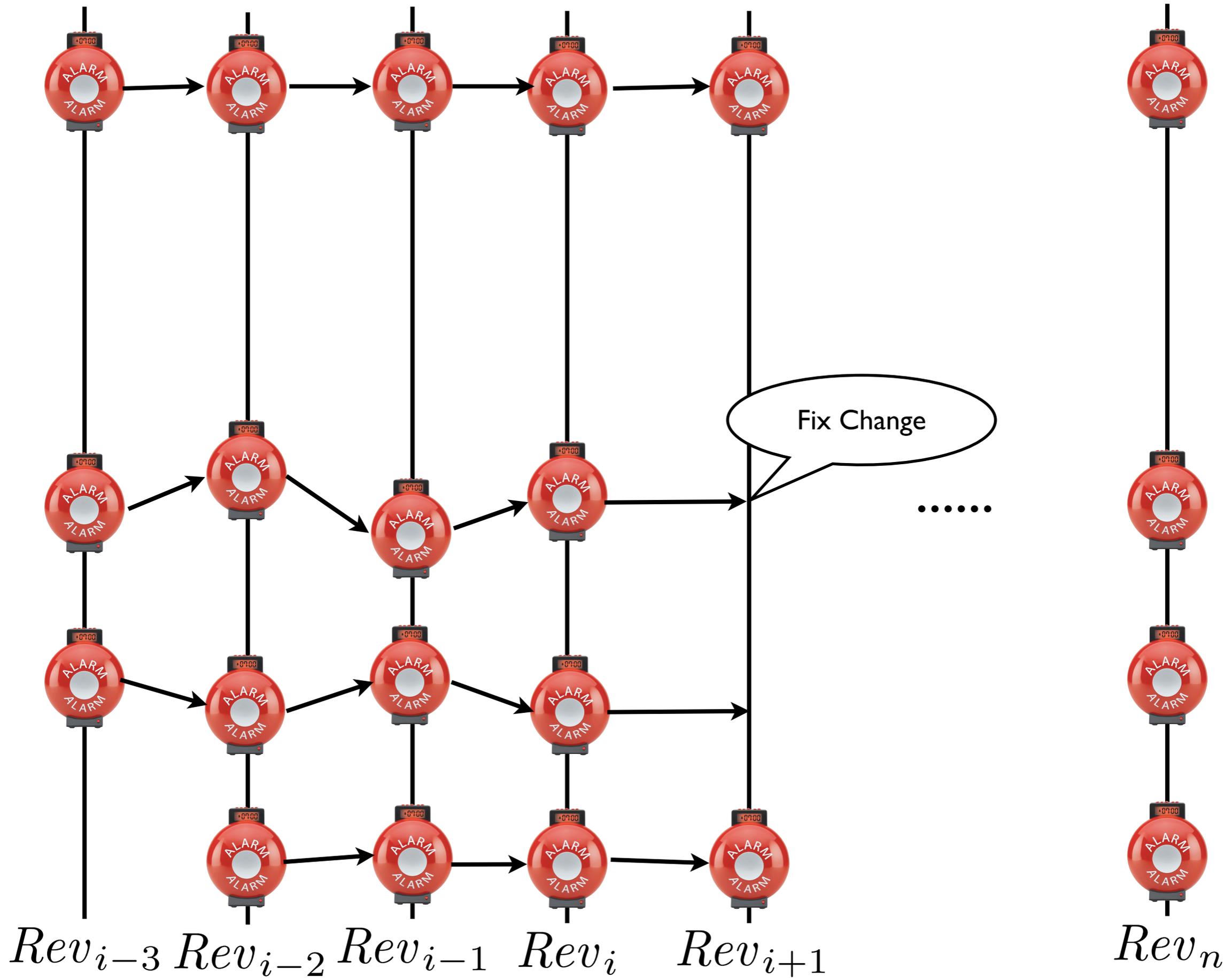
A version control system indicates when each file is changed. A software change can be classified as a fix-change or a non-fix change. A fix-change is a change that fixes a bug or other problem. A non-fix change is a change that does not fix a bug, such as a feature addition or refactoring.

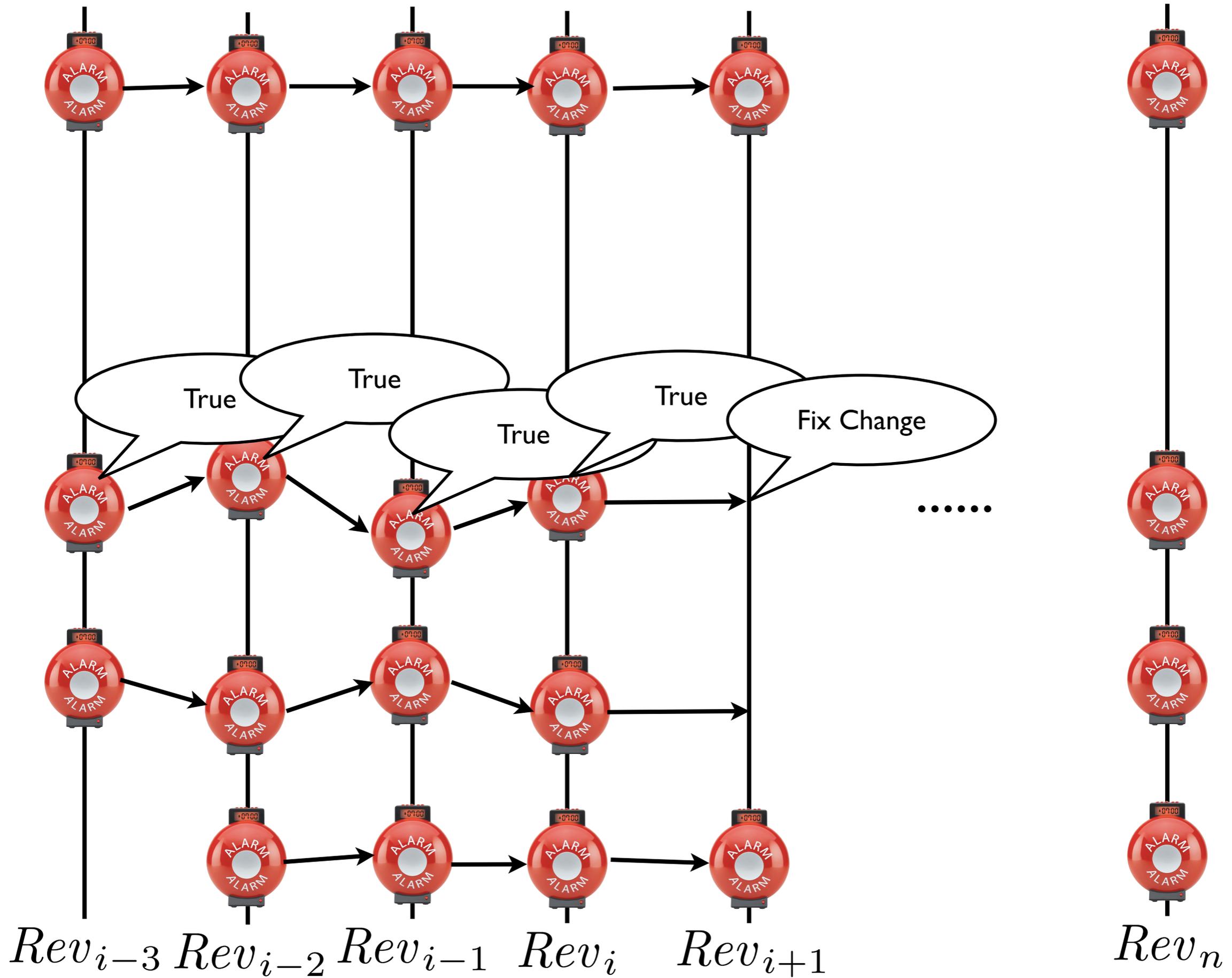
Suppose that during development, a bug-finding tool would issue a warning instance from the *Overflow* category. If a developer finds the underlying problem and fixes it, the warning is probably important. (We do not assume the software developer is necessarily using the bug-finding tool.) On the other hand, if a warning instance is not removed for a long time, then warnings of that category may be neglectable, since the problem was not noticed or was not considered worth fixing.

Using this intuition, we set a weight for each warning category to represent its importance. The weight of a category is proportional to the number of warning instances from that category that are eliminated by a change, with fix-changes contributing more to the

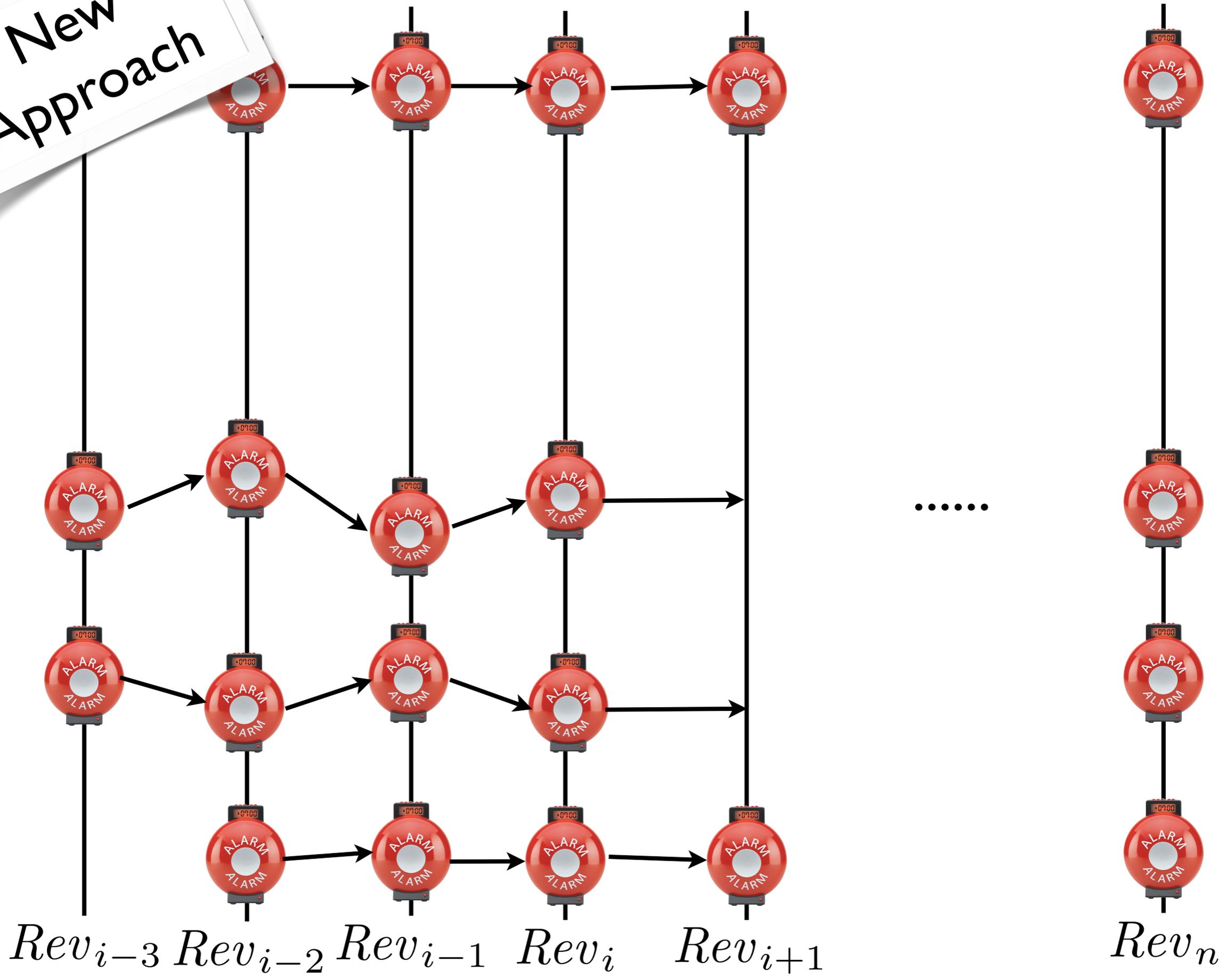
FSE'07



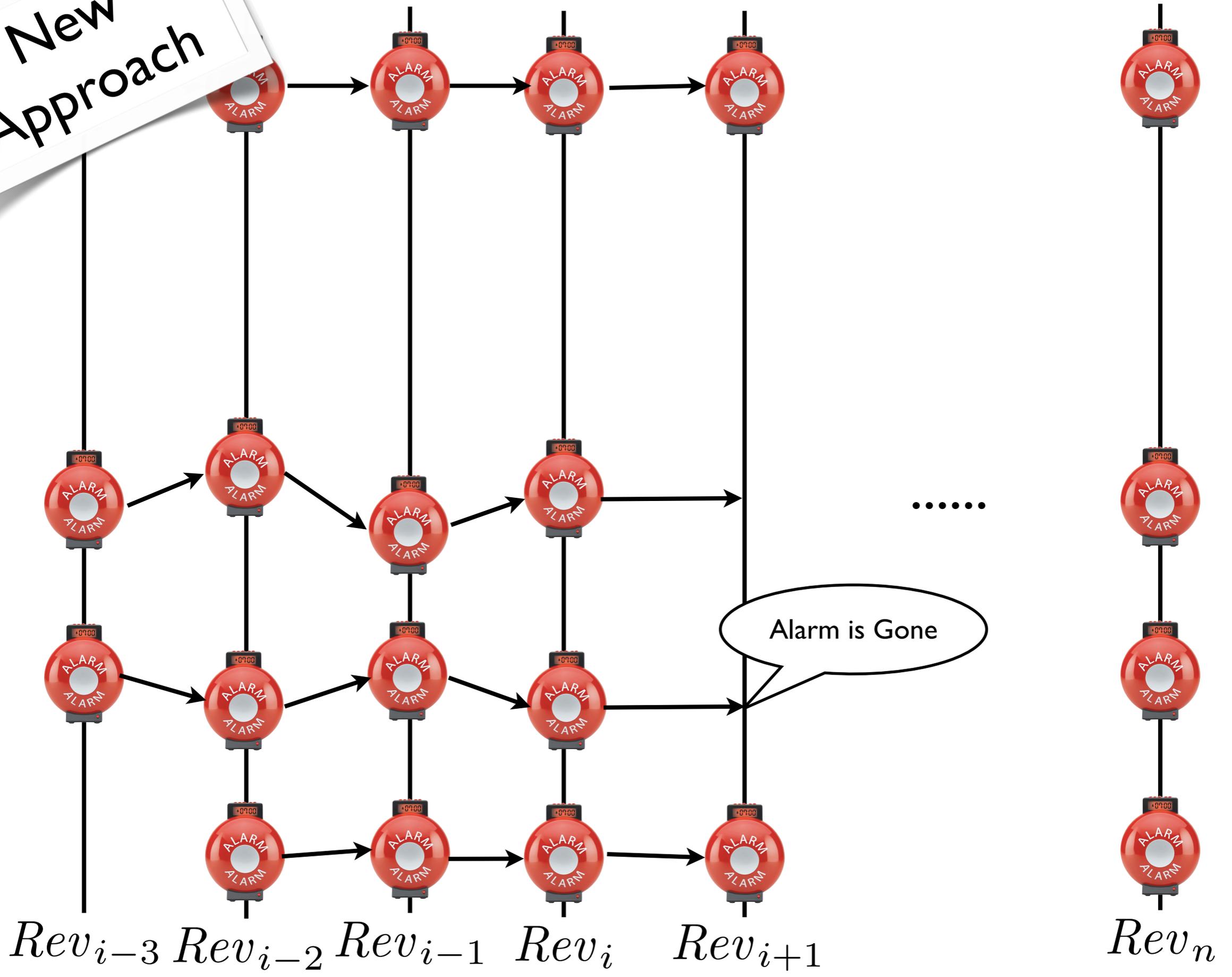




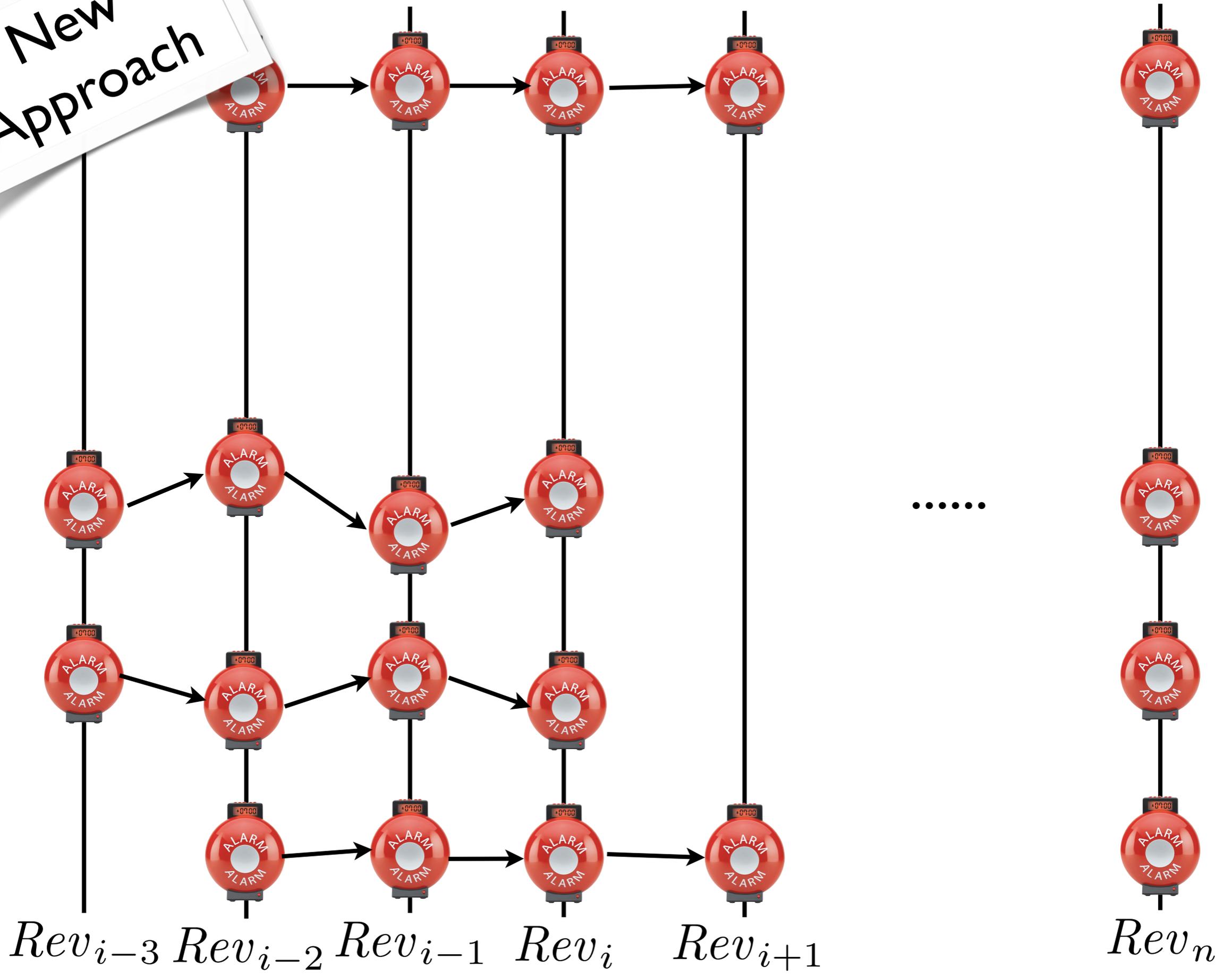
New Approach



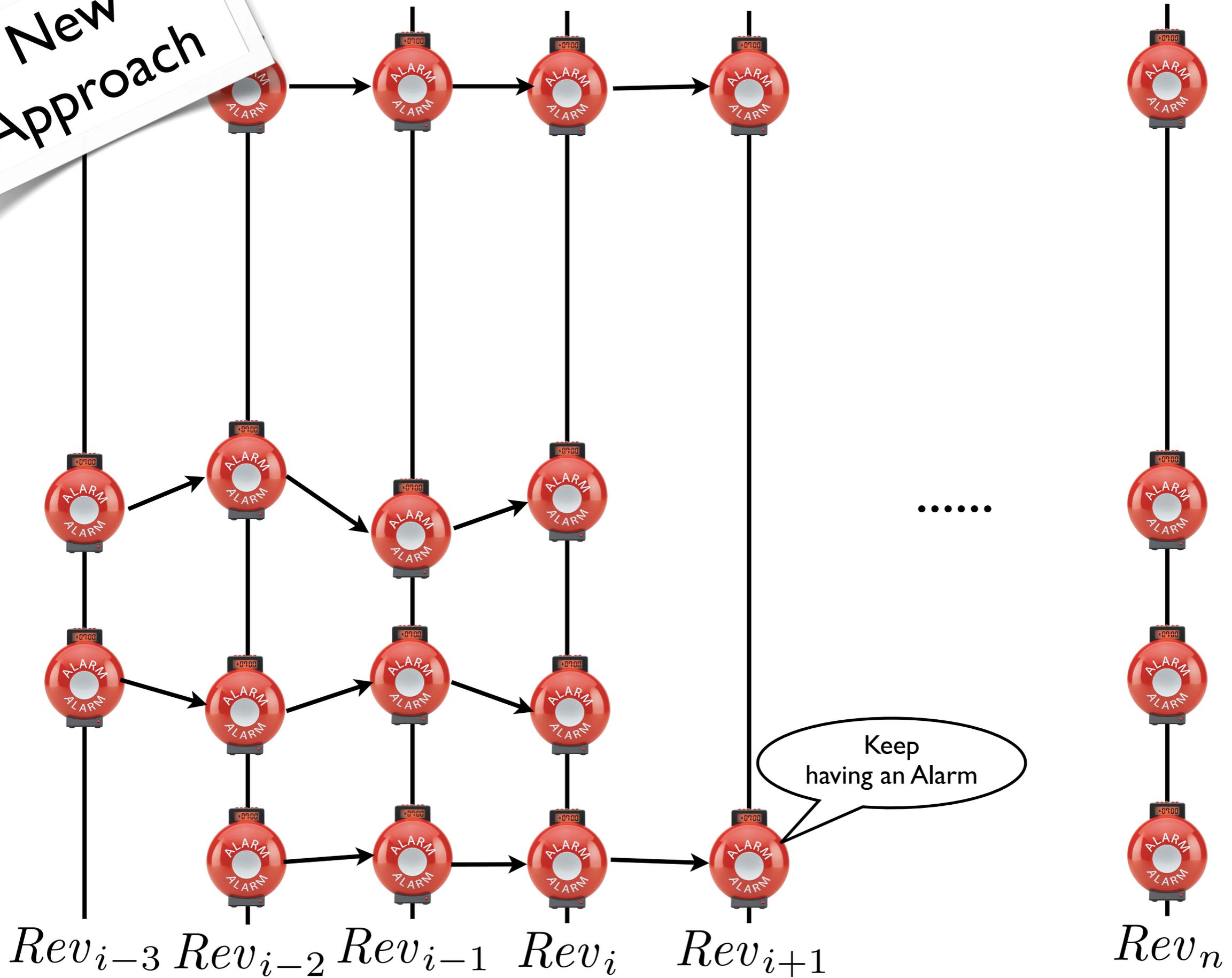
New Approach



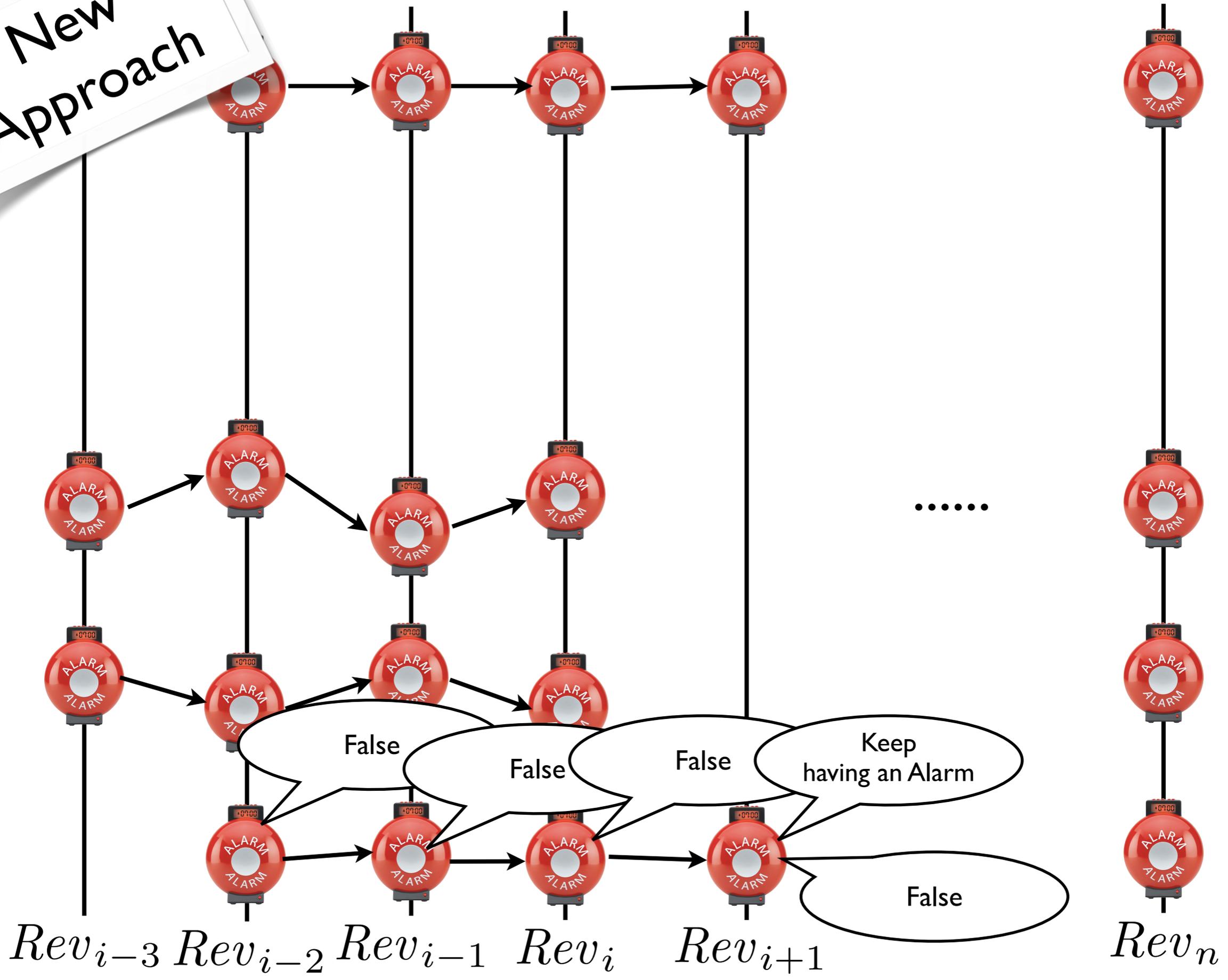
New Approach



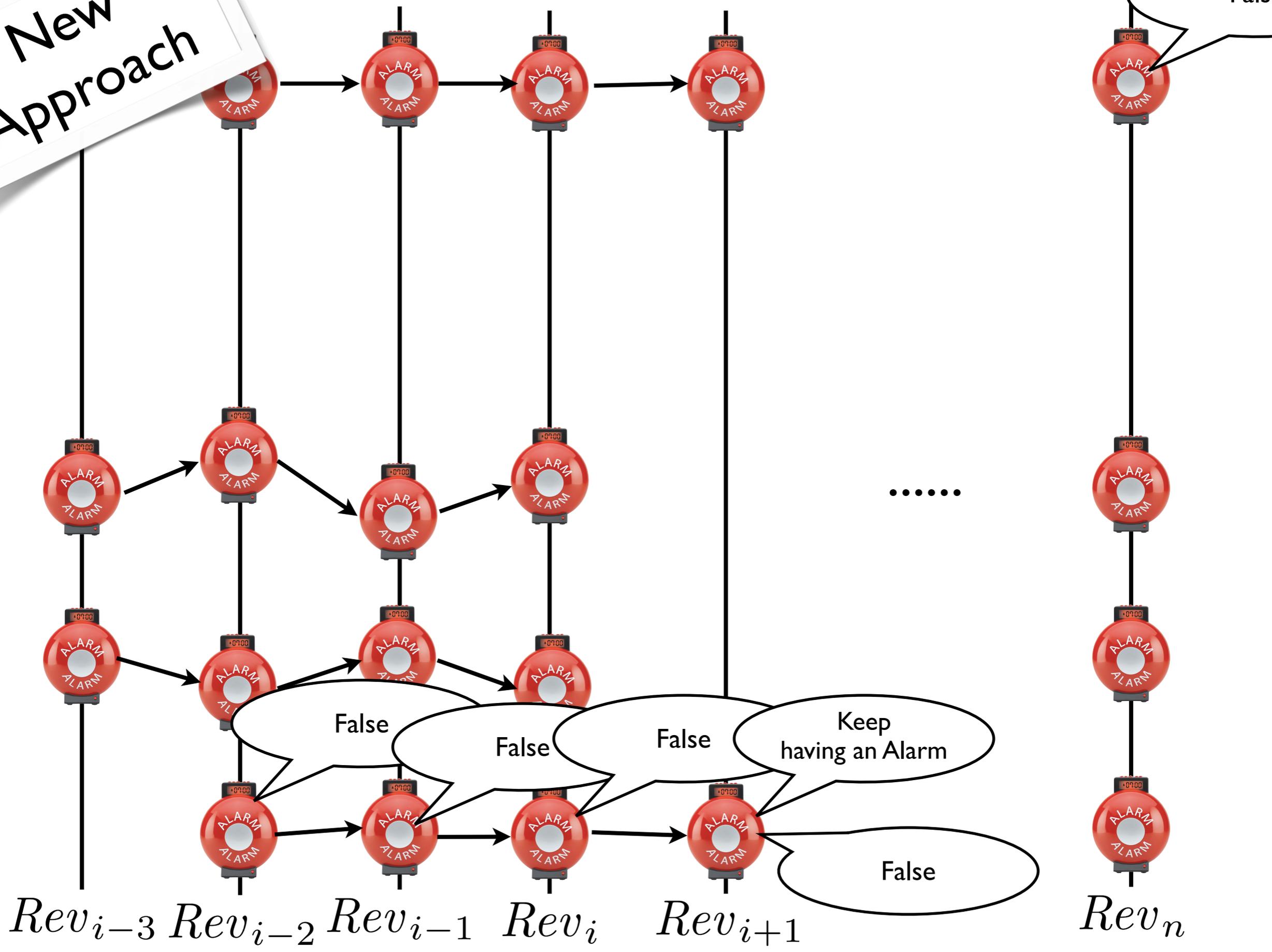
New Approach



New Approach



New Approach



False

False

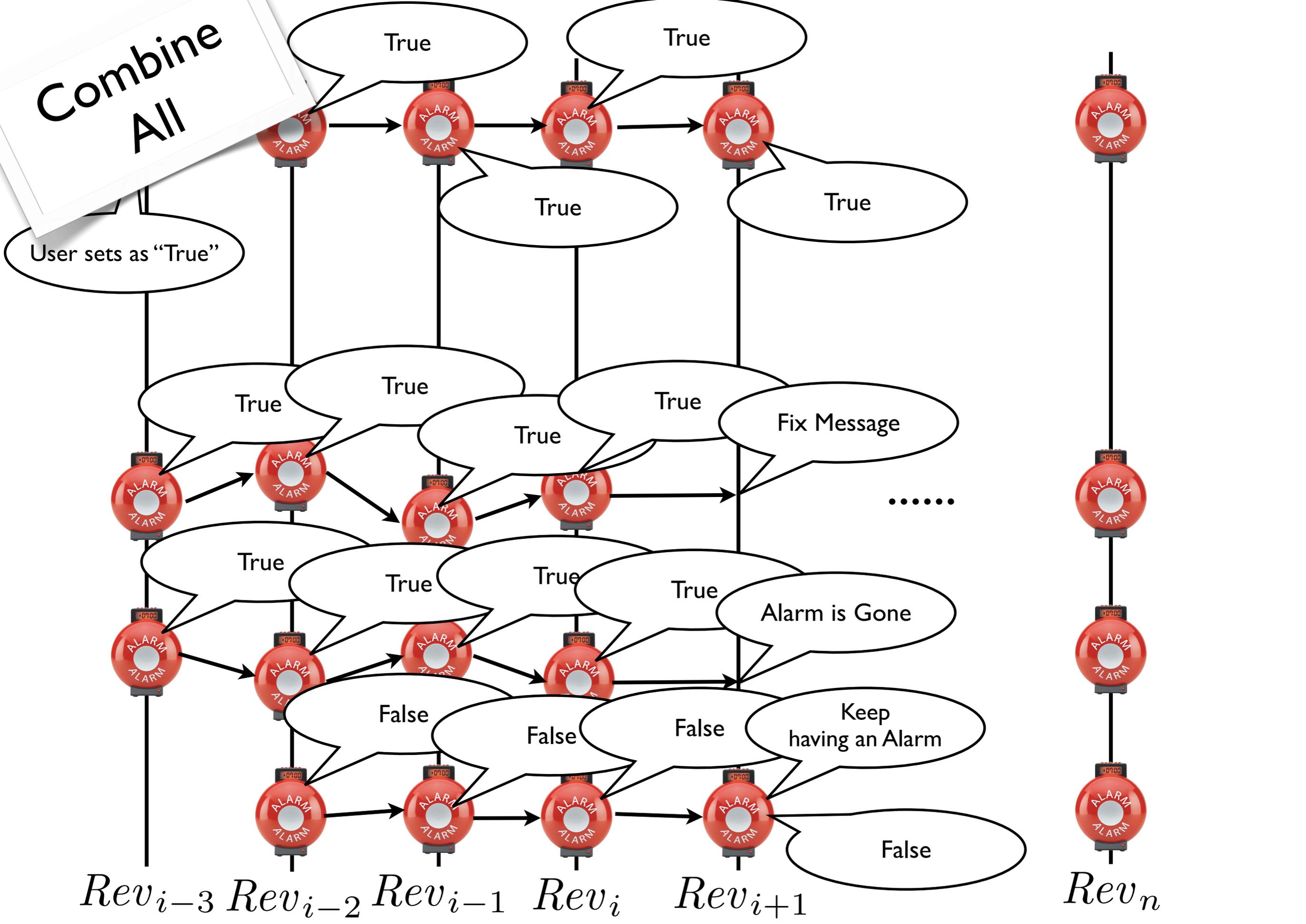
False

False

Keep having an Alarm

False

Combine All



Key Issues

- How to Determine Two Alarms are Similar
- How to Keep Track of an alarm

What I'm
doing, Now

Experiment

- Apache Httpd
 - 787,887 Revisions
 - 200K LOC
 - 194 Buffer Overrun Alarms
 - 308 Memory Leak Alarms
 - Analysis Time : 90 mins
- Subversion
 - 38,283 Revisions
 - 472K LOC
 - 726 Buffer Overrun Alarms
 - 804 Memory Leak Alarms
 - Analysis Time : 50 Hours



Thank you