

INTERPRETING DEEP LEARNING MODELS FINE-TUNED FOR DETECTING VULNERABILITIES RELATED TO MISSING CODE

Claudia Mamede^{1,2}, Claire Le Goues², José Campos^{1,3}, Rui Abreu^{1,4}

¹ Faculty of Engineering, University of Porto, Portugal

² Carnegie Mellon University, USA

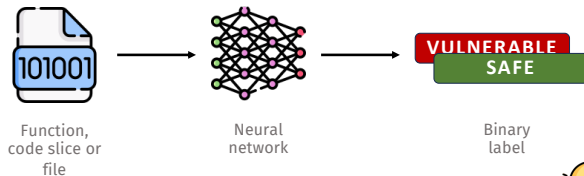
³ LASIGE, Faculty of Sciences, University of Lisbon, Portugal

⁴ INESC-ID, Instituto Superior Técnico, University of Lisbon, Portugal



MOTIVATION

BACKGROUND. LEARNING-BASED VULNERABILITY DETECTION IS TYPICALLY TREATED AS A BINARY CLASSIFICATION PROBLEM.



PROBLEM 1. MODEL'S HIGH ACCURACY DOES NOT ENSURE REAL-WORLD EFFECTIVENESS

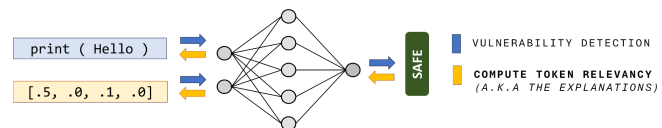
IF THE HAS > 500 LoC, HOW CAN WE FIND AND FIX THE VULNERABILITY?

PROBLEM 2. EXISTING xAI METHODS ONLY EXPLAIN TOKENS PRESENT IN THE INPUT

HOW CAN WE EXPLAIN VULNERABILITIES ARISING FROM THE ABSENCE OF CODE?

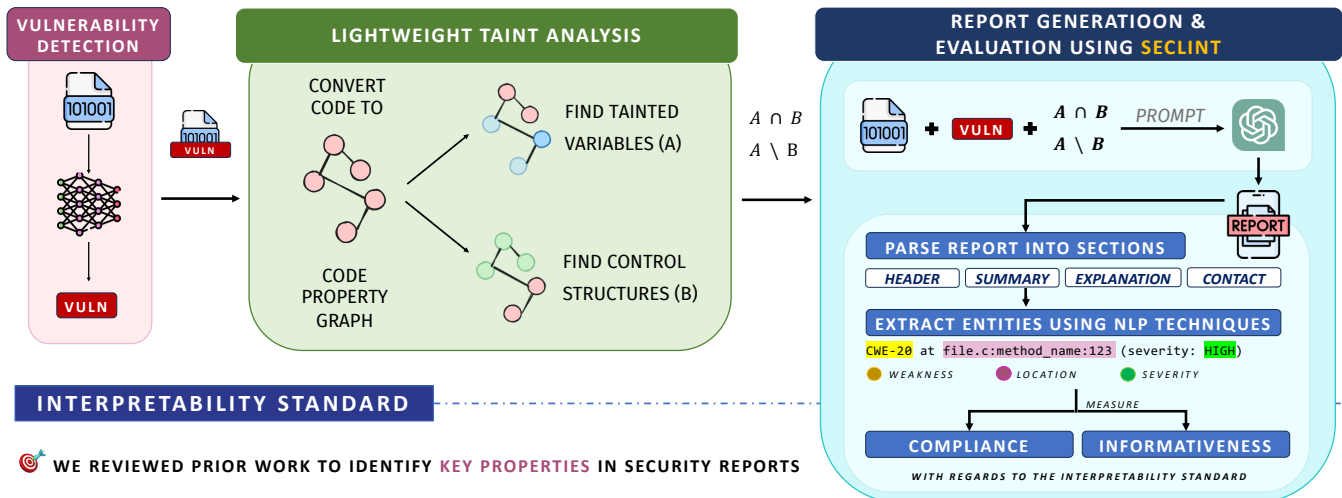


EXPLAINABLE AI (xAI) HELPS DEVELOPERS INTERPRET MODEL OUTPUTS BY ASSIGNING RELEVANCY SCORES TO INPUT TOKENS.



IN VULNERABILITY DETECTION, xAI MAY HELP US LOCATE VULNERABILITIES.

METHODOLOGY



INTERPRETABILITY STANDARD

WE REVIEWED PRIOR WORK TO IDENTIFY KEY PROPERTIES IN SECURITY REPORTS THAT HELP DEVELOPERS UNDERSTAND AND ADDRESS VULNERABILITIES.

1. !VULN-DETECT: <WEAKNESS NAME/ID> AT <LOCATION> (SEVERITY: <LEVEL>)
 2. !WHAT: DESCRIBE THE WEAKNESS/PROBLEM
 3. !WHY: DESCRIBE ITS IMPACT
 4. !HOW: DESCRIBE HOW THE WEAKNESS CAN BE TRIGGERED
 5. !WHEN: DESCRIBE WHEN THE PROBLEM WAS FOUND
 6. !WHERE: DESCRIBE WHERE THE PROBLEM IS LOCATED
- (...)

! – MANDATORY FIELDS

? – OPTIONAL FIELDS

FUTURE WORK

WE ARE STILL REFINING THE INTERPRETABILITY STANDARD



CHECK OUT THE QR CODE BELOW FOR THE UPDATED VERSION & GIVE US FEEDBACK !



INTERPRETING DEEP LEARNING MODELS
FINE-TUNED FOR DETECTING VULNERABILITIES
RELATED TO MISSING CODE

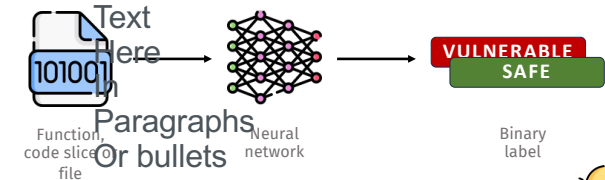
Claudia Mamee^{1,2}, Claire Le Goues², José Campos^{1,3}, Rui Abreu^{1,4}

¹ Faculty of Engineering, University of Porto, Portugal
² Carnegie Mellon University, USA
³ LASIGE, Faculty of Sciences, University of Lisbon, Portugal
⁴ INESC-ID, Instituto Superior Técnico, University of Lisbon, Portugal

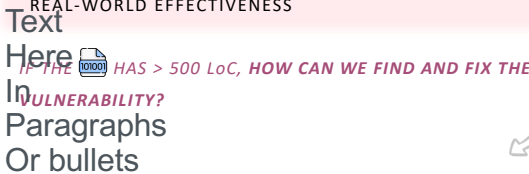


MOTIVATION

BACKGROUND: LEARNING-BASED VULNERABILITY DETECTION IS
TYPICALLY TREATED AS A BINARY CLASSIFICATION PROBLEM.



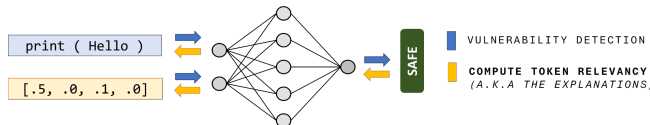
PROBLEM 1. MODEL'S HIGH ACCURACY DOES NOT ENSURE
REAL-WORLD EFFECTIVENESS



EXPLAINABLE AI (xAI) HELPS DEVELOPERS INTERPRET MODEL OUTPUTS
BY ASSIGNING RELEVANCY SCORES TO INPUT TOKENS.

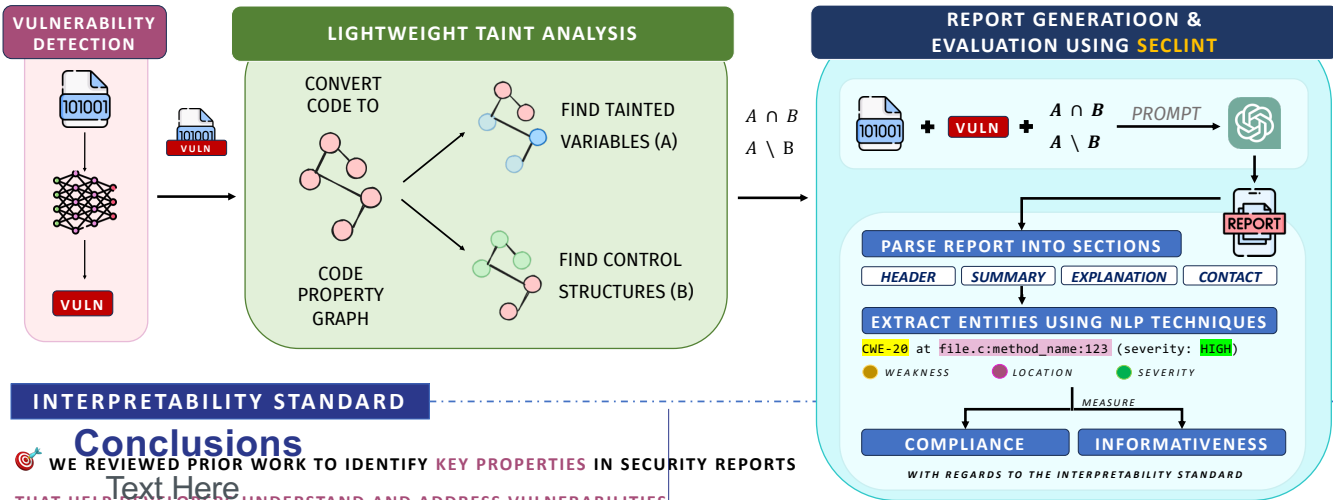
PROBLEM 2. EXISTING xAI METHODS ONLY
EXPLAIN WHAT IS PRESENT IN THE INPUT

HOW CAN WE EXPLAIN VULNERABILITIES ARISING FROM
THE ABSENCE OF CODE?



IN VULNERABILITY DETECTION, xAI MAY HELP US LOCATE VULNERABILITIES.

METHODOLOGY



INTERPRETABILITY STANDARD

WE REVIEWED PRIOR WORK TO IDENTIFY KEY PROPERTIES IN SECURITY REPORTS
THAT HELP DEVELOPERS UNDERSTAND AND ADDRESS VULNERABILITIES.

Text Here	Paragraphs Or bullets
1. !VULN DETECT:	WEAKNESS NAME/ID AT <LOCATION> (SEVERITY: <LEVEL>)
2. !WHAT:	DESCRIBE THE WEAKNESS/PROBLEM
3. !WHY:	DESCRIBE ITS IMPACT
4. !HOW:	DESCRIBE HOW THE WEAKNESS CAN BE TRIGGERED
5. !WHEN:	DESCRIBE WHEN THE PROBLEM WAS FOUND
6. !WHERE:	DESCRIBE WHERE THE PROBLEM IS LOCATED
(...)	

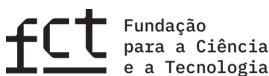
! - MANDATORY FIELDS ? - OPTIONAL FIELDS

Acknowledgements

This work was partially funded by the FCT scholarship or Project's reference/name)

FUTURE WORK

CHECK OUT THE QR CODE BELOW FOR
THE UPDATED VERSION & GIVE US FEEDBACK !



LEARN MORE
ABOUT THIS
RESEARCH