

# 多言語機械翻訳を目指したモンゴル語テキスト文法分析

伊・達瓦<sup>1</sup> 呼斯勤<sup>2</sup>、六月<sup>2</sup>、岳耀明<sup>2</sup> 吾浪<sup>3</sup>、白双成<sup>2</sup> 巴特塞恒<sup>2</sup> 井佐原均<sup>1</sup>

<sup>1</sup>情報通信研究機構(NICT) 日本、<sup>2</sup>MENKsoft CO.,LTD 中国、<sup>3</sup>新疆人民放送局蒙古語編輯部 中国

{idawa,isahara}@nict.go.jp

## 1. 序論

モンゴル語はアルタイ語族に属し、語順や単語の構造及び日本語の「て/に/を/は」にあたる語尾助詞や動詞活用形などは日本語とよく似ている。日本語の文を作る感覚で単語を並べて行けばモンゴル語 - 日本語 (逆にしても良い) の文の翻訳やモンゴル語者テキストへの翻字が可能となる。

機械翻訳に用いるモンゴル語テキストは図1に示したように、使用する国や地域によって、異なる文字システムで書かれ、異なる処理システムを使って処理されている。使用者によって縦書きのものや横書きのものがあり、地方方言によって単語や文の構造も異なる場合が多い。したがって、本研究では、それらのテキスト間との通信や変換及び多言語への翻訳を目指して、まずモンゴル語者テキストをローマ字(Latin)表記の中間テキストに書き換え、モンゴル語のコーパスとして利用している。本論文では現代モンゴル語 NM(New Mongolian, 主に、モンゴル国やロシアカリムク、パイラト連邦国使用)のテキストを処理対象として、自立語語尾に接続される助詞の分割ルール及び伝統的モンゴル語 TM (Traditional Mongolian, 主に、中国内モンゴ地域使用) モンゴル語 Todo(中国新疆地域、ロシア国カリムク地域に居住蒙古民が使用)テキストへの翻字における長母音変換ルールについて検討した結果を報告する。

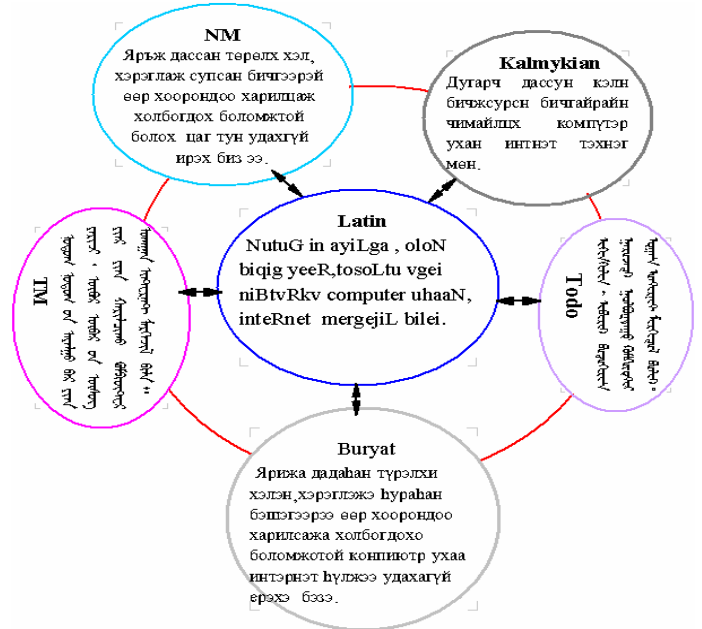


図1 モンゴル語者言語電子化テキスト例

## 2. モンゴル語の構造と文法

### 2.1 モンゴル語者テキストの構造特徴

日本語や韓国語などのアルタイ語族諸言語に見られる「主語、目的語述語」という語順がモンゴル語者テキストにも当てはまり、日本語と英語での語の対応に見られるような語順の入れ替えが、日本語とモンゴル語ではほとんど起こらない。また、モンゴル語の諸テキストは(図2のように)単語ごとに空白によって分かち書きができ、特に、単語の分割は不要であるが、現代モンゴル語の場合は助詞や格助詞などは直前の単語に接続して書く場合が一般的である(図2の実線部分)。単語から助詞を分割しておく、NMから他の言語への翻訳や翻字などの処理には都合がよくなる。特に、TM や Todo テキストへの変換には、これが一つの不可欠な処理であり、研究[1]でも検討されている。

また、分割された部分文字列は直前単語の母音調和によって、例えば「-aar/-oor/-eer/-o' o' r」(日本語格助詞の「で」にあたる)といった4つの交替形を持ち、またNMからTMやTodoでの助詞翻字の場合は、語末文字が母音か子音かによってその表記が異なる。この点は日本語と違う[2]。例えば、NMの文字列「mongoloor」中の「-oor」はそれにあたり、(モンゴル語で)という意味を現しているが、直前の自立語「mongol」の第一音節に陽性母音「o」があるため、母音調和によって「-oor」と書けなければならない。

言語学的及びテキストデータの統一的考察によるNMテキスト自立語に後続する格助詞の分割 変換規則を表1にまとめた[3]。実テキストには本規則に従い、正確に分割されるものもあれば、誤分割されるものも多い。例えば、NMテキストから空白を区切りとして抽出されたある文字列(以後形式単語と呼び、ローマ字表記する)「mongoloor」について、表1のルールを用いて分割すると、前例

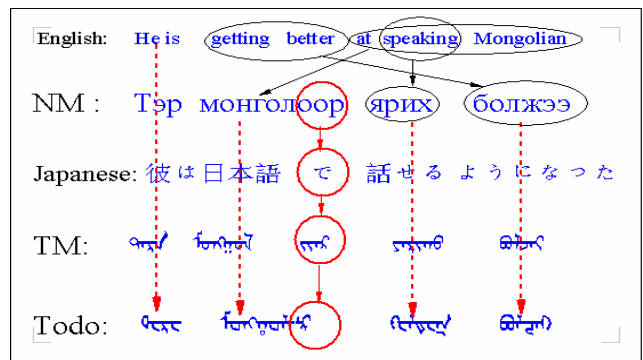


図2 モンゴル語者テキスト文法構造と日本語構造比較

のように、正確な意味を保って分割されるが、同じくNM文字列「Odoovgeer」(今大丈夫)中の形式単語「vgeer」(大丈夫)は「vgeer」方向でまたは「矢」として分割され、元の文字列が正確な意味を失ってしまう。また、NMテキストの「terbawgaasaijzvrhn hagaranvhjee」(彼は熊に怖がらされて心臓が止まった)を、一対一の文字単位でTodoテキストへ翻字した場合は、「tere babugagaasaiji zvrkvi hagarad kvjji」(彼は女房に怖がらされて心臓が止まった)という怪しい意味になってしまう。即ち、言語学的なルールまたは文字単位での翻字だけではかなり限界があると思われる。本研究では、多言語自立語対応辞書[4]、言語的ルール集+DPマッチング手法(2.2節で述べる)及び文字単位変換コーパス照合修正といった総合的な翻字手法を提案する。

表1:NM 形式単語語尾格助詞分割 変更ルール集(表中欄NM final of a word はNM 自立語の語尾文字を、NM suffixは形式単語の語尾文字列を現す)

NM final of a word	NM suffix	TM CP	Todo CP	Jp CP	NM final of a word	NM suffix	TM CP	To do CP	Jp CP
V <sub>i</sub> V <sub>j</sub> V	н	/yin/ ㄣᠨ	/gin/ ᠭᠢᠨ ᠭᠠᠨ ᠮᠠᠨ		V	аар ээр оор еер	/bar/ ᠪᠠᠷ ᠪᠡᠷ	/ber/ ᠪᠡᠷ	で
V	ийн	/yin/ ㄣᠨ	/in/ ᠢᠨ	の	C		/iyar/ ᠶᠠᠷ ᠶᠠᠷ	/yar/ ᠶᠠᠷ ᠶᠠᠷ 1. ㄣ - ㄣ 2. ㄣᠡ, ㄣᠢ	/を
C	хээ	/un/ /vn/ ᠤᠨ / ᠮᠠᠨ	/gin/ ᠭᠢᠨ		V/VV +r		/iyar/ ᠶᠠᠷ	/gar/ /gar/ ᠭᠠᠷ ᠭᠠᠷ ᠭᠠᠷ	
н	ий	/u/ /v/ ᠤ / ᠮᠠ	/ni/ ᠨᠢ		V	аа ээ оо	/ben/ ᠪᠡᠨ ᠪᠠᠨ	/been/ ᠪᠡᠨ ᠪᠡᠨ	を
	ы		/ni/ ᠨᠢ		C	еө жан	/iyen/ ᠶᠡᠨ ᠶᠡᠨ	/yen/ ᠶᠡᠨ ᠶᠡᠨ	/に
	хөө	/yi/ ᠶᠢ	/neen/ ᠨᠡᠨ		VV +r		/iyen/ ᠶᠡᠨ	/gaan/ ᠭᠠᠨ ᠭᠠᠨ	
V	ын	/yin/ ㄣᠨ	/in/ ᠢᠨ		V	аас ээс оос хах	/eje/ ᠡᠵᠡ ᠡᠵᠡ	/eec/ ᠡᠶᠡ ᠡᠶᠡ	~ かゝる よめ
C	ын	/un/ /vn/ ᠤᠨ / ᠮᠠᠨ	/in/ ᠢᠨ		C	хай той гэй гаа	/tai/ ᠲᠠᠢ ᠲᠠᠢ	/tai/ ᠲᠠᠢ ᠲᠠᠢ	と
V	ийг	/yi/ ᠶᠢ	/igi/ ᠶᠢᠭᠢ	を	V	даа доо	/dagan/ ᠳᠠᠭᠠᠨ ᠳᠠᠭᠠᠨ	/daan/ ᠳᠠᠨ ᠳᠠᠨ	に
C	ийг	/i/ ᠶᠢ	/igi/ ᠶᠢᠭᠢ	/のを	C	дээ дөө	/degen/ ᠳᠡᠭᠡᠨ ᠳᠡᠭᠡᠨ	/deen/ ᠳᠡᠨ ᠳᠡᠨ	へ
V	ыг	/yi/ ᠶᠢ	/igi/ ᠶᠢᠭᠢ	/のを	V	гаа тоо	/tagan/ ᠲᠠᠭᠠᠨ ᠲᠠᠭᠠᠨ	/taan/ ᠲᠠᠨ ᠲᠠᠨ	
C	ыг	/i/ ᠶᠢ	/igi/ ᠶᠢᠭᠢ	/のを	C	тээ төө	/tegen/ ᠲᠡᠭᠡᠨ ᠲᠡᠭᠡᠨ	/teen/ ᠲᠡᠨ ᠲᠡᠨ	
н(ng) V <sub>i</sub> V <sub>j</sub>	ыг ийг	/yi/ ᠶᠢ	/gi/ ᠭᠢ		V	гүү гэй	/gei/ ᠭᠡᠢ ᠭᠡᠢ	/gei/ ᠭᠡᠢ ᠭᠡᠢ	ごい
б/г/р с/д/	т	/tu/ /tv/ ᠲᠤ / ᠲᠤ	ᠲᠦ ᠲᠦ	に	C				
V/хөх/д /н(ng)/	д	/du/ /dv/ ᠳᠤ / ᠳᠤ	ᠳᠦ ᠳᠦ		V				
other	нд				C				
н(ng)	ууд	ᠤᠤ / ᠮᠠᠨ	1. ᠤᠤᠣ ᠤᠤᠣ	~ たち	V				
Other C	үүд	ᠤᠤ / ᠮᠠᠨ	2. ᠤᠤᠣ ᠤᠤᠣ	~ ら	C				

NM: New Mongolian; TM: Traditional Mongolian; Todo: Mongolian Todo; V: Vowel; C: Consonant; 1. Positive, 2. Negative; / ~ / : reading; “ - ” : link to pre-word; CP: Case particle; *ij*

2.2 DP マッチングによる単語類似度計算

前述した各文構造特徴によって実際に用いているテキスト形式単語または分割された部分文字列を自立語辞書から直接調べると、マッチングできない場合が多い。本研究では DP(Dynamic Programming) マッチングを適用することによって、高速高類似度単語照合を行う[5]。DP マッチングのアルゴリズムは、次のように定義される。例えば、形式単語や辞書単語それぞれを A, B とし、式(1)で与えられる。(a<sub>i</sub>, b<sub>j</sub> それぞれは A, B 中の i, j 番目の文字を表す)

$$\begin{cases} A = a_1, a_2, \dots, a_I \\ B = b_1, b_2, \dots, b_J \end{cases} \quad (1) \quad (I, J \text{ は単語長となる})$$

次に、文字同士間距離を d<sub>n</sub>(i, j) とし、式(2)で計算する。

$$d_n = \begin{cases} 1 & \text{if } a_i = b_j \\ 0 & \text{else} \end{cases} \quad (2), \quad g_n[0][0] = \begin{cases} 0 & \text{if } A[0] = B[0] \\ 1 & \text{else} \end{cases} \quad (3)$$

$$\text{For } j=1 \text{ to } |B|, \quad g_n[0][j] = g_n[0][j-1] + d_n[0][j]$$

$$\text{For } i=1 \text{ to } |A|, \quad g_n[i][0] = g_n[i-1][0] + d_n[i][0] \quad (4)$$

そこで、単語間正規化最小値関数は(初期値は式(3), (4)とする)下式(5)として定義される。(n は辞書長を示す)

$$g_n(i, j) = \min \begin{cases} g_n(i, j-1) + d_n(i, j) \\ g_n(i-1, j-1) + 2 \times d_n(i, j) \\ g_n(i-1, j) + d_n(i, j) \end{cases} \quad (5)$$

したがって、k 個の候補単語距離を D<sub>min</sub><sup>k</sup> = (1/I+J)g<sub>k</sub>(I, J) とすると、最終的に選択される単語距離は、条件 (C<sub>max\_slope</sub><sup>k</sup> ≥ I/2) によって、下の式(6)として決定される。

$$D_x = \min \{ w_k \times D_{\min}^k(A, B) \} \quad (6)$$

ここで、重み w<sub>k</sub> = C<sub>max\_slope</sub><sup>k</sup> / I は、形式単語長 I 及び k 番目の候補単語が i = j で連続してマッチした部分文字列最大値 C<sub>max\_slope</sub><sup>k</sup> を用いて計算される。

2.3 NM 文字列後続助詞分割 変換

分かち書き NM テキストから空白によって抽出されたある文字列、すなわち、ある形式単語に対して、次の三つのス

トップで分割-翻字処理が行われる。

Step1: ある形式単語に対して、モンゴル語多言語自立語対応電子化辞書 MPEDMCJK (セクション3で紹介する)を用いて、まず、該当文字列は自立語として存在するかどうかを調べる。次に、もし、調べた結果が“yes”なら、該当単語に対応した他言語自立語がそのまま出力される。

Step2: 調べた結果が“no”なら、表1を利用する。そこで、もし、該当する形式単語語尾が表1のルールに従うものであれば、マッチングされた語尾は一旦切り出され、残りの文字列に対して、NM自立語辞書とのDPマッチングが行われ、高類似度のk個の候補単語が選択される。そこで、式(6)を用いて最小距離となる候補が選択され、それに対応した他言語単語が出力される。一方、一旦切り出された語尾は表1に定義された他言語対応格助詞(CP)へ入れ替え、前単語と空白を空けて出力される。例えば、NM形式単語を、“aawd”とすると、Step1では、本文字列が見つからない。そこで、Step2の処理を行うと、“aawd”は“aaw\_d”に書き換え、また、前半部分の文字列“aaw”に対してNM辞書とのDPマッチングを行うと自立語「aaw」が選択され、それに対応したTM単語「abu」が出力される。同時に、単語「abu」の語尾は母音であることから、表1を参照して語尾の「d」は「du」に置き換えられる。

Step3: Step1,2では変換できなかった形式単語においては、日本語への翻訳の場合は該当文字列は未知語として処理される。TMやTodoへの翻字の場合は、次節に述べる長母音変換処理を先に行った後で翻字が行われる。

#### 2.4 モンゴル語諸テキスト長母音変換処理

NMやTodoテキストでは、長母音を表記する専用の音節が存在するが、TM文語では長母音表記音節がない。即ち、諸言語テキスト同士間の変換には長母音翻字は必須となる。

そこで、本研究では、各テキストに対して言語学的ルール考察を行った。NMからTMまたはTodoへの長母音変換ルールを表2,3にまとめた。TM文語では、非規則単語表記(習慣的な表記)が多く見られ、長母音音節を表現しない場合がある。また、NM自立語に後続する助詞は自立語の性(陽、陰性)によって変わることもある。したがって、NMからTMやTodoへの長母音音節変換は必ずしも表2のルール集に従うとは言えない。本研究では、表2,3のルール集、言語コーパス及びDPマッチング手法を同時に参照した処理手法を検討した。長母音変換処理は前節2.3の処理では変換できなかった形式単語において単語ごとに行う。

まず、文節の中に長母音が包含されるかどうかを調べ、もし、包含されるのであれば、文字単位変換テーブル(NMの一文字がTMやTodoの一文字に対応した表、本文中では省略する)と表2,3を用いて、長母音出現位置(Initial)か(medial)かによって他言語文字列へ翻字される。

次に、前処理で翻字された文字列や長母音を包含しない形式単語に対して、翻字先言語大規模テキストから作成した独立単語コーパスとのDPマッチングを行い、類似度が最も高い文字列を最終的な結果として出力する。

#### 2.5 多義語処理対策

日本語からモンゴル語への翻訳の場合、日本語の格助詞「~から」「~より」などは、NMの場合は、母音調りによって「~aas, ~ees, ~oos, ~o`o`s」の4つの形式に対応するが、TMでは、ただ「eqe」、またTodoテキストでは「ecee」

表2 NMからTMまたはTodoテキスト長母音条件付け変換ルール集

long vowel(Initial) in NM		long vowel(Initial) in TM	long vowel(Initial) in Todo
aa	/aa/	/aga/ᠠᠭᠠ, /a/ᠠ	/aa/ᠠᠠ
ee	/ee/	/ege/ᠡᠭᠡ, /e/ᠡ	/ee/ᠡᠡ
oo	/oo/	/ogu/ᠣᠭᠤ	/oo/ᠣᠣ
yy	/yy/	/ugu/ᠤᠭᠤ, /agu/ᠠᠭᠤ	/uu/ᠤᠤ
ий	/ii/	/ii/ᠢᠢ	/ii/ᠢᠢ
ee	/o'o'/	/o'ge/ᠣᠭᠡᠭᠡ	/o'o'/ᠣᠣ
yy	/vvi/	/egvi/ᠡᠭᠦᠢ, /vgvi/ᠦᠭᠦᠢ	/vvi/ᠦᠦ
Eve alphabet	longvowel(medial) in NM	longvowel(medial) in TM	long vowel(medial) in Todo
ж/з	aa	/iga/ᠢᠭᠠ, /iya/ᠢᠶᠠ, /aga/ᠠᠭᠠ	/aa/ᠠᠠ
p/ш		/iga/ᠢᠭᠠ, /iya/ᠢᠶᠠ, /vga/ᠦᠭᠠ	/aa/ᠠᠠ, /ia/ᠢᠠ
other			
ᠠᠠᠨ	ee	/ege/ᠡᠭᠡ, /ebe/ᠡᠪᠡ	/ee/ᠡᠡ, /ei/ᠡᠢ
p/ш		/iyel/ᠢᠶᠡᠯ, /age/ᠠᠭᠡ	
other		/ege/ᠡᠭᠡ, /age/ᠠᠭᠡ, /o'ge/ᠣᠭᠡᠭᠡ, /iyel/ᠢᠶᠡᠯ	
c	ий	/igi/ᠢᠭᠢ	/ii/ᠢᠢ
p/ш	oo	/oga/ᠣᠭᠠ, /aga/ᠠᠭᠠ, /ool/ᠣᠣᠯ, /ogu/ᠣᠭᠤ	/ool/ᠣᠣ
other		/oga/ᠣᠭᠠ, /iga/ᠢᠭᠠ, /ogu/ᠣᠭᠤ, /aga/ᠠᠭᠠ	
c	yy	/agu/ᠠᠭᠤ, /ugu/ᠤᠭᠤ, /uu/ᠤᠤ	/uu/ᠤᠤ, /u/ᠤ
ᠠ	ee	/o'ge/ᠣᠭᠡᠭᠡ, /o'gv/ᠣᠭᠦᠭᠦ, /ege/ᠡᠭᠡ	/o' o' /ᠣᠣ
c		/o'ge/ᠣᠭᠡᠭᠡ, /o'gv/ᠣᠭᠦᠭᠦ, /ege/ᠡᠭᠡ	
p/ш/ш	yy	/egvi/ᠡᠭᠦᠢ, /vgvi/ᠦᠭᠦᠢ, /vgvi/ᠦᠭᠦᠢ	vvi/ᠦᠦ, /v/ᠦ
ж/з		/egvi/ᠡᠭᠦᠢ	
other		/egvi/ᠡᠭᠦᠢ, /vgvi/ᠦᠭᠦᠢ, /vgvi/ᠦᠭᠦᠢ	

表3 NMからTMまたはTodoテキスト無条件変換ルール集

直接変更 NM	TM	Todo
aa	/yaga/ᠶᠠᠭᠠ	/yaa/ᠶᠠᠠ, /a/ᠠ
ee	/yage/ᠶᠠᠭᠡ	/yee/ᠶᠡᠡ, /e/ᠡ
oy	/yagu/ᠶᠠᠭᠤ	/yuu/ᠶᠠᠤ, /u/ᠤ
oy	/yagv/ᠶᠠᠭᠦ	/yvv/ᠶᠠᠦ, /v/ᠦ
ee	/yogel/ᠶᠣᠭᠡᠯ	/yool/ᠶᠣᠣᠯ, /yee/ᠶᠡᠡ

に対応する。また、NMの一つの固有名詞がTMの複数の固有名詞(同音異義語)に対応する場合もある。例えば、図3の例のように、NM固有名詞の「/sara/」(日月の「月」、天体の「月」)はTM固有名詞の「/sara/ (ゲツ)又は/sar\_a/ (ツキ)」に対応する。今後はこのような語に対する翻字において、文字コード(Unicode)の差違を用い、大規模なテキストデータから作成された同音異義語コンテキストコーパスを利用した処理方法を検討したい。

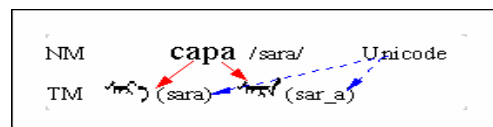


図3 TMテキストに見られる同音異義語とUnicode差違例

### 3. 多言語自立語対応 品詞付け電子化辞書

モンゴル語自然言語処理用大規模コーパスの構築及び多言語機械翻訳システムの実現を目指して、日本独立行政法人情報通信研究機構自然言語グループと中国内蒙古自治区社会科学院モンゴル語情報処理センター(MIT), 内蒙古 MENKsoftCO, .LTD 社及び韓国 KAIST コンピュータ情報処理所は 2005 年から共同研究プロジェクトを実施し、この一年中は現代モンゴル語自立語 5 万個を軸とした“多言語対応品詞付け電子化辞書 MPEMDCJK”を構築した。図 4,5 それぞれは MPEMDCJK 構造及びフォーマット設計を示す。本研究で構築する MPEMDCJK 辞書は 2006 年 4 月に最終的な評価を行い、6 月から公開される予定である[4]。



図 4 MPEMDCJK 構造

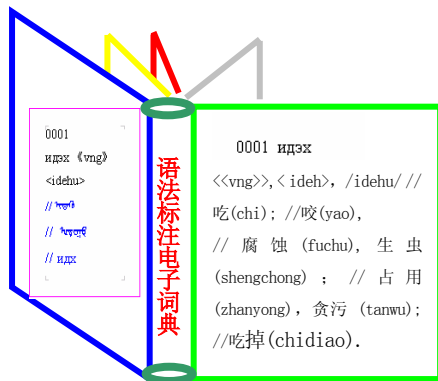


図 5 MPEMDCJK フォーマット設計

### 4. 評価

#### 4.1 実験 1. NM テキスト単語後続助詞分割実験

実験 1 では、MENKsoft 社提供の現代モンゴル語テキスト (19960 語) を利用した。実験は表 4 に示すように、3 つの段階に分けて行った。表 4 の分割単語数はテキスト単語語尾が表 1 の(NM suffix)各項目に存在した形式単語数を、D (有、なし)は MPEMDCJK 辞書を利用した又は利用しない場合を表す。今回の実験評価は本手法によって処理された結果を、MENKsoft 社及び MIT 本プロジェクトメンバーが人手でチェックすることによって行った。

実験表 4 NM テキスト分割実験語数 19960 (\*MPEMDCJK)

	D*なし DP 適用なし	D* 有 DP 適用なし	D* 有 DP 適用
分割単語数	8972	5211	5211
正解率%	64.3	82.6	91.3

4.2 実験 2. NM テキスト→TM または Todo テキスト変換実験 NM→TM, Todo 変換実験は、上述実験 1 では変換が行われなかった全形式単語において実行された。DP マッチング言語情報コーパスとしては、本研究で作成した TM テキストから作成した 28000 語彙、Todo から作成した 34 127 語彙を用いた。データは MENKsoft 社及び MIT が提供した。評価を表 5 に示した。

表 5 NM→TM または Todo 変換実験 (DP マッチング適用)

NM 単語数	T M	Todo
8190 ->	7289	7486
正解率 (%)	89	91.4

4.3 実験結果分析：実験結果の表 4,5 では、辞書及び DP マッチングを適用しない場合と比べ、辞書、コーパス及び DP を適用した場合は高精度な分割や翻字ができることが分かる。特に Todo への変換では最も高い変換結果が得られた。それは、Todo テキストでは、助詞や長母音の表現は NM とあまり差がないためであると思われる。今回の実験では、辞書やコーパスではマッチングできなかった多義語や活用形語尾動詞は文字単位で翻字されている。

### 5. おわりに

本論文では、多言語情報コーパス、言語学文法ルール及び文字列間 DP マッチング最小距離類似度を用いたモンゴル語多文字システムテキスト処理変換手法を提案した。提案手法を利用した今回の実験 (NM テキストから TM 及び Todo テキストへの翻字) では、TM テキストへの変換精度は 89% で、Todo テキストへの場合は 91.4% となり、ほぼ実用レベルに達していることが確認できた。

本システムはモンゴル語を使用する地域や国間での文書情報交換には大きな貢献となる一方、モンゴル語諸テキストの自動編集-修正、品詞情報付け言語コーパスの作成及びモンゴル語自然言語処理、機械翻訳システムの実現に有用と期待される。

モンゴル語の任意のテキスト間の変換、多義語処理対策、実用化システムの実装及び効果的な評価手段の検討などは今後の研究課題となる。

### 参考文献

- [1] 満 都拉、石川 徹也ら、“伝統的モンゴル語と現代モンゴル語の双方向的な翻字”、言語処理学会第 11 回年次大会 発表論文集 pp360-363(2005.3).
- [2] Galcan funcag “Mongol ulassin kiril vsygiin dvrem”, (現代モンゴル語文法)、中国内蒙古人民出版社 2001, Hohhot. (蒙古語)
- [3] Erden Bagan, Gongqig sorung, “蒙古語文法”, 中国内蒙古教育出版社 1998.5
- [4] 伊達瓦、井佐原 均ら、“モンゴル言語 文字自動処理”、中国語情報処理学会誌、2005.11 採録論文 4232.
- [5] Masek, W., and Paterson, M “A faster algorithm computing string edit distance” J. Comput. System Sci. 20(1980), 18-31.