

# Automated Program Repair & Verification

## Coming up

- This Thursday, Beta presentations!
- Next week:
  - Tuesday, Lecture on ethics in software engineering
  - Thursday, no class — at home, online assignment/activity
    - Graded — part of the participation grade (3% of overall class grade)
    - Can optionally opt into sharing anonymized data with researchers
- Also Thursday, Beta due!

## The Cost of Poor Software Quality in the US

Software engineers spend 35-50% of their time validating and debugging software.

Cost of debugging, testing, and verification accounts for 50-75% of the software development budgets.

- Devon H. O'Dell, The Debugging Mindset. ACM QUEUE, <http://doi.org/10.1145/3055301.3068754>

Herb Krasner, The Cost of Poor Software Quality in the US. A 2020 Report. <https://www.it-cisq.org/pdf/CPSQ-2020-report.pdf>

## Machine Learning in Software Engineering

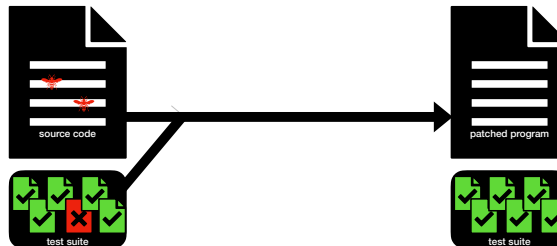
### Your AI pair programmer

GitHub Copilot uses the OpenAI GPT-3 to suggest code and entire functions in real-time, right from your editor.

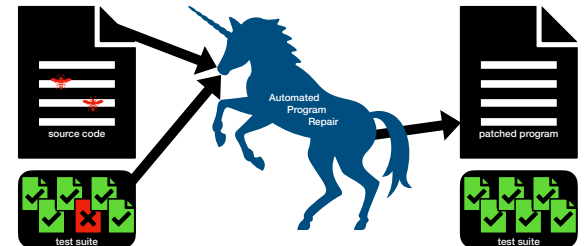
```
1 #!/usr/bin/env ts-node
2
3 import { fetch } from "fetch-h2";
4
5 // Determine whether the content of text is positive
6 // via a web service
7 export function isPositive(text: string): Promise {
8   const response = await fetch("http://localhost:3000/api/assessments", {
9     method: "POST",
10    body: JSON.stringify({
11      content: text,
12      "content-type": "application/json"
13    })
14  });
15   const json = await response.json();
16   return json.is_pos;
17 }
```

<https://github.com/features/copilot>

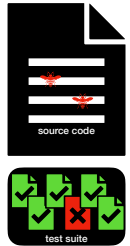
## A simpler problem: Automated program repair



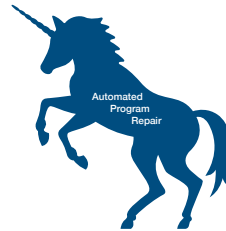
## A simpler problem: Automated program repair



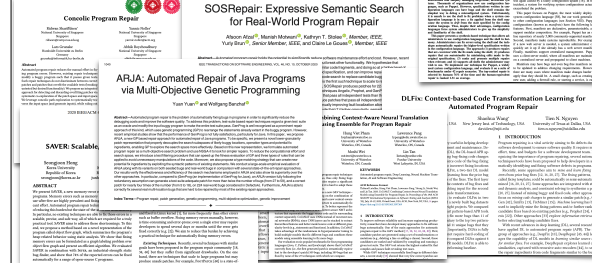
## Program repair techniques



- Tweak the program
- Check if tests pass
- If not, repeat

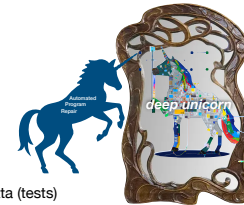


## Program repair techniques



## APR is a form of machine learning

- first, many techniques rely on ML to learn
  - where to edit the code
  - how to edit the code
  - how to decide which patches are good
- second, the underlying problem is learning a function (program) using training data (tests)



## How well does APR work?

- Evaluated 4 techniques
  - GenProg
  - Par
  - TrpAutoRepair
  - SimFix
- Measured patch quality
- Measured what affects patch quality



## Quality vs. quantity

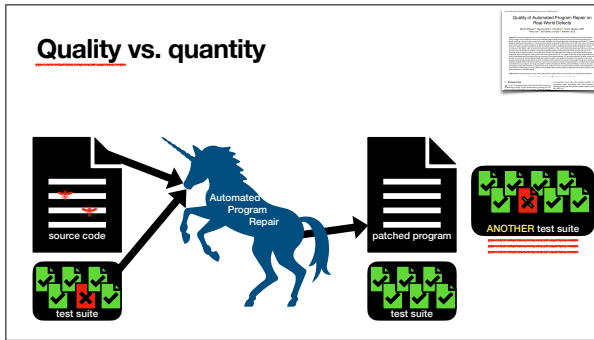
technique	defects patched
GenProg	49 (13.7%)
Par	38 (10.6%)
SimFix	68 (19.0%)
TRPAutoRepair	44 (12.3%)
total	106 (29.7%)

When applied to real-world Java code, APR produces patches for 10.6-19.0% of the defects

## Quality vs. quantity

### Potential problem: Overfitting

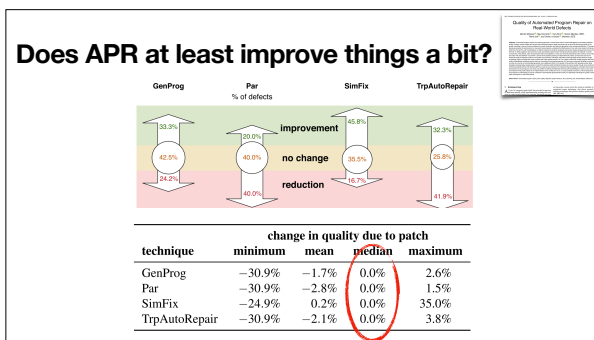
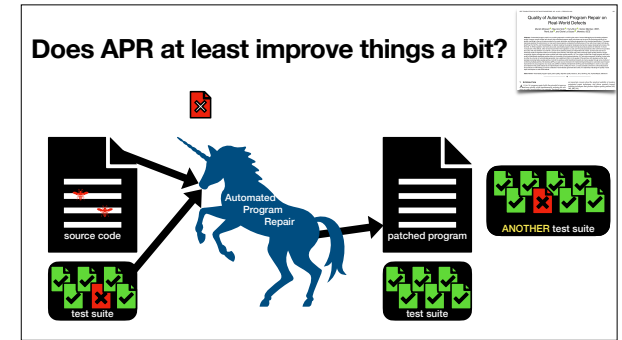
APR uses a set of tests to guide repair. Tests are inherently partial. No way APR can know if a patch captures intended behavioral constraints.



## Quality vs. quantity

technique	minimum	patch quality mean	patch quality median	maximum	100%-quality patches
GenProg	64.8%	95.7%	98.4%	98.4%	24.3%
Par	64.8%	96.1%	98.5%	100.0%	13.8%
SimFix	65.0%	96.3%	99.9%	100.0%	46.1%
TrpAutoRepair	64.8%	96.4%	98.4%	100.0%	19.5%

Less than half (14-46%) of the patches are correct



### Is the Cure Worse Than the Disease? Overfitting in Automated Program Repair

Edward K. Smith<sup>1</sup>, Earl T. Barr<sup>2</sup>, Claire Le Gouar<sup>3</sup>, Yuriy Brun<sup>4</sup>  
<sup>1</sup>University of Massachusetts – University College, London <sup>2</sup>Carnegie Mellon University, Pittsburgh, PA, USA  
<sup>3</sup>INRIA, Inria-France, France <sup>4</sup>University of California, Berkeley, CA, USA

**ABSTRACT**

Automated program repair has shown promise for reducing the significant overhead of patching code. However, we show that APR can sometimes overfit to the test suite, producing patches that pass the test suite but do not address the underlying defect. We show that APR can sometimes overfit to the test suite, producing patches that pass the test suite but do not address the underlying defect. We show that APR can sometimes overfit to the test suite, producing patches that pass the test suite but do not address the underlying defect.

Takeaway: Tests are an imperfect oracle, so APR suffers, producing low-quality patches.

Takeaway: Tests are an imperfect oracle, so APR suffers, producing low-quality patches.

Can we find a domain with better oracles?

## Formal verification allows proving software correct



## Interactive theorem provers for formal verification

Formal verification comes with a built-in oracle:  
The theorem prover



## Industrial impact of theorem proving



## Prohibitively difficult

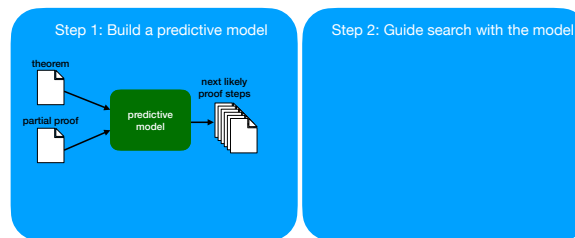
Verified software requires a lot of time and a lot of proofs in proportion to code

Proof is about 8 times bigger than the compiler code  
3 person years of work

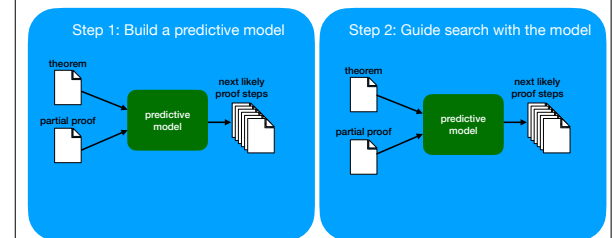
POPL 2006

Virtually all software that ships today is unverified.

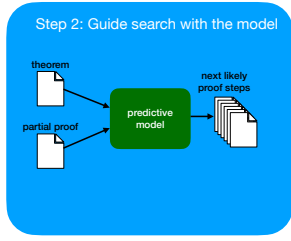
## Proposal: Use APR-style technology to synthesize proofs



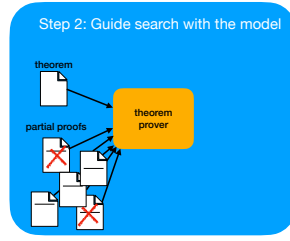
## Proposal: Use APR-style technology to synthesize proofs



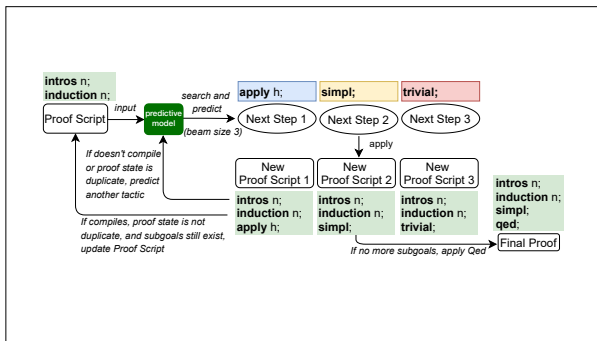
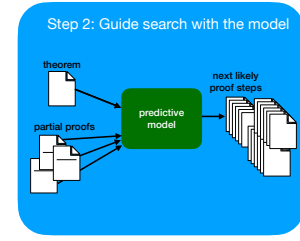
**Proposal: Use APR-style technology to synthesize proofs**



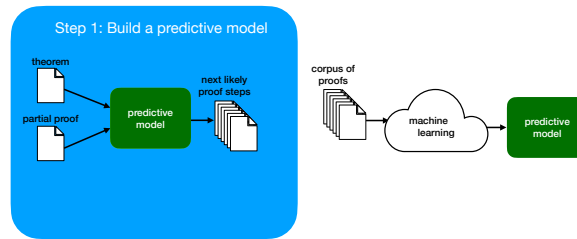
**Proposal: Use APR-style technology to synthesize proofs**



**Proposal: Use APR-style technology to synthesize proofs**



**How to learn a predictive model**



**TacTok (OOPSLA'20)**

**TacTok Semantics-Aware Proof Synthesis**

EMILY FIRSH, University of Massachusetts Amherst, USA  
 YURIY BRUN, University of Massachusetts Amherst, USA  
 ARDEN CHIAI, University of Massachusetts Amherst, USA

Formally verifying software correctness is a highly manual process. However, because verification proof scripts often share structure, it is possible to learn from existing proof scripts to fully automate some formal verification. The goal of this paper is to improve proof script synthesis and enable fully automating some verification. Because theorem provers work in the Coq proof assistant, their programmers to write partial proof scripts, observe the semantics of the proof state that is, and then attempt more progress. Knowing the

**TacTok models partial proof and the current proof state, together**

Training Proofs    Training Instances    AST

Early Proof, Yuriy Brun, and Arden Chiai. 2020. TacTok: Semantics-Aware Proof Synthesis. Proc. ACM Program. Lang. 4(OOPSLA, Article 21), December 2020. 11 pages. <https://doi.org/10.1145/3432029>

ASTatic (Yang and Deng, Learning to Prove Theorems via Interacting with Proof Assistants, IJML'19) modeled just proof state, Hellendoorn, Devanbu, Alipour. On the naturalness of proofs, ESEC/FSE NIER'18) looked at predictability of proof sequences.

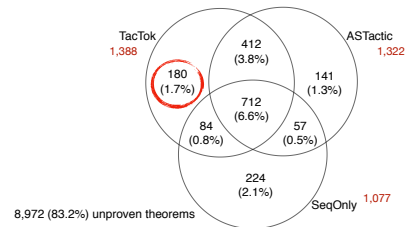
## CoqGym Dataset

- 123 open-source software projects in Coq
- 70,856 theorems
- Broken down into 96 projects (57,719 proofs) for training and 27 projects (13,137 theorems) for testing

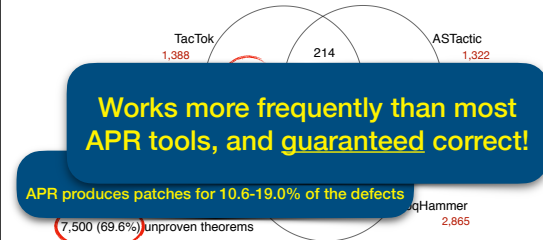
<https://github.com/princeton-vl/CoqGym>

[Yang and Deng, Learning to Prove Theorems via Interacting with Proof Assistants, ICML'19]

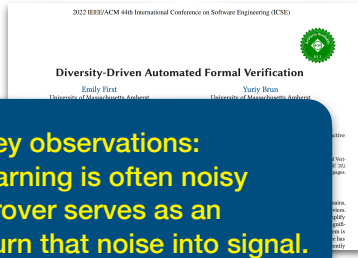
## TacTok vs. ASTactic vs. SeqOnly



## TacTok vs. ASTactic vs. CoqHammer



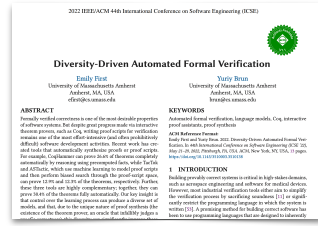
## Diva (ICSE'22)



**2 key observations:**

- Machine learning is often noisy
- Theorem prover serves as an oracle to turn that noise into signal.

## Diva (ICSE'22)



- Vary:
  - proof tactic and token depth
  - learning rate
  - embedding size
  - number of layers
  - training order
  - access to proof state, partial proof, Gallina proof term

## Diva vs. state-of-the-art

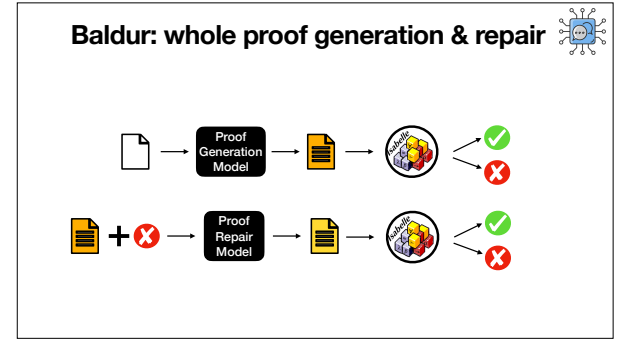
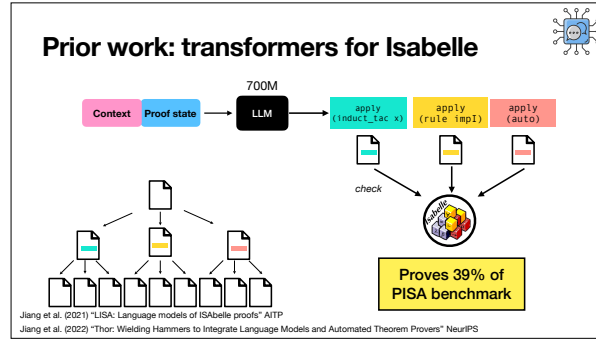
**Diversity inherent in ML increases the proving power 68%-77% over prior search-based synthesis tools, and 27% over CoqHammer.**

<https://github.com/LASER-UMASS/Diva/>

### ChatGPT

Examples	Capabilities	Limitations
"Explain quantum computing in simple terms" →	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?" →	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?" →	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021

Explain formal verification to me



### Evaluation

**Baldur: Whole-Proof Generation and Repair with Large Language Models**

Ioana Firoiu  
University of Pennsylvania  
firoiuf@cis.upenn.edu

Markus N. Rabe  
Google  
mrabe@google.com

Talita Wagner  
University of Black Mountains  
twagner@ubm.ac.uk

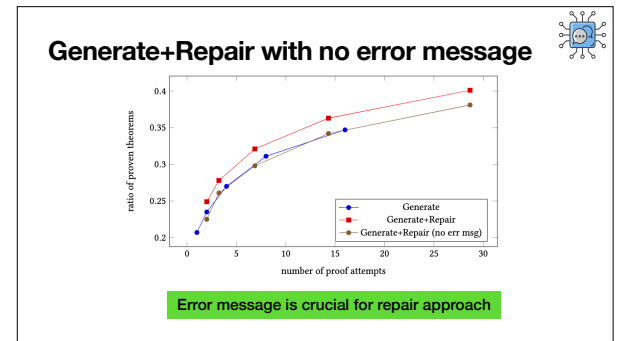
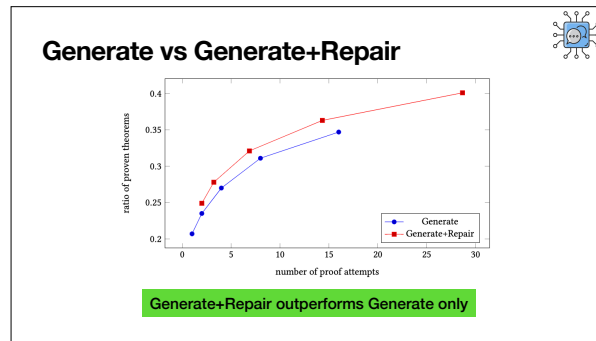
Tarun Brahm  
University of Pennsylvania  
brahm@cis.upenn.edu

**Abstract**  
 Formal verification software packages are highly desirable for their ability to ensure correctness of program code. However, they are often difficult to use, and their use is limited to a small number of domains. In this paper, we propose a new approach to formal verification that leverages large language models (LLMs) to generate and repair proofs. We evaluate our approach on the Isabelle/HOL theorem prover, and show that it can generate and repair proofs for a wide range of problems. Our approach is based on a novel architecture that combines an LLM with a theorem prover. We show that our approach can generate and repair proofs for a wide range of problems, and that it can be used to generate and repair proofs for a wide range of domains. Our approach is based on a novel architecture that combines an LLM with a theorem prover. We show that our approach can generate and repair proofs for a wide range of problems, and that it can be used to generate and repair proofs for a wide range of domains.

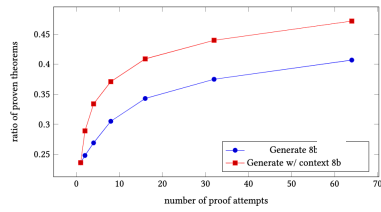
**Introduction**  
 Formal verification is a technique for ensuring the correctness of program code. It is a highly desirable property for many systems, but it is often difficult to use. In this paper, we propose a new approach to formal verification that leverages large language models (LLMs) to generate and repair proofs. We evaluate our approach on the Isabelle/HOL theorem prover, and show that it can generate and repair proofs for a wide range of problems. Our approach is based on a novel architecture that combines an LLM with a theorem prover. We show that our approach can generate and repair proofs for a wide range of problems, and that it can be used to generate and repair proofs for a wide range of domains.

**Baldur + Thor prove nearly 2/3 of PISA test set!**

Jiang et al. (2021) "LISA: Language models of Isabelle proofs" AITP  
 Jiang et al. (2022) "Thor: Welding Hammers to Integrate Language Models and Automated Theorem Provers" NeurIPS 2022



## Generate with context



Proof context helps improve proof generation

## Fully Automated Formal Verification

Machine learning and meta-heuristic search  
can fully automate  
some bug-repair and formal verification.

While APR underperforms because it is driven by an unreliable oracle,  
formal verification is a killer app for APR because  
the theorem prover provides a reliable oracle.