

UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE



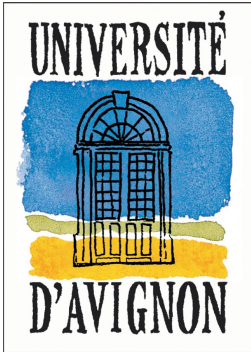
Speech Processing

LIA : Laboratory of Computer Science, University of Avignon

Language Processing group :

- about 12 researchers/18 Phd Students
- Topics : rich transcription, speech analytics, dialogue systems, natural language processing

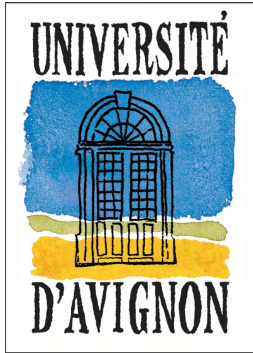
Georges Linarès (georges.linares@univ-avignon.fr)



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speech Processing

- Outline :
 - _ Introduction
 - _ What is Speech ?
 - _ Speech as a part of the Artificial intelligence project
 - _ An historical view of speech processing
 - _ Generalities about statistical speech processing
 - _ Speech recognition Systems : state of the art
 - _ Speaker identification
 - _ Practical work : structuring video database by analysing spoken contents



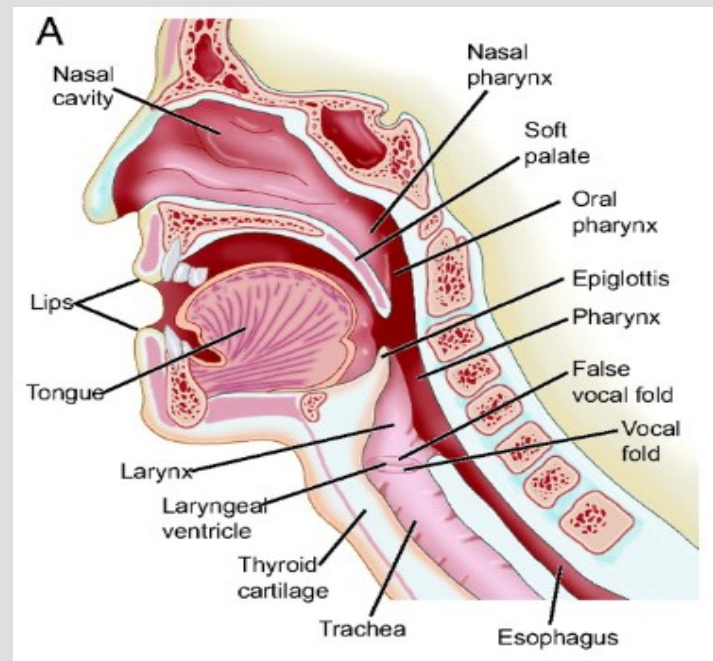
Speech Processing

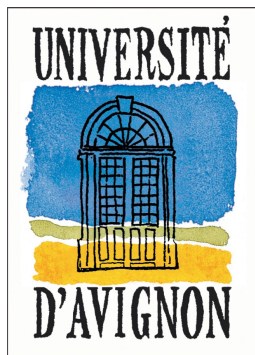
- What is Speech ?
 - Verbal mean of communication
 - Speech is not *written language*
 - Technically :
 - Sounds produced by the vocal folds, the breathing, the articulatory system :
 - Source : vocal folds (pitch)
 - Modulation due to articulators

What is Speech ?

Phonological point of view : Speech is a production of the human vocal system

Articulators :





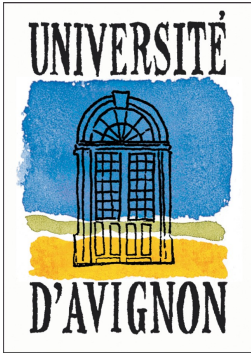
UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

What is Speech ?

***Sociological* point of view :**

Speech is the main communication mean of human communities





UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speech Processing : the engineer point of view

- Speech is a mean to exchange information
- Speech records contain information related to :
 - The semantic contents
 - The speaker identity, emotional state, intents, ...
 - The context
- Speech is useful to :
 - Driving machines by voice commands
 - Extract informations

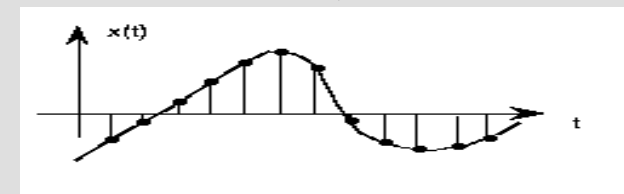
Speech Processing : the engineer point of view

Representation of speech signal

Analogic speech signal (acquisition)

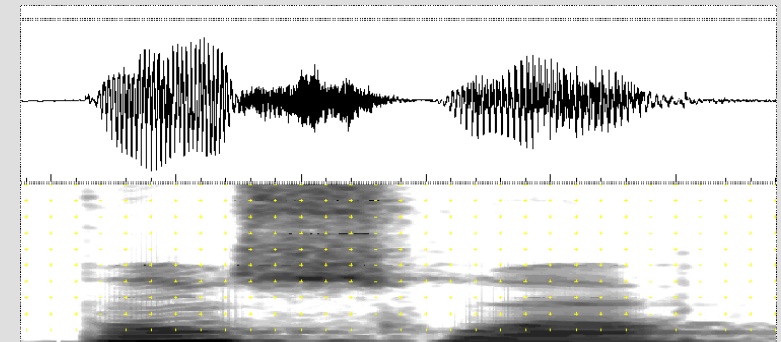


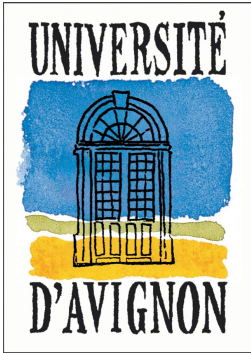
Digitalization, quantification



2D representation :

- Feature extraction
- Analysis in a sliding window





UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

What is Speech ?

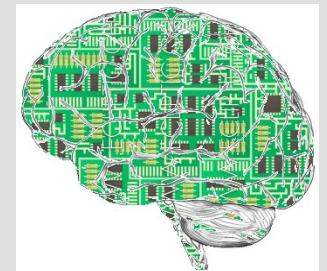
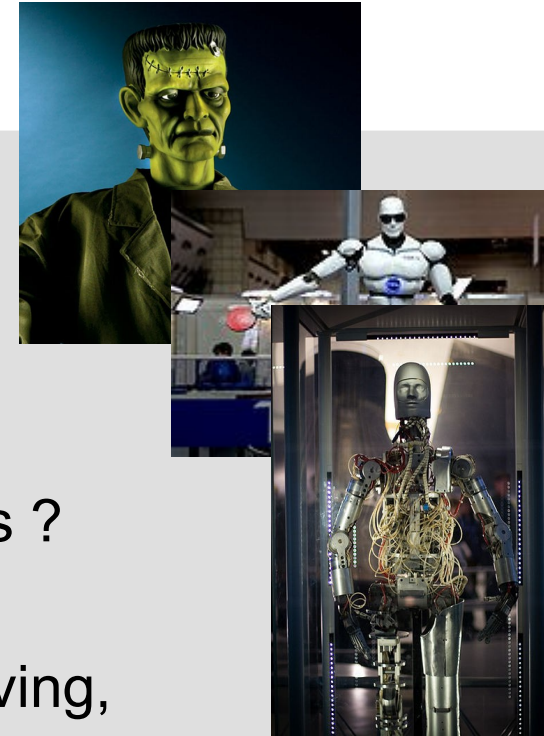
Speech is a complex communication mean :

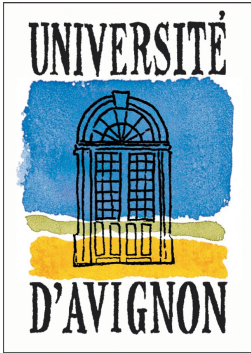
- Human languages are complex (acoustic/linguistic structure)
- Human thinking is complex
- A spoken message results from :
 - the context of the discourse
 - acoustic environments
 - Who are the speakers and the listeners
- Speech understanding relies on :
 - Various level of knowledge (lexical, linguistic, pragmatic, semantic, ...)
- High variability of speech (knowledge sources, media, contents)

Speech processing and A.I.

A.I. project : an *intelligent* machine :

- android?
- a thinking machine ?
- a machine able to perform complex tasks ?
- Simulation of human capacities :
 - Perception (recognition), problem solving, decision making
- **Industrial applications**
- **Building system tractable by computers**

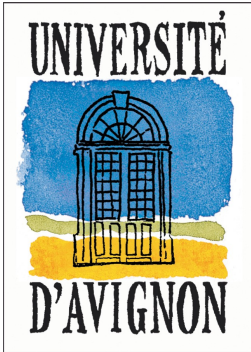




UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

History of Speech processing

- Human-inspired approaches (1970-)
 - Knowledge-based approaches
 - But *humans are still mysterious for the science*
 - Neuromimetic approaches
 - Artificial brains seem easy to build (not so clear...)
 - **machine learning** *versus* **machine knowing**
- Nowadays (1990-):
 - *Machine Learning* but **statistical modeling**



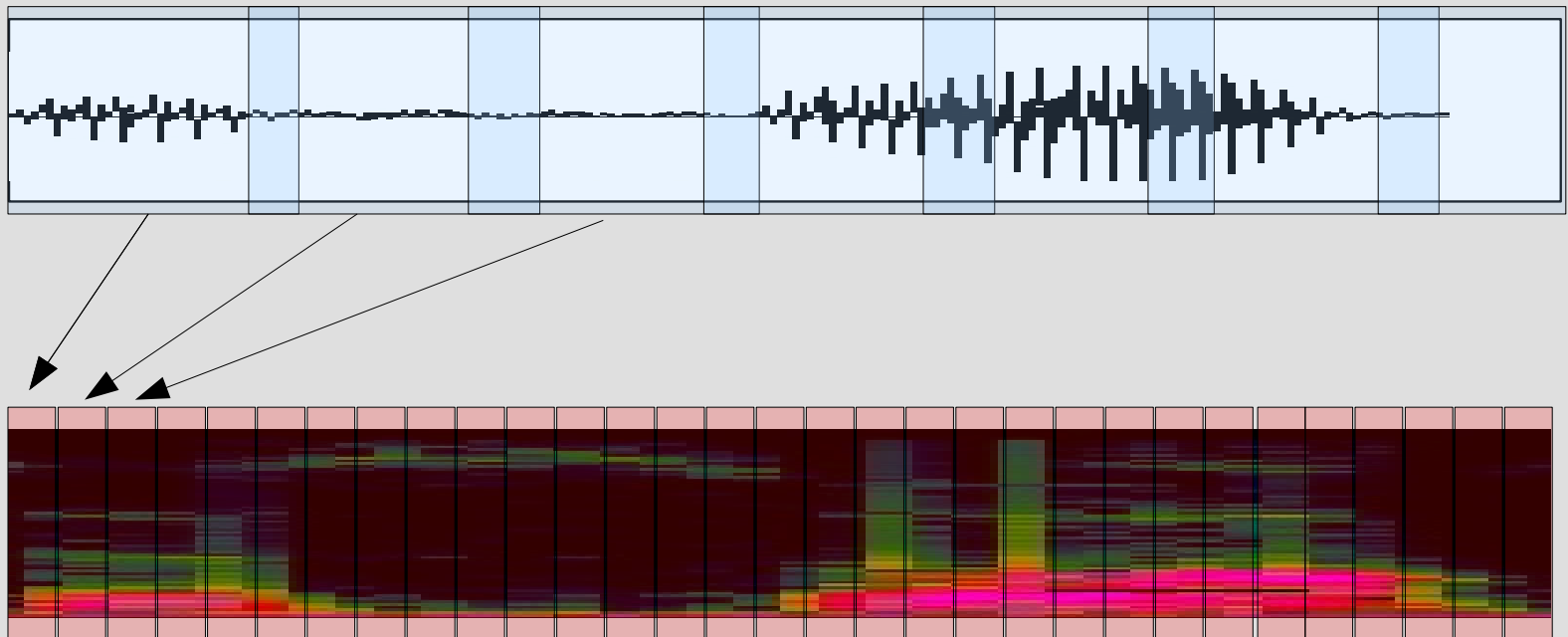
UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

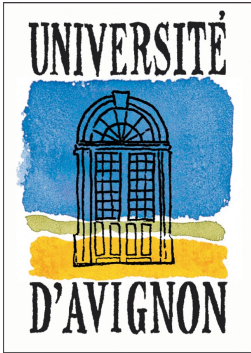
Statistical Speech modeling

- 4 key issues :
 - Features extraction
 - Spectral/cepstral models
 - pattern recognition problems
 - Modeling temporal structures
 - Modeling high level information
 - linguistic, semantic, pragmatic...

Statistical Speech modeling : feature extraction

Speech parametrization





UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : feature extraction

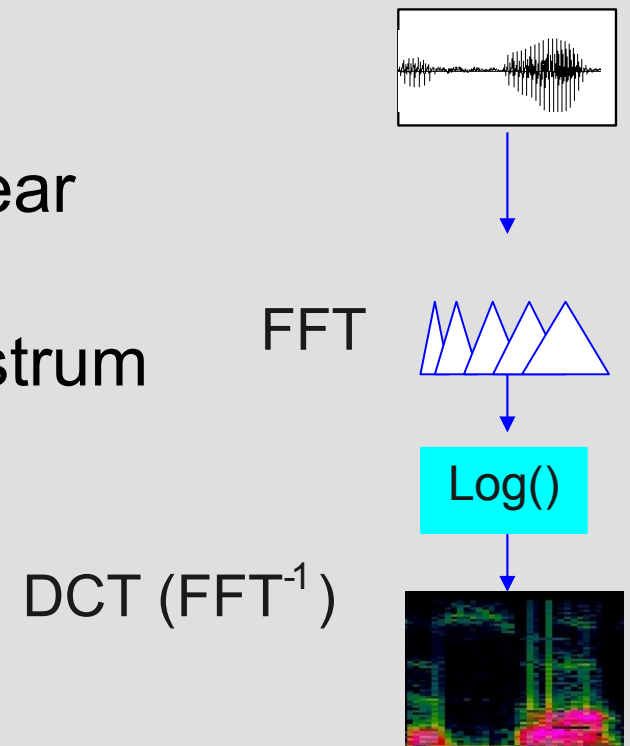
- Speech parametrization :
- LPC : Linear Predictive Coding
 - Principle :
 - to code the prediction errors
 - Auto-regressive models

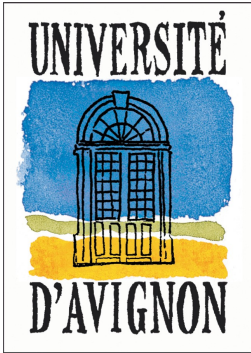
$$s(n) = \sum_{i=1}^P a_i s(n-i)$$

$$e(n) = x(n) - s(n) = x(n) - \sum_{i=1}^P a_i x(n-i)$$

Statistical Speech modeling : feature extraction

- Speech parametrization :
 - PLP : perceptually-based linear prediction
 - MFCC : Mel Frequency Cepstrum coefficients



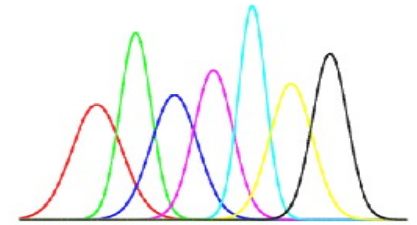


UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : feature extraction

- Speech parametrization : open issues
 - **Robustness**
 - Dimensionality reduction
 - Discriminative/generative approaches
 - LDA : Linear discriminant analysis
 - PCA : principal component analysis
 - ICA : independant component analysis
 - **Combination of audio features**
 - Complementarity of features

Statistical Speech modeling : modeling cepstral patterns

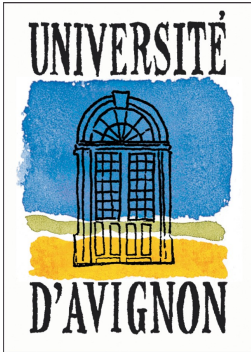


- Modeling cepstral features
 - Consensual approach : Gaussian Mixture Models
 - Principle : approximation of probability density function of cepstral patterns

$$l(x|\mu, \Sigma) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$$l(x|Gmm_k) = \sum_{i=0}^N w_k \cdot l(x|\mu_k, \Sigma_k)$$

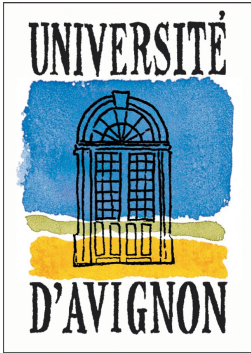
Parameters : $\lambda_k = (\mu_k, \Sigma_k)$



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : modeling cepstral patterns

- Gaussian Mixture Models
 - Allow us to estimate probabilities (likelihood) of observations knowing a model
 - Problems :
 - Estimation (training)
 - Integration to client systems



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : modeling cepstral patterns

- Training Gaussian Mixture Models
 - Principle
 - Criterion: Maximun Likelihood
 - Optimization algorithm :
 - Estimate of the model parameters maximizing the Likelihood
 - Training strategy : Expectation-Maximisation
 - E : estimate of the likelihood
 - M : updating de parameters to maximize Likelihood.

Statistical Speech modeling : modeling cepstral patterns

- EM :
 - Iterative process (until convergence)
 - Updating functions (N-component Gmm) :

$$l(x | Gmm) = \sum_{k=0}^N w_k \cdot l(x | \lambda_k)$$

Likelihoods

$$w'_k = \frac{\sum_X P(X | \lambda_k)}{\sum_X \sum_i P(X, | \lambda_i)}$$

Weights

$$P(X | Gmm_i) = \frac{l(X | Gmm_i)}{\sum_{k=0}^N l(X | Gmm_k)}$$

Probabilities

Statistical Speech modeling : modeling cepstral patterns

- Updating rules :

Means \longrightarrow

$$\mu_k^{t+1} = \frac{\sum_x P(x|\lambda_k^t) \cdot x}{\sum_x \sum_i^N P(x|\lambda_i^t)}$$

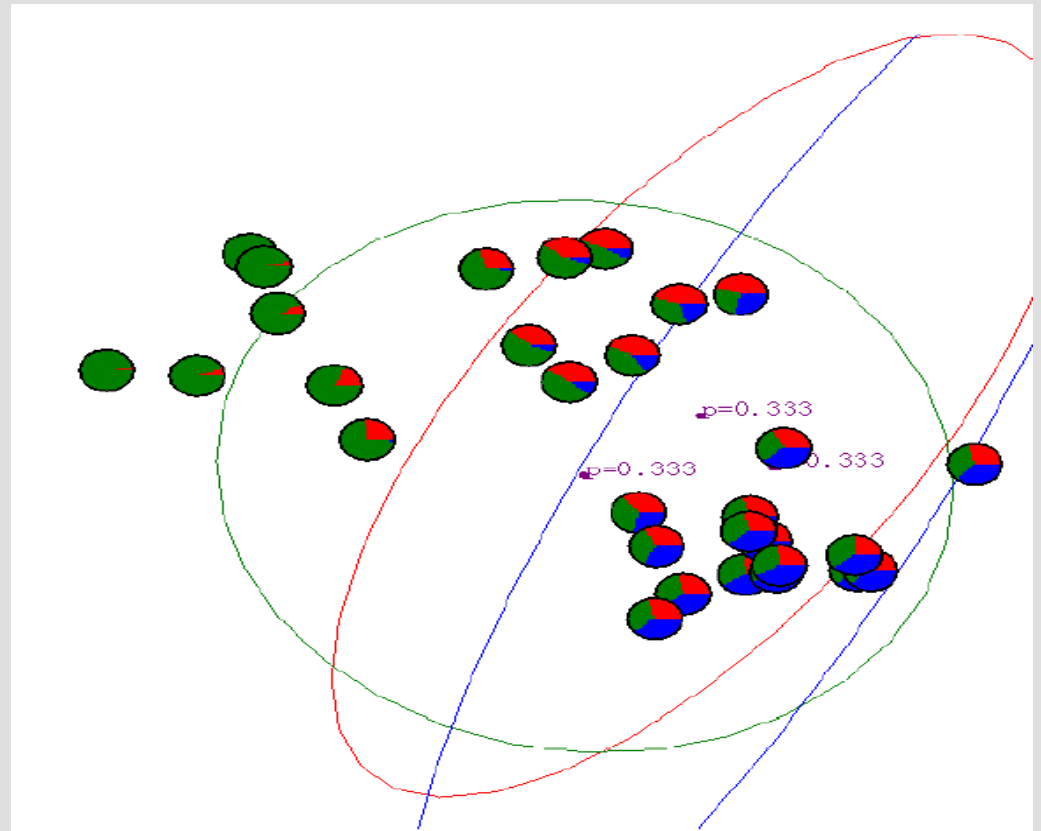
Variance \longrightarrow

$$\Sigma_k^{t+1} = \frac{\sum_x P(x|\lambda_k^t) [x - \mu_k^{t+1}] [x - \mu_k^{t+1}]^T}{\sum_x \sum_i^N P(x|\lambda_i^t)}$$

Statistical Speech modeling : modeling cepstral patterns

EM - Example

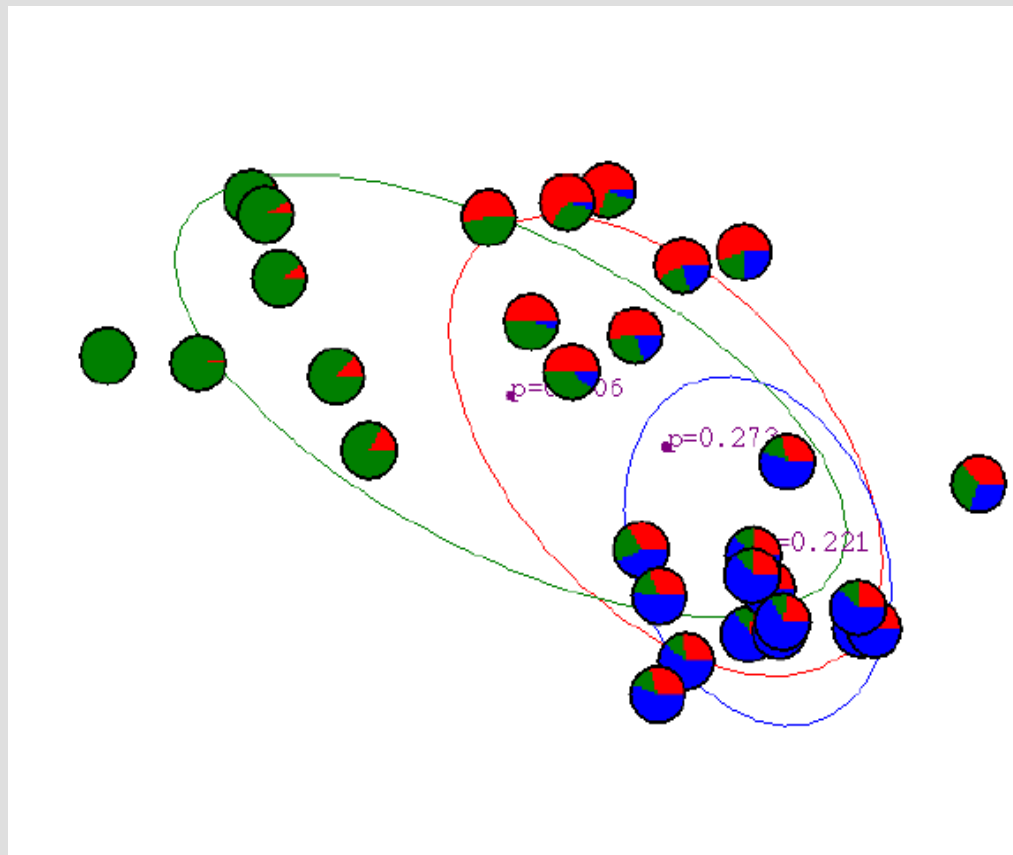
3 Gaussian
components



A.W. Moore

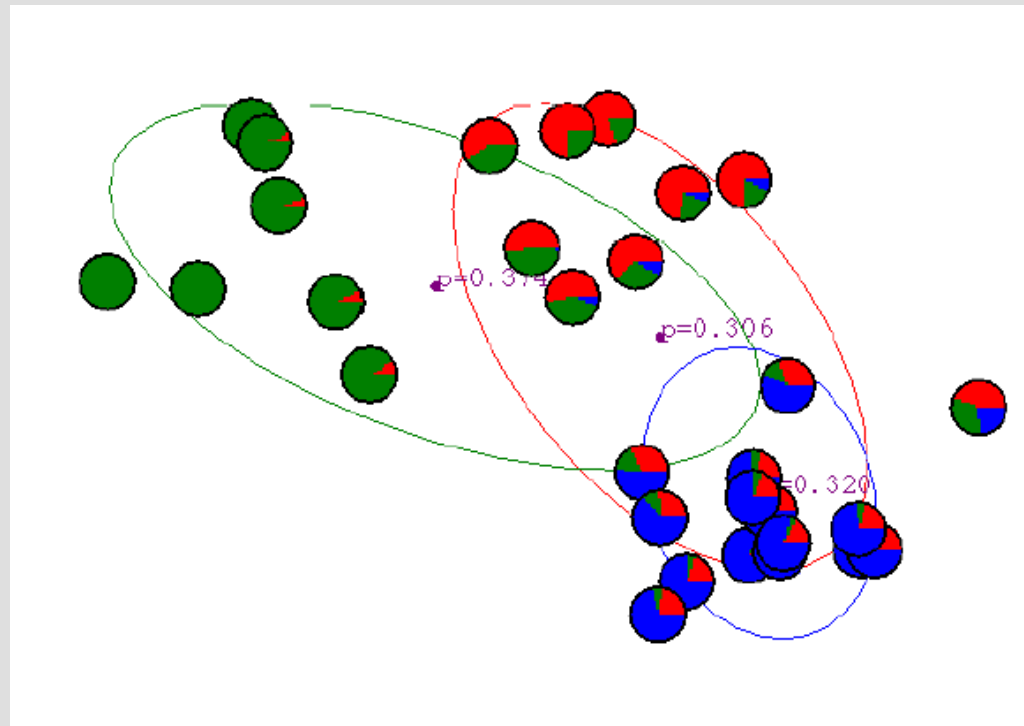
Statistical Speech modeling : modeling cepstral patterns

Iteration 1



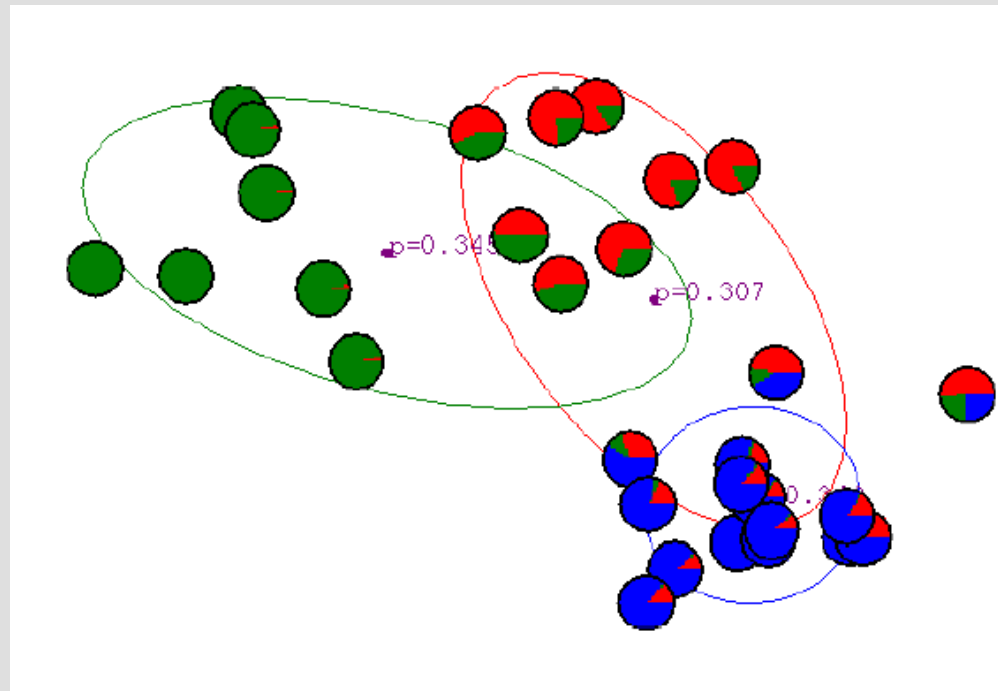
Statistical Speech modeling : modeling cepstral patterns

Iteration 2



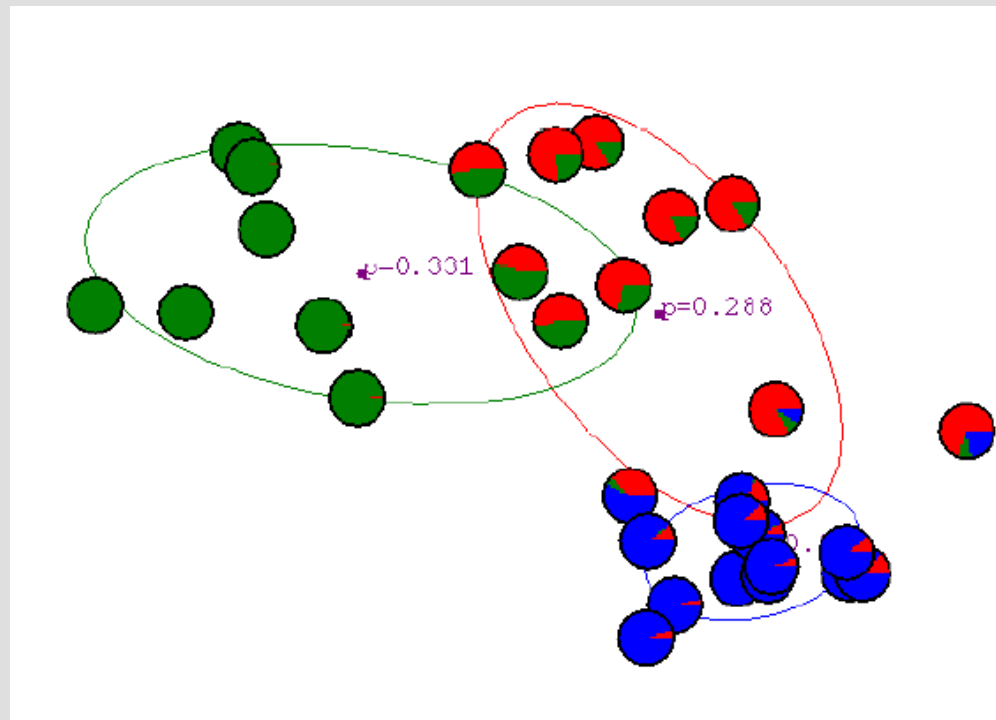
Statistical Speech modeling : modeling cepstral patterns

Iteration 3



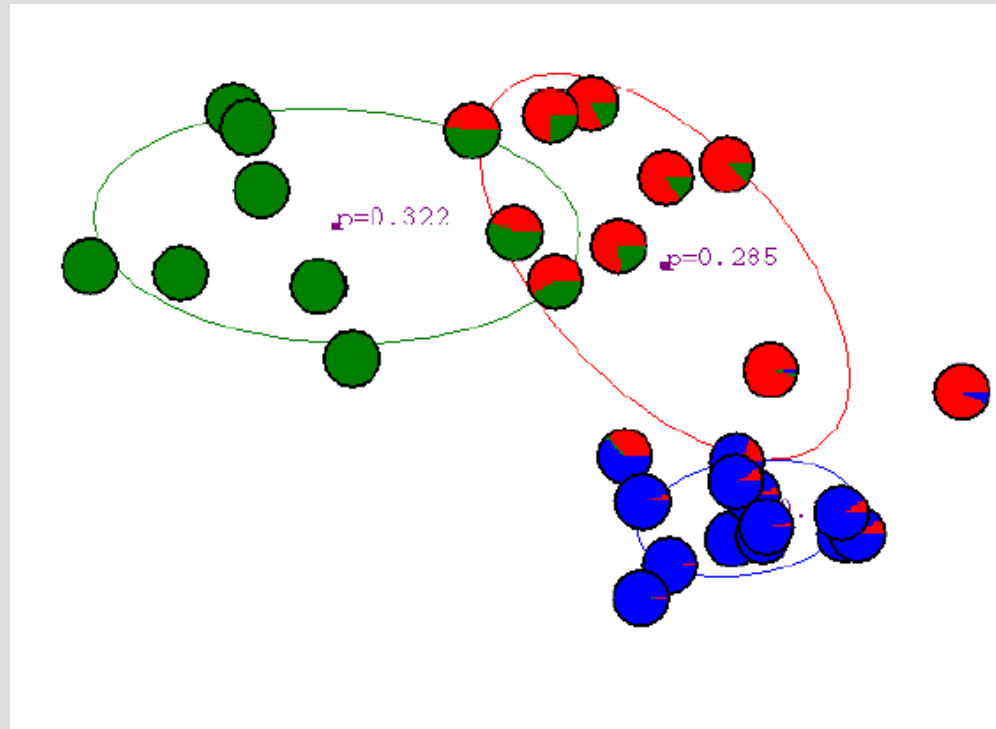
Statistical Speech modeling : modeling cepstral patterns

Iteration 4



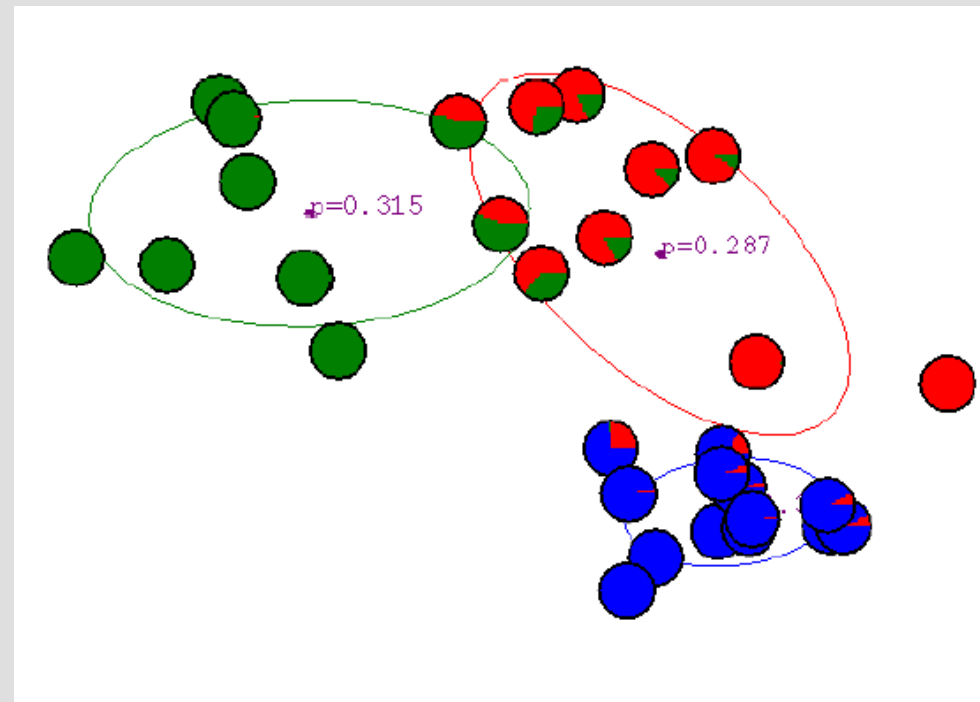
Statistical Speech modeling : modeling cepstral patterns

Iteration 5



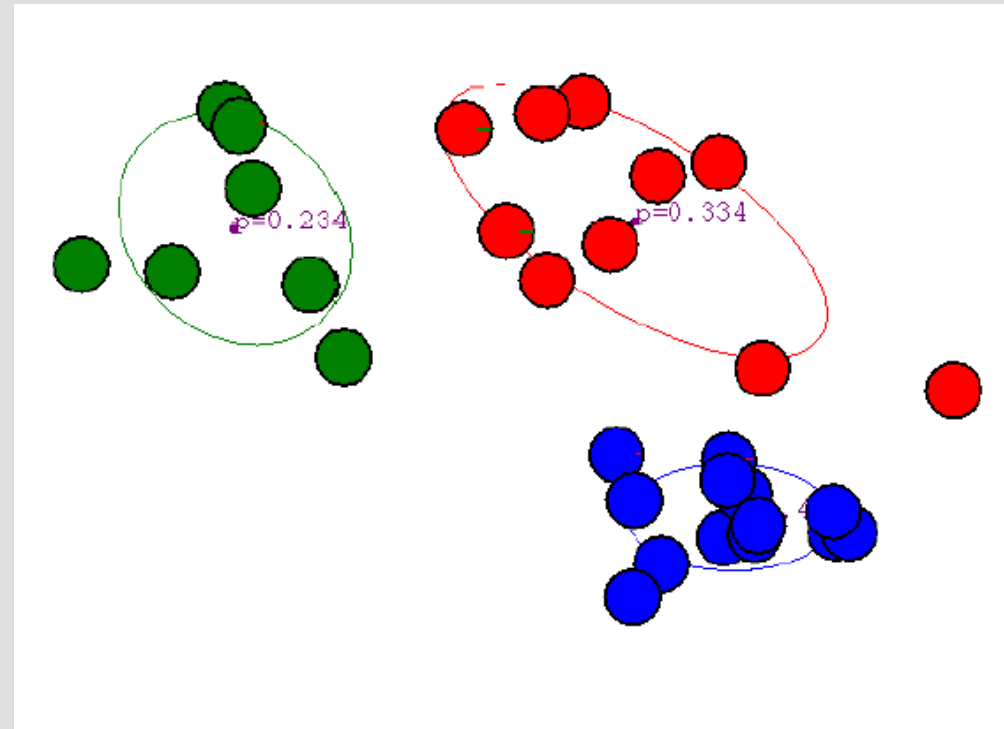
Statistical Speech modeling : modeling cepstral patterns

Iteration 6



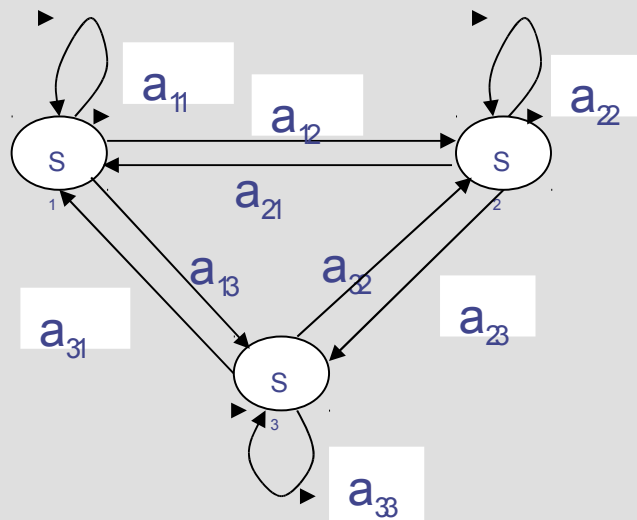
Statistical Speech modeling : modeling cepstral patterns

Iteration 20

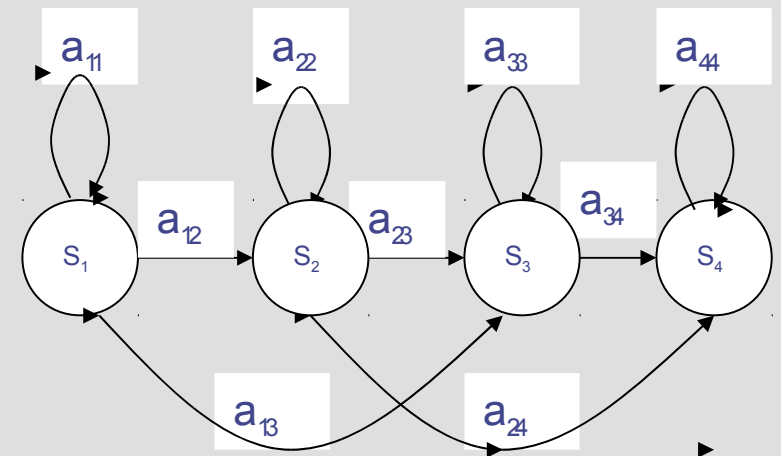


Statistical Speech modeling : modeling **temporal** patterns

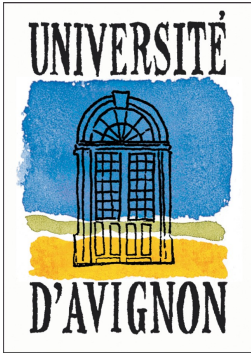
Hidden Markov Models (HMM)



Ergodic HMM



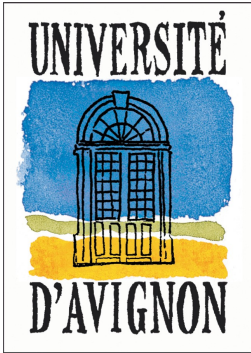
Left-right HMM



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : modeling **temporal** patterns

- Model topology:
 - State number
 - Links (transitions)
- Transition probabilities
- States :
 - density probability functions
 - GMM
 - Neural networks
 - others.....

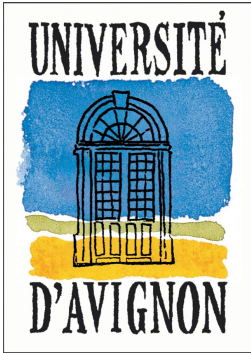


UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical Speech modeling : modeling **temporal** patterns

HMM is well defined tool to :

- estimate $P(O|\lambda)$, for an observation sequence O and a HMM λ ,
- search of the state sequence Q maximizing $P(Q|O, \lambda)$ (decoding stage, Viterbi Algorithm)
- Training stage : find the optimal λ (Baum-Welsh/EM algorithms)



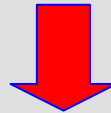
UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speech Processing **systems**

- Automatic Speech Recognition :
 - Extracting the linguistic content of audio/audiovisual documents
 - Many applications
 - Research since 1970

ASR : some applications

Rich transcription

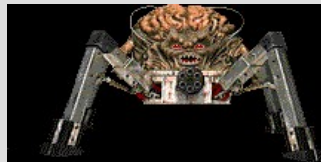
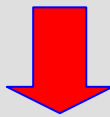


You are listening BBC news...

- Large vocabulary
- Continuous speech
- Dealing with speaking styles
- Extracting meta-data
 - Speaker, topics, etc..

ASR : some applications

Voice command



Voice command :

- small/medium vocabulary
- isolated words
- Embedded systems
 - Hardware constraints

ASR : some applications

Voice command



→ ROBUSTNESS MAY BE CRITICAL!

ASR : some applications

Audio search

Nicolas et Niak
Le Blog de Gringo (S'abonner gratuitement) - 19/04/2009 - Dailymotion
Niak au pré Vidéo envoyée par nicocentaure Nicolas est sans doute un **cheval** réincarné en humain...il va falloir enquêter sur cela ! Indéniablement, il parle "**cheval**". Qui n'a jamais essayé de demander quelque chose à son **cheval** lorsqu'il est au pré ? Et bien, pas si facile, souvent le **cheval**...

Mes premières vidéos en lignes
Bébé crapule (S'abonner gratuitement) - 15/04/2009 - YouTube
Tentatives de sourires : A **cheval** ; Mon premier œuf de Pâques ;

KENZA CHEVAL MERC 080409
treca (S'abonner gratuitement) - 09/04/2009 - Dailymotion
• KENZA au club d'Equit. - LE PETIT MONDE DE NATH

Le retour du père freine tard !!!!
C'est pas bien !!! (S'abonner gratuitement) - hier - Dailymotion
La saison RACE CAR SERIES 2009 est lancée, et ce fut une très belle première sur le circuit de Nogaro. Une voiture de 500 **chevaux**, des pilotes prêts à en découdre, des faits de courses permanents, cette discipline spectaculaire a définitivement sa place dans le sport auto français. La sensation du WE nous vient de François...

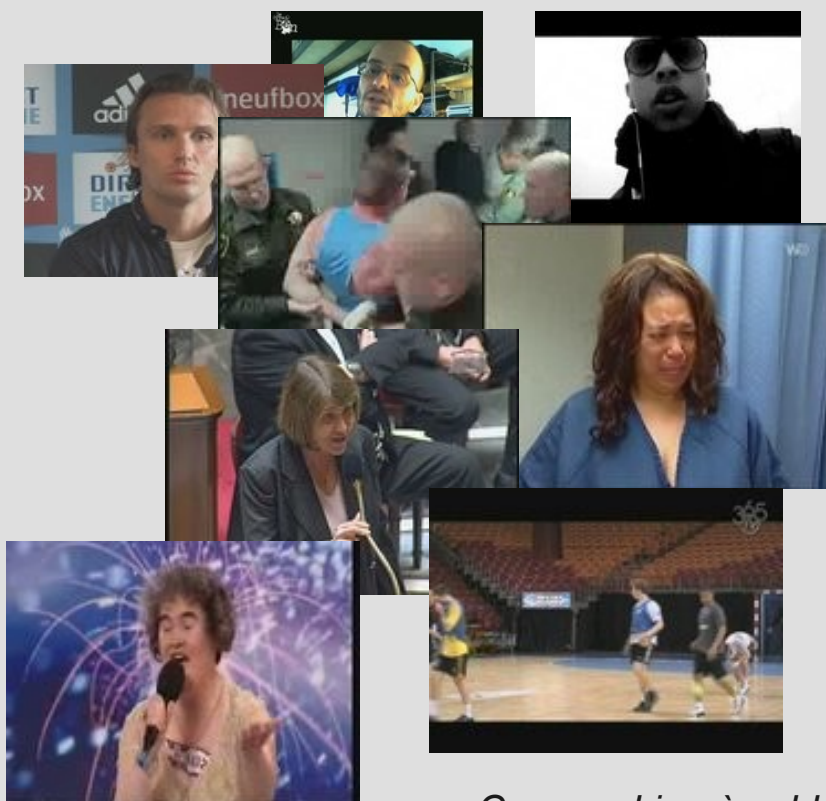
une histoire émouvante - Christian le lion -
nature-boy-79 (S'abonner gratuitement) - 18/04/2009 - Dailymotion
Vegan végétalien végétarien écologie décroissant veganisme choc choquant fourrures mode criminelle specisme race boucherie viande consommation alimentation mépris vie animale sentiment accident panthère lionnes lionne noire accident Afrique europe industrie chacal hygiène jaguar course poursuite police tighn année félin félins...

Antoine (et les autres) à la mer - Pâques 2009
Help !!! (S'abonner gratuitement) - 20/04/2009 - YouTube
Antoine sous toutes les coutures : à pieds, à **cheval**, en voiture, dans le sable, en tenue piscine et à moto ! Antoine et son nouveau camion Plic ploc

- Audio search
- ➔ Spoken term detection
- ➔ Topic detection
- ➔ Entities search ?

ASR : some applications

Structuring audiovisual databases



By-content structuring for efficient
archiving and access

- LVCSR
- Extra-linguistic contents
- Unexpected conditions

ASR : some applications

Example : video genre identification



Need of extra-linguistic features

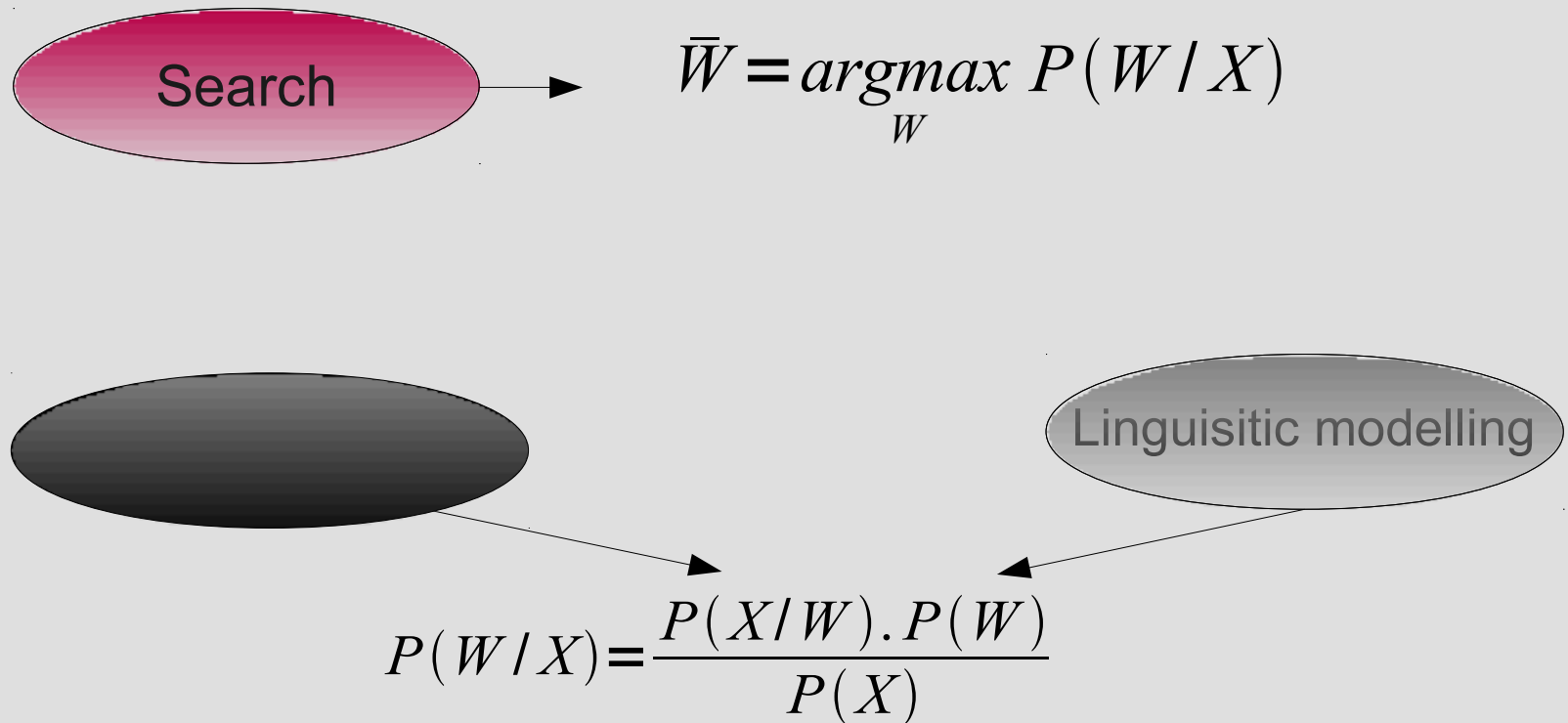
- Spontaneity
- Interactivity
- Emotional load ?

Genre

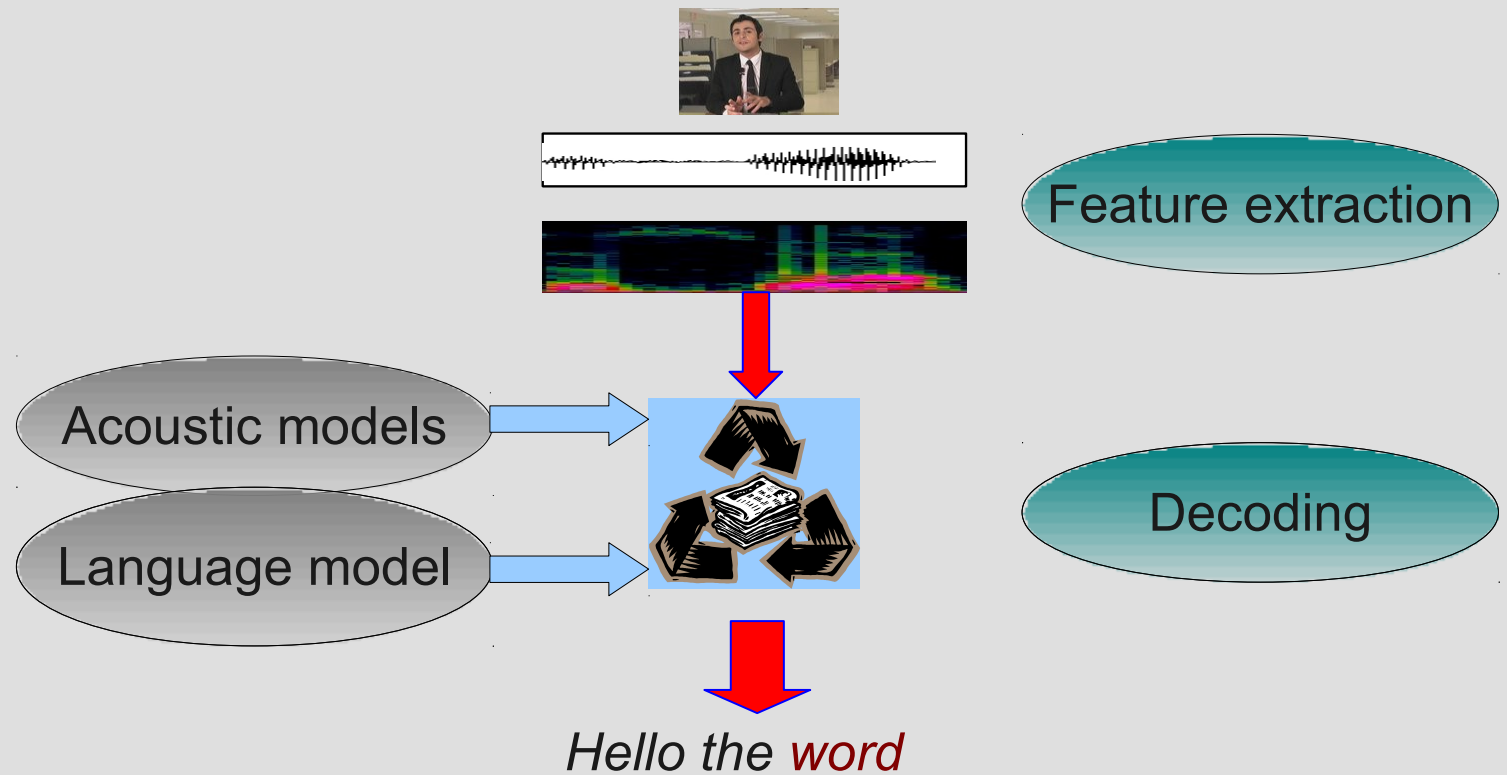
Topics

Environnement

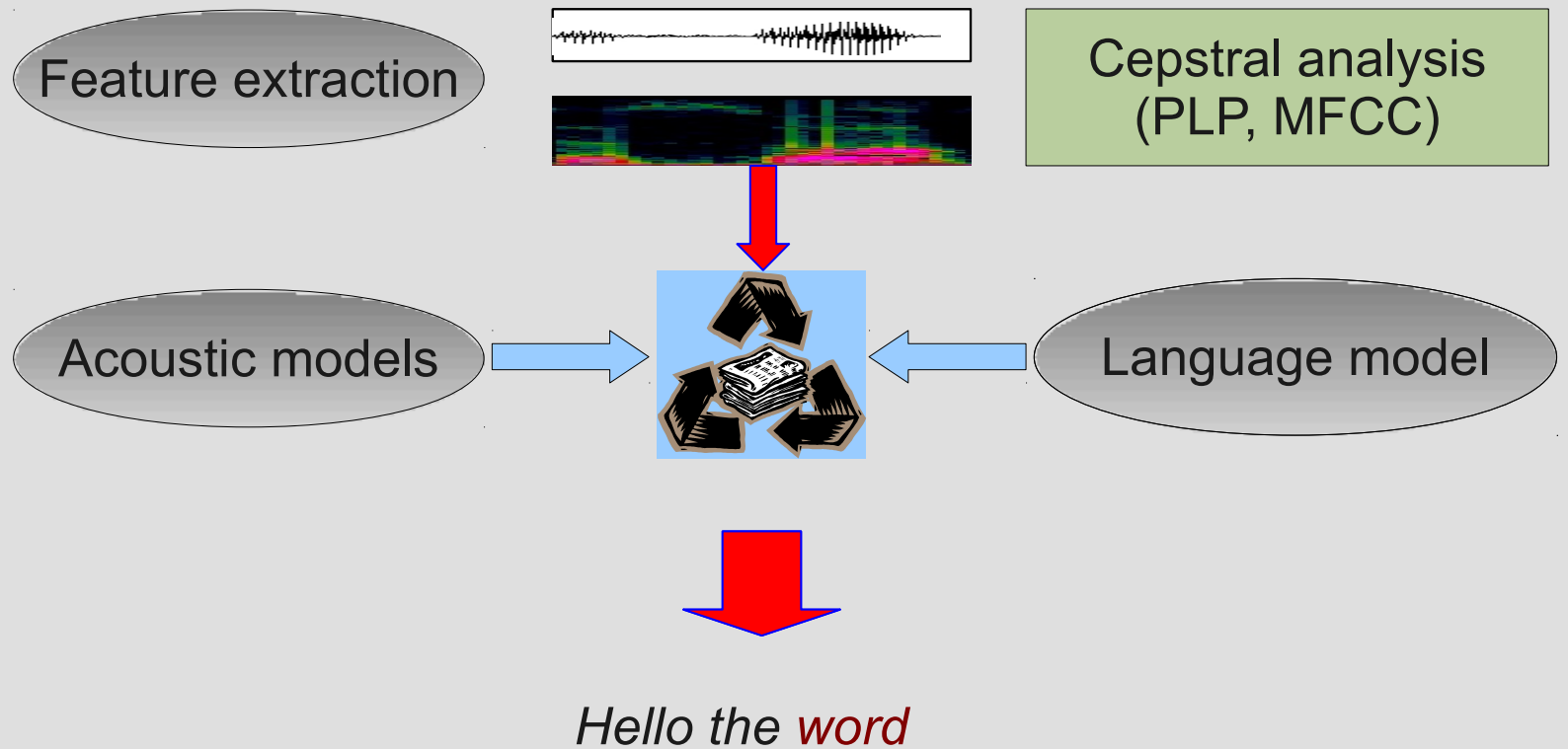
Fundamentals of statistical ASR



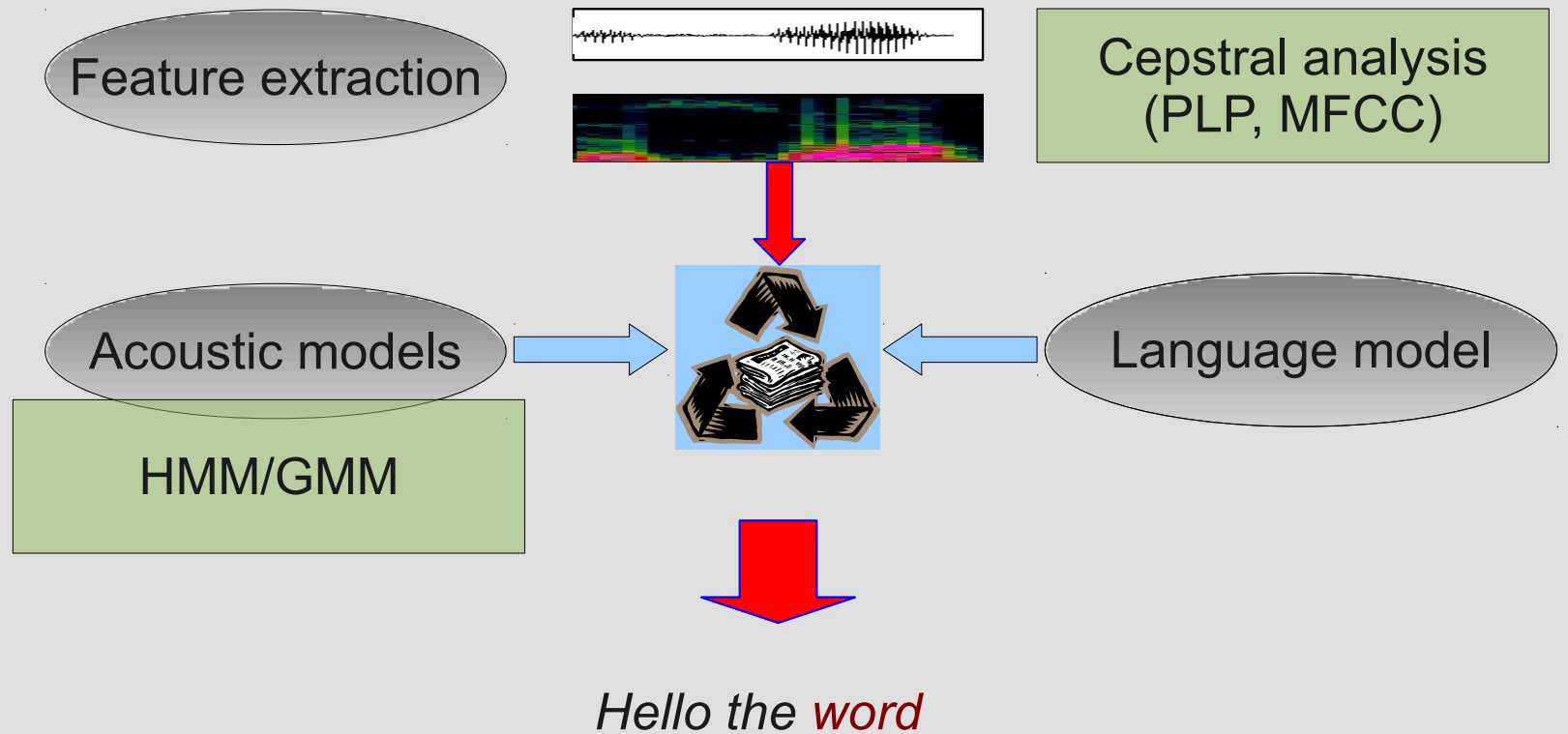
Fundamentals of statistical ASR



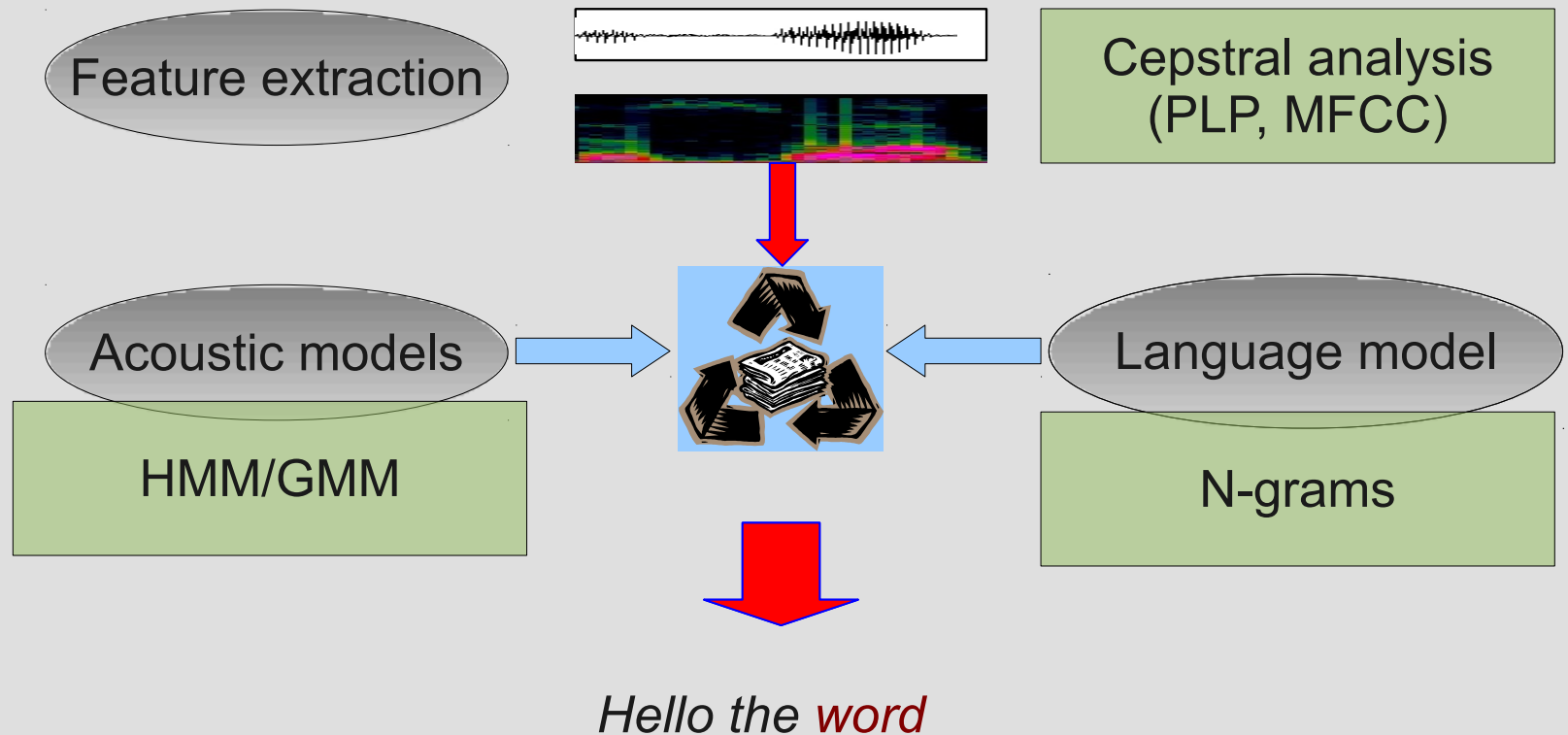
Fundamentals of statistical ASR

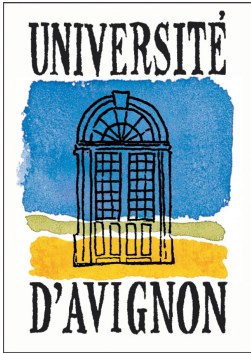


Fundamentals of statistical ASR



Fundamentals of statistical ASR

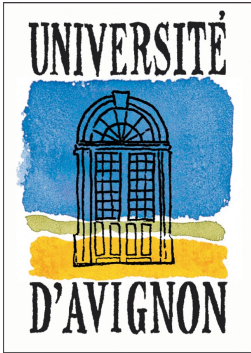




UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Fundamentals of statistical ASR

- Acoustic modeling :
 - Features in cepstral domain
 - Hidden Markov Models (HMM) for phoneme modeling
 - Gaussian Mixture Models (GMM)
 - Training on large annotated corpora
 - Generative models / discriminative learning :
MLE+MMIE/MPE



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Fundamentals of statistical ASR

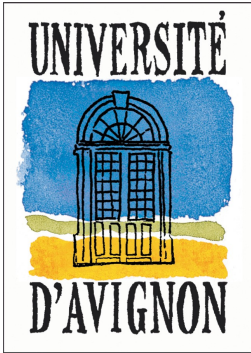
- Main issues in acoustic modelling :
 - Models estimate
 - Cost of the training corpora
 - Tuning the training algorithms
 - Feature/models combination
 - Robustness
 - to acoustic conditions
 - to speaker, speaking styles
(read/spontaneous/conversational speech)

Fundamentals of statistical ASR

- Language models
 - N-grams statistics :

$$P(W) = P(W_i / W_{i-1}, \dots, W_{i-n}) \dots P(W_{i-k} / W_{i-2}, \dots, W_{i-n-k})$$

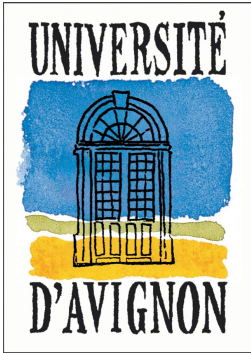
- N=3,4,5
- Require very large corpora (several million words)
- Modelling the unseen events ?
 - interpolation or back-off to (n-1)grams



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Fundamentals of statistical ASR

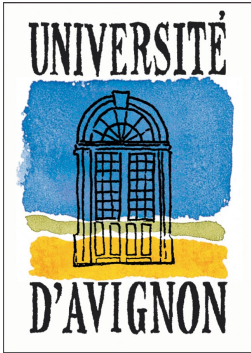
- Main issues in language models :
 - Exhaustive coverage of topics and speaking styles
 - Dealing with unseen events
 - OOV discovering and integration to LMs
 - Lack of semantics
 - Long-term dependencies
 - Semantic relationships



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Fundamentals of statistical ASR

- Search algorithm
 - 2 main approaches :
 - Beam search (Viterbi)
 - Depth-first search (A*)
 - Issues :
 - Dealing with hardware constraints
 - Embedded systems/Large scale ASR
 - Fast decoding

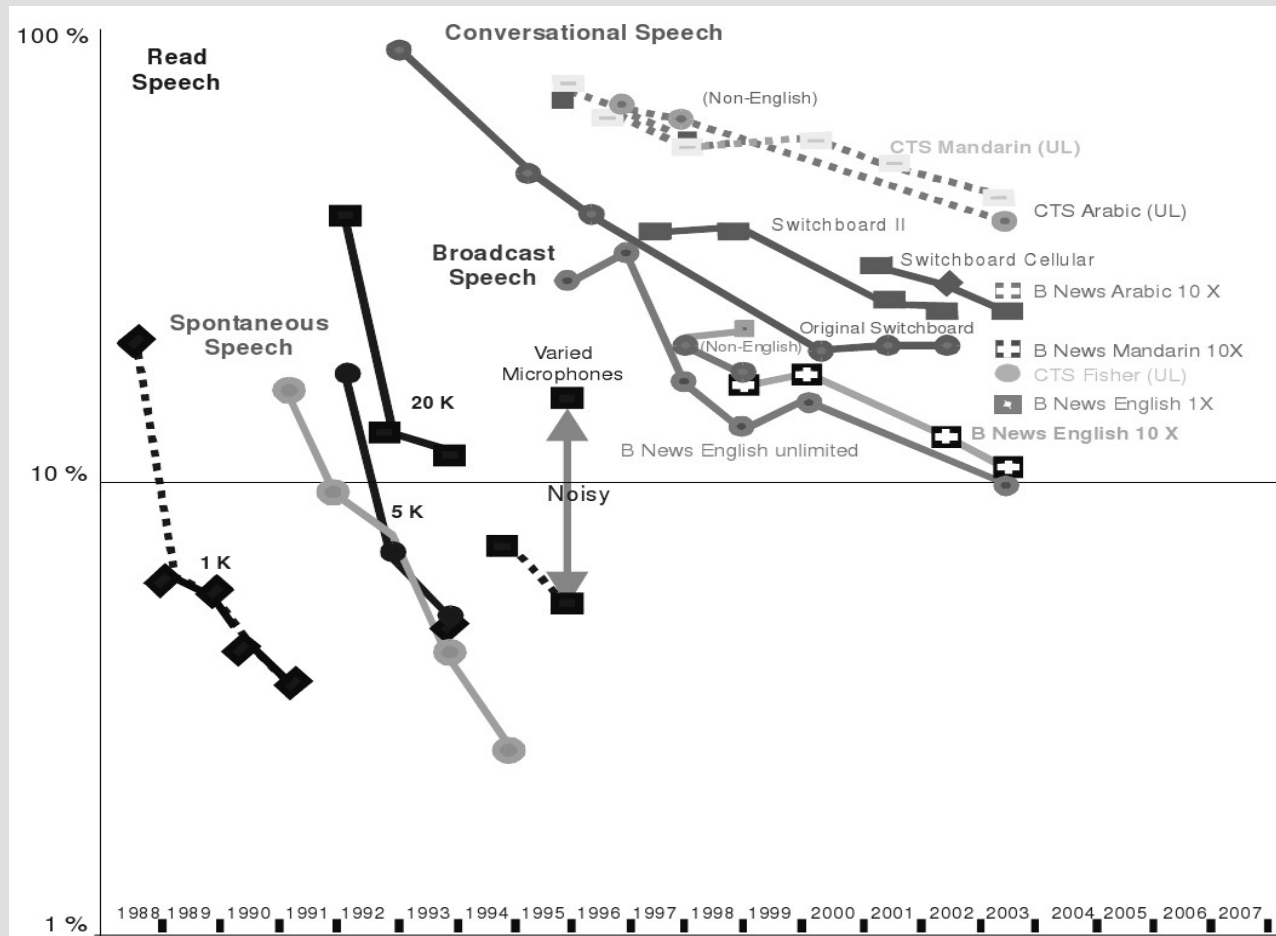


UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Fundamentals of statistical ASR

- Recent advances in Search
 - System combination :
 - A posteriori combination (ROVER)
 - Integrated approaches (Lecouteux & al, 2008)
 - Requires system complementarity, but similar accuracy (Bresdin & Gales, 2007)

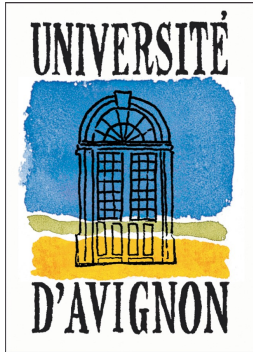
Where we are in LVCSR ?



What about the cost ?

Detail of Improvement	% WER
0. Baseline (RT-03 system)	13.4
1. 843-hour acoustic training	12.1
2. 1700-hour acoustic training	11.3
3. + MMI SAT PTM	11.2
4. + MMI SI PTM, SCTMs	11.0
5. + duration modeling	10.9
6. + online speaker clustering	10.8
7. + longer utterances	10.5
8. + new lexicon, LM	10.4

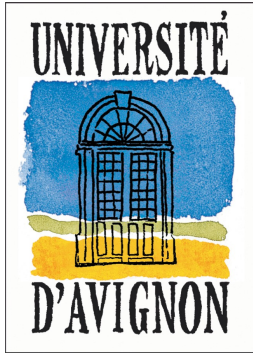
From n-Guyen & al, 2004



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

ASR: conclusion

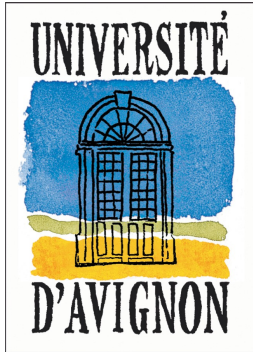
- Is it a success story?
 - Tradeoff between the expected gains and the costs is not so good
 - Are we on the limits of the HMM/N-gram framework ?
 - Is it the good/best paradigm ?
 - Can we do a better usage of ASR ?
 - What is the goal ?



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speech Analytics

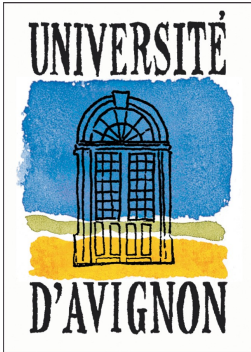
- Analyzing speech to extract high level information :
 - Topics
 - Opinions
 - Roles...
 - speech understanding/interpretation.



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speech Analytics

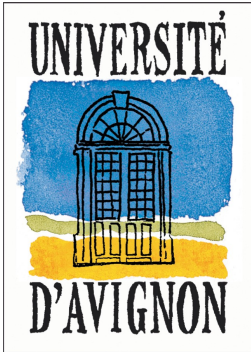
- Methods :
 - 2 steps :
 - (1) automatic transcription
 - (2) ASR outputs processing
 - Critical points :
 - Feature extraction from rich transcription
 - Classification tasks
 - Tools : SVM, Neural Nets, Boosting,...
 - Application-oriented speech processing



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speaker Recognition

- **Motivation :**
 - _ speech used to determine the true identity of the speaker
- **Speaker identification**
 - _ Who is speaking?
- **Speaker verification**
 - _ is the claimed identity true?
- **Speaker segmentation** : speaker turn detection, speaker tracking
- **Constraints :**
 - _ Open/close speaker set
 - _ Text dependent/independent

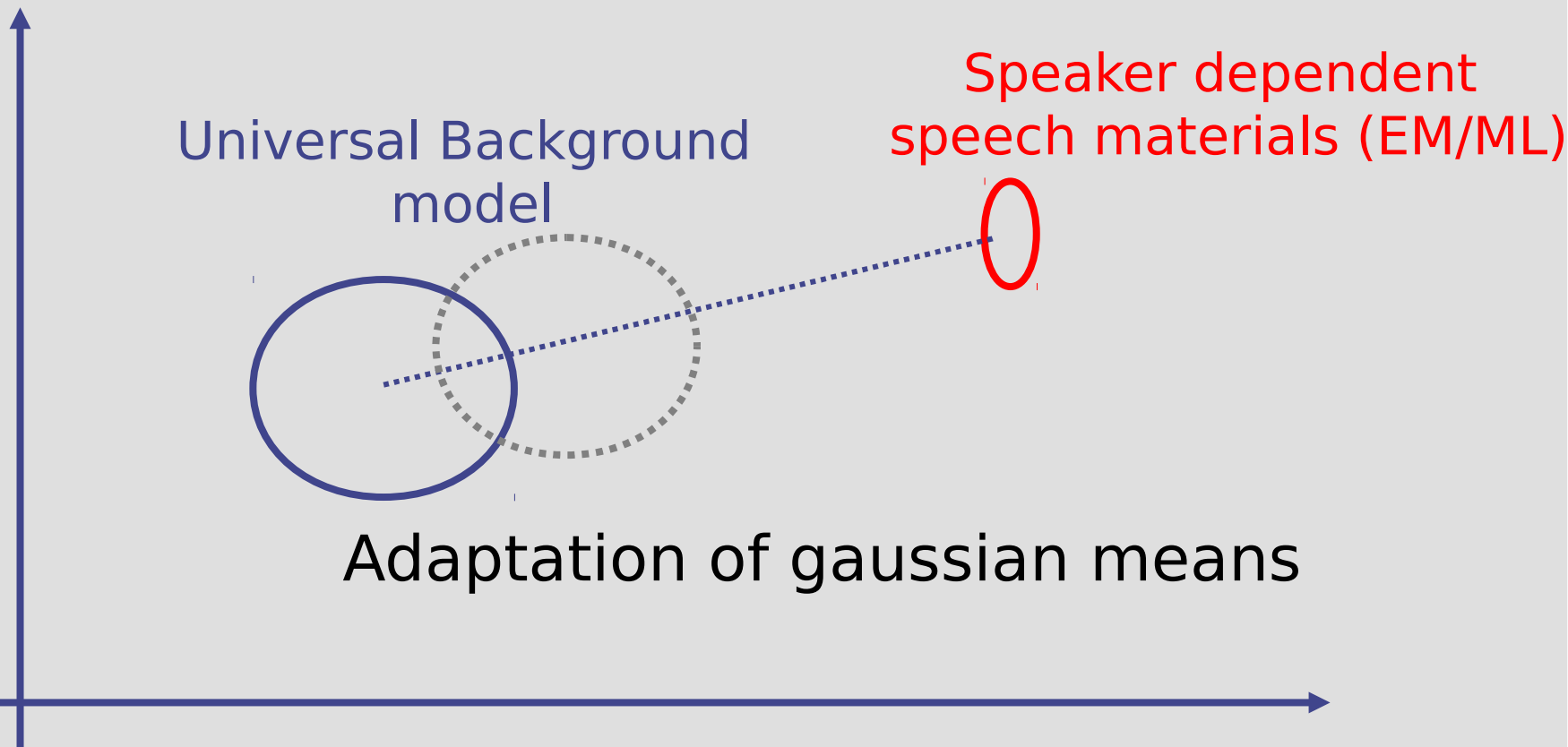


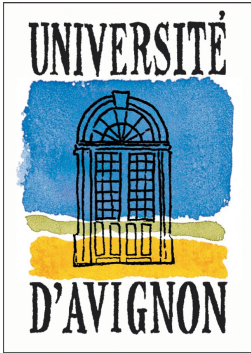
UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Speaker Recognition

- **General approach :**
 - Pattern recognition problem
 - Training phase : speaker-dependent models are enrolled on some speech samples
 - GMM-based approaches
 - Temporal structure of speech is not used

GMM based speaker identification

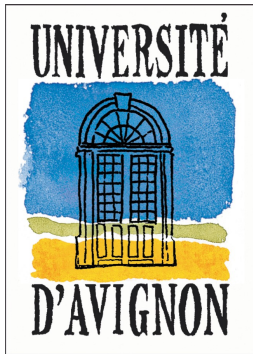




UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Performance of state-of-the art speaker identification systems

- NIST evaluation campaigns (speaker id)
 - American speakers, conversational speech
 - 5% error rate for the best system (2m30s)
- Identification on closed set
 - < 1% on studio data, 630 speakers (6s enrol., 3s test)
- **Error rates increase strongly on spontaneous speech and adverse conditions**
- New advances :
 - factor analysis for variability reduction
 - System combination



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

Statistical speech processing : conclusions

- Statistics is well defined framework to formulate speech processing problems
- ...and to build efficient systems
- Speech processing systems frequently rely on machine learning methods
- Research efforts mainly focused on the best way to apply mathematical tools to SP problems
- Many SP systems require other kind of information/other modeling paradigms to obtain significant improvements