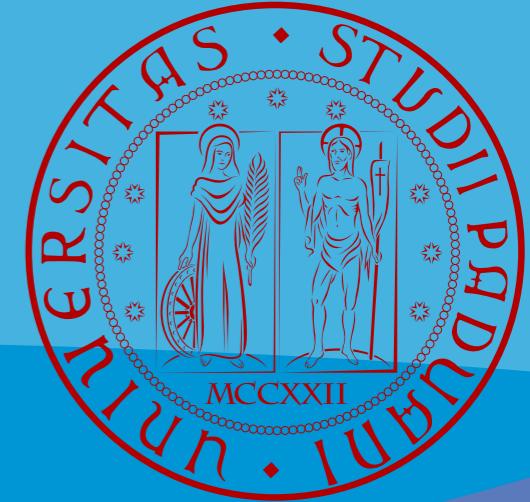




# PROMISE

Participative Research labOratory for Multimedia and  
Multilingual Information Systems Evaluation



# Evaluation Infrastructures

Nicola Ferro

Information Management Systems (IMS) Research Group  
Department of Information Engineering (DEI)  
University of Padua, Italy

PROMISE Winter School 2013 on Bridging between Information Retrieval and Databases  
Bressanone, Italy, 7 February 2013



- Motivations for an evaluation infrastructure
- Modelling
- Architecture and outcomes
- Final remarks

# Why an Evaluation Infrastructure?

**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**

# What is an Infrastructure?

**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**



# What is an Infrastructure?

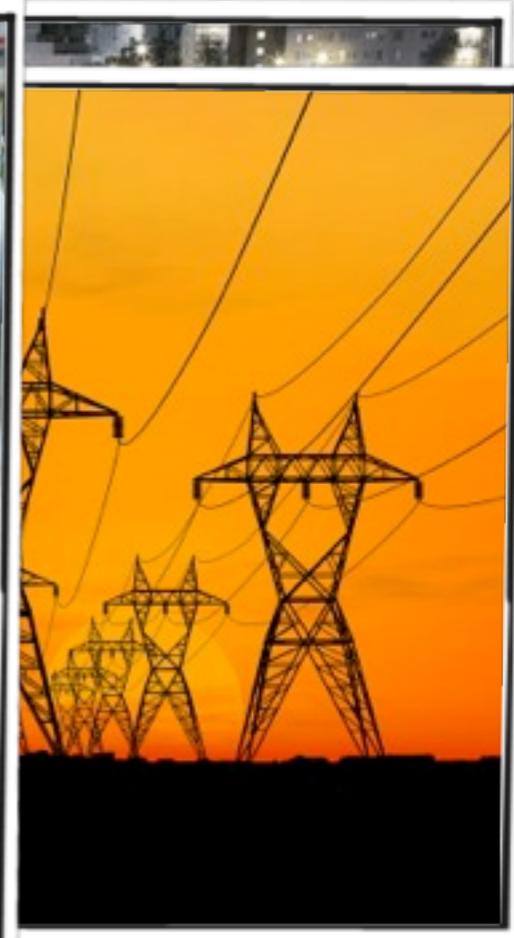
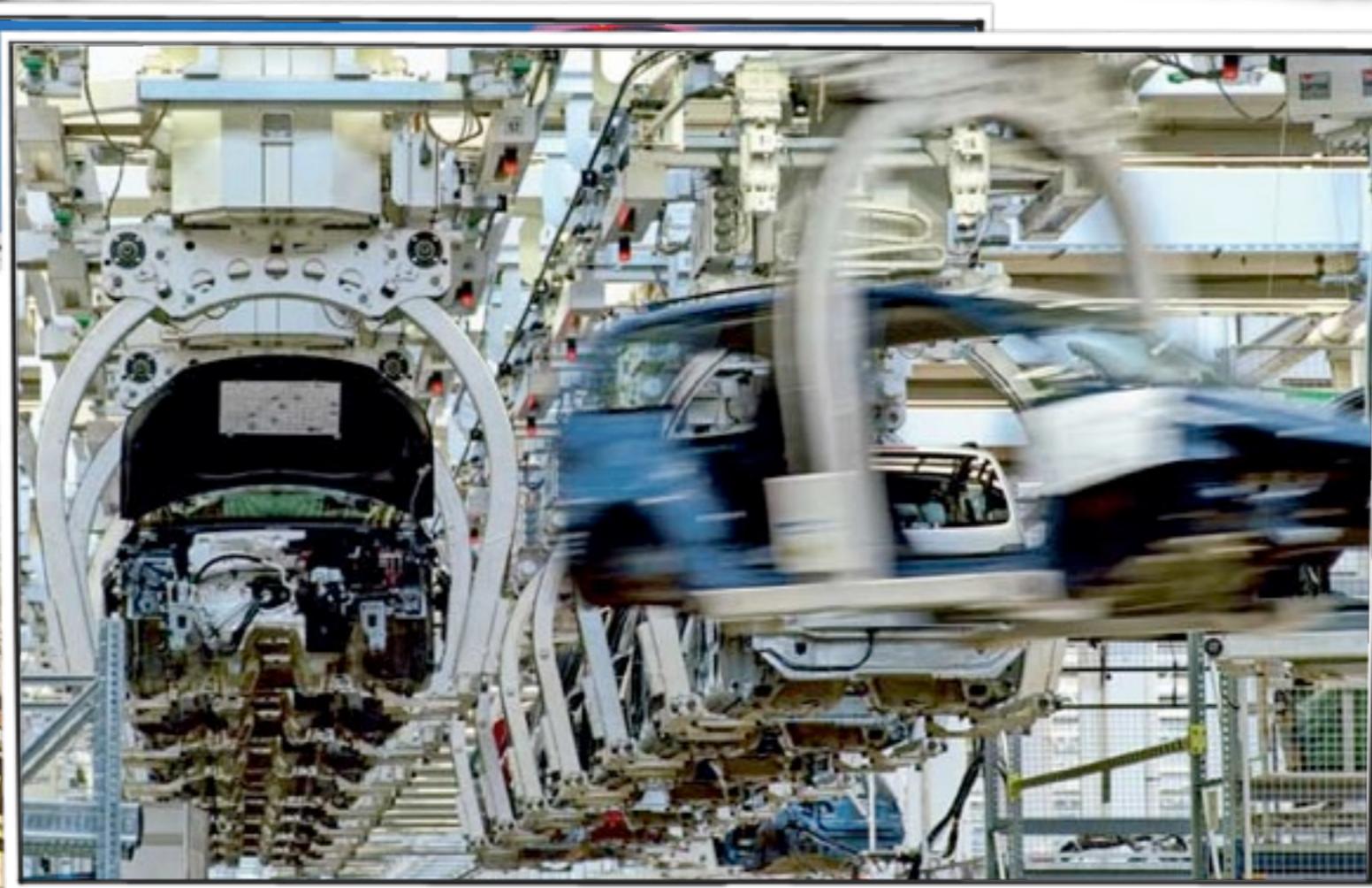
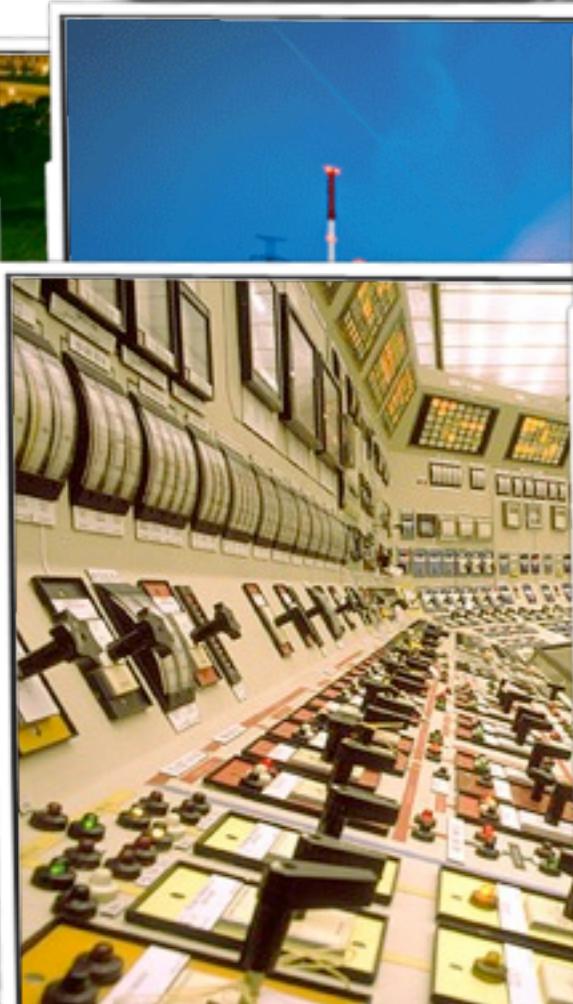
**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**



**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**



**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**



# What is an Infrastructure?

**the basic physical and organizational structures and facilities (e.g., buildings, roads, and power supplies) needed for the operation of a society or enterprise**



# In our case



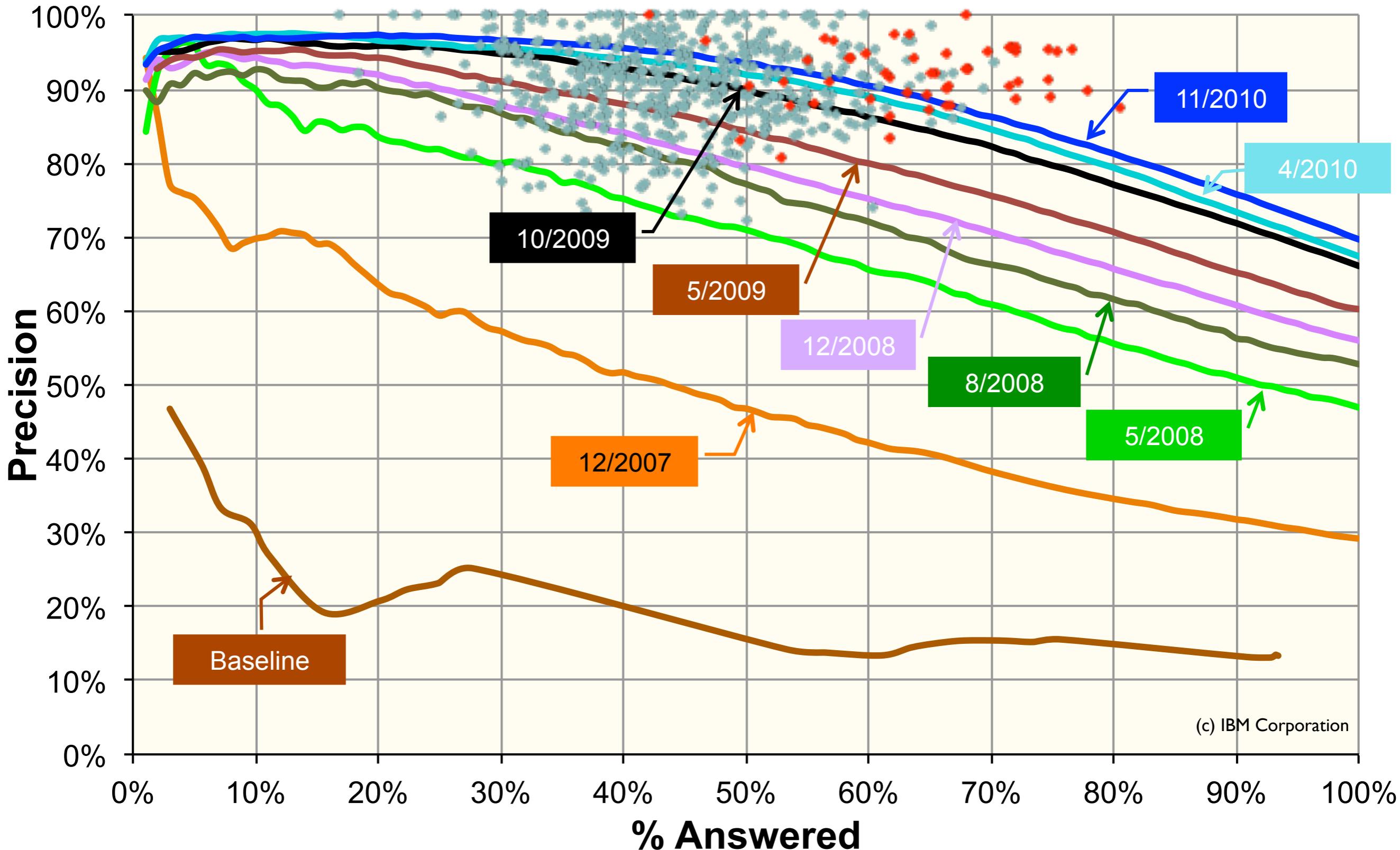
# Why Evaluation?

(c) IBM Corporation. <http://www.youtube.com/watch?v=3G2H3DZ8rNc>



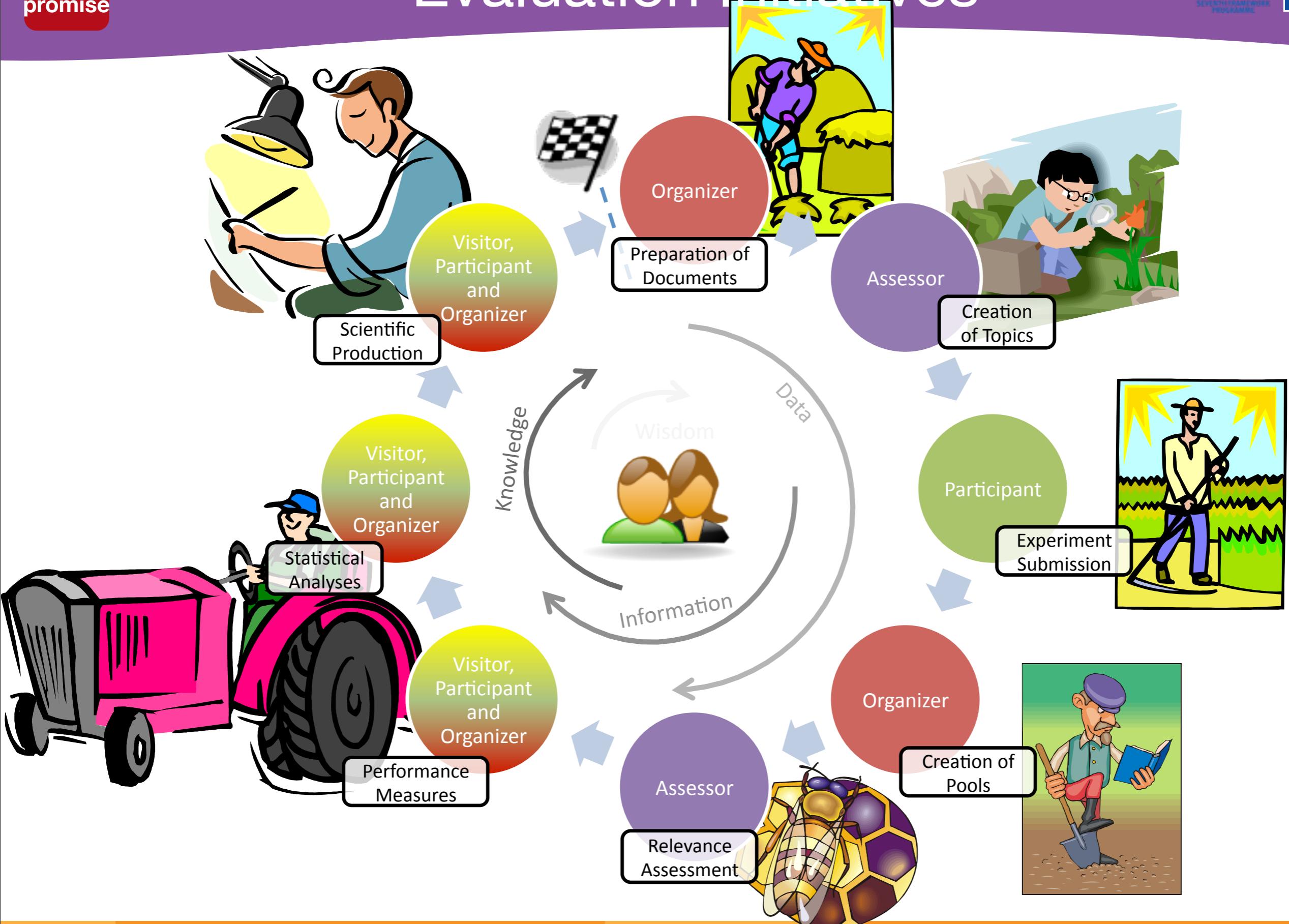
## IBM Watson: Deep QA Project

# Why Evaluation?



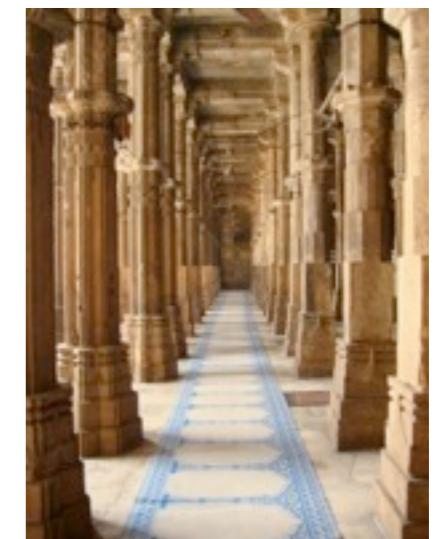
- Evaluation initiatives
- Researchers
- Developers and companies
- Stakeholders and adopters

# Evaluation Initiatives



# How much Valuable are our Data?

- The **TREC 2010 Economic Impact** study estimated in about **30 M\$** the **overall investment** in TREC by NIST
  - probably much much more if we had a means to estimate also the investment by participants in TREC
- The **Reliable Information Access (RIA)** workshop in 2004 estimated in **11 to 40 hours per topic** to carry out failure analysis
- They are the **pillars** for all the subsequent **scientific research and technology development**
  - TREC estimated the return on investment in the range of **3\$-5\$** for each invested dollar



- The TREC 2010 Economic Impact study estimated in about 30 M\$ the overall investment in TREC by NIST



## We need:

- (i) **to guarantee their accessibility and interpretability over time**
- (ii) **to curate and enrich them**
- (iii) **to ease their (re-)use**



- TREC estimated the return on investment in the range of 3\$-5\$ for each invested dollar

```

<topic number="6" type="ambiguous">
    <query>
        kcs
    </query>
    <description>
        Find information on the Kansas City Southern railroac
    </description>
    <subtopic number="1" type="nav">
        Find the homepage for the Kansas City Southern railrc
    </subtopic>
    <subtopic number="2" type="inf">
        I'm looking for a job with the Kansas City Southern r
    </subtopic>
    <subtopic number="3" type="nav">
        Find the homepage for Kanawha County Schools in West
    </subtopic>
    <subtopic number="4" type="nav">
        Find the homepage for the Knox County School system i
    </subtopic>
    <subtopic number="5" type="inf">
        Find information on KCS Energy, Inc., and their merge
    </subtopic>
</topic>

```

```

<session num="1" starttime="08:59:47.258675">
    <topic>
        <title>
            peacecorp
        </title>
        <desc>
            Find information about the peace corp
        </desc>
        <narr>
            When was it started and by whom? What services does it provide and where does
        </narr>
    </topic>
    <interaction num="1" starttime="09:00:04.155323">
        <query>
            peace corp
        </query>
        <results>
            <result rank="1">
                <url>
                    http://www.peacecorps.gov/
                </url>
                <clueweb09id>
                    clueweb09-en0011-60-08003
                </clueweb09id>
                <title>
                    Peace Corps
                </title>
                <snippet>
                    Fighting hunger, disease, poverty, and lack of opportunity.
                </snippet>
            </result>
            ...
        </results>
        <clicked>
            <click num="1" starttime="09:00:09.943356" endtime="09:01:13.434255">
                <rank>

```

## This situation hampers:

- automatic management
- accessibility and interpretability of the data
- ease of (re-)use
- interoperability
- take-up from new comers



```
303 0 APW19980609.1531 2
303 0 APW19980610.1778 1
303 0 APW19980715.1061 2
303 0 APW19980910.1078 0
```

```
1 0 clueweb09-en0120-13-20479 0
1 1 clueweb09-en0120-13-20479 0
1 2 clueweb09-en0120-13-20479 0
```

```
101 0 clueweb09-en0047-33-20039 1
101 0 clueweb09-en0004-66-09322 2
101 0 clueweb09-en0033-30-08382 0
101 0 clueweb09-en0000-45-05740 -2
101 0 clueweb09-en0020-92-11795 1
```

```
20002 0 clueweb09-en0006-85-33170 1 1 10.5
20004 0 clueweb09-en0005-28-20976 1 1 10.5
20006 0 clueweb09-en0010-07-21538 1 1 10.5
```

**ad-hoc**

**diversity**

**ad-hoc with grades**

**relevance feedback**

This situation hampers:

- automatic management
- accessibility and interpretability of the data
- ease of (re-)use
- interoperability
- take-up from new comers



303 0 APW19980609.1531 2  
303 0 APW19980610.1778 1

303 0 APW19980715.1061 2  
303 0 APW19980910.1078 0

## We need:

- (i) **to agree on a common data model which allows for extension**
- (ii) **to provide the basic experimental data with proper metadata (descriptive, administrative, copyright, ...)**

- interoperability
- take-up from new comers



- Explanation of experimental data is usually reported in scientific papers that do not provide direct links to them
- It is often difficult to exactly know which data have been used in a paper and have access to them
- It is ever more difficult to exactly know and reproduce data cleaning and processing operations



## We need:

- (i) to have the possibility of citing experimental data in our papers as any other references**
- (ii) to make our papers actionable and executable providing access to the mentioned experimental data**

- Automatic evaluation of a system
  - increase in the number of experiment
  - more detailed information
  - reduction of the effort to evaluate a system
  - more comparability
- Programmatic access
  - increase in the number of experiment
  - new usages and new applications driven by the community
  - evaluation as part of the development process, like unit testing
  - opportunistic use by third parties with confidentiality requirements
- Improve ground-truth creation
  - reduction of human effort
  - opportunity for more user-centred evaluation



- Lowering the effort (and costs) for carrying out experimental evaluation
- Reducing fragmentation in the experimental evaluation (diverse tasks and metrics, heterogeneous collections, different systems, ...)
- Providing management, curation, accessibility over the time to experimental data as well richer interaction, visualization, and exploration of them
- Supporting continuous evaluation of information access and retrieval systems
- Offering programmatic access to the experimental data and extensibility

# Bridging between IR and DB



- Modelling



- Designing



- Developing

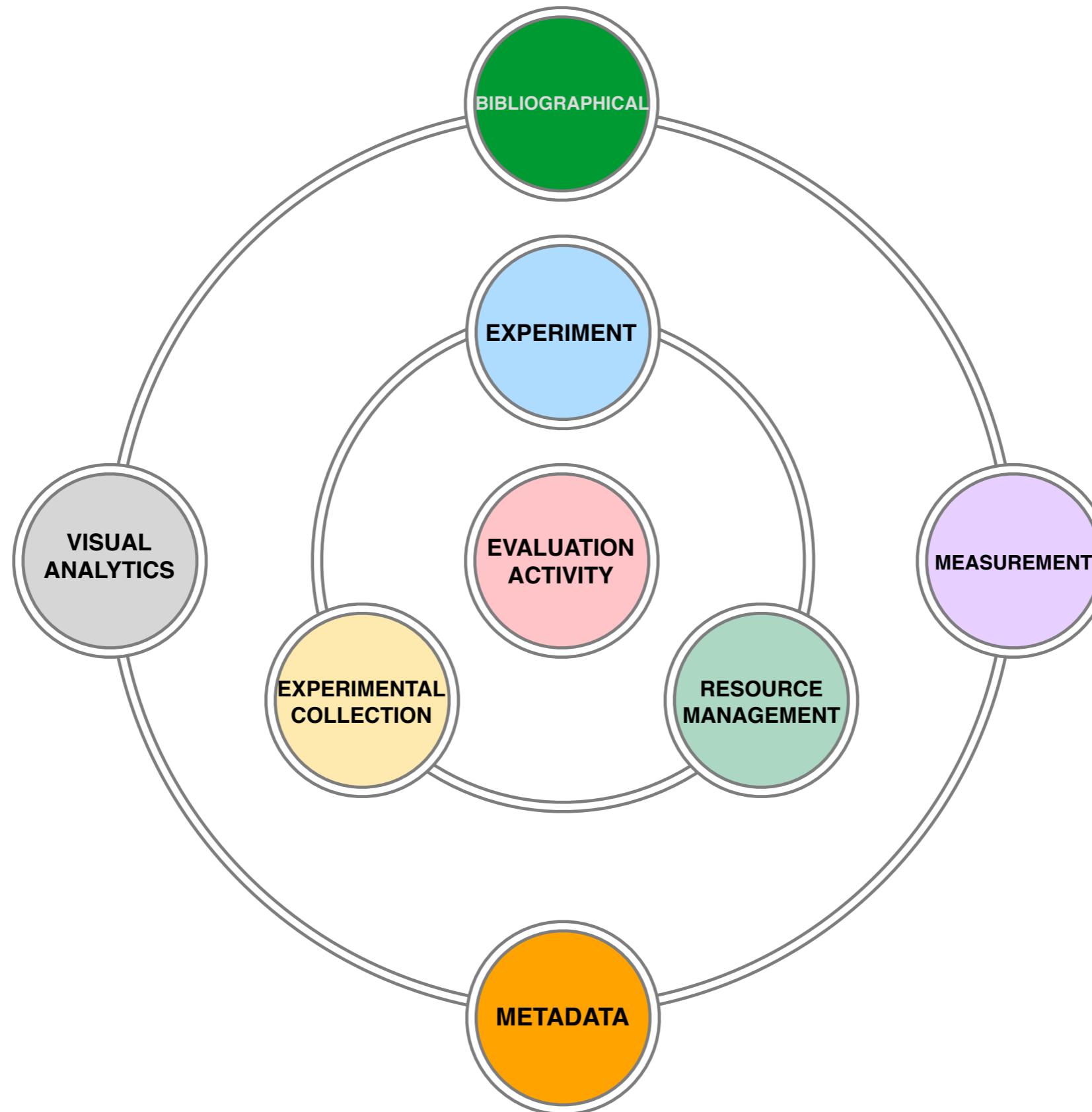


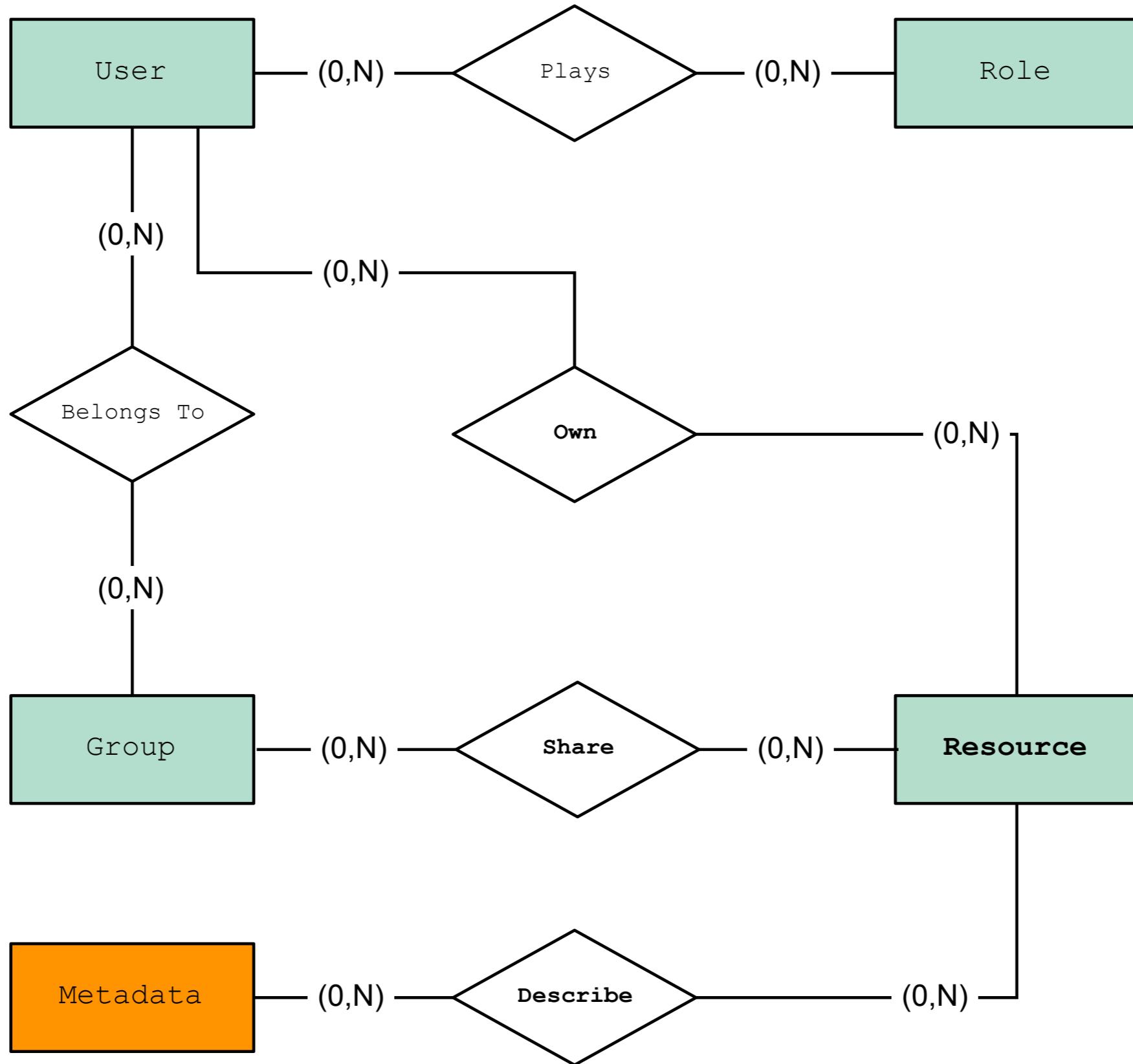
- Sharing

- Extending

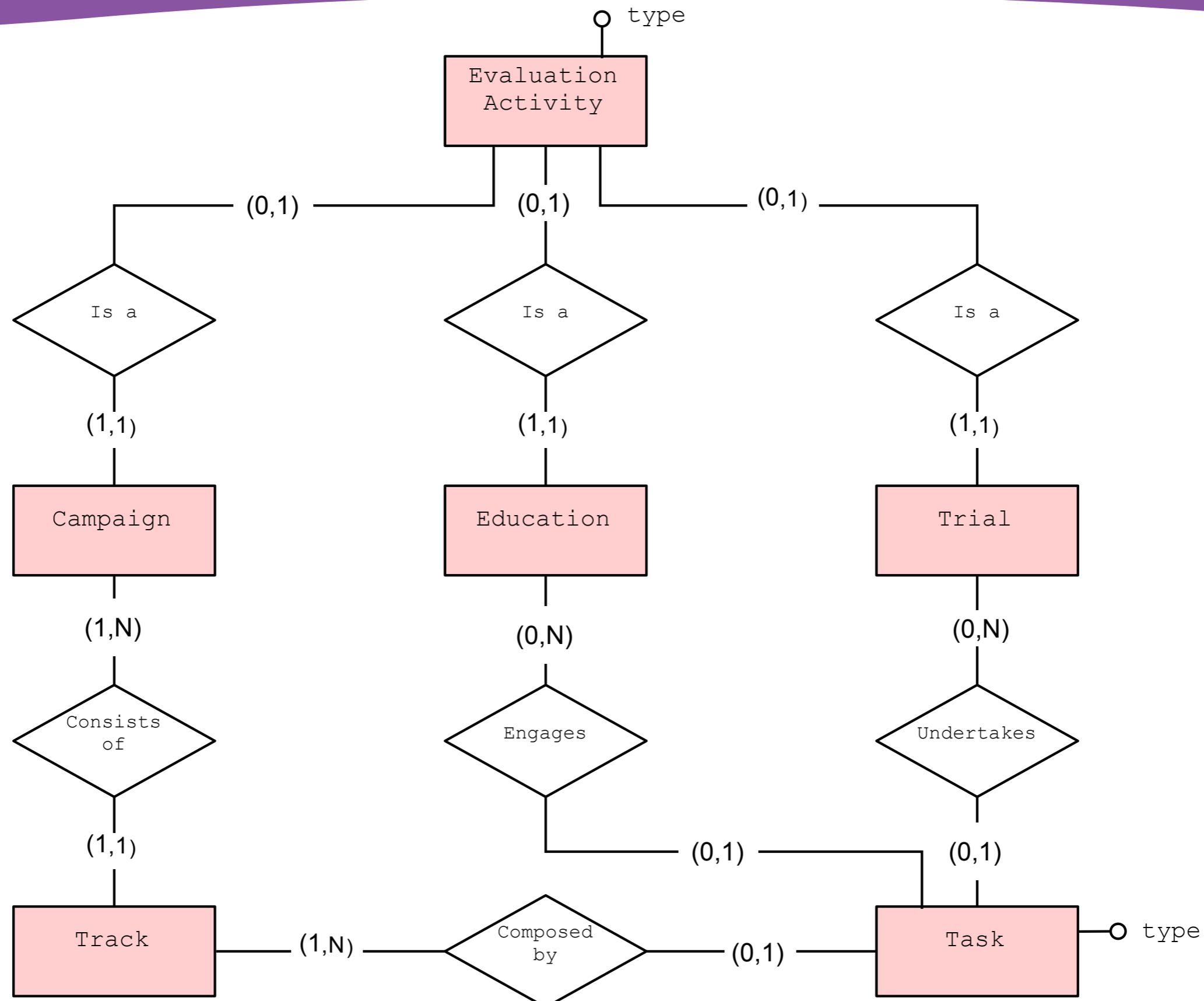
# Modelling

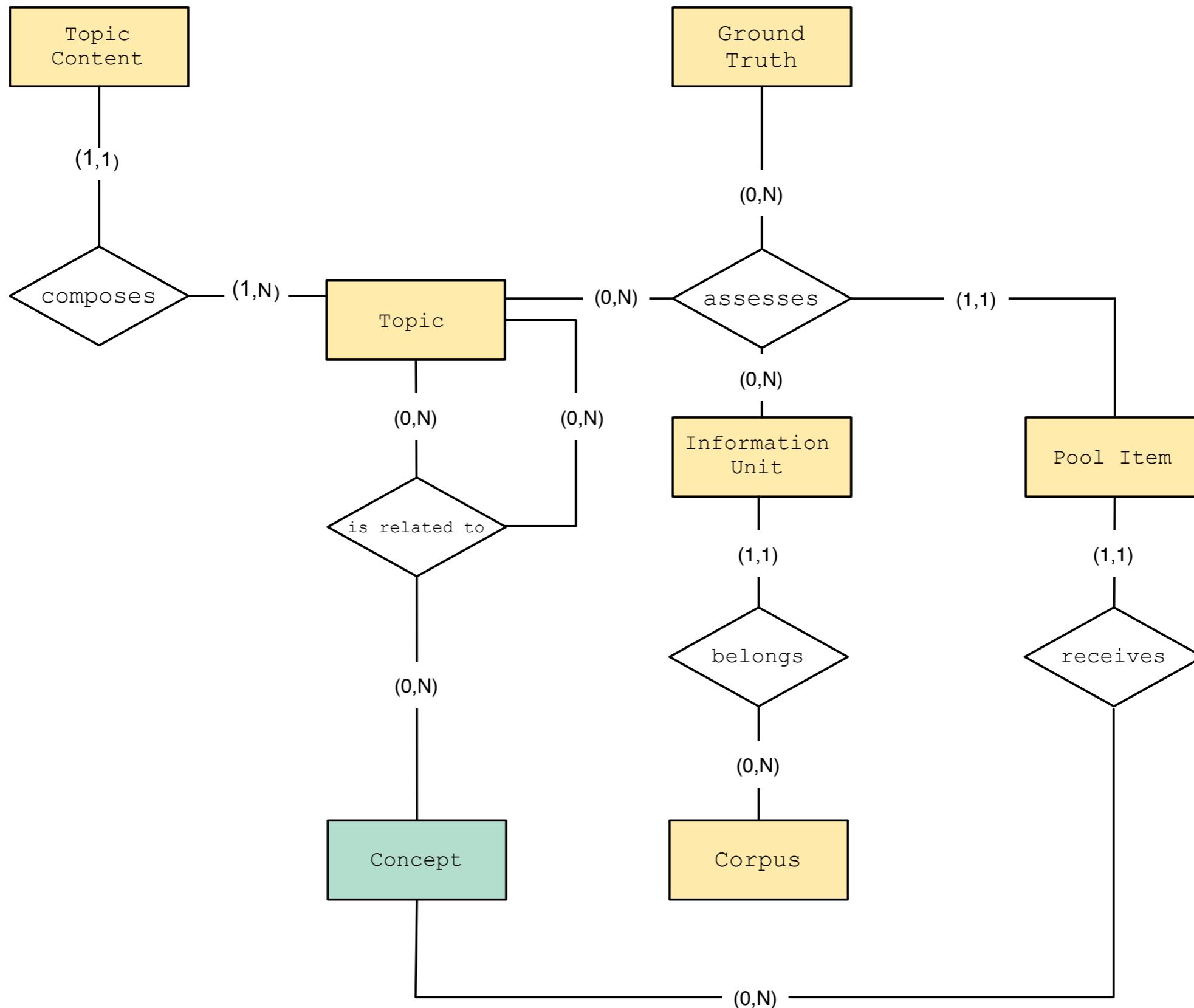
# Modelling Areas



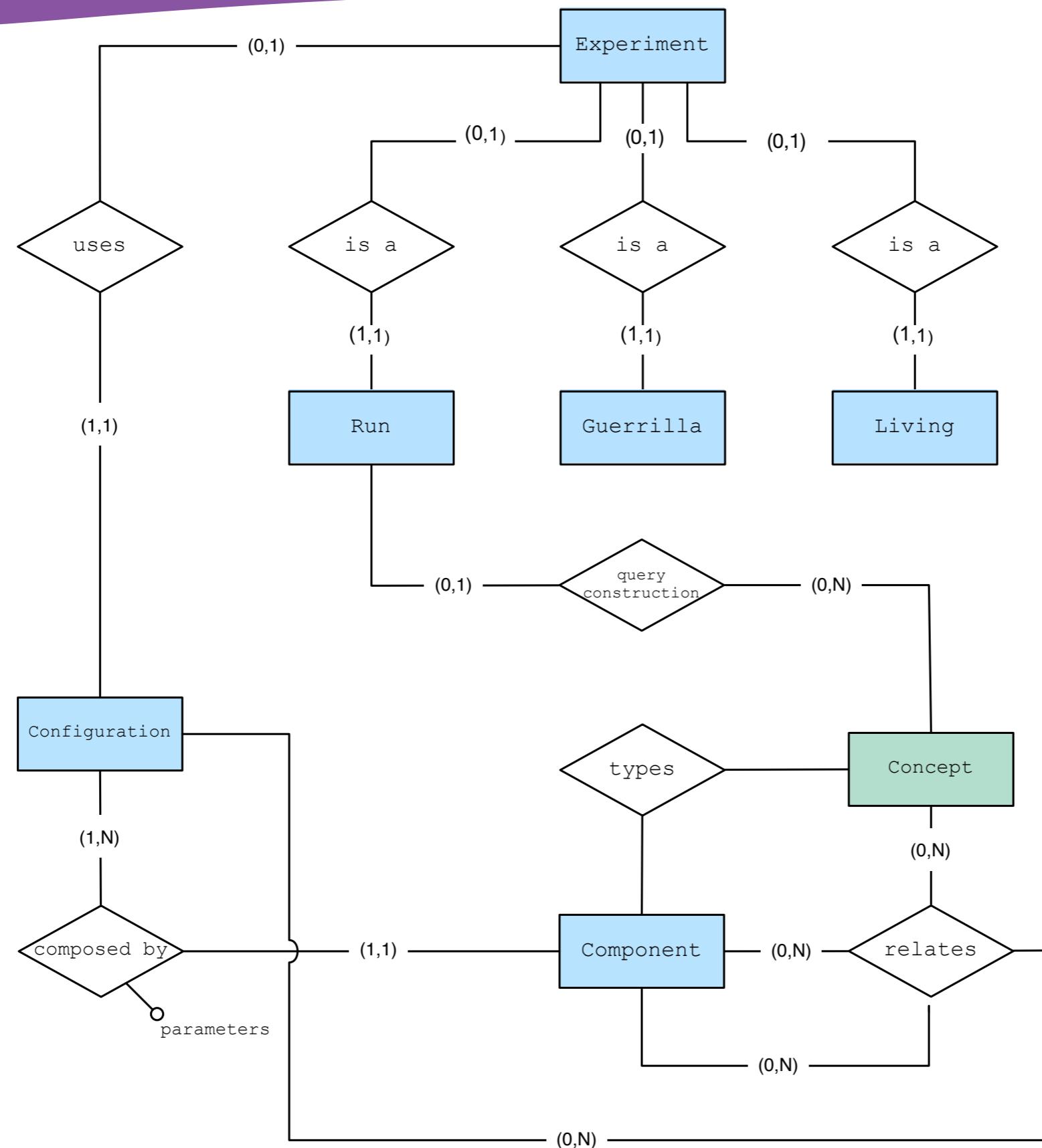


# Evaluation activities

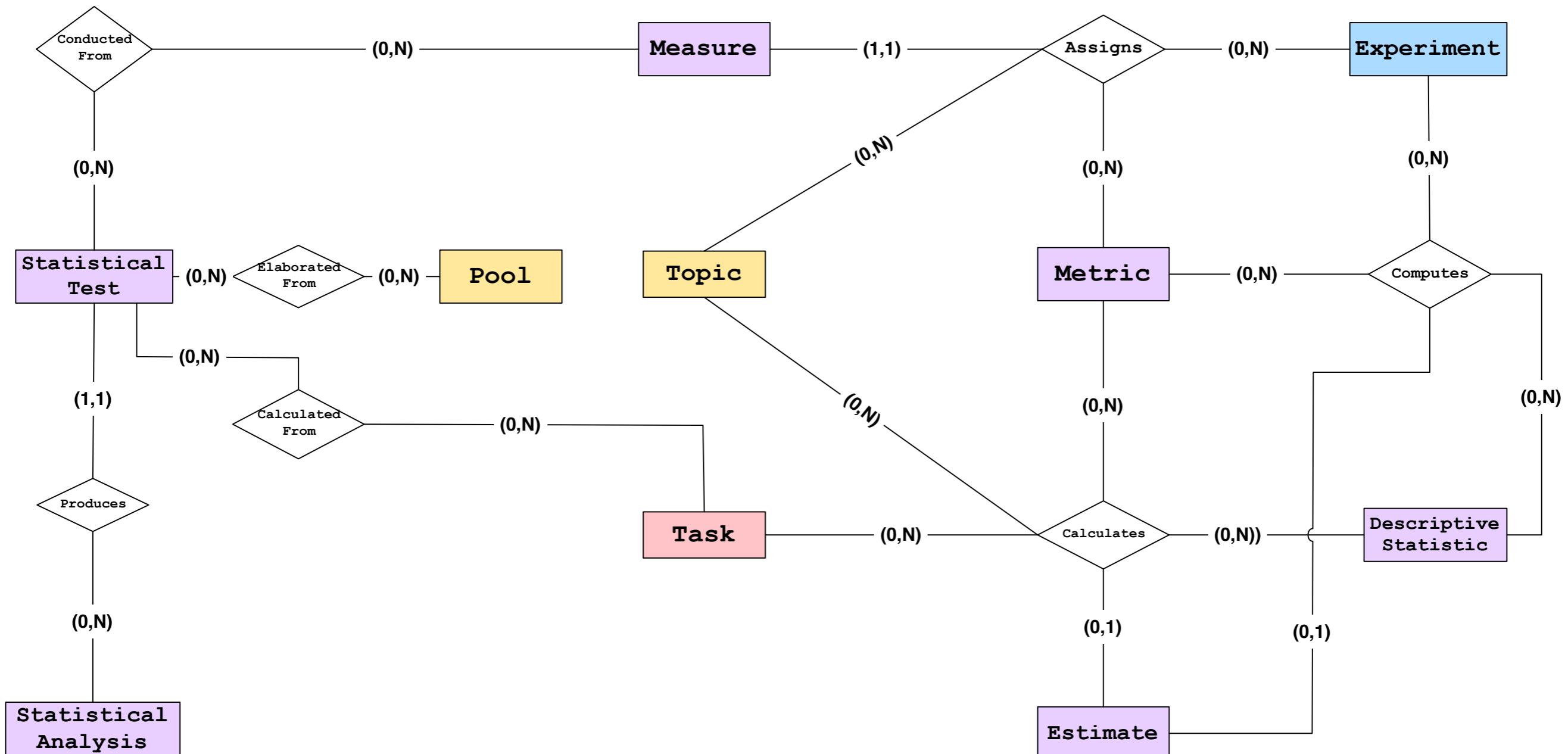




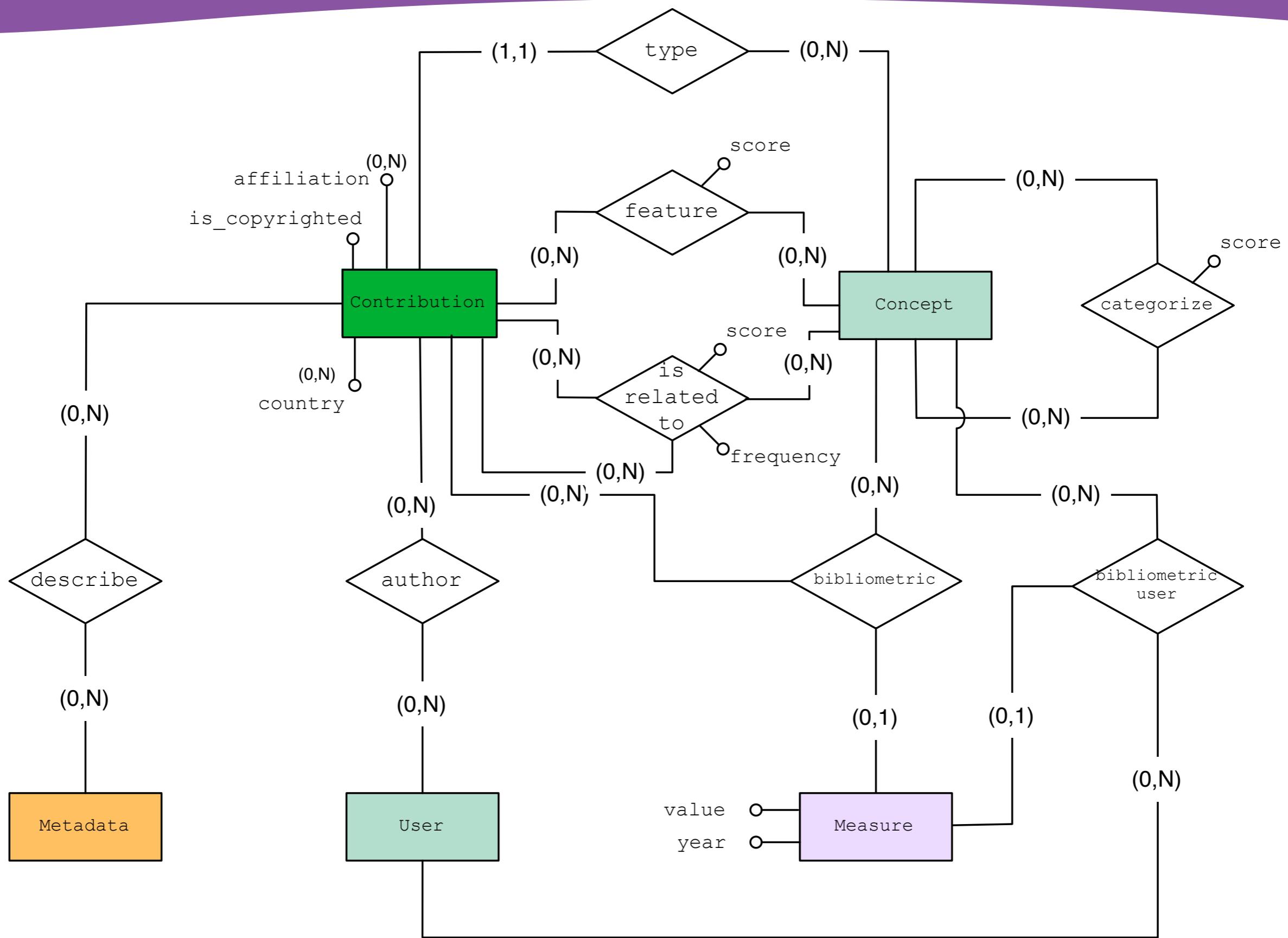
# Experiments



# Measurement



# Bibliography and Impact Analysis



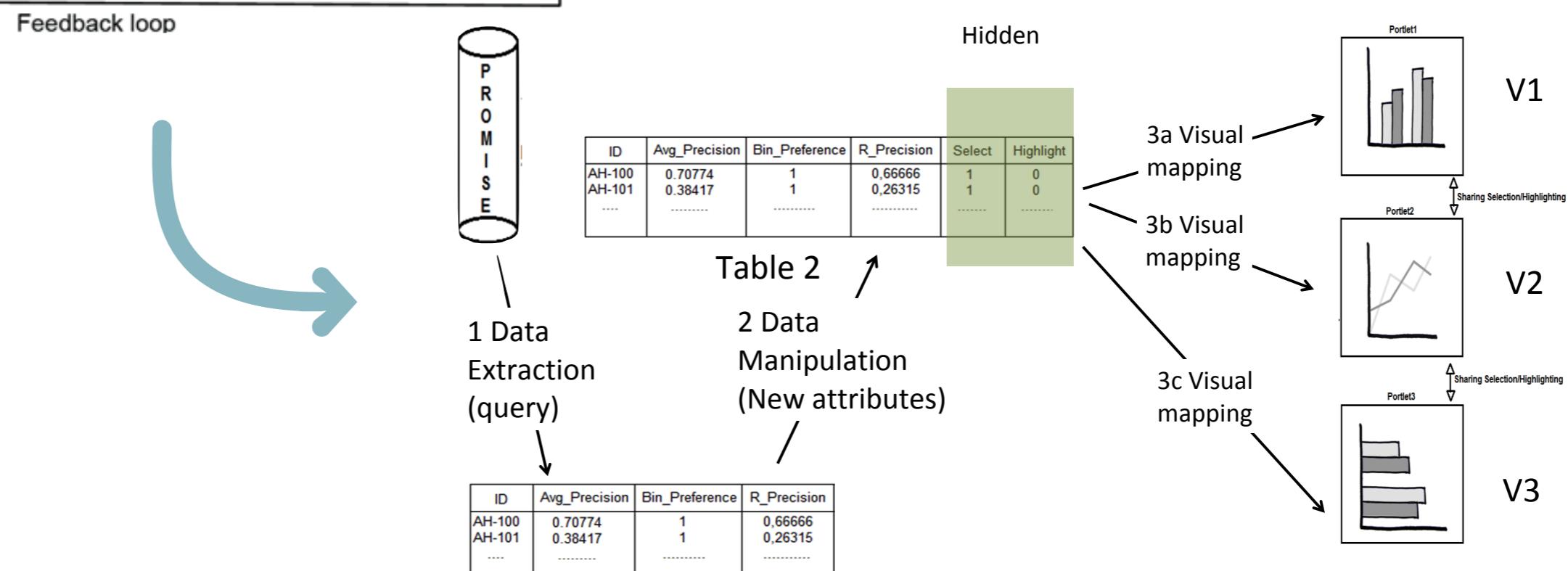
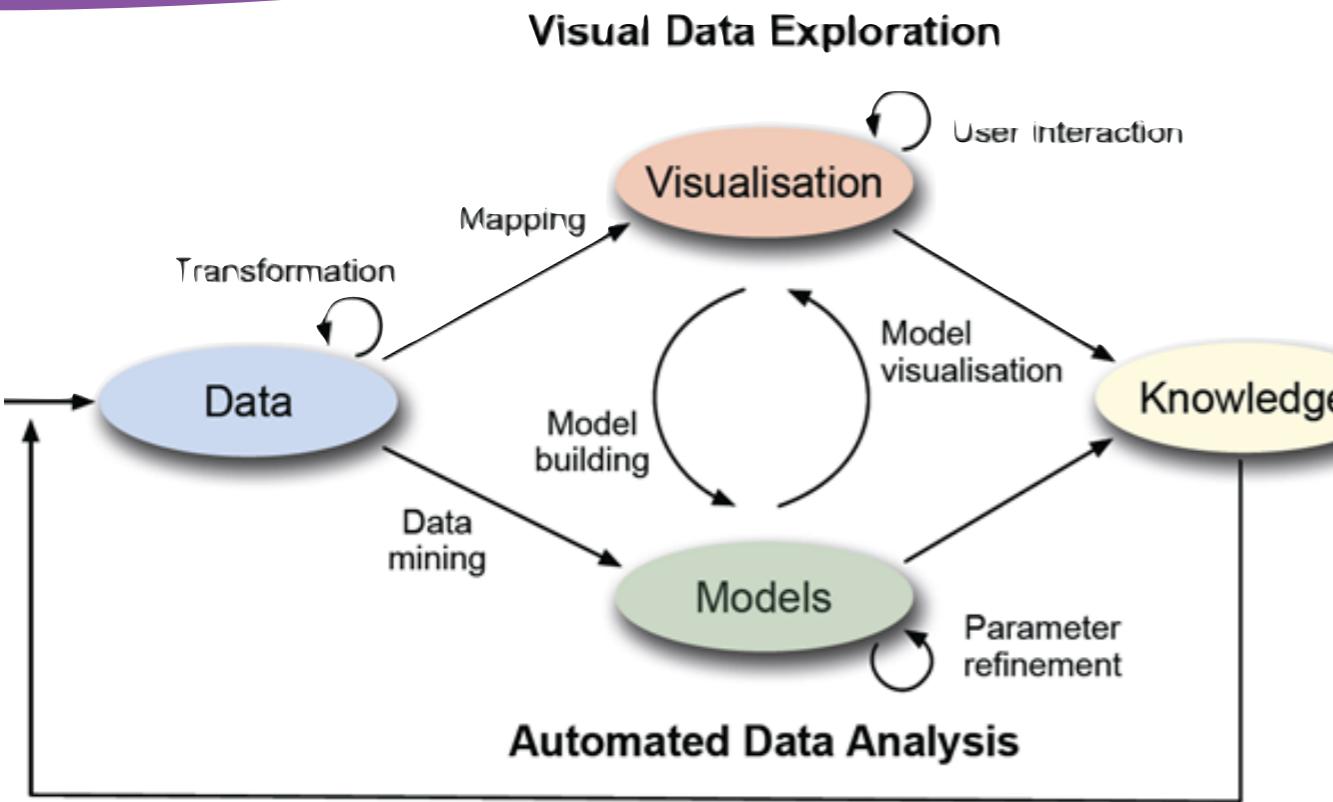


Table 1

## Visual Data Exploration

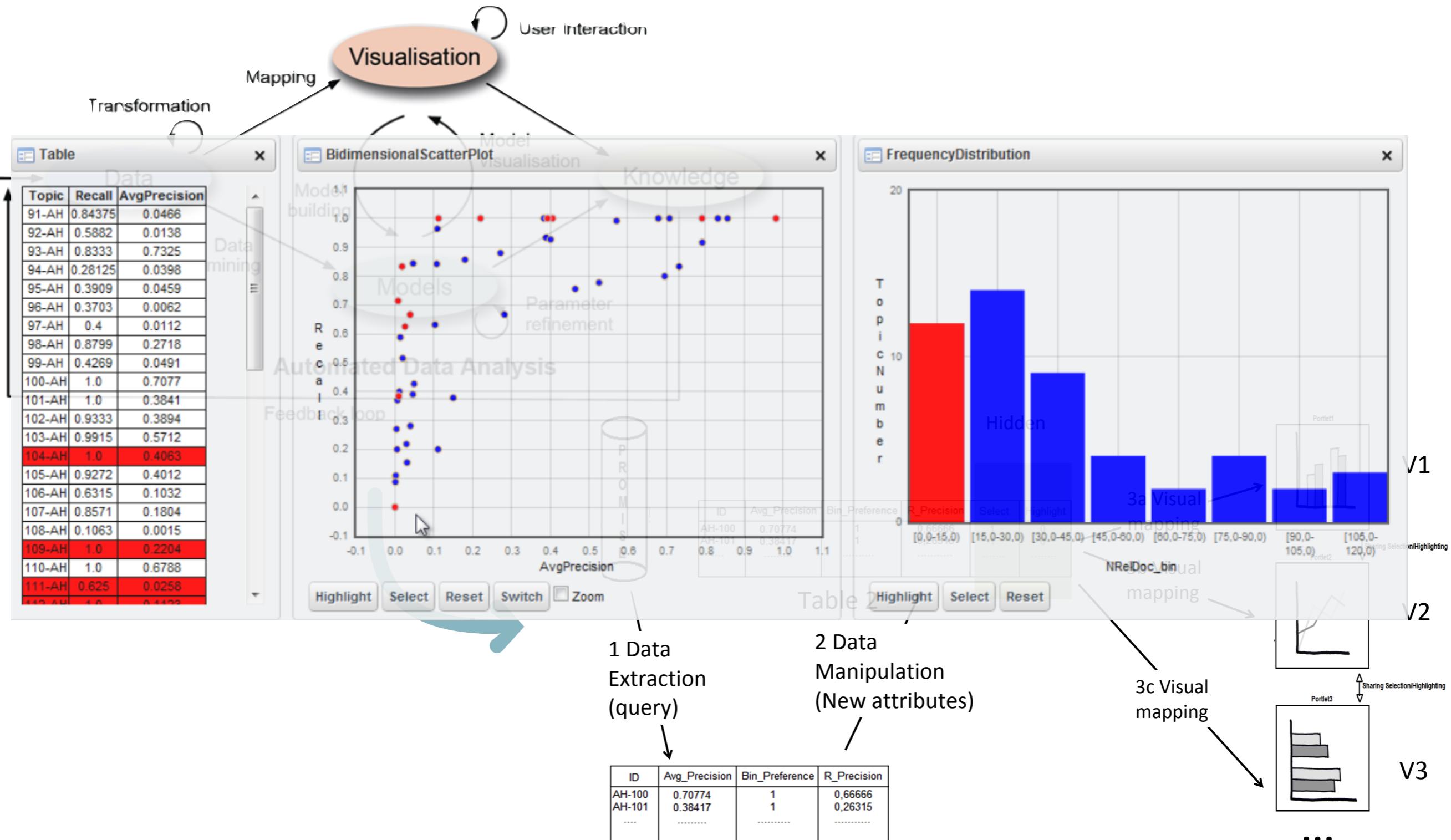
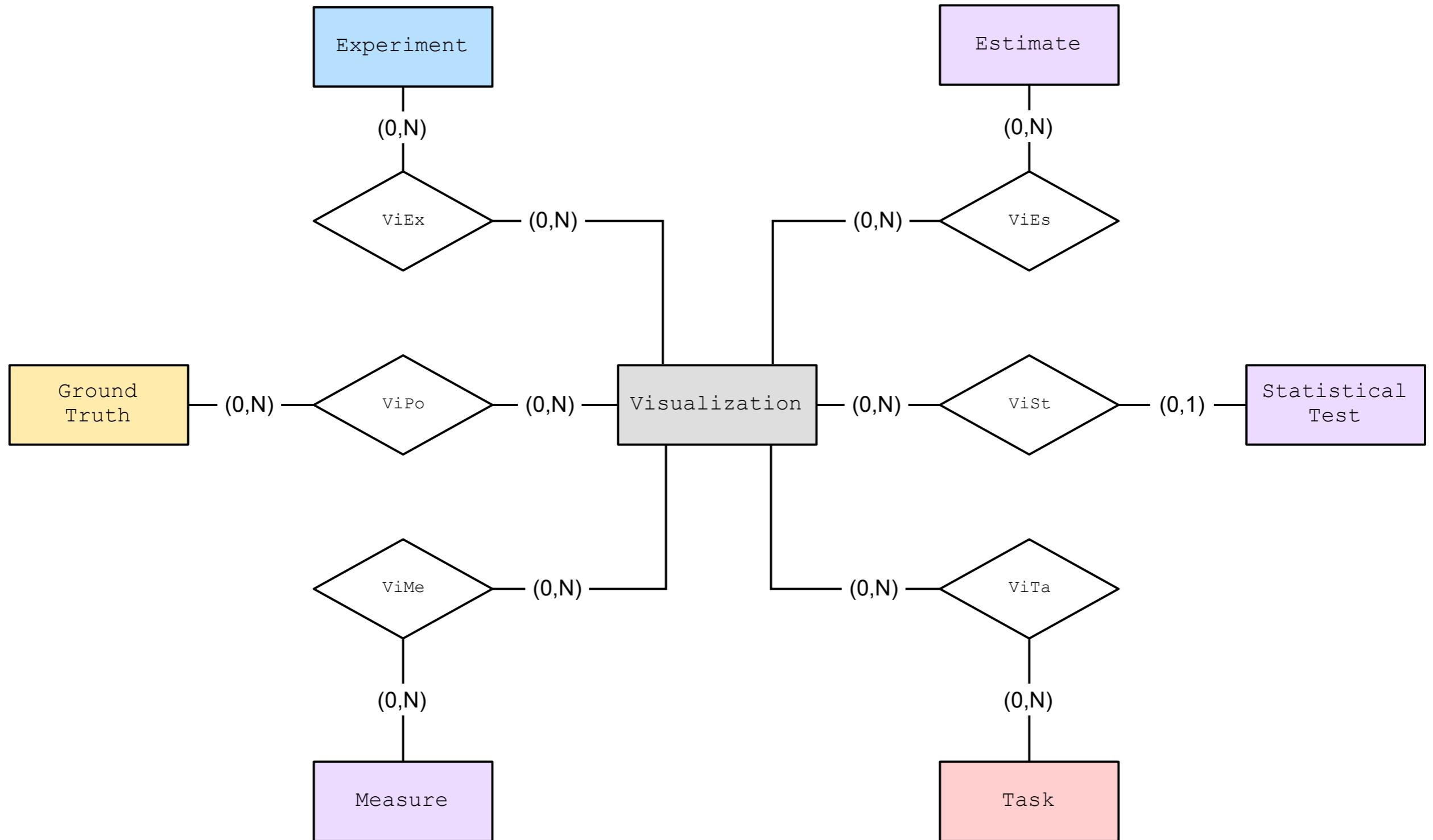
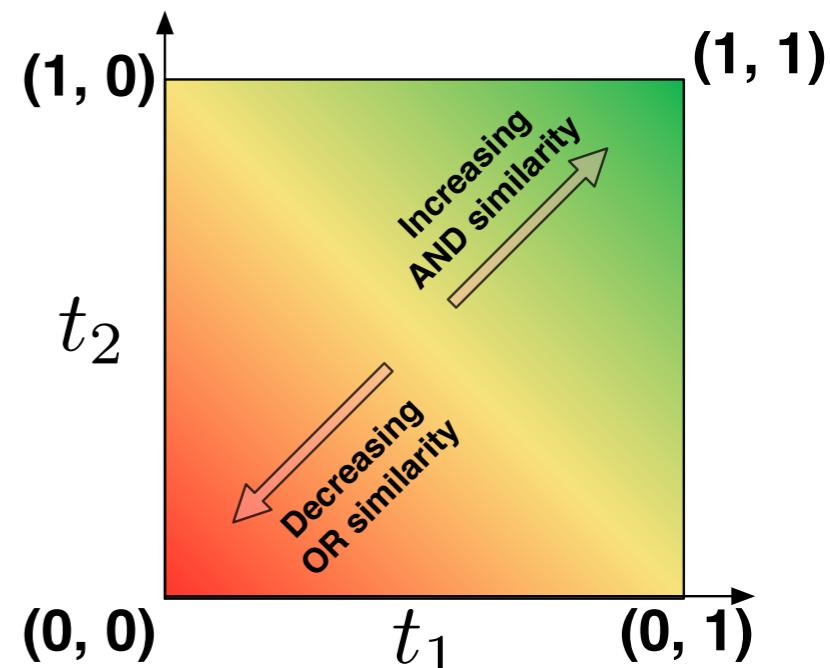


Table 1



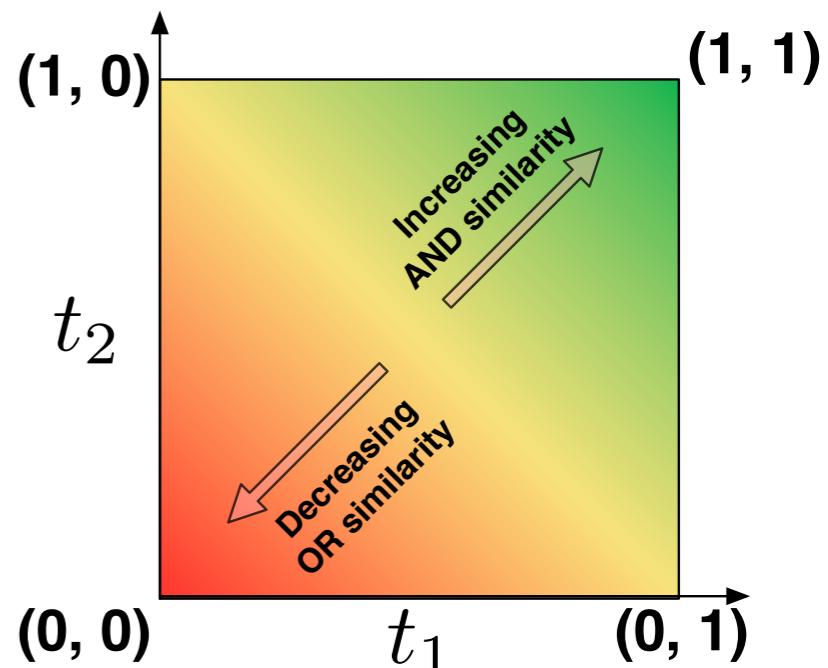
- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance



$$\text{sim}_p^{\text{or}}(r, q) = \left[ \frac{\text{sim}(r, t_1)^p + \dots + \text{sim}(r, t_n)^p}{n} \right]^{\frac{1}{p}}$$

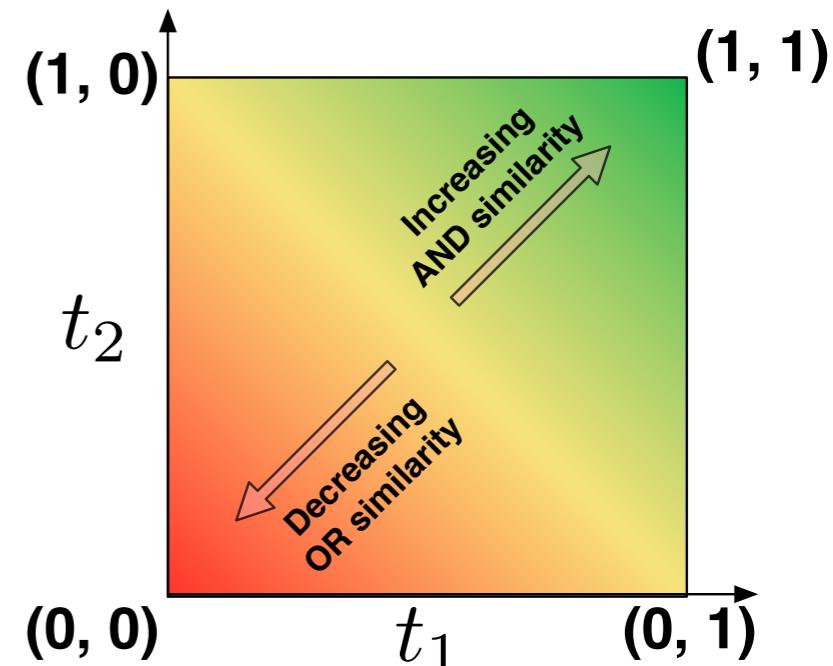
$$\text{sim}_p^{\text{and}}(r, q) = 1 - \left[ \frac{(1 - \text{sim}(r, t_1))^p + \dots + (1 - \text{sim}(r, t_n))^p}{n} \right]^{\frac{1}{p}}$$

- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance
- Best match ( $p = 1$ )



$$\text{sim}_{\text{best}}^{\text{or}}(r, q) = \text{sim}_{\text{best}}^{\text{and}}(r, q) = \frac{\text{sim}(r, t_1) + \text{sim}(r, t_2) + \dots + \text{sim}(r, t_n)}{n}$$

- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance

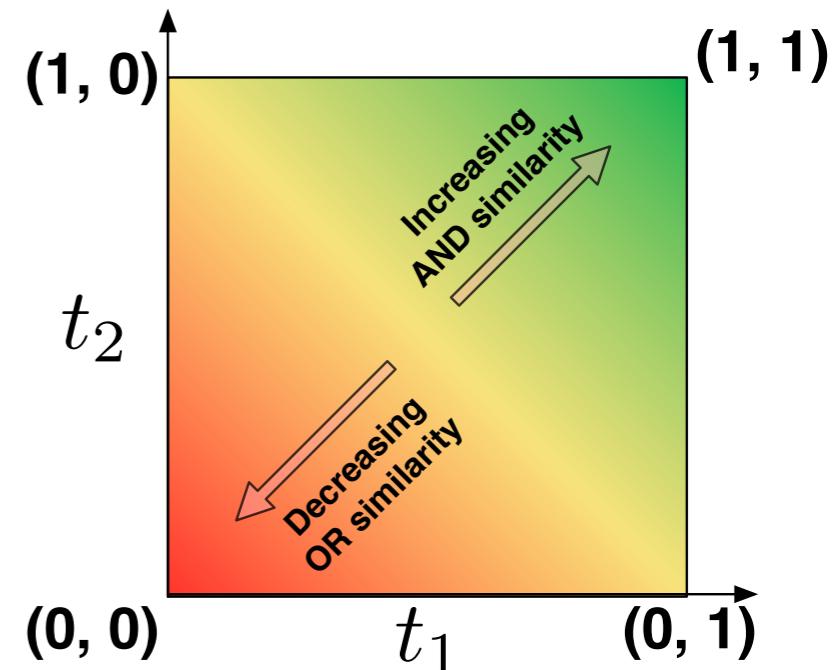


- Loose match ( $p = 5$ )

$$\text{sim}_{\text{loose}}^{\text{or}}(r, q) = \sqrt[5]{\frac{\text{sim}^5(r, t_1) + \text{sim}^5(r, t_2) + \dots + \text{sim}^5(r, t_n)}{n}}$$

$$\text{sim}_{\text{loose}}^{\text{and}}(r, q) = 1 - \sqrt[5]{\frac{(1 - \text{sim}(r, t_1))^5 + (1 - \text{sim}(r, t_2))^5 + \dots + (1 - \text{sim}(r, t_n))^5}{n}}$$

- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance

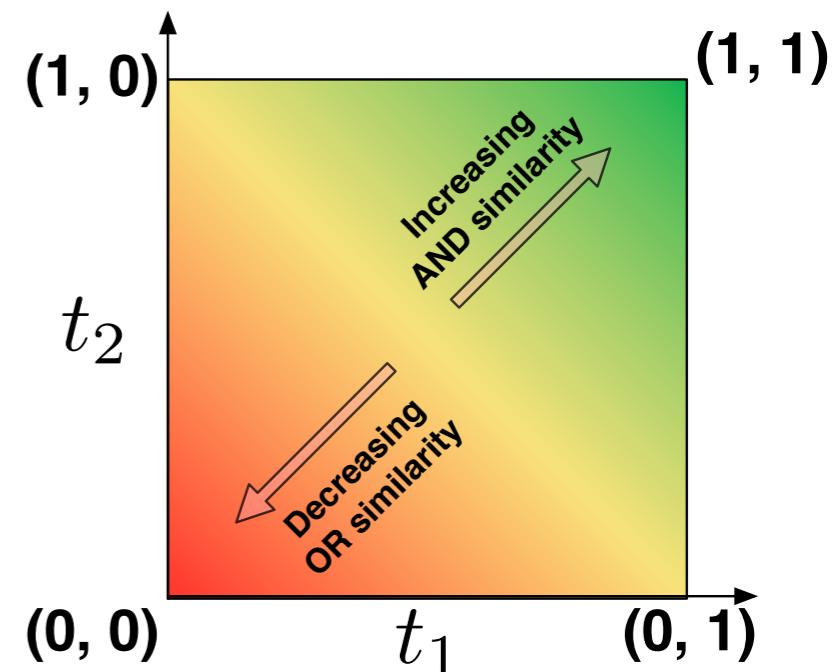


- Fuzzy match ( $p = 2$ )

$$\text{sim}_{\text{fuzzy}}^{\text{or}}(r, q) = \sqrt[2]{\frac{\text{sim}^2(r, t_1) + \text{sim}^2(r, t_2) + \dots + \text{sim}^2(r, t_n)}{n}}$$

$$\text{sim}_{\text{fuzzy}}^{\text{and}}(r, q) = 1 - \sqrt[2]{\frac{(1 - \text{sim}(r, t_1))^2 + (1 - \text{sim}(r, t_2))^2 + \dots + (1 - \text{sim}(r, t_n))^2}{n}}$$

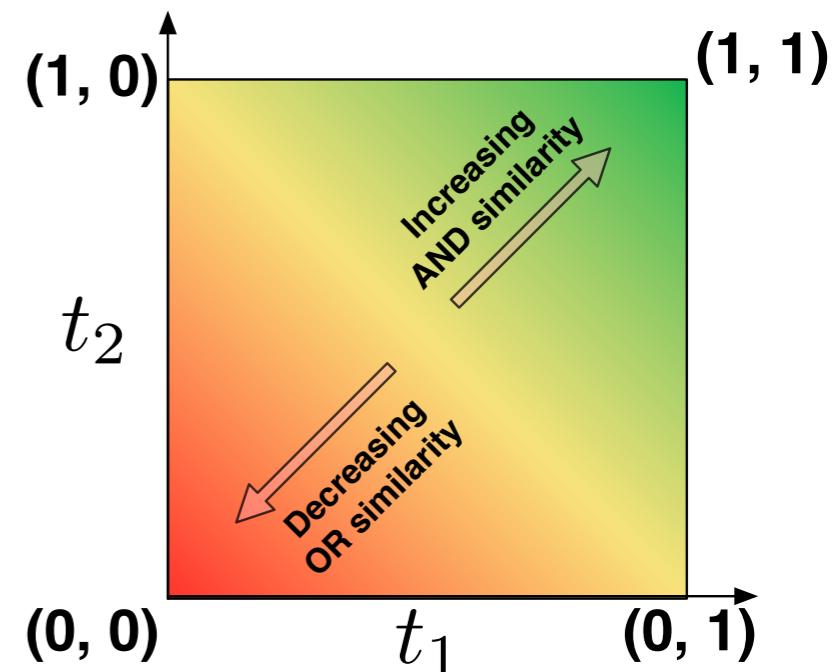
- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance
- Exact match ( $p = \infty$ )



$$\text{sim}_{\text{exact}}^{\text{or}}(r, q) = \max (\text{sim}(r, t_1), \text{sim}(r, t_2), \dots, \text{sim}(r, t_n))$$

$$\text{sim}_{\text{exact}}^{\text{and}}(r, q) = \min (\text{sim}(r, t_1), \text{sim}(r, t_2), \dots, \text{sim}(r, t_n))$$

- Extended boolean retrieval model for supporting exact and best match queries on the managed resources
- P-norm is used to compute distance

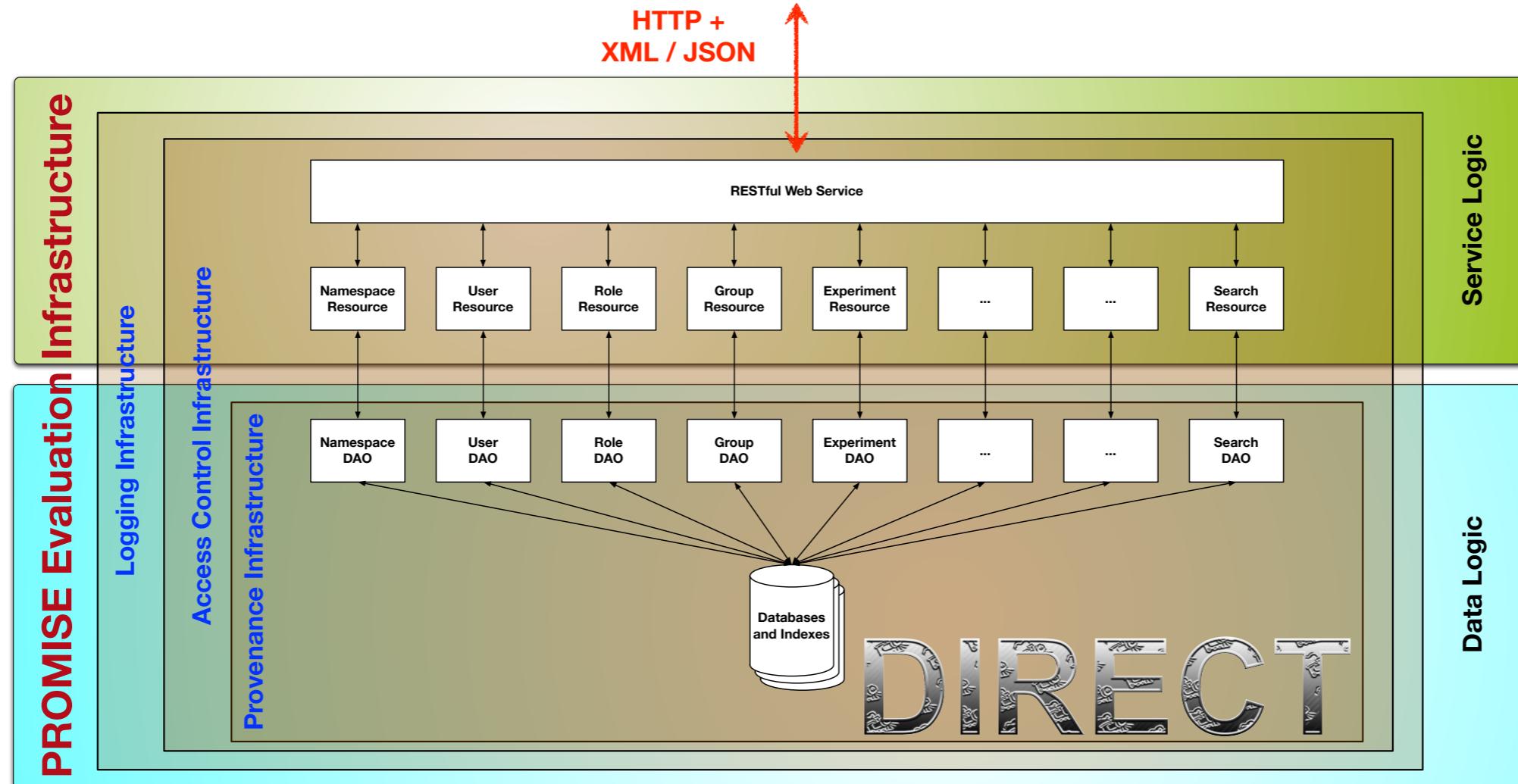


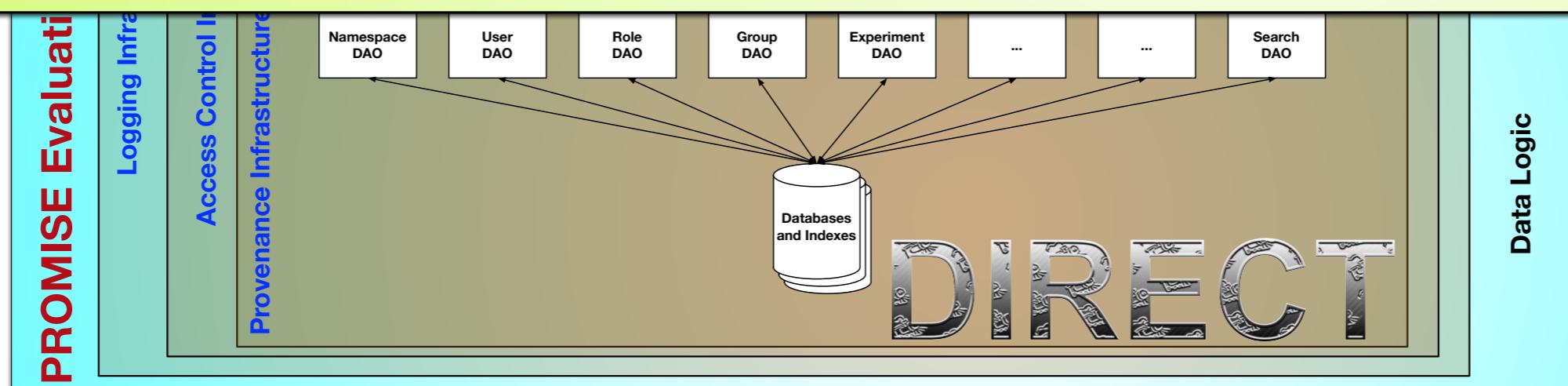
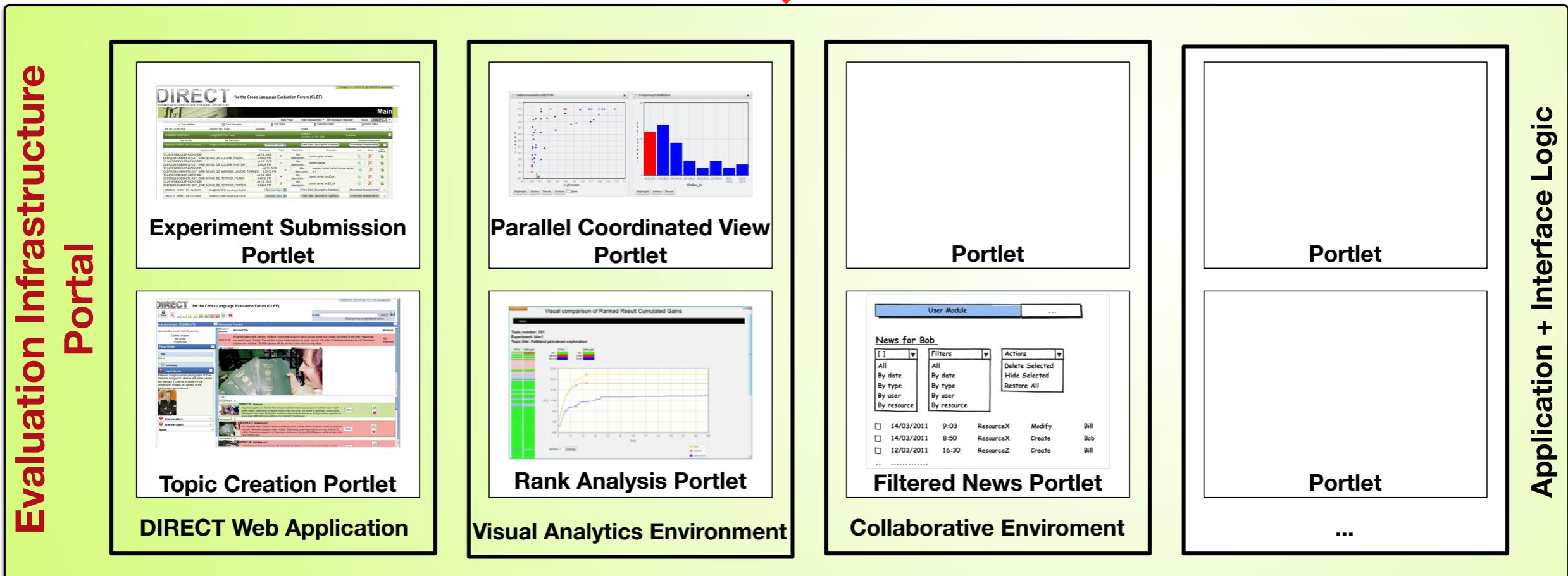
- Query language based on the CQL (Contextual Query Language) syntax
  - It allows for embedding queries in URL and linking to them

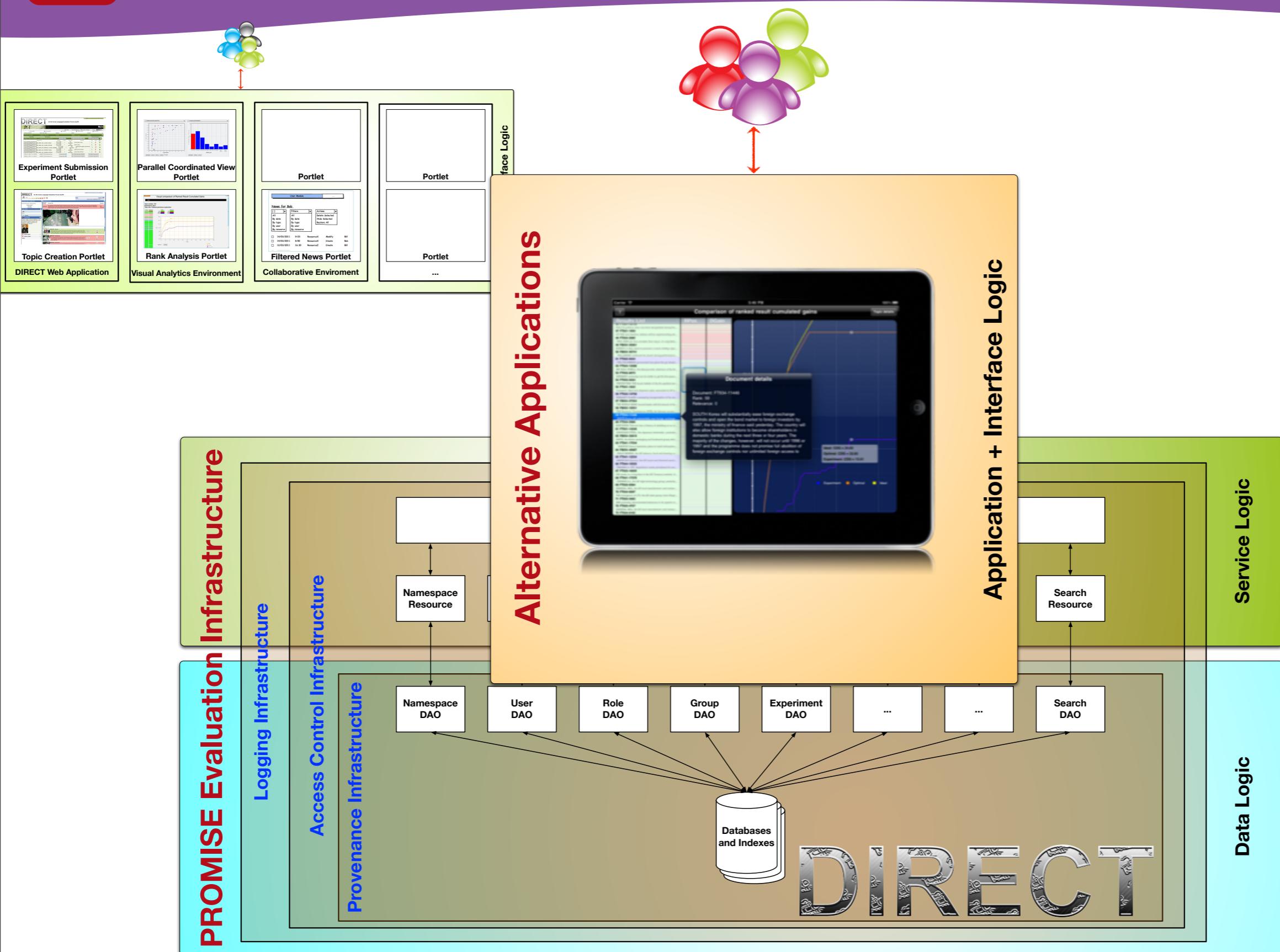
```
direct.experiment.description = "language model"  
and/match=exact
```

```
direct.experiment.user.id = unine
```

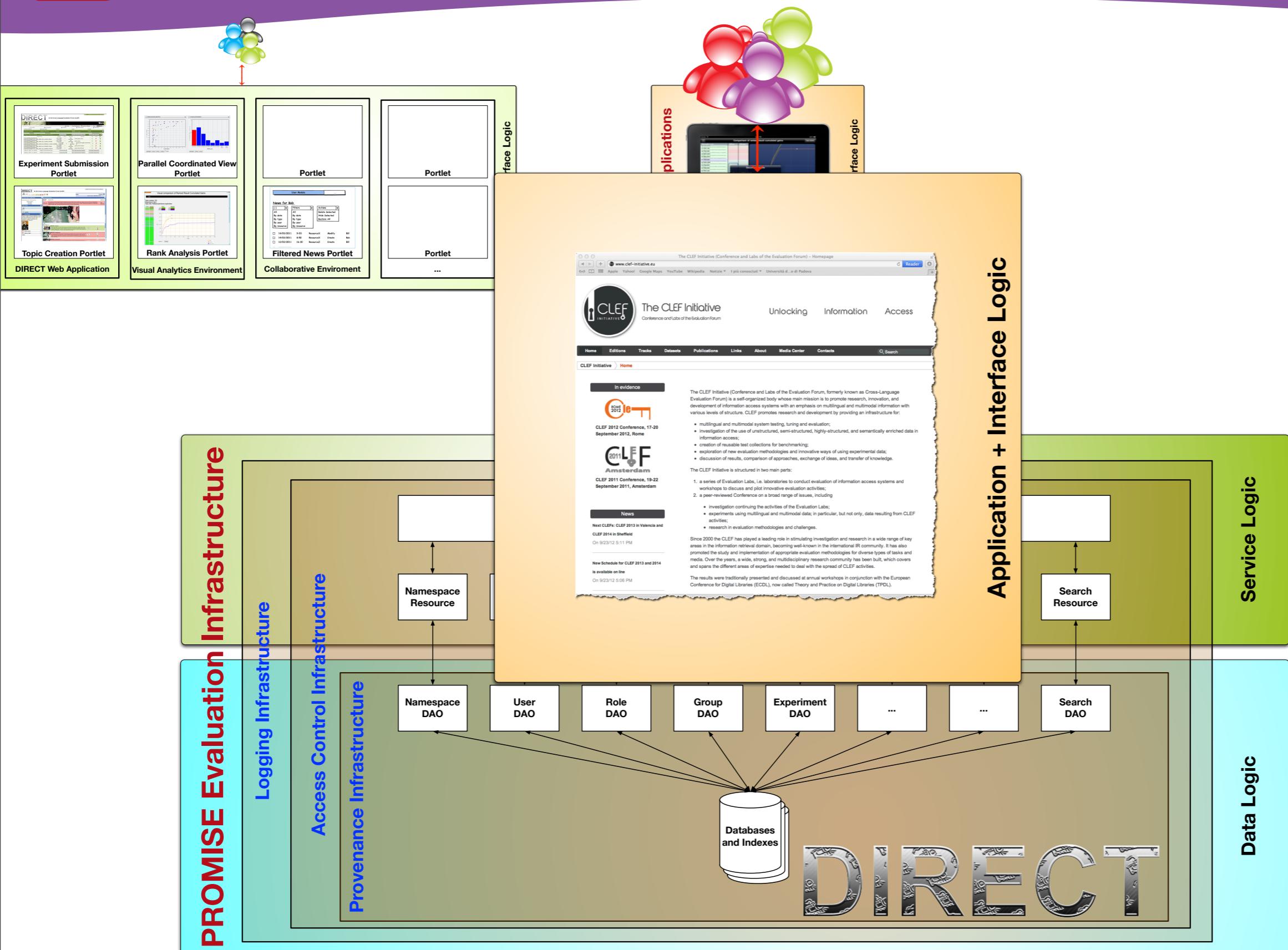
# Designing and Developing



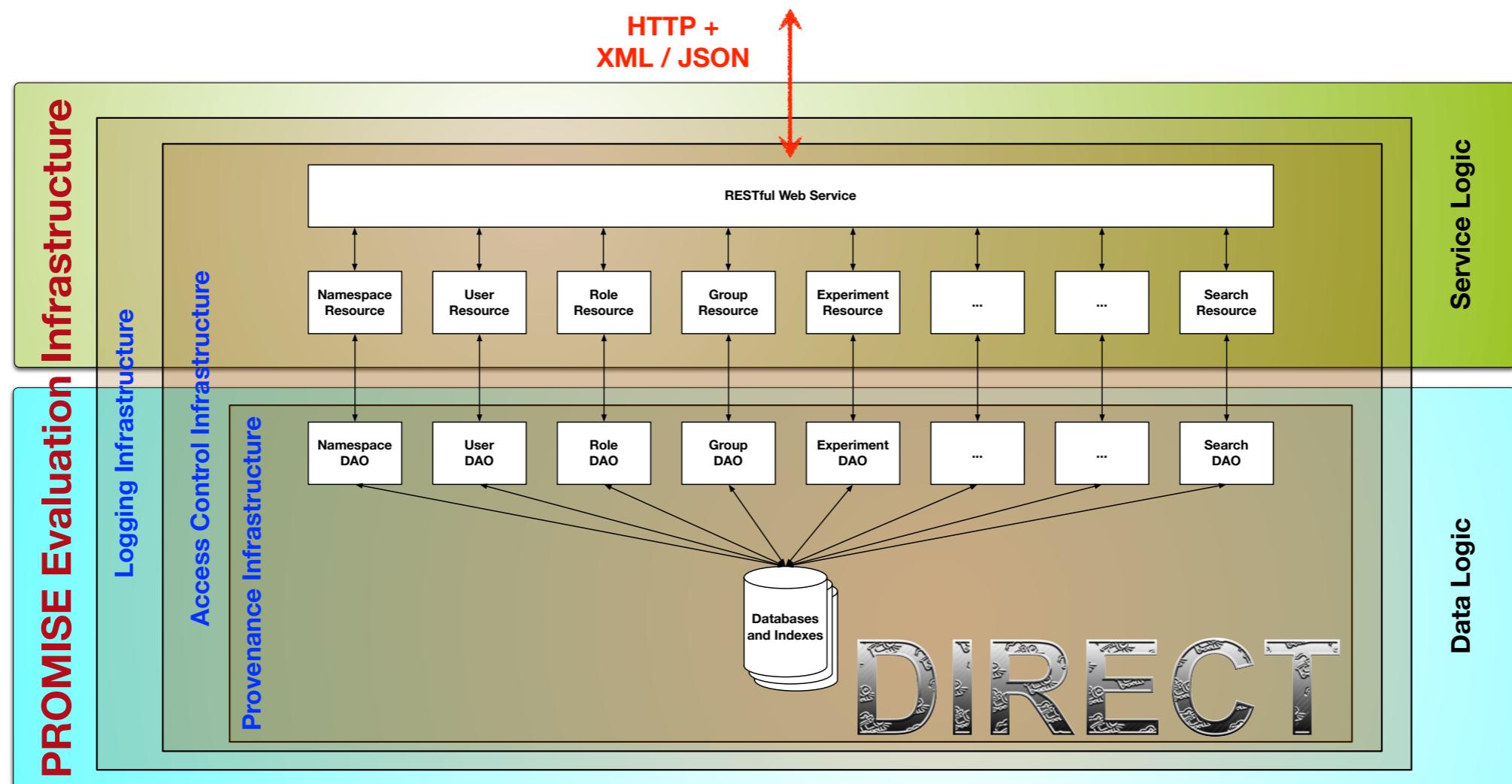




# Architecture



# Architecture



# The DIRECT Application

Logged as testassessor @ CLEFTEST [Logout]

Distributed Information Retrieval Evaluation Campaign Tool (DIRECT) – Experiment View

direct.del.unipd.it direct.dev.unipd.it direct.dev.unipd.it

Distributed Information Retrieval Evaluation Campaign Tool (DIRECT) – Portal Main Page

Reader

Monday 17 December 2012

Distributed Information Retrieval Evaluation Campaign Tool (DIRECT) – Main

Logged as clef@tu-chemnitz.de @ CLEF2012 [Logout] - Monday 17 December 2012

**DIRECT** for Evaluation of Multilingual Information Access in PROMISE NoE

Search

Experiment View

About DIRECT Aboutness

for Palestinian Not Relevant

Download as PNG

Download as PNG

**Portal Main Page**

Translation Manager

Main Page

DIRECT Portal

Campaigns

- + CLEF 2000
- + CLEF 2001
- + CLEF 2002
- + CLEF 2003
- + CLEF 2004
- + CLEF 2005
- + CLEF 2006
- + CLEF 2007
- + CLEF 2008
- CLEF 2009
- Tracks
- Ad-Hoc Persian Track
- Tasks
- + Ad-Hoc TEL Bilingual Persian Task
- Ad-Hoc TEL Monolingual Persian Task
- Download Pool
- + Download Topics
- + Ad-Hoc Persian Topic Creation Track
- + Ad-Hoc Robust Track
- + Ad-Hoc TEL Track
- + Ad-Hoc TEL Topic Creation Track
- + Grid@CLEF Pilot Track
- + ImageCLEF Track
- + CLEF Test

Select all Unselect all Download Selected

< prev (1 of 1) next >

Show 20 rows

Identifier	Participant	Description	Query Construction	Source Language	Is Pooled	View	Download
AH-PERSIAN-MONO-FA-CLEF2009.JHU-APL.JHUFA4R100TD	jhu-apl	4-grams; 100 rf terms	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.JHU-APL.JHUFA4R100TDN	jhu-apl	4-grams; 100 rf terms	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.JHU-APL.JHUFA5R100TD	jhu-apl	5-grams; RF 100 terms	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.JHU-APL.JHUFASK4R400TD	jhu-apl	skip 4-grams (1 skip); 400 RF terms	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.JHU-APL.JHUFATR5R50TD	jhu-apl	truncated word forms; 50 RF terms	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.OPENTEXT.OTFA09T	opentext	title-only, Neuchatel stemmer	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.OPENTEXT.OTFA09TD	opentext	same as otFA09t except both title and description fields used	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.OPENTEXT.OTFA09TDE	opentext	blind feedback based on top-3 rows of otFA09td	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.OPENTEXT.OTFA09TDZ	opentext	depth-10000 sampling run (orig ranks in first 5 decimal places of rsv)	AUTOMATIC	fa	true		
AH-PERSIAN-MONO-FA-CLEF2009.QAZVINIAU.IAUPERFA1	qazviniau	(Title+Desc), Stemmed Collection, Stemmed Queries, Query Stops, PRF(5,10)	AUTOMATIC	fa	true		

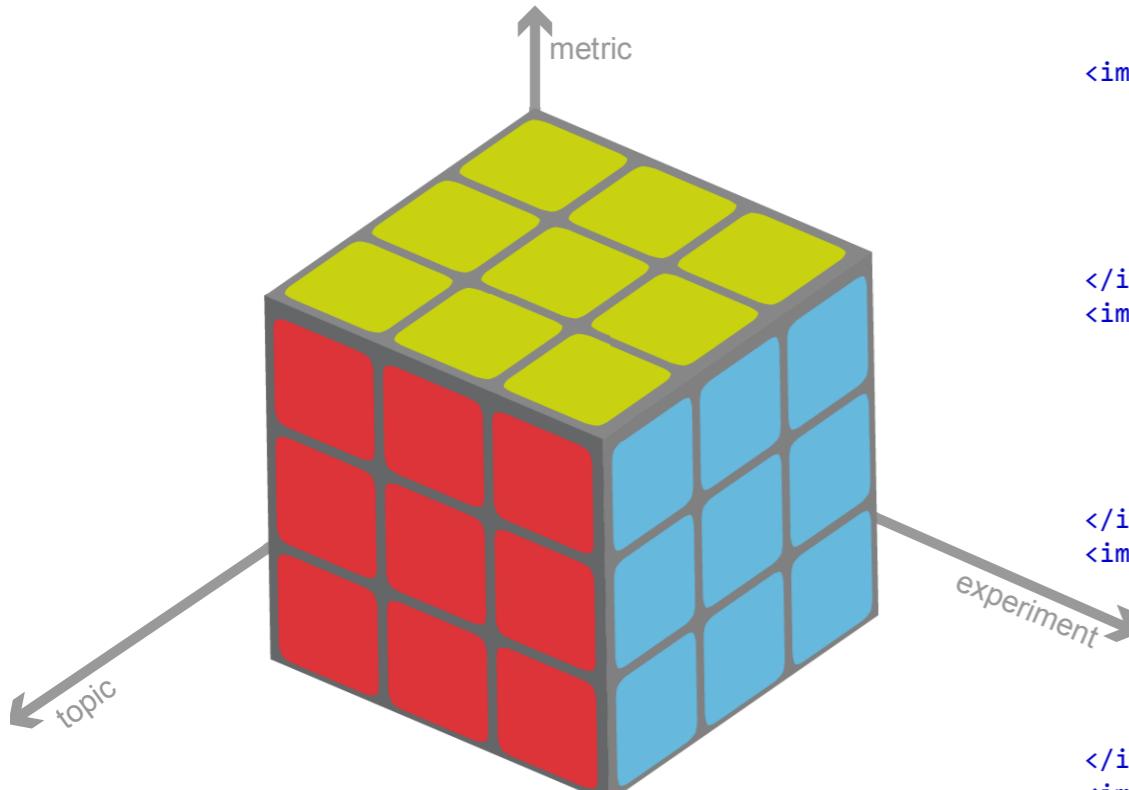
```

1 while ($ranking = shift) {
2   $rankname[$i] = $ranking;
3   my $runcount;
4
5   open(IN, $ranking) or die "Can't read $ranking: $!\n";
6   while (<IN>) {
7     [...]
8
9   $rank[$i]{$run} = $score;
10 }
11 close(IN);
12 }
13 }
14 #
15 #
16
17 for $i (0 .. $numrank - 2) {
18   for my $j ($i + 1 .. $numrank - 1) {
19     my $t = tau($i, $j);
20     print "$t\n";
21     $sum += $t;
22
23   [...]
24 }
25 }
26
27 sub tau {
28   my ($i, $j) = @_;
29   my ($con, $dis, $tie_i, $tie_j, $n, @pairs);
30   my (@runs, %runstat);
31
32   for my $run (sort { $rank[$i]{$b} <=> $rank[$i]{$a} } keys %{$rank[$i]}) {
33     push @pairs, [ $rank[$i]{$run}, $rank[$j]{$run} ];
34     push @runs, $run;
35   }
36
37
38   [...]
39
40   return $t;
41 }
42
43
44
45
46
47
1 while ($ranking = shift) {
2   $rankname[$i] = $ranking;
3   my $runcount;
4
5   my $ua = LWP::UserAgent->new();
6   my $response = $ua->get(sprintf(
7     'http://direct.dei.unipd.it/run/' . $ranking . '/items',
8     uri_escape($text),
9     uri_escape($key),
10    ));
11
12   while (<response>) {
13     [...]
14   }
15   $rank[$i]{$run} = $score;
16 }
17
18
19 }
20
21 #
22
23 for $i (0 .. $numrank - 2) {
24   for my $j ($i + 1 .. $numrank - 1) {
25     my $t = tau($i, $j);
26     print "$t\n";
27     $sum += $t;
28
29   [...]
30 }
31 }
32
33 sub tau {
34   my ($i, $j) = @_;
35   my ($con, $dis, $tie_i, $tie_j, $n, @pairs);
36   my (@runs, %runstat);
37
38   for my $run (sort { $rank[$i]{$b} <=> $rank[$i]{$a} } keys %{$rank[$i]}) {
39     push @pairs, [ $rank[$i]{$run}, $rank[$j]{$run} ];
40     push @runs, $run;
41   }
42
43
44   [...]
45
46   return $t;
47 }

```

The screenshot shows a web browser window displaying the CLEF Initiative portal. The title bar reads "The CLEF Initiative (Conference and Labs of the Evaluation Forum) - CLEF2012 Working Notes". The page header includes the CLEF logo, the text "The CLEF Initiative Conference and Labs of the Evaluation Forum", and the words "Unlocking Information Access". The navigation menu at the top includes links for Home, Editions, Tracks, Datasets, Publications, Links, About, Media Center, and Contacts, along with a search bar. A breadcrumb navigation path shows "CLEF Initiative > Editions > CLEF2012 > Working Notes". On the left, a sidebar for "CLEF2012" lists "Agenda" and "Working Notes" (which is currently selected). The main content area is titled "CLEF 2012 Working Notes" and describes working notes for the CLEF 2012 Conference held in Rome, Italy. It includes a link to view papers for all tracks. Below this, sections include "Contents", "Cultural Heritage in CLEF (CHiC)", "Cultural Heritage in CLEF (CHiC) Overview 2012" (with authors Vivien Petras, Nicola Ferro, Maria Gädé, Antoine Isaac, Michael Kleineberg, Ivano Masiero, Mattia Nicchio, and Juliane Stiller and a link to the paper), "The Sheffield and Basque Country Universities Entry to CHiC: Using Random Walks and Similarity to Access Cultural Heritage" (with authors Eneko Agirre, Paul Clough, Samuel Fernando, Mark Hall, Aranba Otegi, and Mark Stevenson and a link to the paper), "UniNE at CLEF 2012" (with authors Mitra Akasereh, Nada Naji, and Jacques Savoy and a link to the paper), "Chemnitz at the CHiC Evaluation Lab 2012: Creating an Xtrieval Module for Semantic Enrichment" (with authors Jens Kürsten, Thomas Wilhelm, Daniel Richter, and Maximilian Eibl and a link to the paper), "Dealing with Sparse Document and Topic Representations: Lab Report for CHiC 2012" (with authors Philipp Schaefer, Daniel Hienert, Frank Sawitzki, Andras Wira-Alam, and Thomas Lüke and a link to the paper), "Query Expansion Using Wikipedia and Dbpedia" (with authors Nitish Aggarwal and Paul Buitelaar and a link to the paper), "Retrieval in the Intellectual Property Domain (CLEF-IP)" (with a link to the paper), and "Image Retrieval in CLEF (ImageCLEF)" (with a link to the paper).

## GET /task/CHIC-AH-MONO-EN-CLEF2012/measure

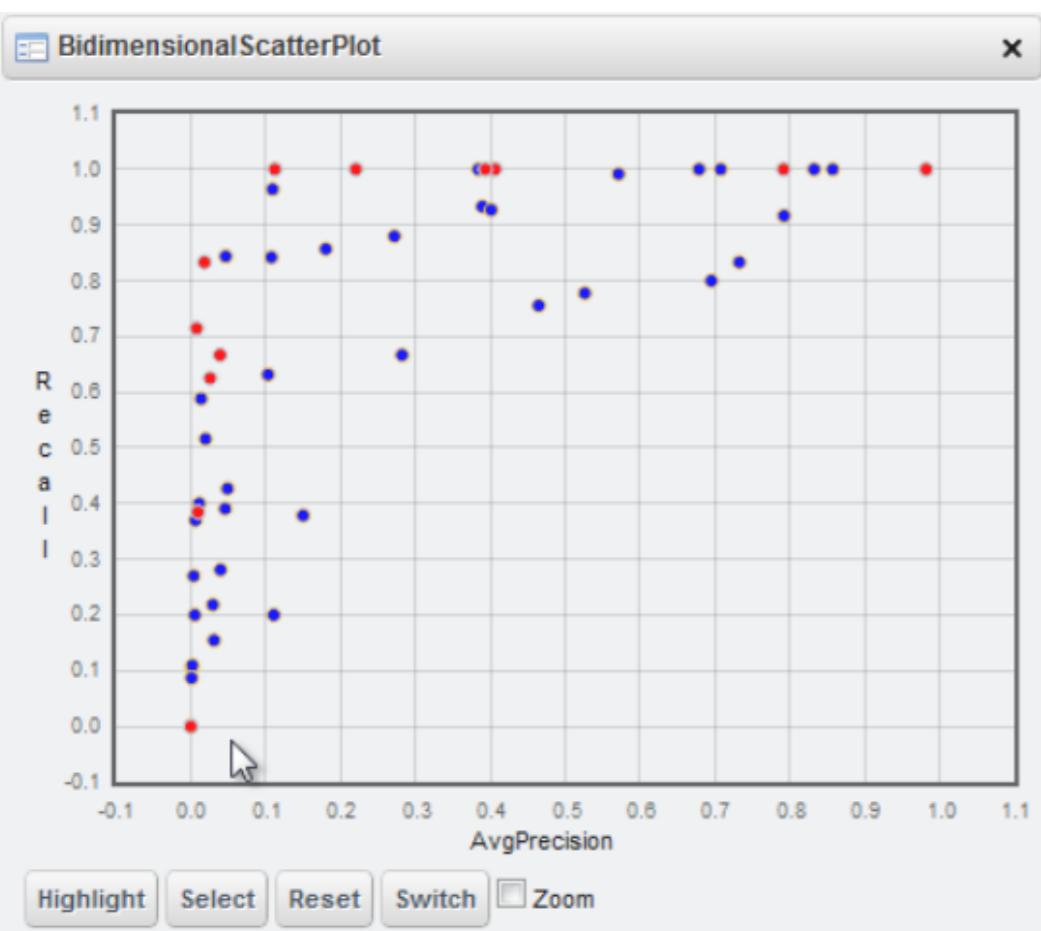


```

<ims:measure ims:identifier="177bcef2-00a0-4f59-b781-f285610f1c6f" ims:value="3.59523803E-1"
    ims:created="2012-10-11T16:23:31.011+02:00" ims:last-modified="2012-10-11T16:23:31.011+02:00">
    <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="00681943-2249-4a56-9c1e-5dbf82dcd1b5" ims:value="0.0E0"
    ims:created="2012-10-11T16:23:30.698+02:00" ims:last-modified="2012-10-11T16:23:30.698+02:00">
    <ims:concept ims:identifier="R_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="0339409c-2baa-44e1-b977-48e4d79dc0ed" ims:value="5.0E-1"
    ims:created="2012-10-11T16:23:30.815+02:00" ims:last-modified="2012-10-11T16:23:30.815+02:00">
    <ims:concept ims:identifier="PRECISION_AT_0_INTERPOLATED_RECALL_LEVEL"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="006df09d-157a-466c-bf2e-c0b3b10b6bc0" ims:value="4.4742262363E-1"
    ims:created="2012-10-11T16:23:30.882+02:00" ims:last-modified="2012-10-11T16:23:30.882+02:00">
    <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>
<ims:measure ims:identifier="164ef66a-aa48-413f-a229-540bb19f2d2e" ims:value="5.9090906382E-1"
    ims:created="2012-10-11T16:23:30.760+02:00" ims:last-modified="2012-10-11T16:23:30.760+02:00">
    <ims:concept ims:identifier="R_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>
<ims:measure ims:identifier="29a91449-61ae-49ef-96dc-2ae67377ef86" ims:value="1.0E0"
    ims:created="2012-10-11T16:23:30.594+02:00" ims:last-modified="2012-10-11T16:23:30.594+02:00">
    <ims:concept ims:identifier="PRECISION_AT_0_INTERPOLATED_RECALL_LEVEL"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>

```

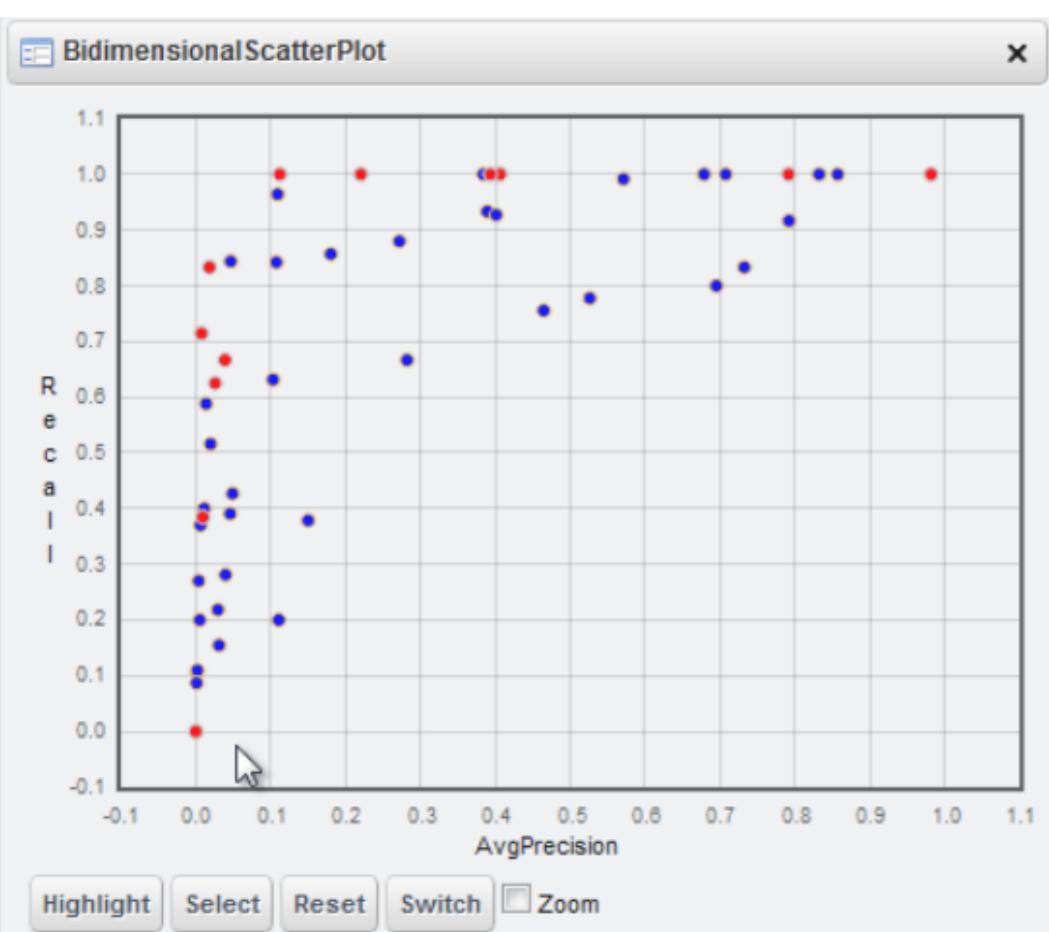
# Visualizations



```

<ims:measure ims:identifier="177bcef2-00a0-4f59-b781-f285610f1c6f" ims:value="3.59523803E-1"
    ims:created="2012-10-11T16:23:31.011+02:00" ims:last-modified="2012-10-11T16:23:31.011+02:00">
    <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="00681943-2249-4a56-9c1e-5dbf82dcd1b5" ims:value="0.0E0"
    ims:created="2012-10-11T16:23:30.698+02:00" ims:last-modified="2012-10-11T16:23:30.698+02:00">
    <ims:concept ims:identifier="R_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="0339409c-2baa-44e1-b977-48e4d79dc0ed" ims:value="5.0E-1"
    ims:created="2012-10-11T16:23:30.815+02:00" ims:last-modified="2012-10-11T16:23:30.815+02:00">
    <ims:concept ims:identifier="PRECISION_AT_0_INTERPOLATED_RECALL_LEVEL"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-001" />
</ims:measure>
<ims:measure ims:identifier="006df09d-157a-466c-bf2e-c0b3b10b6bc0" ims:value="4.4742262363E-1"
    ims:created="2012-10-11T16:23:30.882+02:00" ims:last-modified="2012-10-11T16:23:30.882+02:00">
    <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>
<ims:measure ims:identifier="164ef66a-aa48-413f-a229-540bb19f2d2e" ims:value="5.9090906382E-1"
    ims:created="2012-10-11T16:23:30.760+02:00" ims:last-modified="2012-10-11T16:23:30.760+02:00">
    <ims:concept ims:identifier="R_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>
<ims:measure ims:identifier="29a91449-61ae-49ef-96dc-2ae67377ef86" ims:value="1.0E0"
    ims:created="2012-10-11T16:23:30.594+02:00" ims:last-modified="2012-10-11T16:23:30.594+02:00">
    <ims:concept ims:identifier="PRECISION_AT_0_INTERPOLATED_RECALL_LEVEL"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
    <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
    <ims:topic ims:identifier="CHIC-002" />
</ims:measure>
```

## POST /visualization/



```

<ims:visualization ims:identifier="3dedf5b5-5026-4054-9fc6-827512791960" ims:scope="PUBLIC">
  <ims:concept ims:identifier="SCATTERPLOT"
    ims:namespace="http://direct.dei.unipd.it/visualizations/" />
  <ims:parameters>
    <ims:parameter>
      <ims:concept ims:identifier="X_AXIS"
        ims:namespace="http://direct.dei.unipd.it/visualizations/" />
      <ims:value>
        <ims:concept ims:identifier="AVERAGE_PRECISION"
          ims:namespace="http://direct.dei.unipd.it/metrics/" />
      </ims:value>
    </ims:parameter>
    <ims:parameter>
      <ims:concept ims:identifier="Y_AXIS"
        ims:namespace="http://direct.dei.unipd.it/visualizations/" />
      <ims:value>
        <ims:concept ims:identifier="RECALL"
          ims:namespace="http://direct.dei.unipd.it/metrics/" />
      </ims:value>
    </ims:parameter>
    <ims:parameter>
      <ims:concept ims:identifier="HIGHLIGHT"
        ims:namespace="http://direct.dei.unipd.it/visualizations/" />
      <ims:value>
        <ims:measure ims:identifier="177bcef2-00a0-4f59-b781-f285610f1c6f"/>
        <ims:measure ims:identifier="006df09d-157a-466c-bf2e-c0b3b10b6bc0"/>
      </ims:value>
    </ims:parameter>
  </ims:parameters>
  <ims:measures>
    <ims:measure ims:identifier="177bcef2-00a0-4f59-b781-f285610f1c6f" ims:value="3.59523803E-1" >
      <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
      <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
      <ims:topic ims:identifier="CHIC-001" />
    </ims:measure>
    <ims:measure ims:identifier="006df09d-157a-466c-bf2e-c0b3b10b6bc0" ims:value="4.4742262363E-1">
      <ims:concept ims:identifier="AVERAGE_PRECISION"
        ims:namespace="http://direct.dei.unipd.it/metrics/" />
      <ims:experiment ims:identifier="CHIC-AH-MONO-EN-CLEF2012.ARANTZA.OTEGI@EHU.ES.EXP_UKB_WN100" />
      <ims:topic ims:identifier="CHIC-002" />
    </ims:measure>
  </ims:measures>
  <ims:snapshots>
    <ims:snapshot ims:identifier="017c333a-4b7c-4267-926d-f15fe3554efd" ims:media-type="application/pdf" >
      <ims:content ims:content-transfer-encoding="base64">
        c2lnbiBjb250ZW50
      </ims:content>
    </ims:snapshot>
  </ims:snapshots>
</ims:visualization>

```

CHIC2012-Overview-Paper-updated-2012-09-03.htm

## Cultural Heritage in CLEF (CHiC) Overview 2012

Vivien Petras, Nicola Ferro, Maria Gäde, Antoine Isaac, Michael Kleineberg, Ivano Masiero, Mattia Nicchio and Juliane Stiller

### 1 Introduction

Cultural heritage content is often multilingual and multimedia (e.g. text, photographs, images, audio recordings, and videos), usually described with metadata in multiple formats and of different levels of complexity....

### 2 Results Analysis

**Table 1.** Best monolingual Experiments and Performance Difference between best and last (up to 5) Experiment (in MAP)

Track	Rank	Part.	Experiment Identifier	MAP
<b>Monolingual English</b>	1 <sup>st</sup>	<b>UPV</b>	<a href="#">EXP_UKB_WN100</a>	<b>51.61%</b>
	2 <sup>nd</sup>	Chemnitz	QE0X20NO	48.60%
	3 <sup>rd</sup>	Neuchatel	UNINEENEN1	44.87%
	4 <sup>th</sup>	Gesis	GESIS_WIKI_ENTITY_EN_EN	43.96%
	5 <sup>th</sup>	Berkeley	MONO_EN_TD_T2FB	36.40%
	Dif.			

Figures 7 to 9 show the interpolated recall vs. average precision for the top groups of the monolingual tasks.

CHiC Ad-Hoc Monolingual English Task Top 5 Participants - Standard Recall Levels vs Mean Interpolated Precision

Standard Recall Level	UPV (Precision)	Chemnitz (Precision)	Neuchatel (Precision)	Gesis (Precision)	Berkeley (Precision)
0.1	51.61%	48.60%	44.87%	43.96%	36.40%
0.2	50.00%	47.00%	43.00%	41.00%	35.00%
0.3	48.00%	45.00%	41.00%	39.00%	33.00%
0.4	46.00%	43.00%	39.00%	37.00%	31.00%
0.5	44.00%	41.00%	37.00%	35.00%	29.00%
0.6	42.00%	39.00%	35.00%	33.00%	27.00%
0.7	40.00%	37.00%	33.00%	31.00%	25.00%
0.8	38.00%	35.00%	31.00%	29.00%	23.00%
0.9	36.00%	33.00%	29.00%	27.00%	21.00%
1.0	34.00%	31.00%	27.00%	25.00%	20.00%

< a href="http://direct.dei.unipd.it/user/UPV">UPV</a>

CHIC2012-Overview-Paper-updated-2012-09-03.htm

## Cultural Heritage in CLEF (CHiC) Overview 2012

Vivien Petras, Nicola Ferro, Maria Gäde, Antoine Isaac, Michael Kleineberg, Ivano Masiero, Mattia Nicchio and Juliane Stiller

### 1 Introduction

Cultural heritage content is often multilingual and multimedia (e.g. text, photographs, images, audio recordings, and videos), usually described with metadata in multiple formats and of different levels of complexity....

### 2 Results Analysis

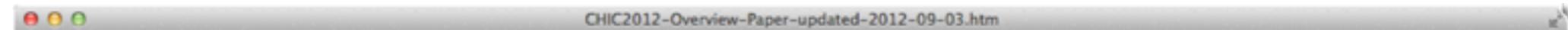
**Table 1.** Best monolingual Experiments and Performance Difference between best and last (up to 5) Experiment (in MAP)

Track	Rank	Part.	Experiment Identifier	MAP
<b>Monolingual English</b>	1 <sup>st</sup>	<b>UPV</b>	<a href="#">EXP_UKB_WN100</a>	<b>51.61%</b>
	2 <sup>nd</sup>	Chemnitz	QEWX20NO	48.60%
	3 <sup>rd</sup>	Neuchatel	UNINEENEN1	44.87%
	4 <sup>th</sup>	Gesis	GESIS_WIKI_ENTITY_EN_EN	43.96%
	5 <sup>th</sup>	Berkeley	MONO_EN_TD_T2FB	36.40%
	Diff.			

Figures 7 to 9 show the interpolated recall vs. average precision for the top groups of the monolingual tasks.

**CHiC Ad-Hoc Monolingual English Task Top 5 Participants - Standard Recall Levels vs Mean Interpolated Precision**

[EXP\\_UKB\\_WN100](http://direct.dei.unipd.it/experiment/EXP_UKB_WN100)



Vivien Petras, Nicola Ferro, Maria Gäde, Antoine Isaac, Michael Kleineberg, Ivano Masiero, Mattia Nicchio and Juliane Stiller

## 1 Introduction

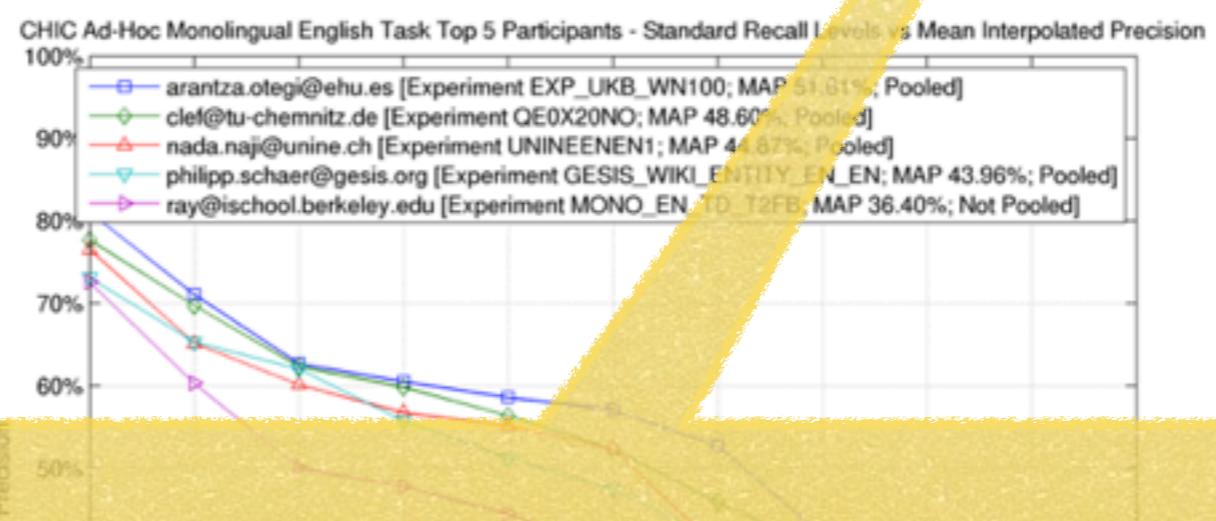
Cultural heritage content is often multilingual and multimedia (e.g. text, photographs, images, audio recordings, and videos), usually described with metadata in multiple formats and of different levels of complexity....

## 2 Results Analysis

**Table 1.** Best monolingual Experiments and Performance Difference between best and last (up to 5) Experiment (in MAP)

Track	Rank	Part.	Experiment Identifier	MAP
<b>Monolingual English</b>	1 <sup>st</sup>	<a href="#">UPV</a>	<a href="#">EXP_UKB_WN100</a>	<b>51.61%</b>
	2 <sup>nd</sup>	Chemnitz	QE0X20NO	48.60%
	3 <sup>rd</sup>	Neuchatel	UNINEENEN1	44.87%
	4 <sup>th</sup>	Gesis	GESIS_WIKI_ENTITY_EN_EN	43.96%
	5 <sup>th</sup>	Berkeley	MONO_EN_TD_T2FB	36.40%
	Diff.			41.78%

Figures 7 to 9 show the interpolated recall vs. average precision for the top groups of the monolingual tasks.



<[51.61%](http://direct.dei.unipd.it/estimate/017c333a-4b7c-4267-926d-f15fe3554efd)>



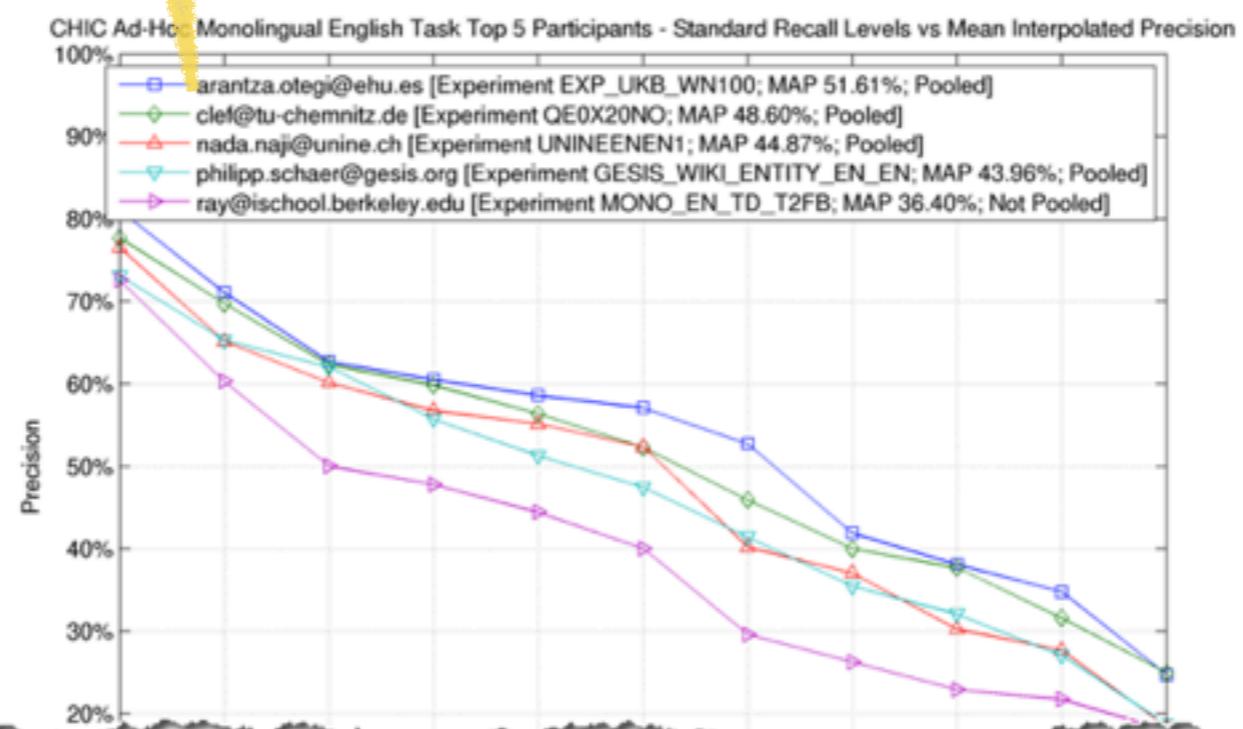
Cultural heritage content is often multilingual and multimedia (e.g. text, photographs, images, audio recordings, and videos), usually described with metadata in multiple formats and of different levels of complexity...

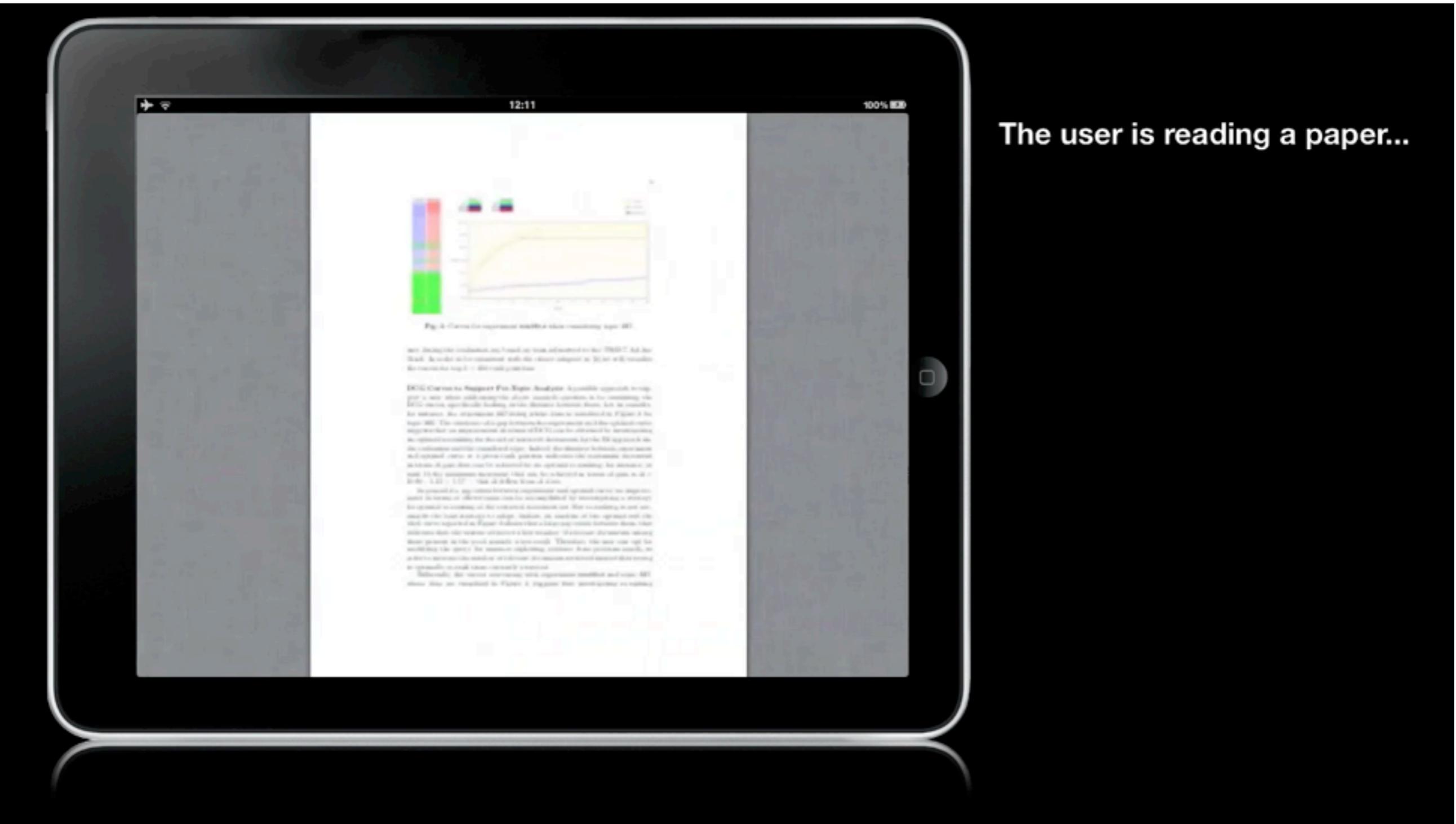
## 2 Results Analysis

**Table 1.** Best monolingual Experiments and Performance Difference between best and last (up to 5) Experiment (in MAP)

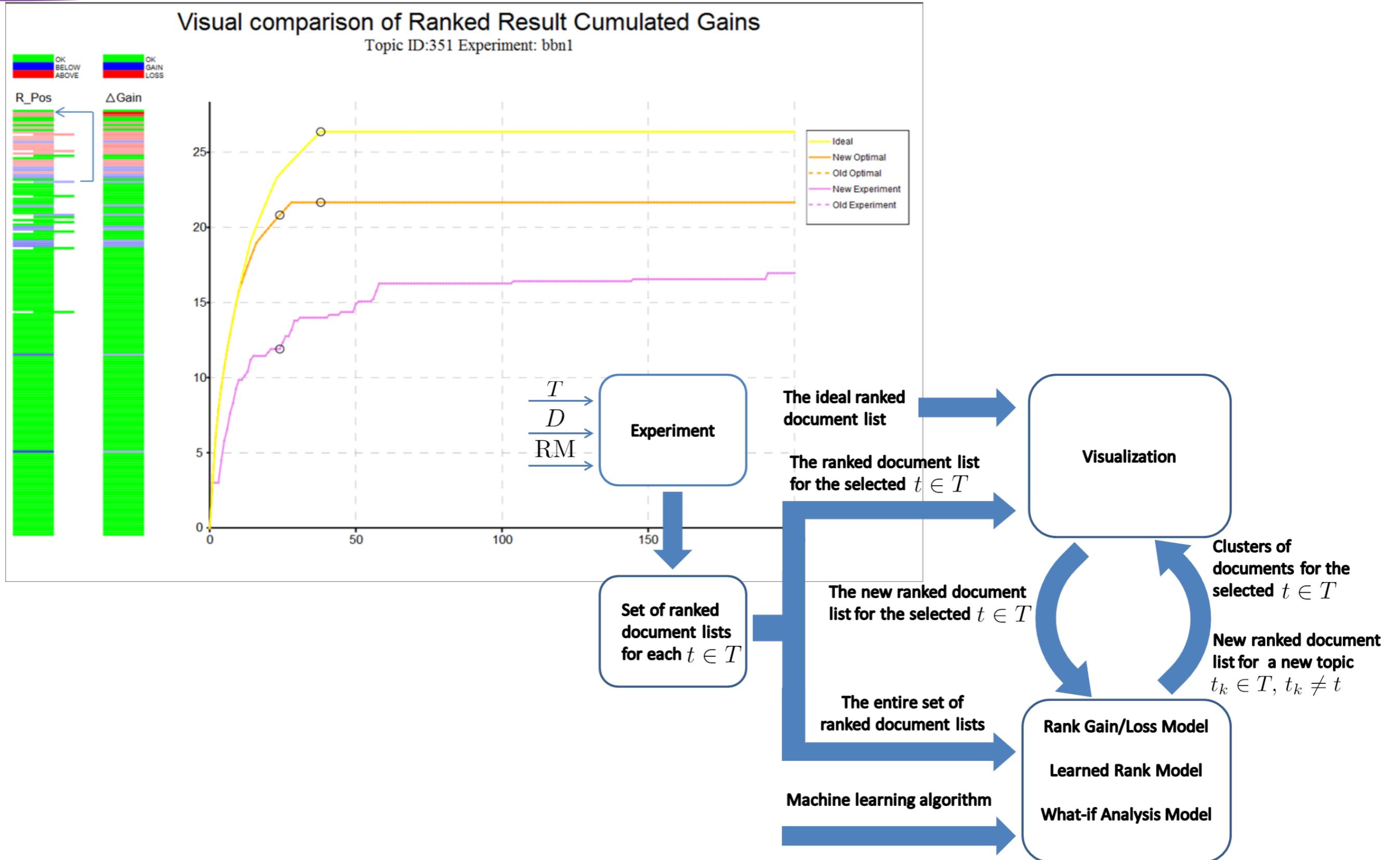
Track	Rank	Part.	Experiment Identifier	MAP
Monolingual English	1 <sup>st</sup>	UPV	EXP_UKB_WN100	51.61%
	2 <sup>nd</sup>	Chemnitz	QE0X20NO	48.60%
	3 <sup>rd</sup>	Neuchatel	UNINEENEN1	44.87%
	4 <sup>th</sup>	Gesis	GESIS_WIKI_ENTITY_EN_EN	43.96%
	5 <sup>th</sup>	Berkeley	MONO_EN_TD_T2FB	36.40%
	Diff.			41.78%

Figures 7 to 9 show the interpolated recall vs. average precision for the top groups of the monolingual tasks.



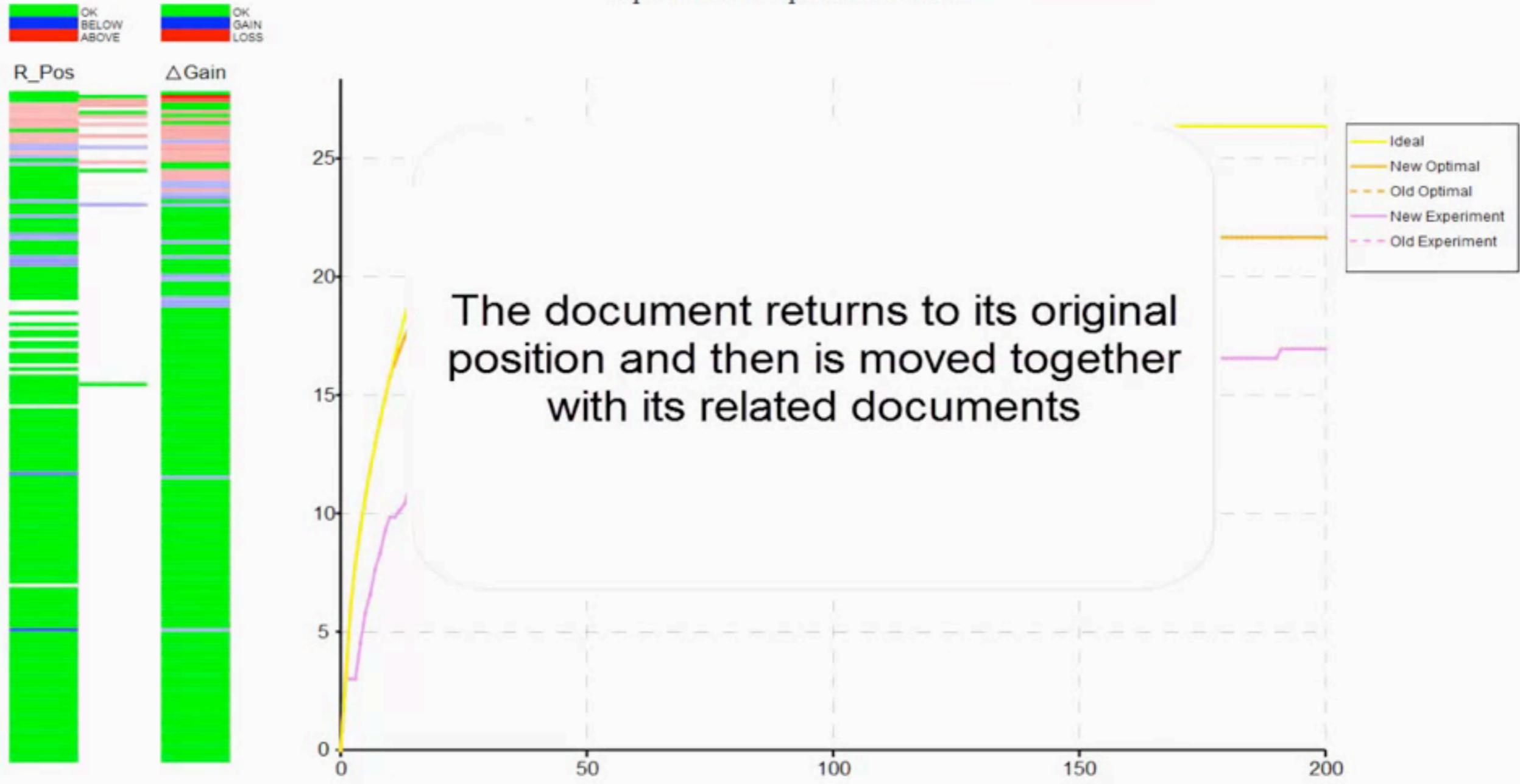


The user is reading a paper...



## Visual comparison of Ranked Result Cumulated Gains

Topic ID:351 Experiment: bbn1



# User Annotations

**User not logged into Direct**

**DIRECT** for the Cross Language Evaluation Forum (CLEF)  
DISTRIBUTED INFORMATION RETRIEVAL EVALUATION CAMPAIGN TOOL

**Portal Main Page**

Username  Password  Login

Campaigns

- + CLEF 2000
- + CLEF 2001
- + CLEF 2002
- + CLEF 2003
- + CLEF 2004
- + CLEF 2005
- + CLEF 2006
- + CLEF 2007
- Tracks
- Ad-Hoc Track
- + Tasks
- + Ad-Hoc Bilingual Bulgarian Task
- + Ad-Hoc Bilingual Czech Task
- + Ad-Hoc Bilingual English Task
- Download Pool
- Download Topics
- am
- bn
- bg
- es
- zh
- cs
- fr
- hi
- hu
- id
- it
- mr
- om
- ta
- te
- + Ad-Hoc Bilingual Hungarian Task
- + Ad-Hoc Monolingual

Select all  Unselect all  Download Selected

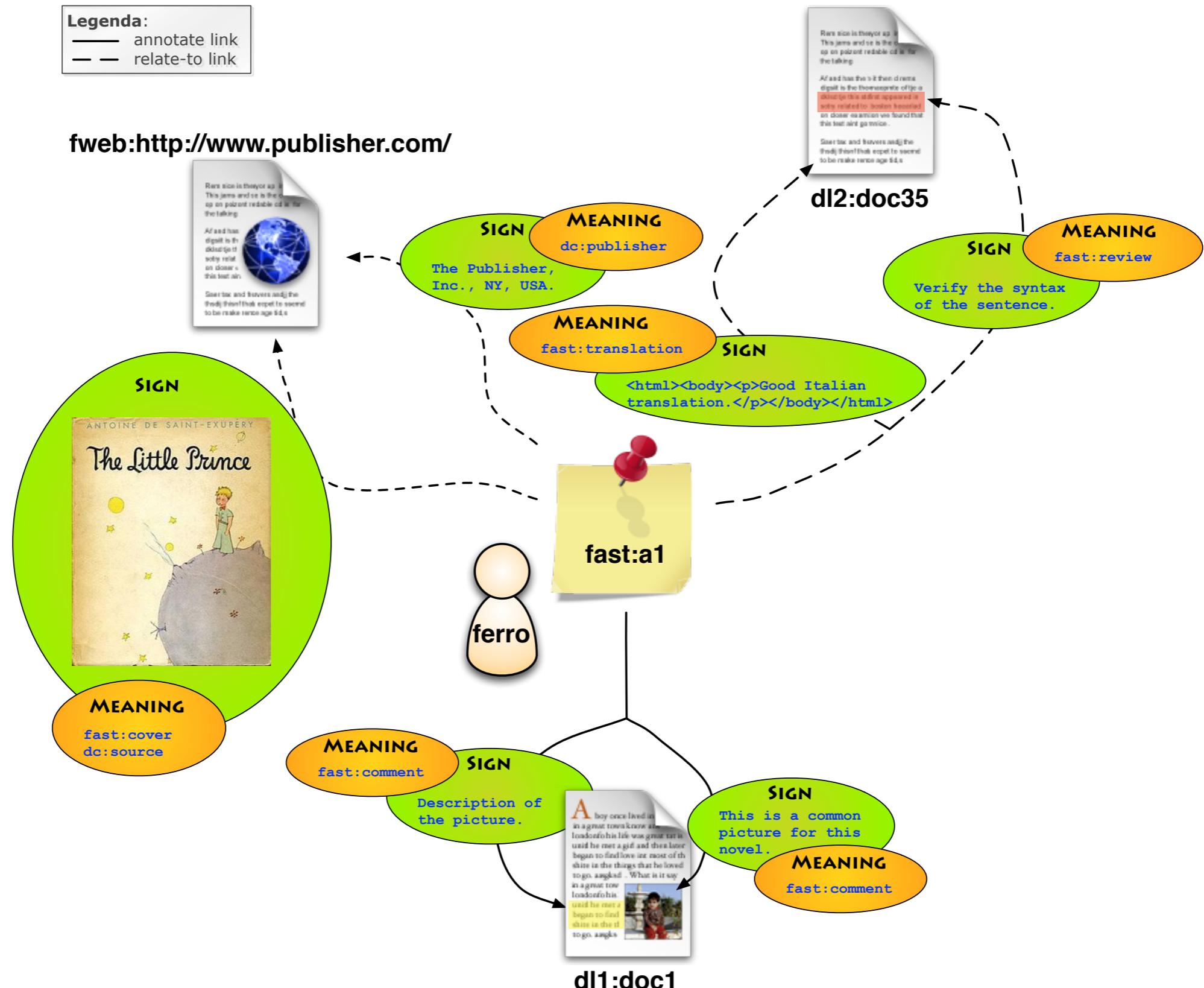
Identifier	Parameter	Dice Coefficient for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_HINDI_TITLEDESC_DICE	bombay-ltrc	Point-wise Mutual Info (PMI) for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_HINDI_TITLEDESC_PMI	bombay-ltrc	Point-wise Mutual Info (PMI) for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_HINDI_TITLE_DICE	bombay-ltrc	Dice Coefficient for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_HINDI_TITLE_PMI	bombay-ltrc	Point-wise Mutual Info (PMI) for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_MAR_TITLE_DICE	bombay-ltrc	Dice Coefficient for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BOMBAY-LTRC.IITB_MAR_TITLE_PMI	bombay-ltrc	Point-wise Mutual Info (PMI) for Similarity	AUTOMATIC	hi	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_1	budapest-acad	simple dictionary, best supposed params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_2	budapest-acad	simple dictionary, 2nd best	AUTOMATIC	hu	true		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_3	budapest-acad	simple dictionary, 3rd best	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_4	budapest-acad	simple dictionary, 4th best	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_5	budapest-acad	simple dictionary, 5th best	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.ACAD.BILING_6	budapest-acad	simple dictionary, 6th best	AUTOMATIC	hu	false		
	budapest-acad	using wikipedia, best supposed params	AUTOMATIC	hu	true		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.BILING_WIKI2	budapest-acad	with Wikipedia, 2nd best params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.BILING_WIKI3	budapest-acad	with Wikipedia, 3rd best params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.BILING_WIKI4	budapest-acad	with Wikipedia, 4nd best params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.BILING_WIKI5	budapest-acad	with Wikipedia, 5th best params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-CLEF2007.BUDAPEST-ACAD.BILING_WIKI6	budapest-acad	with Wikipedia, 6th best params	AUTOMATIC	hu	false		
10.2415/AH-BILI-X2EN-		Transitive translation for title and description query using German as pivot language and using #svn	AUTOMATIC	..	..		

**This experiment is interesting**

**This experiment should have been pooled: it brings a lot of new relevant docs**

**Check the spelling of this topic**

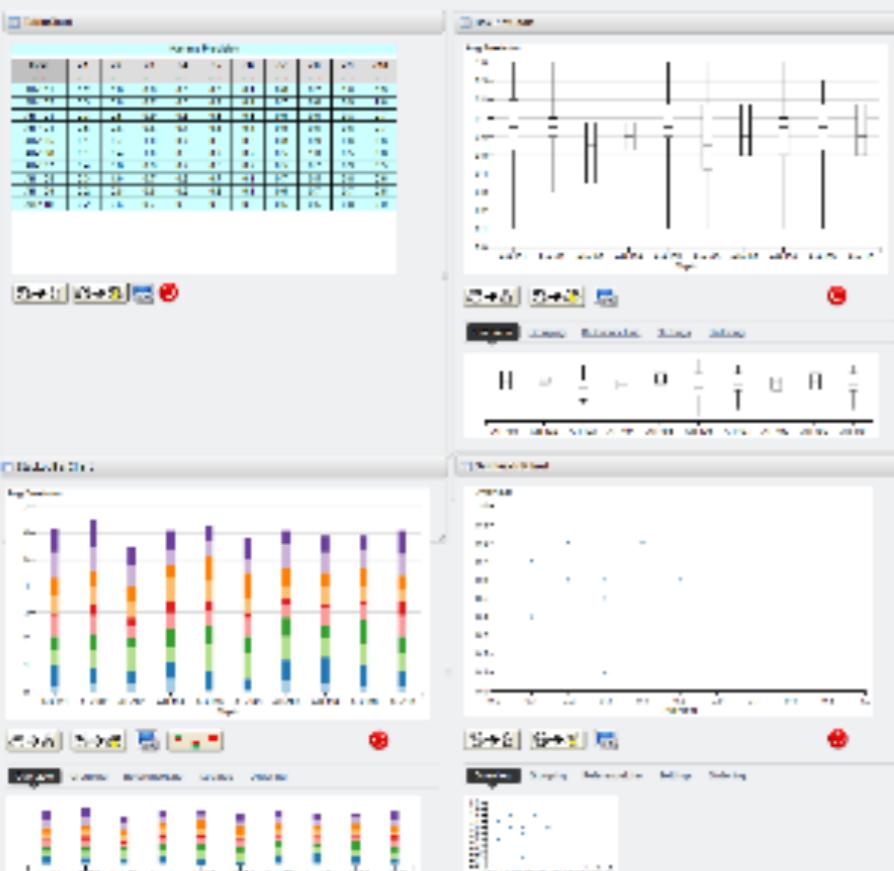
# User Annotations



# User Annotations



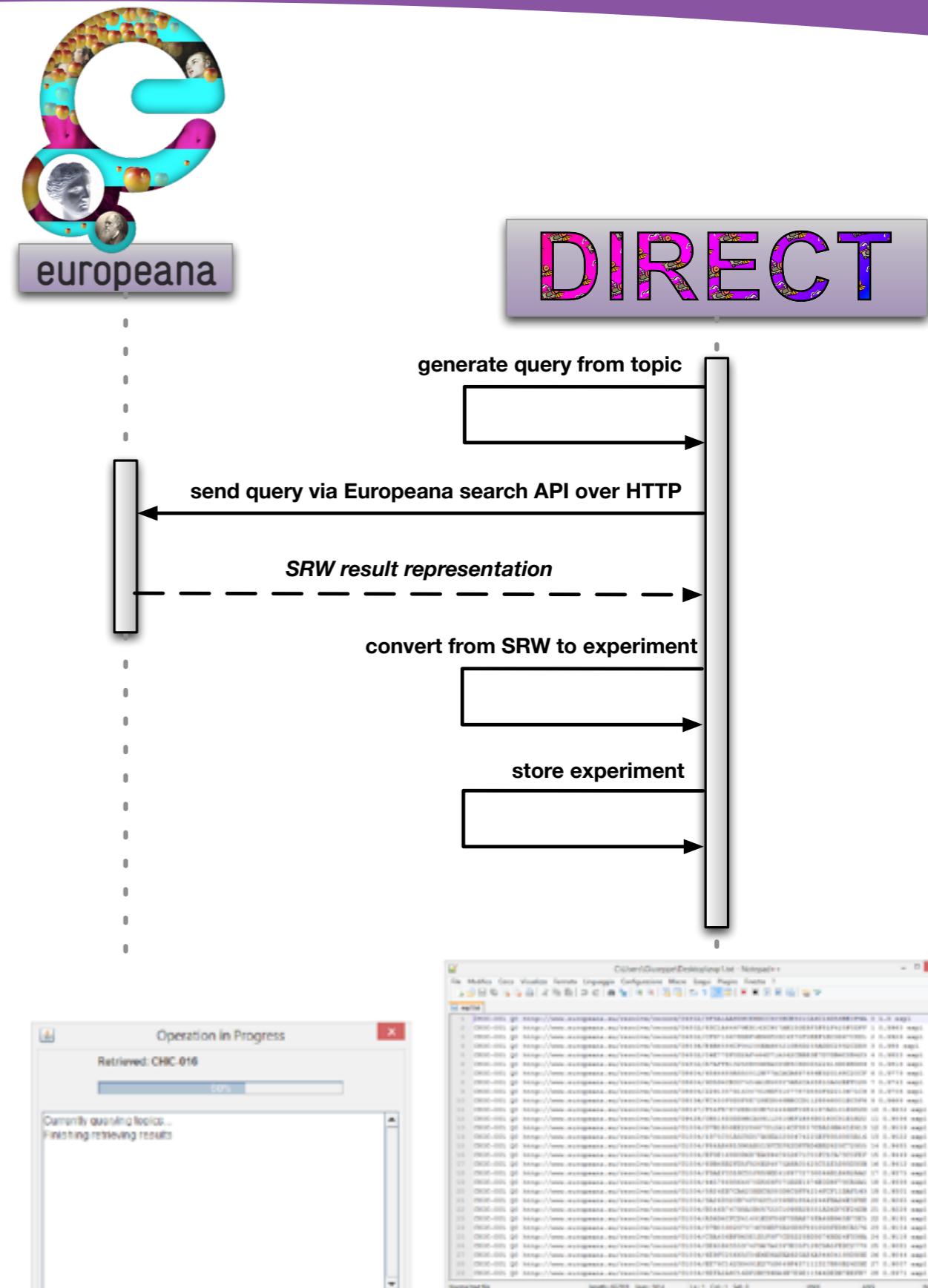
# User Annotations



- good results for experiments on topic 351
- please check correlation trend on scatter plot

 Table chart BoxPlot chart Stacked Bar chart Scatter Plot chart**SUBMIT**

- Programmatically evaluate a system under examination
  - somewhat similar to unit testing in software engineering
- Currently experimented with Europeana
  - based on the CHiC test collections



# Sharing and Extending

# The Community



- We have discussed why and how to design an evalution infrastructure
  - bridging between IR and DB is a key factor
- The need for an evaluation infrastructure is raising in the community
  - DESIRE 2011, SWIRL II, PROMISE Retreat
  - EvaluIR, Apache Open Relevance Project
- To be successful, enablers are needed
  - community adoption, modification, and extension
  - “political” willingness
  - economical sustainability

# Thank You

