

A ROBUST SPEECH COMMAND RECOGNIZER FOR EMBEDDED APPLICATIONS

Alexandre Maciel

Arlindo Veiga

Carla Lopes

Cláudio Neves

Fernando Perdigão

José David Lopes

Luís de Sá

INSTITUIÇÕES ASSOCIADAS:



INSTITUTO
SUPERIOR
TÉCNICO



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

universidade
de aveiro



Inovação

SIEMENS
Communications



instituto de
telecomunicações

creating and sharing knowledge for telecommunications

© 2005, it - instituto de telecomunicações. Todos os direitos reservados.

Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API
- Results
- Demo
- Conclusions

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

2 | Gaia 29th July 2008



instituto de
telecomunicações

Summary

🔗 **Human-Machine Interaction**

🔗 Noise Reduction Techniques

🔗 Model Training

- ⋮ Audio Database

- ⋮ Acoustic Modelling

🔗 Decoder

🔗 API

🔗 Results

🔗 Demo

🔗 Conclusions

INSTITUIÇÕES ASSOCIADAS:

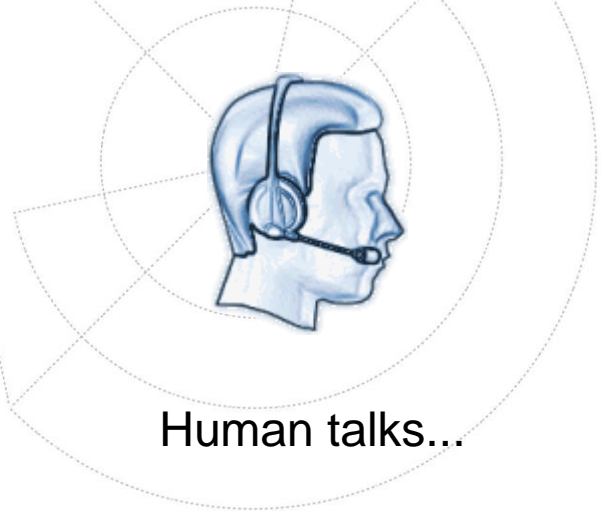


Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

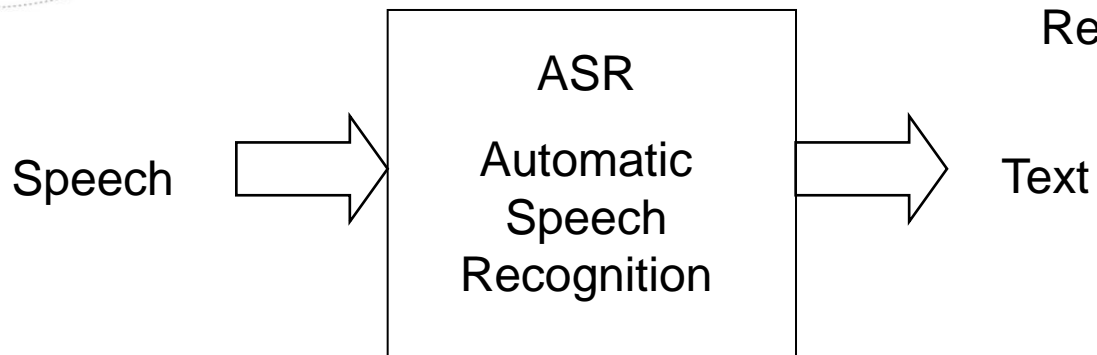


Human-Machine Interaction

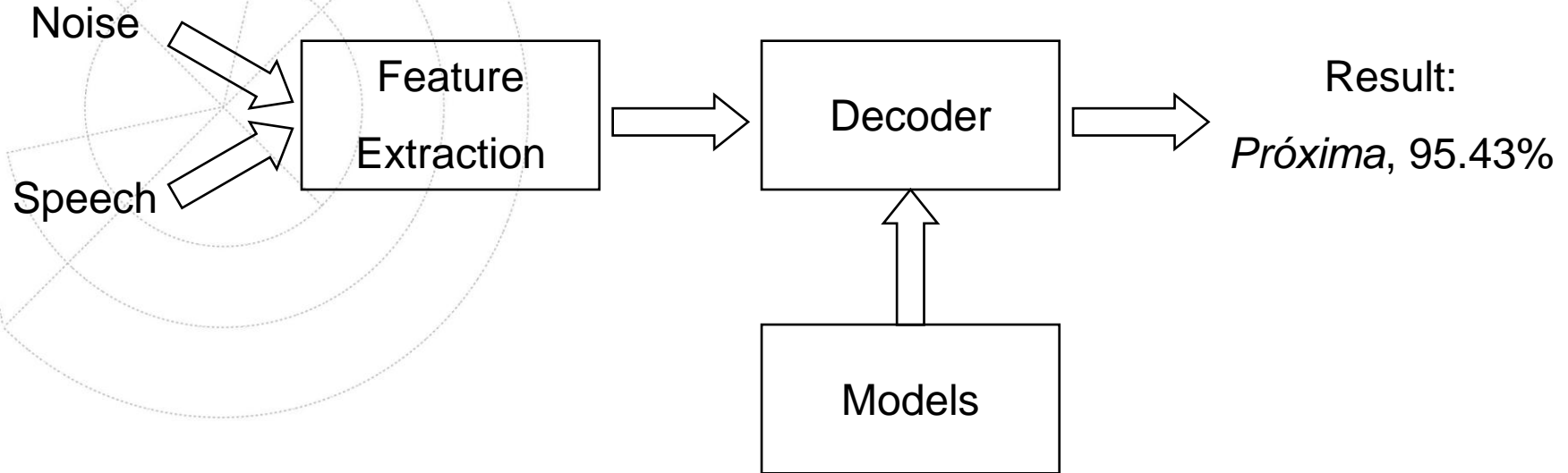


The machine performs:

Automatic Speech
Recognition



Human-Machine Interaction



🔗 Feature Extraction

- : Noise reduction techniques
- : Voice Activity Detection (VAD)

🔗 Models (previously trained)

🔗 Decoder ⇒ result (textual information or action and a confidence measure)

Summary

- Human-Machine Interaction
- Noise Reduction Techniques**
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API
- Results
- Demo
- Conclusions

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

6 | Gaia 29th July 2008



instituto de
telecomunicações

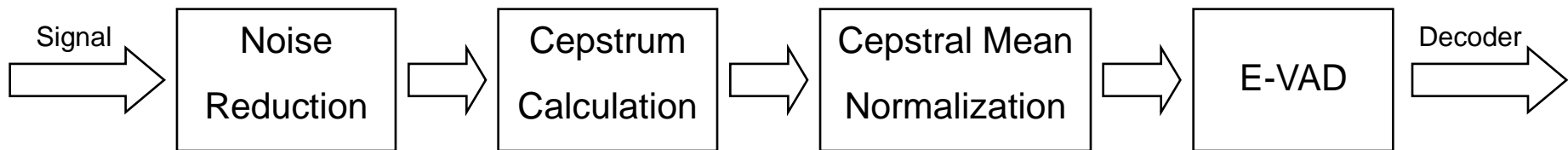
Noise Reduction Techniques

Algorithm E-AFE:

- based on an ETSI standard (ETSI ES 202 050);
- performs noise-robust feature extraction;
- Neves, C., Veiga, A., Sá, L. and Perdigão, F., 2008. *Efficient Noise-Robust Speech Recognition Front-End based on the ETSI Standard*. In proc. of ICSP'2008. Beijing, China.

Algorithm E-VAD:

- based on logarithm energy.



Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training**
 - Audio Database**
 - Acoustic Modelling**
- Decoder
- API
- Results
- Demo
- Conclusions

Model Training: Speech Database

✂ Speaker distribution
by gender and acoustical
environment

✂ 184 hours of
recordings

✂ 232,000 audio files

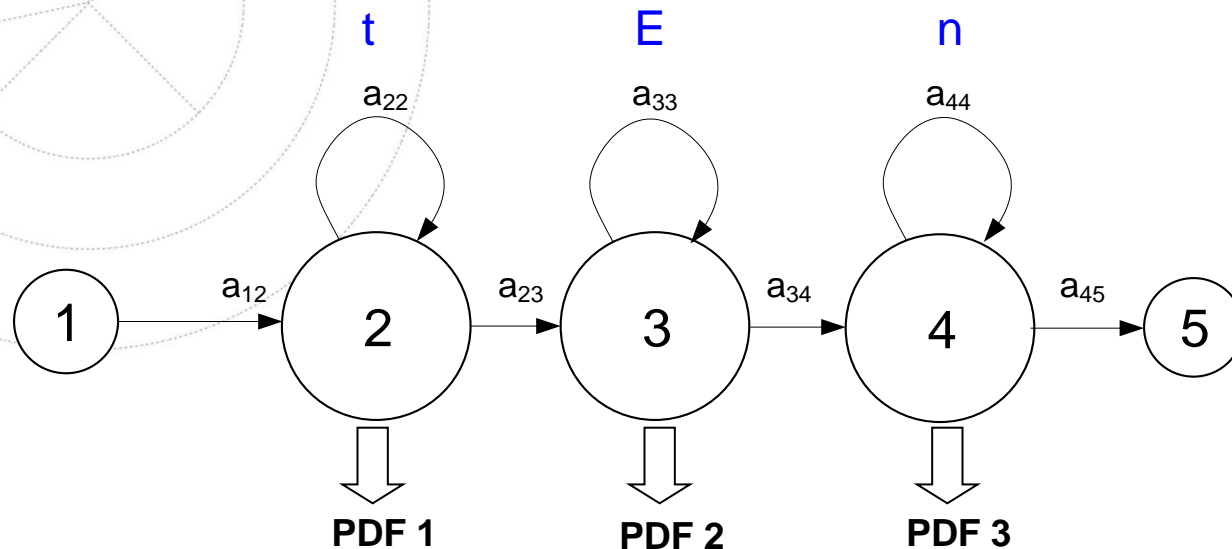
✂ DB splitting:
:75% for training
:20% for testing
:5% for development

<i>Gender</i>	<i>TVFL</i> <i>(clean)</i>	<i>TVF</i> <i>(factory)</i>	<i>TVV</i> <i>(vehicle)</i>
Female	103	20	9
Male	197	16	23
Total	300	36	32

Model Training: Acoustic Modelling

Acoustic models based on Hidden Markov Models

Example of a model for the word “ten” (t E n)



- Number of states
- Transition probabilities
- Probability density function

Model Training: Acoustic Modelling

- Whole-word models
- Monophone models
- Triphone Models

abril sp 6 b r i l ~ sp

acrescentar sp 6 k r @ S s e ~ t a r sp

portagem sp p u r t a Z e ~ j ~ sp

próximo sp p r O s i m u sp

abril sp 6+b 6-b+r b-r+i r-i+l~ i-l~ sp

acrescentar sp 6+k 6+k-r k-r+S ... t-a+r a-r sp

portagem sp p+u p-u+r ... Z-e~+j~ e~j~ sp

próximo sp p+r p-r+O ... i-m+u m-u sp

Summary

✂ Human-Machine Interaction

✂ Noise Reduction Techniques

✂ Model Training

⋮ Audio Database

⋮ Acoustic Modelling

✂ **Decoder**

✂ API

✂ Results

✂ Demo

✂ Conclusions

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

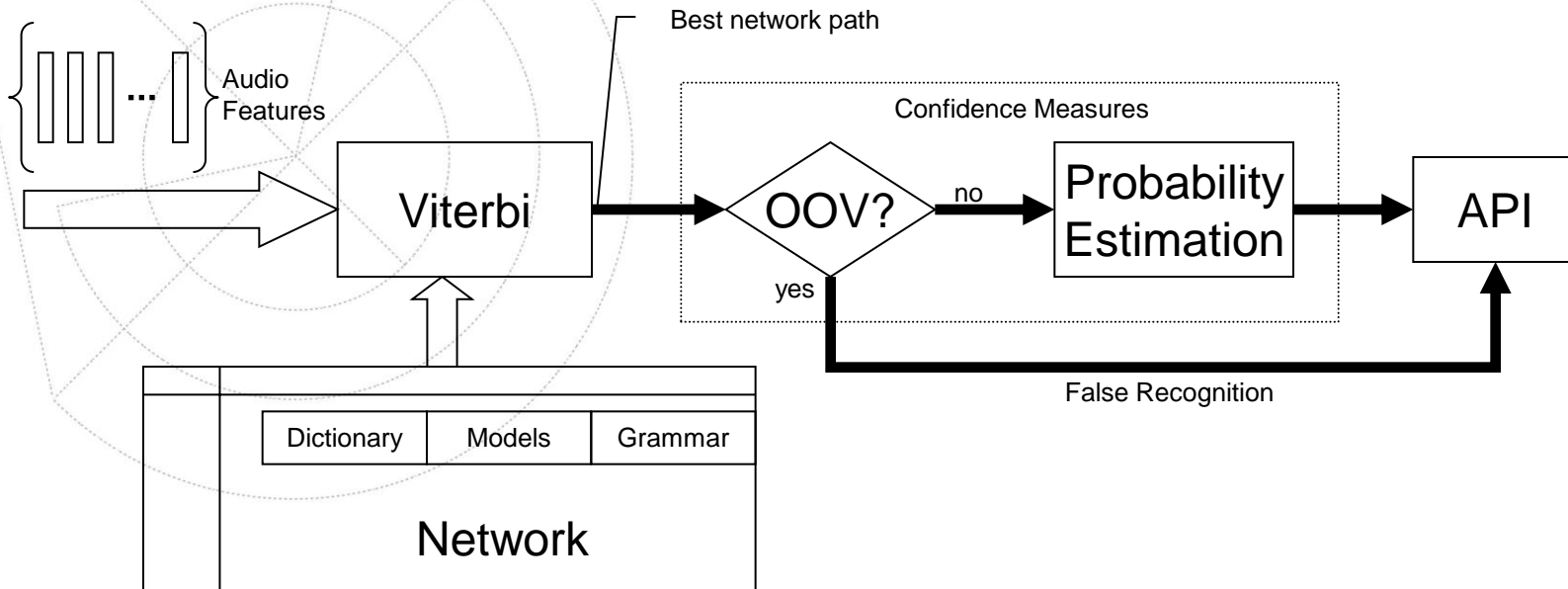
SIGMAP 2008

12 | Gaia 29th July 2008



instituto de
telecomunicações

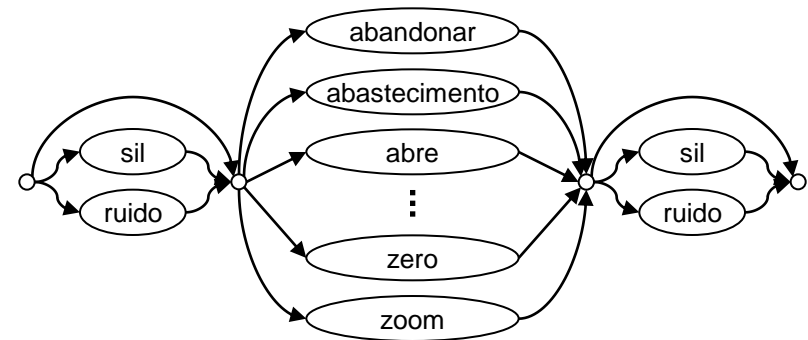
Decoder



Decoder Step

- 1: Built network;
- 2: Wait for VAD notification;
- 3: Process all audio features;
- 4: Select winner path;
- 5: Test if OOV;
 - If OOV, notify "False Recognition" to API;
 - Else, estimate result probability and send to API
- 6: Go to 2:

Commands' Grammar



Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API**
- Results
- Demo
- Conclusions

INSTITUIÇÕES ASSOCIADAS:



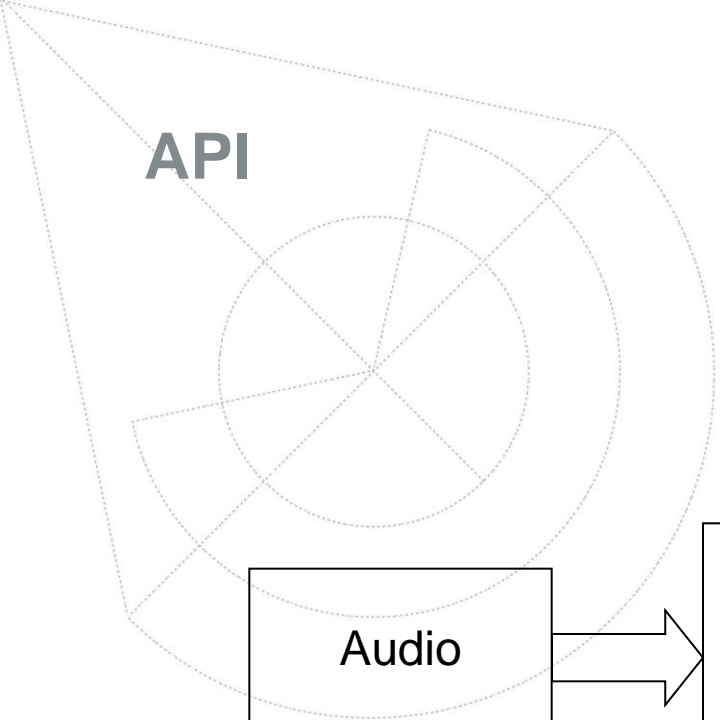
Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

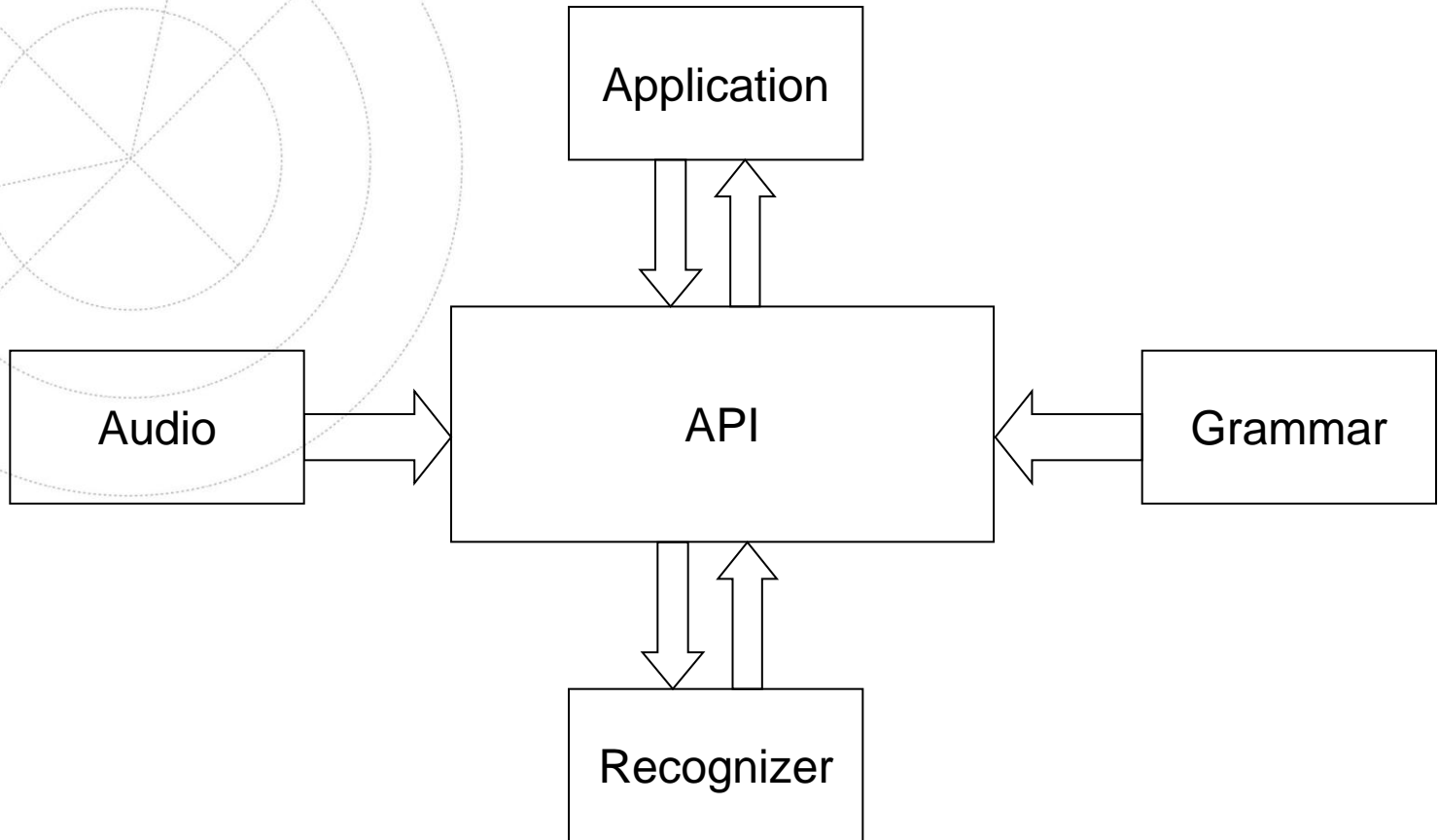
14 | Gaia 29th July 2008



instituto de
telecomunicações



API



Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API
- Results**
- Demo
- Conclusions

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

16 | Gaia 29th July 2008



instituto de
telecomunicações

Results

How to measure performance?

Original test file

“*/TVFL_FLI_LIM_F46_024@8.lab”

sil
opcoes
sil
.

“*/TVFL_FLI_LIM_MG6_050@8.lab”

sil
frente
sil
.

“*/TVFL_FLI_LIM_MH8_219@8.lab”

sil
vai
sil
.

sil
deletion



Recognized file

“*/TVFL_FLI_LIM_F46_024@8.rec”

opcoes
sil
.

“*/TVFL_FLI_LIM_MG6_050@8.rec”

sil
frente
sil
.

Successful
recognition

Word
misrecognized



“*/TVFL_FLI_LIM_MH8_219@8.rec”

sil
dezoito
sil
.

Results

Front-end evaluation:

- ETSI Front-end: 94.88 %
- Efficient Front-End: 96.88%

Acoustic modelling evaluation

- Whole-word models: 96.76%
- Monophone models: 89.28%
- Triphone models: 97.03%

Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API
- Results
- Demo**
- Conclusions

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

19 | Gaia 29th July 2008



instituto de
telecomunicações

Demo

 Voice controlling presentation

 Recognizer working in real-time in a vehicle environment

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

20 | Gaia 29th July 2008



instituto de
telecomunicações

Summary

- Human-Machine Interaction
- Noise Reduction Techniques
- Model Training
 - Audio Database
 - Acoustic Modelling
- Decoder
- API
- Results
- Demo
- Conclusions**

INSTITUIÇÕES ASSOCIADAS:



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra

SIGMAP 2008

21 | Gaia 29th July 2008



instituto de
telecomunicações

Conclusions

- ✎ The recognizer works in real-time over low performant hardware
- ✎ Triphone models are more likely to be used
- ✎ The ETSI noise reduction front-end might be biased to the databased used in development



Questions?

:Thank You

INSTITUIÇÕES ASSOCIADAS:



INSTITUTO
SUPERIOR
TÉCNICO



Faculdade de Ciências
e Tecnologia da
Universidade de Coimbra



universidade
de aveiro



Inovação



End



instituto de
telecomunicações

creating and sharing knowledge for telecommunications

© 2005, it - instituto de telecomunicações. Todos os direitos reservados.