

# Metrics in Audio Security & Privacy

Andreas Nautsch  
EURECOM

EAB-RPC — RESPECT project  
2020-09-15 — Virtual Conference

# Outline

- Audio Security

- Automatic Speaker Verification Anti-Spoofing Challenge 2019 (ASVspoof 2019)  
<https://www.asvspoof.org/>
- Kinnunen et al.: “Tandem Assessment of Spoofing Countermeasures and Automatic Speaker Verification: Fundamentals,” IEEE/ACM TASLP 2020, DOI: 10.1109/TASLP.2020.3009494

- Audio Privacy

- VoicePrivacy 2020 Challenge  
<https://www.voiceprivacychallenge.org/>
- Nautsch et al.: “The Privacy ZEBRA: Zero Evidence Biometric Recognition Assessment,” Proc. Interspeech 2020, pre-print arxiv:2005.09413

# Audio Security Metric

## tandem Decision Cost Function (t-DCF)

Speech  $\Leftrightarrow$  SC37 dictionary :)

tar ~ mated, bona fide, ...

non ~ non-mated, non-attack, ....

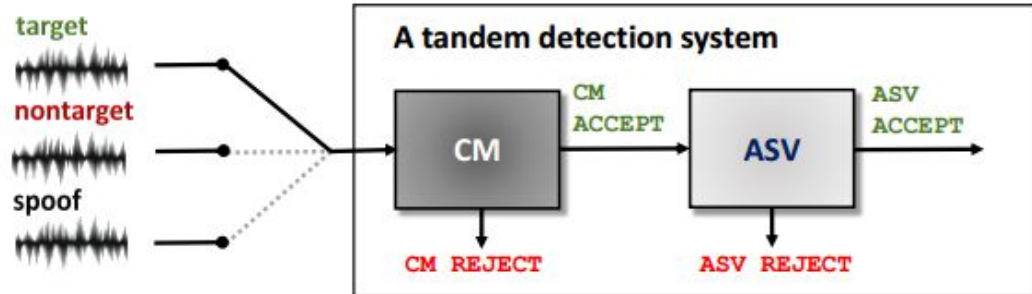
spoof ~ presentation attack,  
logical access spoof, ...

miss ~ FRR, FNMR, BPCER, ...

false alarm (fa) ~ FAR, FMR, APCER, ...

# Audio Security: The Setting

- Anti-Spoofing
  - “Physical Access”      Replay attacks      (see Voice PAD)
  - “Logical Access”      Voice synthesis/morphing/conversion attacks      (not PAD)
- In-Scope
  - Tandem operation of countermeasure (CM) and ASV sub-systems
  - Throughout formalised assessment
- Out-Scope
  - Informal descriptors by error rates
  - Purely CM-focused performance



# Audio Security: Expected Cost as Metric

- Quantification of beliefs
  - What is the impact of a decision outcome?
  - How likely is a decision outcome?
- Expected class discrimination risk
  - $\mathbb{E}$  [ risk | costs, class priors, classification rates ]
  - Sweep thresholds, take minimum

	Actual class	Tandem decision	Unit cost	Actual class	Asserted prior
a.	Target	REJECT (by ASV)	$C_{miss}$	Target	$\pi_{tar}$
b.	Nontarget	ACCEPT	$C_{fa}$	Nontarget	$\pi_{non}$
c.	Spoof	ACCEPT	$C_{fa,spoo}$	Spoof	$\pi_{spoo}$
d.	Target	REJECT (by CM)	$C_{miss}$		$\Sigma = 1$

$$t\text{-DCF} = C_{miss} \cdot \pi_{tar} \cdot P_a + C_{fa} \cdot \pi_{non} \cdot P_b + C_{fa,spoo} \cdot \pi_{spoo} \cdot P_c + C_{miss} \cdot \pi_{tar} \cdot P_d$$



# Audio Security: Tandem Classification Rates

$$C_{fa,spoof} \cdot \pi_{spoof} \cdot P_c$$

(CM ACCEPT, ASV ACCEPT)

Spoof

$$P_c(\tau_{cm}, \tau_{asv}) = P_{fa}^{cm}(\tau_{cm}) \times P_{fa,spoof}^{asv}(\tau_{asv})$$

$$C_{fa} \cdot \pi_{non} \cdot P_b$$

(CM ACCEPT, ASV ACCEPT)

Nontarget

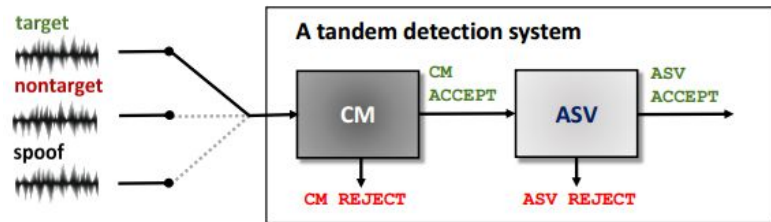
$$P_b(\tau_{cm}, \tau_{asv}) = (1 - P_{miss}^{cm}(\tau_{cm})) \times P_{fa}^{asv}(\tau_{asv})$$

$$C_{miss} \cdot \pi_{tar} \cdot P_d$$

(CM REJECT)

Target

$$P_d(\tau_{cm}, \tau_{asv}) = P_{miss}^{cm}(\tau_{cm})$$



$$C_{miss} \cdot \pi_{tar} \cdot P_a$$

(CM ACCEPT, ASV REJECT)

Target

$$P_a(\tau_{cm}, \tau_{asv}) = (1 - P_{miss}^{cm}(\tau_{cm})) \times P_{miss}^{asv}(\tau_{asv})$$

# Audio Security: Metric Normalisation

- Better comparability (other costs/priors)
- What are the extreme actions?

- CM & ASV: all-pass

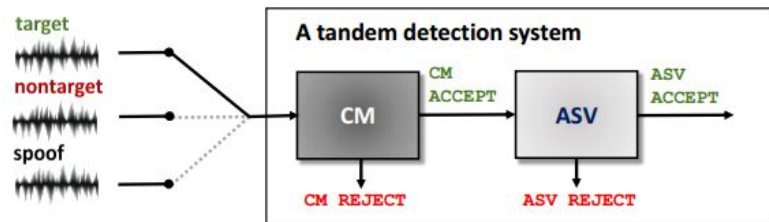
$$C_{fa} \cdot \pi_{non} \cdot \boxed{1} + C_{fa,spoo} \cdot \pi_{spoo} \cdot \boxed{1}$$

- CM: no-pass

$$C_{miss} \cdot \pi_{tar} \cdot \boxed{1}$$

- CM: all-pass & ASV: no-pass

$$C_{miss} \cdot \pi_{tar} \cdot \boxed{1}$$



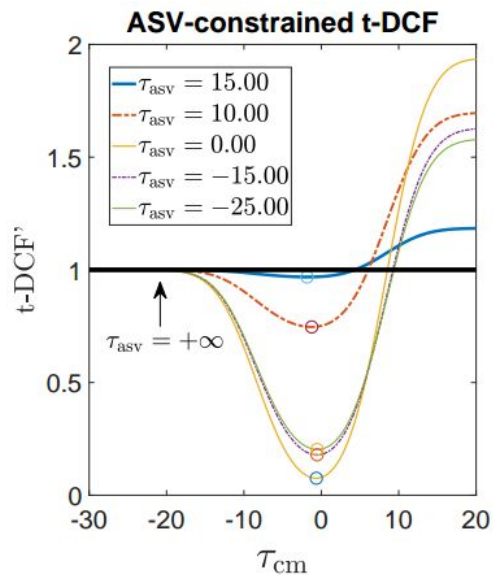
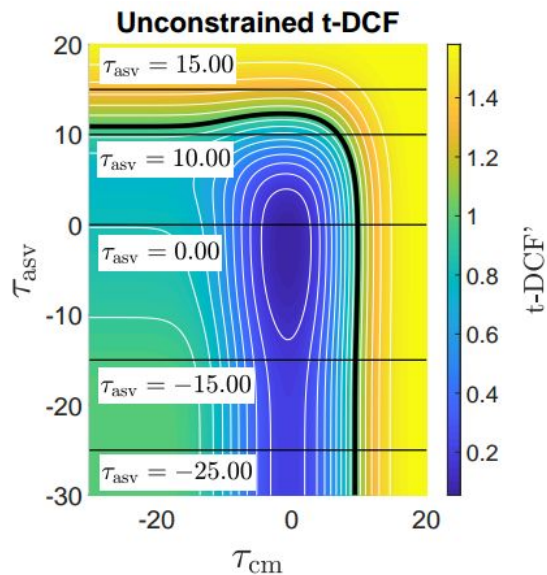
$$t\text{-DCF}'(\tau_{cm}, \tau_{asv}) = \frac{t\text{-DCF}(\tau_{cm}, \tau_{asv})}{t\text{-DCF}_{\text{default}}}$$

$$t\text{-DCF}'_{\min} = \frac{t\text{-DCF}_{\min}}{t\text{-DCF}_{\text{default}}} \leq \frac{t\text{-DCF}_{\min}}{t\text{-DCF}_{\min}} = 1$$

$$t\text{-DCF}_{\text{default}} = \min \{ C_{fa} \cdot \pi_{non} + C_{fa,spoo} \cdot \pi_{spoo}, C_{miss} \cdot \pi_{tar} \}$$

# Audio Security: t-DCF Examples

- ASVspoof 2019 Challenge
  - Cost & prior parameters as per challenge
  - Synthetic ASV & CM scores



## ASV-constrained t-DCF

$$t\text{-DCF}(\tau_{cm}) = C_0 + C_1 P_{miss}^{cm}(\tau_{cm}) + C_2 P_{fa}^{cm}(\tau_{cm})$$

$$C_0 = \pi_{tar} C_{miss} P_{miss}^{asv} + \pi_{non} C_{fa} P_{fa}^{asv}$$

$$C_1 = \pi_{tar} C_{miss} - (\pi_{tar} C_{miss} P_{miss}^{asv} + \pi_{non} C_{fa} P_{fa}^{asv})$$

$$C_2 = \pi_{spoo} C_{fa,spoo} P_{fa,spoo}^{asv}$$



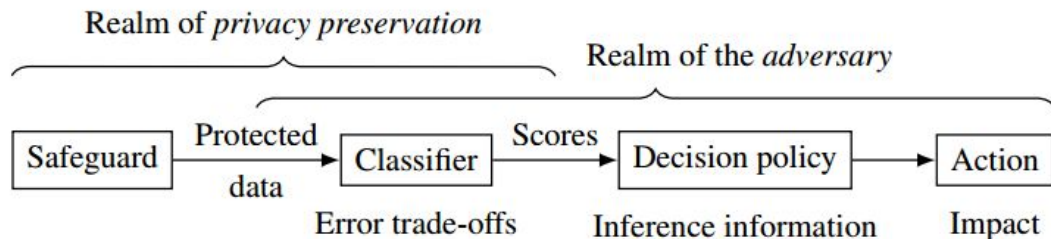
# Audio Privacy Metric

## Zero Evidence Biometric Recognition Assessment (The Privacy ZEBRA)



# Audio Privacy: The Setting

- Pseudomise audio speech data
- Decoupling layers & taking the perspective of an adversary

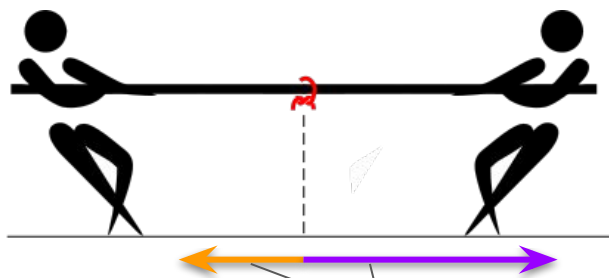
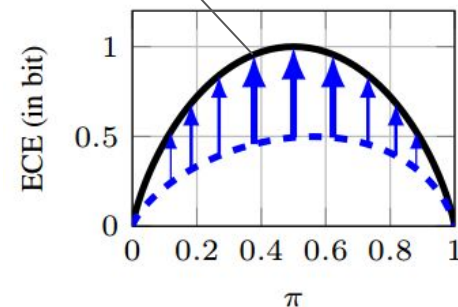


- Existing metrics do not suffice!
  - Zero-knowledge proofs are unavailable.
  - EER is the worst possible decision policy that an adversary can take for herself.
  - Unlinkability (not devised for this setting) — identity confirmation but not short-listing.
  - Any fixed error rate/cost metric prejudices privacy disclosure impacts to an individual.

# Audio Privacy: Zero Evidence as Metric

- Population level: Empirical Cross-Entropy (ECE)
  - Idea: prior entropy  $\Rightarrow$  evidence  $\Rightarrow$  posterior entropy
  - Cross-entropy of classification by scores from ground truth
  - Zero evidence: prior ECE = posterior ECE, regardless of prior  $\pi$
- Individual level: Zero Strength of Evidence
  - Forensic sciences: likelihood ratio
  - Who is stronger: prosecutor or defendant?

Coin-tossing simulation:  
all scores are equal  
“prior = posterior ECE”

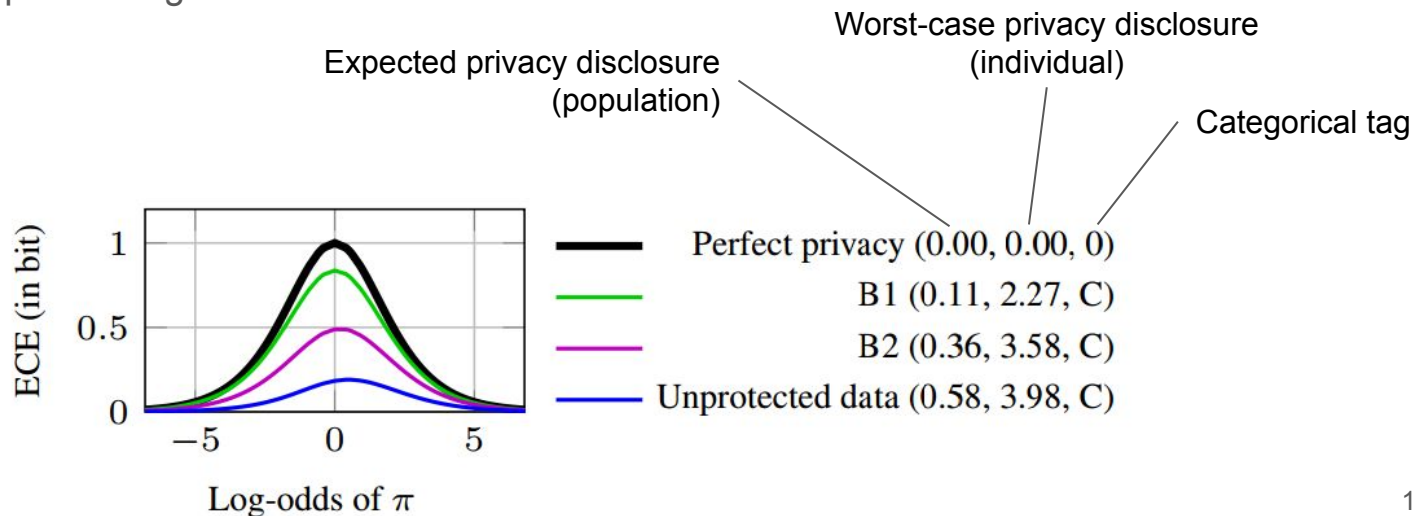


Tag	Category	Posterior odds ratio (flat prior)
0	$l = 1 = 10^0$	50 : 50 (flat posterior)
A	$10^0 < l < 10^1$	more disclosure than 50 : 50
B	$10^1 \leq l < 10^2$	one wrong in 10 to 100
C	$10^2 \leq l < 10^4$	one wrong in 100 to 10 000
D	$10^4 \leq l < 10^5$	one wrong in 10 000 to 100 000
E	$10^5 \leq l < 10^6$	one wrong in 100 000 to 1 000 000
F	$10^6 \leq l$	one wrong in at least 1 000 000

# Audio Privacy: ZEBRA Examples

- VoicePrivacy 2020 Challenge

- Task: speech recognition should work — voice biometrics not  $\Rightarrow$  modification of raw audio
- ASV: pre-trained kaldi x-vector recipe
- B1: DNN baseline
- B2: signal processing baseline



# Summary & Conclusion

- Summary
  - Audio security: cost-based approach for expected risk minimization
  - Audio privacy: relative information & strength of evidence approach
- Conclusion
  - Constrained cost as a guide for the CM optimization given a biometric system  
⇒ taking a holistic perspective
  - Audio privacy must achieve privacy for every single one; *are we a marginalising society?*  
⇒ expectation & worst-case estimates

Cheers.