



中科院计算所
INSTITUTE OF COMPUTING TECHNOLOGY, CAS

ICTCAS's ICTGrasper at TAC 2008: Summarizing Dynamic Information with Signature Terms Based Content Filtering

Jin Zhang

zhangjin@software.ict.ac.cn

**Key Laboratory of Network Science and Technology,
Institute of Computing Technology, P. R. China**

TAC Update Task Results

- **Automatic Evaluation (71 peers)**

Criterion	Rank	Score
BE	1*	0.06480
ROUGE-2	3	0.09776
ROUGE-SU4	5	0.13295

- **Manual Evaluation (64 peers)**

Criterion	Rank
Pyramid	2*

ICTGrasper obtained top 2 ranks on BE, and placed the 2nd and 3rd on Pyramid

- **Did pretty well on both measures**

Update Summary Introduction

- Temporal content can be divided into sub-collections corresponding to time intervals

- Update summary focuses on the dynamic information between current interval and its previous ones

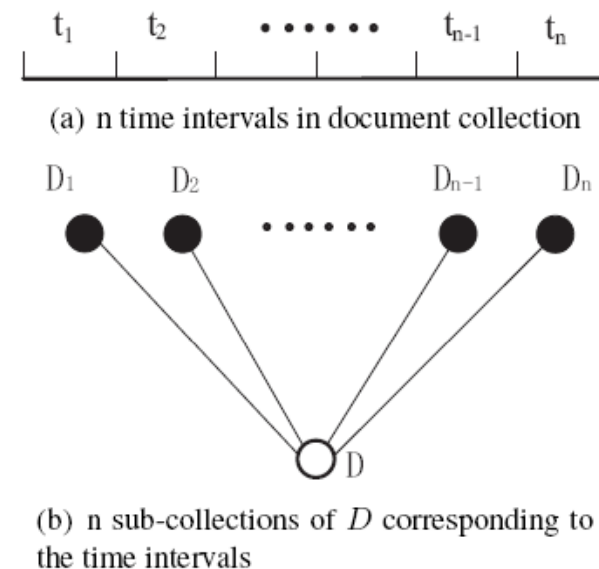


Figure 2: Formalization of Dynamic Summarization

Divide temporal content into current information I_c and history information I_h

For TAC 2008, document set A can be denoted as I_h , B can be denoted as I_c .

Problem Formalization

- **Update Summary of TAC 2008**

- Divide temporal content into current information I_c and history information I_h
 - document set A: I_h
 - document set : I_c

- **Task of Update Summary**

- Identify dynamic information
 - Based on the relationship between I_c and I_h
- Evaluate importance of dynamic information
 - Sentence ranking criterion
- Select important dynamic information

Overview of ICTGrasper

- Content Filtering Model
- Rerank sentences based on signature terms

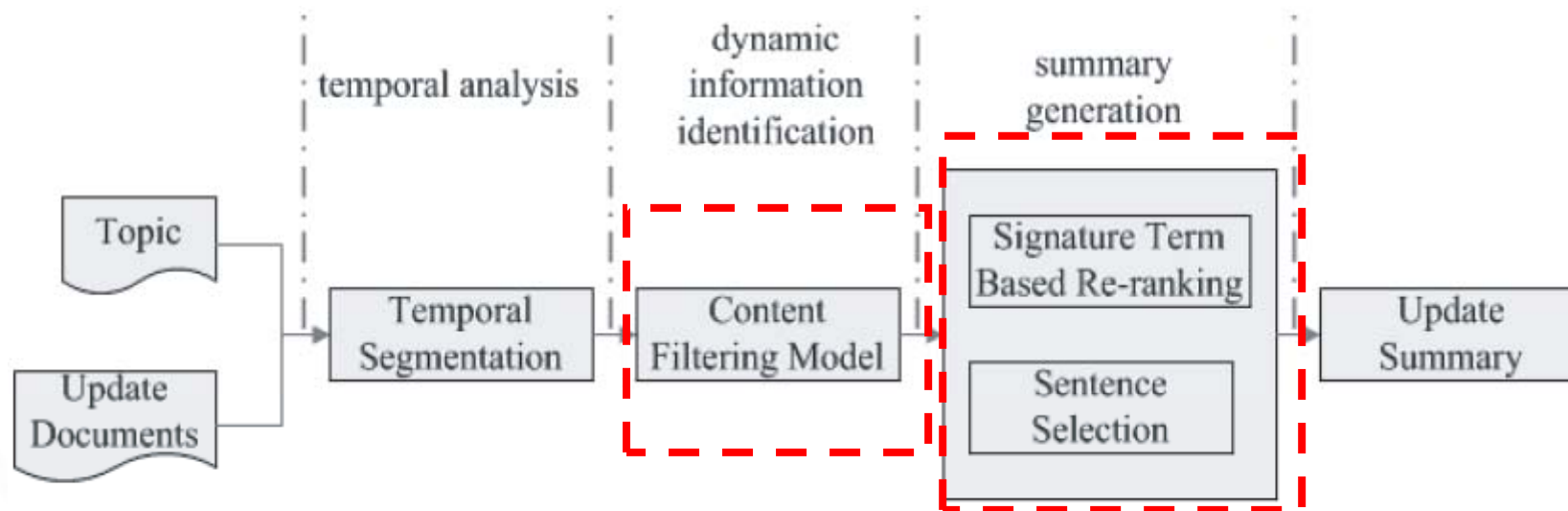


Figure 1: The framework of ICTGrasper system for dynamic summarization

Content Filtering Model

- **Content filtering model**

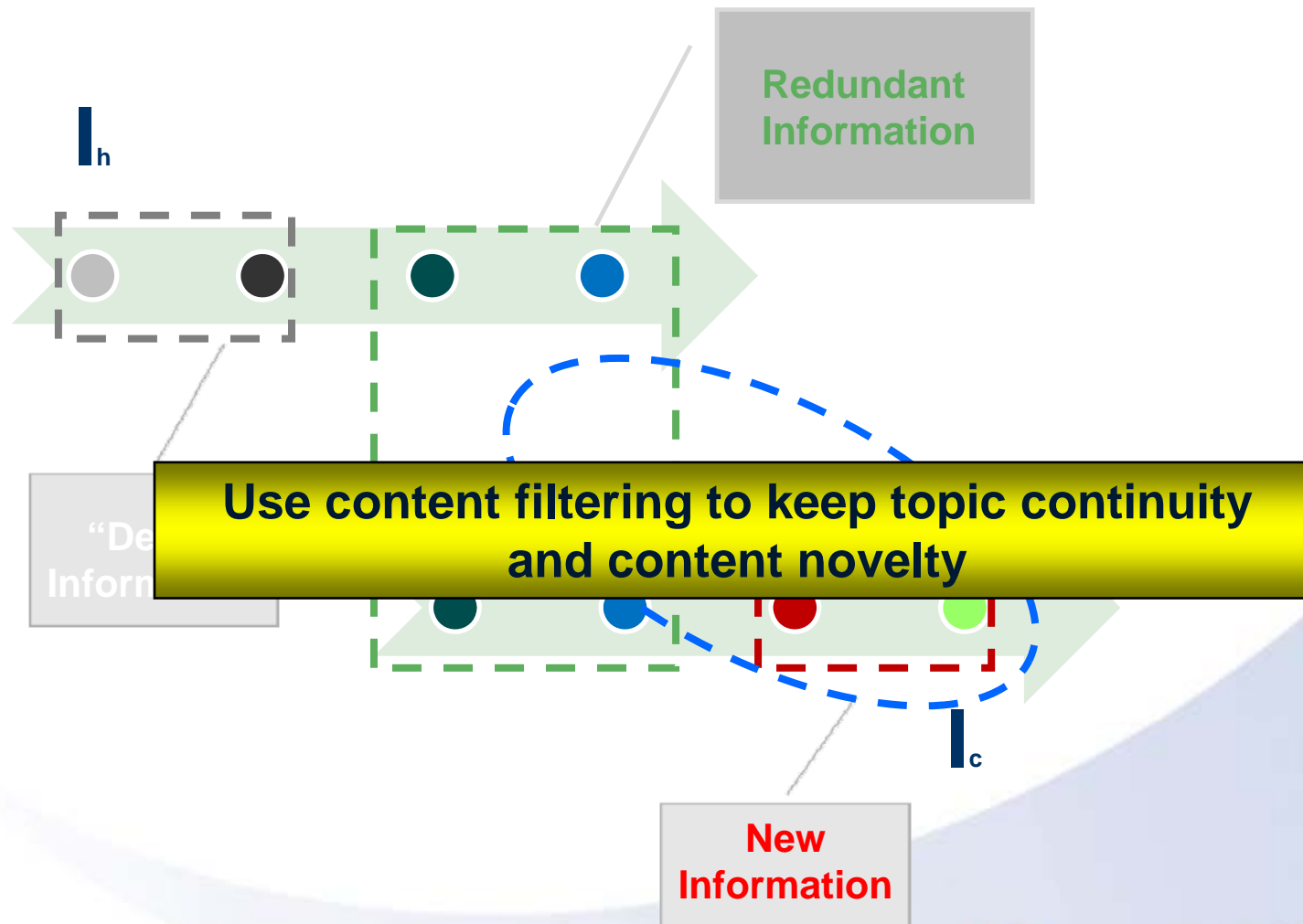
- **Dynamic information identification at sentence level**

- **Document Filtering Model**
- **Summary Filtering Model**
- **Union Filtering Model**

Employ content filtering model for dynamic information identification at sentence level



Content Filtering Model



Content Filtering Model

● Solution method

– Static summarization algorithm

● GSPS

$$f(t+1) = \lambda r + (1 - \lambda) \hat{P}f(t)$$

– DFM

$$f(I_c - I_h) = GSPS(I_c - I_h) = GSPS(\{s | s \in I_c, \{s\} \not\subseteq I_c \tilde{\cap} I_h\})$$

– **S** **f** is a summarization function, and the minus symbol “-” denotes the filtering operation.

– UFM

$$f(I_c + I_h) - I_h = GSPS(I_c + I_h) - I_h$$



Topic Signature Reranking

● Topic Signatures

- signature terms are those terms which occur significantly more than expected "at large"

Temporal Topic Signature	AV	Tf	Temporal Topic Signature	AV	Tf
european central bank	5	6	single currency	8	11
single currency	4	11	central bank	6	7
currency euro	4	6	national currency	4	5
new single currency	4	4	exchange rate	4	4
central bank	3	8	foreign currency	4	4
european single currency	3	4	adopt euro	3	4
european union	3	4	currency euro	3	4
member state	3	3	san marino	3	4

Terms importance varies with time intervals

Topic Signature Reranking

- **Adapt a linear to optimize importance measure of a sentence**

$$Rank(s) = \alpha * cen(s) + \beta * sim(s, T) + \gamma * score(s|t)$$

- ***cen(s)* is centrality of *sentence, s, with LexRank***
- ***sim(s, T)* is the similarity between *s* and topic**
- ***score(s|t)* is the score decided by topic signatures in a specified time interval**



Experiments

- **Conduct experiments on DUC 2007 Dataset**
- **Set up two experiments**
 - The effectiveness of content filtering model
 - The improvement of signature terms
- **Evaluate with ROUGE Metric**
 - ROUGE-2
 - ROUGE-SU4



Experiments

● Performance of Content Filtering Models

Table 1: The results of three models for dynamic summarization, where the degree denotes the degree of membership in set, R-2 and R-SU4 are the scores of ROUGE-2 and ROUGE-SU4 on dataset of DUC 2007 update task.

Degree	DFM1		DFM2		SFM1		SFM2		UFM1		UFM2	
	R-2	R-SU4	R-2	R-SU4	R-2	R-SU4	R-2	R-SU4	R-2	R-SU4	R-2	R-SU4
1.0	0.1114	0.1450	0.1142	0.1482	0.1142	0.1482	0.1142	0.1482	0.1163	0.1492	0.10268	0.13847
0.4	0.1126	0.1467	0.1146	0.1478	0.1165	0.1495	0.1142	0.1482	0.1180	0.1506	0.1050	0.1398

Union content filtering model can be chosen as the optimal model for dynamic information identification

Experiments

● Performance of Content Filtering Models

Table 2: Performance comparison with state-of-the-art systems of DUC 2007 update task, where *UFM1* represents the performance of our proposed union filtering model, *LCC*, *IIIT*, and *NUS* are the top performing systems.

System	ROUGE-2	ROUGE-SU4
UFM1(Degree=0.4)	0.1180	0.1506
LCC (Rank 1)	0.1119	0.1431
IIIT (Rank 2)	0.0985	0.1352
NUS (Rank 3)	0.0962	0.1325
Generic Baseline	0.0850	0.1225

Performance of UFM1 outperforms DUC 2007 top systems

Experiments

● Performance of Signature Terms

Table 3: The performances of signature terms, where UFM1-N is the system with UFM1 but without consideration signature terms, and UFM1-S is the system with signature term based UFM1.

System	ROUGE-2	ROUGE-SU4
UFM1-S	0.1219	0.1581
UFM1-N	0.1180	0.1506
LCC	0.1119	0.1431

Performance improvement for both ROUGE-2 and ROUGE-SU4 is very obvious



TAC 2008 Evaluation Results

● Automatic Evaluation

Table 4: The evaluation results of TAC 2008 top performing systems, where $R - 2$ and $R - SU4$ stand for the ROUGE-2 and ROUGE-SU4 scores in ROUGE evaluation. For convenience, only four runs with stable performance are illustrated in this table, and Run 14 and Run 65 are two runs of ICTGrasper.

Run	ROUGE				BE	
	R2	Rank	R-SU4	Rank	BE	Rank
14	0.09776	3	0.13295	5	0.06480	1
65	0.09559	5	0.13151	9	0.06293	2
43	0.10395	1	0.13646	1	0.06267	3
60	0.09449	6	0.13583	3	0.06203	4

TAC 2008 Evaluation Results

- **Manual Evaluation**

- **ICTGrasper placed 2nd and 3rd in the average modified (pyramid) score over 64 peers**



TAC 2008 Evaluation Results

● Evaluation Results on Document Set B

Table 5: The performances of Grasper's best run (Run 14) on the Document Set *B*.

Metric	Score	Rank
mod Pyramid score - B	0.344	1
numScus - B	4.063	2
ROUGE-2 Recall - B	0.101	1
ROUGE-SU4 Recall - B	0.137	1
BE Recall - B	0.076	1



Conclusion

- **Introduced a signature terms based content filtering approach for update summarization**
 - Identify dynamic information at sentence level with content filtering models
- **Proposed to rerank the importance of filtered content based on topic signatures**
- **The results show that our proposed approach works very well in update task**



中科院计算所
INSTITUTE OF COMPUTING TECHNOLOGY, CAS

Thanks !