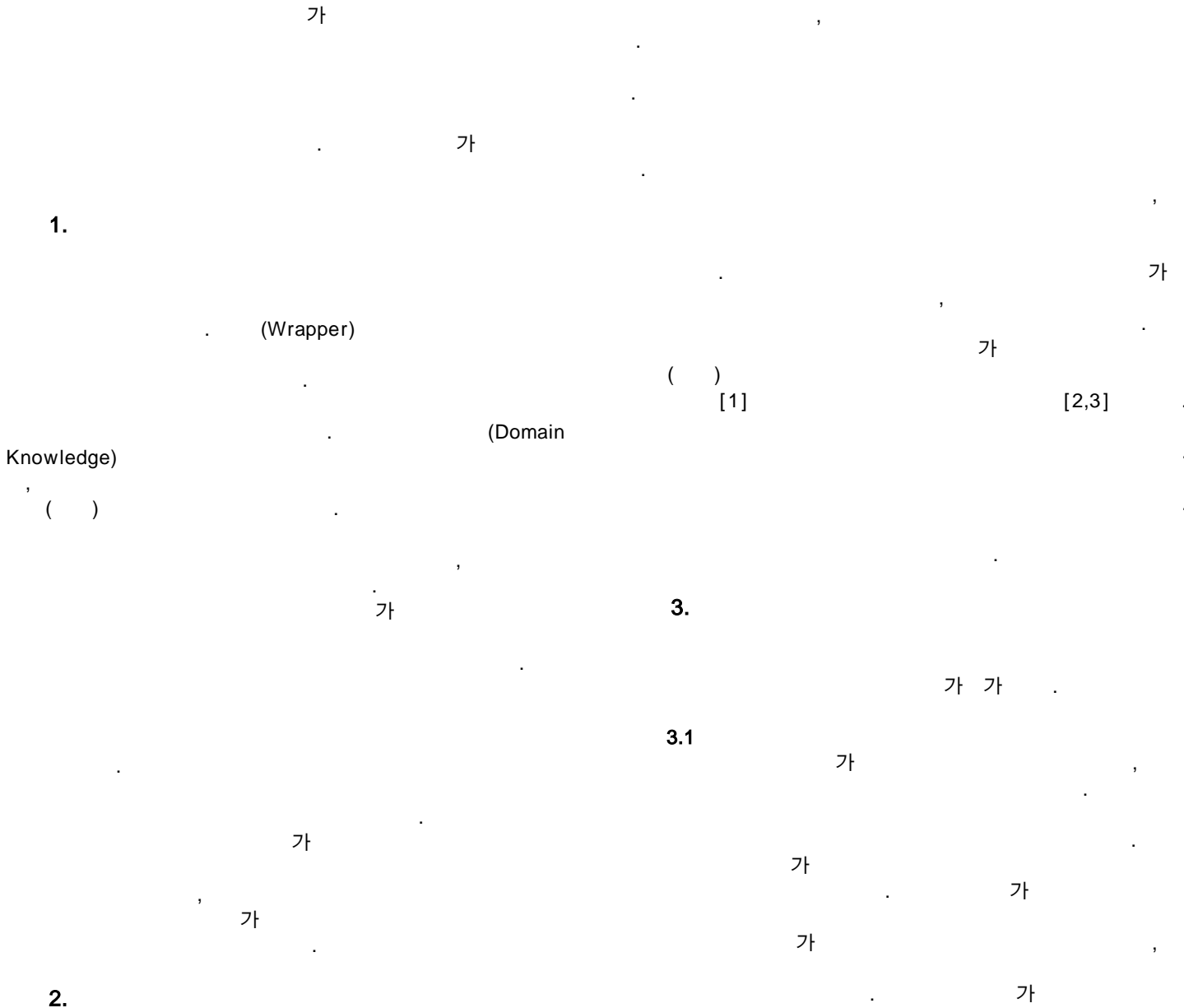


{chjeong, jerry, ybj}@kisti.re.kr,
myaeng@icu.ac.kr

Improving Rule Generation Precision for Wrappers using Domain Knowledge

Chang-Hoo Jeong^{*}, Jerry Hyeon Seo^{*}, Beom-Jong You^{*}, Sung-Hyon Myaeng^{**}
Korea Institute of Science and Technology Information^{*},
Information and Communications University^{**}



$$P(e'_i | t'_j) = \frac{P(t'_j | e'_i) * P(e'_i)}{P(t'_j)} \dots$$

$$\cong P(e'_i) * P(t'_j | e'_i)$$

$$P(e'_i | T'_j) \cong P(e'_i) * P(T'_j | e'_i) \dots$$

$$\cong P(e'_i) * \frac{1}{v} \sum_{k=1}^v P(t'_{jk} | e'_i)$$

-
-
-
- 가
- 가
-
- 가
- 가
- 가
-
- 가
-
- 가
-
- (Front page)
- (Back end page)

가

가

가

가

3.2

$$P(e'_i = t'_j | \{e_1 = t_1 \& e_2 = t_2 \& \dots \& e_n = t_n\})$$

$$= \{P(e'_i = t'_j | e_1 = t_1) * P(e_1 = t_1)\}$$

$$+ \{P(e'_i = t'_j | e_2 = t_2) * P(e_2 = t_2)\}$$

$$+ \dots$$

$$+ \{P(e'_i = t'_j | e_n = t_n) * P(e_n = t_n)\}$$

$$= \sum_{h=1}^n P(e'_i = t'_j | e_h = t_h) * P(e_h = t_h)$$

$$P(e'_i = T'_j | \{e_1 = T_1 \& e_2 = T_2 \& \dots \& e_n = T_n\})$$

$$\cong \frac{1}{v} \sum_{k=1}^v \sum_{h=1}^n P(e'_i = t'_{jk} | e_h = t_{hk}) * P(e_h = t_{hk})$$

t e

, t' e'

t

T T'

가

가

가

가

가

가

()

가

가

가

가

가

P1,

P2,

() α

P1) + ((1 - α) * P2)가

α

(α *

4.1

가 , 가 , 가
 가 , 가
 가

$$(P) = (\quad / \quad) * 100$$

가 가

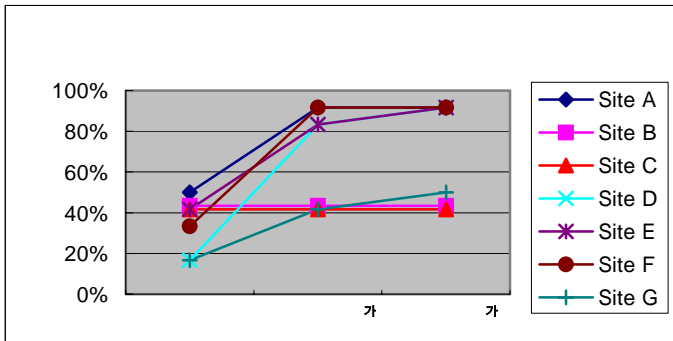
$$(Global\ Precision) = (\quad / \quad)$$

가

가

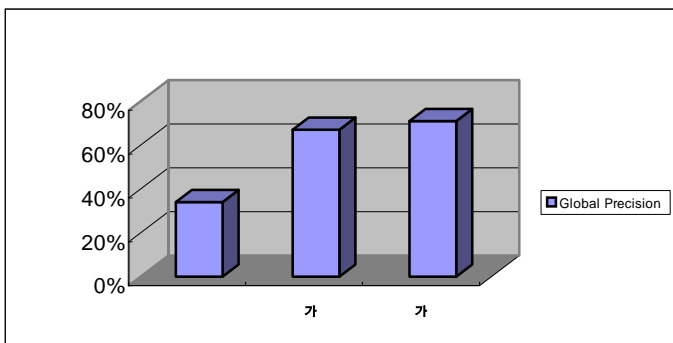
4.2

[1]



[1]

[2]



[2]

가 가

가 가 가

5.

가 가
 가 가
 가 가

6.

[1] N. Kushmerick, D. Weld, and R. Doorenbos, "Wrapper Induction for information extraction", International Joint Conference on Artificial Intelligence (IJCAI), Nagoya, Japan, 1997.

[2] H. Seo, J. Yang, and J. Choi, "Knowledge-based Wrapper Generation by Using XML", IJCAI-2001 Workshop on Adaptive Text Extraction and Mining (ATEM 2001), pp. 1-8, Seattle, USA, 2001.

[3] , "Wrapper", : , 29 1-2 , pp. 42-52, 2002.