

ISS World 2007

DUBAI, UAE – February 27, 2007

**THE SEARCH FOR RESULTS  
IN VOICE ANALYSIS:  
how different identification  
technologies can work together  
effectively**

**Luciano Piovano**

Government Intelligence Solutions, V.P.

# Loquendo Voice Technologies for COMINT



Forensics



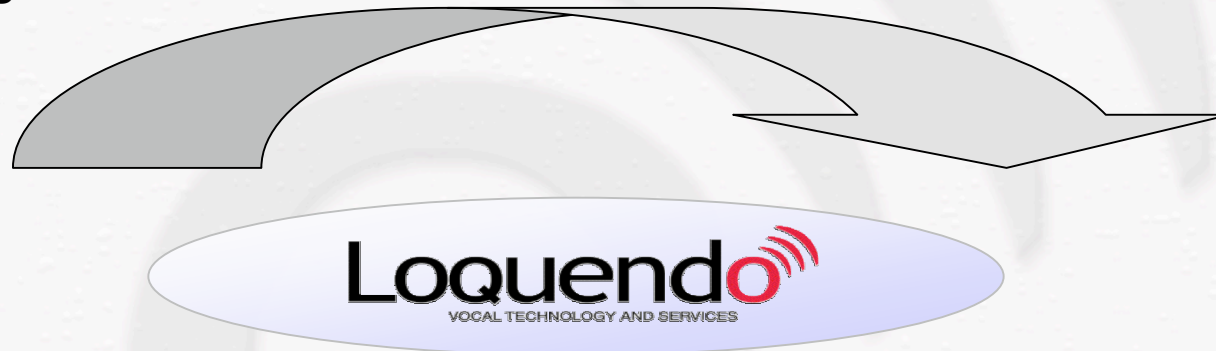
LEA  
investigation



Counter Terrorism  
Intelligence



Battlefield



- **Speaker Recognition** through Voice-Print comparison of free speech
- **Language Identification** – also for dialect/accent recognition
- **Keyword spotting** to detect words of special interest to investigators

# Different scenarios for Speaker Identification applications

## Intelligence/CounterTerrorism

- Huge volume of intercepts
- Various targets (sometimes several hundred)
- Different languages spoken
- Emphasis on spotting targets as calls come in
- Limited accuracy usually sufficient
- Strict time constraints
- Usually no need to gather evidence

## Intelligence Agencies

## Criminal Investigation

- Limited number of intercepted calls
- Fewer targets
- Spoken language generally known in advance
- Each call can be analyzed
- High accuracy required
- Looser time constraints
- Intercepts may have to be produced as evidence

## Law Enforcement Agencies

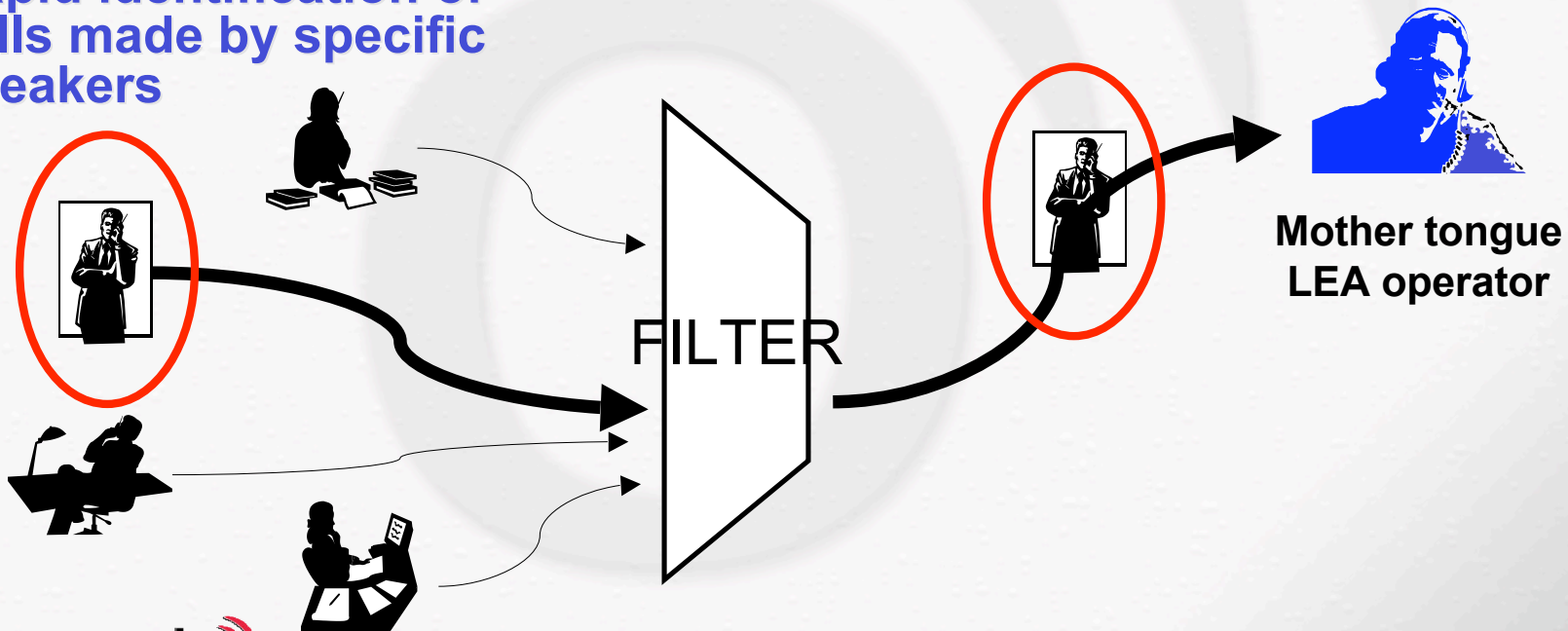
# Intelligence / Counter-Terrorism



- Huge volume of telephone intercepts
- Hundreds of target speakers
- Different languages spoken
- Spotting of targets as calls come in
- Multiple investigation scenarios

## Objective:

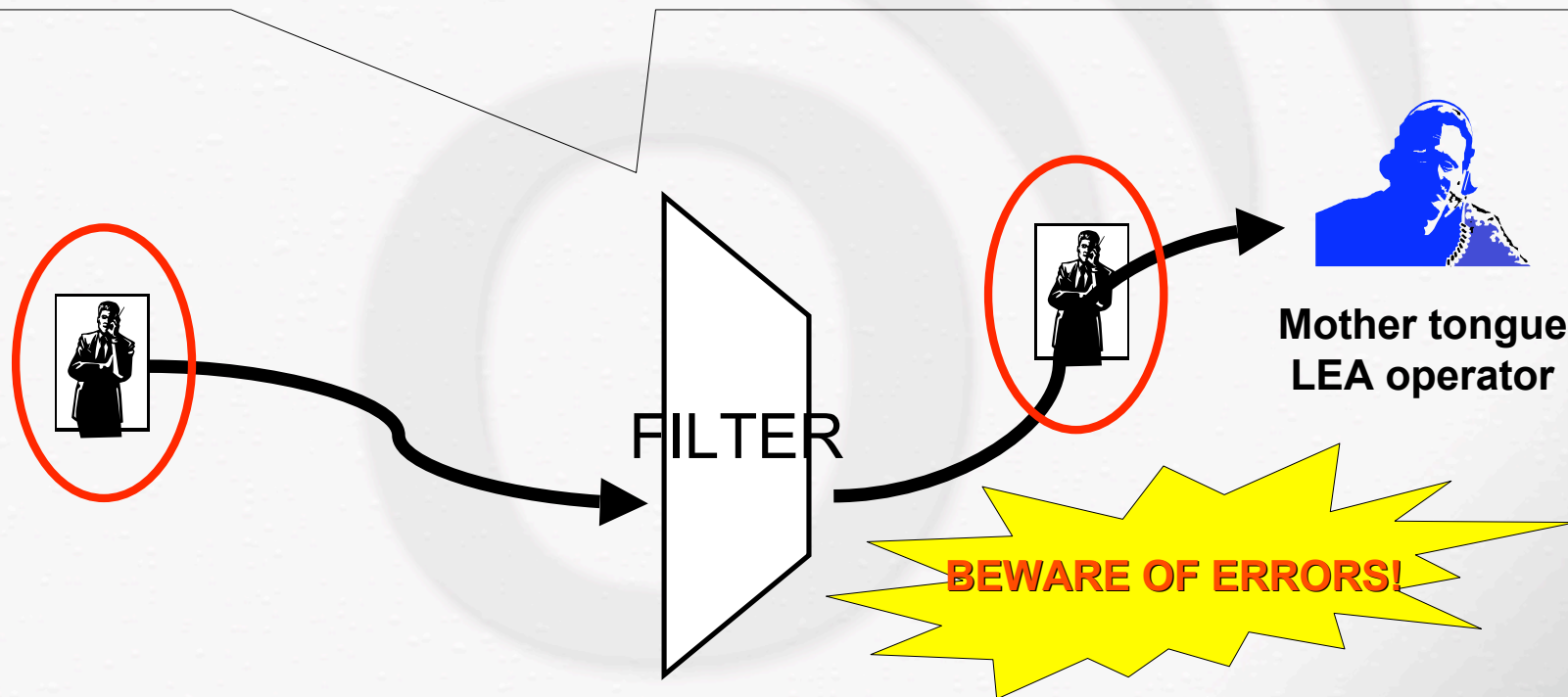
Rapid identification of calls made by specific speakers





# Elements used for Filtering

- 1) Investigative knowledge
- 2) Network parameters (CLI, DN, IMEI code,...)
- 3) Speech content (spoken language, keywords,...)
- 4) Speaker features (biometrics, gender, emotion, ...)



# LEA Investigations – An example

## Finding for a phone call in an international trunk traffic



- Int'l trunk
- ...
- PABX

How can I spot the right calls without infringing other people's privacy?

**Automatic real-time extraction of calls matching target Voice Prints**

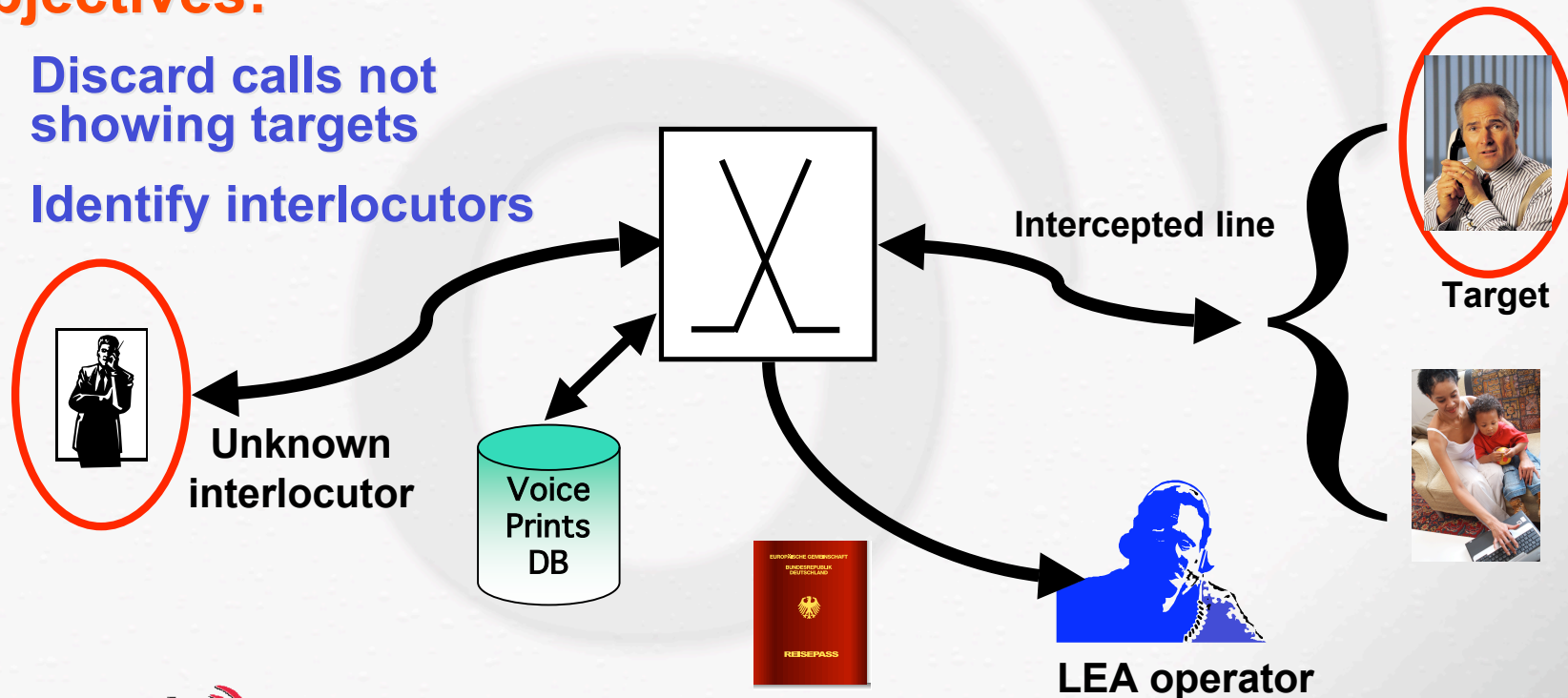
# Criminal Investigations



- Limited volume of telephone intercepts
- Dozens of target speakers
- Spoken languages known in advance
- Ranking of intercepted calls
- Usually narrow investigation scenarios

## Objectives:

1. Discard calls not showing targets
2. Identify interlocutors



# Speaker Identification through Biometrics

- Every voice contains acoustic-phonetic features that can be extracted, amplified, stored and used to build Voice Prints (VPs)
- VPs are based on “certified” audio recordings
- Like fingerprints, VPs can also be used for comparison with elements gathered in the field
- Accuracy scores are intrinsically statistical ( $P_{\text{Err}} > 0$ )
- In telephone intercepts, voice is the only “signature” that can be assessed



**Each individual can be assigned a Voice Print to determine his/her identity**

## **LFSI – Loquendo Free Speech Identification**

- **Software technology allows the identification of speakers in natural speech telephone calls**
- **Phonetic GMM recognition**
- **Search for several targets at the same time**
- **Real time processing of audio files**
- **Provides normalized scores for every “voice print – audio file” pairing**
- **Language independent**
- **Channel independent (mobile, fixed, VoIP)**
- **Excellent accuracy results (obtained at NIST '06 SRE)**



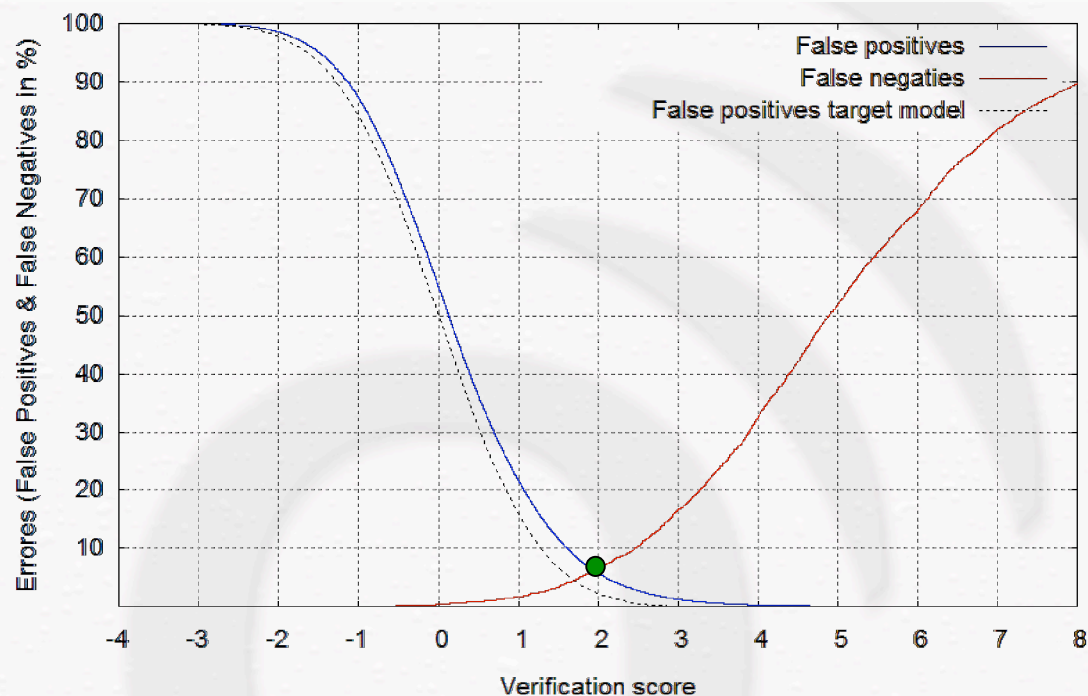
# What about the accuracy?

Elements to consider:

- 1) **A *a priori* probability of correct target interception**
- 2) **False Alarms (False Positives) FA**
  - 1) Should tend to zero in authentication applications
  - 2) May be more acceptable in Intelligence applications
- 3) **False Miss (False Negatives) FM**
  - 1) Normally unacceptable in Intelligence
  - 2) More acceptable in authentication applications
- 4) **Impossibility of optimizing both error rates (FA and FM) at the same time**

# System Characterization (1)

## LFSI Error Rate Plot



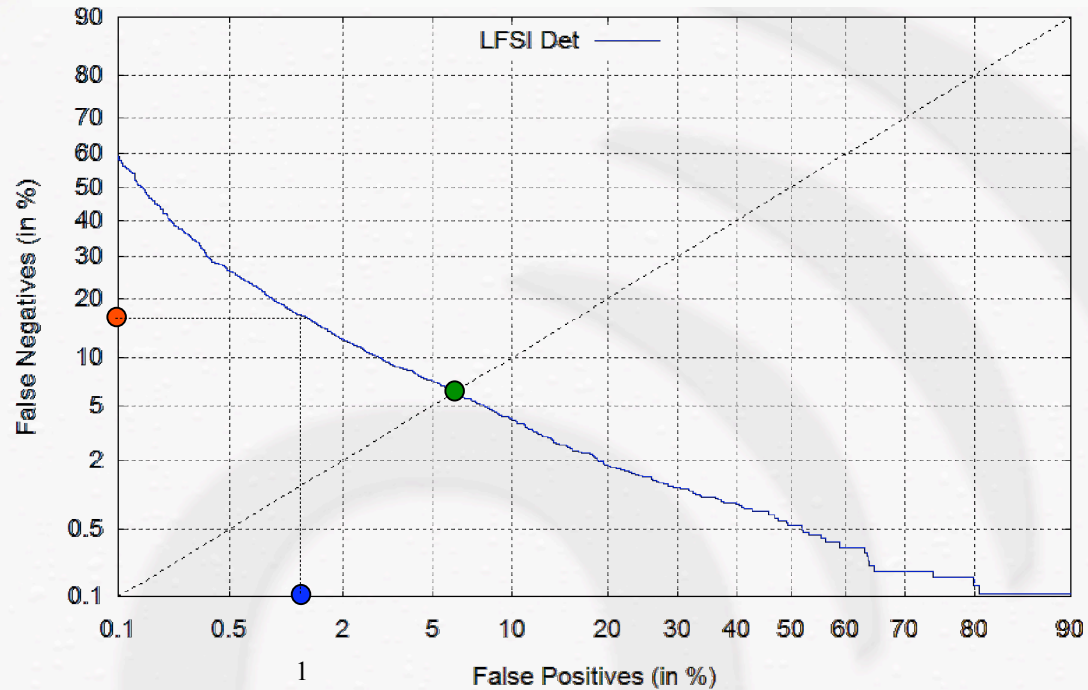
**False Positives = False Alarms**

**False Negatives = False Miss**

**Equal Error Rate**

# System Characterization (2)

## LFSI Detection Error Tradeoff Plot



**False Positives = False Alarms**

**False Negatives = False Miss**

**Equal Error Rate**

## Enough accuracy? An example

a) Working Point where  $P_{\text{FA}|1\text{target}} = 1\%$

⇔ then an average of 1 call out of 100 will be wrong  
with reference to each specific target

If you look for 100 targets

$$P_{\text{FA}|100\text{targets}} = 1 - P_{\text{right}} = 1 - (0,99)^{100} = 63\%$$

**USUALLY UNACCEPTABLE**

b) Working point where  $P_{\text{FA}|1\text{target}} = 0,1\%$

$$P_{\text{FA}|100\text{targets}} = 9\%$$

**MUCH BETTER**

## How to improve accuracy

# What's next?

We have only considered point 4): **Voice Prints comparison**

- 1) Investigative knowledge
- 2) Network parameters (CLI, DN, IMEI code,...)
- 3) Speech content (**Spoken Language**, keywords,...)
- 4) **Speaker features** (VP biometrics, **gender**, emotion, ...)

So now let's consider point 3): **Spoken Language**  
and 4) **Gender**



## Language Identification (L2I)

- A model of each individual language can be made using its characteristic features
- A likelihood score can be calculated from comparing speech recordings to language models
- The likelihood scores indicate which language is being spoken
- Based on sufficient speech recordings in a specific language coming from a variety of speakers, the language identification engine can be trained to recognize new languages
- Also suitable for dialects (may be less precise)
- Suitable for Accent Identification (development in progress)

# Gender Identification

- **A model of each gender (male/female) can be made using general voice features**
- **A likelihood score can be calculated from comparing speech recordings to gender models**
- **Suitable for filtering calls (men are often targets)**

## Example of combinations of different filters (1/2)

### Investigative assumptions

Example involves an Italo-American company

One branch in the US, one in Italy

Drug-trafficking involved

Bad guys are Italian (could be located in Italy and USA)

1000 calls a day on that link

50% involve women

### Voice Print library knowledge/assumptions

100 targets related to drug trafficking:

10 women

90 men, of which

30 Americans

60 Italians

## Example of combinations of different filters (2/2)

### Technology assumptions

$FA_{\text{Gender Id}} \cong FA_{\text{Speaker Id}} \cong FA_{\text{Language Id}}$

**Then** the comparison will be made between:

60 VPs belonging to Italian men involved in drug trafficking

The percentage of the 1000 calls/day where only men are present

The system will first perform a comparison to check gender and then if only men are involved in the call it will perform the Italian male VPs comparison

**Therefore:**

60 VPs instead of 100  $\Rightarrow FA_{\text{total}} = 5,8\%$  (instead of 9%)

Applied to 500 calls instead of 1000 per day

Without any classification there would be an average of 90 FA/day

**WITH THE FILTERS  $\Rightarrow$  29 FA/day**

# CONCLUSIONS

**Intelligent adoption of different filtering criteria may improve the chances of a successful search and reduce time wasted on analysis of irrelevant material**

**The search for specific targets (based on Voice Print comparison) can be enhanced if individuals are also grouped according to the languages they speak/ their gender**

**Loquendo provides solutions combining Speaker Identification and Language Identification as well as Gender Identification**



# CONTACTS

**LOQUENDO booth  
at ISS World exhibition**

[security@loquendo.com](mailto:security@loquendo.com)

**THANK YOU !**