




## GikiCLEF overview

Diana Santos  
Luís Miguel Cabral  
Nuno Cardoso  
Linguateca






SINTEF Information and Communication Technologies



## GikiCLEF in a nutshell

- Asking open list questions to Wikipedia
- A difficult task for people and for computers
- A realistic situation where crosslingual and multilingual skills may be of real interest

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

- A merger of QA and GIR, a follow-up of the GikiP pilot
- 10 Wikipedia collections, 50 culturally motivated topics
- 8 participants in 2009

SINTEF Information and Communication Technologies

## Topic titles in GikiP 2008: 3 different language cultures, but maybe artificial

ID	English topic title
GP1	Which waterfalls are used in the film "The Last of the Mohicans"?
GP2	Which Vienna circle members or visitors were born outside the <u>Austria-Hungarian</u> empire or Germany?
GP3	<u>Polish-speaking</u> rivers that flow through cities with more than 150,000 inhabitants
GP4	Which <u>Swiss</u> cantons border Germany?
GP5	Name all wars that occurred on Greek soil.
GP6	Which <u>Australian</u> mountains are higher than 2000 m?
GP7	African capitals with a population of two million inhabitants or more
GP8	Suspension bridges in <u>Brazil</u>
GP9	Composers of Renaissance music born in <u>Germany</u>
GP10	Polynesian islands with more than 5,000 inhabitants
GP11	Which plays of <u>Shakespeare</u> take place in an Italian setting?
GP12	Places where <u>Goths</u> lived
GP13	Which navigable rivers in Afghanistan are longer than 1000 km?
GP14	<u>Reverend</u> architects who designed buildings in Europe
GP15	French bridges which were in construction between 1980 and 1990

SINTEF Information and Communication Technologies

## GikiCLEF task to the topic group

- Please find realistic questions that make sense in your language / in your culture (9 non-English languages) and hopefully have a better coverage in that particular Wikipedia
- Hopefully even difficult to translate or render in other languages

Thanks to the topic group (14): Corina Forascu, Pamela Forner, Danilo Giampiccolo, Katrin Lamm, Thomas Mandl, Petya Osenova, Anselmo Peñas, Erik Tjong Kim Sang, Diana Santos, Julia Schulz, Yvonne Skalban, Alvaro Rodrigo Yuste (+Paula Carvalho, Christian-Emil Ore)

Thanks to Fred Gey for English nativization

SINTEF Information and Communication Technologies

## Examples of GikiCLEF-like systems in use

- In a mountain trip in the middle of Norway, an Italian family...



- In which European countries is the bidet usual?
  - Portugal
  - Spain
  - Andorra
  - Italy
  - ...



SINTEF Information and Communication Technologies

## Examples of GikiCLEF-like systems in use

- In the middle of a football match on TV in a foreign station



- What south American countries have yellow in their football team?
  - Brazil
  - Colombia
  - Ecuador

SINTEF Information and Communication Technologies

## Examples of GikiCLEF-like systems in use

- German youth gathering (or chatting) discussing the future: where should we move to study and still have fun?



- Which German cities have more than one university?

- Berlin
- Bonn
- Cologne
- Aachen
- Bremen
- Augsburg
- Hamburg
- ...

## Examples of GikiCLEF-like systems in use

- Where was this American museum which had a Picasso painting on which we saw that program last Winter?



- Which American museums have Picasso works?

- Museum of Modern Art (MOMA)
- Museum of Fine Arts (Boston)
- Solomon R. Guggenheim Museum
- Metropolitan Museum of Art
- National Gallery of Art
- Art Institute of Chicago
- Denver Art Museum
- (...)

## Examples of GikiCLEF-like systems in use

- Which Romanian poets published volumes with ballads before 1931? ... which may have influenced Mircea Eliade? Preliminary research for a MA Thesis in Romanian literature...



Dimitrie Bolintineanu



Vasile Alecsandri



George Coșbuc



Elisabeta de Neuwied

## Results from the topic group work

- We gathered a lot of more or less culturally-laden topics (70 -> 50)
  - Some of them **cross-cultural**, in fact (like the influence of Italy in Hemingway or the spread of Picasso's influence in the American continent)
  - Most around themes such as thinkers/writers/famous people, and places
- Not really too geographical for the most part
  - Still, Alps are **subdivided differently** in different countries/languages
  - Flemish towns** turn out to be a not clear concept internationally
  - We allowed **visual** clues to decide (for snow, colours and/or left in a map)
- A lot of little-developed and inaccurate Wikipedia pages were found a posteriori, showing that it may not be so obvious as initially thought that searching Wikipedia crosslingually is a good idea

## The task as seen by a participant system

- In order to provide a correct answer, a system had to produce **justification** in at least one of the languages returned
- Even if correct, an answer would not be rewarded in GikiCLEF if it were not possible to assess by a human reader
- Justification could be in the page itself, or in a set of other pages provided together
  - Example question: Name places where Goethe fell in love.
  - One correct answer is Leipzig. But in the article about Leipzig this is not stated ©
  - Where to find the justification? In "Johann Wolfgang von Goethe" page we find the following excerpt: "In **Leipzig**, Goethe fell in love with Käthchen Schönkopf and wrote cheerful verses about her in the Rococo genre."
  - Correct and justified GikiCLEF answer: **Leipzig**, Johann Wolfgang von Goethe

## GikiCLEF evaluation measures computation

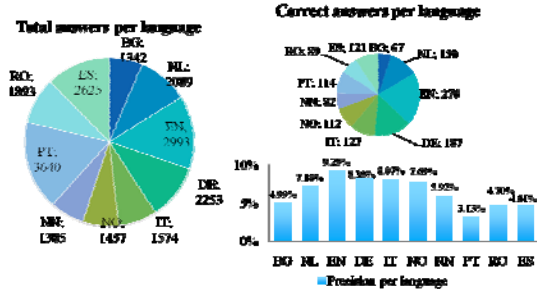
- C: Number of correct answers
- N: Total number of answers provided by the system
- GP: score per language: precision vs. C:  $C * C / N$  (so the score of Romanian will be  $C_{ro} * C_{ro} / N_{ro}$ )
- Total score: Sum of all scores per languages:  $\sum GP_{lang}$





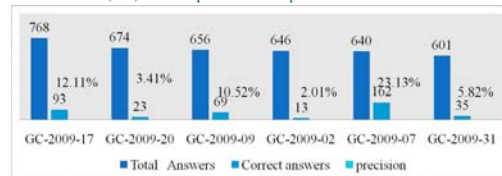


## Results overview per language



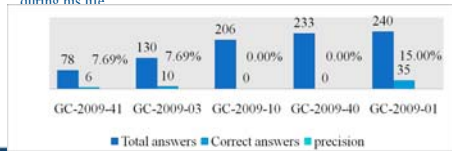
## Topics with most answers (by the participant set)

- GC-2009-17 (768): List the 5 Italian regions with a special statute
- GC-2009-20 (674): List the name of the sections of the North-Western Alps
- GC-2009-09 (656): Name places where Goethe fell in love.
- GC-2009-02 (646): Which countries have the white, green and red colors in their national flag?
- GC-2009-07 (640): What capitals of Dutch provinces received their town

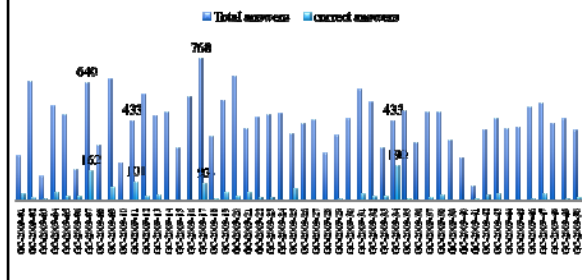


## Topics with least answers (by the participant set)

- GC-2009-41 (78): Chefs born in Austria who have received a Michelin Star.
- GC-2009-03 (130): In which countries outside Bulgaria are there opinions on Petar Dunov's (Beinsa Duno's) ideas?
- GC-2009-10 (206): Which Flemish towns hosted a restaurant with two or three Michelin stars in 2008?
- GC-2009-40 (233): Which rivers in North Rhine Westphalia are approximately 10km long?
- GC-2009-01 (240): List the Italian places which Ernest Hemingway visited during his life.



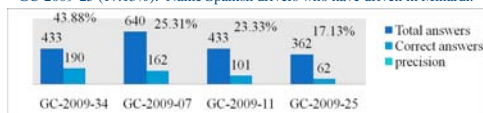
## Number of answers per topic



## Topics result overview: by precision

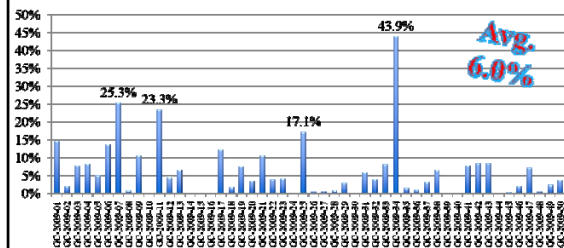
### Highest precision

- GC-2009-34 (43.9%): What eight thousanders are at least partially in Nepal?
- GC-2009-07 (25.3%): Which capitals of Dutch provinces received their town privileges before the fourteenth century?
- GC-2009-11 (23.3%): What Belgians won the Ronde van Vlaanderen exactly twice?
- GC-2009-25 (17.13%): Name Spanish drivers who have driven in Minardi.



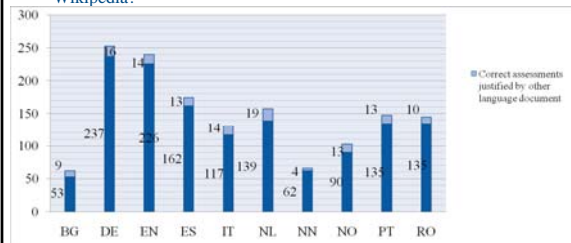
- 9 topics without any correct answer by the systems

## Precision by topic



## Authoritative languages: where is the justification

- For each language, was the justification found in the corresponding Wikipedia?



## Public resources

<http://www.lingueca.pt/GikiCLEF/>

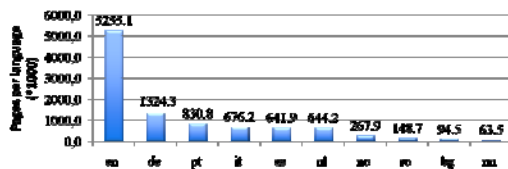
- Topic lists per language, in XML and text format
- Topic descriptions (in English, together with assessment decisions)
- Lists of correct answers
  - Correct and justified in some language: GikiCLEF\_answers\_correct\_justified.txt
  - Correct (because we know): GikiCLEF\_answers\_correct.txt
- Results (global, per language) and general statistics
- SIGA system, open source
- GikiCLEF collections

## Final discussion: Was GikiCLEF worthwhile?

- Recurrent comment: **too difficult!**
- Often not clear which (Wikipedia) pages were sought
  - Flag pages or country pages?
  - Team pages or country pages?
- Most (organizers) were probably not expecting so much work
- It is possible to create a system capable of answering many languages
- It pays to process English! The amount provided by each language, even in culturally-aware topics, is negligible
- The quality of other languages' Wikipedias is in strong need of improvement

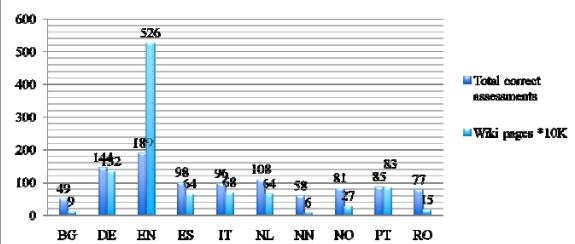
## Further details

## Size of collections (pages \* 1k)

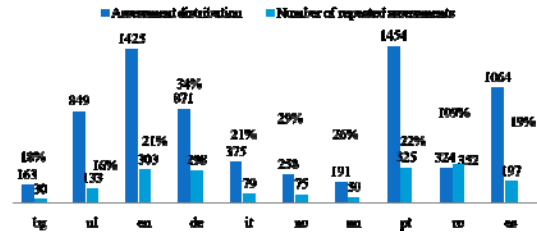


All collections made available were dated from June 2008, except for the English collection SQL dump, from May and July instead (since the SQL dump of June was not available)

## Number of authoritative answers per wiki size

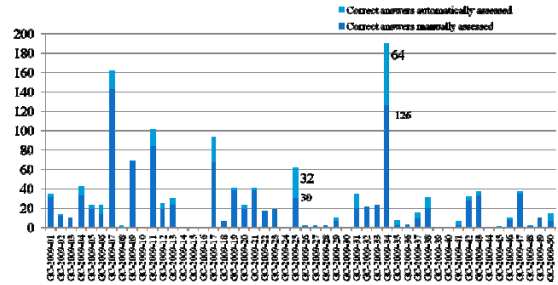


## Assessment processing

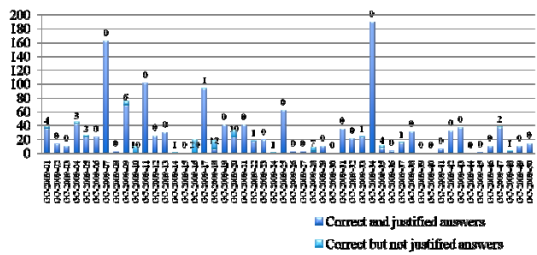


\*Overlap\*: ratio of repeated vs. total, but some answers were assessed by three or even four assessors, since repetition was automatically computed by SIGA.

## Number of correct answers automatically assessed and manually assessed



## Correct and justified answers VS correct not justified answers



## Use of justification

- Only 2 systems made use of the justification field (4 runs) with a total of 229 answers with another justification document

