

А.Б. Романюк, В.М. Макар  
Національний університет “Львівська політехніка”,  
кафедра систем автоматизованого проектування

## РОЗРОБЛЕННЯ ДВОРІВНЕВОГО МОРФОЛОГІЧНОГО АНАЛІЗАТОРА УКРАЇНСЬКОЇ МОВИ У ПРОГРАМНОМУ СЕРЕДОВИЩІ РС-КІММО

© Романюк А.Б., Макар В.М., 2008

**Останнім часом у лінгвістичному забезпеченні САПР спостерігається тенденція створення діалогових мов, наближених до природної мови. Використання природної мови зумовлює необхідність розроблення морфологічних, синтаксичних та семантичних аналізаторів. У статті описано розроблення морфологічного аналізатора української мови на основі системи РС-КІММО.**

**In the last years in the CAD linguistic providing there is a tendency of creation of dialog languages of close to the human language. The use of human language predetermines the necessity of development of morphological, syntax and semantic analyzers. In this article development of morphological analyzer of Ukrainian is described on the basis of the system РС-KIMMO.**

### Вступ

Серед відомих розробок морфологічних аналізаторів найпоширеніша розроблена фінським дослідником Кіммо Коскенніємі система КІММО, яку також називають дворівневою моделлю. На базі цієї моделі створено як інструментальні засоби, які дають змогу без додаткового програмування розробляти описи морфологічних структур природних мов, так і бібліотеки процедур, які уможливають використання морфологічних процедур у прикладних [1].

У системі реалізований погляд на опис морфології природної мови як на формально-математичну задачу, що дало змогу використовувати виключно математичний апарат при розробленні моделі та оцінці її обчислювальної ефективності. Морфологічний аналізатор і синтезатор у КІММО-моделі – це перетворювачі, які працюють за принципом скінченного автомата.

До переваг КІММО-моделі насамперед належить її високий ступінь мовної незалежності за більша, ніж в інших моделях, простота складання правил перетворень. Мовна незалежність та засновані на строгих і добре відомих програмістам математичних конструкціях механізми реалізації сприяли розробленню різних версій ефективного програмного забезпечення системи програмістами з різних країн і різними мовами програмування. Наявність ефективних програмних засобів сприяла розробленню різних версій лінгвістичного забезпечення для багатьох мов світу. Отже, на базі моделі КІММО виник унікальний для лінгвістичної практики феномен поширеної стандартної мовнезалежної морфологічної системи із розвинутим набором програмних засобів і добре вивченим корпусом лінгвістичних описів. Тому вивчення можливостей КІММО-моделі для створення морфологічного аналізатора української мови є актуальною науковою задачею, яка має також значну практичну цінність [5].

Експериментальні реалізації морфологічних процесорів виконані з використанням КІММО-моделі у програмних середовищах РС-КІММО, РУ-КІММО для японської, арабської, шведської, німецької, фінської, англійської, аккадської, румунської, російської, церковнослов'янської, французької, іспанської, іврити, турецької та інших мов [1]. Для деяких мов реалізовано описи тільки окремих морфологічних підсистем [4].

### РС-КІММО: дворівневий морфологічний аналізатор

РС-КІММО – це програма, яка здійснює аналіз та/або синтез слів, використовуючи дворівневу модель, в якій словоформа подається як відповідність між її формою на лексичному рівні та на поверхневому рівні [2]. Поверхнева форма – це запис словоформи згідно із орфографічними правилами, тобто формування словоформи з основи та аломорфів. Лексична форма словоформи – це узагальнений запис словоформи послідовністю морфем.

У програмному середовищі РС-КІММО опис морфології природної мови здійснюється створенням двох файлів:

1. Файла правил, в якому визначається абетка та фонологічні (орфографічні) правила.
2. Файла лексику, в якому містяться лексичні одиниці (слова і морфеми), правила їхньої сполучуваності та відповідні їм морфологічні характеристики.

Наприклад, для англійської мови поверхнева форма словоформи “spies” повинна відповідати її поверхневій формі “`spru+s” так ( ` – символ наголосу, + символ морфемного шва, 0 – символ нульового елемента):

Лексична форма:	`	s	p	y	+	0	s
Поверхнева форма:	0	s	p	i	0	e	s

Правила повинні описувати відповідність між символами двох рівнів `:0, y:i, +:0, та 0:e. Наприклад, дворівневе правило для встановлення відповідності між y:i буде таким:

$$y:i \Rightarrow @:C\_+ :0 \quad (1)$$

Генератор та аналізатор є основними функціональними частинами РС-КІММО. Генератор на вході отримує лексичну форму, застосовує до неї фонологічні (орфографічні) правила і повертає відповідну поверхневу форму. Аналізатор на вході одержує поверхневу форму, застосовує до неї фонологічні (орфографічні) правила, переглядає лексикон і повертає відповідну лексичну форму з набором відповідних морфологічних [3].

Правила і лексикон опрацьовують, використовуючи скінченні автомати. Наприклад, дворівневе правило (1) для його використання в РС-КІММО повинно бути перетворене на таку таблицю переходів скінченного:

	@	y	+	@
	C	i	0	@
1:	2	0	1	1
2:	2	3	2	1
3:	0	0	1	0

### Розроблення файлів правил та лексику для опису морфології української мови.

Для реалізації опису морфології української мови в інструментальному середовищі РС-КІММО необхідно розробити файл правил та файл лексику.

Потрібно передбачити, що символи кирилиці не підтримуються в РС-КІММО і буде використовуватись така відповідність між символами:

a-a, б-b, в-v, г-g, ґ – G, д-d, е-e, є-E, ж-Z, з-z, и-y, і-i, ї-I, й-j, к-k, л-l, м-m, н-n, о-o, п-p, р-r, с-s, т-t, у-u, ф-f, х-h, ц-c, ч-C, ш-S, щ-H, ю-U, я-A, ь-B.





```

RULE " 0:'=> __ Apost" 2 3
    0 Apost @
    ' Apost @
1: 2 1 1
2. 0 1 0

RULE " t:C => __u" 2 3
    t u @
    C u @
1: 2 1 1
2. 0 1 0

RULE " t:C => __Ss" 2 3
    t Ss @
    C Ss @
1: 2 1 1
2. 0 1 0

RULE " k:c => __i" 2 3
    k i @
    c i @
1: 2 1 1
2. 0 1 0

RULE " e:i => __l#" 3 4
    e l # @
    i l # @
1: 2 1 1 1
2. 0 3 0 0
3. 0 0 1 0

RULE " i:e => __C Ie" 3 4
    i C Ie @
    e C Ie @
1: 2 1 1 1
2. 0 3 0 0
3. 0 0 1 0

RULE " 0:C => __U" 2 3
    0 U @
    C U @
1: 2 1 1
2. 0 1 0
END

```

Файл лексикону в PC-KIMMO складається з двох частин. У першій частині описуються стани скінченного автомата, а в другій частині – дуги переходів між цими станами. Стани скінченного автомата описуються в частині, яка називається Alternations, а його переходи в частині Lexicons [2]. В Alternations описуються всі назви станів скінченного автомата, який дає змогу

побудувати словоформу. Перші складові цієї частини описують кореневі поняття дерева лексикону. Кожен рядок файла лексикону має такий формат:

1. Ключове слово ALTERNATION.
2. Назва класу Alternation (стан скінченного автомата).
3. Одна або більше назва стану скінченного автомата. Наприклад:

ALTERNATION      Begin    <назви станів скінченного автомата >

Один з станів скінченного автомата повинен здійснювати перехід до спеціального стану

End. Наприклад:      ALTERNATION      Foo    End

Перший рядок файла лексикону для морфологічного аналізатора української мови буде таким:

ALTERNATION      Begin prefixVV1 prefixVV2 prefixVV3 prefixVV4 verb/var1 verb/var2  
verb/var3 verb/var4 noun/var1/tv noun/var2/tv noun/var3 noun/var4 noun/var1/mi noun/var2/mi  
noun/var1/mj noun/var2/mj

Після опису можливих станів скінченного автомата необхідно описати переходи між цими станами. Опис кожного з переходів містить три основні складові: (1) символ переходу, рядок символів лексичної форми; (2) назву можливого стану скінченного автомата; (3) будь-який рядок символів, який необхідно подати на вихід при переході між станами скінченного автомата (набір морфологічних характеристик). Фрагменти файла лексикону з описом переходів наведено в табл. 2.

Таблиця 2

#### Фрагменти файла лексикону

Приклад опису префікса 'на':	Приклад опису кореня 'смаж' дієслова:
<pre>\f na+ \x prefixVV1 \alt verb/var1 \fea \gl Dokonanyj vyd_</pre>	<pre>\f smaZ \x verb/var4 \alt suffixVV4 \fea \gl smaZ</pre>
Приклад опису кореня 'фермер' іменника:	Приклад опису суфікса 'ити':
<pre>\f fermer \x noun/var2/tv \alt endingNV2/tv \fea \gl fermer</pre>	<pre>\f yty \x suffixVV3 \alt endingVInfV3 \fea \gl +yty_</pre>
Приклад опису флексії 'меш':	Приклад опису флексії 'ла':
<pre>\f meS \x endingVInfV1 \alt End \fea \gl +meS_sg.2 pers. Future</pre>	<pre>\f la \x endingVPastV4 \alt End \fea \gl +la_sg.fem.past</pre>

Морфологічний аналізатор української мови, розроблений в програмному середовищі PC-KIMMO, дає змогу здійснювати морфологічний аналіз та синтез певного набору іменників та дієслів. У табл. 3 та 4 наведені флексії та суфікси, які дають змогу встановити такі морфологічні характеристики іменників, як група, відміна, відмінок та число [3].

Таблиця 3

Група	Відмінок	Закінчення однини	Закінчення множини	Закінчення однини	Закінчення множини
		Перша відміна		Друга відміна	
тверда	Н.	-а	-и	-0/--	-и/а
	Р.	-и	--	-а	-ів/--
	Д.	-і	-ам	-у	-ам
	З.	-у	--	-о	-ів/и/а
	О.	-ою	-ами	-ом	-ами
	М.	-і	-ах	-і/ові	-ах
	Кл.	-о	-и	-е/у/о	-и/а
м'яка	Н.	-я	-і	-е/я/--	-і/ї/я
	Р.	-і	--	-я/ю	-ів/їв/ь
	Д.	-і	-ям	-ю	-ям
	З.	-ю	--	-я/е/--	-ів/ї/я
	О.	-ею	-ями	-ем/єм/ям	-ями
	М.	-і	-ях	-і/ї/ю	-ях
	Кл.	-е	-і	-е/ю/я	-і/ї/я
мішана	Н.	-а	-і	--	-і/а
	Р.	-і	-ей	-а	-ів/--
	Д.	-і	-ам	-у	-ам
	З.	-у	-ей	-а/е/--	-і/ів/а
	О.	-ею	-ами	-ем	-ами
	М.	-і	-ах	-і/у	-ах
	Кл.	-е	-і	-е	-і/а

Таблиця 4

Відмінок	Закінчення однини	Закінчення множини	Закінчення однини	Закінчення множини
	Третя відміна		Четверта відміна	
Н.	--	-і	-а/я	-а
Р.	-і	-ів/ей	-и	--
Д.	-і	-ям/ам	-і	-ам
З.	-ю	-ів/і	-я	-а
О.	-ом	-ами/ями	-ам/ям	-ами
М.	-і	-ах/ях	-і	-ах
Кл.	-е/и	-і	-я/а	-а
Суфікси	-	-	-ат, - ен, -ят	-ат, - ен, -ят

Загалом дієслова, аналіз яких забезпечує цей морфологічний аналізатор, складаються з основи і закінчення. У табл. 5 наведено перелік флексій, суфіксів та префіксів, які визначають час, особу та число дієслів [3].

Таблиця 5

Час	Особа Число	Дієслова із закінченням 'ять'	Дієслова із закінченням 'уть'	Дієслова із закінченням 'ють'	Дієслова із закінченням 'ать'
1	2	3	4	5	6
Теперішній	1р, sg	-у	-у	-ю	-у
	2р, sg	-иш	-иш	-еш	-иш
	3р, sg	-ить	-ить	-е	-ить

1	2	3	4	5	6
	1р, pl	-имо	-имо	-ємо	-имо
	2р, pl	-ите	-ите	-єте	-ите
	3р, pl	-ять	-уть	-ють	-ать
Минулий	1р, sg	-в	-в	-в	-в
	2р, sg	-в	-в	-в	-в
	3р, sg	-в(М),-ла(Ф), -ло(N)	-в(М),-ла(Ф), -ло(N)	-в(М),-ла(Ф), -ло(N)	-в(М),-ла(Ф), - ло(N)
	1р, pl	-ли	-ли	-ли	-ли
	2р, pl	-ли	-ли	-ли	-ли
	3р, pl	-ли	-ли	-ли	-ли
Майбутній	1р, sg	-му	-му	-му	-му
	2р, sg	-меш	-меш	-меш	-меш
	3р, sg	-ме	-ме	-ме	-ме
	1р, pl	-мемо	-мемо	-мемо	-мемо
	2р, pl	-мете	-мете	-мете	-мете
	3р, pl	-муть	-муть	-муть	-муть
Інфінітив		-ити	-ати	-ити	-ити
Суфікси		-ити, и	-ати, и	-ити, и	-ити, и
Префікси		На -	На -	Ви-	За-

Таблиця 6

**Приклади використання дворівневого морфологічного  
аналізатора української мови**

Словоформа	Результати морфологічного аналізу
‘пишу’ (зміна приголосної ‘с’- >’ш’ перед голосною ‘у’)	<pre>PC-KIMMO&gt;recognize pysu pysu      [pys_u_1 sg present] PC-KIMMO&gt;</pre>
‘піччю’ (подвоєння приголосної ‘чч’->’ч’)	<pre>PC-KIMMO&gt;recognize piCCU piCU      [piC+U_Possesive case,sg]</pre>
‘смажитимемо’	<pre>PC-KIMMO&gt;recognize smaZytymemo smaZytymemo [smaZ+yty_+memo_pl.1 pers. future] PC-KIMMO&gt;</pre>
‘вип’ю’ (префікс ‘ви’, корінь ‘п’, флексія ‘ю’)	<pre>PC-KIMMO&gt;recognize vyp'U vyp+pU    [Dokonanyj vyd_p_u_1 sg present]</pre>
‘випити’	<pre>PC-KIMMO&gt;recognize vypyty vyp+pyty [Dokonanyj vyd_p+yty__infinityv]</pre>
‘телятах’	<pre>PC-KIMMO&gt;recognize telAtah telAtah   [tel+At_+ah_locative case.pl]</pre>
‘лошати’	<pre>PC-KIMMO&gt;recognize loSaty loSaty    [loS+at_+y_Genitive case,sg]</pre>
‘печі’ (чергування голосної ‘і’- >’е’)	<pre>PC-KIMMO&gt;recognize peCi piCi      [piC+il]</pre>
‘сіл’ (чергування голосної ‘е’->’і’)	<pre>PC-KIMMO&gt;recognize sil sel       [sel_sg]</pre>



## Висновки

Розроблення морфологічного аналізатора української мови в програмному середовищі PC-KIMMO є можливим.

Розроблення морфологічного аналізатора полягає у створенні файла правил і файла лексикону.

Доцільно використовувати розроблений морфологічний аналізатор у навчальному процесі.

Доцільно продовжити розроблення морфологічного аналізатора української мови у програмному середовищі PC-KIMMO.

1. Гельбух А. Ф. *Эффективно реализуемая модель морфологии флективного естественного языка / Диссертация на соискание ученой степени кандидата технических наук.* – М., 1994.
2. Antworth E.L. *PC-KIMMO: a two-level processor for morphological analysis. Occasional Publications in Academic Computing No. 16. Dallas: Summer Institute of Linguistics, 1990, 273 p.*
3. Плющ М.Я., Грунас Н.Я. *Грамматика української мови в таблицях.* – К.: Вища школа. – 2004. – С. 5–55 с.
4. Bear J. *A morphological recognizer with syntactic and phonological rules // Proceedings of Coling'86, Association for Computational Linguistics, 1986, p. 272–276.*
5. Jurafsky D., Martin J. *Speech and Language Processing.* – Prentice-Hall, 2000.