## When in Doubt Throw It Out: Building on Confident Learning for Vulnerability Detection





### Yuanjun Gong · Fabio Massacci

2025 ACM/IEEE 47th International Conference on Software Engineering: New Ideas and Emerging Results







Funded by the European Union

## What Confident learning is all about

- IF dataset not so good  $\rightarrow$  use model to re-label it
- IF Score(x,1) < Confid(Score(x,1))  $\rightarrow$  label(x,0)
- and so on for 0

BUT

For vulnerabilities, model aren't good either...

- Wen et al. ICSE'23  $\rightarrow$  Best precision < 55%
- Croft et al. ICSE'23  $\rightarrow$  Wrong labels > 28%

## Our Idea: when in doubt throw them out

- Relabelling is too risky  $\rightarrow$  too many wrong elements
- Drop elements we are not sure  $\rightarrow$  they are confusing the model

Confidence level = weighted average of the predicted vulns

$$t_{\ell} = \frac{1}{\sum_{k} label(k, \ell)} \sum_{k} pred(k, \ell) \cdot label(k, \ell)$$
  
Well predicted  $\rightarrow pred(k, \ell) > t_{\ell} > \frac{1}{2}$   
Poorly predicted  $\rightarrow pred(k, \ell) \le \frac{1}{2}$   
The rest  $\rightarrow$  the one confusing the models

## **Two possible conditions**



## Two possible conditions $\rightarrow$ confidence of 0s



## And it works....

Running vulnerability dataset with 8K samples (50-50) on a BERT Model

	Data kept	Precision	Recall
Original Dataset	100%	63%	48%
Invert Mislabelled	100%	60%	70%
Remove Mislabelled	72%	59%	77%
Remove Confusing 0	89%	60%	75%
Remove Confusing 1	99%	59%	77%

DROP 1% of datapoints and Recall increases Without big change on Precision

# When in Doubt Throw It Out: Building on Confident Learning for Vulnerability Detection

#### What Confident learning is all about

IF dataset not so good  $\rightarrow$  use model to re-label it

- IF Score(x,1) < Confid(Score(x,1))  $\rightarrow$  label(x,0)

- and so on for 0

#### BUT

For vulnerabilities, model aren't good either...

- Wen et al. ICSE'23  $\rightarrow$  Best precision < 55%
- Croft et al. ICSE'23  $\rightarrow$  Wrong labels > 28%



#### Our Idea: when in doubt throw them out

- Relabelling is too risky  $\rightarrow$  too many wrong elements
- Drop elements we are not sure  $\rightarrow$  they are confusing the model

Confidence level = weighted average of the predicted vulns

$$t_{\ell} = \frac{1}{\sum_{k} label(k, \ell)} \sum_{k} pred(k, \ell) \cdot label(k, \ell)$$

Well predicted  $\rightarrow pred(k, \ell) > t_{\ell} > \frac{1}{2}$ Poorly predicted  $\rightarrow pred(k, \ell) \leq \frac{1}{2}$ The rest  $\rightarrow$  the one confusing the models

#### And it works ....

Running vulnerability dataset with 8K samples (50-50) on a BERT Model

	Data kept	Precision	Recall
Original Dataset	100%	63%	48%
Invert Mislabelled	100%	60%	70%
Remove Mislabelled	72%	59%	77%
Remove Confusing 0	89%	60%	75%
Remove Confusing 1	99%	59%	77%



#### DOI 10.5281/zenodo.14744047



#### Yuanjun Gong



#### Fabio Massacci