# A Proposal to Extend and Enrich the Scientific Data Curation of Evaluation Campaigns

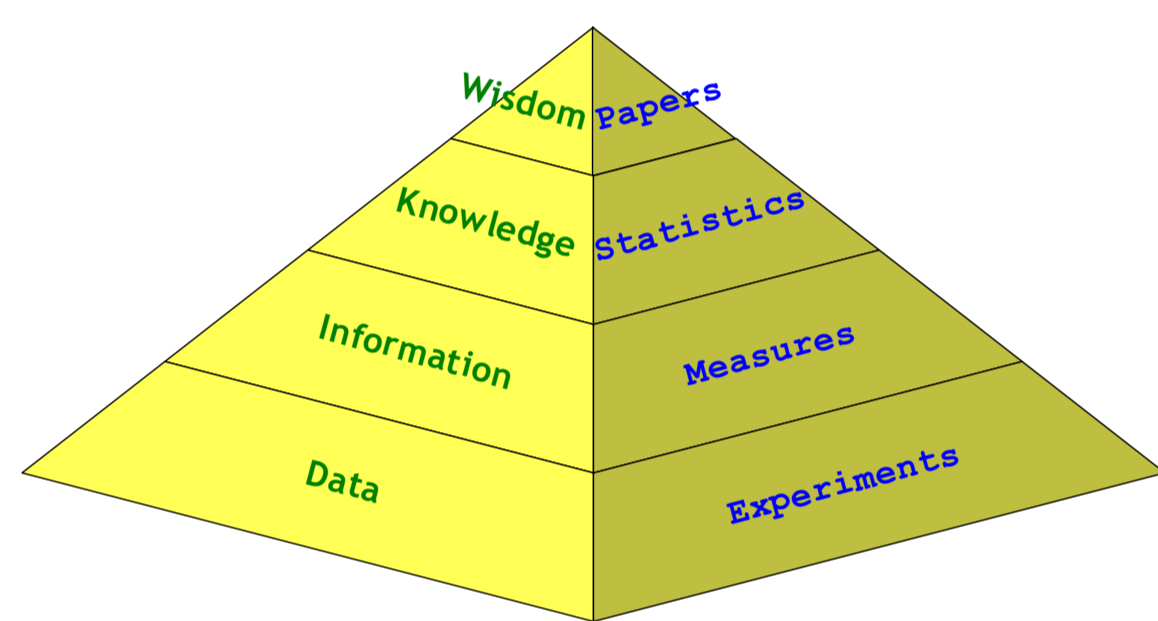## Maristella Agosti, Giorgio M. Di Nunzio, and Nicola Ferro

Information Management Systems (IMS) Research Group
Department of Information Engineering – University of Padova
{agosti, dinunzio, ferro}@dei.unipd.it

## Methodology

The current approach for laboratory evaluation of information access systems relies on the Cranfield methodology, which makes use of *experimental collections*.

▶ An experimental collection $\mathcal{C}$ allows the comparison of information access systems according to some measurements which quantify their performances;

▶ If we reasoning about this evaluation paradigm, a first step is to point out that the experimental evaluation in the *Information Retrieval (IR)* field is a scientific activity and, as such, its outcomes are different kinds of valuable *scientific data*

▶ Using the experimental data, we produce different performance measurements, such as precision and recall, that are standard measures that are used to evaluate the performances of an *Information Retrieval System (IRS)* for a given experiment. Starting from these performance measurements, we can compute descriptive statistics, such as mean or median, used to summarize the overall performances achieved by an experiment or by a collection of experiments. Finally, we can perform hypothesis tests and other statistical analyses to conduct an in-depth analysis and comparison over a set of experiments.



▶ *data*: the *experimental collections* and the *experiments* correspond to the "data level" in the hierarchy, since they are the raw, basic elements needed for any further investigation and they would have little meaning by themselves. In fact, an experiment and the list of results obtained conducting it are almost useless without a relationship with the experimental collection with respect to which the experiment has been conducted and the list of results produced; those data constitute the basis for any subsequent computation;

▶ *information*: the *performance measurements* correspond to the "information level" in the hierarchy, since they are the result of computations and processing on the data, so that we have associated a meaning to the data by way of some kind of relational connection. For example, precision and recall measures are obtained by relating the list of results contained in an experiment with the relevance judgements $J$;

▶ *knowledge*: the *descriptive statistics* and the *hypothesis tests* correspond to the "knowledge level" in the hierarchy, since they are a further elaboration of the information carried by the performance measurements and provide us with some insights about the experiments;

▶ *wisdom*: *theories*, *models*, *algorithms*, *techniques*, and *observations*, which are usually communicated by means of papers, talks, and seminars, correspond to the "wisdom level" in the hierarchy, since they provide interpretation, explanation, and formalization of the content of the previous levels.

## Infrustructure

▶ The experimental evaluation is usually carried out in important international evaluation campaigns which bring research groups together, provide them with the means for measuring the performances of their systems, discuss and compare their results.

▶ There are issues raised at international level which suggest that the IR experimental evaluation as a source of scientific data requests and the evaluation methodology itself need to be reconsidered to be properly supported by an organizational, hardware, and software infrastructures which allow for management, search, access, curation, enrichment, and citation of the produced scientific data.

▷ The EC in the i2010 Digital Library Initiative clearly states that "digital repositories of scientific information are essential elements to build European eInfrastructure for knowledge sharing and transfer, feeding the cycles of scientific research and innovation up-take" **?**.

▷ The US National Scientific Board points out that "organizations make choices on behalf of the current and future user community on issues such as collection access; collection structure; technical standards and processes for data curation; ontology development; annotation; and peer review" **?**.

▷ The Australian Working Group on Data for Science suggests to "establish a nationally supported long-term strategic framework for scientific data management, including guiding principles, policies, best practices and infrastructure" **?**.

## Extending Evaluation

▶ Scientific data, their curation, enrichment, and interpretation are essential components of scientific research. These issues are better faced and framed in the wider context of the *curation of scientific data*, which plays an important role on the systematic definition of a proper methodology to manage and promote the use of data.

▶ Therefore, we have to take into consideration the possibility of information enrichment of scientific data, meant as archiving and preserving scientific data so that the experiments, records, and observations will be available for future research, as well as provenance, curation, and citation of scientific data items.

▶ Furthermore, the importance of some of the many possible reasons for which keeping data is important are for example:

▷ re-use of data for new research, including collection based research to generate new science;

▷ retention of unique observational data which is impossible to re-create;

▷ retention of expensively generated data which is cheaper to maintain than to re-generate;

▷ enhancing existing data available for research projects;

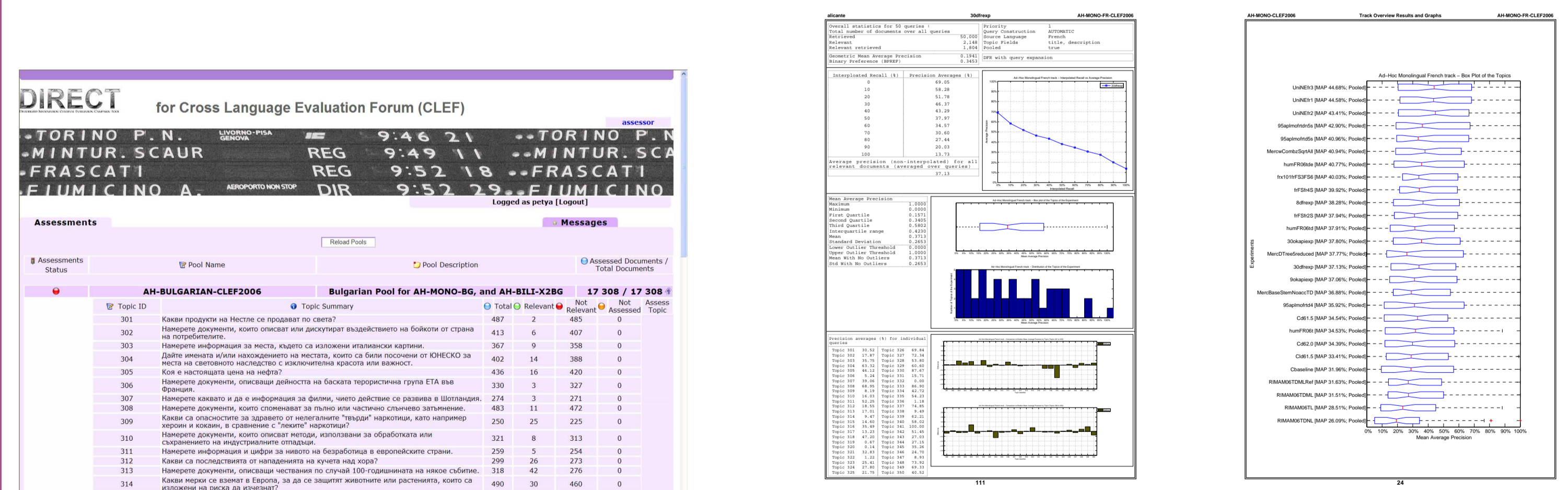▷ validating published research results.

## Key Points

**Conceptual model and metadata** the information space implied by an evaluation campaign needs an appropriate conceptual model which takes into consideration and describes all the entities involved by the evaluation campaign. From a conceptual model we can derive also appropriate data formats for exchanging information among organizers and participants.

**Unique Identification Mechanism** the lack of a conceptual model also implies that there is no common mechanism for uniquely identify the different digital objects involved in an evaluation campaign. Indeed, the possibility of citing scientific data and their further elaboration is an effective way for making scientists and researchers an active part of the digital curation process.

**Statistical Analyses** in developing an infrastructure, it is advisable to add some form of support and guide to participants for adopting a more uniform way of performing statistical analyses on their own experiments. If this support is added, participants can not only benefit from standard experimental collections which make their experiments comparable, but they can also exploit standard tools for the analysis of the experimental results, which would make the analysis and assessment of their experiments comparable too.

## Running System



## References

[1] Agosti, M., Di Nunzio, G. M., and Ferro, N. (2006a). A Data Curation Approach to Support In-depth Evaluation Studies. In *Proc. International Workshop on New Directions in Multilingual Information Access (MLIA 2006)*, pages 65–68. http://ucdata.berkeley.edu/sigir2006-mlia.htm.

[2] Agosti, M., Di Nunzio, G. M., and Ferro, N. (2006b). Scientific Data of an Evaluation Campaign: Do We Properly Deal With Them? In *Working Notes for the CLEF 2006 Workshop*. http://www.clef-campaign.org/2006/working_notes/workingnotes2006/agostiCLEF2006.pdf.

[3] Di Nunzio, G. M. and Ferro, N. (2005). DIRECT: a System for Evaluating Information Access Components of Digital Libraries. In *Proc. 9th European Conference on Research and Advanced Technology for Digital Libraries (ECDL 2005)*, pages 483–484. Lecture Notes in Computer Science (LNCS) 3652, Springer, Heidelberg, Germany.

[4] Di Nunzio, G. M. and Ferro, N. (2006). Scientific Evaluation of a DLMS: a service for evaluating information access components. In *Proc. 10th European Conference on Research and Advanced Technology for Digital Libraries (ECDL 2006)*, pages 536–539. Lecture Notes in Computer Science (LNCS) 4172, Springer, Heidelberg, Germany.