

# Anytime Approximate Formal Feature Attribution

---

Jinqiang Yu<sup>1,2</sup>, Graham Farr<sup>1</sup>, Alexey Ignatiev<sup>1</sup>, Peter J. Stuckey<sup>1,2</sup>

1. Department of Data Science and AI, Monash University, Australia
2. Australian Research Council OPTIMA ITTC, Australia



MONASH  
University



OPTiMA

ARC TRAINING CENTRE IN  
OPTIMISATION TECHNOLOGIES  
INTEGRATED METHODOLOGIES  
AND APPLICATIONS

# Table of Contents

Introduction

Background

Approximate Formal Feature Attribution

Experimental Results

Conclusions

# Explainable Artificial Intelligence (XAI)

Rapid advances in Artificial Intelligence (AI) and Machine Learning (ML) algorithms.

**Issue:** opaque models → lack of trust.

**Rise:** Explainable Artificial Intelligence (XAI).

**Solution:** Post-hoc explanations.

**Post-hoc explanations** answer '*why?*' questions and '*how?*' questions

**Heuristic approach:**

- \* Feature selection: Anchor.
- \* Feature Attribution: LIME, SHAP, etc.
- \* Issue:
  - \* Explanation quality.

**Post-hoc explanations** answer '*why?*' questions and '*how?*' questions

**Formal approach:**

- \* Correct and minimal.
- \* Feature selection: abductive explanation (AXp) and contrastive explanation (CXp).
- \* Feature attribution: formal feature attribution (FFA).

## Formal feature attribution (FFA):

Provide the importance of each feature, i.e. the proportion of AXp's in which it appears.

### FFA approach:

- \* Make use of the *hitting set duality* between AXp's and CXp's.
- \* Collect AXp's as a *side effect* of CXp enumeration algorithm.

## Observations from the FFA approach:

- \* Usually find many AXp's before finding the first CXp.
- \* AXp's are diverse → good approximation of FFA.
- \* AXp enumeration can get *exact* FFA faster.

**Issue:** *Exact* FFA is hard to compute.

# Table of Contents

Introduction

Background

Approximate Formal Feature Attribution

Experimental Results

Conclusions



## Boolean Satisfiability (SAT)

- \* Decision problem for propositional logic.
- \* Formula  $\varphi$ : Conjunctive normal form (CNF).
  - \* Clause: a disjunction of literals.
  - \* Literal: a Boolean variable  $b$  or  $\neg b$ .
  - \* Example:  $(a \vee \neg c) \wedge (b \vee c)$ .
- \* Satisfiable: there exists an assignment  $\mu$  satisfying the formula.

**Maximum Satisfiability (MaxSAT):** maximize the number of satisfied clauses.

**Partial Unweighted MaxSAT:**  $\phi = \mathcal{H} \wedge \mathcal{S}$ .

- \*  $\mathcal{H}$ : hard clauses, which *must* be satisfied.
- \*  $\mathcal{S}$ : soft clauses, which represent a *preference* to satisfy those clauses.
- \* Aim: maximize the number of satisfied soft clauses.

Let  $\phi = \mathcal{H} \cup \mathcal{S}$  and  $\phi \models \perp$ .

## Minimal Unsatisfiable Subset (MUS):

A subset of clauses  $\mu \subseteq \mathcal{S}$  is a *MUS* iff  $\mathcal{H} \cup \mu \models \perp$  and  $\forall \mu' \subsetneq \mu$  it holds that  $\mathcal{H} \cup \mu' \not\models \perp$ .

## Minimal Correction Subset (MCS):

A subset of clauses  $\sigma \subseteq \mathcal{S}$  is a *MCS* iff  $\mathcal{H} \cup \mathcal{S} \setminus \sigma \not\models \perp$  and  $\forall \sigma' \subsetneq \sigma$  it holds that  $\mathcal{H} \cup \mathcal{S} \setminus \sigma' \models \perp$ .

# Minimal Hitting Set Duality

**Minimal hitting set (MHS)** duality relationship between MUSes and MCSes, i.e.

$$\mathbb{U}_\phi = \text{MHS}(\mathbb{C}_\phi) \text{ and } \mathbb{C}_\phi = \text{MHS}(\mathbb{U}_\phi)$$

where

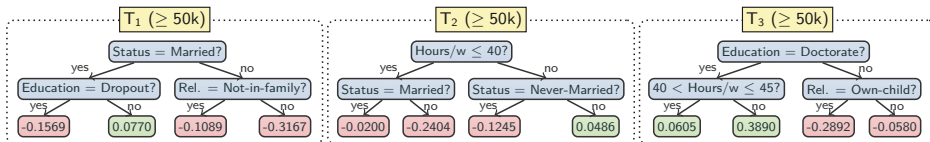
- \*  $\mathbb{U}_\phi$ : MUSes.
- \*  $\mathbb{C}_\phi$ : MCSes.
- \*  $\text{MHS}(S)$  returns the minimal hitting sets of  $S$ .
- \* *Minimal* sets that share an element with each subset in  $S$ .

# Classification Example

A classification function  $\tau: \mathcal{F} \rightarrow K$ .

- \*  $\mathcal{F}$  : complete feature space.
- \*  $K$ : a set of classes.

# Classification Problem



**Figure 1:** Boosted tree model trained on the *adult* classification dataset.

## Instance:

{Education=Bachelors, Status=Separated, Occupation=Sales, Relationship=Not-in-family, Sex=Male, Hours/w  $\leq 40$ }

**Score:**  $-0.4073 = (-0.1089 - 0.2404 - 0.0580) < 0 \rightarrow$  prediction  $< 50k$

## Formal Explanation

**Abductive explanation (AXp)**  $\mathcal{X}$ : subset-minimal set of features sufficing to explain the prediction .

$$\forall(\mathbf{x} \in \mathbb{F}). \left[ \bigwedge_{i \in \mathcal{X}} (x_i = v_i) \right] \rightarrow (\kappa(\mathbf{x}) = c)$$

**Contrastive explanation (CXp)**  $\mathcal{Y}$ : subset-minimal set of features that are necessary to change the prediction.

$$\exists(\mathbf{x} \in \mathbb{F}). \left[ \bigwedge_{i \notin \mathcal{Y}} (x_i = v_i) \right] \wedge (\kappa(\mathbf{x}) \neq c)$$

**Minimal Hitting Set Duality:** CXps minimally hit every AXp, and vice-versa.

# AXp and CXp Examples

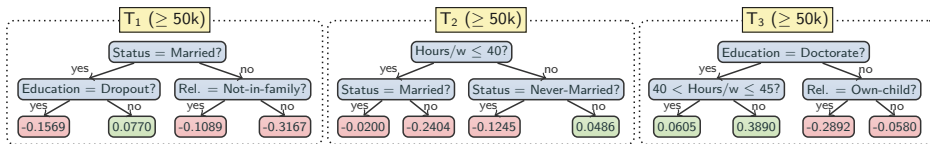


Figure 2: Boosted tree model trained on the *adult* classification dataset.

---

$\mathcal{X}_1 = \{$	Education, Hours/w $\}$
IF	Education = Bachelors
AND	Hours/w $\leq 40$
THEN	Target < 50k

---

(a) AXp  $\mathcal{X}_1$ .

---

$\mathcal{X}_2 = \{$	Education, Status $\}$
IF	Education = Bachelors
AND	Status = Separated
THEN	Target < 50k

---

(b) AXp  $\mathcal{X}_2$ .



# AXp and CXp Examples

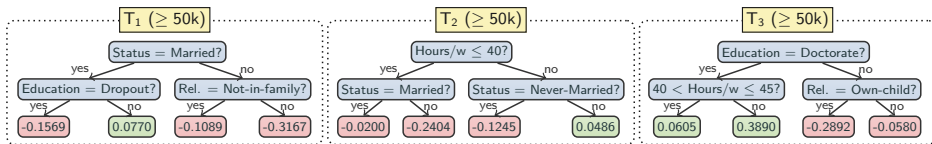


Figure 2: Boosted tree model trained on the *adult* classification dataset.

	$\mathcal{Y}_1 = \{ \text{Education} \}$
IF	Education $\neq$ Bachelors
THEN	Target can be changed to $\geq 50k$

(a) CXp  $\mathcal{Y}_1$ .

	$\mathcal{Y}_2 = \{ \text{Hours/w, Status} \}$
IF	Hours/w $\leq 40$
AND	Status $\neq$ Separated
THEN	Target can be changed to $\geq 50k$

(b) CXp  $\mathcal{Y}_2$ .

# Formal Feature Attribution (FFA)

Inspired by the implicit hitting set based algorithm eMUS/MARCO.

$$\text{ffa}_{\kappa}(i, (\mathbf{v}, c)) = \frac{|\{\mathcal{X} \mid \mathcal{X} \in \mathbb{A}_{\kappa}(\mathbf{v}, c), i \in \mathcal{X}\}|}{|\mathbb{A}_{\kappa}(\mathbf{v}, c)|} \quad (1)$$

where

- \*  $\mathbf{v}$ : instance.
- \*  $c$ : prediction.
- \*  $\mathcal{X}$ : AXp.
- \*  $\mathbb{A}_{\kappa}(\mathbf{v}, c)$ : the set of AXp's for  $\mathbf{v}$ .
- \*  $i$ : feature.

---

**Algorithm 1** Anytime Explanation Enumeration

---

**Input:** Classifier:  $\kappa$ , instance:  $\mathbf{v}$ , prediction:  $c$

**Output:** AXP's:  $\mathbb{A}$ , CXP's:  $\mathbb{C}$

```
1:  $(\mathbb{A}, \mathbb{C}) \leftarrow (\emptyset, \emptyset)$   $\triangleright$ Sets of AXP's and CXP's to collect.
2: while resources available do
3:    $\mathcal{Y} \leftarrow \text{MINIMALHS}(\mathbb{A}, \mathbb{C})$   $\triangleright$ Get a new MHS of  $\mathbb{A}$  subject to  $\mathbb{C}$ .
4:    $\triangleright \mathcal{Y}$  is the CXP candidate.
5:   if  $\mathcal{Y} = \perp$  then break  $\triangleright$ Stop if none is computed.
6:    $\triangleright$ Check CXP condition for  $\mathcal{Y}$ .
7:   if  $\exists (\mathbf{x} \in \mathbb{F}). \left[ \bigwedge_{i \notin \mathcal{Y}} (x_i = v_i) \right] \wedge (\kappa(\mathbf{x}) \neq c)$  then
8:      $\mathbb{C} \leftarrow \mathbb{C} \cup \{\mathcal{Y}\}$   $\triangleright \mathcal{Y}$  appears to be a CXP.
9:   else
10:     $\mathcal{X} \leftarrow \text{EXTRACTAXP}(\mathcal{F} \setminus \mathcal{Y}, \kappa, \mathbf{v}, c)$ 
11:     $\mathbb{A} \leftarrow \mathbb{A} \cup \{\mathcal{X}\}$   $\triangleright$ There must be a missing AXP  $\mathcal{X} \subseteq \mathcal{F} \setminus \mathcal{Y}$ .
return  $\mathbb{A}, \mathbb{C}$ 
```

---

---

**Algorithm 2** ExtractAXP

---

**Input:** Candidate:  $\mathcal{X}$ , classifier:  $\kappa$ , instance:  $\mathbf{v}$ , prediction:  $c$

**Output:** AXP:  $\mathcal{X}$

```
1: for  $j \in \mathcal{X}$  do
2:   if  $\forall (\mathbf{x} \in \mathbb{F}). \left[ \bigwedge_{i \in \mathcal{X} \setminus \{j\}} (x_i = v_i) \right] \rightarrow (\kappa(\mathbf{x}) = c)$  then
3:      $\mathcal{X} \leftarrow \mathcal{X} \setminus \{j\}$ 
return  $\mathcal{X}$ 
```

---

# FFA Example

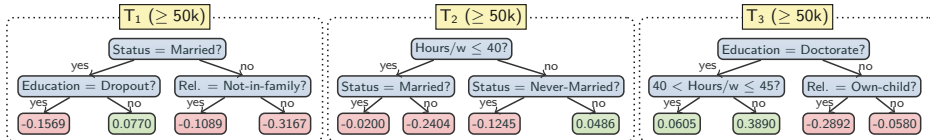


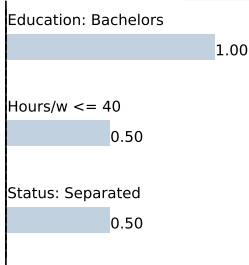
Figure 4: Boosted tree model trained on the *adult* classification dataset.

$\mathcal{X}_1 = \{ \text{Education, Hours/w} \}$
IF Education = Bachelors AND Hours/w ≤ 40 THEN Target < 50k

(a) AXp  $\mathcal{X}_1$ .

$\mathcal{X}_2 = \{ \text{Education, Status} \}$
IF Education = Bachelors AND Status = Separated THEN Target < 50k

(b) AXp  $\mathcal{X}_2$ .



(c) FFA.

# LIME, SHAP and FFA Examples

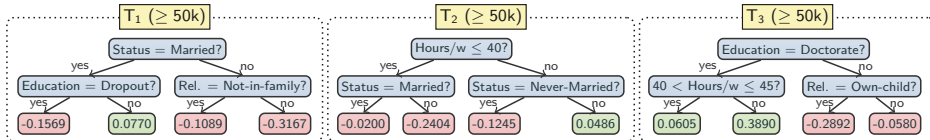
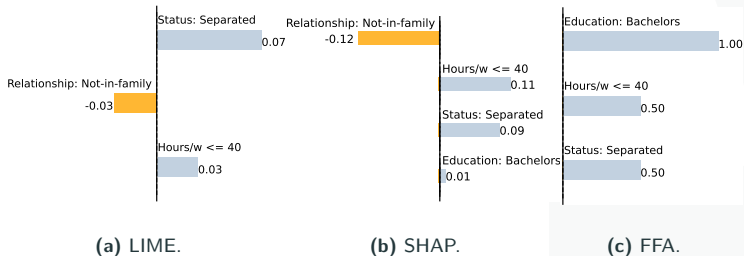


Figure 6: Boosted tree model trained on the *adult* classification dataset.



Issue of LIME and SHAP: Some *irrelevant* features have non-zero attribution,

# Table of Contents

Introduction

Background

Approximate Formal Feature Attribution

Experimental Results

Conclusions

# Switching from CXp to AXp Enumeration

**Issue in the FFA approach:** *Exact* FFA hard to compute.

**Inspirations:**

- \* AXp enumeration → getting *exact* FFA faster.
- \* Diverse AXp's → good FFA approximations.
- \* AXp enumeration → AXp's are not diverse.
- \* CXp enumeration → diverse AXp's → quick convergence.

**Proposed approach:**

- \* Anytime approach to computing approximate FFA
- \* Start with CXp enumeration
- \* Switch to AXp enumeration at some point.

# Algorithm with Switching

---

## Algorithm 3 Adaptive Explanation Enumeration

---

- 1:  $(\mathbb{E}_0, \mathbb{E}_1) \leftarrow (\emptyset, \emptyset)$   $\triangleright$  CXp's and AXp's to collect
  - 2:  $\rho \leftarrow 0$   $\triangleright$  Target phase of enumerator, initially CXp
  - 3: **while** true **do**
  - 4:    $\mu \leftarrow \text{MINIMALHS}(\mathbb{E}_{1-\rho}, \mathbb{E}_\rho, \rho)$   $\triangleright$  MHS of  $\mathbb{E}_{1-\rho}$  s.t.  $\mathbb{E}_\rho$ .
  - 5:   **if**  $\mu = \perp$  **then break**  $\triangleright$  Stop if none is computed.
  - 6:    $\triangleright$  Check CXp condition for  $\mathcal{Y}$ .
  - 7:   **if**  $\text{ISTARGETXP}(\mu, \tau, \mathbf{v}, \mathbf{c})$  **then**
  - 8:      $\mathbb{E}_\rho \leftarrow \mathbb{E}_\rho \cup \{\mu\}$   $\triangleright$  There must be a missing AXp  $\mathcal{X} \subseteq \mathcal{F} \setminus \mathcal{Y}$ .
  - 9:   **else**
  - 10:      $\triangleright$  Collect target expl.  $\mu$
  - 11:      $\nu \leftarrow \text{EXTRACTDUALXP}(\mathcal{F} \setminus \mu, \tau, \mathbf{v}, \mathbf{c})$
  - 12:      $\mathbb{E}_{1-\rho} \leftarrow \mathbb{E}_{1-\rho} \cup \{\nu\}$   $\triangleright$  Collect dual expl.  $\nu$
  - 13:      $\triangleright$  The difference!
  - 14:     **if**  $\text{ISSWITCHNEEDED}(\mathbb{E}_\rho, \mathbb{E}_{1-\rho}, w, \alpha, \varepsilon)$  **then**
  - 15:        $\rho \leftarrow 1 - \rho$   $\triangleright$  Flip phase of MINIMALHS
  - return**  $\mathbb{E}_1, \mathbb{E}_0$   $\triangleright$  Result AXp's and CXp's
-



# Switching Criteria

**Criterion 1:** Switch when CXp's on average are *much* smaller than AXp's, i.e. when

$$\frac{\sum_{\mathcal{X} \in \mathbb{A}^w} |\mathcal{X}|}{\sum_{\mathcal{Y} \in \mathbb{C}^w} |\mathcal{Y}|} \geq \alpha, \quad (2)$$

**Criterion 2:** Switch when the average CXp size “stabilizes”.

$$\left| |\mathcal{Y}_{\text{new}}| - \frac{\sum_{\mathcal{Y} \in \mathbb{C}^w} |\mathcal{Y}|}{w} \right| \leq \varepsilon, \quad (3)$$

**Rational:**

- \* Normally  $|\mathcal{X}| > |\mathcal{Y}|$
- \* CXp extraction: check satisfiable  $\rightarrow$  cheap.
- \* AXp extraction: check unsatisfiable  $\rightarrow$  expensive.
- \* Before switching: ensure AXp's diverse.
- \* After switching: single call for AXp, multiple calls for CXp extraction.

# Switching Criteria

**Criterion 1:** Switch when CXp's on average are *much* smaller than AXp's, i.e. when

$$\frac{\sum_{\mathcal{X} \in \mathcal{A}^w} |\mathcal{X}|}{\sum_{\mathcal{Y} \in \mathcal{C}^w} |\mathcal{Y}|} \geq \alpha, \quad (2)$$

**Criterion 2:** Switch when the average CXp size “stabilizes”.

$$\left| |\mathcal{Y}_{\text{new}}| - \frac{\sum_{\mathcal{Y} \in \mathcal{C}^w} |\mathcal{Y}|}{w} \right| \leq \varepsilon, \quad (3)$$

Switch when meeting *either* of the two criteria!

# Table of Contents

Introduction

Background

Approximate Formal Feature Attribution

**Experimental Results**

Conclusions

# Experimental Setup

**Datasets:** 3 Images and 2 text data

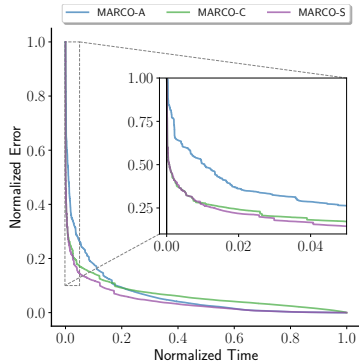
**Metrics:**

- \* **Errors:** Manhattan distance, i.e. the sum of absolute differences across all features.
- \* **Kendall's Tau:** Similarity of two rankings. Ranging  $[-1, 1]$ . The higher the closer.
- \* **Rank-biased overlap (RBO):** Similarity of two rankings. Ranging  $[0, 1]$ . The higher the closer.
- \* **Kullback–Leibler (KL) divergence:** Statistical distance between two probability distributions. Ranging from 0 to  $\infty$ .
- \* **Number of AXp's**

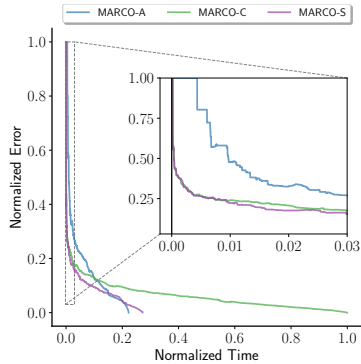
## Average Runtime:

- \* MARCO-S (Our approach): 3509.50s (9.26s – 30881.55s)
- \* MARCO-A (AXp enumeration): 3255.30s (2.15s – 29191.42s)
- \* MARCO-C (CXp enumeration): 19311.87s (9.39s – 55951.57s)

# Error Results



(a) Mean

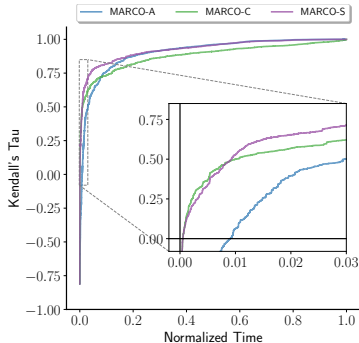


(b) Median

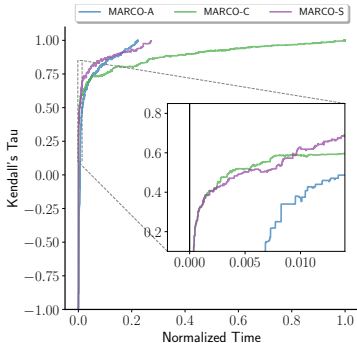
**Figure 8:** FFA approximation error over time.

- \* **MARCO-S:** Propose approach
- \* **MARCO-A:** AXp enumeration
- \* **MARCO-C:** CXp enumeration

# Kendall's Tau Results



(a) Mean

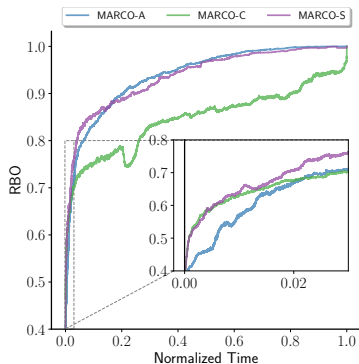


(b) Median

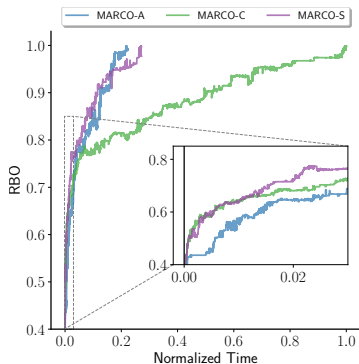
Figure 9: Kendall's Tau over time.

- \* MARCO-S: Propose approach
- \* MARCO-A: AXp enumeration
- \* MARCO-C: CXp enumeration

# RBO Results



(a) Mean



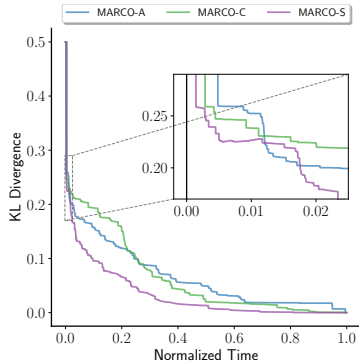
(b) Median

Figure 10: RBO over time.

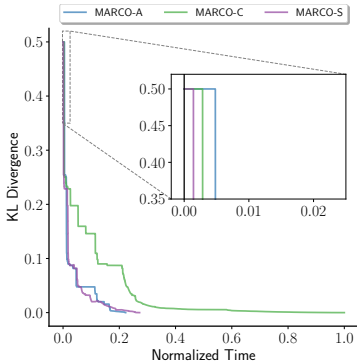
- \* MARCO-S: Propose approach
- \* MARCO-A: AXp enumeration
- \* MARCO-C: CXp enumeration



# KL Divergence Results



(a) Mean



(b) Median

**Figure 11:** KL divergence over time.

- \* **MARCO-S:** Propose approach
- \* **MARCO-A:** AXp enumeration
- \* **MARCO-C:** CXp enumeration

# Number of AXp's Results

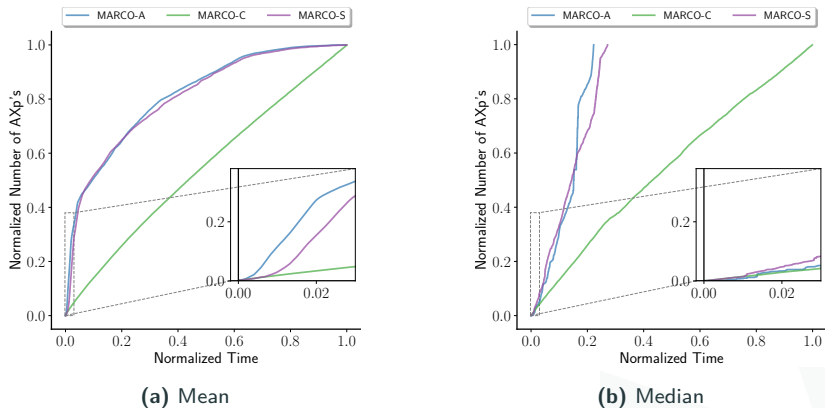


Figure 12: Number of AXp's over time.

- \* MARCO-S: Propose approach
- \* MARCO-A: AXp enumeration
- \* MARCO-C: CXp enumeration

# Number of AXp's Examples

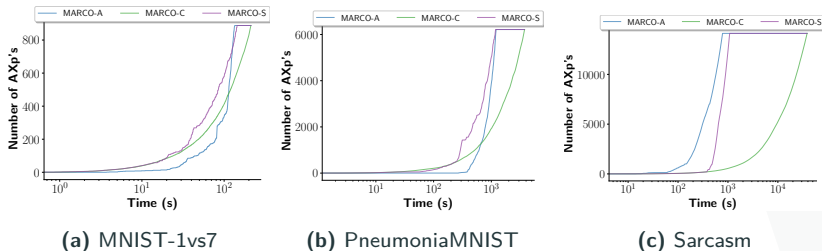


Figure 13: Number of AXp's over time in example instances.

- \* MARCO-S: Propose approach
- \* MARCO-A: AXp enumeration
- \* MARCO-C: CXp enumeration

# Table of Contents

Introduction

Background

Approximate Formal Feature Attribution

Experimental Results

Conclusions

# Conclusions

- \* The proposed approach can replicate the behavior of the superior competitor → efficient and good approximation of FFA.
- \* Start with CXp enumeration → diverse AXp's.
- \* Switching to AXp enumeration → extracting AXp's faster.
- \* The proposed mechanism can be readily adapted to a multitude of other problems, e.g. in over-constrained systems or model-based diagnosis (MBD)

Thank you!