

工學碩士 學位論文

Identification of Numerical Expressions using
Finite-State Automata

指導教授 朴 侏 讚

2002年 2月

韓國海洋大學校 大學院

工 學 科

白 億 種

本 論 文 白 億 種 工 學 碩 士 學 位 論 文 認 准

委 員 長 工 學 博 士 孫 周 永 印

委 員 工 學 博 士 辛 沃 根 印

委 員 工 學 博 士 朴 然 讚 印

2001年 12月

韓 國 海 洋 大 學 校 大 學 院

工 學 科 白 億 種

1	1
2	3
2.1	3
2.2	6
2.3	9
3	11
3.1	11
3.2	12
3.3	14
4	18
4.1	18
4.2	19
4.3	20
4.4	21
4.5	22
4.6	22
4.7	가	23
5	가	26
5.1	26
5.2	가	26
5.3	가	27
5.4	30
6	33
	35

Identification of Numerical Expressions using Finite-State Automata

Ock-Jong Baek

Department of Computer Engineering, Korea Maritime University, Pusan, Korea

Abstract

There have been many trials to parse sentences in text to search complete and exact parses, but it is very hard because of unavoidable incompleteness of lexicon and grammar. Recently, to alleviate these difficulties, partial parsing appears as an alternative in the field. Partial parsing aims to recover syntactic information efficiently and reliably from unrestricted text, by sacrificing completeness and depth of analysis.

As a part of partial parsing, the identification of Korean numerical expressions in text is described in this paper. Numerical expressions are required in several systems such as information extraction systems and question-answering systems. One of desired characteristics of these systems is the fastness. To achieve this goal, we use a finite-state automaton, for which we could use a tool like lex. So that we could rapidly implement the system. We observed that the system is fast and correct through several experiments. To evaluate our system, we used newspaper as test collection. We achieved the recall of 90.8%, and the precision of 86.9%. Experiments show that our system is comparatively correct.

1

(ambiguity)

(partial parsing)

가

[1-4].

가 가

[5-7].

가

(numerical

expressions)

(temporal expression)

(quantity

expression)

(date), (time),

(duration)

(money), (percent),

(measure), (cardinal)

[8].

(regular expression)¹⁾

(finite-state

automata, FSA)

가

1)

가 [12,13,14].

lex

[9- 11].

가

가

[5, 7]

[12]

2

3

4

5

가

6

2

2.1

(regular expression)

(regular grammar)

(regular language)

가

1

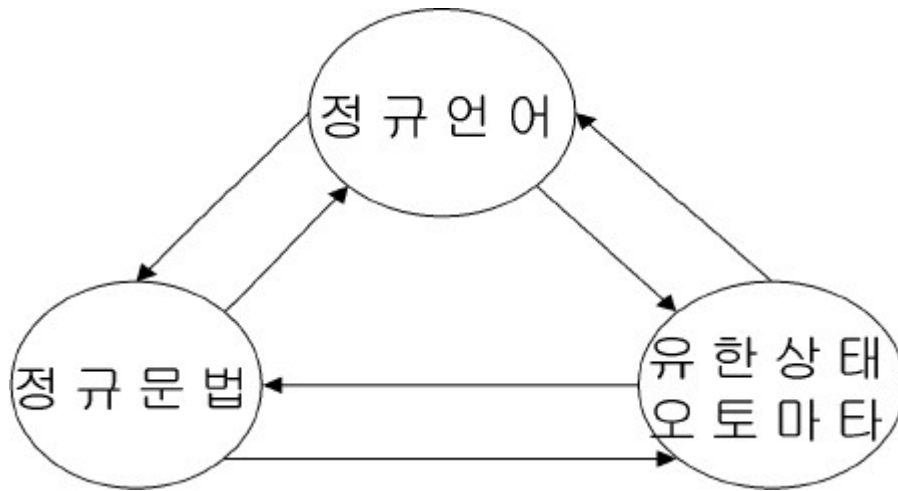
가 [9].

(lex) 1975

(Lesk, M. E)

가

C



1

Fig. 1 Regular expression, regular grammar, and finite-state automata

(lex source)

```

<      >
%%
<      >
%%
<      >
%%가      ,
<      >   %{  }%      C
      ,      ,      .      %{  }%
lex.yy.c
<      >
      <      >
lex.yy.c
  
```


(text character)

(operator character)

, *

a* a

1

Table 1 Meta characters for patterns and their meanings

.	newline character (\n) (single)
*	0
[]	Bracket [] character class
^	[] negation
\$.
{ }	{ }
\	disable , C escape sequence
+	+ 1
?	? 0 1
"..."	quotation mark
/	slash "/"
()	()

2.2

[1-3]. (cascaded finite-state automata)
가

가
가
가

2.2.1 Cascaded Analysis of Syntactic Structure(CASS)

가 . CASS

[1].

- 1 : NP -> D? A* N+ | Pron
- VP -> Md Vb | Vz | Hz Vbn | Bz Vbn | Bz | Vbg
- 2 : PP -> P NP
- 3 : SV -> NP VP
- 4 : S -> (Adv | PP)? SV NP? (Adv | PP)*

D : , A : , N : , Pron : , P : , Adv :
NP : , VP : , PP :

-0 , L L-1

(longest match) . CASS

, 95%

2.2.2 FASTUS (Finite-State Automation Text Understanding System)

Hobbs et al.[6]

(information extraction system) . FASTUS

1. .
2. .
3. .
4. .
5. .

FASTUS

, MUC(message understanding conference) 가 100 12
가 .

2.2.3 Fidditch

Fidditch Hindle[13] , 가

. Fidditch

(unrestricted text)

. Fidditch Marcus (punt) 가

(skip and fit),

(, that, which),

Fidditch

가

2.2.4

[14]

[5]
가 .

[14]
, (local grammar) , FST (finite state transducer)

2.3

, 가

2.3.1

가

가

(template)

[5, 7].

가가

DB

가

가 .

2.3.2

가

가

가

[12].

3.2

가

가

가

가

가

“ 5 ” “ ’ ”

가

3.2.1

가

가

(.)가

(.)

가

(-)

(:)

(\$)

(%)

()

가

(3-1) 2.7%

(3-2) 2,000,000

(3-3) 2001-3-20 (2001 3 20)

(3-4) 12:30 (12 30)

3.2.1

가
(countable),
(uncountable)

[15]. “ : ‘ ’ , ‘ ’)
가 ” 3

3

Table 3 Subcategory of unit bound nouns

		, , , ,
		, ,
		, 가 , , , ,
		, , ,
		, , , , ,
		, ,
	
		, , , ,
		, , , ,
		, , , (), (), (, 가) (), (), (,), () (), (,), (), ()
		, , (), (), (), (), (,)

3.3

(DTE), (TME), (DUR)
(MNY), (PCT), (MSR), (CRD) [8].

“2001 12 25 ”

“ ” “ ” “ ” “ ” “ ” “ ”

가

3.3.1 (Date)

가

, ,

가

, “2001 3 2 ”

“ ”

가

가

(3-5) 1999 12 25 [redacted]
 (3-6) 76 3 [redacted]
 (3-7) 76 [redacted]
 (3-8) [redacted] 7 가 가
 가 가
 10 12 가

(3-9) [redacted]
 (3-10) [redacted]
 (3-11) [redacted] 5 [redacted]
 (3-12) [redacted] 7 17 [redacted]
 (3-13) [redacted]
 (3-14) [redacted] [redacted]
 (3-15) [redacted] 9 11 가
 (3-16) 3 [redacted] 가

3.3.2 (Time)

가 가
 가

(3-14) 11 11 11 [redacted]
 (3-15) 10 30 [redacted]
 (3-16) 7 30 가 [redacted]
 (3-17) [redacted] 9 [redacted]
 (3-18) [redacted] 12 가 [redacted]
 (3-19) [redacted] 4 [redacted]

3.3.3 (Duration)

“2 3 ” () 가

“ ”

“ ” “ ”

- (3-20) 12
- (3-21) 10 12 가
- (3-22) 1 3
- (3-23)
- (3-24) : 1994 2000
- (3-25) : 14 16

3.3.4 (Money)

가 , , , , ,

\$, ₩ 弗 ,

가

가 가 ,

가

가

- (3-26) 2
- (3-27) 2,000
- (3-28) 17 가
- (3-29) 1\$ 1200.36 가 가

3.3.5 (Percent)

%

가

- (3-30) 5% 가
- (3-31) 7.6%
- (3-32) 가

3.3.6 (Measure)

- 가 , , , , , , ,
- 가 가
- [16] , , , , 5 14
- [17] , , , , , , , , ,
- 10가 ,

- (3-33) 79 1
- (3-34) 18
- (3-35) 100
- (3-36)
- (3-37)
- (3-38)
- (3-39)

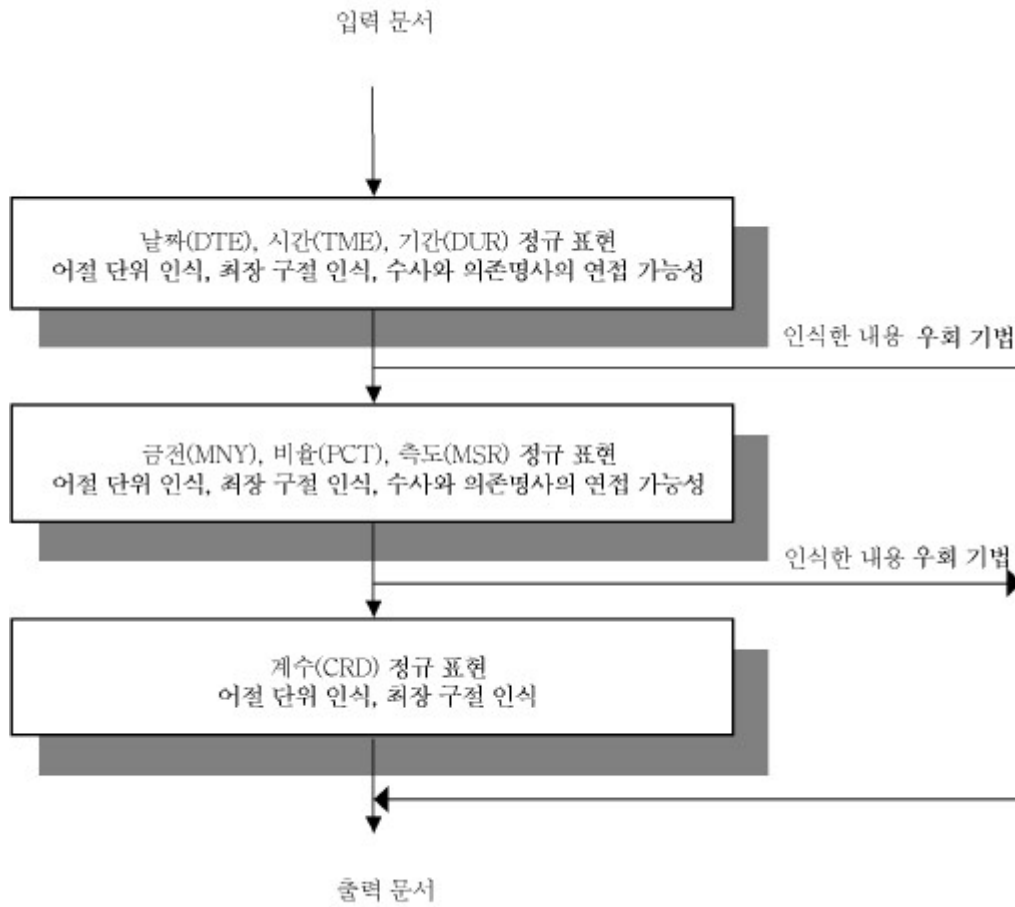
3.3.7 (Cardinal)

- (3-40) 2 , 2.7
- (3-41) 3
- (3-42) 1 25

4

4.1

가
가 , 가
가 , 가
1 , 2
3
가 , 가
가 , 가 , 가



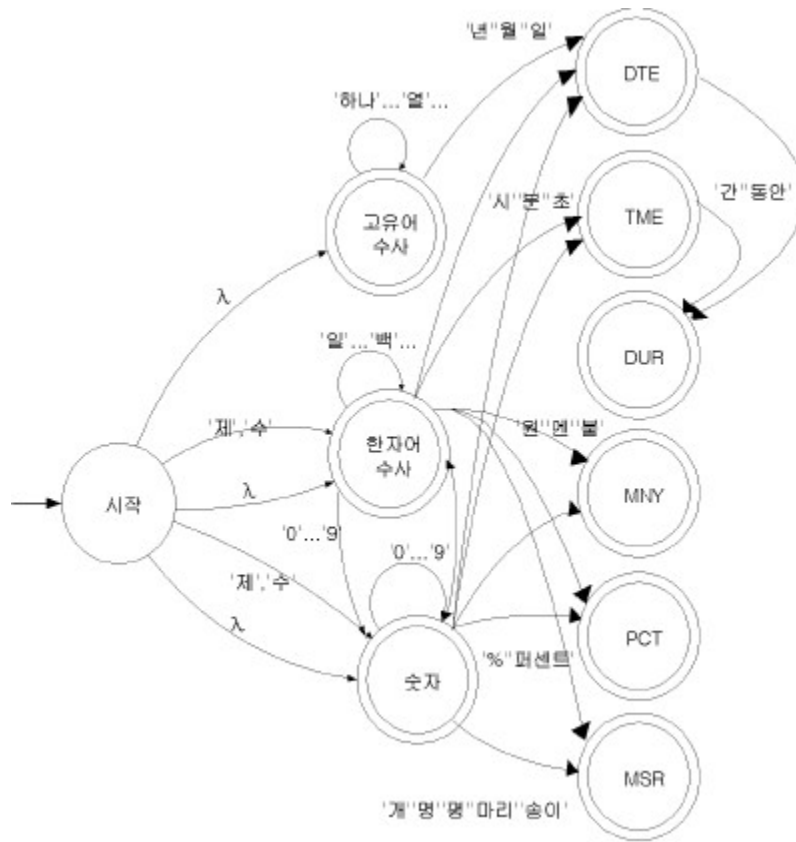
2

Fig. 2 Overview of a numerical expression identifier

4.2

가

3



3

Fig. 3 A finite automata for recognizing numerals

3

가

4.3

:

가

가

lex

NNN ([0-9])+
 NNC (“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”)+
 NNK (“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”|“ ”)+

NDy ({NNN}|{NNN}|{NNC})*{SP}?+“ ”?“ ”
 NDm ({NNN}|{SP}?)“ ”
 NDd ({NNN})+“ ”?“ ”

NTh ({NNN}|{NNK})+“ ”
 NTm ({NNN}|{SP}?)+“ ”?“ ”
 NTs ({NNN}|{SP}?)+“ ”?“ ”

NDy_{md} ({TENTWELVE}|{NDy}){SP}?{NDm}{SP}?{NDd}
 NDy_m {NDy}{SP}?{NDm}
 NDm_d {NDm}{SP}?{NDd}

NTh_{ms} {NTh}{SP}?{NTm}{SP}?{NTs}
 NTh_m {NTh}{SP}?{NTm}
 NTm_s {NTm}{SP}?{NTs}

4.4

가

‘ ’

‘ ’

가

가

가 .

: 7 60 가
 : <MNY>7 60 </MNY> 가
 : <MNY><MSR>7 </MSR>60 </MNY> 가
 : <MNY>7 60 </MNY> 가

4.5

가 , , 가 , 가 . , , , . , , .

- (4-1) <MSR> </MSR>
- (4-2) <MSR> </MSR>
- (4-3) <MSR> </MSR>
- (4-4) <MSR> </MSR>
- (4-5) <CRD> </CRD>
- (4-6) <MSR> </MSR>
- (4-7) <MSR> </MSR>
- (4-8) <MSR> </MSR>

가

- (4-9) <MSR>4 </MSR>

4.6

가

“2001 12 30 ”

가

가

가

Lex

가

, 가

(longest match) ,

(rule given first)

가

,

4.7

가

가

가

가 가

“10 10 ”

“

”

10

가

가

가

가

가

가

가

가

[18, 19].

가

가 , 가 20 30,
 100 가
 가 , 가 100 가
 가 가
 가 가 ,
 가

“ , , , , ” 가

[20, 21].

가 , 가 ,
 가 가 , “ ” , “ ” , “ ”

가

() : UNIT

5 가

5.1

가 가 94
 가 , 73,547 5,657
 . 가 ,
 가 .
 SUN Sparc 10 Workstation .

5.2 가

가 (Recall) R
 (Precision) P ,
 가 $Fscore$, (1), (2) (3)
 [22].

$$R = \frac{N_R}{N_C} \tag{1}$$

$$P = \frac{N_R}{N_S} \tag{2}$$

$$Fscore = \frac{(\times + 1)PR}{\times (P + R)} \tag{3}$$

N_S , N_R 가
 . N_C 가

$R > 1$, $P < 1$ 가 1, 2, 5

5.3 가

5.3.1

4 가 (*Recall*), (*Precision*), *Fscore*

, 가

가 가

4
 Table 4 System performances

가	3,862
	4,036
	3,507
<i>Recall</i>	90.8%
<i>Precision</i>	86.9%
<i>Fscore</i>	88.8%

5.3.2

5

90% (DUR), (CRD)

80%

가 (TME), (CRD)

5

Table 5 Performance for each tag

	가			<i>Recall</i>	<i>Precision</i>	<i>Fscore</i>
TME	73	96	72	98.6	75.0	85.2
DTE	1317	1258	1243	94.4	98.8	96.5
DUR	206	168	154	74.8	91.7	82.4
MNY	261	267	242	92.7	90.6	91.7
PCT	212	207	206	97.2	99.5	98.3
MSR	1372	1449	1291	94.1	89.1	91.5
CRD	421	591	311	73.9	52.6	61.5

5.4

530 가 , 가
 356 가 .
 501 .
 가 , “56 44” , 가 <PCT>56
 44</PCT> , <MSR>56 </MSR><CRD>44</CRD>
 , “4 1” 가 <CRD>4 1</CRD>,
 <TME>4 </TME> <CRD>1</CRD>
 가 .

오류 점유율



4

Fig. 4 Error Analysis

5.4.1

- (5-1) <MSR> </MSR>
- (5-2) <MSR> </MSR>
- (5-3) <MSR> </MSR>
- (5-4) <CRD> </CRD>

5.4.2

39%

가

(5-5) <CRD> </CRD>

(5-6) <CRD> </CRD>

(5-7) <CRD> </CRD>

(5-8) <CRD> </CRD>

(5-9) <CRD> </CRD>

(5-10) <CRD> </CRD>

5.4.3

가

가

(5-11) <MSR>25 </MSR>

(5-12) <MS>1 </MSr>

(5-13) <MSR> 1 </MSR>

(5-14) <MSR> 2 </MSR>

(5-15) <DUR>3 </DUR>

5.4.4

가

가

” , +
 , +
 , 가

- (5-16) <DTE>4 · 19</DTE>
- (5-17) <CRD>4</CRD> · <CRD> 19</CRD>
- (5-18) <PCT>56 44</PCT>
- (5-19) <MSR>56 </MSR><CRD>44</CRD>
- (5-20) B<MSR>777 </MSR>
- (5-21) ISO<MSR>9001 </MSR>
- (5-22) <MSR> </MSR>
- (5-23) <MSR> </MSR>
- (5-24) <MNY> </MNY>
- (5-25) <MSR> </MSR>
- (5-26) <MNY> </MNY>
- (5-27) <CRD> </CRD>

5.4.5

가

가

가

가

(5-28) 10 9 7
(5-29) 10 9 7

가 .

가

, (5-28)

, (5-29)

가 ,

가 .

가 .

가 .

6

가 , 가

가

가 가 94 가 가
가 가 , 가
86.9%

90.8%

가

가

[23-25]

[15]

- [1] Steven Abney, "Parsing by Chunks," Kluwer Academic Publisher, 1991.
- [2] Salah Ait-Mokhtar, J. P Chanod, "Incremental finite state parsing," In ANLP'97, 1997.
- [3] , " , " , 7 , 6 , pp. 83-96, 2000.
- [4] , " , " , , , 1998.
- [5] , , , , " , " 10 pp. 3-8, 1998.
- [6] Douglas E. Appelt, Jerry R. Hobbs, John Bear, David Israel, and Mabry Tyson, "FASTUS : A finite-state processor for information extraction from real-word text," in Proceedings IJCAI '93, Chambéry, France, August 1993.
- [7] , , , , " , " , 19 , 10 , pp. 27-39, 2001.
- [8] N. Chinchor, P. Robinson, and E. Brown, Hub4 "Named Entity Task Definition (version 4.8)," SAIC, August 1998.
- [9] Lesk, M.E. & Schmidt, E. "Lex : A Lexical Analyzer Generator. Anonymous (Ed.), Unix Research System Papers, Tenth Edition. Murray Hill, NJ: AT&T Bell Laboratories," 1990.

- [10] , , , 1996.
- [11] , , , 1994.
- [12] , , , , “
 ,” 4
 pp.469-477, 1992.
- [13] Donald Hindle. “User manual for Fidditch. Technical Report 7590-142,
 Naval Research Laboratory,” 1998.
- [14] , , , “ FST
 ,” 11
 pp.231-236, 1999.
- [15] , , “ ,” 12
 pp.395-401, 2000.
- [16] , “ ,” 89, , 1983.
- [17] , , , 1994.
- [18] , , , 1986
- [19] , , , 2000.
- [20] , , , “ / /
 ,” 13 pp.341-347,
 2001.
- [21] , “ ,” 9
 pp.136-142, 1997.
- [22] Frakes, W. B. and Baega-Yates, R., Information Retrieval : Data
 Structures & Algorithms, Prentice Hall, 1992.
- [23] , “ ,”

(B), 23, 9, pp.991- 1000, 1996.

[24], “ ,” 10
 , pp. 137- 142, 1998.

[25], “ bigram ,” 12
 , pp. 85- 88, 2000.

ICQ ,

가

가

가

가

3

(?)

(?)

가

, 가 ,

. , ,

, 가 ,

,

.

가 가 ,

.

가

,

,

,

2

,
가