# The promise and perils of using machine learning when engineering software.
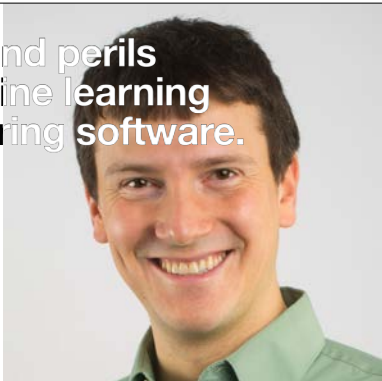
Yuriy Brun

MaLTeSQuE 2022

---

## Machine learning today

### Forbes
How Artificial Intelligence And Machine Learning Are Transforming The Future Of Renewable Energy

3 Ways Artificial Intelligence Automatically Improves Your Pictures

---

### Resilient cities  Cities
## Predicting crime, LAPD-style

Cutting edge data-driven analysis directs Los Angeles patrol officers to likely future crime scenes - but critics worry that decision-making by machine will bring 'tyranny of the algorithm'

● Join our live Q&A with Homicide Watch this Friday

https://www.theguardian.com/cities/2014/jun/25/predicting-crime-lapd-los-angeles-police-data-analysis-algorithm-minority-report

---

### ACLU
## The Government Is Blacklisting People Based on Predictions of Future Crimes

By Hina Shamsi, Director, ACLU National Security Project

**Modern software uses machine learning to influence critical decisions**

https://www.aclu.org/blog/national-security/discriminatory-profiling/government-blacklisting-people-based-predictions

---

### Forbes
The Algorithm That Beats Your Bank Manager

HAAS NEWS > NEWS CATEGORIES > RESEARCH NEWS

## Minority homebuyers face widespread statistical lending discrimination, study finds

By Laura Counts | NOVEMBER 13, 2018

Face-to-face meetings between mortgage officers and homebuyers have been rapidly replaced by online applications and algorithms, but lending discrimination hasn't gone away.

A new University of California, Berkeley study has found that both online and face-to-face lenders charge higher interest rates to African American and Latino borrowers, earning 11 to 17 percent higher profits on such loans. All told, those homebuyers pay up to half a billion dollars more in interest every year than white borrowers with comparable credit scores do, researchers found.

The findings raise legal questions about the rise of statistical discrimination in the fintech era, and point to potentially widespread violations of U.S. fair lending laws, the researchers say. While lending discrimination has historically been caused by human prejudice, pricing disparities are increasingly the result of algorithms that use machine learning to target applicants who might shop around less for higher-priced loans.

"The mode of lending discrimination has shifted from human bias to algorithmic bias," said study co-author Adair Morse, a finance professor at UC Berkeley's Haas School of Business. "Even if the people writing the

---
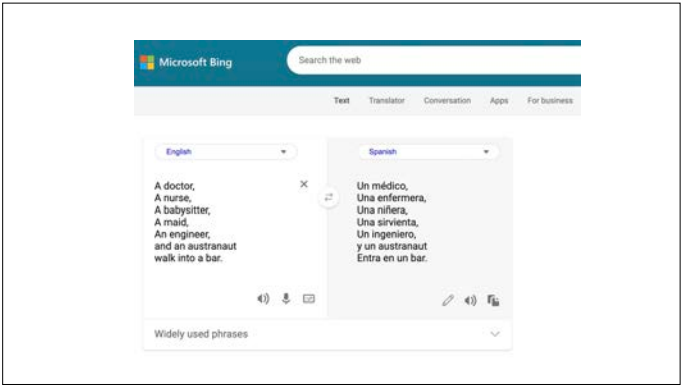
### HEALTH IT ANALYTICS
Machine-Learning Mode... Detect Early-Stage Canc...

A new study suggests that machine-learni... occult nodal metastasis in patients with a... cavity cancer with more accuracy than sta...

### MIT News
ON CAMPUS AND AROUND THE WORLD

## Using AI to predict breast cancer and personalize care

MIT/MGH's image-based deep learning model can predict breast cancer up to five years in advance.

Adam Conner-Simons and Rachel Gordon | CSAIL
May 7, 2019

**Software can make bad decisions. Software can discriminate!**

Machine learning has great promise,
but with that promise, come risks.

Today's goal:
Identifying and addressing the risks

Part I
Automated program repair

Part II
Software discrimination
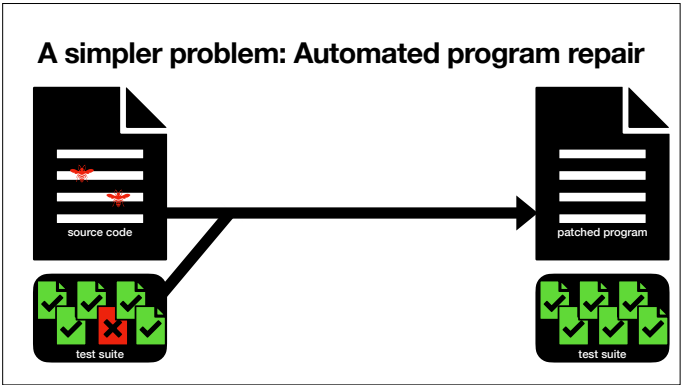
## Machine Learning in Software Engineering



https://github.com/features/copilot

## A simpler problem: Automated program repair



source code → patched program

test suite

## A simpler problem: Automated program repair

source code

test suite

Automated Program Repair

patched program

test suite

## Program repair techniques

source code

test suite

Automated Program Repair

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

## Program repair techniques

Concolic Program Repair

SOSRepair: Expressive Semantic Search for Real-World Program Repair

Alison Afzal, Manish Motwani, Kathryn T. Stolee, Member, IEEE, Yuriy Brun, Senior Member, IEEE, and Claire Le Goues, Member, IEEE

ARJA: Automated Repair of Java Programs via Multi-Objective Genetic Programming

Yuan Yuan and Wolfgang Banzhaf

SAVER: Scalable

Tortoise: Interactive System Configuration Repair

DLFix: Context-based Code Transformation Learning for Automated Program Repair

Mining Context-Aware Neural Translation

Ensemble for Program Repair

## 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)

Automated Program Repair

deep unicorn

## How well does APR work?

Quality of Automated Program Repair on Real-World Defects

Manish Motwani, Mauricio Soto, Yuriy Brun, Senior Member, IEEE, René Just, and Claire Le Goues, Member, IEEE

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

- Measured what affects patch quality

## Quality vs. quantity

Quality of Automated Program Repair on Real-World Defects

| technique | | | |
|---|---|---|---|
| GenProg | | | |
| Par | | | |
| SimFix | | | |
| TRPAutoRepair | | | |
| total | | | |

ESEC/FSE'19

Empirical Review of Java Program Repair Tools: A Large-Scale Experiment on 2,141 Bugs and 23,551 Repair Attempts

Thomas Durieux
INESC-ID and IST, University of Lisbon, Portugal
thomas@durieux.me

Fernanda Madeiral
Federal University of Uberlândia, Brazil
fer.madeiral@gmail.com

Matias Martinez
Université Polytechnique Hauts-de-France, France
matias.martinez@uphf.fr

Rui Abreu
INESC-ID and IST, University of Lisbon, Portugal
rui@computer.org

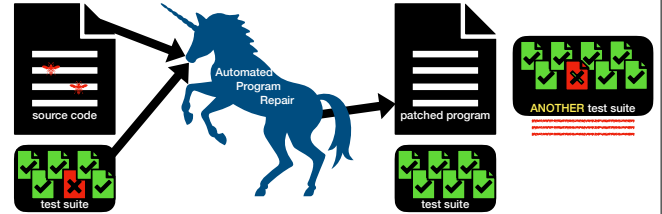When applied to rea
APR produces patches for

## Slide 1: Quality vs. quantity

**Potential problem: Overfitting**

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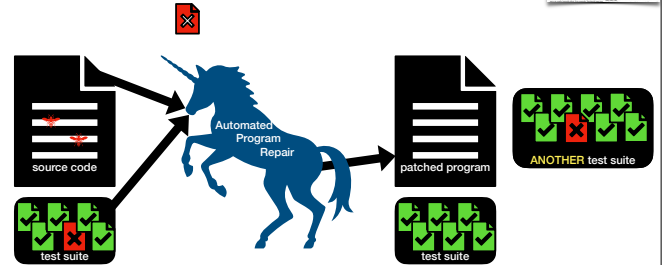No way APR can know if a patch captures intended behavioral constraints.

## Slide 2: Quality vs. quantity



source code → Automated Program Repair → patched program
test suite
ANOTHER test suite
test suite

## Slide 3: Quality vs. quantity

| technique | minimum | patch quality mean | median | maximum | 100%-quality patches |
|---|---|---|---|---|---|
| GenProg | 64.8% | 95.7% | 98.4% | 100.0% | 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**

## Slide 4: Does APR at least improve things a bit?



source code → Automated Program Repair → patched program
test suite
ANOTHER test suite
test suite

## Slide 5: Does APR at least improve things a bit?

GenProg, Par (% of defects), SimFix, TrpAutoRepair
improvement
no change
reduction

| technique | change in quality due to patch minimum | mean | median | maximum |
|---|---|---|---|---|
| GenProg | −30.9% | −1.7% | 0.0% | 2.6% |
| Par | −30.9% | −2.8% | 0.0% | 1.5% |
| SimFix | −24.9% | 0.2% | 0.0% | 35.0% |
| TrpAutoRepair | −30.9% | −2.1% | 0.0% | 3.8% |

## Slide 6

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

Edward K. Smith    Earl T. Barr    Claire Le Goues    Yuriy Brun
University of Massachusetts · University College, London · Carnegie Mellon University
Amherst, MA, USA    London, UK    Pittsburgh, PA, USA
{tedks, brun}@cs.umass.edu, e.barr@ucl.ac.uk, clegoues@cs.cmu.edu

Smith, et al., Is the Cure Worse That the Disease?
Overfitting in Automated Program Repair, ESEC/FSE 2015.

**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

Proof script

```
Proof.
intros.
induction n.
.
.
Qed.
```

Agda

---

**Industrial impact of theorem proving**

AIRBUS

BEDROCK
Systems Inc

galois

aws    CERTORA

android

---

**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**

Step 1: Build a predictive model

theorem

partial proof

predictive model

next likely proof steps

Step 2: Guide search with the model

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

Step 1: Build a predictive model

theorem → predictive model → next likely proof steps
partial proof →

Step 2: Guide search with the model

theorem → predictive model → next likely proof steps
partial proof →

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

Step 2: Guide search with the model

theorem → predictive model → next likely proof steps
partial proof →

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

Step 2: Guide search with the model

theorem → theorem prover
partial proofs →

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

Step 2: Guide search with the model

theorem → predictive model → next likely proof steps
partial proofs →

**intros** n;
**induction** n;

Proof Script → input → predictive model → search and predict (beam size 3) → Next Step 1 / Next Step 2 / Next Step 3

**apply** h; / **simpl;** / **trivial;**

apply

New Proof Script 1 / New Proof Script 2 / New Proof Script 3

*If doesn't compile or proof state is duplicate, predict another tactic*

*If compiles, proof state is not duplicate, and subgoals still exist, update Proof Script*

**intros** n; **induction** n; **apply** h;

**intros** n; **induction** n; **simpl;**

**intros** n; **induction** n; **trivial;**

**intros** n; **induction** n; **simpl; qed;**

*If no more subgoals, apply Qed* → Final Proof

# How to learn a predictive model

Step 1: Build a predictive model

theorem → predictive model → next likely proof steps
partial proof →

corpus of proofs → machine learning → predictive model

## Slide 1: TacTok (OOPSLA'20)

**TacTok (OOPSLA'20)**

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

Training Proofs

Training Instances

Tactic

AST

proof

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

## Slide 2: CoqGym Dataset

**CoqGym Dataset**

- 123 open-source software projets 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]

## Slide 3: TacTok vs. ASTactic vs. SeqOnly

**TacTok vs. ASTactic vs. SeqOnly**

TacTok
1,388

ASTactic
1,322

412 (3.8%)

141 (1.3%)

180 (1.7%)

712 (6.6%)

84 (0.8%)

57 (0.5%)

224 (2.1%)

SeqOnly 1,077

8,972 (83.2%) unproven theorems

## Slide 4: TacTok vs. ASTactic vs. CoqHammer

**TacTok vs. ASTactic vs. CoqHammer**

TacTok
1,388

ASTactic
1,322

214

**Works more frequently than most APR tools, and <u>guaranteed</u> correct!**

APR produces patches for 10.6-19.0% of the defects

CoqHammer
2,865

7,500 (69.6%) unproven theorems

## Slide 5: Diva (ICSE'22)

**Diva (ICSE'22)**

2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)

Diversity-Driven Automated Formal Verification

Emily First
University of Massachusetts Amherst

Yuriy Brun
University of Massachusetts Amherst

**2 key observations:**
- **Machine learning is often noisy**
- **Theorem prover serves as an oracle to turn that noise into signal.**

## Slide 6: Diva (ICSE'22)

**Diva (ICSE'22)**

2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)

Diversity-Driven Automated Formal Verification

Emily First
University of Massachusetts Amherst
Amherst, MA, USA
efirst@cs.umass.edu

Yuriy Brun
University of Massachusetts Amherst
Amherst, MA, USA
brun@cs.umass.edu

- 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/

## 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.**

---

…let's talk about a different peril of machine learning that verification might help with.

**Part II
Software discrimination**

---

**Part II
Software discrimination**

Data-driven systems can exhibit undesirable properties.

Can we build systems to be safe and fair?

---

Testing systems for bias

---

## Themis
automated test-suite generator

**How much does my software discriminate with respect to …?**

**Does my software discriminate more than 10% of the time, and against what?**

http://fairness.cs.umass.edu

Galhotra, Brun, and Meliou, Fairness Testing: Testing Software for Discrimination. ESEC/FSE 2017

## Debugging bias

## Debugging



## fairkit-learn



**Fairkit, Fairkit, on the Wall, Who's the Fairest of Them All?
Supporting Data Scientists in Training Fair Models**

Brittany Johnson, Jesse Bartola, Rico Angell, Katherine Keith,
Sam Witty, Stephen J. Giguere, Yuriy Brun

University of Massachusetts Amherst, USA

## Can we verify systems to be safe and fair?

## How would that work?

**User specifies a definition of safe or fair behavior.**

| training | testing | safety |

**Train classifiers,
selects one to satisfy fairness,
verify safety on held-out suite.**

Example scenario:

Suppose a university wants to train a model to predict student success from entrance exam scores, while ensuring the model is fair:
roughly the same fraction of men and women are predicted to be successful.
(This is called Disparate Impact.)

## Slide 1: Disparate Impact

# Disparate Impact



Seldonian | Standard | Fairlearn | Fairness Constraints

Fairlearn: Agarwal et al. A reductions approach to fair classification. ICML 2018.
Fairness Constraints: Zafar et al., Fairness Constraints: A Mechanism for Fair Classification. FATML 2015.

## Slide 2: Disparate Impact

# Disparate Impact



Seldonian | Standard | Fairlearn | Fairness Constraints

Fairlearn: Agarwal et al. A reductions approach to fair classification. ICML 2018.
Fairness Constraints: Zafar et al., Fairness Constraints: A Mechanism for Fair Classification. FATML 2015.

## Slide 3: Equalized Odds

# Equalized Odds



Thomas, Castro da Silva, Barto, Giguere, Brun, and Brunskill.
"Preventing Undesirable Behavior of Intelligent Machines", Science 366 (6468), Nov 22, 2019

## Slide 4

# And this approach is very versatile: …works for policy selection



Metevier, Giguere, Brockman, Kobren, Brun, Brunskill, Thomas. Offline Contextual Bandits with High Probability Fairness Guarantees.
NeurIPS 2019.

## Slide 5

Example scenario:

One source of ML bias comes from deploying a model on data that is fundamentally different from the data the model was trained on.

## Slide 6

# What if software is deployed on data fundamentally different from training data?



Giguere, Metevier, Brun, Castro da Silva, Thomas, and Niekum.
Fairness Guarantees under Demographic Shift, ICLR 2022.

Machine learning can result in unexpected, unintended behavior.

But machine learning can be leveraged to produce verified safe and fair models, avoiding such behavior.

## Contributions