



# Evaluation is not just a score

Olivier Galibert, Juliette Kahn, Ilya Oparin

Firstname.lastname@lne.fr

**MESURES  
& RÉFÉRENCES**

Clés de la COMPÉTITIVITÉ  
et d'un MONDE PLUS SÛR

Laboratoire national de métrologie et d'essais

## Modular

- ▶ Modules are evaluated independently
- ▶ Traditional and “simple” metrics
  - WER, TER, BLEU...

## Evaluation is about computing an overall score

- ▶ Little to no performance and error analysis
- ▶ Test data selected to get a global score
  - Application is the primary guide



## Evaluations that help improving systems

- ▶ Provide clues to system performance
  - Weak and strong points
  - Factors of variation
- ▶ Comprehensive evaluation report

## Evaluations that are not fully blackbox

- ▶ Inter-related modules
- ▶ Enhanced metrics and evaluation protocols
- ▶ Coherent data selection and annotation

## In-depth evaluations : explain and help



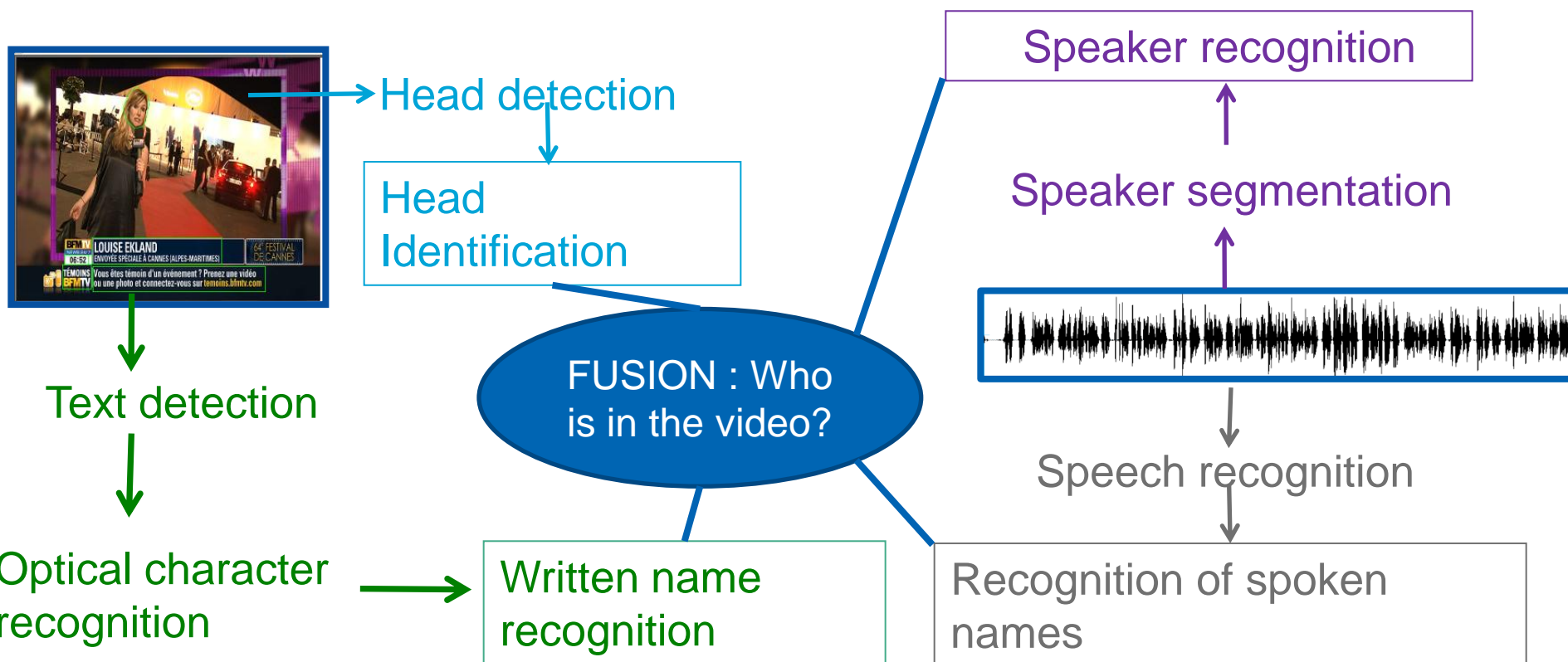
## How to build the perfect evaluation (ha!)

- ▶ Task definition and decomposition
- ▶ Data definition
- ▶ Metrics definition
- ▶ Results analysis



# Défi REPERE – Building a complex task

- ▶ Bring together several consortia to develop systems for person recognition in audiovisual data



## Corpus representative for the task

- Classify the difficulties of the corpus
- Composed from documents corresponding to applicative tasks

## Annotations

- Identify phenomena to be found
- Annotate potential explicative factors of system errors



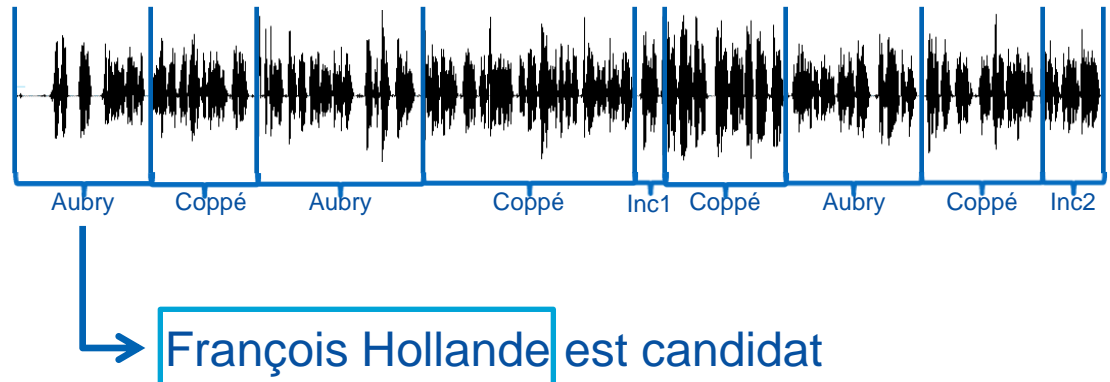
## Annotations for the tasks

- ▶ Speech transcription
- ▶ Annotation of persons names in the speech transcription, text
- ▶ Identification of the speakers and the heads

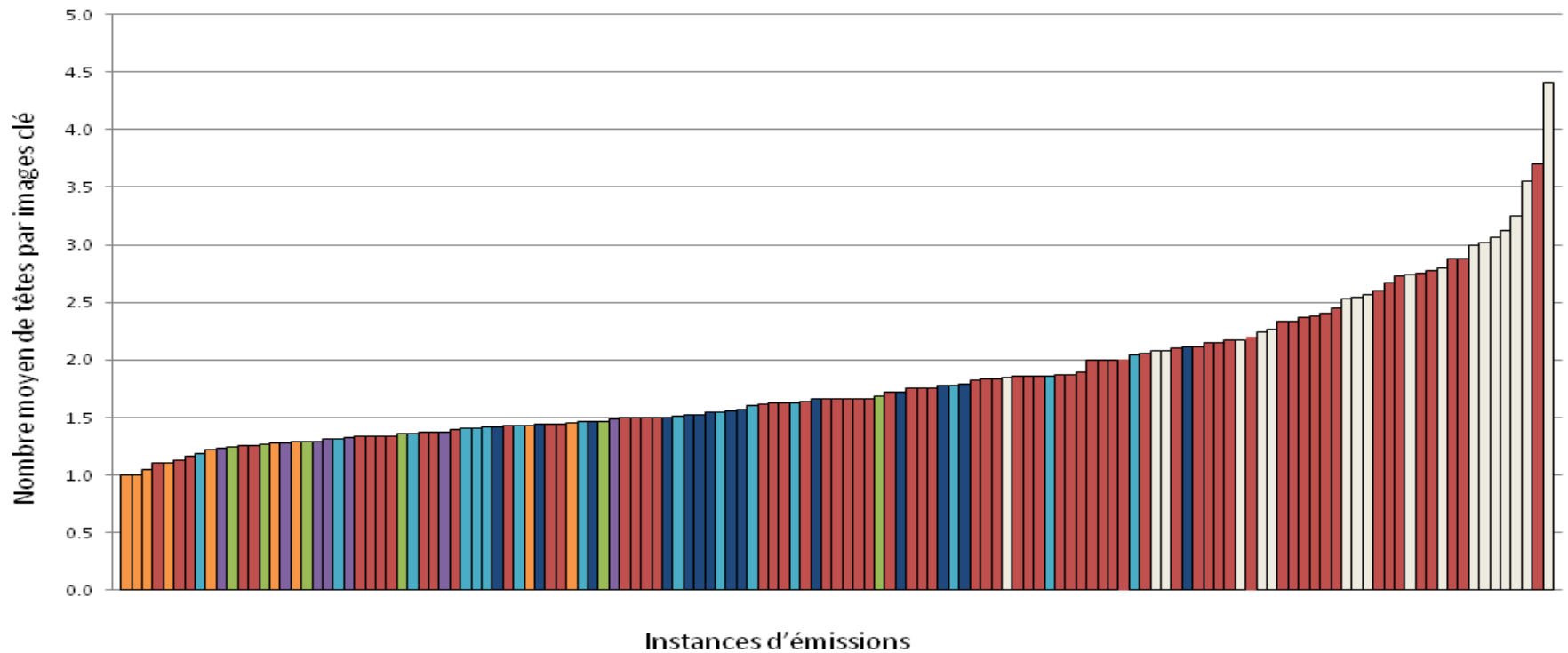


## Annotations for the analysis

- ▶ Gender
- ▶ Persons' roles
- ▶ Head descriptions



## Number of heads per annotated image

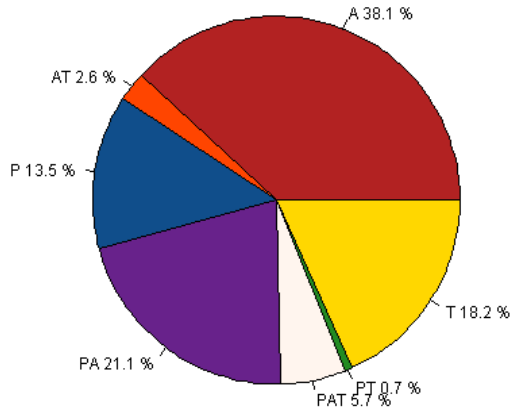


[Kahn J. et al, 2012]





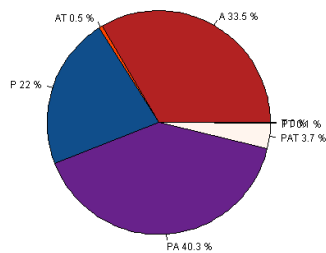
Répartition des modalités permettant de reconnaître les personnes



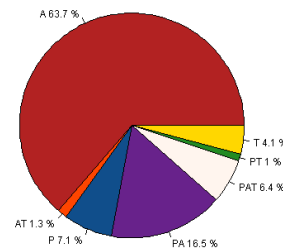
## Modalities that can be used

A : Seen  
 P : Spoken  
 T : His/her name shown on screen

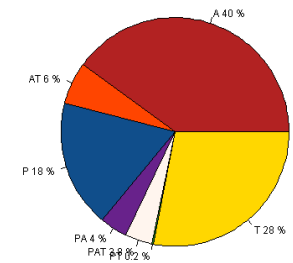
Répartition des modalités permettant de reconnaître les personnes Entre\_les\_lignes



Répartition des modalités permettant de reconnaître les personnes Top\_Questions



Répartition des modalités permettant de reconnaître les personnes Planete\_Showbiz

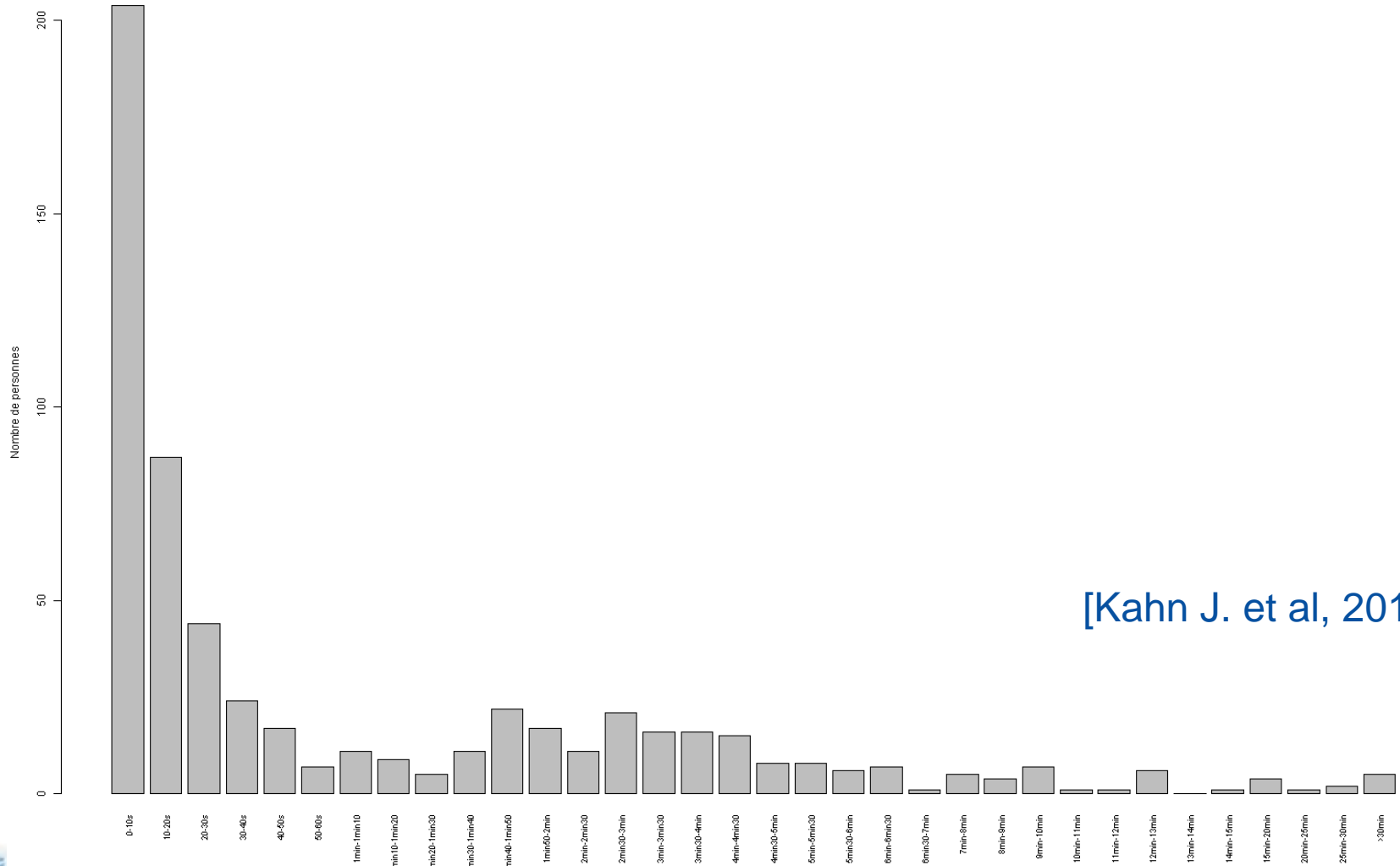


[Kahn J. et al, 2012]



# New data, known task or not ?

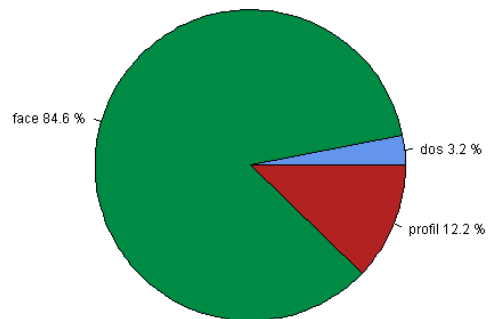
Répartition du temps de parole par personne indépendamment des émissions et des fichiers



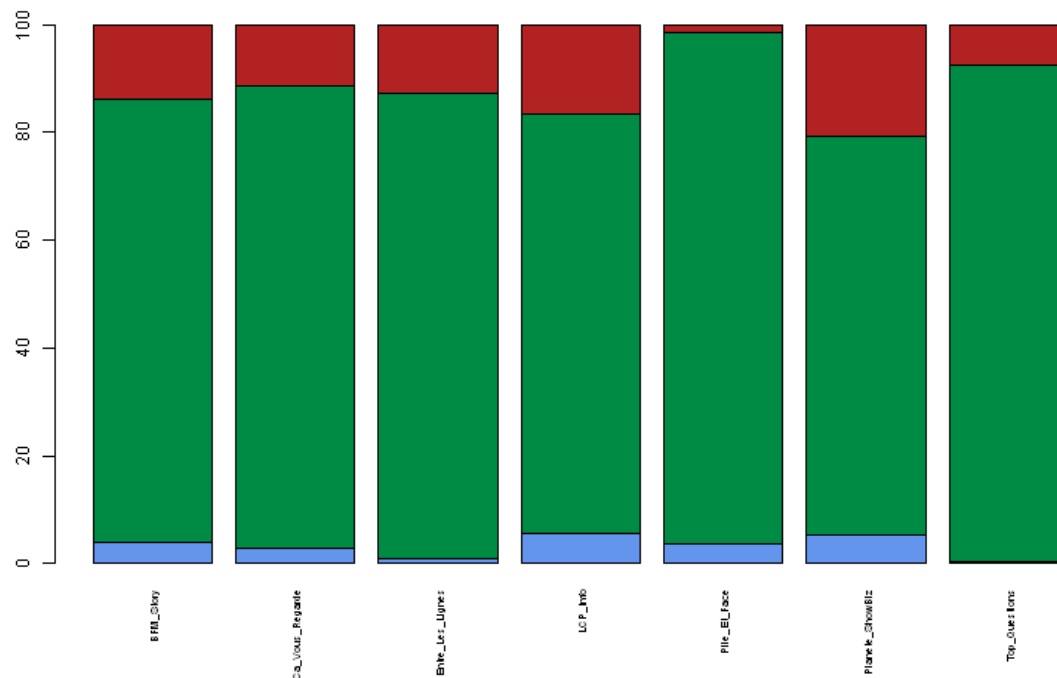
[Kahn J. et al, 2012]



Répartition des orientations des têtes dans le corpus d'apprentissage



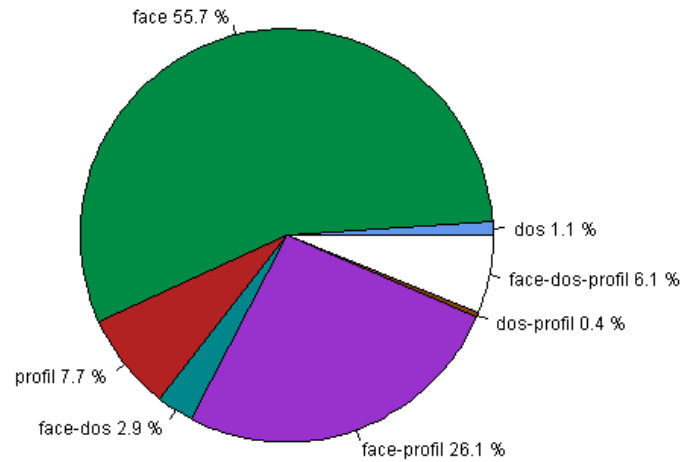
Head orientation : Face, profile, back



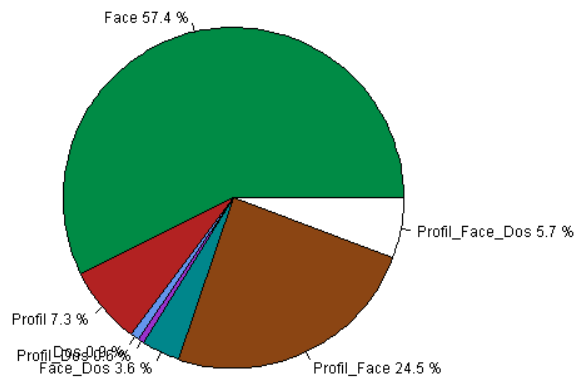
[Kahn J. et al, 2012]



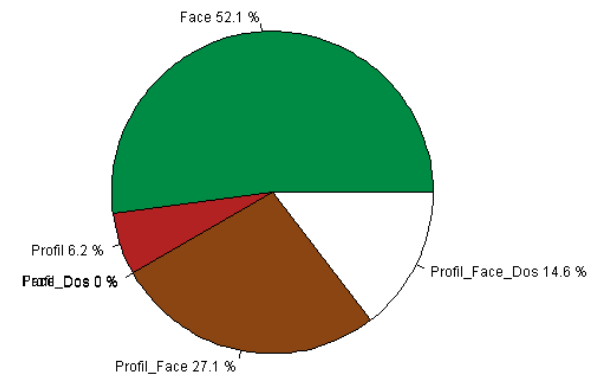
Répartition des personnes selon l'orientation des têtes



BFMStory



CaVousRegarde



[Kahn J. et al, 2012]

Marginal phenomenon except for the factor of glasses

Caractéristique binaire	NON	OUI	Pourcentage de OUI
Synthétique	13190	4	0.03%
Lunettes	9 674	3520	26.68%
Coiffes	12987	207	1.57%
Moustache	12981	213	1.61%
Barbe	12940	254	1.93%
Piercing	12434	760	5.76%
Autre	13165	29	0.22%

[Kahn J. et al, 2012]



- ▶ Some metrics compute a global opaque score, some try to enumerate errors
  - BLEU vs. TER
  - F-measure vs. SER
  - Jaccard vs. DetEval/ZoneMap
- ▶ Global scores are often correctness scores, which people often feel are more understandable
- ▶ Error-enumeration based metrics usually give error rates, which enable classification and importance tweaking
  - Classes are often a little large though

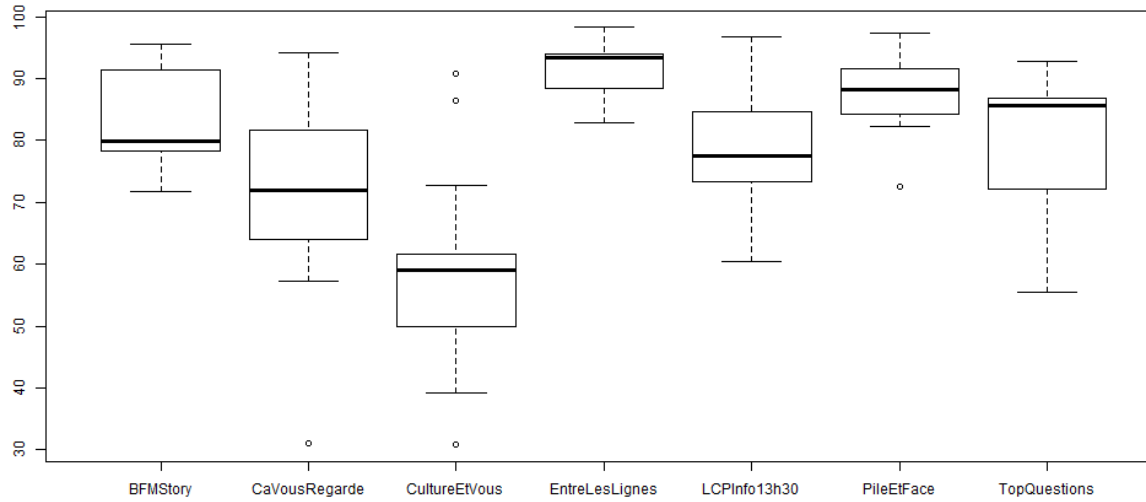


- ▶ Try to classify errors
  
- ▶ Try to extract explanation factors, discriminative factors between system outputs
  
- ▶ Representation is key... but also very hard
  - Information overload vs. oversimplification
  
- ▶ And data sparseness is quickly a killer
  - You don't do statistics on a couple of examples...
  - ...but building confidence intervals is an open problem

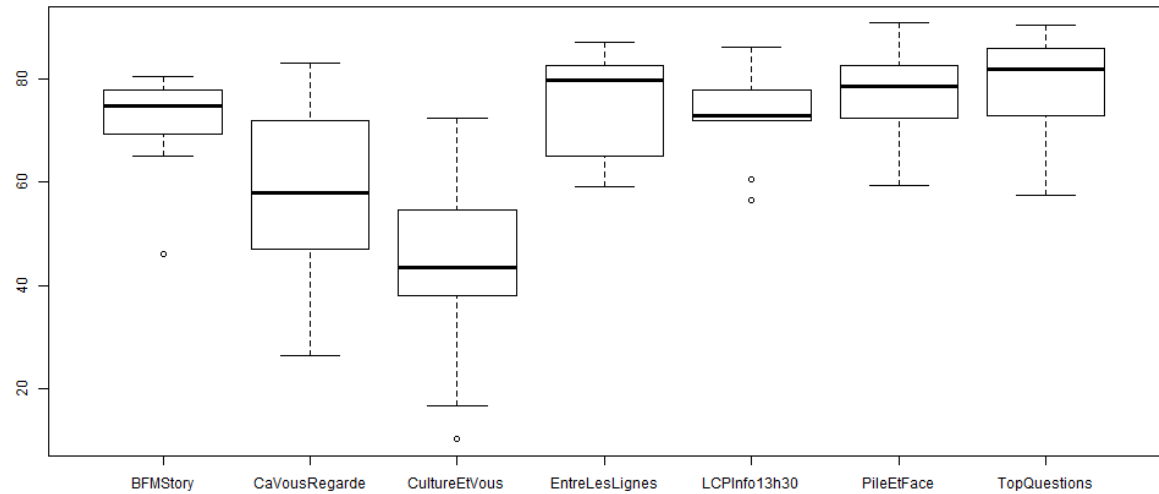


# Box plot vs...

Precision Who is visible or is speaking ?

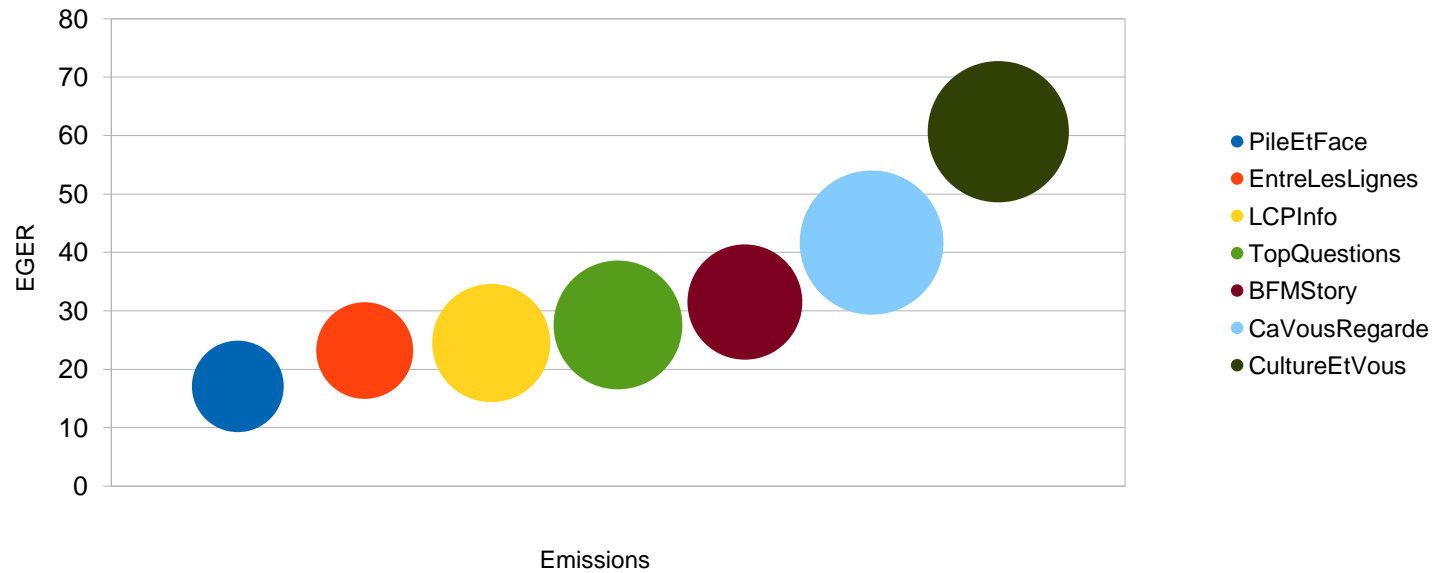


Recall Who is visible or is speaking ?



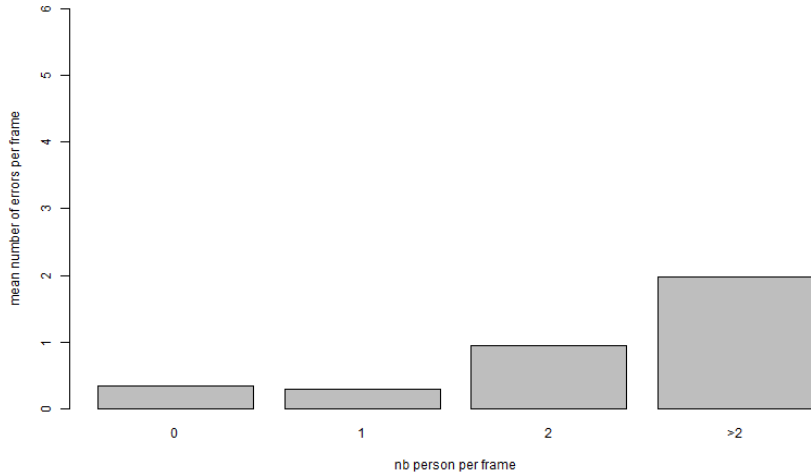


Qui est visible ou audible ?

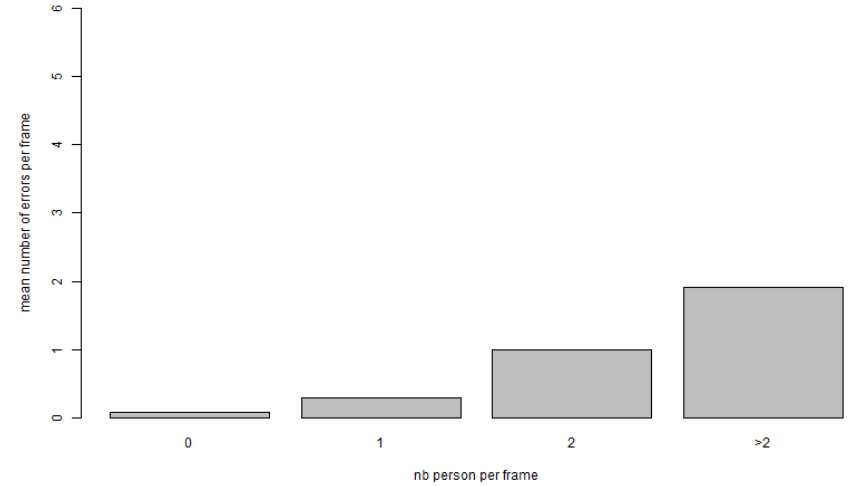


## Influence of the number of heads by frame

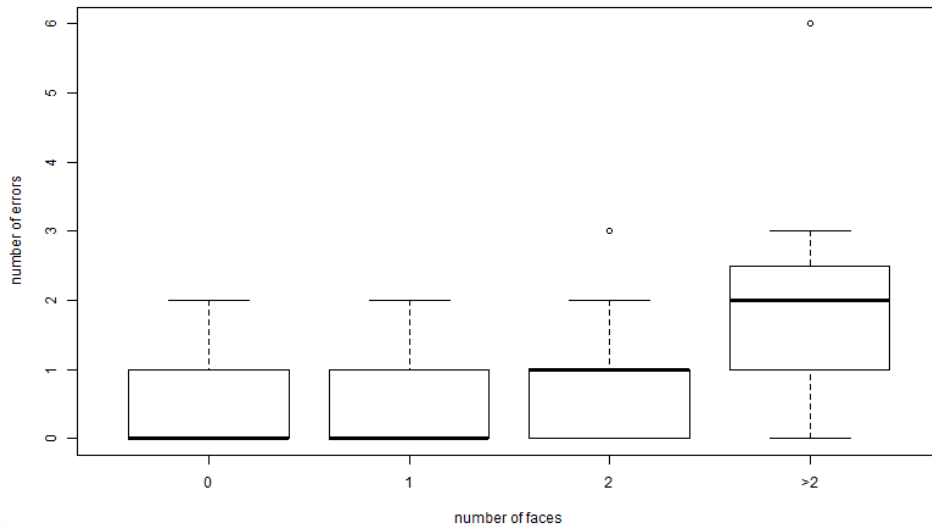
TPS\_PERCOL\_Primary\_Bouillabaisse



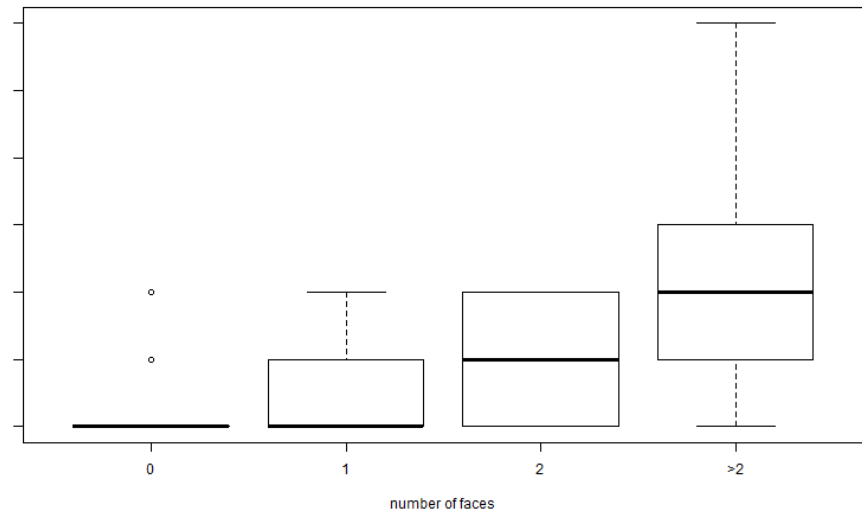
TPS\_SODA\_Primary



TPS\_PERCOL\_Primary\_Bouillabaisse



TPS\_SODA\_Primary



- ▶ Corpus definition
  - Factor-based structure vs. usefulness ratio
  - Annotation for analysis... if there is enough data
  
- ▶ Metrics definition
  - Metrics that enumerate errors
  - Tweakability
  
- ▶ Evaluations
  - Factor analysis and sparsity
  - Pertinence of representations





Thanks for your  
attention

[Firstname.lastname@lne.fr](mailto:Firstname.lastname@lne.fr)



**MESURES  
& RÉFÉRENCES**

Clés de la COMPÉTITIVITÉ  
et d'un MONDE PLUS SÛR

Laboratoire national de métrologie et d'essais