# Metamorphic Testing and Debugging of Tax Preparation Software

**Saeid Tizpaz-Niari**
Computer Science Department
University of Texas at El Paso

Verya Monjezi
Computer Science Department
University of Texas at El Paso

Morgan Wagner
Psychology Department
University of Texas at El Paso

Shiva Darian
Information Science Department
University of Colorado Boulder

Krystia Reed
Psychology Department
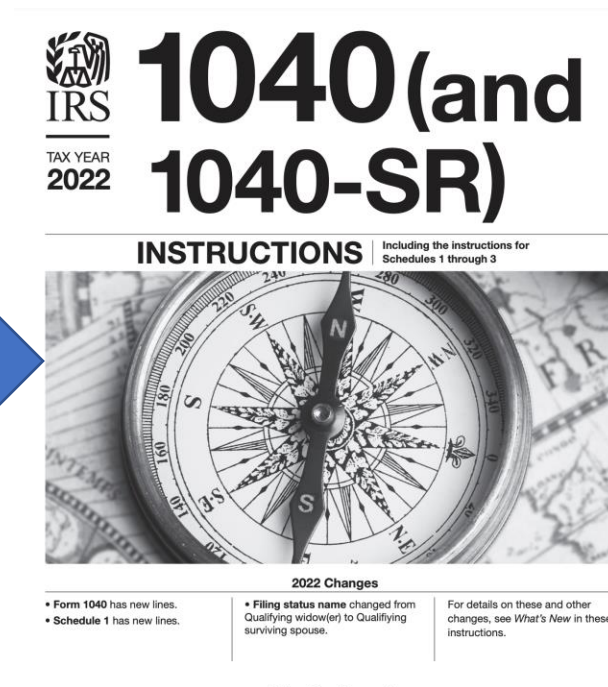University of Texas at El Paso

Ashutosh Trivedi
Computer Science Department
University of Colorado Boulder

" *Our new Constitution is now established, and has an appearance that promises permanency; but in this world nothing can be said to be certain, except death and taxes.*"

— *Benjamin Franklin, in a letter to* [Jean-Baptiste Le Roy](#)*, 1789*

# U.S. Tax 101: Manual Tax Filling



**Publication 596 (EITC)**

# Tax Preparation Software (US-based)

- 72 million tax returns via software

- 11.2 billion dollars industry

- Free (Open-source) options for low-income

**Open-source Tax Solver (OTS)**

Metamorphic Testing and Debugging
of Tax Preparation Software (ICSE-SEIS 2023)

**Langley v. Comm'r, T.C. Memo. 2013-22.** *The misuse of tax preparation software, even if unintentional or accidental, is no defense to accuracy-related penalties under section 6662.*

# Accountable Tax Software

- Comply with laws, regulations, or public policies as they evolve over time.

- Approaches for Accountability of Software

  - Formal verification to ensure compliance;

  - Methodologies for software design, development, and maintenance; and

  - Specification and reasoning about software compliance and accountability.

See solicitation for *Designing Accountable Software Systems* (DASS): https://www.nsf.gov/pubs/2022/nsf22512/nsf22512.htm

# Challenges

- **Absence of Oracle**
  - Given a taxpayer profile, the ground truth for the tax returns, eligibilities, and credits are not known a prior even for the tax experts;

- **Lack of Trustworthy Dataset**
  - Due to obvious privacy and legal concerns; and

- **Computationally difficult**
  - Finding similar tax profiles is hard (scale, notion of similarity, etc).

# Differential Debugging of Tax Software

- **Observation 1**:
  - Tax law adheres to the principles of ``common'' law;
  - It implements the legal doctrine of precedent; hence,
  - Similar cases must follow similar rulings.

- **Observation 2**:
  - *Horizontal equity* in taxation: relation between similarly situated tax-payers;
  - *Vertical equity* in taxation: relation between taxpayers in different income buckets

# Equity in Tax Domain Goes Beyond Software



**Black Americans Face More Audit Scrutiny, IRS Acknowledges**

Black taxpayers were three to five times more likely than taxpayers who are not Black to be audited, research published this year found.

May 15, 2023

Racial Bias in IRS Tax Audits

# Metamorphic Specifications

- Validation of software correctness by comparing inputs/outputs

- Example 1: Search Engine
  - $\forall\, q1, q2.\ q1 \subseteq q2 \Rightarrow Items(q1) \geq Items(q2)$

- Example 2: Numerical Software
  - $\forall\, \theta_1, \theta_2.\ \theta_2 = 2 * \pi + \theta_1 \Rightarrow Sin(\theta_1) == Sin(\theta_2)$

- Example 3: Tax Software
  - $\forall\, x_1, x_2.\ x_2 \equiv_{age} x_1 \wedge x1.age \geq x2.age \Rightarrow Return(x_1) \geq Return(x_2)$

# TenForty

**Tax Law and Policy**

IRS 1040

&

Related Publications

**Scenarios**

(1) Eligibility for senior
and disability benefits;

(2) Eligibility for EITC
benefits;

...

**Metamorphic Specification**

$\forall \mathbf{x}, \mathbf{y} \ (\mathbf{x} \equiv_{age} \mathbf{y}) \wedge (\mathbf{x}.age \geq 65) \wedge$
$(\mathbf{y}.age < 65) \implies (\mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y}))$

$\forall \mathbf{x} (\mathbf{x}.sts = MFJ) \wedge (\mathbf{x}.AGI > 56,844)$
$\implies \forall \mathbf{y} (\mathbf{x} \equiv_{EITC} \mathbf{y} \wedge \mathbf{x}.EITC > 0.0 \wedge$
$\mathbf{y}.EITC = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$

...

---

**Algorithm 1:** RANDOMTESTCASEGENERATION

**Input:** Tax preparation software $\mathcal{P}$, initial input seeds $I$, metamorphic property $p$, a tolerance threshold $\delta$, a Bayes factor $B$, a lower-bound on the confidence $\theta$, and timeout $T$.

**Output:** Passed/Failed, test cases, decision tree

1   $(x_m, \Delta_{FTR}, res) \leftarrow$ SAMPLE$(I)$, 0, True
2   **while** $time() - start\_time < T$ **do**
3     $k \leftarrow 0$
4     $x_1 \leftarrow$ UNIFORMPERTURB$(x_m, p)$
5     $x_2 \leftarrow$ UNIFORMPERTURB$(x_1, p)$
6     $\Delta \leftarrow$ DISTANCE$(\mathcal{P}(x_1), \mathcal{P}(x_2))$
7     **if** $\Delta > \delta$ **then**
8       $I$.ADD$((x_1, x_2),$ 'failed'$)$
9       **if** $\Delta > \Delta_{FTR}$ **then**
10         $x_m \leftarrow x_1$
11         $\Delta \leftarrow \Delta_{FTR}$
12       $res \leftarrow$ False
13     **else**
14       $I$.ADD$((x_1, x_2),$ 'passed'$)$
15       $k \leftarrow k + 1$
16       **if** $k < \frac{-\log B}{\log \theta}$ **then**
17         Go to 5

# Research Questions

- RQ1: Are **metamorphic relation (MR)** useful to capture the legal requirements of tax preparation software?

- RQ2: Can randomized algorithm with Bayesian guarantees be effective in **testing tax preparation** software against the MR?

- RQ3: Could data-driven fault localization help **pinpoint the root of failures** in the internal and input spaces?

# RQ1: Suitability of MR for Tax Law and Policy

| Id | Domain | Metamorphic Property |
|---|---|---|
| 1 | Disability | $\forall \mathbf{x}, \mathbf{y}((\mathbf{x} \equiv_{age} \mathbf{y}) \wedge (\mathbf{x}.age \geq 65) \wedge (\mathbf{y}.age < 65)) \vee ((\mathbf{x} \equiv_{blind} \mathbf{y}) \wedge (\mathbf{x}.blind \wedge \neg \mathbf{y}.blind)) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 2 | Disability | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \implies \forall \mathbf{y}((\mathbf{x} \equiv_{s\_age} \mathbf{y}) \wedge (\mathbf{x}.s\_age \geq 65) \wedge (\mathbf{y}.s\_age < 65)) \vee ((\mathbf{x} \equiv_{s\_blind} \mathbf{y}) \wedge (\mathbf{x}.s\_blind \wedge \neg \mathbf{y}.s\_blind)) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 3 | EITC | $\forall \mathbf{x}(\mathbf{x}.sts = MFS) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{L27} \mathbf{y} \wedge \mathbf{x}.L27 > 0.0 \wedge \mathbf{y}.L27 = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 4 | EITC | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \wedge (\mathbf{x}.AGI > 56,844) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{L27} \mathbf{y} \wedge \mathbf{x}.L27 > 0.0 \wedge \mathbf{y}.L27 = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 5 | EITC | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{AGI} \mathbf{y} \wedge \mathbf{x}.AGI \leq 56,844 \wedge \mathbf{y}.AGI > 56,844) \vee (\mathbf{x} \equiv_{L27} \mathbf{y} \wedge \mathbf{x}.L27 > 0.0 \wedge \mathbf{y}.L27 = 0.0) \vee (\mathbf{x} \equiv_{QC} \mathbf{y} \wedge \mathbf{x}.QC \geq \mathbf{y}.QC) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 6 | EITC | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \wedge (\mathbf{x}.AGI \leq 56,844) \implies \forall \mathbf{y}((\mathbf{x} \equiv_{L27} \mathbf{y}) \wedge \mathbf{x}.L27 \geq \mathbf{y}.L27) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 7 | CTC | $\forall \mathbf{x}(\mathbf{x}.sts = MFS) \wedge (\mathbf{x}.AGI < 200k) \forall \mathbf{y}((\mathbf{x} \equiv_{L19} \mathbf{y}) \wedge (\mathbf{x}.L19 > \mathbf{y}.L19)) \implies \mathcal{F}(\mathbf{x}) > \mathcal{F}(\mathbf{y})$ |
| 8 | CTC | $\forall \mathbf{x}, \mathbf{x}'(\mathbf{x}.sts = \mathbf{x}'.sts = MFJ) \wedge (\mathbf{x}.AGI < 400k) \wedge (\mathbf{x}'.AGI \geq 400k) \wedge \lceil \mathbf{x}'.AGI - 400k \rceil_{1k} * 0.05 < \mathbf{x}'.QC * 2k + \mathbf{x}.OD * 0.5k \implies \forall \mathbf{y}, \mathbf{y}'(\mathbf{x} \equiv_{\{QC, OD\}} \mathbf{y}) \wedge (\mathbf{x}' \equiv_{\{QC, OD\}} \mathbf{y}') \wedge (0 \leq \mathbf{y}.QC = \mathbf{y}'.QC \leq \mathbf{x}.QC = \mathbf{x}'.QC \leq 10) \wedge (0 \leq \mathbf{y}.OD = \mathbf{y}'.OD \leq \mathbf{x}.OD = \mathbf{x}'.OD \leq 10) \implies (\mathcal{F}(\mathbf{x}) - \mathcal{F}(y)) \geq (\mathcal{F}(x') - \mathcal{F}(y'))$ |
| 9 | ETC | $\forall \mathbf{x}(\mathbf{x}.sts = MFS) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{L29} \mathbf{y} \wedge \mathbf{x}.L29 > 0.0 \wedge \mathbf{y}.L29 = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 10 | ETC | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \wedge (\mathbf{x}.AGI \geq 180k) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{L29} \mathbf{y} \wedge \mathbf{x}.L29 > 0.0 \wedge \mathbf{y}.L29 = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 11 | ETC | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \wedge (\mathbf{x}.AGI \leq 160k) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{L29} \mathbf{y} \wedge \mathbf{x}.L29 \geq \mathbf{y}.L29) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 12 | ETC | $\forall \mathbf{x}, \mathbf{x}'(\mathbf{x}.sts = \mathbf{x}'.sts = MFJ) \wedge (\mathbf{x}.AGI \leq 160k) \wedge (160k < \mathbf{x}'.AGI < 180k) \implies \forall \mathbf{y}, \mathbf{y}'((\mathbf{x} \equiv_{L29} \mathbf{y}) \wedge (\mathbf{x}' \equiv_{L29} \mathbf{y}') \wedge (\mathbf{x}.L29 = \mathbf{x}'.L29 \geq \mathbf{y}.L29 = \mathbf{y}'.L29)) \implies (\mathcal{F}(\mathbf{x}) - \mathcal{F}(y)) \geq (\mathcal{F}(x') - \mathcal{F}(y'))$ |
| 13 | ID | $\forall \mathbf{x}, \mathbf{y}(\mathbf{x} \equiv_{MDE} \mathbf{y}) \wedge (\mathbf{x}.MDE \leq \mathbf{x}.AGI * 7.5\%) \wedge (\mathbf{y}.MDE = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 14 | ID | $\forall \mathbf{x}(\neg \mathbf{x}.iz) \implies \forall \mathbf{y}(\mathbf{x} \equiv_{MDE} \mathbf{y} \wedge \mathbf{x}.MDE > 0.0 \wedge \mathbf{y}.MDE = 0.0) \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ |
| 15 | ID | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \implies \forall \mathbf{y}((\mathbf{x} \equiv_{iz, L12} \mathbf{y}) \wedge (\mathbf{x}.iz \wedge \neg \mathbf{y}.iz) \wedge (\mathbf{x}.L12 \leq 24.8k \wedge \mathbf{y}.L12 = 0.0)) \implies \mathcal{F}(\mathbf{x}) \leq \mathcal{F}(\mathbf{y})$ |
| 16 | ID | $\forall \mathbf{x}(\mathbf{x}.sts = MFJ) \implies \forall \mathbf{y}((\mathbf{x} \equiv_{iz, L12} \mathbf{y}) \wedge (\mathbf{x}.iz \wedge \neg \mathbf{y}.iz) \wedge (\mathbf{x}.L12 > 24.8k \wedge \mathbf{y}.L12 = 0.0)) \implies \mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |

# RQ1: Suitability of MR for Tax Law and Policy

| Id | Year 2018 | Year 2019 | Year 2021 |
|---|---|---|---|
| 1,2 | No Change | No Change | No Change |
| 3 | No Change | No Change | $\mathcal{F}(\mathbf{x}) \geq \mathcal{F}(\mathbf{y})$ |
| 4 | $\mathbf{x}.AGI > 54,884$ | $\mathbf{x}.AGI > 55,952$ | $\mathbf{x}.AGI > 57,414$ |

**Answer RQ1**:
❖ Metamorphic relations are suitable to specify the correctness requirements in tax software.
❖ These relations allow us to update the requirements as the tax policies evolve over time.

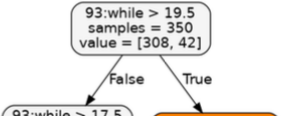| Id | Year 2018 | Year 2019 | Year 2021 |
|---|---|---|---|
| 14 | Not Possible | Not Possible | No Change |
| 15 | $\mathbf{x}.L8 \leq 24.0k \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ | $\mathbf{x}.L9 \leq 24.4k \implies \mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{y})$ | $\mathbf{x}.L12 \leq 25.1k \implies \mathcal{F}(\mathbf{x}) \leq \mathcal{F}(\mathbf{y})$ |
| 16 | $\mathbf{x}.L8 > 24.0k$ | $\mathbf{x}.L9 > 24.4k$ | $\mathbf{x}.L12 > 25.1k$ |

# RQ2: Testing Software against MR requirements

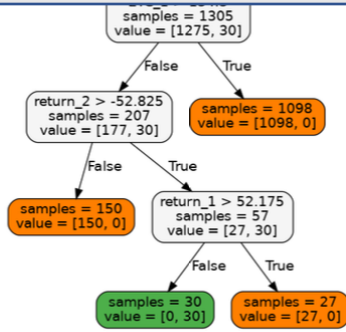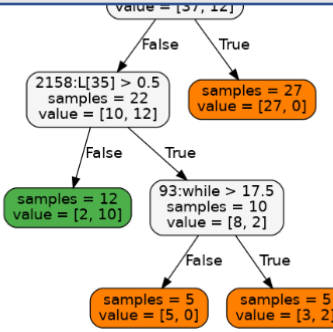| Property ID | OpenTaxSolver 2018 | | | | OpenTaxSolver 2019 | | | | OpenTaxSolver 2020 | | | | OpenTaxSolver 2021 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *#test cases* | *#fail* | *#pass* | *$T_F(s)$* | *#test cases* | *#fail* | *#pass* | *$T_F(s)$* | *#test cases* | *#fail* | *#pass* | *$T_F(s)$* | *#test cases* | *#fail* | *#pass* | *$T_F(s)$* |
| Disability (1) | 36,558 | 0 | 36,558 | *N/A* | 35,970 | 0 | 35,970 | *N/A* | 36,255 | 0 | 36,255 | *N/A* | 32,456 | 0 | 32,456 | *N/A* |
| Disability (2) | 36,369 | 0 | 36,369 | *N/A* | 36,780 | 0 | 36,780 | *N/A* | 35,790 | 0 | 35,790 | *N/A* | 32,355 | 0 | 32,355 | *N/A* |
| EITC (3) | 3 | | | | | | | | | | | | | | 32,343 | *N/A* |
| EITC (4) | 3 | | | | | | | | | | | | | | 0 | 0.05 |
| EITC (5) | 3 | | | | | | | | | | | | | | 32,883 | *N/A* |
| EITC (6) | 3 | | | | | | | | | | | | | | 32,962 | *N/A* |
| CTC (7) | 3 | | | | | | | | | | | | | | 32,388 | *N/A* |
| CTC (8) | 1 | | | | | | | | | | | | | | 16,346 | *N/A* |
| ETC (9) | 3 | | | | | | | | | | | | | | 1,102 | 0.05 |
| ETC (10) | 3 | | | | | | | | | | | | | | 34 | 0.05 |
| ETC (11) | 1 | | | | | | | | | | | | | | 16,459 | 29.02 |
| ETC (12) | 1 | | | | | | | | | | | | | | 14,636 | *N/A* |
| ID (13) | 36,801 | 0 | 36,801 | *N/A* | 36,210 | 0 | 36,210 | *N/A* | 36,160 | 15 | 36,145 | 70.09 | 27,348 | 5,508 | 21,840 | 0.06 |
| ID (14) | — | — | — | — | — | — | — | — | 36,405 | 0 | 36,405 | *N/A* | 31,916 | 0 | 31,916 | *N/A* |
| ID (15) | 36,926 | 0 | 36,926 | *N/A* | 36,630 | 0 | 36,630 | *N/A* | 36,315 | 0 | 36,315 | *N/A* | 32,793 | 0 | 32,793 | *N/A* |
| ID (16) | 36,846 | 0 | 36,846 | *N/A* | 36,570 | 0 | 36,570 | *N/A* | 36,235 | 10 | 36,225 | 46.02 | 32,363 | 8 | 32,355 | 44.34 |

**Answer RQ2**:
- ❖ Updated software is no longer satisfying the correctness requirements.
- ❖ Multiple weakness areas relate to married filing separately status.
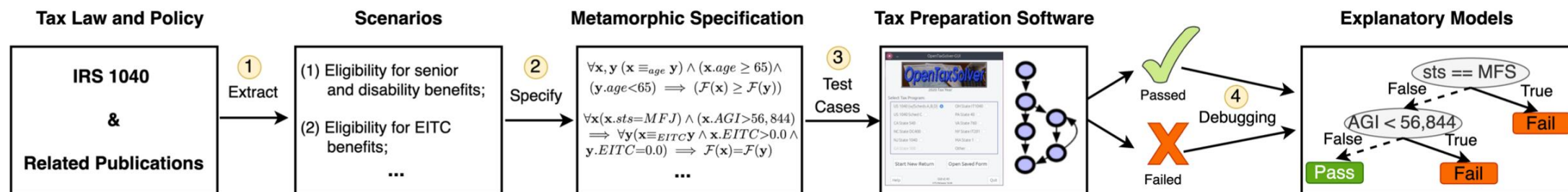
# RQ3: Data-Driven Root Cause Identification



| Id | Debugging Input Space | Debugging Internal Space |
|---|---|---|
| | | 93:while > 19.5<br>samples = 350<br>value = [308, 42] |

**Answer RQ3**:

❖ Decision trees are useful artifacts to explain failing circumstances.

❖ Our experiences show that the software might completely miss an eligibility condition.

❖ Our results also showed unexpected errors due to finite precision in the computation.

**Forensic DNA Software**
- New York City's Office of Chief Medical Examiner (OCME) for thousands of criminal cases between 2011 and 2017

- Undisclosed data dropping method CheckFrequencyForRemoval()

- Falsely skew results toward false inclusion for individuals whose DNA was not present.

**"Do I Qualify?" Screening Software**
- Poverty management systems in Pennsylvania (Check Eligibility)

- Comparative implementation of benefit eligibility handbook

- Errors in the eligibility checking: Exclude the most vulnerable families from receiving the essential aids