

XBenchMatch: a Benchmark for XML Schema Matching Tools



Fabien Duchateau¹, Zohra Bellahsene¹ and Ela Hunt²
¹LIRMM, Univ. Montpellier 2-CNRS, ²ETH Zurich

duchatea@lirmm.fr, bella@lirmm.fr, hunt@inf.ethz.ch



XBenchMatch uses as input the result of a schema matching algorithm (set of mappings and/or an integrated schema) and generates statistics about the quality of this input and the performance of the matching tool. A demo version of the prototype is available at <http://www.lirmm.fr/duchatea/XBenchMatch>.

GOALS: extensibility, portability, simplicity (ease of use), scalability, genericity, completeness

XBenchMatch FEATURES

- **Extensibility.** The benchmark should be able to be extended to include new measures and new format
- **Portability.** The benchmark should be OS-independent,
- **Simplicity.** since both end-users and schema matching experts are targeted by this benchmark tool.
- **Scalability** on two aspects: creating new benchmark scenarii is an easy task. And a benchmark composed of many scenarii should be easy to build and evaluate.
- **Genericity.** It should work with most of the available matchers.

KIND OF EVALUATION

1. **Quality of Mappings and Quality of Integrated Schema**
 - based on the use of the metrics
2. **Performance of Matching Algorithms (time).**

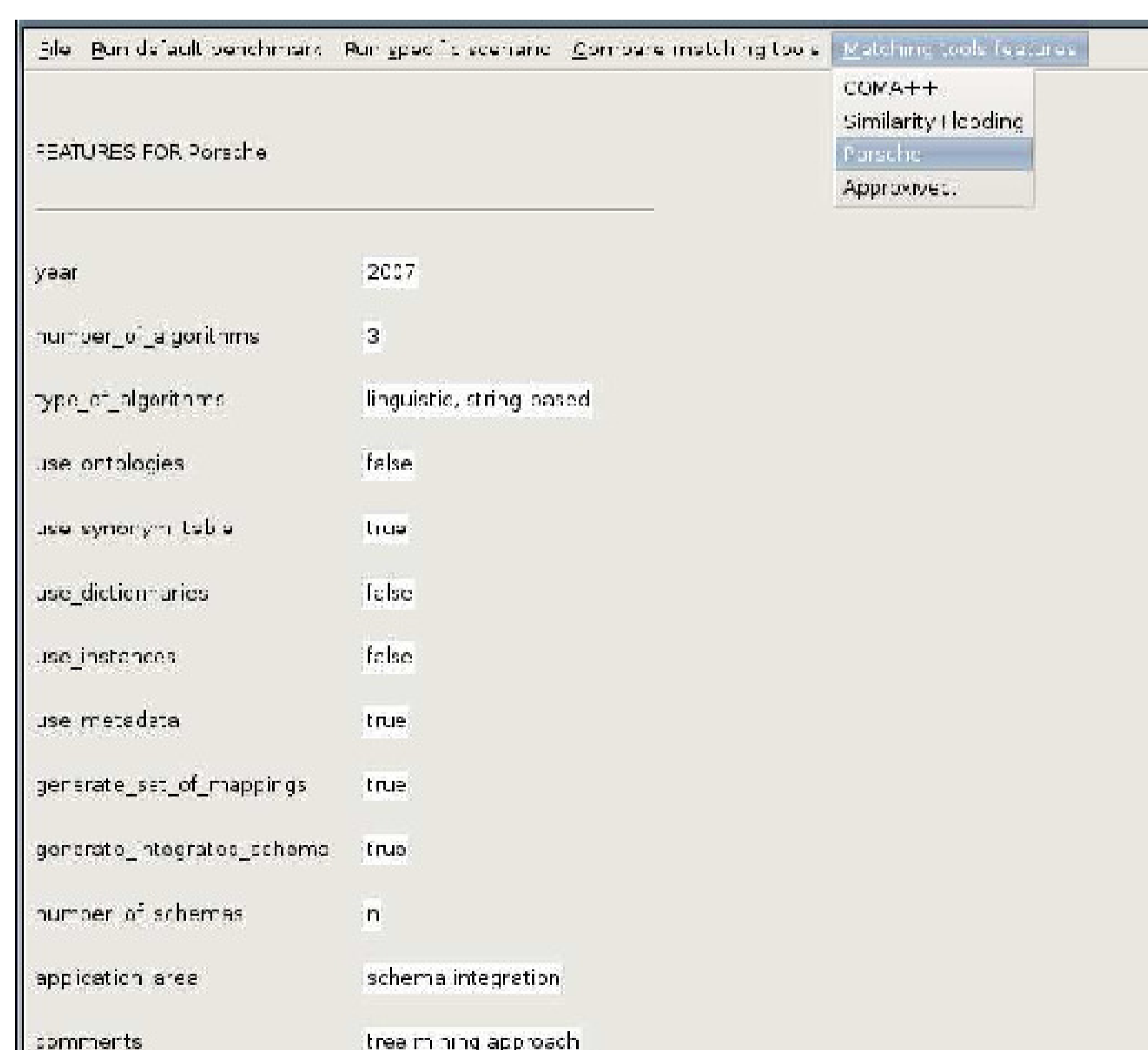
MAPPING QUALITY MEASURES, T_{map} are derived mappings, T_{ex} are expert mappings, expressed as paths:

1. Precision = $|T_{map} \cap T_{ex}| / |T_{map}|$
2. Recall = $|T_{map} \cap T_{ex}| / |T_{ex}|$
3. Fmeasure = $(2 \cdot \text{precision} \cdot \text{recall}) / (\text{precision} + \text{recall})$

INTEGRATED SCHEMA QUALITY MEASURES

for an integrated schema S_i , and an input schema S_g :

1. **Backbone measure, BM,** corresponds to the size of the largest common subtree of S_g and S_i (measured in nodes), seen against the background of the integrated schema S_i . **BM** = $|LCS_{sub}(S_i, S_g)| / |S_i|$
2. **Structural overlap** corresponds to the number of nodes shared by S_i and S_g and included in a common subtree. **Sub** is the set of all disjoint subtrees (each containing a minimum of two nodes) common to S_i and S_g . **kSub** is the total number of elements of all subtrees in **Sub**. **StructuralOverlap** = $kSub / |S_i|$
3. **Structural proximity** considers the number of subtrees common to S_i and S_g . **o** is the number of elements in S_i that are not included in any common subtree, **o** = $|S_i| - kSub$. **StructuralProximity** = $kSub / \sqrt{(|S_i| \times |Sub| + o)}$



XBenchMatch: a Benchmark for XML Schema Matching Tools

duchatea@lirmm.fr, bella@lirmm.fr, hunt@inf.ethz.ch

Experiments

SCHEMAS

1. Person schemas are small and strongly heterogeneous.
2. Purchase orders, XCBL collection 3, demonstrate matching of a large schema to a smaller one.
3. University course schemas are from Thalia [4].
4. Biological schemas correspond to Uniprot protein DB, and GeneCards integrate data from over 100 databases.

TESTED MATCHERS

Porsche, COMA++ and Similarity Flooding.

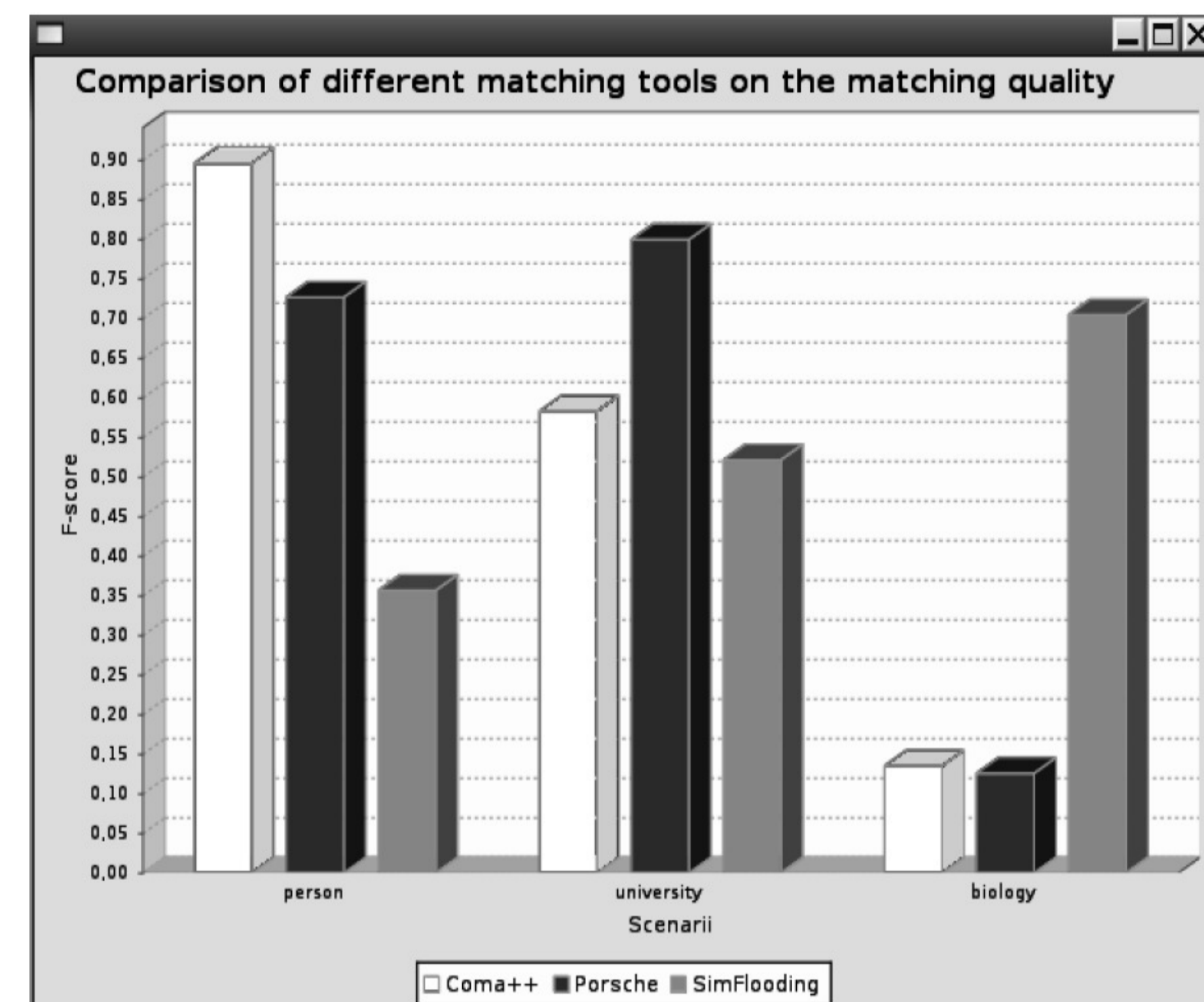


Figure 2. Matching precision on the three scenarios for three schema matchers.

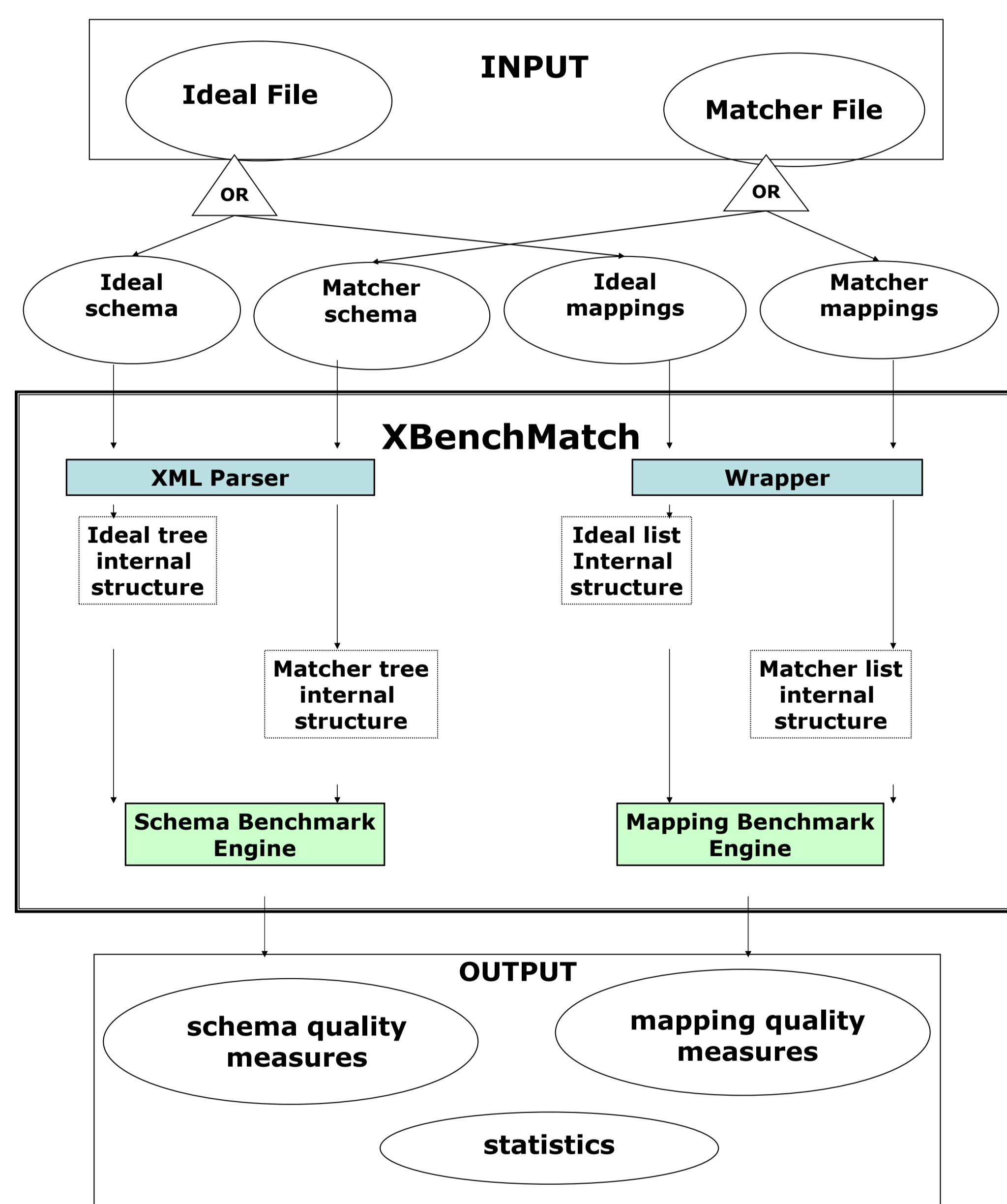


Figure 1. Architecture of XBenchMatch

	Person	University	Order	Biology
NB nodes (S_1 / S_2)	11 / 10	18 / 18	20 / 844	719 / 80
Avg NB of nodes	11	18	432	400
Max depth (S_1 / S_2)	4 / 4	5 / 3	3 / 3	7 / 3
NB of Mappings	5	15	10	57

Table 1: Summary of four evaluation scenarios.

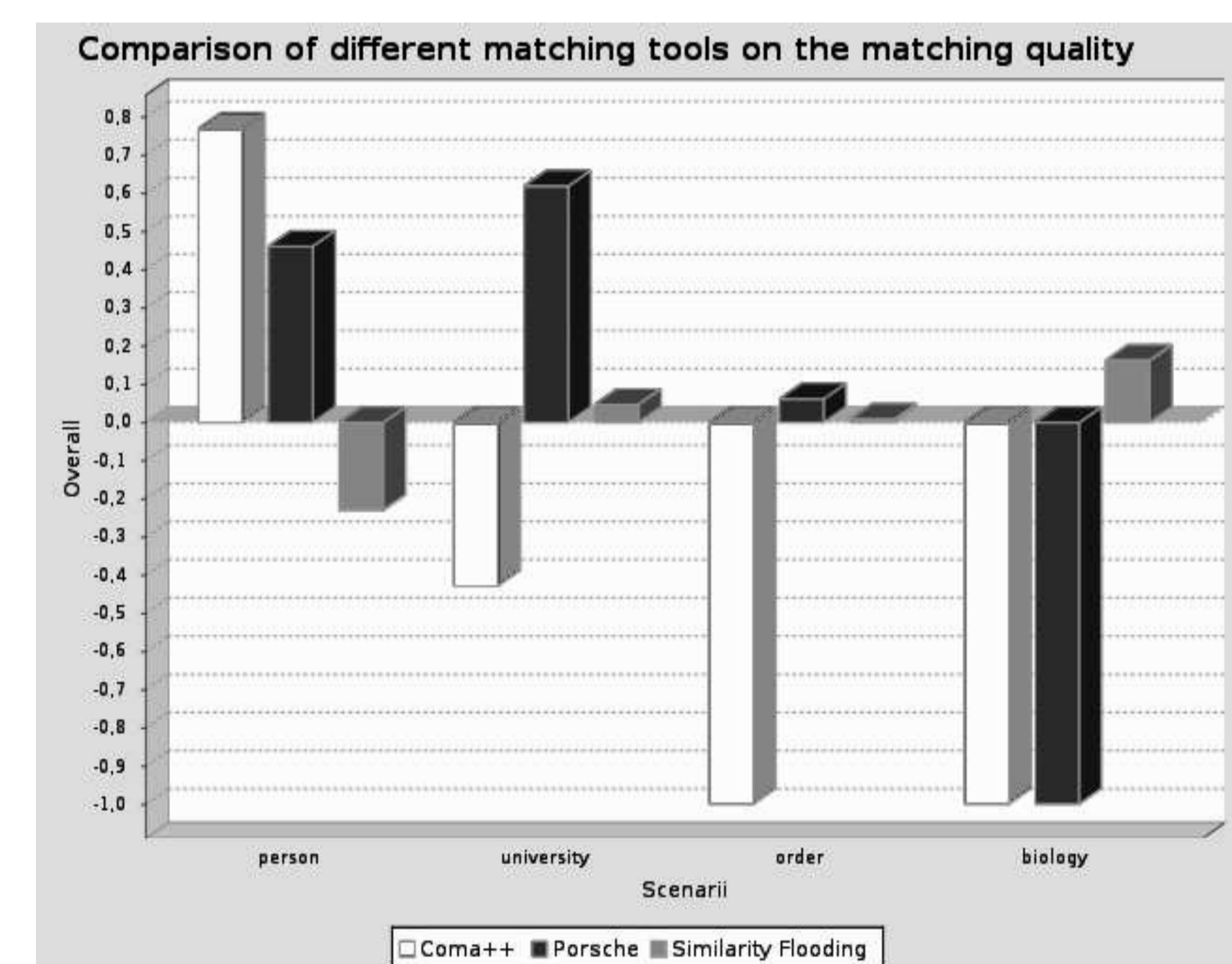


Figure 3. Matching quality of the four scenarios for three schema matchers.