

## Automating the Identification of Safety Events Involving ML-Enabled Medical Devices



MACQUARIE  
University

AUSTRALIAN INSTITUTE  
OF HEALTH INNOVATION



Dr Ying WANG  
Macquarie University  
@YingWang\_CHI



Dr David LYELL  
Macquarie University  
@David\_Lyell



Prof Enrico COIERA  
Macquarie University  
@EnricoCOIERA

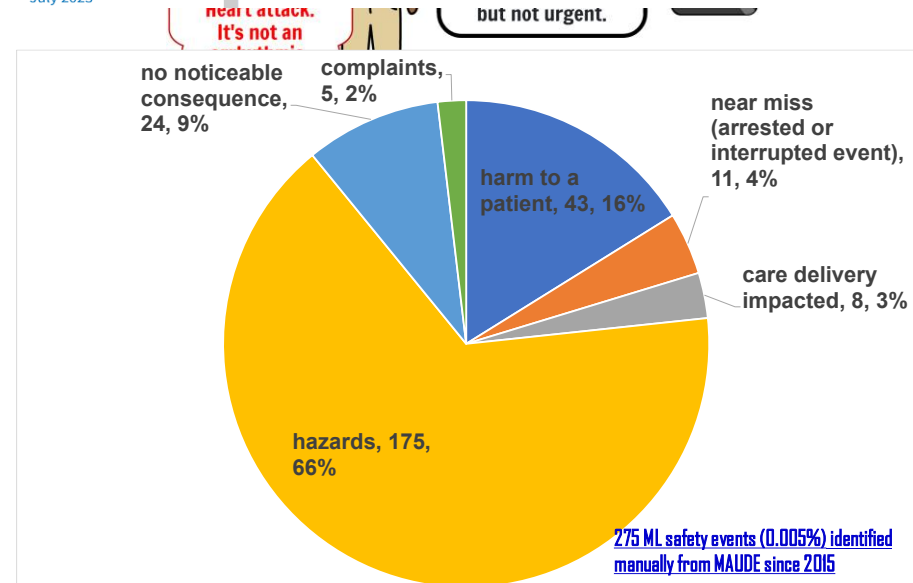
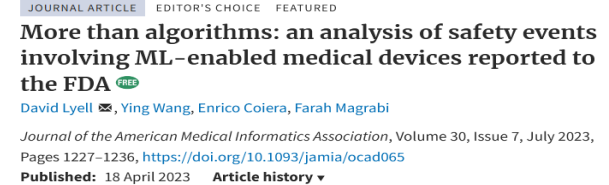


Prof Farah MAGRABI  
Macquarie University  
@Farah MAGRABI



# Background

1. Machine learning (ML) enabled systems have the potential to improve healthcare delivery.
2. With use in clinical settings and by consumers, safety events involving ML-enabled medical devices are emerging.
3. FDA MAUDE database is a valuable resource to understand patient safety risks associated with ML medical devices in real worlds.
4. Considering over 2 million events were reported to MAUDE in 2021, manual review on an ongoing basis is not feasible.



# Machine Learning driven solution



## Data collection

- Data processing: bag of words
- Numeric representation: term frequency-inverse document frequency (TF-IDF)

## Feature representation

- Event reports including updates and follow-up investigations
- Brand names plus reports
- Manufacture names plus reports

## Model training and testing

- Binary discriminative classifiers of SVM with radial-basis function kernel
- 10-fold subsampling cross-validation
- Four stratified classifiers with varied feature sets above

## Error analysis and verification

- Performance measures: precision, recall and F1 score
- Verification of classifier identified positives

# Training and testing datasets

Datasets	Event type	Stratified dataset (n=)
Training set (75%) <u>From 1 January 2018 to 31 October 2021</u>	ML safety events	207
	General safety events	4,044,796
	<b>Total</b>	<b>4,045,003</b>
Testing set (25%)	ML safety events	68
	General safety events	1,348,266
	<b>Total</b>	<b>1,348,334</b>
External testing set <u>From 1 November 2021 to 31 March 2022</u>		895,627

Data pre-processing



Feature representation



Model training



Model testing

# Results: performance of four classifiers

Dataset		Report classifier	Generic type classifier	Brand name classifier	Combined classifier
Testing set	False negative (n)	27	22	23	4
	False positive (n)	4	2	5	0
	<b>Total (n)</b>	<b>31</b>	<b>24</b>	<b>28</b>	<b>4</b>
	Precision (%)	94.4	97.1	93.2	100
	Recall (%)	71.6	75.6	74.7	93.9
	F1 score (%)	81.4	<b>85</b>	82.9	<b>96.9</b>
External testing set	Classifier identified MLSEs (n)	43	38	38	7,707
Verified as	MLSEs (n)	31	38	38	/
	Data input issues attributed to non-ML	5	0	0	/
	False positives (n)	7	0	0	/

Report classifier

Generic type classifier

Brand name classifier

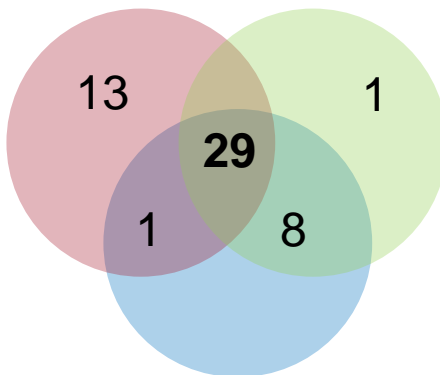
Combined classifier

Table: Performance of report classifier, generic type classifier, brand name classifier and combined classifier on testing and external datasets.

# Results: error analysis

## Classifier identifications

- **52** ML safety events identified by three classifiers on external testing set.
- **29** (56%) identified by all three involving known ML devices.
- **9** events identified by two classifiers
  - 8 events by brand name and generic type
  - 1 event by report and generic name
- **13** events identified by report alone and one by brand name



## True and false positives

- **40** of 52 events (77%) manually verified as true positives (TP).
- **5 false positives** (10%) did not involve ML devices but associated with device input problems, such as poor data acquisition.
- **7 false positives** (13%) neither HIT problems nor ML devices.
- All 29 events overlapped by the three classifiers were verified as TPs.
  - **Data acquisition problems:** 90% (n=26) from imaging systems.
  - **Algorithm errors:** two TPs (6.9%) with an ECG device designed for detecting normal sinus rhythm and severe arrhythmias.

## Feasibility

Model performance

- Stratified text classifiers can be used to identify rare ML events from large collections of data.

## Generalizability

Real-world use of AI

- Current classifiers don't generalize well to identify rare class events from highly imbalanced dataset.

## Device input problems

Common safety risks

- Mostly identified by classifiers and made up the bulk of the training set.

## Algorithm errors

Critical safety risk for ML devices

- One user's device suggested atrial fibrillation during heart attack.
- Another user experienced atrial fibrillation, but the device indicated a normal rhythm.

# Future work

## Active learning

Data annotation

- iteratively ask safety experts to verify events identified by the classifier and label more ML safety events.

## Large Language Models

Generative AI

- pre-trained or fine-tuned on a dataset of safety events to identify ML safety events.

## SafeHealth.ai

Knowledge sharing

- Insights on medical AI devices
  - Safety risk analysis of ML events



**MACQUARIE**  
University

AUSTRALIAN INSTITUTE  
OF HEALTH INNOVATION

**Thank you and comments?**

**Dr Ying WANG**

Macquarie University

@YingWang\_CHI

E: [ying.wang@mq.edu.au](mailto:ying.wang@mq.edu.au)

W: [safehealth.ai](http://safehealth.ai)

