# Robust Performance Metrics for Authentication Systems
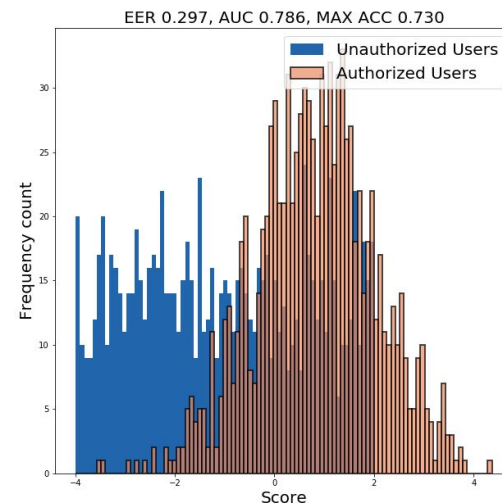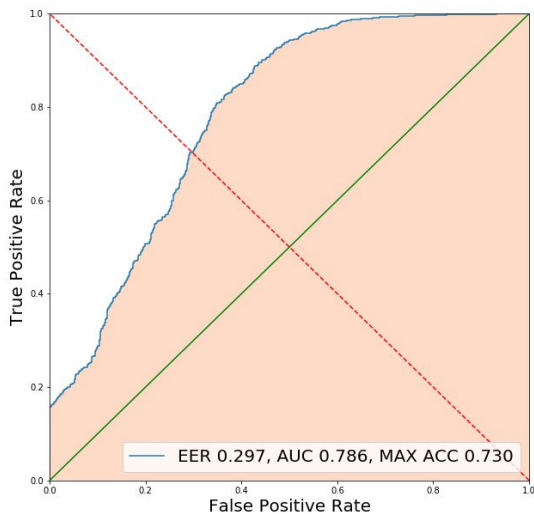
**Shridatt Sugrim**, Can Liu, Meghan C. McLean, Janne Lindqvist
Rutgers University
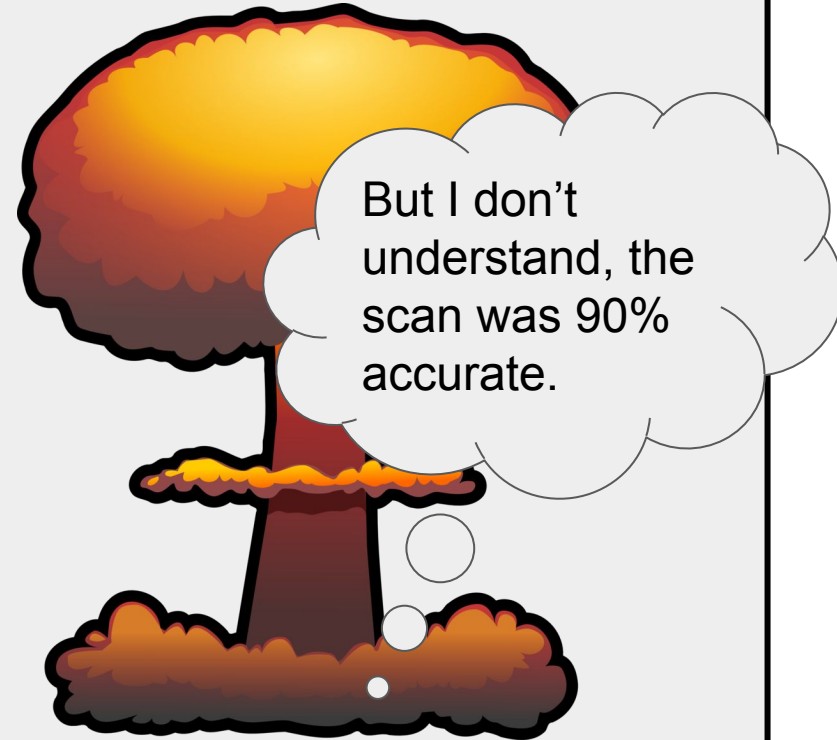
# TL; DR:
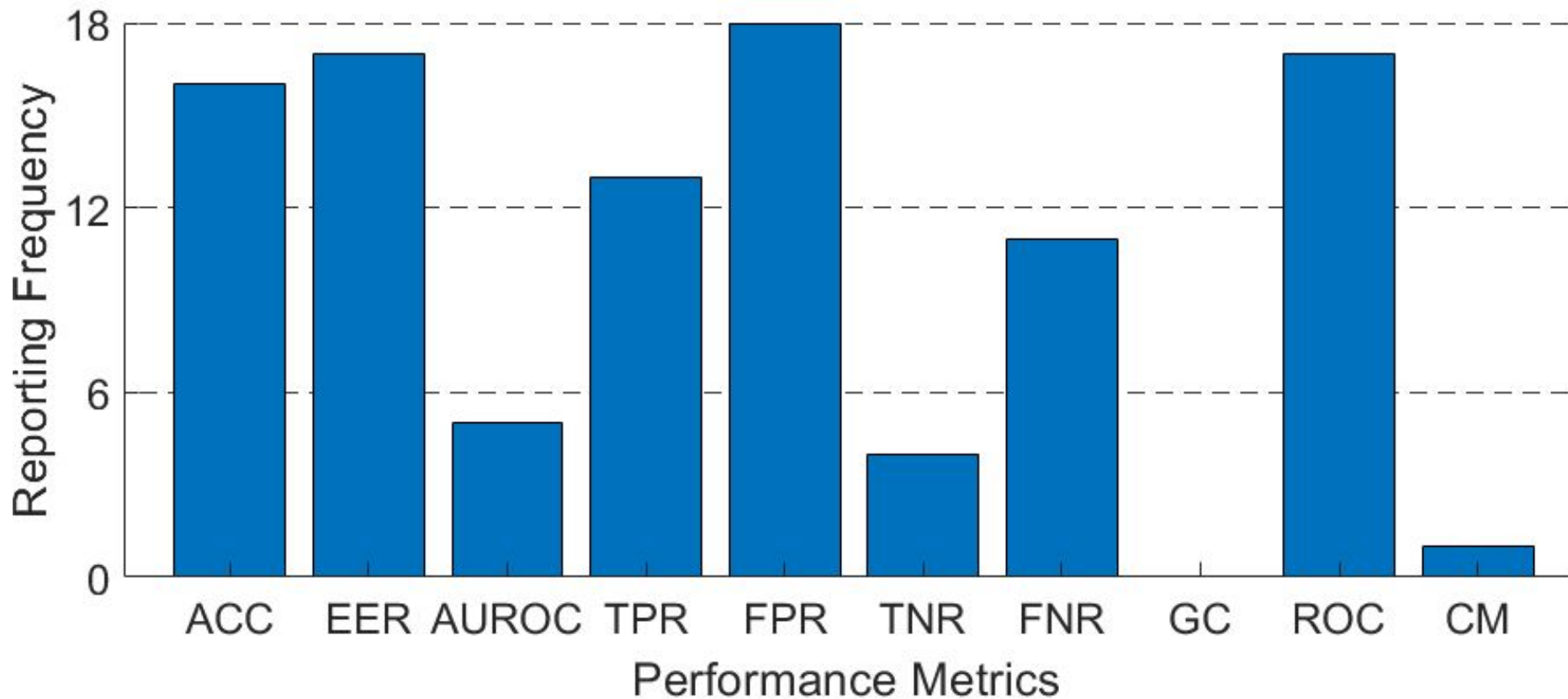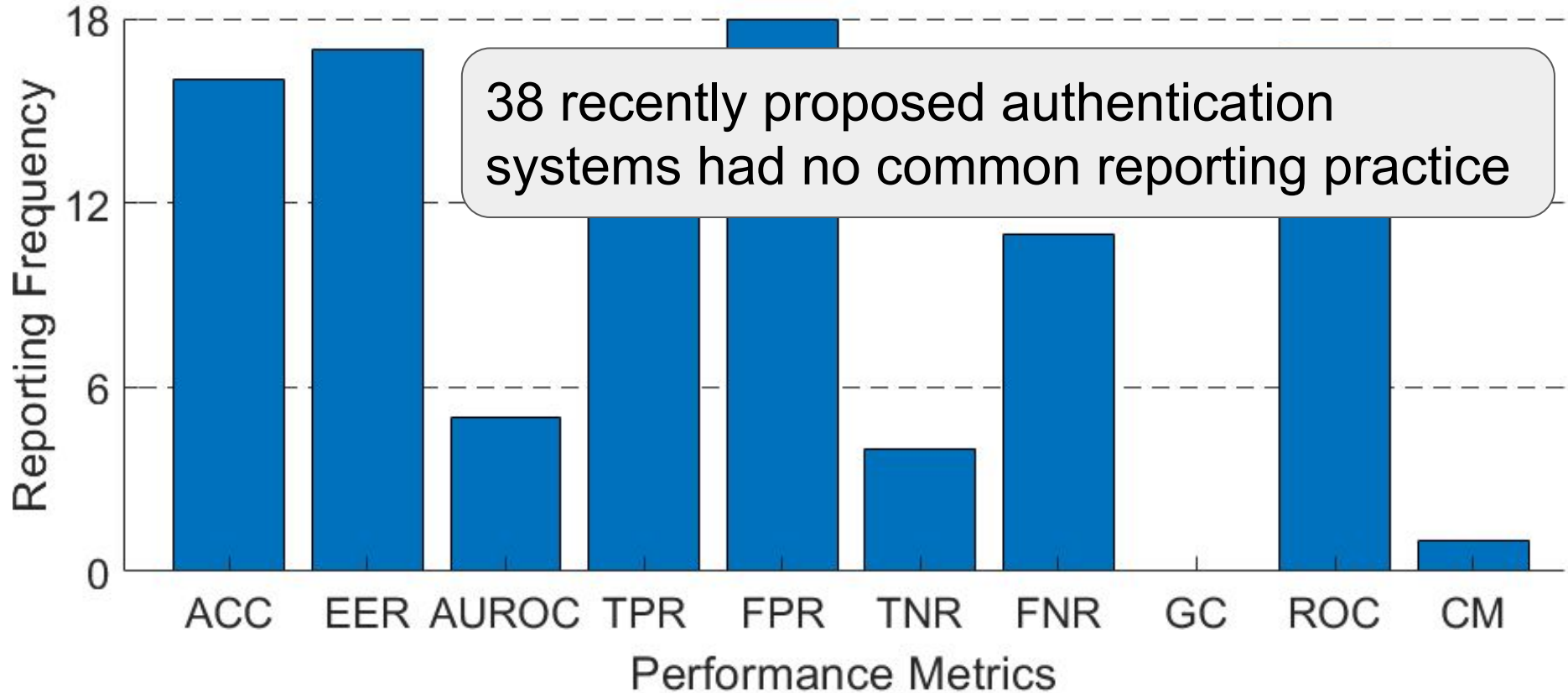## Performance reporting might surprise you

# TL; DR:
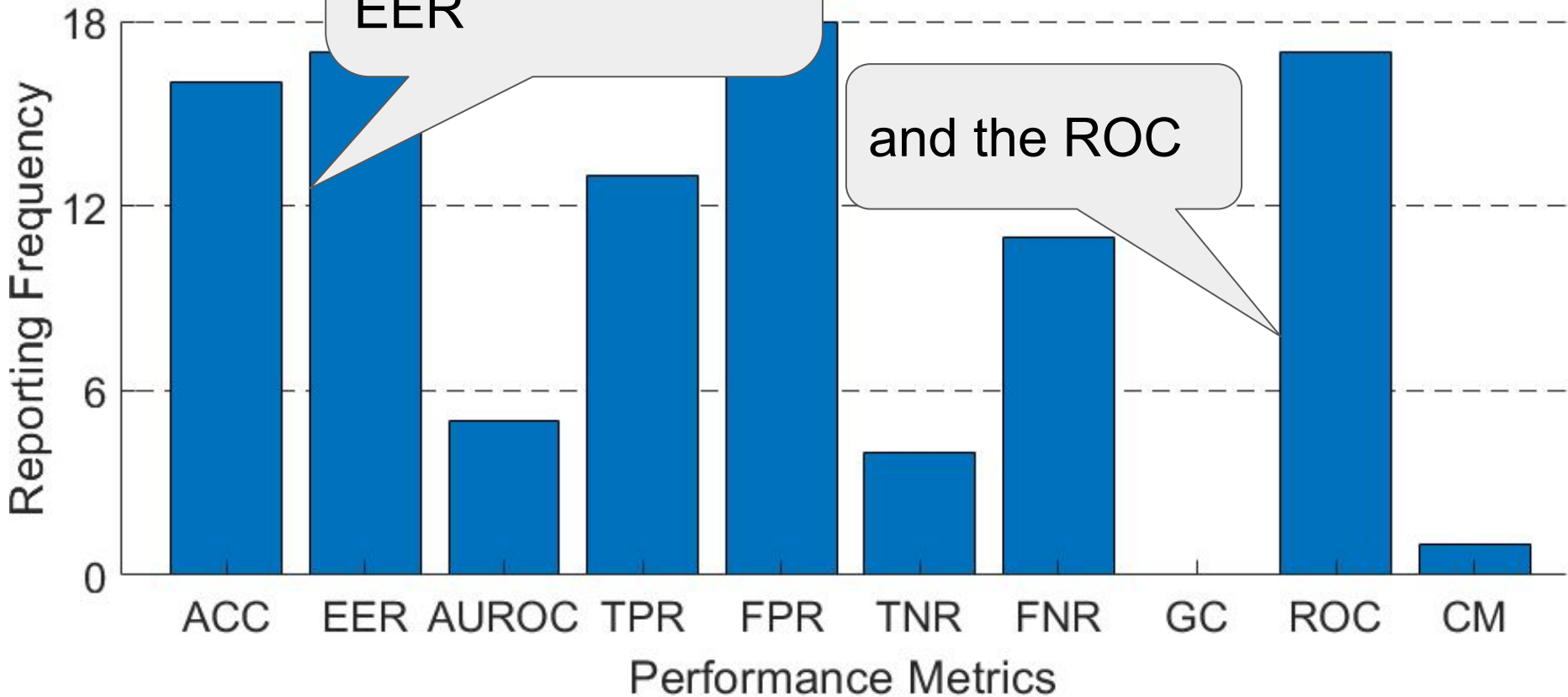# Performance reporting might surprise you

# Why is evaluation of authentication systems hard?
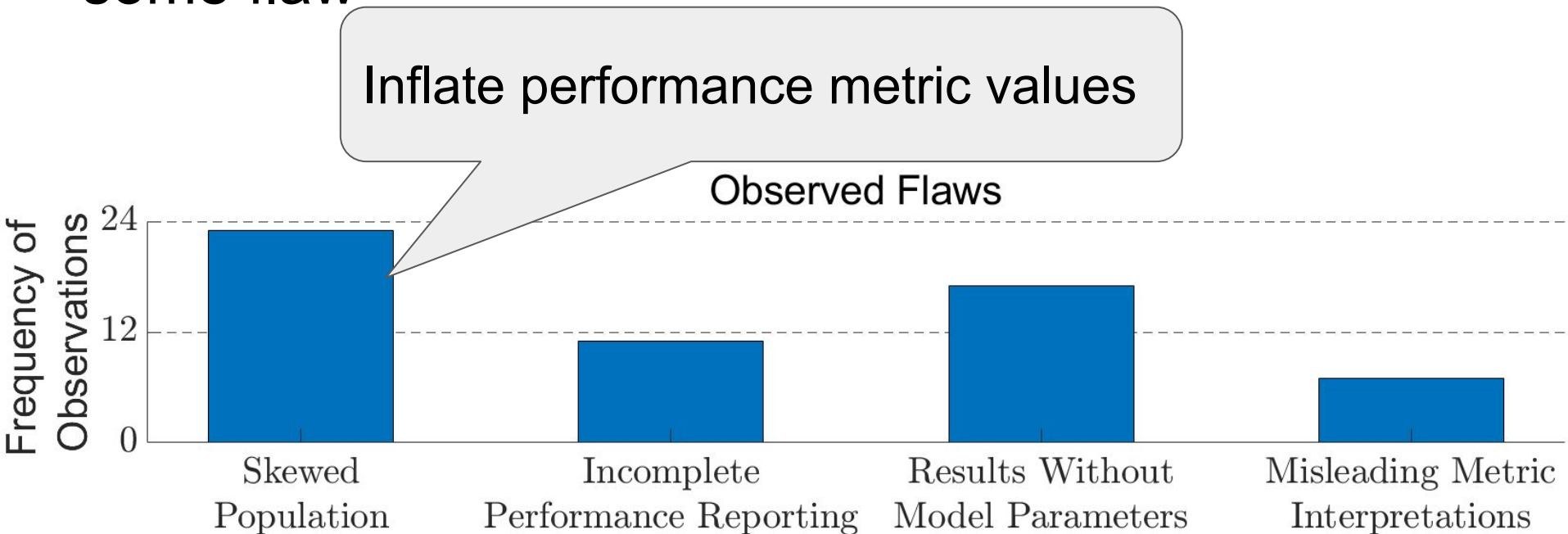
# Why is evaluation of authentication systems hard?



38 recently proposed authentication systems had no common reporting practice

# Why is ev̲̲̲̲̲̲̲̲̲̲tication systems hard?



6

# Most (36 of 38) of the performance reporting had some flaw

# Most (36 of 38) of the performance reporting had some flaw



Inflate performance metric values

Observed Flaws

# Most (36 of 38) of the performance reporting had some flaw

# Most (36 of 38) of the performance reporting had some flaw

Prevents reproducing results

### Observed Flaws

Frequency of Observations

24

12

0

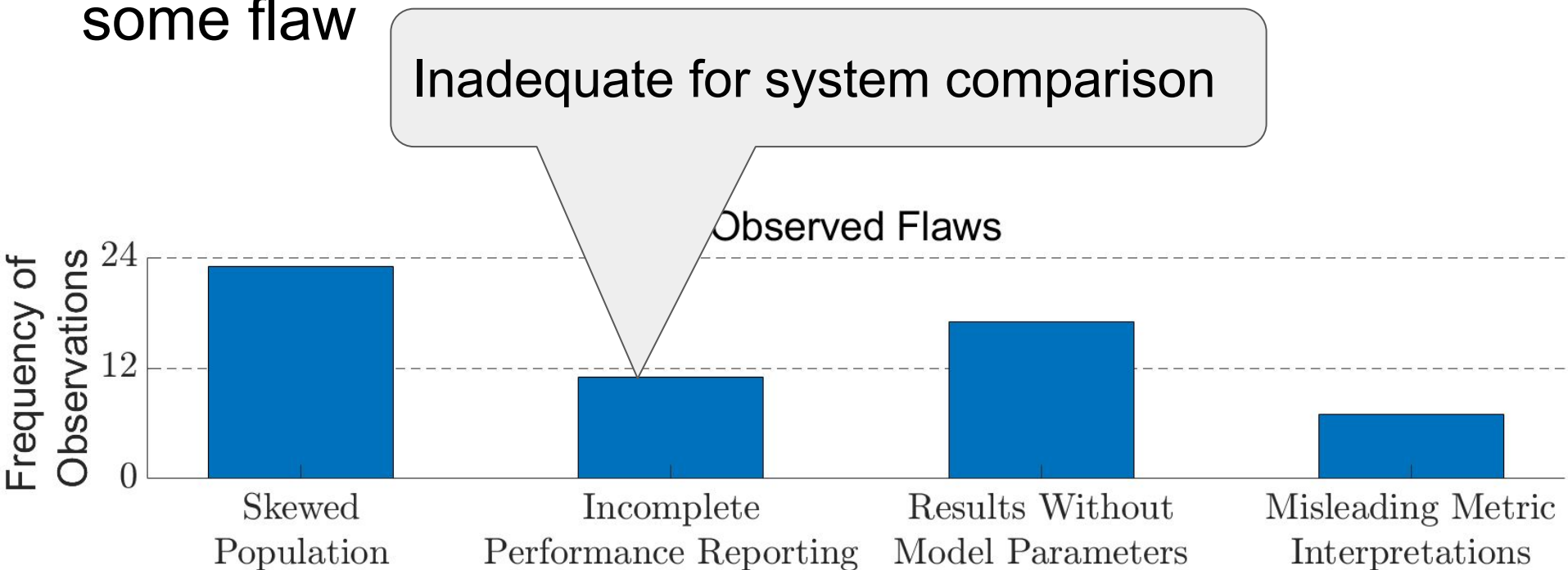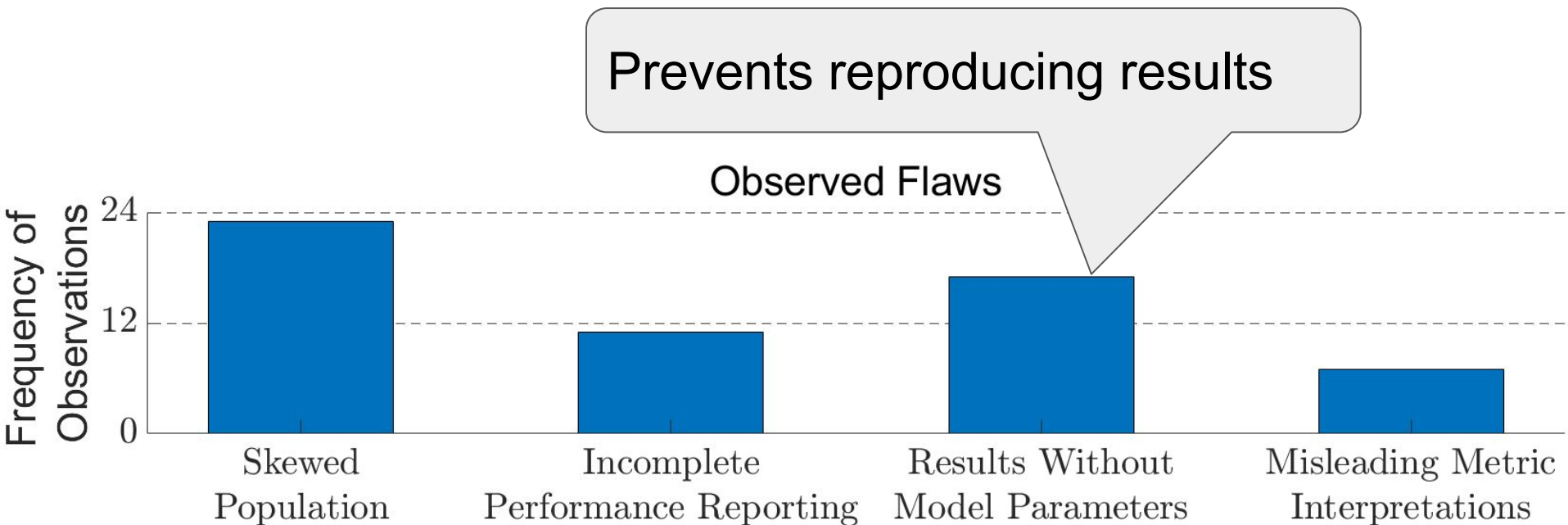Skewed Population | Incomplete Performance Reporting | Results Without Model Parameters | Misleading Metric Interpretations
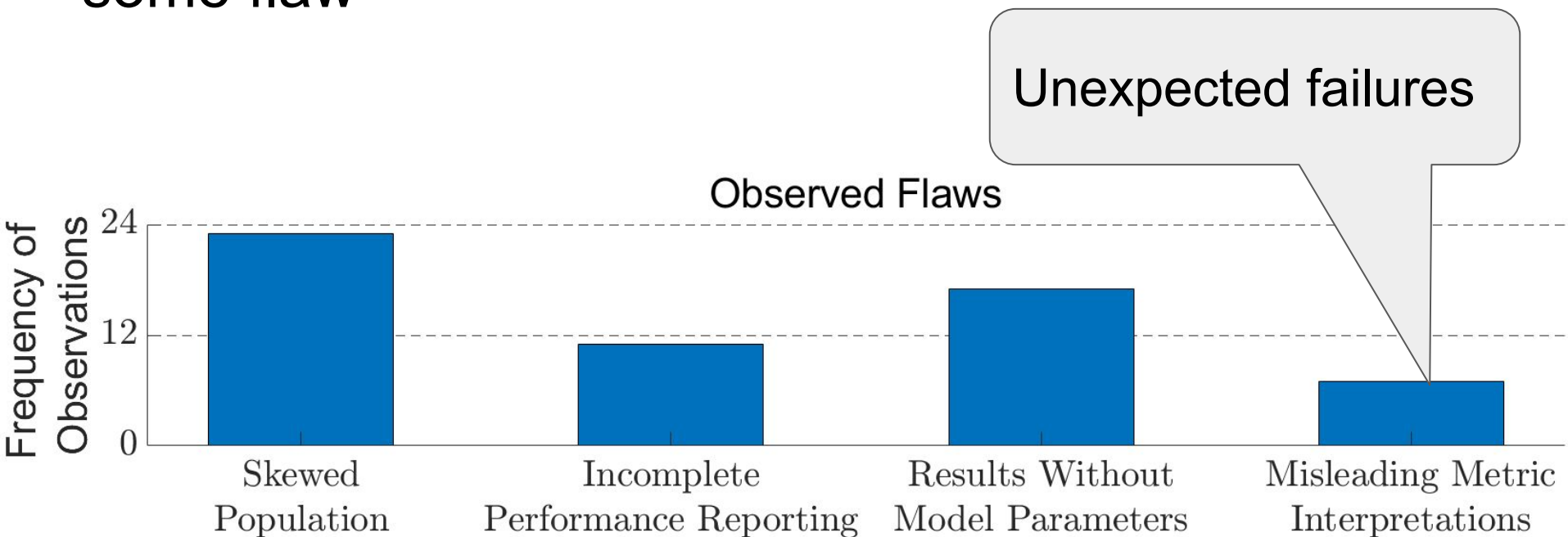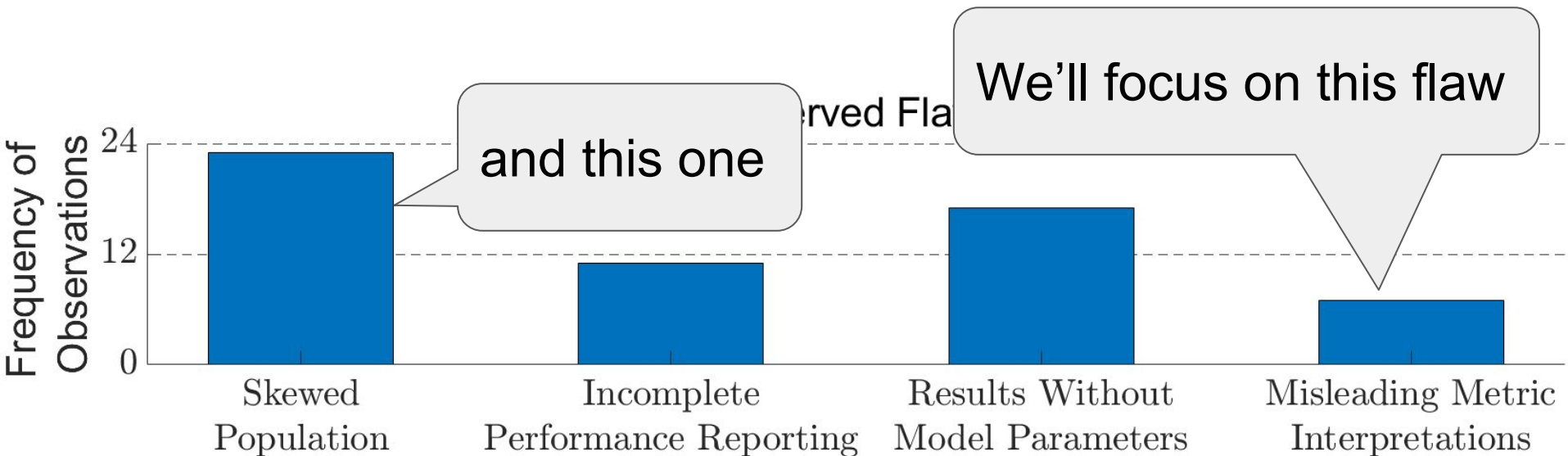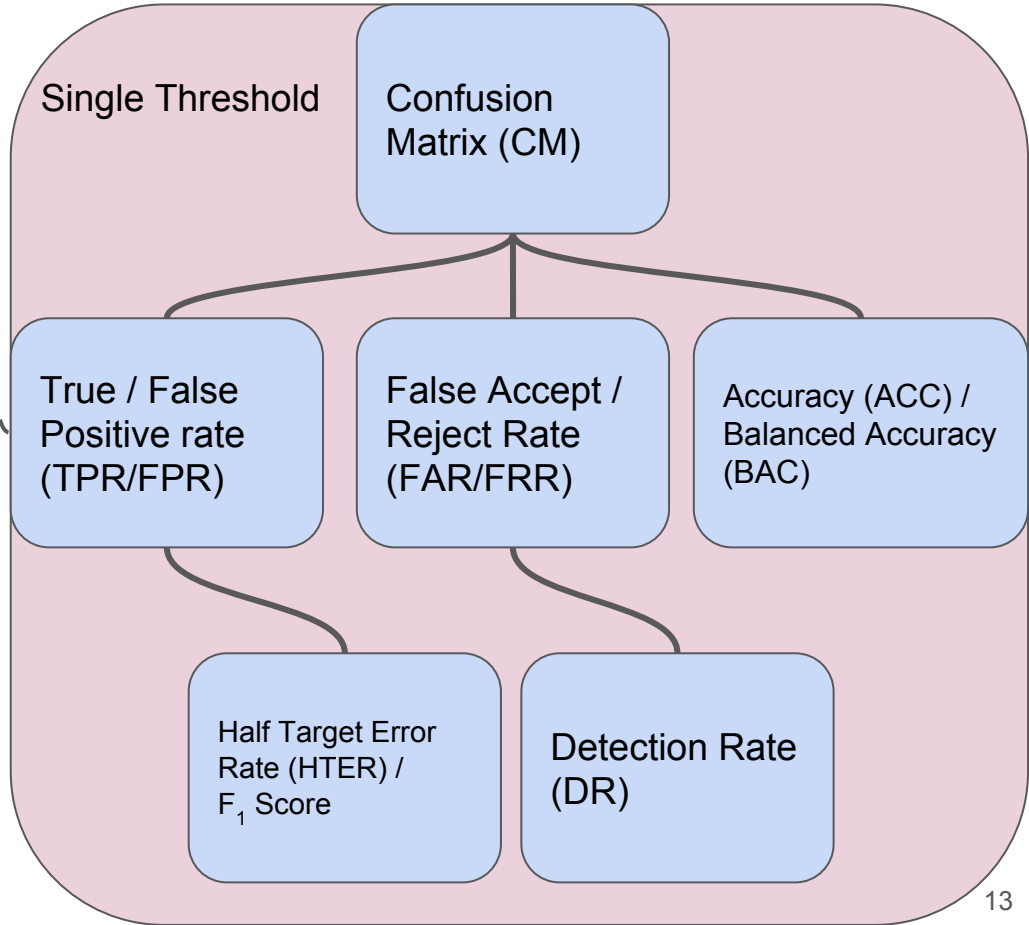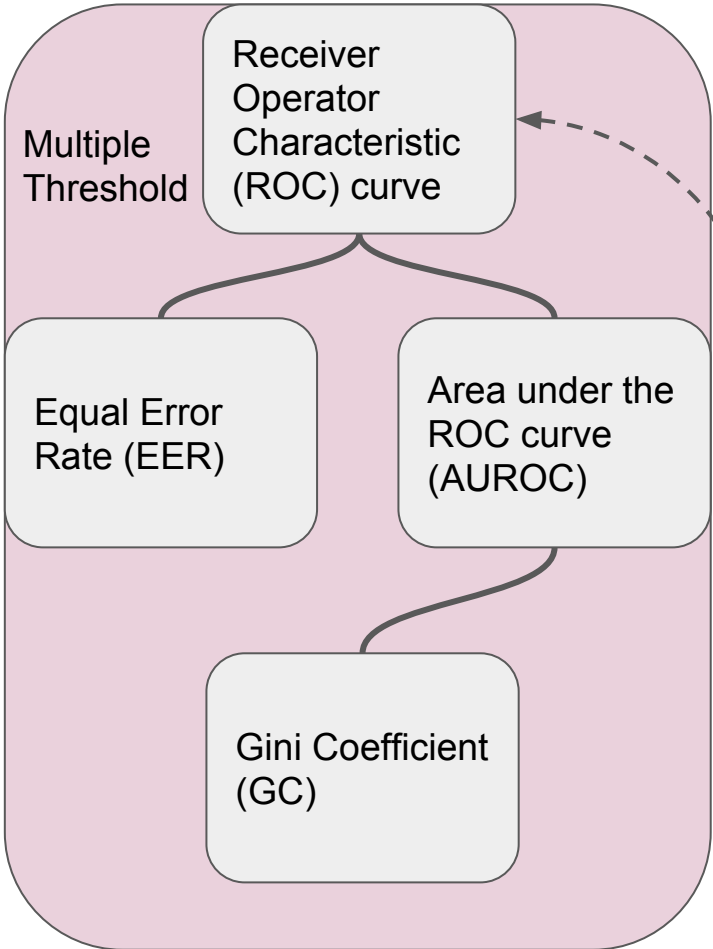
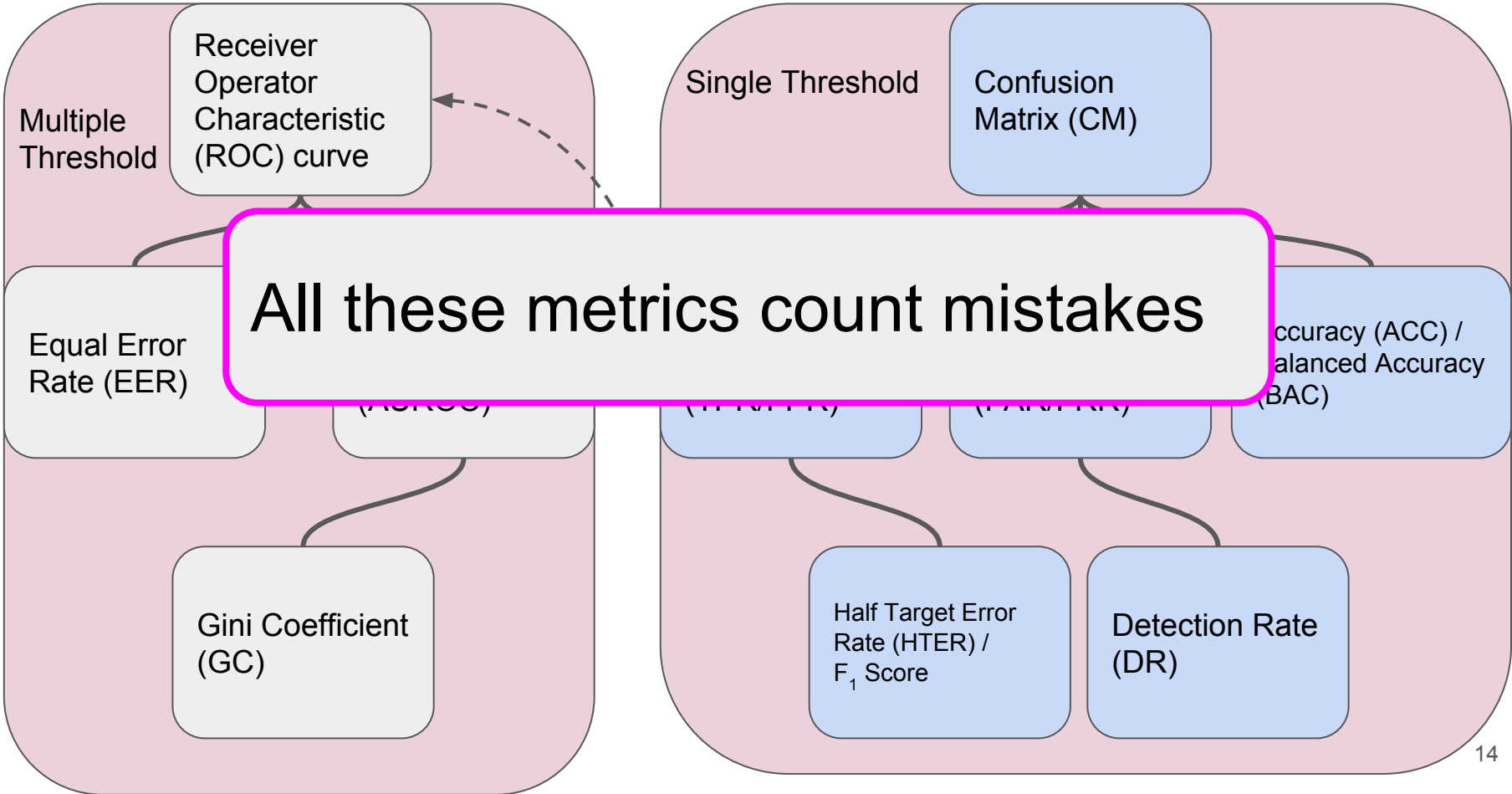# Most (36 of 38) of the performance reporting had some flaw

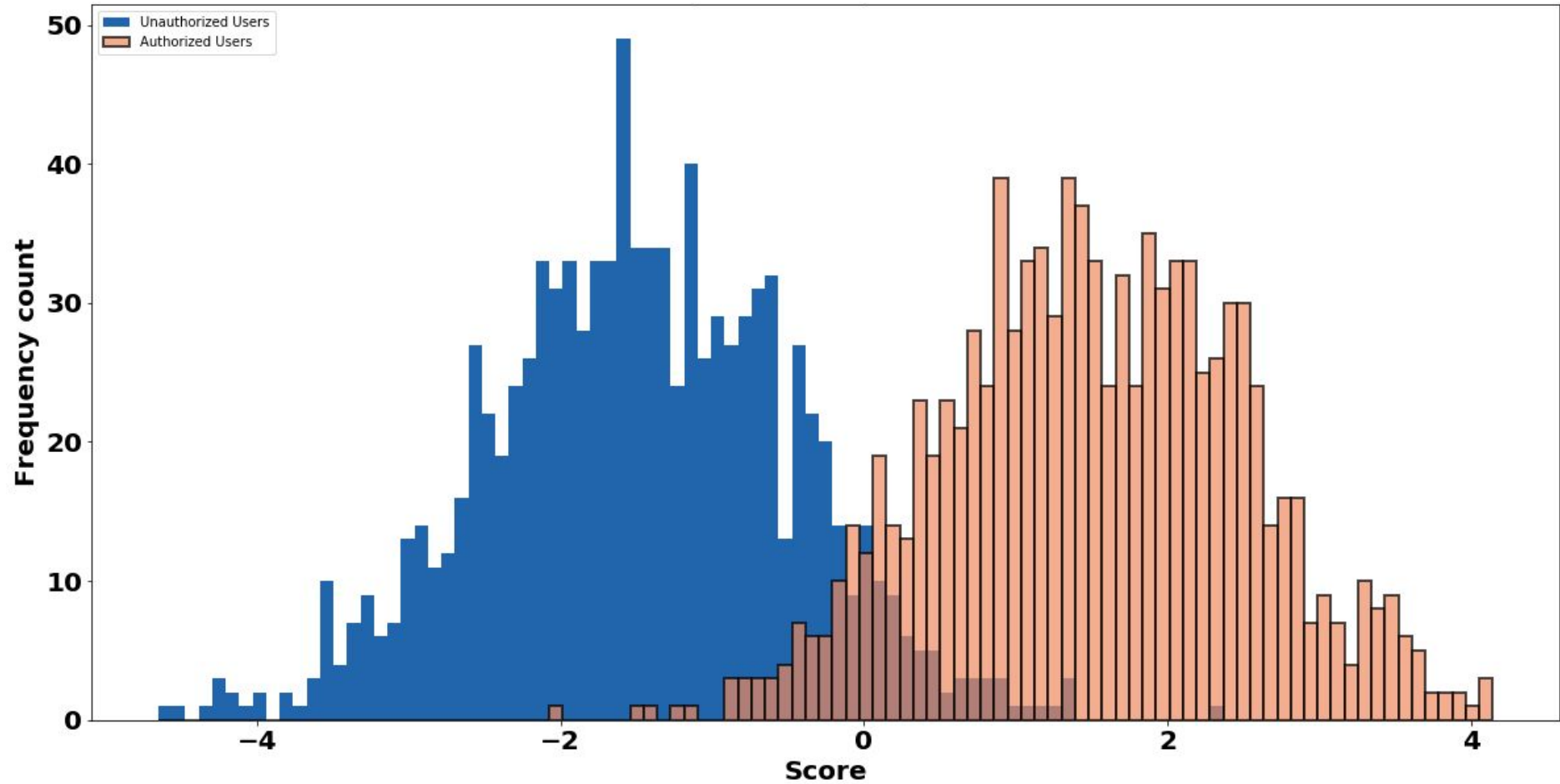# Most (36 of 38) of the performance reporting had some flaw

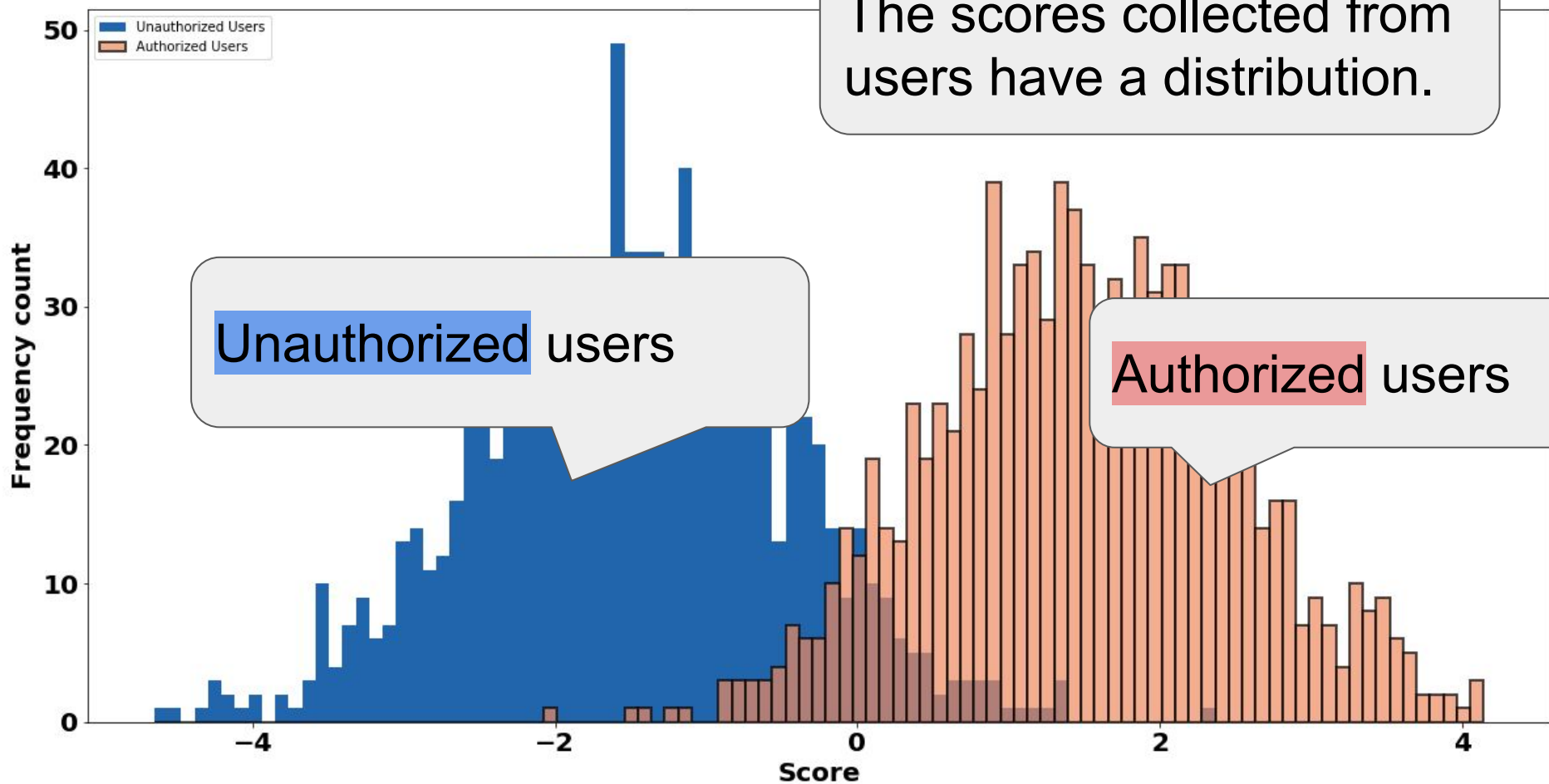# Related metrics have similar properties

# Related metrics have similar properties.



Multiple Threshold

Receiver Operator Characteristic (ROC) curve

Equal Error Rate (EER)

(AuROC)

Gini Coefficient (GC)

Single Threshold

Confusion Matrix (CM)

(TPR/FPR)

(FAR/FRR)

Accuracy (ACC) / Balanced Accuracy (BAC)

Half Target Error Rate (HTER) / $F_1$ Score

Detection Rate (DR)

## All these metrics count mistakes

# How authentication systems work
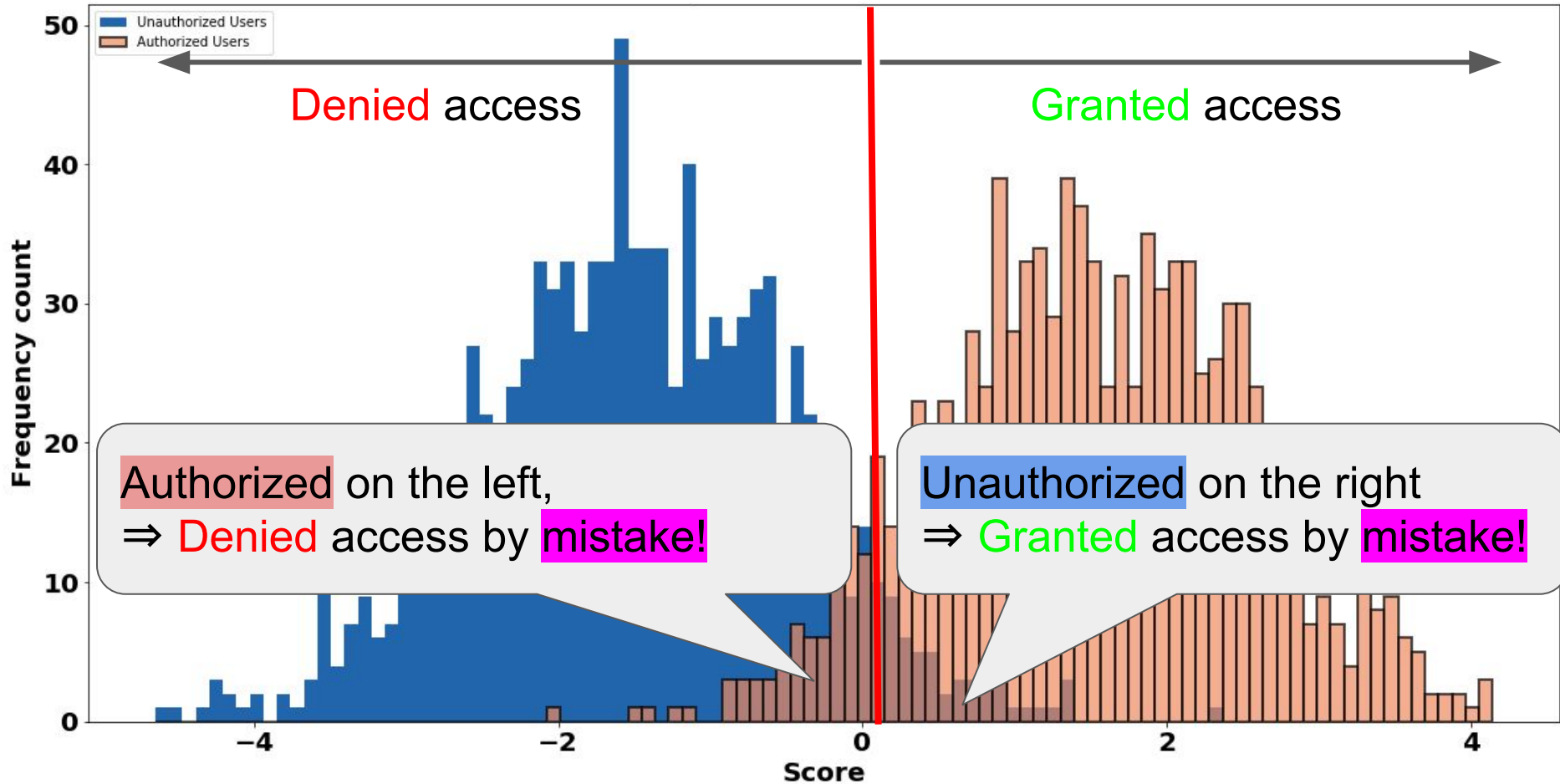
# How authentication systems work

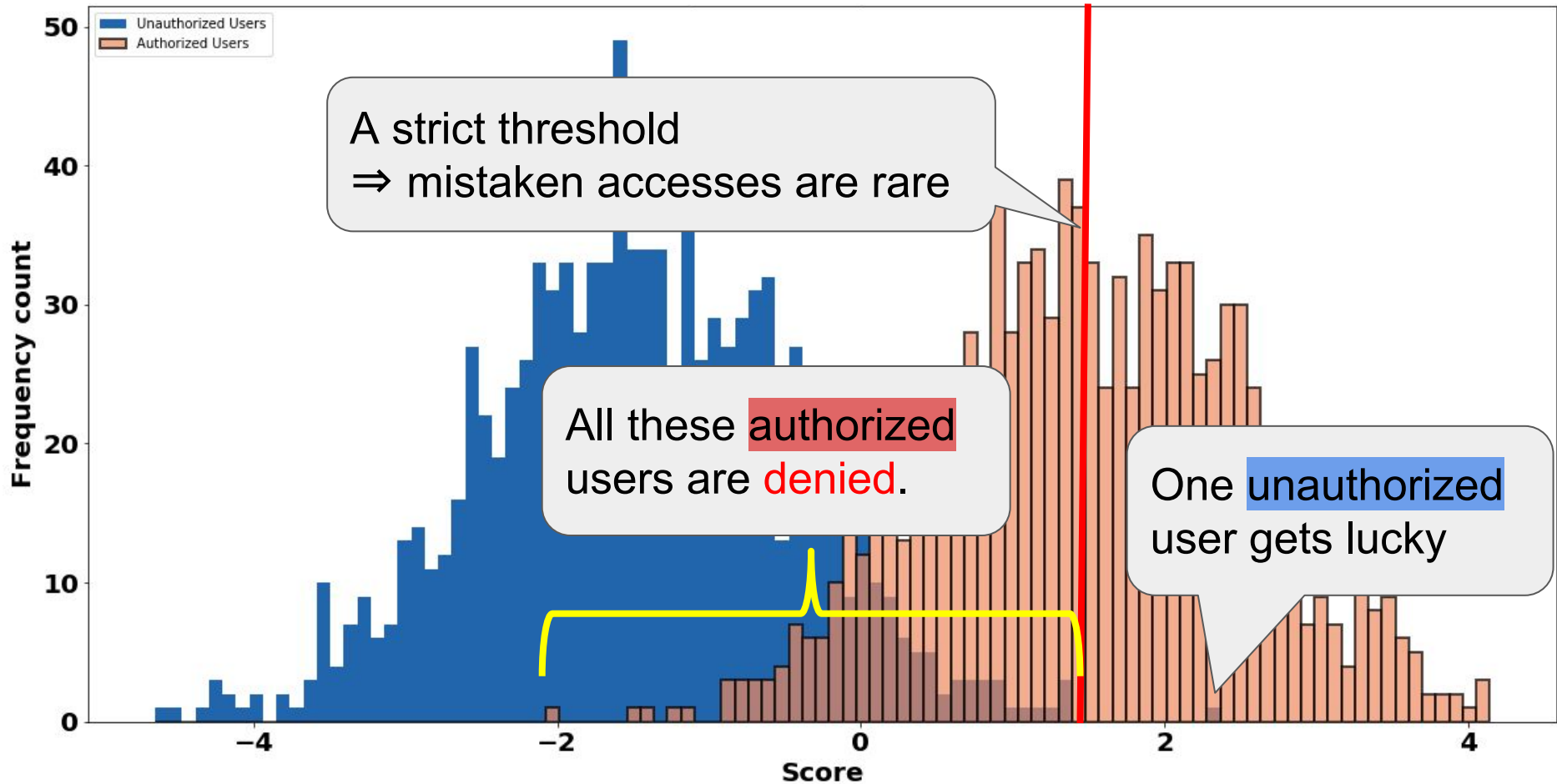# How authentication systems work
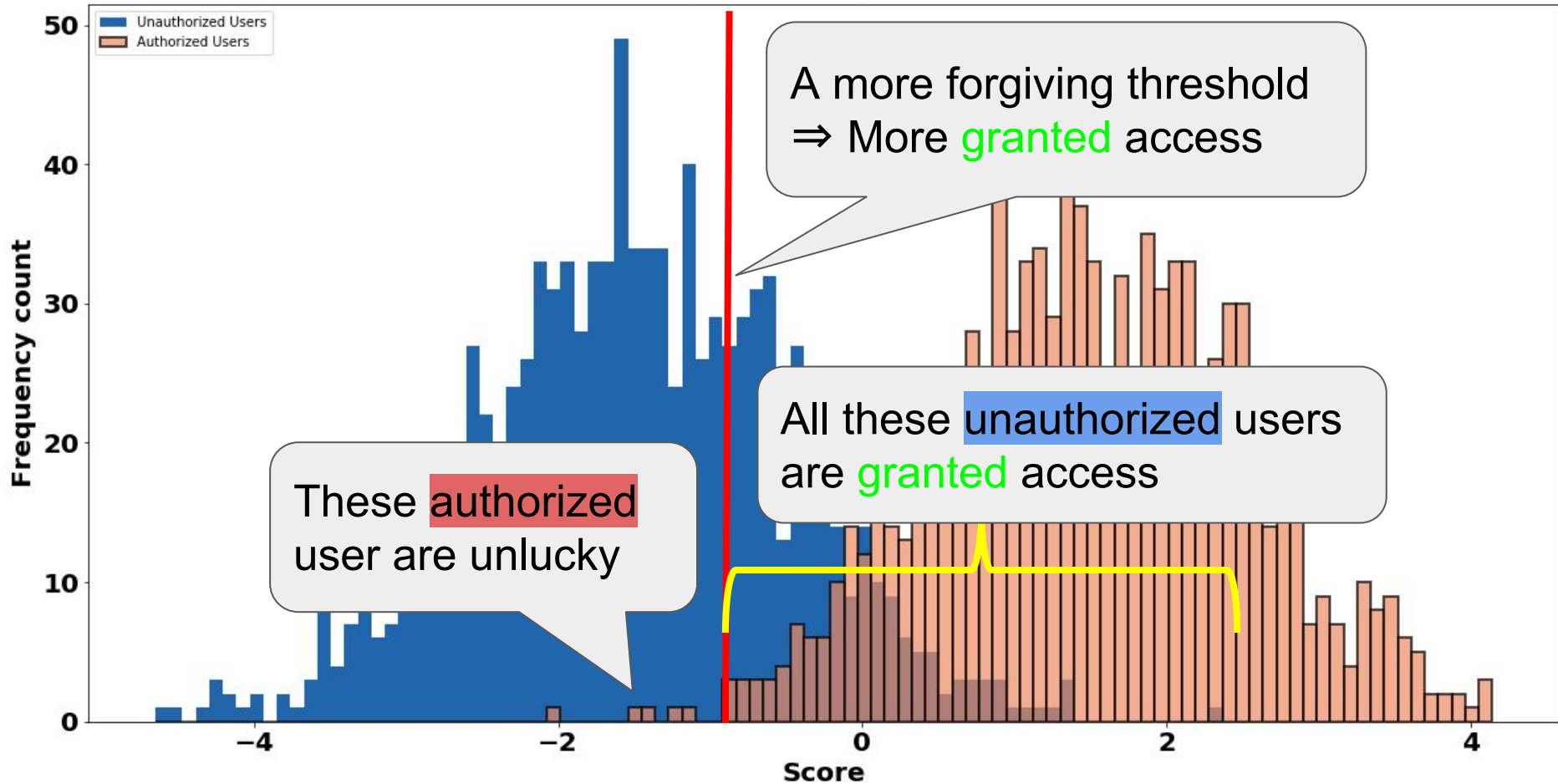
# How authentication systems work

# Where do authentication mistakes come from?

# Thresholds matter

# Thresholds matter

# Thresholds matter



Reporting metrics with one threshold
⇒ Inadequate for system comparisons
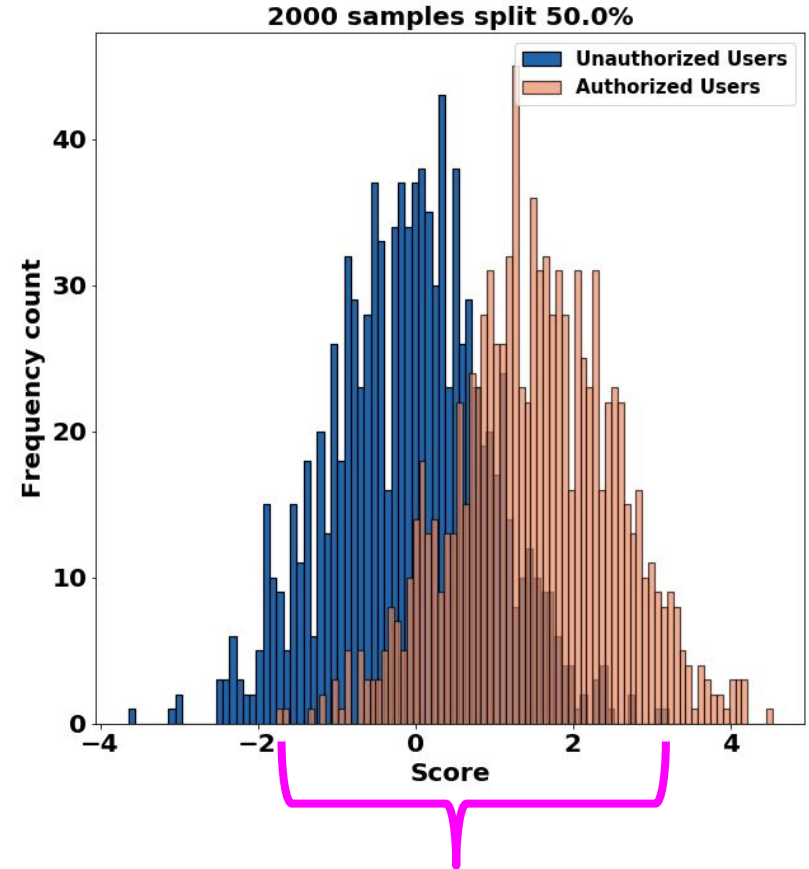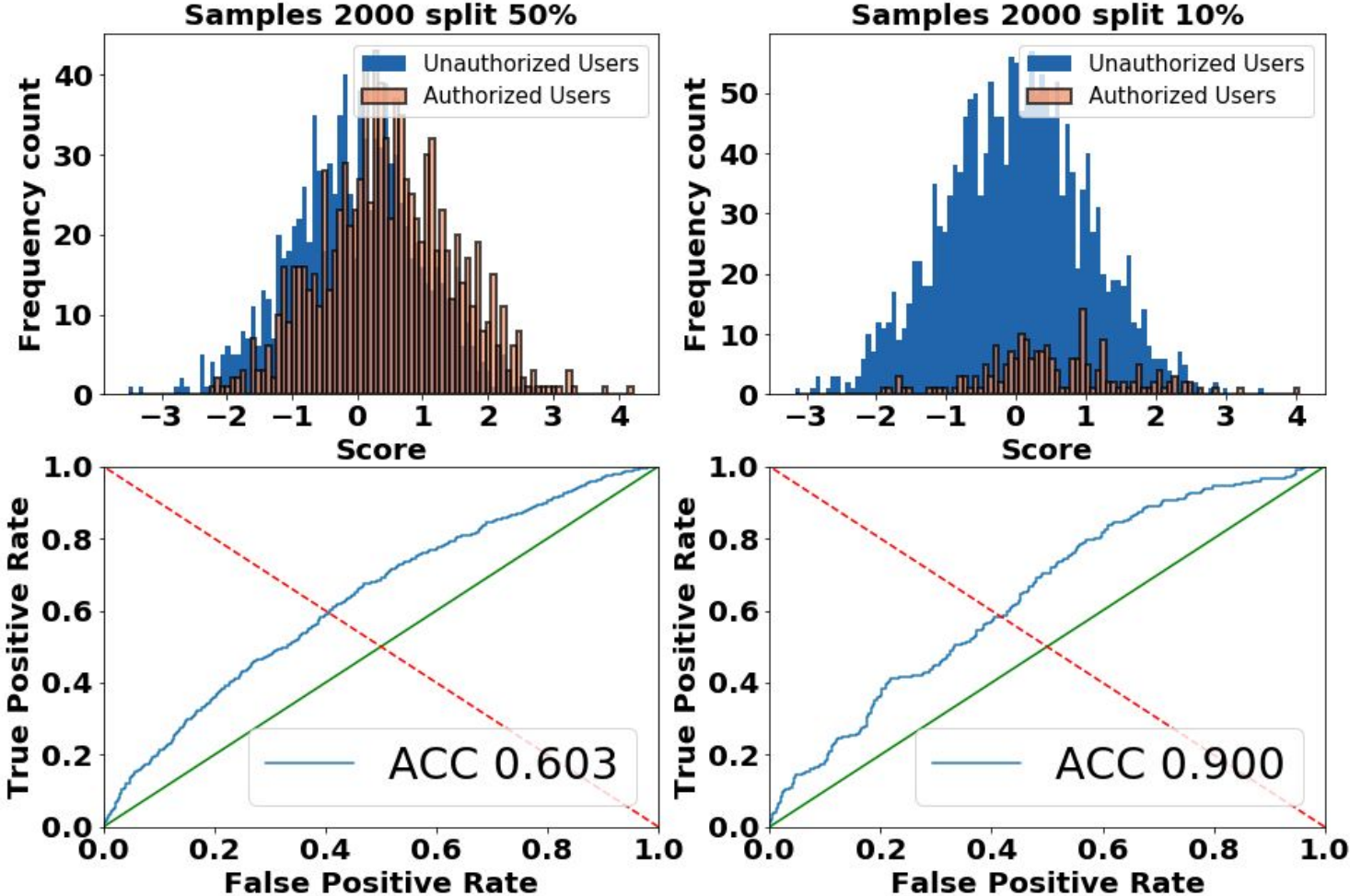
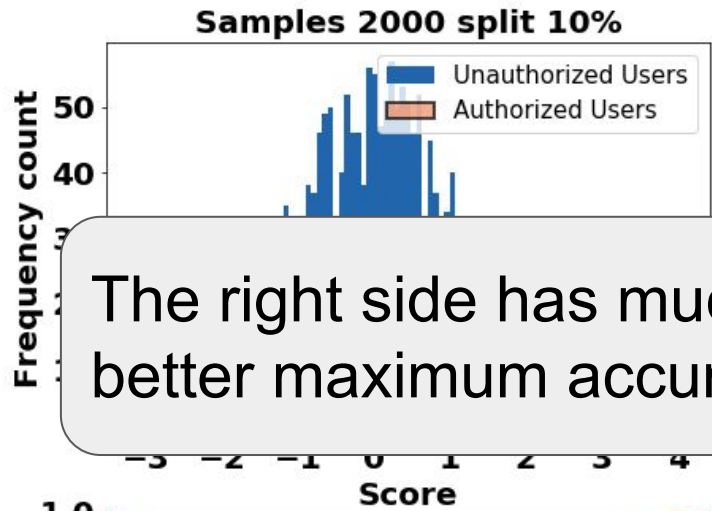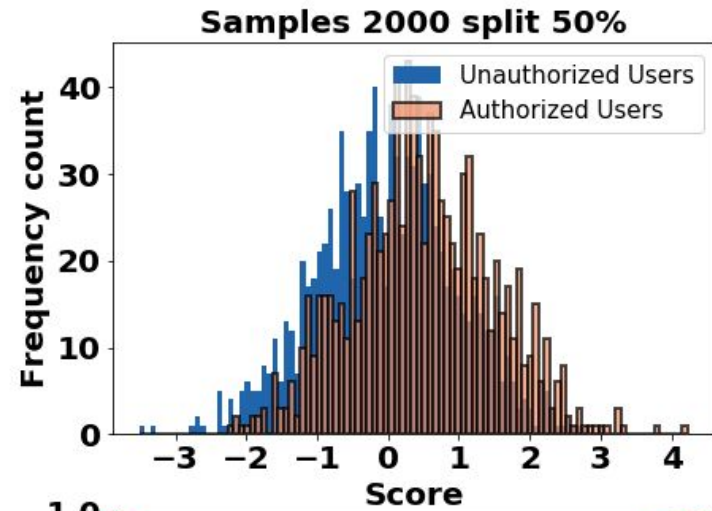# Frequency Count of Scores (FCS) helps visualize problems

- The distribution of scores plays an important role in the system performance
- **The potential for error is directly proportional to the width of the score overlap**
- The FCS can be used to identify problems with scoring
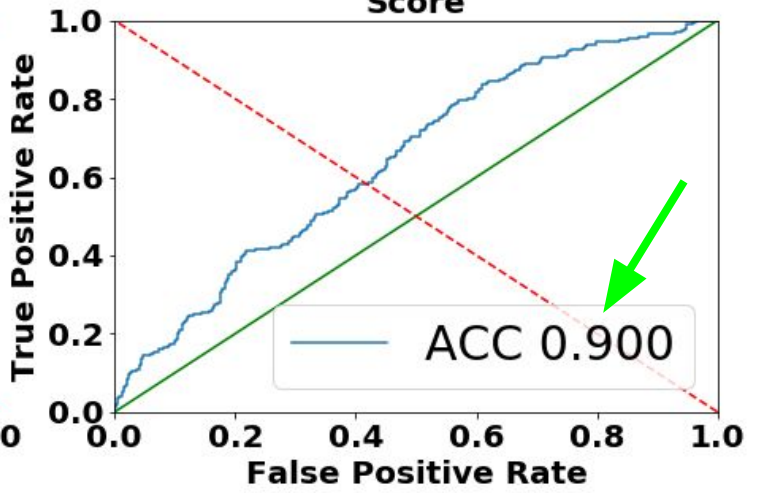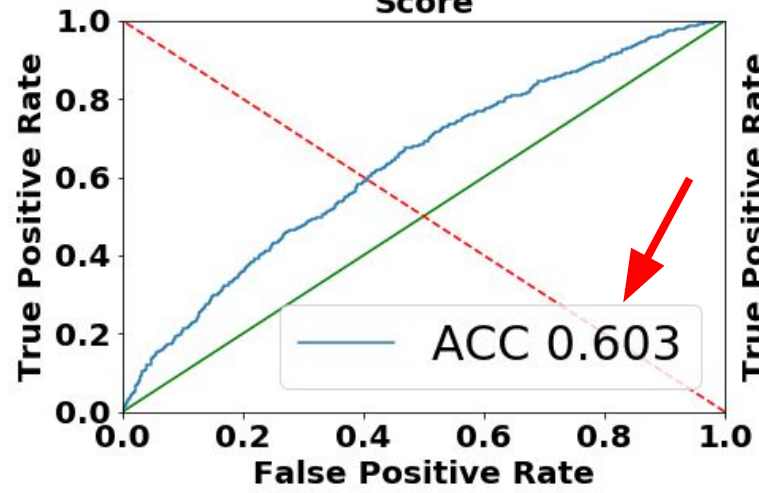


2000 samples split 50.0%

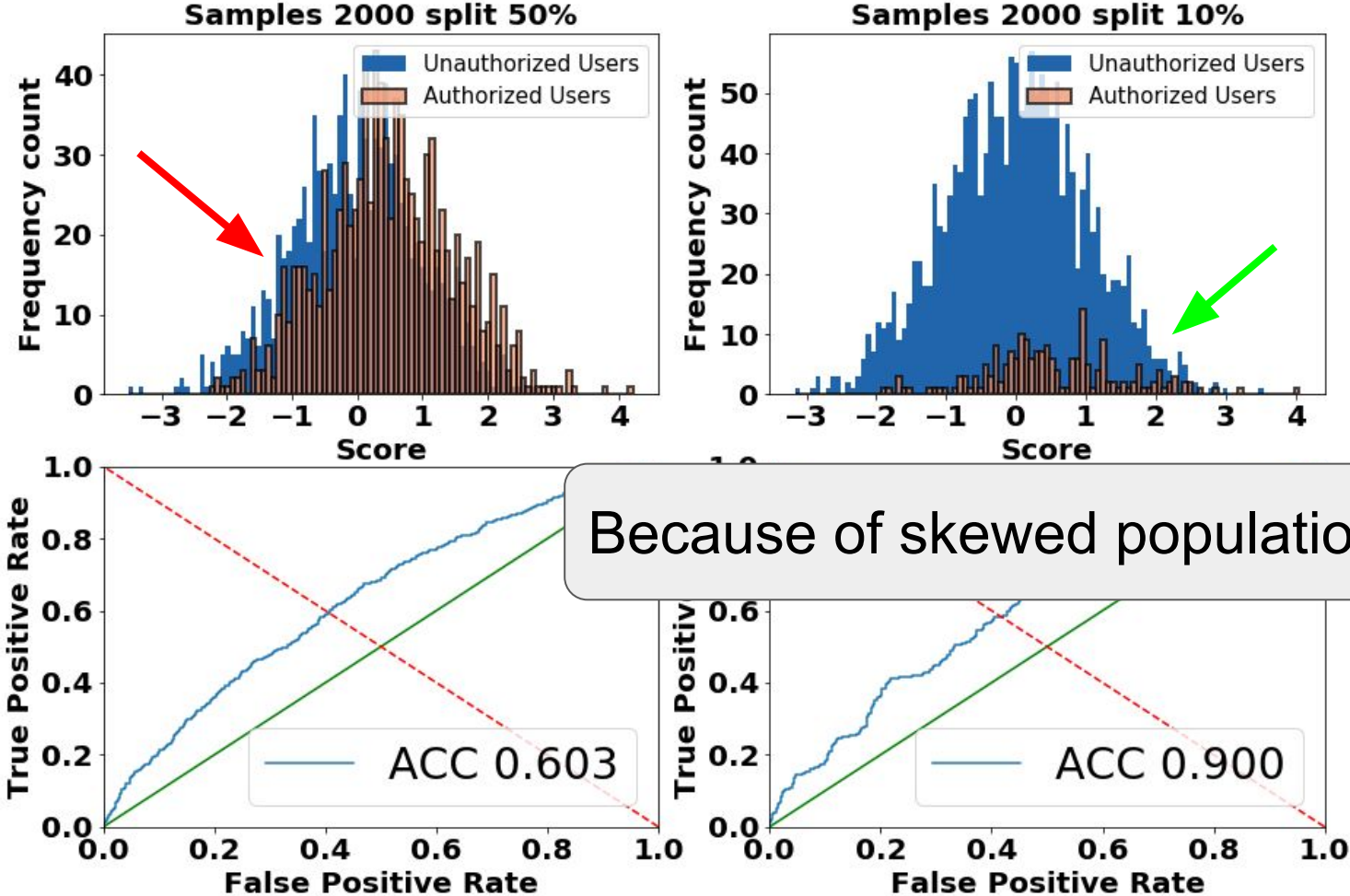# Skewed populations make accuracy unreliable

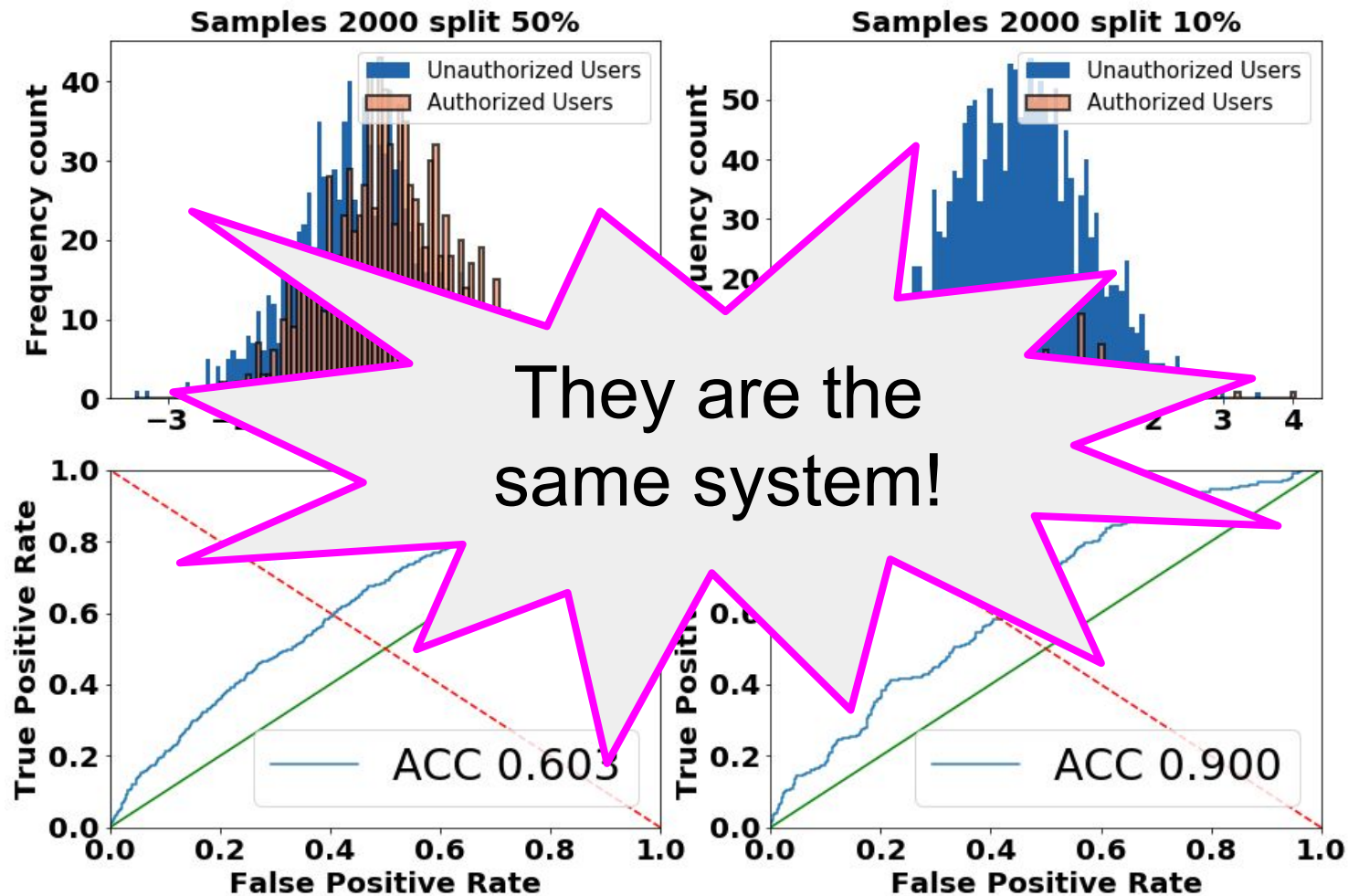# Skewed populations make accuracy unreliable



The right side has much better maximum accuracy!

# Skewed populations make accuracy unreliable
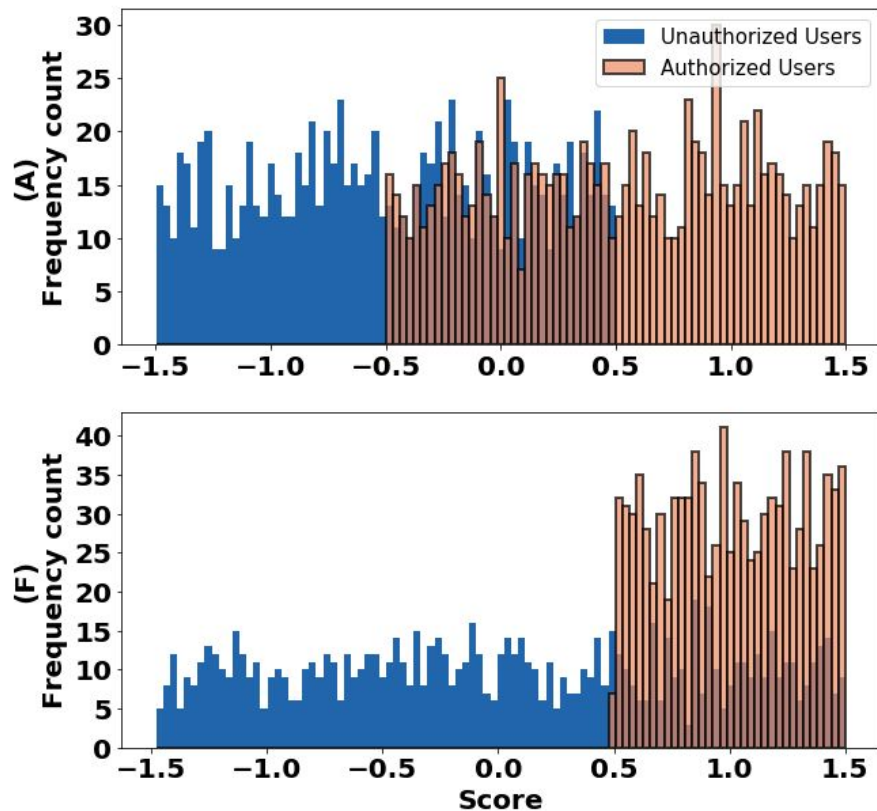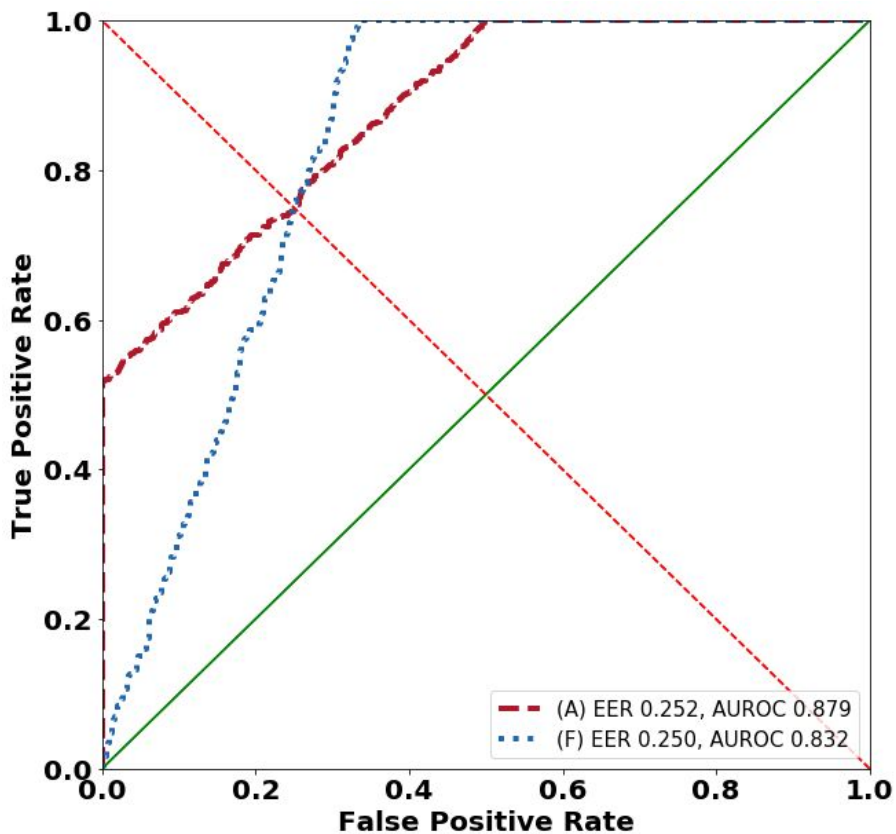


Samples 2000 split 50% — ACC 0.603

Samples 2000 split 10% — ACC 0.900

Because of skewed population

# Skewed populations make accuracy unreliable
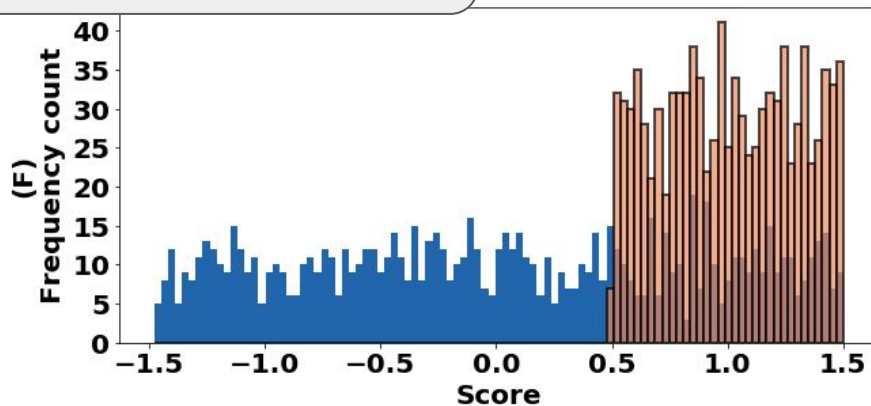


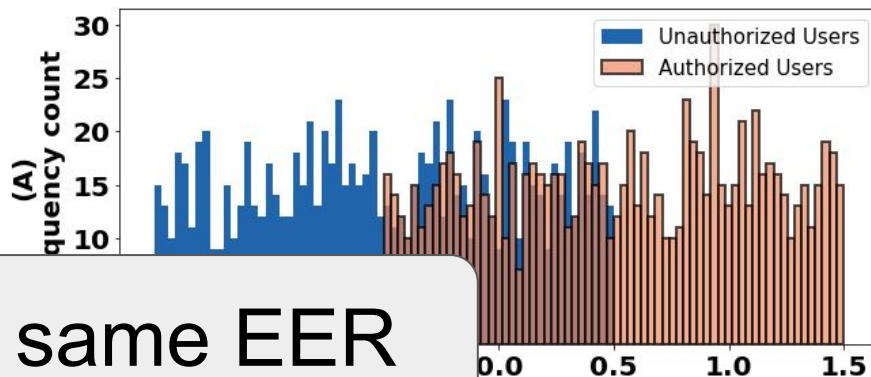They are the same system!

# Similar EER does not mean similar performance

# Similar EER does not mean similar performance



About the same EER

(A) EER 0.252, AUROC 0.879
(F) EER 0.250, AUROC 0.832

# Similar EER does not mean similar performance



Very different TPR at FPR 0.1

(A) EER 0.252, AUROC 0.879
(F) EER 0.250, AUROC 0.832

# Similar EER does not mean similar performance



This difference in performance is due to how the scores overlap

# EERs hides performance tradeoffs

. Compared 3 systems: SVC2004, Keystroke, Touchalytics

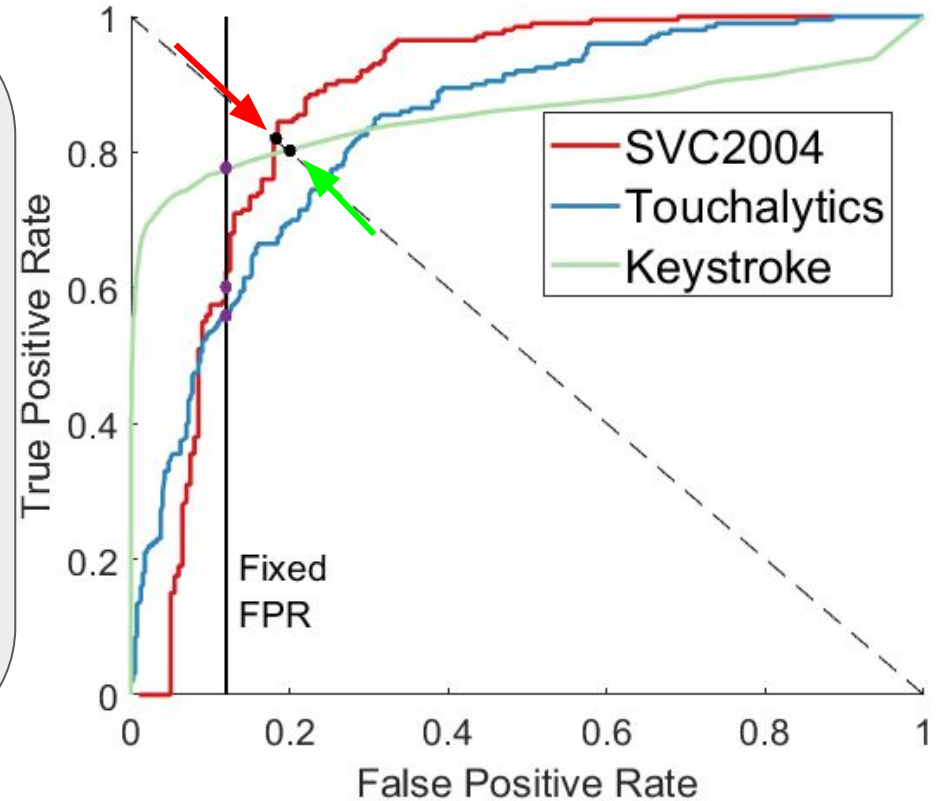. EER is misleading: SVC2004 and Keystroke have similar EERs – Keystroke has FPR of 0.1

# EERs hides performance tradeoffs

SVC2004 has a lower EER making it look better than Keystroke

SVC2004 EER = 0.185
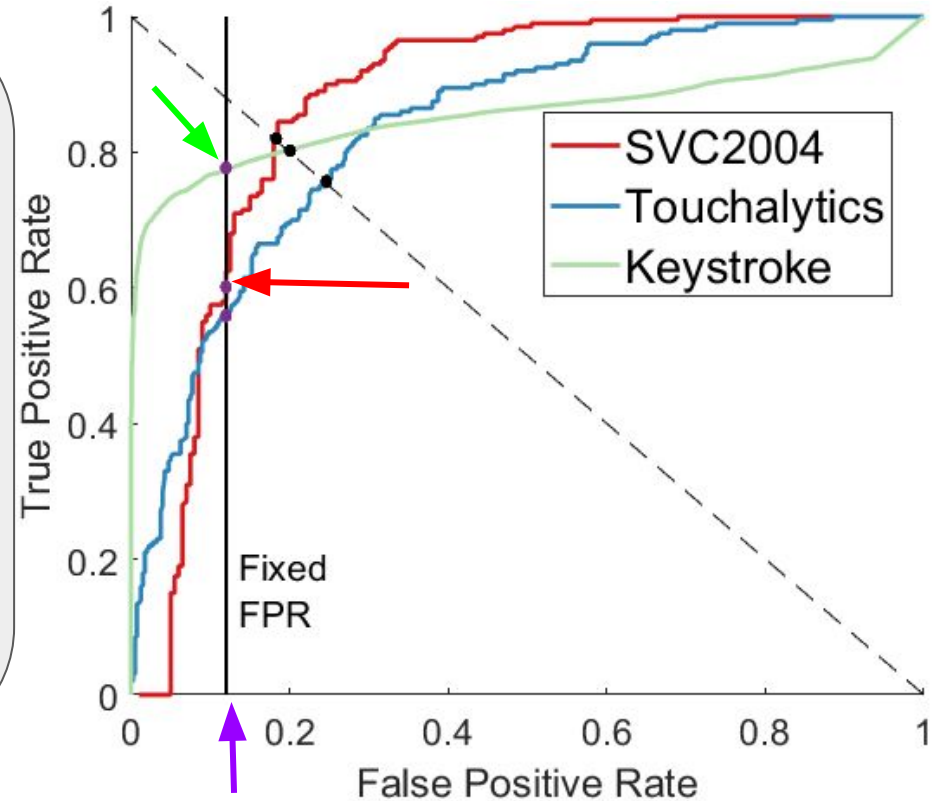Keystroke EER = 0.198

# EERs hides performance tradeoffs
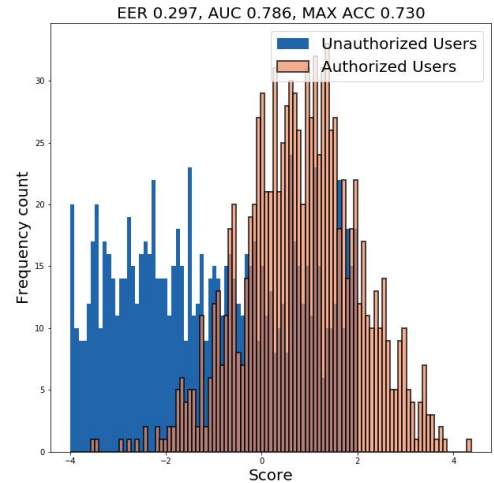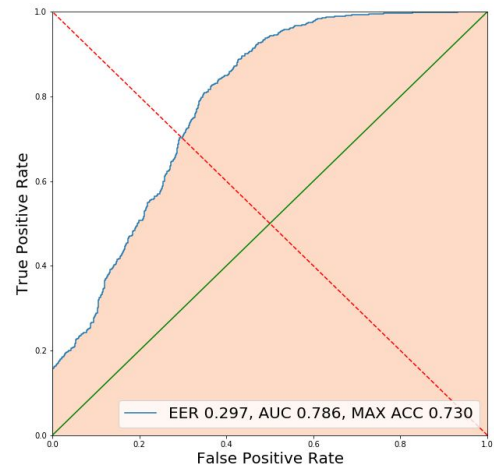
However,
Keystroke performs better
at FPR is 0.1:

SVC2004 TPR = 0.600
Keystroke TPR = 0.776

# Summary

- We propose reporting ROC and FCS to increase transparency
- No common reporting practice across surveyed systems
  - 36 out of 38 proposed systems had flaws in reporting
- Poor performance reporting impedes system comparison and replication
- Common metrics (e.g. accuracy, EER) can be misleading and hide performance tradeoffs

# Questions?

Please visit our websites for more details:

[lindqvistlab.org](lindqvistlab.org)
[scienceofsecurity.science](scienceofsecurity.science)