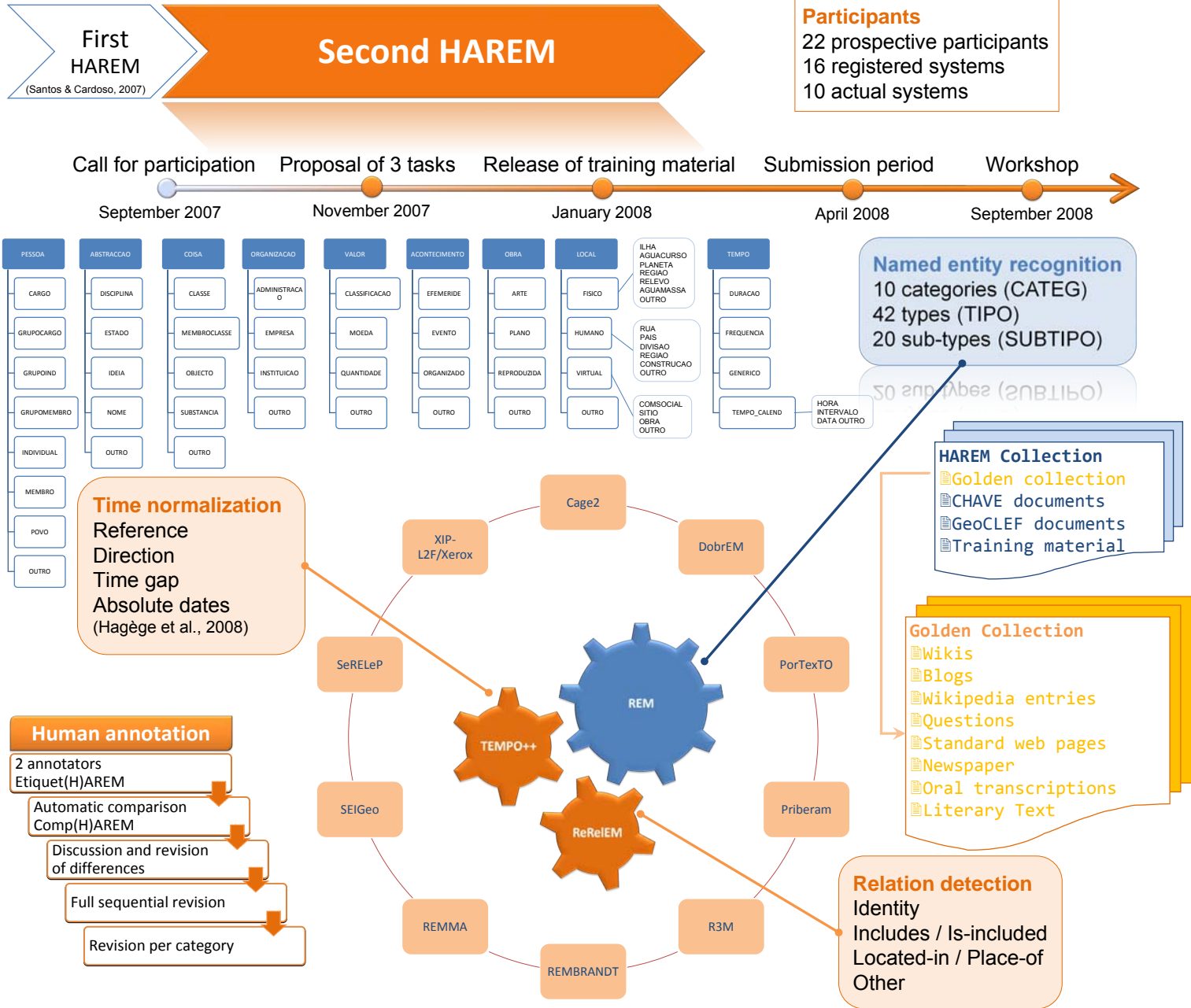


# Second HAREM: new challenges and old wisdom

<http://www.linguateca.pt/HAREM>

Diana Santos [Diana.Santos@sintef.no](mailto:Diana.Santos@sintef.no) Cláudia Freitas [maclaudia.freitas@gmail.com](mailto:maclaudia.freitas@gmail.com)  
 Hugo Gonçalo Oliveirahroliv@dei.uc.pt Paula Carvalho [pqfcarvalho@gmail.com](mailto:pqfcarvalho@gmail.com) \*\*



## Improvements

- ✓ No separation between identification and classification
- ✓ Finer geographic distinctions
- ✓ Finer grained classification of time expressions (Hagège et al., 2008)
- ✓ Not knowing <EM> ≠ Knowing it is not <EM CATEG="OUTRO">
- ✓ Text collections in XML
- ✓ Open source scoring programs
- ✓ Both GC and submissions may have classification (I) and delimitation vagueness (ALT)
- ✓ Partial identification doesn't score unless covered by ALT
- ✓ New evaluation measure

Hagège, Caroline, Jorge Baptista & Nuno Mamede. "Proposta de anotação e normalização de expressões temporais da categoria TEMPO para o HAREM II". 2008.

Santos, Diana & Nuno Cardoso (eds.). *Reconhecimento de entidades mencionadas em português: Documentação e actas do HAREM, a primeira avaliação conjunta na área*. 2007.

## "Old", persistent features

- ✓ Semantic model
- ✓ Selective scenarios

$$\sum_{n=1}^N (1 - 1/num_{cat}) * cat_{certo} * \alpha + (1 - 1/num_{tipos}) * tipo_{certo} * \beta + (1 - 1/num_{sub}) * sub_{certo} * \gamma - \sum_{n=0}^M (1/num_{cat}) * cat_{certo} * \delta + (1/num_{tipos}) * tipo_{certo} * \epsilon + (1/num_{sub}) * sub_{certo} * \phi$$