Symposium on calibration in forensic science
Aston Institute for Forensic Linguistics

# Calibration in Forensic
# Voice Comparison

**Daniel Ramos**

daniel.ramos@uam.es

Audias – Audio, Data Intelligence and Speech
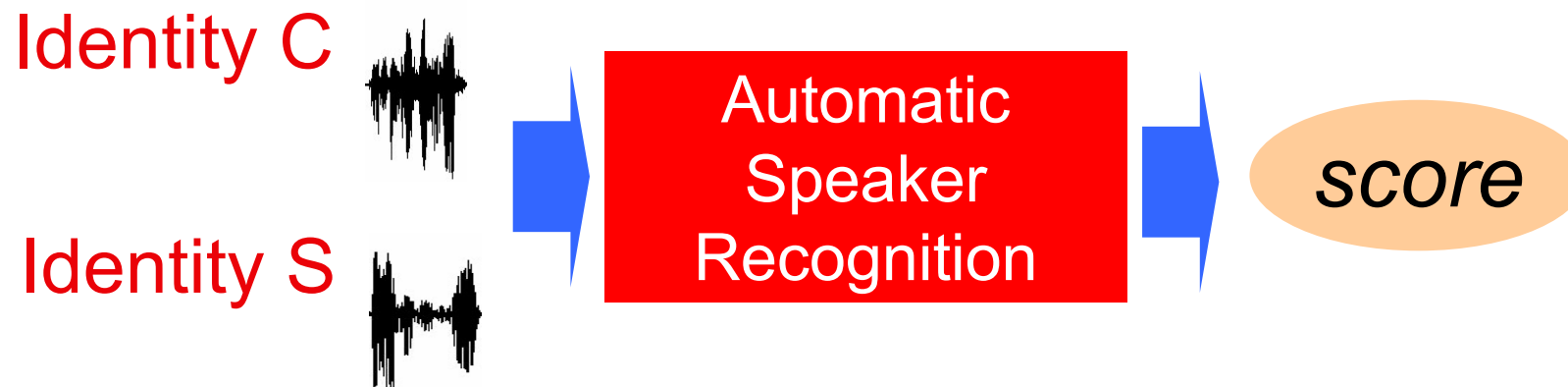Universidad Autónoma de Madrid

http://audias.ii.uam.es

# Automatic Speaker Recognition

- Standard architechture: *scores*

Identity C 

Identity S 

 Automatic Speaker Recognition → *score*

- Ideally:
  - If C y S are same identity (same-source), high**er** score
  - If C y S are different identities (different-source), low**er** score
- Thus, a score allows discrimination
- Not enough in forensics: a likelihood ratio (LR) is needed

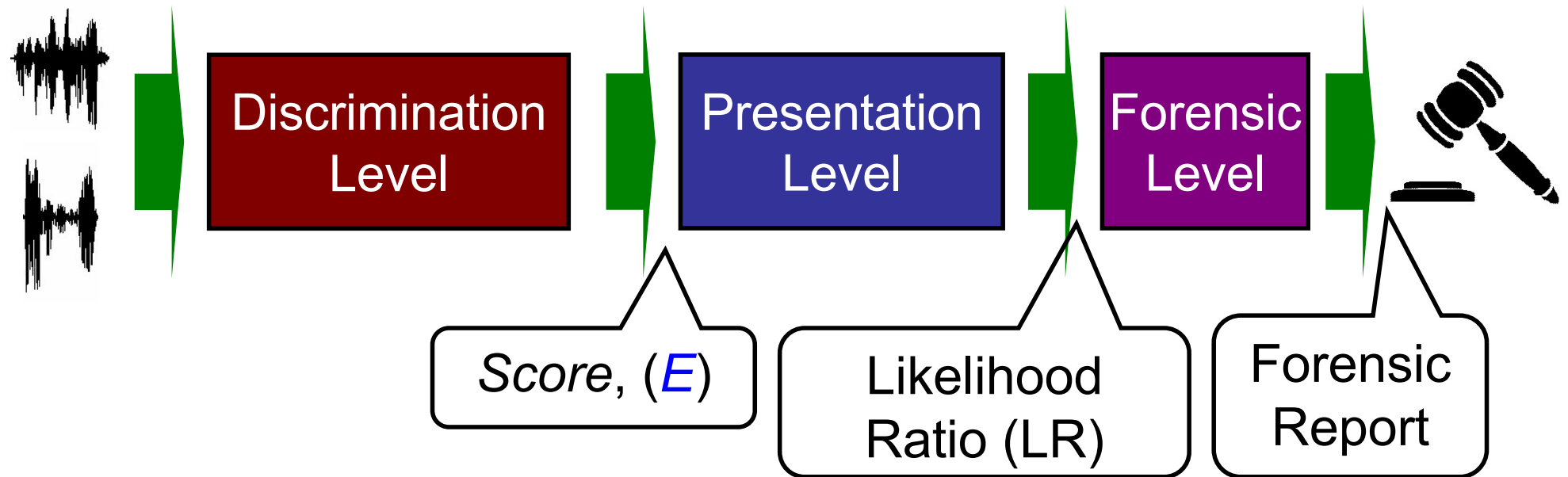# LR with Automatic Speaker Recognition



| Discrimination Level | Presentation Level | Forensic Level |

*Score*, (*E*)

Likelihood Ratio (LR)

Forensic Report

‹aud*i*as›

# LR with Automatic Speaker Recognition

| Discrimination Level | Presentation Level | Forensic Level |
|---|---|---|

- **Objective:** discriminating scores
  - Score-based architecture
  - Improve discrimination
  - So-called automatic speaker recognition

‹audias›

UAM

# LR with Automatic Speaker Recognition



- ■ Objective: transform scores into likelihood ratios

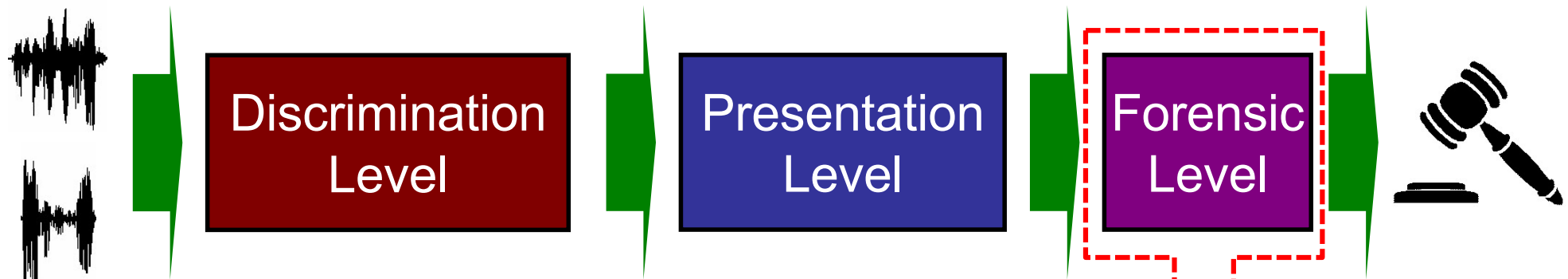  - ❑ Warning: LRs *must be well calibrated*

    **Forensic Science International**

    Reliable support: Measuring calibration of likelihood ratios[☆]

    Daniel Ramos[*], Joaquin Gonzalez-Rodriguez

  - ❑ $C_{llr}$: a popular measure of performance (the lower the better)

⟨aud*i*as⟩

# LR with Automatic Speaker Recognition



**Discrimination Level** → **Presentation Level** → **Forensic Level**

- Objective: adequate forensic reports
  - Probabilistic weight of the evidence
  - Following recommendations (ENFSI)
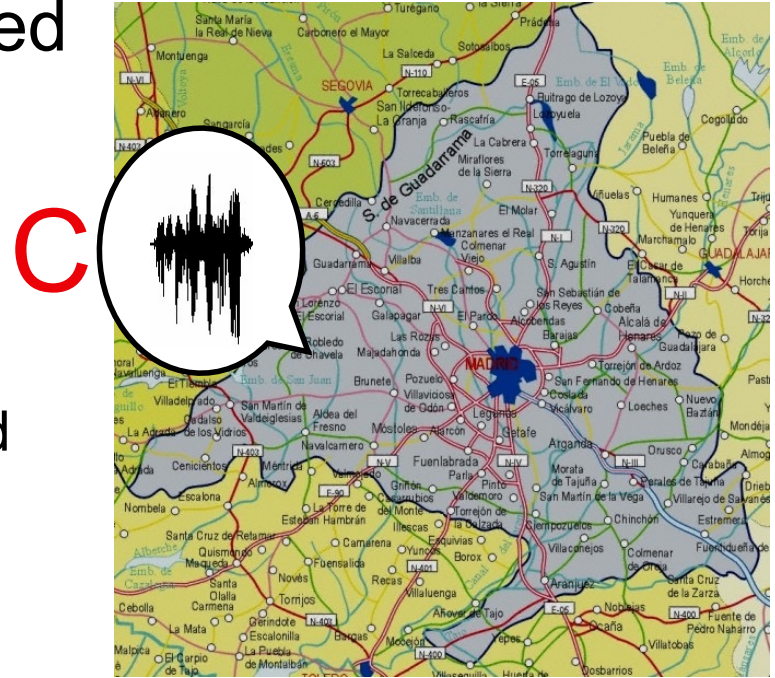  - Validation
  - Accreditation

Methodological Guidelines for Best Practice in Forensic Semiautomatic and Automatic Speaker Recognition

ENFSI GUIDELINE FOR EVALUATIVE REPORTING IN FORENSIC SCIENCE

ENFSI — EUROPEAN NETWORK OF FORENSIC SCIENCE INSTITUTES

‹aud*i*as›

# A Very Simplified
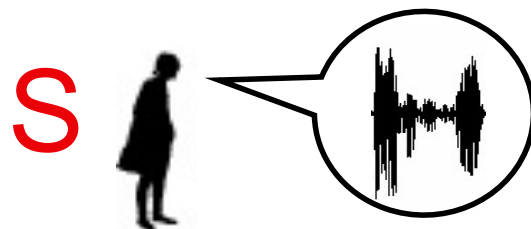# (Yet Illustrative)
# Example

# Simulated Case

- **Incriminating recordings wire-tapped in the Madrid region (trace)**

  - Population: potential sources of the speech

    - Speakers from Madrid region, with similar characteristics with questioned speech

      - Language
      - Accent
      - ...

  - Digital Wire-Tapping (SITEL, Spanish nationwide system)

- **Police investigations lead to a suspect**

C

S

‹ audias ›

# Simulated Case

■ **Recordings are taken from the suspect (reference speech)**



❑ Typically, controlled recordings

   ■ But very different conditions as for the questioned speech

❑ Could be previous wire-tappings where authorship is accepted

   ■ Similar conditions as questioned speech

Example: Ahumada III Database (Real Cases)

S       C 

(reference)            (trace)

D. Ramos, J. Gonzalez-Rodriguez, J. Gonzalez-Dominguez and J. J. Lucena-Molina, "Addressing database mismatch in forensic speaker recognition with Ahumada III: a public real-case database in Spanish", in Proceedings of Interspeech 2008, pp. 1493-1496, September 2008.

# LR Computation

- Step 1: the automatic system computes a score
  - No meaning on itself
    - 10 with respect to what?
  - In general, non-interpretable
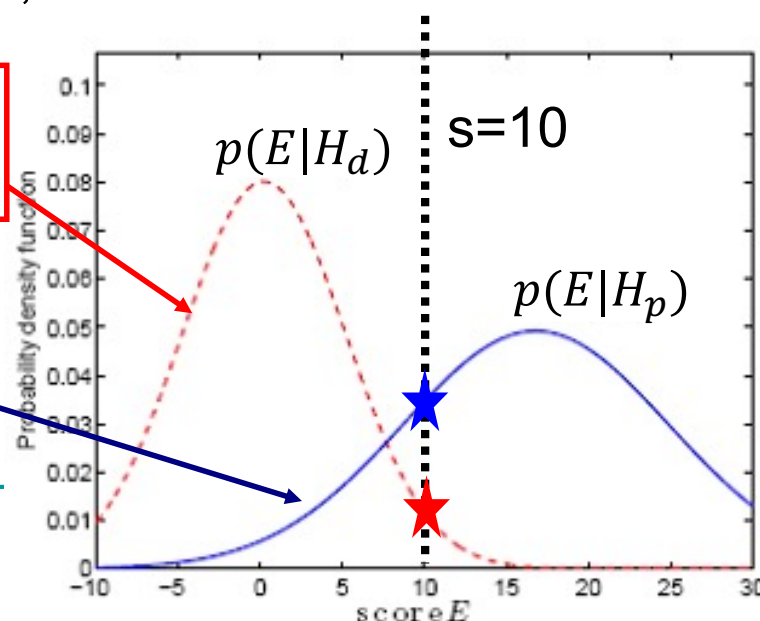    - Its range of variation is not known *a priori*

C ~~~~~

S ~~~~~ → [ Automatic Speaker Recognition ] → s=10

- Step 2: compute the LR
  - In this example, we use a Gaussian model

Different-Speaker Scores

Same-Speaker Scores

$$LR = \frac{0.35}{0.15} = 2.33$$

Weak support to prosecutor proposition ("same-source")



$p(E|H_d)$  s=10

$p(E|H_p)$

ison

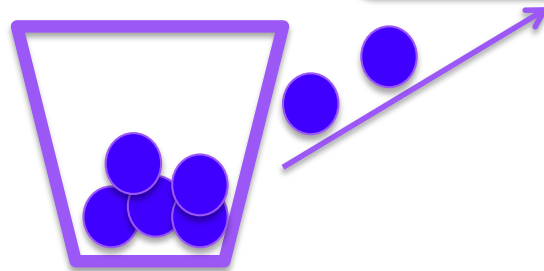<audias>

UAM

# Data for LR Computation



**Hd training scores (for the given case)**

**Hp training Scores (for the given case)**

LR Model

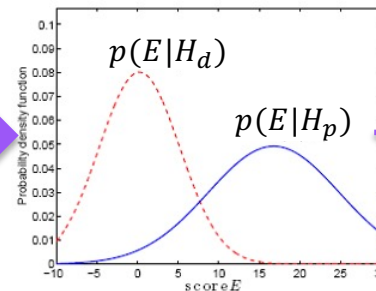Score prob. dentisies

$p(E|H_d)$

$p(E|H_p)$

C

S

Automatic Speaker Recognition

s  **Score**

**Score-LR transformation**

**LR from the case**

<audias>

UAM

# Data for LR Computation

**Hd training scores (for the**

**Selection of traning data is fundamental for LR computation!**

IEEE TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING, VOL. 15, NO. 7, SEPTEMBER 2007

## Emulating DNA: Rigorous Quantification of Evidential Weight in Transparent and Testable Forensic Speaker Recognition

Joaquin Gonzalez-Rodriguez, *Member, IEEE*, Phil Rose, Daniel Ramos, *Student Member, IEEE*,
Doroteo T. Toledano, *Member, IEEE*, and Javier Ortega-Garcia, *Member, IEEE*

## Protocol for the collection of databases of recordings for forensic-voice-comparison research and practice

Geoffrey Stewart Morrison ✉, Philip Rose & Cuiling Zhang

Pages 155-167 | Received 11 Jun 2011, Accepted 05 Oct 2011, Published online: 25 Jan 2012

Australian Journal of **Forensic Sciences**

**Score**

**-LR ...ation**

**...matic ...ker ...tion**

**Hp training Scores (for the g...**

**...m the case**

‹aud*i*as›

# A Guideline for the Validation of LR Methods
# (With Emphasis in Automatic Methods)

# Guideline: Validation of Forensic LR Methods

Contents lists available at ScienceDirect

## Forensic Science International

**ELSEVIER**

## A guideline for the validation of likelihood ratio methods used for forensic evidence evaluation

CrossMark

Didier Meuwly [a,b,*], Daniel Ramos [c], Rudolf Haraksim [d]

[a] Netherlands Forensic Institute, Laan van Ypenburg 6, 2497GB The Hague, The Netherlands
[b] University of Twente, Drienerlolaan 5, 7522NB Enschede, The Netherlands
[c] ATVS – Biometric Recognition Group, Escuela Politecnica Superior, Universidad Autonoma de Madrid, C/Francisco Tomas y Valiente 11, 28049 Madrid, Spain
[d] LTS5 – Signal Processing Laboratory, École Polytechnique Fédérale de Lausanne, Faculty of Electrical Engineering, Station 11, CH-1015 Lausanne, Switzerland

‹aud*i*as›

UAM

# Guideline: Validation of Forensic LR Methods

- Objective

  - Determine if a LR method is valid to be used in casework

  - All the validation process should be documented for transparency

  - Towards standardization of procedures (for biometrics)

- Validation process

  - Based on Empirical Testing

    - Data: still an issue

  - Performance assessment

    - Performance *characteristics*

      - *What aspect of performance should be measured?*

    - Performance *metrics*

      - *How to measure a characteristic?*

    - Performance graphical representations

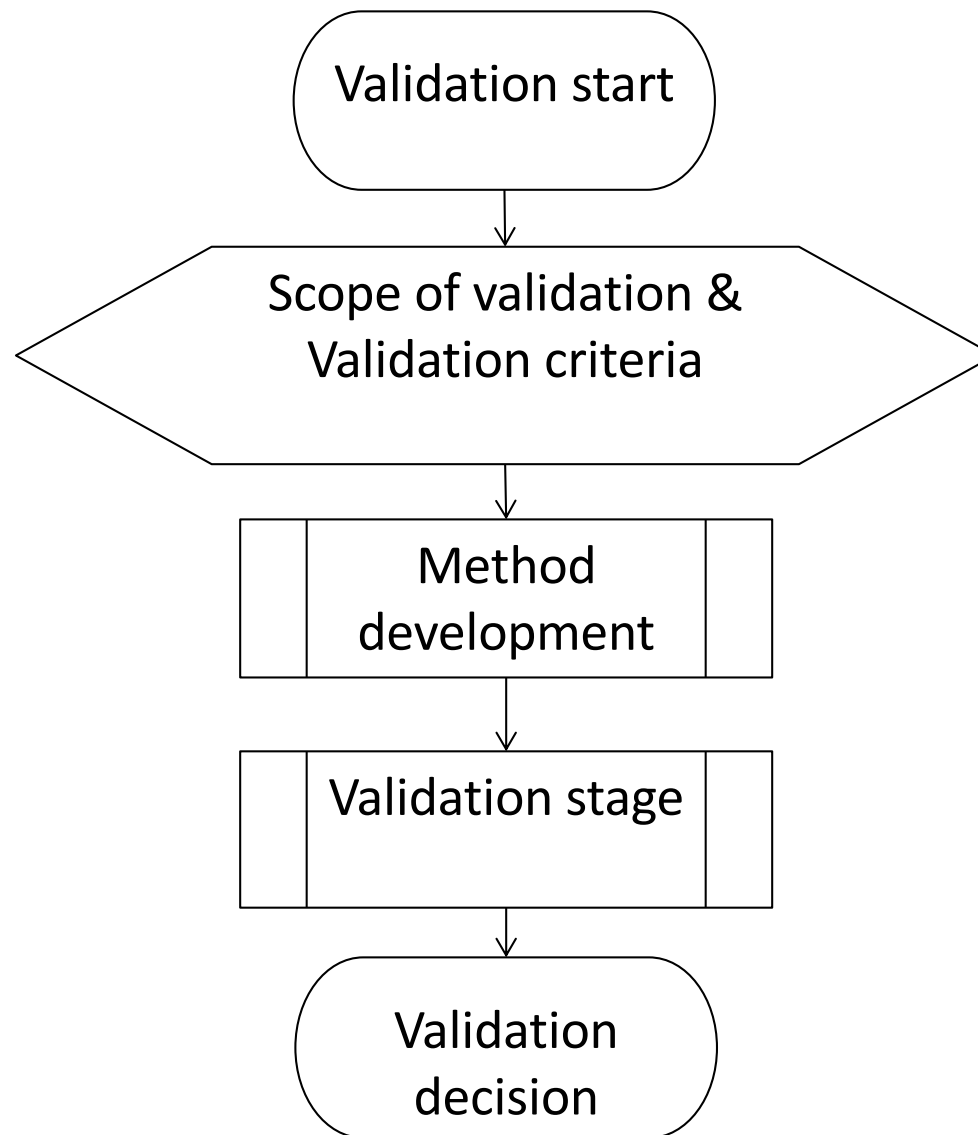      - *Ok, show me an illustrating plot*

Emphasis on Forensic Data
* Lab (development) performance
* Followed by forensic performance

ISO
Terminology

‹audias›

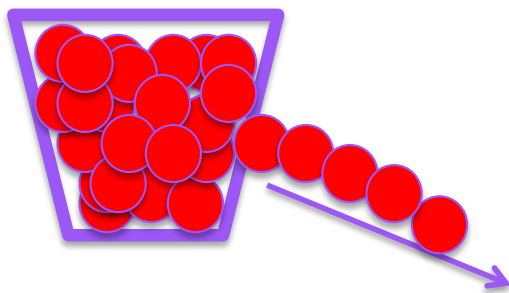# Validation – a process

# Performance Characteristics and Measures

| Performance Characteristic | Performance Metric | Graphical Representation |
|---|---|---|
| Accuracy | Cllr, EER | ECE plot<br>DET plot |
| Discriminating power | $Cllr^{min}$ | $ECE^{min}$ plot |
| Calibration | $Cllr^{cal}$ | Tippett plot |
| Robustness | Cllr, EER<br>LR range | ECE plot<br>DET plot<br>Tippett plot |
| Coherence | Cllr, EER | ECE plot<br>DET plot<br>Tippett plot |
| Generalization | Cllr, EER | ECE plot<br>DET plot |

‹audias›

UAM

# Performance Characteristics and Measures

| Performance Characteristic | Performance Metric | Graphical Representation |
|---|---|---|
| Accura... | | ECE plot |
| Discrim... power | | ...plot |
| Calibra... | | ...lot |
| Robustnes... | Cllr, EER | ECE plot<br>DET plot<br>...ett plot |
| Coherence | | ...plot<br>...plot<br>...ett plot |
| Generalization | Cllr, EER | ...plot<br>DET plot |

All measures require a validation set of LR values (computed automatically from a validation database)

Not restricted to these! The Guideline is thought to be open to modifications

‹audias›

# Experimental Set-Up

**Hd training scores (for the given case)**

**Simulated Case Scores**

Hd true    Hp true

Hd true score

LR Model

Score prob. dentisies

$p(E|H_d)$    $p(E|H_p)$

Score-LR transformation

**Hp training Scores (for the given case)**

**Hd-true LR**

# Experimental Set-Up



**Hd training scores (for the given case)**

**Hp training Scores (for the given case)**

LR Model

Score prob. dentisies

$p(E|H_d)$

$p(E|H_p)$
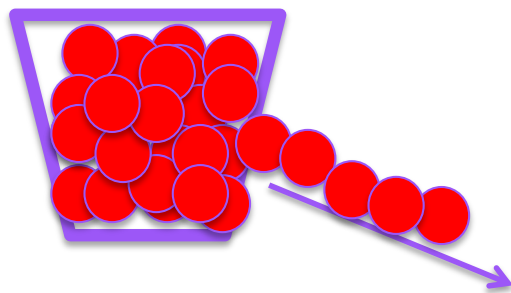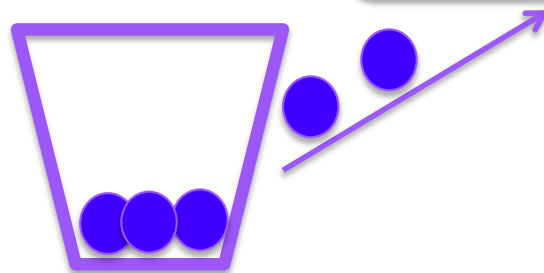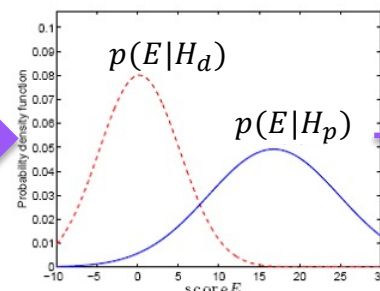
**Simulated Case Scores**

Hd true
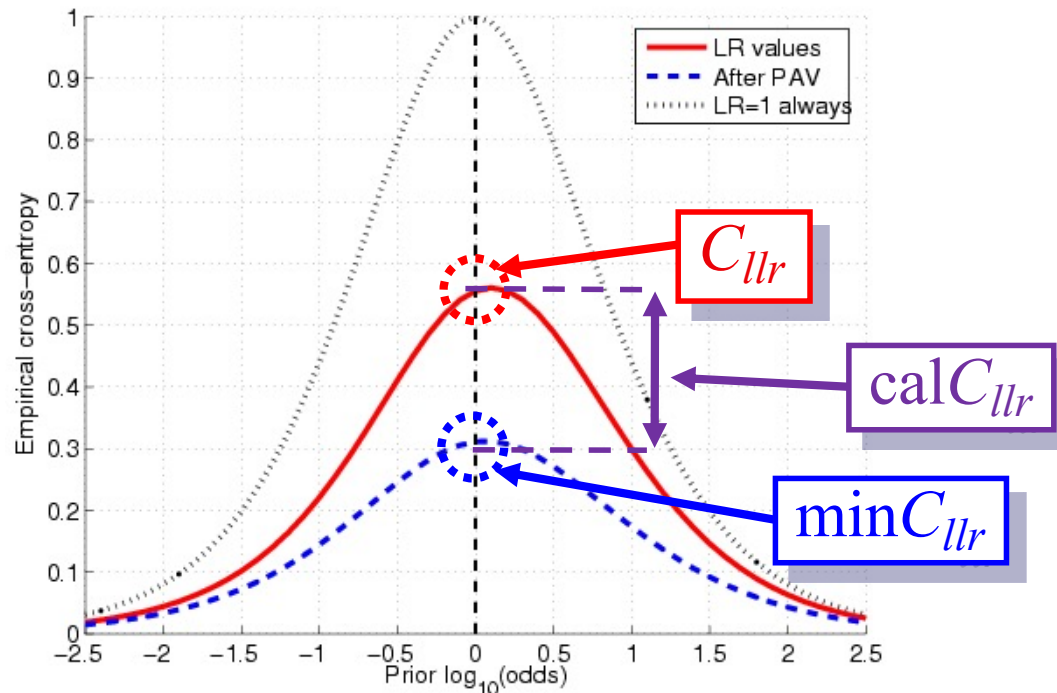
Hp true

**Hd true score**

**Score-LR transformation**

**Hp-true LR**

‹audias›

# Some Performance Metrics and Representations (Included in the Guideline)

# Guideline for Validation: Performance

**Primary**

**Secondary**
(measure behavior of the primary)

| Performance Characteristic | Performance Metric | Graphical Representation |
|---|---|---|
| Accuracy | $C_{llr}$ | ECE plot |
| Discriminating power | $C_{llr}^{min}$<br>EER | $ECE^{min}$ plot<br>DET plot |
| Calibration | $C_{llr}^{cal}$ | Tippett plot |
| Robustness | $C_{llr}$<br>EER<br>LR range | ECE plot<br>DET plot<br>Tippett plot |
| Coherence | $C_{llr}$<br>EER | ECE plot<br>DET plot<br>Tippett plot |
| Generalization | $C_{llr}$<br>EER | ECE plot<br>DET plot |

‹aud*i*as›

# Empirical Cross-Entropy Plots and $C_{llr}$



**ECE curve: Accuracy (the lower the better)**

Calibration (red – blue)

Discrimination (blue curve)

D. Ramos, J. Gonzalez-Rodriguez, G. Zadora and C. Aitken. "Information-theoretical Assessment of the Performance of Likelihood Ratios". Journal of Forensic Sciences (under minor revision)
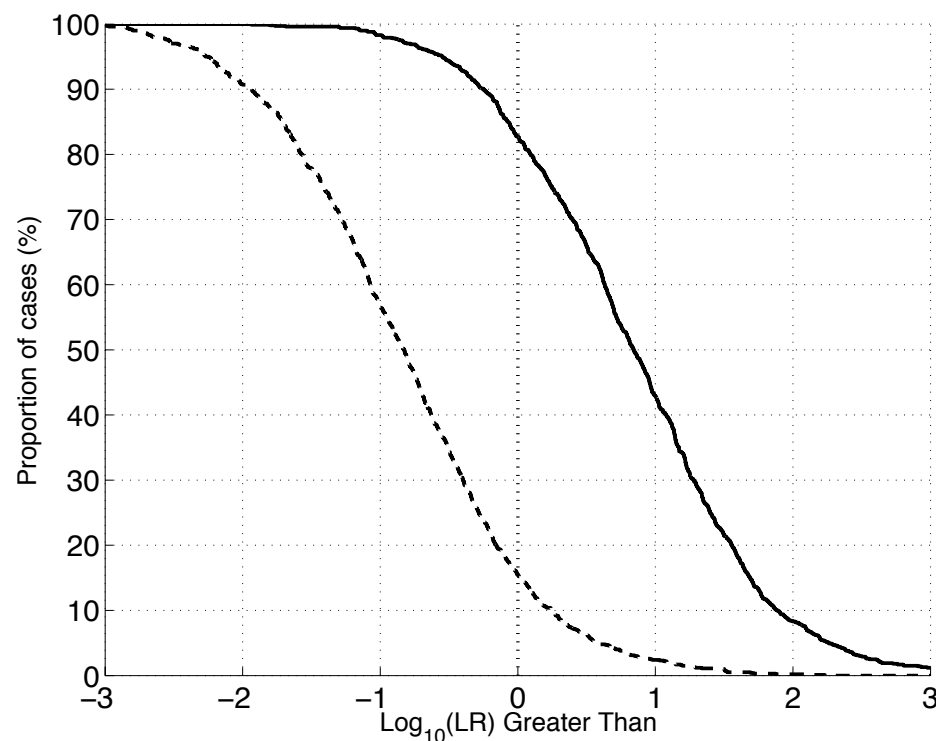
- http://arantxa.ii.uam.es/~dramos/software.html

- **Summarizing metric: $C_{llr}$**

Niko Brümmer [a,b,*], Johan du Preez [b]
Application-independent evaluation of speaker detection
Computer Speech and Language 20 (2006) 230–275

⟨aud*i*as⟩

UAM

# Tippett Plots

- Cumulative distribution of LR values
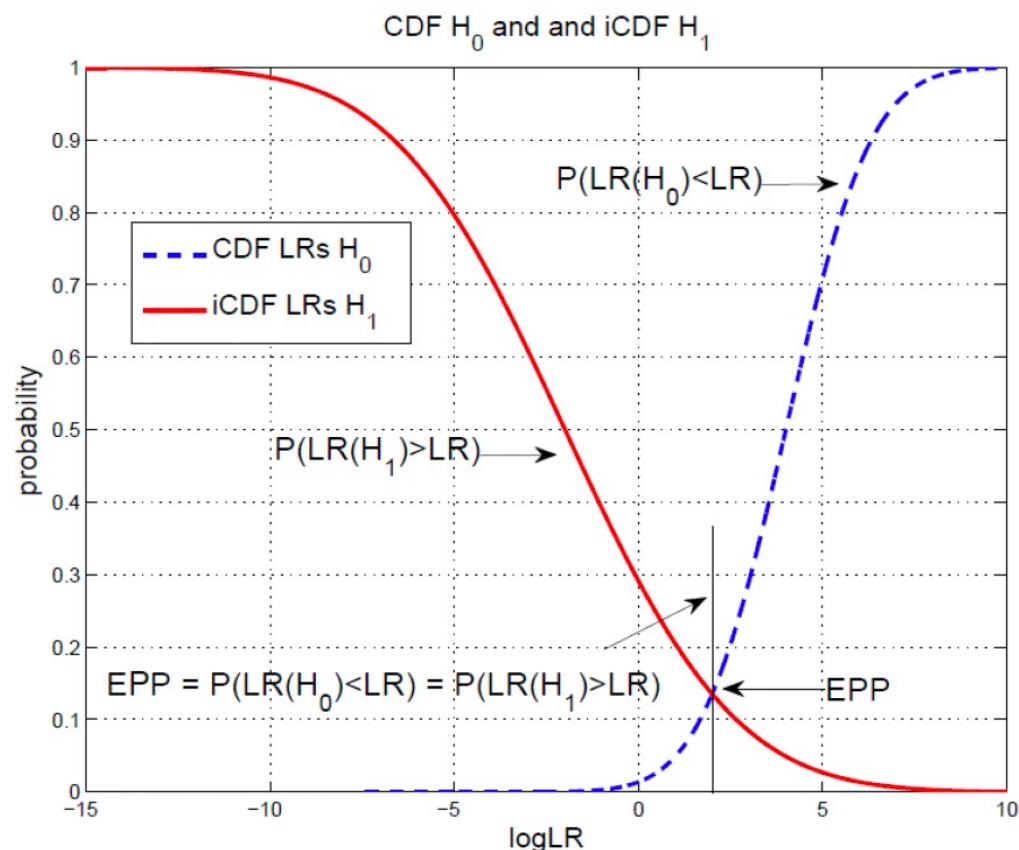


Accuracy

Calibration

Discrimination

WARNING: Tippett plots do not measure them explicitly!

- Summarizing metric:

  - Rates of misleading evidence

# Tippett Plots (Type II representation)

- Cumulative distribution of LR values



Accuracy

Calibration

Discrimination

WARNING: Tippett plots do not measure them explicitly!

**Methodological Guidelines for Best Practice in Forensic Semiautomatic and Automatic Speaker Recognition**

ENFSI
EUROPEAN NETWORK OF FORENSIC SCIENCE INSTITUTES

# Acknowledgements

# Acknowledgements

- People who contributed in one way or another to this presentation:

  - Didier Meuwly (NFI)

  - Rudolf Haraksim (European Comission Joint Research Center)

  - Charles Berger (NFI)

‹audias›

UAM