# Evaluating Focused Retrieval Tasks

Jovan Pehcevski<sup>1</sup> and James A. Thom<sup>2</sup>

<sup>1</sup>AxIS Project Team INRIA Rocquencourt, France jovan.pehcevski@inria.fr

<sup>2</sup>School of Computer Science and Information Technology RMIT University, Melbourne, Australia james.thom@rmit.edu.au

#### Overview

- Focused retrieval
- Taxonomy of text retrieval tasks
- Evaluation framework
- Fidelity tests
- Discussion
- Q&A session



- Focused retrieval, including question answering, passage retrieval, and XML element retrieval, investigates ways to provide users with direct access to relevant information in retrieved documents
- Evaluating focused retrieval is a challenging task
  - different retrieval techniques typically produce answers of various sizes and granularity
  - there is a need for common evaluation framework where different aspects of focused retrieval can be consistently measured and compared

• The INitiative for the Evaluation of XML retrieval (INEX) has studied different aspects of focused retrieval since 2002

▲□▶▲□▶▲□▶▲□▶ □ のへで

・ロン ・ 御 と ・ 注 ・ ・ 注 ・

2

- by mainly considering XML element retrieval techniques that can effectively retrieve information from structured document collections
- by introducing different focused retrieval tasks, such as the in context tasks
- by using a highlighting assessment procedure to gather relevance assessments for the retrieval topics

- \* ロ > \* @ > \* 目 > \* 目 > 「目 ~ のへぐ

INEX In context tasks Evaluation

#### In context tasks

- <u>Relevant in context</u>: retrieve relevant documents, and identify the set of non-overlapping document parts representing the relevant information within each document
- Best in context: retrieve relevant documents, and identify the best entry point for starting to read the relevant information within each document
- The *in context* tasks correspond to end-user tasks, where focused retrieval answers are grouped per document, in their original document order
  - interactive experiments and user studies carried out within and outside INEX provide support for these tasks
  - the tasks loosely correspond to the INEX highlighting assessment procedure



#### Robertson's compatibility argument

"[...] there is a strong compatibility argument for researchers to use the same methods as each other unless there is very good reason to depart from the norm."

S.E. Robertson. Evaluation in information retrieval. In *ESSIR Proceedings*, p. 81–92, 2001.

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

INEX In context tasks Evaluation

#### **Evaluation**

- How to evaluate the in context tasks of focused retrieval?
- Two main requirements
  - the score should reflect the ranked list of documents inherent in the result list
  - the score should also reflect how well the retrieved information per document corresponds to the relevant information
- We want to use measures that directly exploit the INEX highlighting assessment procedure, and that are:
  - simple and easy to interpret
  - natural extensions of the well-established measures used in traditional information retrieval



- We present a taxonomy of text retrieval tasks based on the structure of the *answers* required by a task
- We also discuss some *assumptions* associated to a task, which model what users actually prefer
- These assumptions, together with the answer structure, define a retrieval task and influence how it should be evaluated

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト 一臣 - のへで

8

**Evaluating Focused Retrieval Tasks** 

retrieval tasks son framework Answers Fidelity tests Assumptions Discussion Q&A Session

#### Answers



- We present an evaluation framework for the *in context* tasks of focused retrieval
- The framework focuses on the compound answers shown in the taxonomy
- The evaluation of the *in context* tasks:
  - calculates scores for ranked lists of documents, where
  - the score per document reflects how well the retrieved information corresponds to the relevant information in the document

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Answers Assumptions

## Assumptions

- Basic assumption: Users want to see as much relevant information as possible with as little irrelevant information as possible
- This basic assumption is not sufficient to determine how best to evaluate most text retrieval tasks
- We need to make further assumptions about what users actually prefer. For example:
  - A1: Users consider all answers to be equally useful
  - A2: Users consider longer more detailed answers to be more useful than shorter answers



- Three scores per document S(d) could be calculated, depending on whether a single answer part, a set of answer parts, or a ranked list of answer parts are retrieved from the document
- We focus on the case where a set of non-overlapping answer parts is retrieved
- For a returned document, the text identified by the selected set of retrieved parts is compared to the text highlighted by the assessor

・ロト・4日・4日・4日・日・少へで

11

Score per document Scores for ranked list of documents

#### Score per document ...

- We calculate the following:
  - Precision P(d), as the fraction of retrieved text (in characters) that is highlighted for the document
  - Recall *R*(*d*), as the fraction of highlighted text (in characters) that is retrieved for the document
  - F-Score *F*(*d*), as the combination of precision and recall using their harmonic mean
- We use the F-score as an appropriate document score for the case where a set of non-overlapping answer parts is retrieved:

 SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007
 Evaluating Focused Retrieval Tasks
 13

 Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session
 Score per document
 Scores for ranked list of documents

• generalized Recall *gR*[*r*], as the number of relevant documents retrieved up to a document-rank *r*, divided by the total number of relevant documents (modelling assumption A1):

$$gR[r] = \frac{\sum_{j=1}^{r} rel(d_j)}{Nrel}$$
(2)

• Average generalized Precision *AgP* (modelling assumption A1):

SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007

$$AgP = \sum_{r=1}^{|\mathcal{D}|} \frac{1}{Nrel} \cdot rel(d_r) \cdot gP[r]$$
(3)

**Evaluating Focused Retrieval Tasks** 

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト 一臣 - のへで

15

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Score per document Scores for ranked list of documents

## Scores for ranked list of documents

- Over the ranked list of documents, we calculate the following:
  - generalized Precision gP[r], as the sum of document scores up to a document-rank r, divided by the rank r:

$$gP[r] = \frac{\sum_{j=1}^{r} S(d_j)}{r}$$
(1)



#### Scores for ranked list of documents ...

generalized Recall gR'[r], as the amount of relevant text retrieved up to a document-rank r, divided by the total amount of relevant text highlighted for the topic (modelling assumption A2):

$$gR'[r] = \frac{\sum_{j=1}^{r} rsize(d_j)}{Trel}$$
(4)

• Average generalized Precision *AgP*' (modelling assumption A2):

$$AgP' = \sum_{r=1}^{|\mathcal{D}|} \frac{rsize(d_r)}{Trel} \cdot rel(d_r) \cdot gP[r]$$
(5)

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 → のへで

Score per document Scores for ranked list of documents

## Scores for ranked list of documents ...

- Traditional information retrieval (IR) measures:
  - Precision *P*[*r*], as the fraction of retrieved relevant documents up to a document-rank *r*:

$$P[r] = \frac{\sum_{j=1}^{r} rel(d_j)}{r}$$
(6)

 Average Precision AP, as the average of the precisions calculated at natural recall levels:

$$AP = \sum_{r=1}^{|\mathcal{D}|} \frac{1}{Nrel} \cdot rel(d_r) \cdot P[r]$$
(7)



- For the *in context* retrieval tasks, there are two dimensions that we need to consider within the overall space of possible runs:
  - runs with different amounts of relevant and non-relevant information in the set of passages/elements returned for each document (dimension S)
  - runs with different rankings of the documents (dimension R)
- For a given evaluation measure these two dimensions may interact in unexpected ways

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Simulated runs Expected orderings Experimental results

# Fidelity tests

- Fidelity tests are designed to assess whether evaluation measures indeed measure what they are supposed to measure
- Simulated runs constructed in a controlled way are typically used to determine the fidelity of an evaluation measure
- Depending on the retrieval task, the best retrieval performance should be achieved by using the right (and desired) answer granularity, while preserving a reasonable relative ordering of the other simulated runs





▲□▶ ▲□▶ ▲ □▶ ★ □▶ ▲ □ ● の()

◆□ ▶ ◆□ ▶ ◆ 三 ▶ ◆ 三 ● ○ ○ ○ ○

Simulated runs Expected orderings Experimental results

#### Expected orderings for the S-R space



### Investigating the two dimensions

- Expected run orderings for the S dimension (different sets of answer parts returned for a document)
  - correctly captured by both AgP and AgP', but not by AP
  - information is lost in the abstraction toward the document level needed for AP
- Expected run orderings for the R dimension (different document rankings)
  - correctly captured by AgP
  - the swap of the first two document ranks without inserting a non-relevant document at the top is not captured by *AP*
  - the swap of the first two document ranks after inserting a non-relevant document at the top is not captured by AgP'

「日マネロマネロマネロマン」

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Simulated runs Expected orderings Experimental results

#### **Experimental results**

- We use version 5.0 of the INEX 2006 relevance assessments, which contains a set of judgements for 114 topics from INEX 2006
- We analyse the run performances to separately investigate the expected orderings on each of the two dimensions, as well as on the complete S-R space
- We compare scores obtained with the three overall performance measures (*AgP*, *AgP*', and *AP*)



- Expected run orderings for the S-R space
  - correctly captured by AgP
  - four notable disagreements between *AgP* and *AgP*' when comparing run pairs that insert non-relevant document at the top of their rankings
  - there are cases where the mean absolute performance differences obtained by AgP' are much larger than those obtained by AgP

Simulated runs Expected orderings Experimental results

# Investigating the S-R space ...

			ΔαΡ					ΔαP'		
Run ordering	Diff (%)	>	==	<	р	Diff (%)	>	==	<	р
SR->S <sub>L</sub> R	+20	112	2	0	2.2e-16	+17	112	2	0	2.2e-16
SR->SSR	+13	112	2	0	2.2e-16	+8	112	2	0	2.2e-16
SR->SR	0	0	114	0	_	0	0	114	0	_
SR->SR	+10	114	0	0	2.2e-16	+24	114	0	0	2.2e-16
S <sub>L</sub> R->S <sub>LD</sub> R	+26	113	1	0	2.2e-16	+16	113	1	0	2.2e-16
S <sub>L</sub> R->S <sub>L</sub> R <sub>S</sub>	+0.07	52	13	49	0.6023	+0.5	52	13	49	0.2962
S <sub>L</sub> R->S <sub>L</sub> R	+9	114	0	0	2.2e-16	+20	114	0	0	2.2e-16
S <sub>S</sub> R->S <sub>ST</sub> R	+41	114	0	0	2.2e-16	+45	114	0	0	2.2e-16
SsR->SsRs	+0.07	43	29	42	0.4146	+0.5	43	29	42	0.0963
SsR->SsRI	+9	114	0	0	2.2e-16	+22	114	0	0	2.2e-16
SR <sub>S</sub> ->S <sub>L</sub> R <sub>S</sub>	+20	112	2	0	2.2e-16	+17	112	2	0	2.2e-16
SRs->SsRs	+13	112	2	0	2.2e-16	+9	112	2	0	2.2e-16
SR <sub>S</sub> ->SR <sub>SI</sub>	+10	114	0	0	2.2e-16	+22	114	0	0	2.2e-16
SRI->SI RI	+18	112	2	0	2.2e-16	+14	112	2	0	2.2e-16
SRI->SSRI	+12	112	2	0	2.2e-16	+7	112	2	0	2.2e-16
SRI−>SŘ <sub>SI</sub>	0	0	114	0	—	-2	0	0	114	5.9e-13

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ → ≣ → のへで

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト 一臣 - のへで

27

SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007	Evaluating Focused Retrieval Tasks	25	
Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests <b>Discussion</b> Q&A Session	Passage versus element retrieval Focused versus traditional document retrieval Modelling evaluation assumptions		
Discussion			

- We use our findings to motivate a discussion about the following research topics:
  - the comparison between passage and element retrieval
  - the usefulness of focused and traditional document retrieval in identifying relevant information
  - the importance of modelling appropriate evaluation assumptions for a retrieval task

Focused retrieval Taxonomy of text retrieval tasks **Evaluation framework** Fidelity tests Discussion

Q&A Session

Simulated runs Expected orderings Experimental results

## Investigating the S-R space ...

			AgP					AgP'		
Run ordering	Diff (%)	>	==	<	р	Diff (%)	>	==	<	р
S <sub>LD</sub> R->S <sub>LD</sub> R <sub>S</sub>	+0.7	67	8	39	0.0004	+3	67	8	39	5.9e-0
S <sub>LD</sub> R->S <sub>LD</sub> R	+7	114	0	0	2.2e-16	+18	114	0	0	2.2e-1
S <sub>L</sub> R <sub>S</sub> ->S <sub>LD</sub> R <sub>S</sub>	+27	113	1	0	2.2e-16	+18	113	1	0	2.2e-1
S <sub>L</sub> R <sub>S</sub> ->S <sub>L</sub> R <sub>SI</sub>	+9	114	0	0	2.2e-16	+18	114	0	0	2.2e-1
$S_LR_I -> S_{LD}R_I$	+24	113	1	0	2.2e-16	+14	113	1	0	2.2e-1
S <sub>L</sub> R <sub>I</sub> ->S <sub>L</sub> R <sub>SI</sub>	+0.03	52	13	49	0.6023	-1	25	0	89	2.4e-0
S <sub>ST</sub> R->S <sub>ST</sub> R <sub>S</sub>	+0.1	60	0	54	0.4904	+1	60	0	54	0.214
S <sub>ST</sub> R->S <sub>ST</sub> R	+5	114	0	0	2.2e-16	+11	114	0	0	2.2e-1
S <sub>S</sub> R <sub>S</sub> ->S <sub>ST</sub> R <sub>S</sub>	+41	114	0	0	2.2e-16	+45	114	0	0	2.2e-1
S <sub>S</sub> R <sub>S</sub> ->S <sub>S</sub> R <sub>SI</sub>	+9	114	0	0	2.2e-16	+20	114	0	0	2.2e-1
S <sub>S</sub> R <sub>I</sub> ->S <sub>ST</sub> R <sub>I</sub>	+36	114	0	0	2.2e-16	+34	114	0	0	2.2e-1
S <sub>S</sub> R <sub>I</sub> ->S <sub>S</sub> R <sub>SI</sub>	+0.03	43	29	42	0.4146	-1	12	0	102	1.9e-0
SR <sub>SI</sub> ->S <sub>L</sub> R <sub>SI</sub>	+18	112	2	0	2.2e-16	+14	112	2	0	2.2e-1
SR <sub>SI</sub> ->S <sub>S</sub> R <sub>SI</sub>	+12	112	2	0	2.2e-16	+7	112	2	0	2.2e-1
S <sub>LD</sub> R <sub>S</sub> ->S <sub>LD</sub> R <sub>SI</sub>	+6	114	0	0	2.2e-16	+15	114	0	0	2.2e-1
S <sub>LD</sub> R <sub>I</sub> ->S <sub>LD</sub> R <sub>SI</sub>	+0.4	67	8	39	0.0004	+0.05	46	0	68	0.879
S <sub>L</sub> R <sub>SI</sub> ->S <sub>LD</sub> R <sub>SI</sub>	+24	113	1	0	2.2e-16	+15	113	1	0	2.2e-1
S <sub>ST</sub> R <sub>S</sub> ->S <sub>ST</sub> R <sub>SI</sub>	+5	114	0	0	2.2e-16	+10	114	0	0	2.2e-1
S <sub>ST</sub> R <sub>I</sub> ->S <sub>ST</sub> R <sub>SI</sub>	+0.05	60	0	54	0.4896	-1	48	0	66	0.018
SsRsi->SstRsi	+36	114	0	0	2.2e-16	+35	114	0	0	2.2e-1

SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007

**Evaluating Focused Retrieval Tasks** 26

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion

Q&A Session

Passage versus element retrieval Focused versus traditional document retrieval Modelling evaluation assumptions

#### Passage versus element retrieval

• Perfect retrieval for the relevant in context task can only be achieved when retrieving all the highlighted passages within a document, in their exact size



#### Passage versus element retrieval

Focused versus traditional document retrieval Modelling evaluation assumptions

#### Passage versus element retrieval ...

- The absolute performance difference between the passage run and our best simulated element run was 13%, which shows that no element run can achieve perfect retrieval
- One explanation for this could be that there is an inherent bias of the INEX highlighting assessment procedure towards passage retrieval
- How can passage and element retrieval be sensibly compared?



## Focused versus traditional document retrieval

- Traditional IR measures, such as AP, cannot fully capture the level of detail required by focused retrieval
- More specifically, the AP measure partially captures the different ordering of documents in the result list, but it does not capture how well the retrieved information per document corresponds to the relevant information
- The average generalized precision *AgP* measure is able to fully capture both evaluation aspects, which makes it more useful than *AP* in measuring the retrieval performance
- On the INEX 2006 test collection, AgP is able to distinguish more significant performance differences than AP

31

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Passage versus element retrieval Focused versus traditional document retrieval Modelling evaluation assumptions

#### Passage versus element retrieval ...

- If there is an inherent bias towards passages, then this should be taken into account when comparing these two types of retrieval
- Two different sub-tasks could be identified to allow for sensible comparison:
  - Passage retrieval sub-task, where the retrieval answers are passages and it makes sense to compare whether element retrieval techniques help in identifying more relevant passages
  - Element retrieval sub-task, where the retrieval answers are XML elements and it makes sense to compare whether passage retrieval techniques help in identifying more relevant elements.



#### Modelling evaluation assumptions

- Assumptions A1 and A2 are of particular importance for *in* context retrieval tasks, as it is not entirely clear which of the two assumptions should be preferred for evaluation
- Our fidelity tests demonstrate that the AgP' measure (based on assumption A2) is not entirely measuring what it is supposed to measure, and that the AgP measure (based on assumption A1) correctly captures the expected run orderings
- An argument for assumption A2 is that it motivates the preference given to more exhaustive answers in the evaluation

Passage versus element retrieval Focused versus traditional document retrieval Modelling evaluation assumptions

#### Modelling evaluation assumptions ...

- It may be possible that the current *AgP*' definition (shown in Equation 5) is not correctly modelling assumption A2
- Fixing this definition requires further investigation, which might be solved in one of these two ways:
  - interpolated average generalized precision could be used instead of the current non-interpolated definition
  - the current non-interpolated *AgP*' definition could be re-defined as follows:

$$AgP' = gR'[|\mathcal{D}|] \cdot \frac{\sum\limits_{r=1}^{|\mathcal{D}|} rel(d_r) \cdot gP[r]}{\sum\limits_{r=1}^{|\mathcal{D}|} rel(d_r)}$$
(8)

r=1

ロト 4 同 ト 4 三 ト 4 三 ト 9 오 오

33

35

SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007 Evaluating Focused Retrieval Tasks

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

## **Questions?**



Greetings from Versailles!

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

Passage versus element retrieval Focused versus traditional document retrieval Modelling evaluation assumptions

### Modelling evaluation assumptions ...

- A more fundamental challenge, however, relates to the user preference of the two evaluation assumptions
- Would users regard a focused and more concise answer as more useful than a lengthy exposition?
- Or would they indeed perceive the answer that contains more relevant (and possibly repeating) information as more useful?
- We believe that it may be possible to determine the answers to these and similar questions either via user experiments or by questioning assessors about how they valued the answers for their topics



- Since 2005, a highlighting assessment procedure is used at INEX to gather relevance assessments
- Assessors are asked to highlight sentences representing the relevant information in a pooled set of documents
- The assessment tool automatically computes the relevance of the judged document parts (including full documents) as the ratio of highlighted to fully contained text
- The relevance values are drawn from a continuous [0,1] relevance scale

#### INEX highlighting assessment procedure ...



annication during testing. Thus, testing must cover various kinds of applications and provide very

 Let us assume that two systems, System A and System B, respectively retrieve the following ranked lists of elements:

Rank	System A	System B
1	/article[1]/bdy[1]/sec[1]	/article[1]/bdy[1]/sec[1]/p[1]
2	/article[1]/bdy[1]/sec[2]	/article[1]/bdy[1]/sec[1]/p[2]
3	/article[1]/bdy[1]/sec[3]	/article[1]/bdy[1]/sec[1]/p[3]

Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

#### Appendix B: HiXEval evaluation scenarios

- An XML retrieval system may return a series of smaller elements that belong to a larger fully highlighted element, with the goal to boost the performance scores (overall and at selected rank cutoffs)
- We use two scenarios that allow us to perform a more detailed analysis of this (possibly undesirable) evaluation behaviour



• The recall-base contains only one fully highlighted section, which consists of three fully highlighted paragraphs:

<element path="/article[1]/bdy[1]/sec[1]"
size="99" rsize="99"/>
<element path="/article[1]/bdy[1]/sec[1]/p[1]"
size="33" rsize="33"/>
<element path="/article[1]/bdy[1]/sec[1]/p[2]"
size="33" rsize="33"/>
<element path="/article[1]/bdy[1]/sec[1]/p[3]"</pre>

element path="/article[i]/bdy[i]/sec[i]/p]

```
size="33" rsize="33"/>
```

#### Scenario 2

• The recall-base contains two fully highlighted sections, and the first section consists of three highlighted paragraphs:

```
<element path="/article[1]/bdy[1]/sec[1]"
size="99" rsize="99"/>
<element path="/article[1]/bdy[1]/sec[1]/p[1]"
size="33" rsize="33"/>
<element path="/article[1]/bdy[1]/sec[1]/p[2]"
size="33" rsize="33"/>
<element path="/article[1]/bdy[1]/sec[1]/p[3]"
size="33" rsize="33"/>
```

<element path="/article[1]/bdy[1]/sec[2]"</pre>

size="99" rsize="99"/>

#### Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session

## HiXEval performance scores

	Hi	iXEval	measu	re
System	P@3	R@3	F@3	AP
Scenario	1			
A	0.33	1.00	0.50	1.00
В	1.00	1.00	1.00	1.00
<u>Scenario</u>	2			
А	0.67	1.00	0.80	1.00
В	1.00	0.50	0.67	0.50

Table: HiXEval performance scores for the two evaluation scenarios, obtained with rank cutoff and overall performance measures. Best scores under each HiXEval measure are shown in bold.

	《日》《聞》《臣》《臣》	≣ ୬୍ର୍େ
SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007	Evaluating Focused Retrieval Tasks	41
Focused retrieval Taxonomy of text retrieval tasks Evaluation framework Fidelity tests Discussion Q&A Session		
HiXEval performance score	es	
<ul> <li>The desired trade-off betwee information as possible and amount of non-relevant info AP and F@r (to some exter</li> </ul>	en retrieving as much relevant not retrieving a substantial rmation is correctly captured nt), but not by <i>P</i> @ <i>r</i> and <i>R</i> @ <i>r</i>	ant d by
<ul> <li>We currently normalise ove retrieved (the rank cutoff r), required to reach that rank</li> </ul>	r the number of elements and not over the user effort cutoff	
<ul> <li>Future work: normalise ove instead of over the number calculate Precision/Recall a</li> </ul>	r the amount of text returned of elements retrieved (that is at number of characters reac	d s, l)
How to determine the exact	cutoff values?	-
	<ul> <li>&lt; □&gt; &lt; □&gt; &lt; □&gt; &lt; □&gt; &lt; □&gt;</li> </ul>	き うくで
SIGIR 2007 Workshop on Focused Retrieval, 27/07/2007	Evaluating Focused Retrieval Tasks	43