# Revisiting and Improving Scoring Fusion for Spoofing-aware Speaker Verification Using Compositional Data Analysis
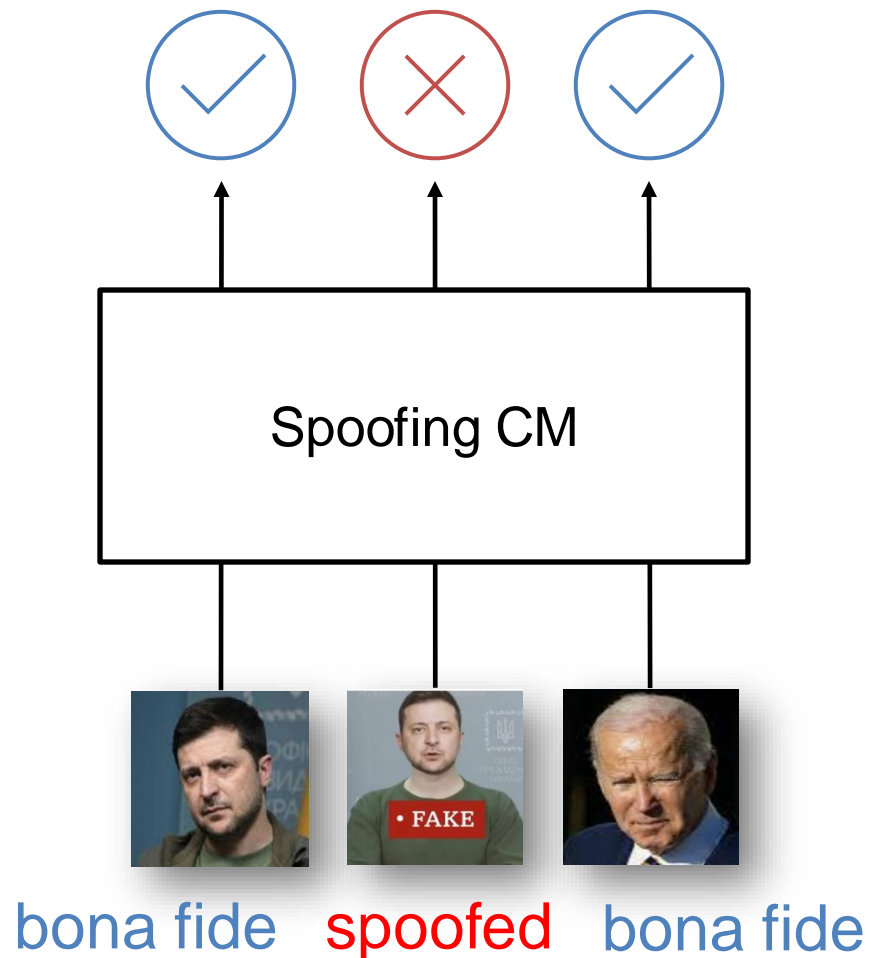
Xin Wang🎤, Tomi Kinnunen, Kong Aik Lee,
Paul-Gauthier Noe, Junichi Yamagishi

NII, JST PRESTO, UEF, PolyU, Inria

wangxin@nii.ac.jp

# Summary in one slide

❑ Question: how ASV and spoofing countermeasure (CM) should be fused **theoretically**?

❑ Message: fusing ASV and CM != fusing ASVs (or CMs)

❑ Methods

  – Linear fusion of log likelihood ratios (LLRs)   } Bayesian

  – Non-linear fusion of LLRs   decision theory

❑ Results: both better than baseline, non-linear the best

# Background: spoofing CM



protect human listeners

protect ASV

bona fide  spoofed  bona fide

# Background: spoofing CM protecting ASV



bona fide
matched

bona fide
not matched

# Background: spoofing-robust ASV (SASV)

# Background: spoofing-robust ASV (SASV)

❑ **Approach 1: end-to-end**



✓ easy to get hands on
x no extra explanation

A single deep neural network (DNN)

enroll

# Background: spoofing-robust ASV (SASV)

❏ **Approach 2: fusion-based**



x   technically demanding
✓ re-use CM & ASV
✓ extra explanations

Fusion

Spoofing CM

ASV

enroll

• FAKE

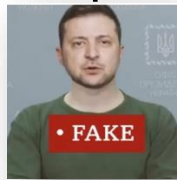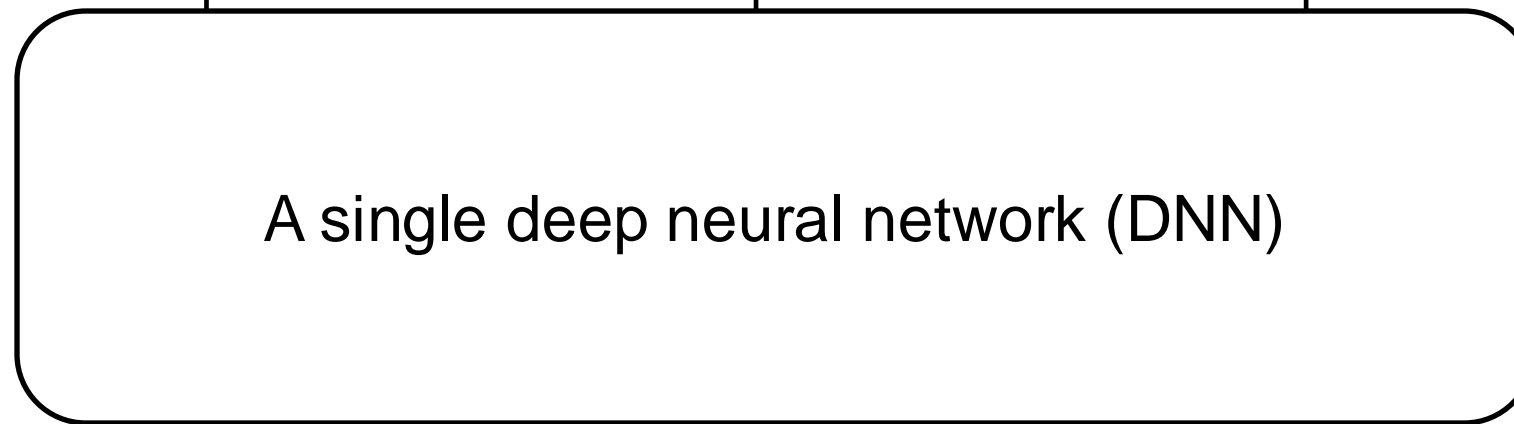# Question: how to *properly* fuse ASV and CM

# Question: how to *properly* fuse ASV and CM

❑ **baseline approach** (Jung 2022)



$$s_{\mathrm{sasv}} = s_{\mathrm{cm}} + s_{\mathrm{asv}}$$

$$s_{\mathrm{cm}} \in \mathbb{R} \qquad\qquad s_{\mathrm{asv}} \in \mathbb{R}$$

Spoofing CM      ASV

$$\boldsymbol{x}^{(r)}$$

$$\boldsymbol{x}^{(p)}$$

Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. 2022. SASV 2022: The first spoofing-aware speaker verification challenge. In *Proc. Interspeech*, 2022. 2893–2897.

# Question: how to *properly* fuse ASV and CM

☐ **baseline approach** (Jung 2022)

$$s_{\text{sasv}} = s_{\text{cm}} + s_{\text{asv}}$$

$$s_{\text{cm}} \in \mathbb{R} \qquad\qquad s_{\text{asv}} \in \mathbb{R}$$

| Spoofing CM | ASV |
|---|---|

$$\boldsymbol{x}^{(r)}$$

? What to do if, say, $s_{\text{cm}} \in [-100, 100]$ $s_{\text{asv}} \in [-1, 1]$

Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. 2022. SASV 2022: The first spoofing-aware speaker verification challenge. In *Proc. Interspeech*, 2022. 2893–2897.

# Question: how to *properly* fuse ASV and CM

☐ **baseline approach** (Jung 2022)

$$s_{\mathrm{sasv}} = s_{\mathrm{cm}} + s_{\mathrm{asv}}$$



? What to do if, say, $s_{\mathrm{cm}} \in [-100, 100]$  $s_{\mathrm{asv}} \in [-1, 1]$
? Why not normalize both, why summation …

*Any thoery to support the good pratice?*

Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. 2022. SASV 2022: The first spoofing-aware speaker verification challenge. In *Proc. Interspeech*, 2022. 2893–2897.

# Answers by this work

❑ **Fusion in SASV != fusion in ASV (or CM) ensemble (sec.2.1)**

▪ Spoofing CM and ASV are dealing with different pairs of hypotheses

▪ A different theory is needed

$$\log \frac{P(H_{\text{tar}}|\mathbf{X})}{1 - P(H_{\text{tar}}|\mathbf{X})} = \log \frac{\pi_{\text{tar}}}{1 - \pi_{\text{tar}}} + \boxed{\sum_{k=1}^{K} \text{llr}_{\text{non}}^{\text{tar}}(\boldsymbol{x}_k),}$$



$$s_{\text{sasv}} = s_{\text{cm}} + s_{\text{asv}}$$

$s_{\text{cm}} \in \mathbb{R}$

tanh

$s_{\text{asv}} \in \mathbb{R}$

Spoofing CM

ASV

$\boldsymbol{x}^{(p)}$

Logistic regression

$s_{\text{asv}} \in \mathbb{R}$

$s_{\text{asv}} \in \mathbb{R}$

ASV subsystem

ASV subsystem

$\boldsymbol{x}^{(p)}$

# Answers by this work

❑ **Fusion in SASV != fusion in ASV (or CM) ensemble (sec.2.1)**

  ▪ Spoofing CM and ASV are dealing with different pairs of hypotheses

  ▪ A different theory is needed

  We explain the practice in this talk

❑ **Linear summation (Sec.2.2 – 2.4)**

  ▪ Bayesian decision theory + compositional data analysis

  ▪ In practice: calibration + sum of CM and ASV LLRs

❑ **Non-linear fusion (Sec.2.5)**

  ▪ Bayesian decision theory (arxiv appendix)

    • the "optimal" solution to minimize a decision cost

  ▪ In practice: calibration & non-linear fusion

# Method 1: linear fusion in good practice

☐ **Score calibrations are needed**

$$s_{\mathrm{sasv}}$$

$$\mathsf{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x}) \qquad \mathsf{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

| Calibration | | Calibration |

$$s_{\mathrm{cm}} \in \mathbb{R} \qquad s_{\mathrm{asv}} \in \mathbb{R}$$

| Spoofing CM | | ASV | $\boldsymbol{x}^{(r)}$ |

$$\boldsymbol{x}^{(p)}$$

# Method 1: linear fusion in good practice

☐ **Score calibrations are needed**
☐ **LLRs should be summed**

? ~~Why normalize $s_{\mathrm{cm}}$, not $s_{\mathrm{asv}}$~~

? summation, product

$$s_{\mathrm{sasv}}$$

$$\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x}) \qquad \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

| Calibration | | Calibration |

$$s_{\mathrm{cm}} \in \mathbb{R} \qquad\qquad s_{\mathrm{asv}} \in \mathbb{R}$$

| Spoofing CM | | ASV |

$$\boldsymbol{x}^{(r)}$$

$$\boldsymbol{x}^{(p)}$$

# Method 1: linear fusion in good practice

☐ **Score calibrations are needed**
☐ **LLRs should be summed**

? ~~Why normalize~~ $s_{\mathrm{cm}}$ ~~, not~~ $s_{\mathrm{asv}}$
? summation, product

$s_{\mathrm{sasv}}$

$\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$

$\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$

Calibration

Calibration

*Three data classes but binary decisions! (sec 2.2 and appendix)*

$$s_{\mathrm{sasv}} = \mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x}) + \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

$$s_{\mathrm{cm}} = \mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x}) - \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

**NII**

# Method 1: linear fusion in good practice

☐ **Score calibration – nothing new**



$$\text{llr}_{\text{spf}}^{\text{tar.bf}}(\boldsymbol{x})$$

Calibration

$$s_{\text{cm}} \in \mathbb{R}$$

Spoofing CM

Logistic regression (Morrison 2013)

$$\text{llr}_{\text{non.bf}}^{\text{tar.bf}}(\boldsymbol{x}) = \underline{a} s_{\text{cm}} + \underline{b}$$

estimate {*a,b*} on using hold-out data

Geoffrey Stewart Morrison. 2013. Tutorial on logistic-regression calibration and fusion: converting a score to a likelihood ratio. Australian Journal of Forensic Sciences 45, 2 (2013), 173–197.
Scikit-learn: https://scikit-learn.org/stable/modules/calibration.html

# Method 1: linear fusion in good practice

❑ **Score calibration – nothing new**

$$\text{llr}_{\text{spf}}^{\text{tar.bf}}(\boldsymbol{x})$$

| Calibration |
|---|

$$s_{\text{cm}} \in \mathbb{R}$$

| Spoofing CM |
|---|

Logistic regression

Generative calibration (Brummer 2014)

$$p(s_{\text{cm}}|\text{spf}) \qquad p(s_{\text{cm}}|\text{tar.bf})$$



1. choose a parametric distribution
2. estimate distribution para. on dev. set
3. compute $\text{llr}_{\text{spf}}^{\text{tar.bf}}(\boldsymbol{x}) = \log \frac{p(s_{\text{cm}}|\text{tar.bf})}{p(s_{\text{cm}}|\text{spf})}$

Niko Brummer, Albert Swart, and David Van Leeuwen. 2014. A comparison of linear and non-linear calibrations for speaker recognition. In *Proc. Odyssey*, 2014. 14–18.
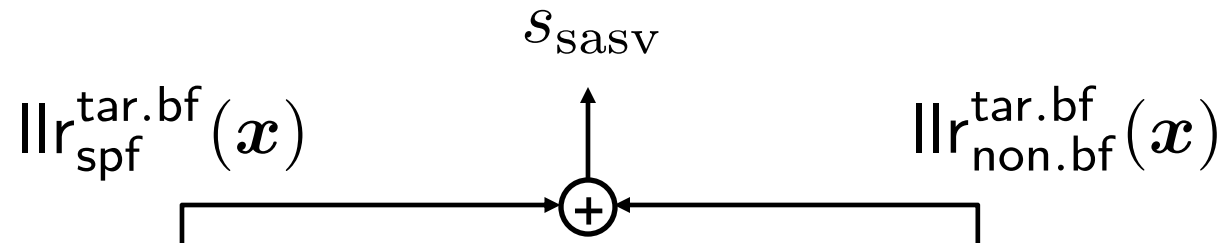
# Method 1: linear fusion in good practice

❑ **Score calibration – nothing new**



$$\text{llr}_{\text{spf}}^{\text{tar.bf}}(\boldsymbol{x})$$

Calibration

$$s_{\text{cm}} \in \mathbb{R}$$

Spoofing CM

Logistic regression

Generative calibration

Many other methods exist
(Ferrer 2022, Leeuwen 2013)

Luciana Ferrer, "Analysis and Comparison of Classification Metrics", arXiv:2209.05355, https://github.com/luferrer/CalibrationTutorial
David A. van Leeuwen and Niko Brümmer. 2013. The distribution of calibrated likelihood-ratios in speaker recognition. In *Proc. Interspeech*, 2013. 1619–1623.

# Method 1: linear fusion in good practice

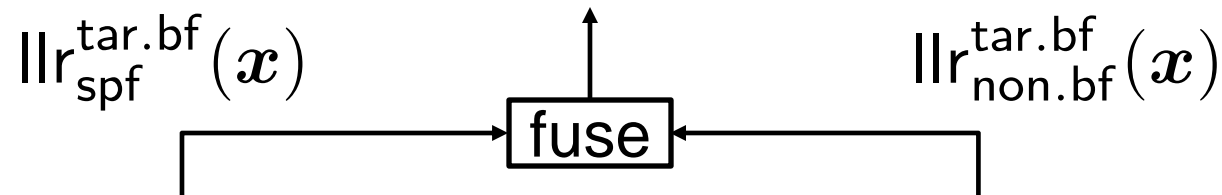☐ **Is linear fusion optimal for decision making?**

- No

$$s_{\mathrm{sasv}}$$

$$\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

$$\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

$+$

*See more in Sec2.5 & Appendix*

| Cost | ✓ | ✕ |
|---|---|---|
| Bona fide matched | 0 | Cmiss |
| Bona fide unmatched | Cfa | 0 |
| Spoofed | Cfa | 0 |

# Method 2: non-linear fusion is better

❑ **Non-linear fusion minimizes the cost**

$$s_{\mathrm{sasv}} = -\log \left[ (1-\rho)e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}}} + \rho e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}} \right] \qquad \text{for Cfa=Cmiss}$$

$$\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}(\boldsymbol{x}) \qquad\qquad \mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}}(\boldsymbol{x})$$

fuse

*See more in Sec2.5 & Appendix*

| Cost | ✓ | ✗ |
|---|---|---|
| Bona fide matched | 0 | Cmiss |
| Bona fide unmatched | Cfa | 0 |
| Spoofed | Cfa | 0 |

# Method 2: non-linear fusion is better

☐ **Non-linear fusion minimizes the cost**

$$s_{\text{sasv}} = -\log\left[(1-\rho)e^{-\text{llr}_{\text{non.bf}}^{\text{tar.bf}}} + \rho e^{-\text{llr}_{\text{spf}}^{\text{tar.bf}}}\right]$$

for Cfa=Cmiss

$\text{llr}_{\text{spf}}^{\text{tar.bf}}(\boldsymbol{x})$       $\text{llr}_{\text{non.bf}}^{\text{tar.bf}}(\boldsymbol{x})$

Asserted spoofing prior [Kinnuen 2023]

| | fuse | |

| Calibration | | Calibration |

$s_{\text{cm}} \in \mathbb{R}$       $s_{\text{asv}} \in \mathbb{R}$

| Spoofing CM | | ASV | $\boldsymbol{x}^{(r)}$ |

$\boldsymbol{x}^{(p)}$

# Method 2: non-linear fusion is better

☐ **Non-linear fusion minimizes the cost**

$$s_{\mathrm{sasv}} = -\log\left[(1-\rho)e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}}} + \rho e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}}\right]$$

for Cfa=Cmiss

$\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}(\boldsymbol{x})$

$\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}}(\boldsymbol{x})$

fuse

Asserted spoofing prior (Kinnuen 2023)

| Calibration | | Calibration |
|---|---|---|

$s_{\mathrm{cm}} \in \mathbb{R}$

$s_{\mathrm{asv}} \in \mathbb{R}$

A general form of ASV ($\rho = 0$) or CM ($\rho = 1$)

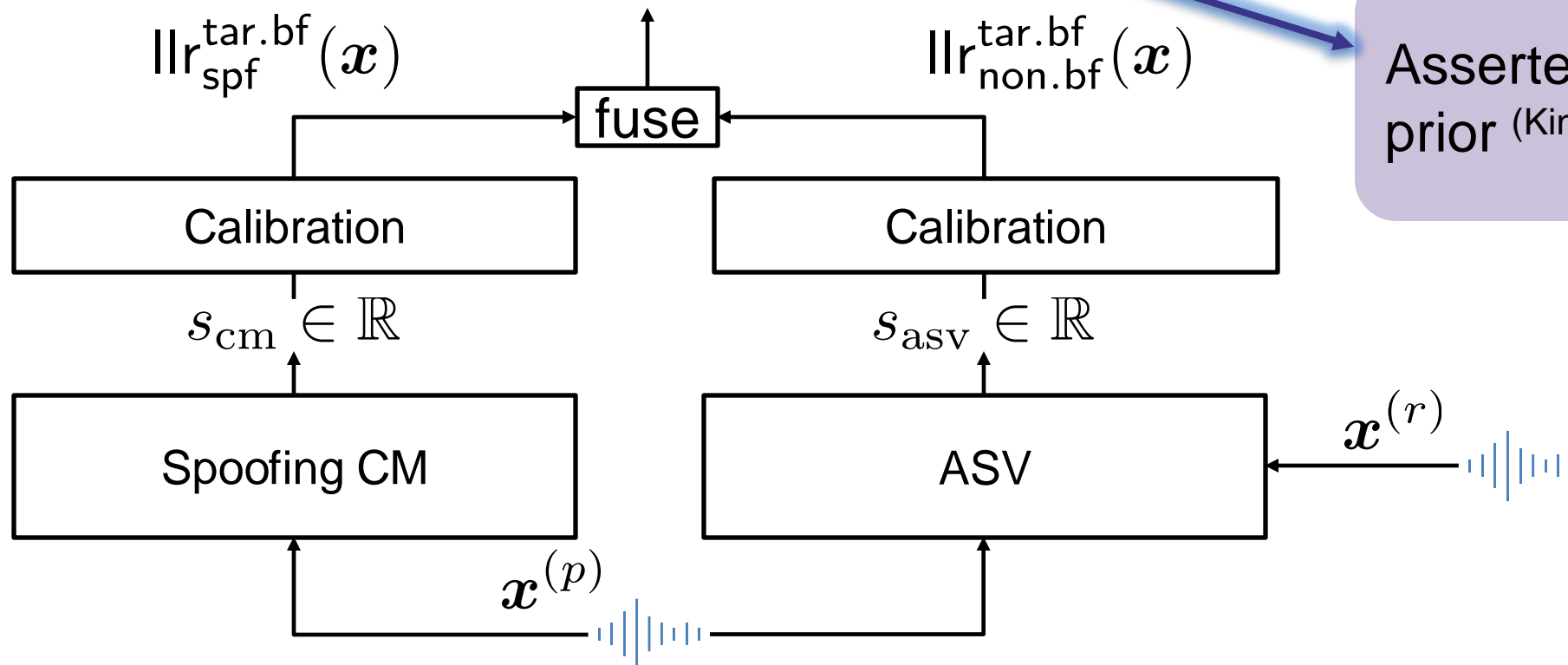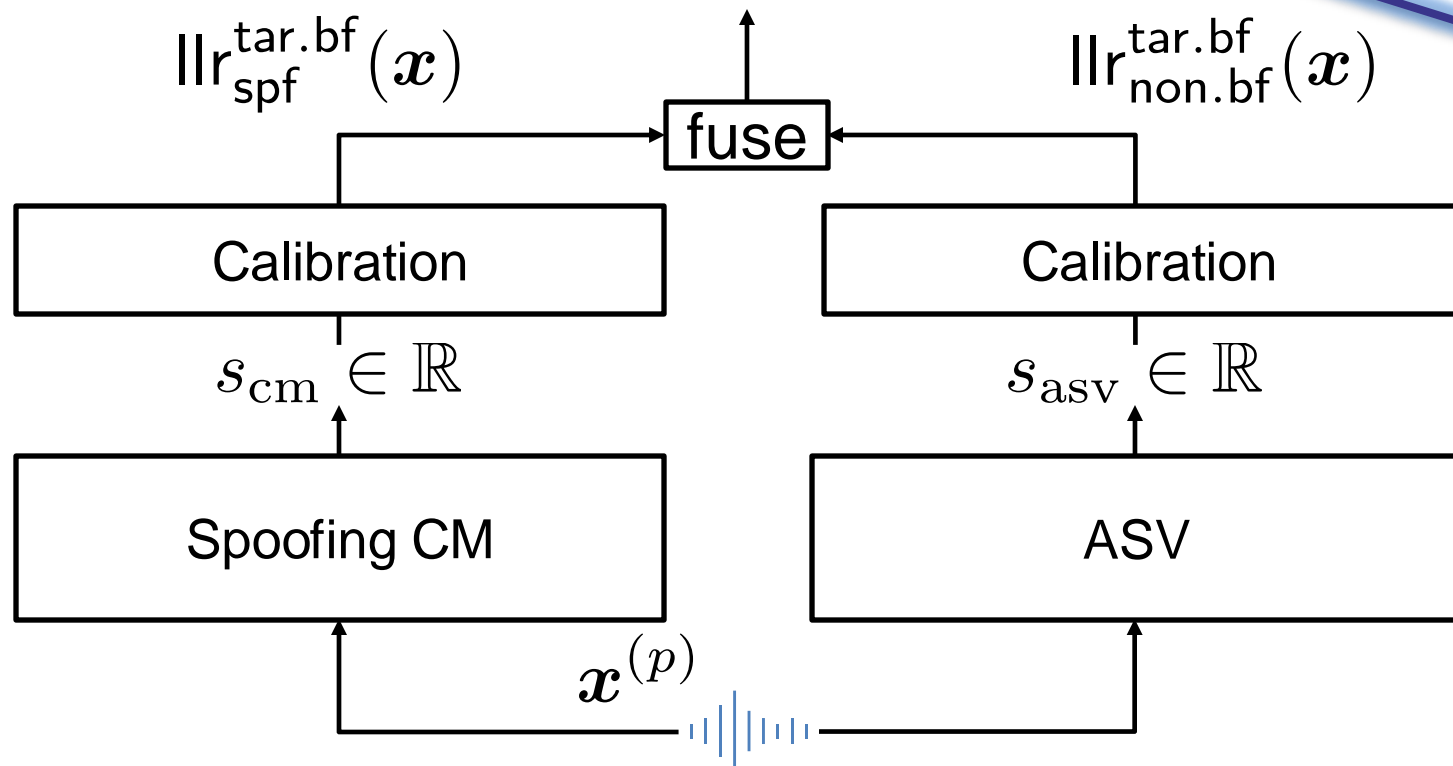| Spoofing CM | | ASV |
|---|---|---|

$\boldsymbol{x}^{(p)}$

# Method 2: non-linear fusion is better

☐ **Non-linear fusion minimizes the cost**

$$s_{\mathrm{sasv}} = -\log\left[(1-\rho)e^{-\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}} + \rho e^{-\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}}\right]$$

for Cfa=Cmiss

$\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$

$\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$

fuse

Calibration

Calibration

$s_{\mathrm{cm}} \in \mathbb{R}$

$s_{\mathrm{asv}} \in \mathbb{R}$

Spoofing CM

ASV

$\boldsymbol{x}^{(p)}$

Asserted spoofing prior (Kinnuen 2023)

A general form of ASV ($\rho = 0$) or CM ($\rho = 1$)

A general form of Gaussian fusion (Todisco 2018)

Massimiliano Todisco, Héctor Delgado, Kong Aik Lee, Md Sahidullah, Nicholas Evans, Tomi Kinnunen, and Junichi Yamagishi. 2018. presentation attack detection and automatic speaker verification: Common features and gaussian back-end fusion. In Proc. Interspeech, 2018. 77–81.

# Demo on toy data set



$$s_{\text{sasv}} = \text{llr}_{\text{non.bf}}^{\text{tar.bf}} + \text{llr}_{\text{spf}}^{\text{tar.bf}}$$

# Demo on toy data set



$$s_{\mathrm{sasv}} = -\log\left[(1-\rho)e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}}} + \rho e^{-\mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}}\right]$$

$$s_{\mathrm{sasv}} = \mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{non.bf}} + \mathrm{llr}^{\mathrm{tar.bf}}_{\mathrm{spf}}$$

# Demo on toy data set



$$s_{\mathrm{sasv}} = -\log\left[(1-\rho)e^{-\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}} + \rho e^{-\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}}\right]$$

different $\rho$

$$s_{\mathrm{sasv}} = \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}} + \mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}$$

# Recap the practices

Linear fusion

Non-linear fusion

$$s_{\mathrm{sasv}} = \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}} + \mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}$$

$$s_{\mathrm{sasv}} = -\log\left[(1-\rho)e^{-\mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}} + \rho e^{-\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}}\right]$$

$$\mathrm{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x}) \qquad \mathrm{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

| fuse |

| Calibration | | Calibration |

$$s_{\mathrm{cm}} \in \mathbb{R} \qquad\qquad s_{\mathrm{asv}} \in \mathbb{R}$$

| Spoofing CM | | ASV | $\boldsymbol{x}^{(r)}$ |

$$\boldsymbol{x}^{(p)}$$

# Experiments

❑ **Data**

- SASV 2022 challenge database, official protocols [Jung 2022]

❑ **Systems**

- All use *pre-trained* ASV and CM from SASV 2022 B1 [Jung 2022]
- Systems differ in score calibration & fusion

❑ **Misc**

- Training & evaluation in six rounds
- Averaged results are reported

Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. 2022. SASV 2022: The first spoofing-aware speaker verification challenge. In Proc. Interspeech, 2022. 2893–2897.

# Experiments

☹ worse ▉▒░ ☺ better

| ID | B1 | B1c | L2 | L2c | L3 | L3c | B1v2 | Post |
|---|---|---|---|---|---|---|---|---|
| Fusion | linear | | linear | | non-linear | | | |
| Calibration | × | ✓ | × | ✓ | × | ✓ | × | × |
| SASV-EER (%) | 20.46 | 2.73 | 3.31 | 1.56 | 1.44 | 1.43 | 1.60 | 1.55 |
| conf. ($\alpha = 5\%$) | ±0.40 | ±0.27 | ±0.31 | ±0.23 | ±0.23 | ±0.23 | ±0.22 | ±0.24 |
| Cllr | 2.17 | 1.09 | 1.04 | 0.14 | 0.18 | 0.16 | 0.96 | 0.84 |
| $\text{Cllr}_{\min}$ | 0.52 | 0.11 | 0.13 | 0.07 | 0.06 | 0.07 | 0.08 | 0.07 |
| $\text{Cllr}_{\text{calib}}$ | 1.64 | 0.98 | 0.91 | 0.07 | 0.11 | 0.10 | 0.88 | 0.78 |
| t-EER (%) | 2.10 | 2.10 | 1.68 | 1.68 | 1.68 | 1.68 | 2.19 | 2.21 |

SASV-EER (Jung2022)

other metrics

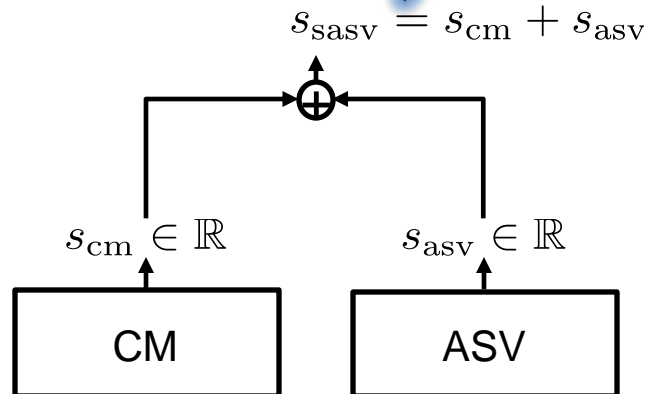Systems with different fusion & calibration methods

From other papers

**NII**

# Experiments

| ID | B1 | B1c | L2 | L2c | L3 | L3c | B1v2 | Post |
|---|---|---|---|---|---|---|---|---|
| Fusion | linear | | | linear | | | | |
| Calibration | × | ✓ | × | ✓ | × | ✓ | × | × |
| SASV-EER (%) | 20.46 | 2.73 | 3.31 | 1.56 | 1.44 | 1.43 | 1.60 | 1.55 |
| conf. $(\alpha = 5\%)$ | ±0.40 | ±0.27 | ±0.31 | ±0.23 | ±0.23 | ±0.23 | ±0.22 | ±0.24 |

no
calibration

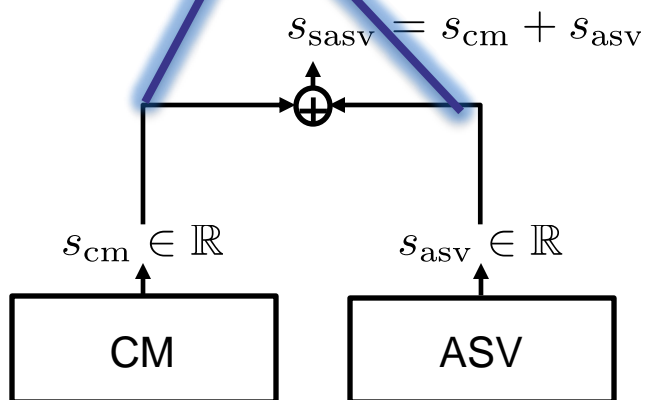log.reg.
calibration

log.reg. + Gaussian
calibration

$s_{\text{sasv}} = s_{\text{cm}} + s_{\text{asv}}$

$\oplus$

$s_{\text{cm}} \in \mathbb{R}$        $s_{\text{asv}} \in \mathbb{R}$

CM        ASV

$s_{\text{sasv}}$

$\oplus$

logistic reg.
calibration        logistic reg.
calibration

$s_{\text{cm}} \in \mathbb{R}$        $s_{\text{asv}} \in \mathbb{R}$

CM        ASV

$s_{\text{sasv}}$

$\oplus$

Gaussian +
logistic reg.        Gaussian +
logistic reg.

$s_{\text{cm}} \in \mathbb{R}$        $s_{\text{asv}} \in \mathbb{R}$

CM        ASV

# Experiments

bona fide matched
bona fide unmatched
spoofed



$$s_{\text{sasv}} = s_{\text{cm}} + s_{\text{asv}}$$

$s_{\text{cm}} \in \mathbb{R}$    $s_{\text{asv}} \in \mathbb{R}$

CM    ASV

$s_{\text{sasv}}$

logistic reg. calibration    logistic reg. calibration

$s_{\text{cm}} \in \mathbb{R}$    $s_{\text{asv}} \in \mathbb{R}$

CM    ASV

$s_{\text{sasv}}$

Gaussian + logistic reg.    Gaussian + logistic reg.

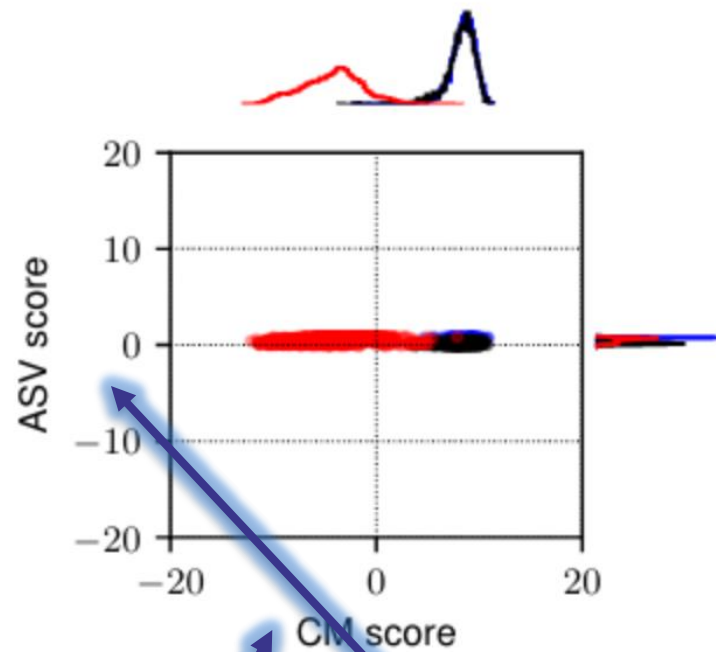$s_{\text{cm}} \in \mathbb{R}$    $s_{\text{asv}} \in \mathbb{R}$
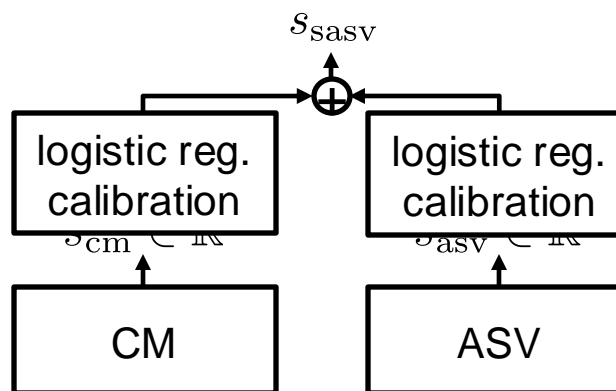
CM    ASV

baseline        good linear fusion        good linear fusion  32
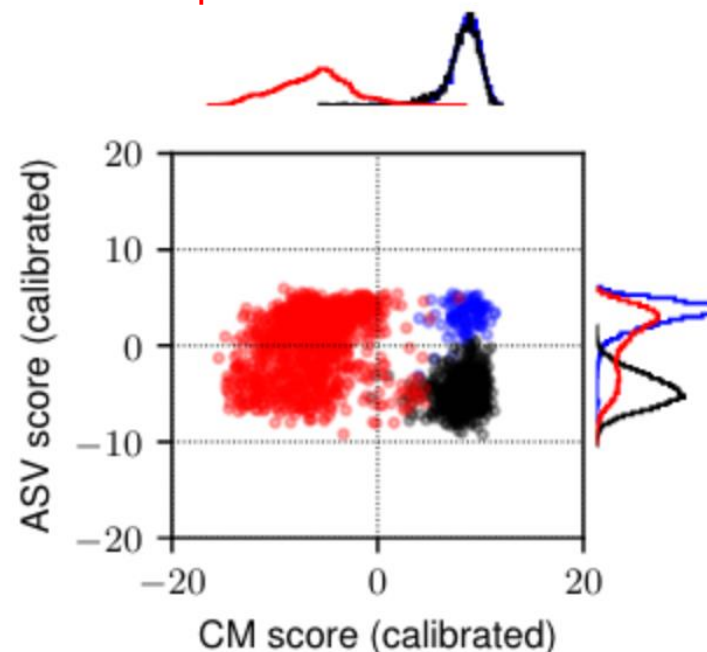
# Experiments

bona fide matched
bona fide unmatched
spoofed
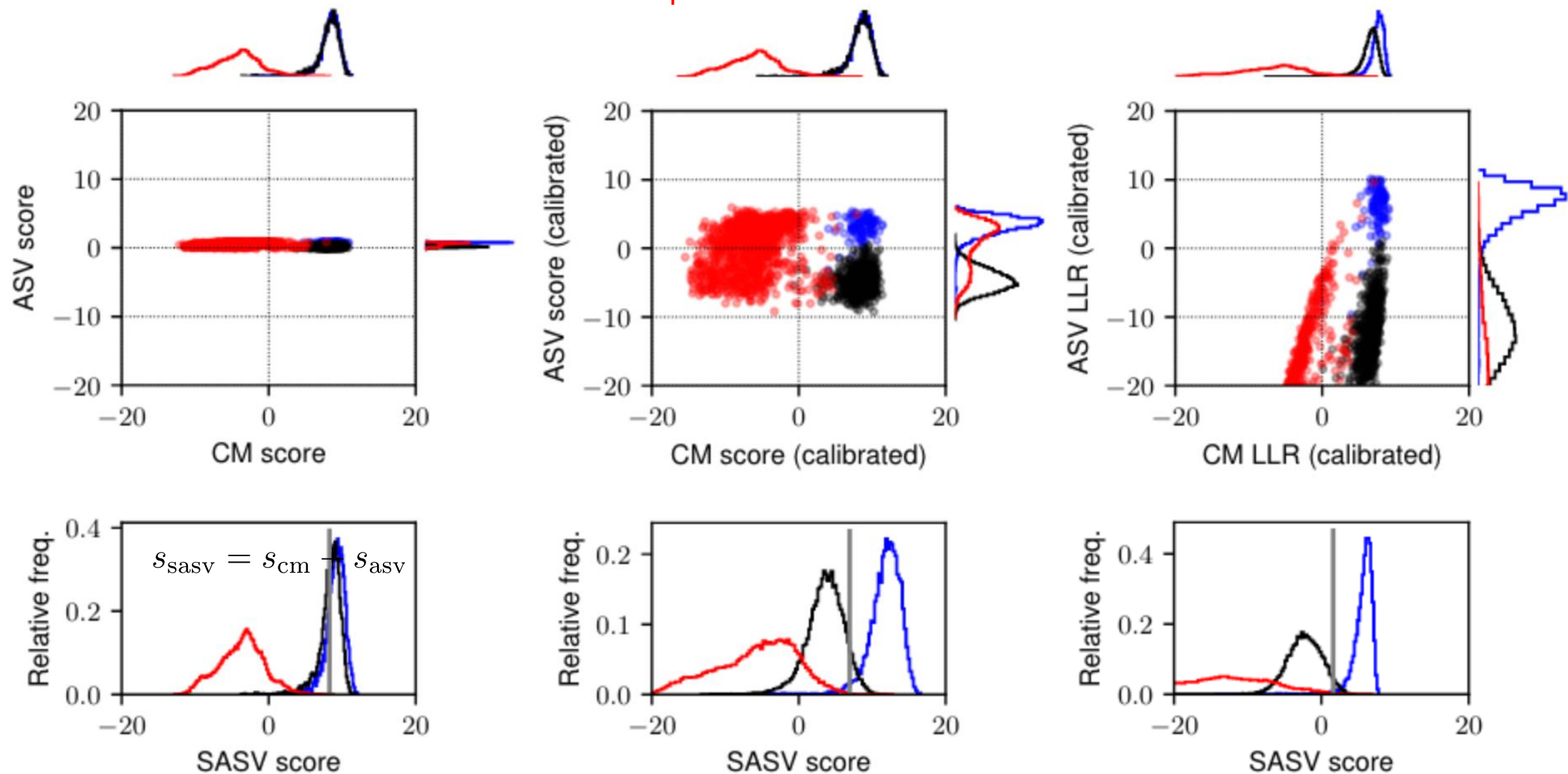
$s_{\mathrm{sasv}} = s_{\mathrm{cm}} + s_{\mathrm{asv}}$

NII

# Experiments

| ID | B1 | B1c | L2 | L2c | L3 | L3c | B1v2 | Post |
|---|---|---|---|---|---|---|---|---|
| Fusion | | | | linear | | non-liear | (Jung 2022) | (Zhang 2022) |
| Calibration | × | ✓ | × | ✓ | × | ✓ | × | × |
| SASV-EER (%) | 20.46 | 2.73 | 3.31 | 1.56 | 1.44 | 1.43 | 1.60 | 1.55 |
| conf. ($\alpha = 5\%$) | ±0.40 | ±0.27 | ±0.31 | ±0.23 | ±0.23 | ±0.23 | ±0.22 | ±0.24 |

The difference is small on this database

good linear fusion

good non-linear fusion

Jee-weon Jung, Hemlata Tak, Hye-jin Shim, Hee-Soo Heo, Bong-Jin Lee, Soo-Whan Chung, Ha-Jin Yu, Nicholas Evans, and Tomi Kinnunen. 2022. SASV 2022: The first spoofing-aware speaker verification challenge. In Proc. Interspeech, 2022. 2893–2897.
You Zhang, Ge Zhu, and Zhiyao Duan. 2022. A Probabilistic Fusion Framework for Spoofing Aware Speaker Verification. In *Proc. Odyssey*, June 28, 2022. ISCA, 77–84.

# Main messages

- ❑ **Fusion SASV != fusion of ASV or CM ensemble**
- ❑ **Linear and non-linear can be suppored by theory**
- ❑ **Calibration affects discrimination**

$$s_{\mathrm{sasv}}$$

$$\mathsf{llr}_{\mathrm{spf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

$$\mathsf{llr}_{\mathrm{non.bf}}^{\mathrm{tar.bf}}(\boldsymbol{x})$$

fuse

| Calibration | | Calibration |

$$s_{\mathrm{cm}} \in \mathbb{R}$$

$$s_{\mathrm{asv}} \in \mathbb{R}$$

| Spoofing CM | | ASV |

$$\boldsymbol{x}^{(r)}$$

$$\boldsymbol{x}^{(p)}$$

# Pointers

❑ **Evaluation using the same Bayes decision cost**

> Hye-jin Shim, Jee-weon Jung, Tomi Kinnunen, Nicholas Evans, Jean-Francois Bonastre, and Itshak Lapidot. 2024. **a-DCF: an architecture agnostic metric with application to spoofing-robust speaker verification**. In Proc. Odyssey, 2024. 158–164. https://doi.org/10.21437/odyssey.2024-23

❑ **SOTA ASV is not robust to spoofing attacks**

> Jee-weon Jung, Xin Wang, Nicholas Evans, Shinji Watanabe, Hye-jin Shim, Hemlata Tak, Sidhhant Arora, Junichi Yamagishi, and Joon Son Chung. 2024. **To what extent can ASV systems naturally defend against spoofing attacks?** In Proc. Interspeech, 2024. .

A4-05.5

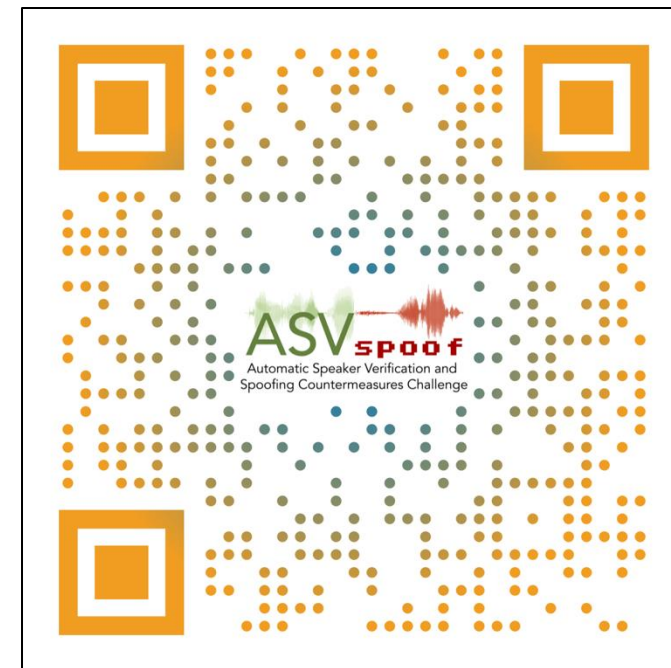❑ **The non-linear fusion has been used by many teams in ASVspoof 5 challenge**

# Thank you


Code & Jupyter notebook
*step-by-step explanation*


Appendix
*theory in details*


ASVspoof

# Fusing CM & ASV is special

❑ **A single ASV**

$$\log \frac{p(\boldsymbol{x}|H_{\mathrm{tar}})}{p(\boldsymbol{x}|H_{\mathrm{non}})} \gtrless -\log \frac{P(H_{\mathrm{tar}})}{P(H_{\mathrm{non}})}$$

match

not match

$$\frac{\textit{decision}}{\textit{scoring}}$$

$$\log \frac{p(\boldsymbol{x}|H_{\mathrm{tar}})}{p(\boldsymbol{x}|H_{\mathrm{non}})}$$

ASV ← $\boldsymbol{x}^{(r)}$

$\boldsymbol{x}^{(p)}$

$\boldsymbol{x} = (\boldsymbol{x}^{(p)}, \boldsymbol{x}^{(r)})$

# Fusing CM & ASV is special

❑ **Fusing ASV, face recognition, and other biometrics**

$$\sum_k \log \frac{p(\boldsymbol{x}_k|H_{\text{tar}})}{p(\boldsymbol{x}_k|H_{\text{non}})} \gtrless -\log \frac{P(H_{\text{tar}})}{P(H_{\text{non}})}$$

*decision*

─────────────────────────────────────

*scoring*

$$\log \frac{p(\boldsymbol{x}_1|H_{\text{tar}})}{p(\boldsymbol{x}_1|H_{\text{non}})}$$

$$\log \frac{p(\boldsymbol{x}_2|H_{\text{tar}})}{p(\boldsymbol{x}_2|H_{\text{non}})}$$

| Face recognition | | ASV |

# Fusing CM & ASV is special

❑ **CM and ASV are dealing with different hypotheses**

$$\sum_k \log \frac{p(\boldsymbol{x}_k|H_{\mathrm{tar}})}{p(\boldsymbol{x}_k|H_{\mathrm{non}})} \gtrless -\log \frac{P(H_{\mathrm{tar}})}{P(H_{\mathrm{non}})}$$

*decision*

*scoring*

$$\log \frac{p(\boldsymbol{x}|H_{\mathrm{bon}})}{p(\boldsymbol{x}|H_{\mathrm{spf}})}$$

$$\log \frac{p(\boldsymbol{x}|H_{\mathrm{tar}})}{p(\boldsymbol{x}|H_{\mathrm{non}})}$$
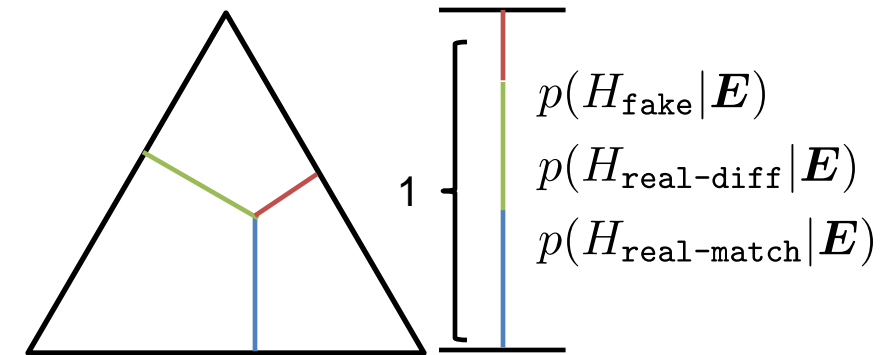
| CM | ASV |

# Fusing CM & ASV is special

❑ **We have three classes of data in two separate hypothesis testings**

$$\{H_{\text{fake}}, H_{\text{real-diff}}, H_{\text{real-match}}\}$$



Simplex ⟹

$$p(H_{\text{fake}}|\boldsymbol{E})$$
$$p(H_{\text{real-diff}}|\boldsymbol{E})$$
$$p(H_{\text{real-match}}|\boldsymbol{E})$$

What we need ⟵ $\tilde{p}_2$

Bayes' rule & Isometric-log-ratio

$\tilde{p}_1$ ⟵

**Optimal way using ternary hypothesis testing**

$$p(\boldsymbol{E}|H_{\text{fake}}) \quad p(\boldsymbol{E}|H_{\text{real-diff}}) \quad p(\boldsymbol{E}|H_{\text{real-match}})$$

# Fusing CM & ASV is special

❑ **We have three classes of data in two separate hypothesis testings**

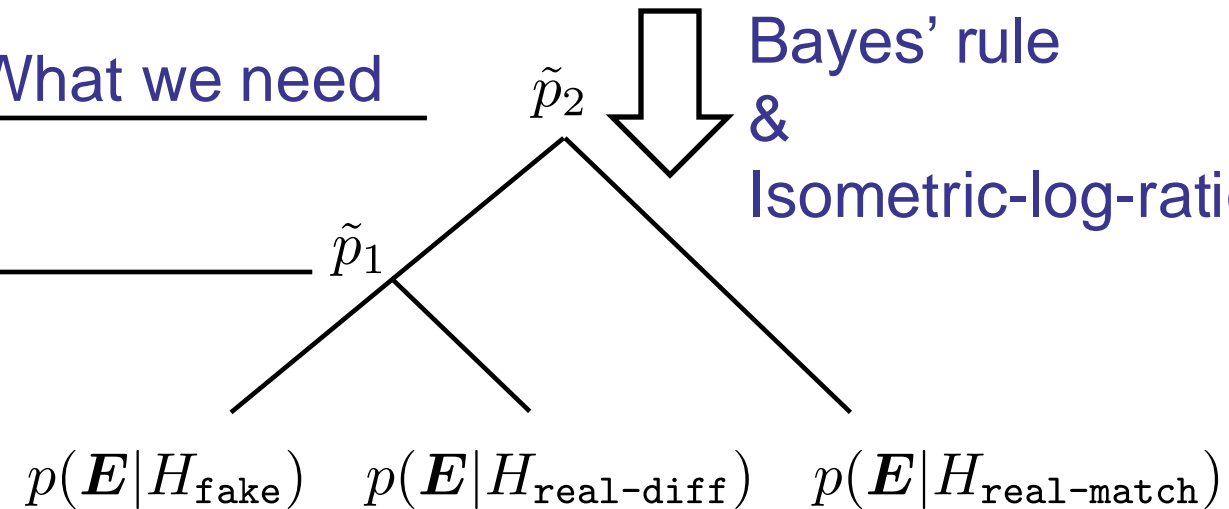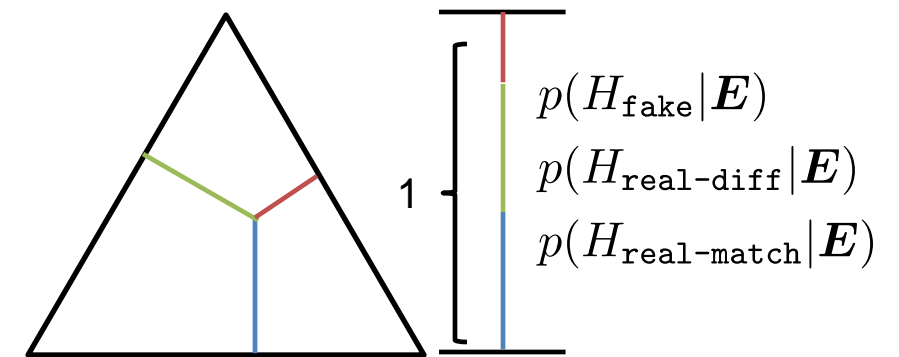$$\{H_{\text{fake}}, H_{\text{real-diff}}, H_{\text{real-match}}\}$$



Simplex

$$p(H_{\text{fake}}|\boldsymbol{E})$$
$$p(H_{\text{real-diff}}|\boldsymbol{E})$$
$$p(H_{\text{real-match}}|\boldsymbol{E})$$

$$\tilde{p}_2 = \frac{1}{\sqrt{6}}\left[\log\frac{p(\boldsymbol{E}|H_{\text{real-match}})}{p(\boldsymbol{E}|H_{\text{fake}})} + \log\frac{p(\boldsymbol{E}|H_{\text{real-match}})}{p(\boldsymbol{E}|H_{\text{real-diff}})}\right] \leftarrow \tilde{p}_2$$

Bayes' rule
&
Isometric-log-ratio

$$\tilde{p}_1$$

vs

log likelihood ratio

vs

log likelihood ratio

$$p(E|H_{\text{fake}}) \quad p(E|H_{\text{real-diff}}) \quad p(E|H_{\text{real-match}})$$