

ローマ字・カタカナ・キリル文字による アイヌ語 Universal Dependencies の可能性

安岡 孝一 (京都大学)[†]

安岡 素子 (京都大学・京都外国語大学)

Universal Dependencies for Ainu Language in Latin Alphabet, Katakana, and Cyrillic

Koichi YASUOKA (Kyoto University)

Motoko YASUOKA (Kyoto University / Kyoto University of Foreign Studies)

要旨・既発表の有無

書写言語としてのアイヌ語は、ローマ字 (ラテンアルファベット)・カタカナ・キリル文字など、多彩な文字と記法によって記述されてきた。その一方、抱合語としてのアイヌ語は、日本語や欧米諸語とは全く異なる言語構造を持つことから、これらの言語向けの言語処理手法は、そのままではアイヌ語に適用できない。ならば Universal Dependencies は、どうだろう。言語横断的な文法構造記述として設計された Universal Dependencies は、書写言語としてのアイヌ語を、どの程度ちゃんと記述できるのだろうか。『アイヌ神謡集』、『アイヌ語會話字典』、アイヌ語訳『五倫名義解』、『Аинско-русский словарь』を Universal Dependencies コーパスとして記述していく中で、われわれは、われわれの見積りが甘かったことを痛感すると同時に、それでも、アイヌ語 Universal Dependencies が、アイヌ語の言語処理に寄与することを確認した。本発表では、その一端について述べる。

本発表は、人文科学とコンピュータ研究会 (第 131 回) の発表「ローマ字・カタカナ・キリル文字併用アイヌ語 RoBERTa・DeBERTa モデルの開発」(c) 情報処理学会を拡張したものである。

1. アイヌ語 Universal Dependencies の概要

アイヌ語 Universal Dependencies は、Senuma and Aizawa (2017, 2018) がローマ字 (ラテンアルファベット) 版を開発⁽¹⁾し、安岡 (2021a,b) がカタカナ⁽²⁾・キリル文字への拡張をおこなった。土台となった Universal Dependencies (UD) は、書写言語における品詞・形態素属性・依存構造 (係り受け関係) を、言語に関わらず記述する手法である [Marneffe et al. (2021)]。句構造を考慮せずに係り受け関係を記述することで、言語横断性を高めており、全ての文法構造を単語間のリンクで記述するのが特徴である。

依存構造解析それ自体は、Tesnière (1959) の構造的統語論に源を発し、Mel'čuk (1988) の有向グラフ記述によって、一応の完成を見た手法である。その最大の特長は、いわゆる動詞中

[†] yasuka@kanji.zinbun.kyoto-u.ac.jp

⁽¹⁾ 瀬沼らによる公開は、知里 (1923) 『アイヌ神謡集』の「ホテナオ」のみだった。

⁽²⁾ われわれがカタカナ表記アイヌ語のデジタル化にかかわったのは、佐藤 (1996) が発端であり、それは JIS X 0213:2000 へのアイヌ語表記用カタカナ追加として結実した。

表 1 CoNLL-U の各フィールド

1. ID: 単語ごとに付与されたインデックスで、文ごとに 1 から始まる整数。縮約語に対しては、単語の範囲を示すのも可。
2. FORM: 語、または、句読記号。
3. LEMMA: 基底形、語幹。
4. UPOS: UD で規定された言語普遍的な品詞タグ (表 2)。
5. XPOS: 言語固有の品詞タグ。
6. FEATS: UD で規定された言語普遍的な形態素属性のリスト。言語固有の拡張も可。
7. HEAD: 当該の単語の係り受け元 ID。係り受け元が無い場合は 0 とする。
8. DEPREL: UD で規定された言語普遍的な係り受けタグ (表 3)。HEAD が 0 の場合は root とする。言語固有の拡張も可。
9. DEPS: 複数の係り受け元を持つ場合、全ての HEAD:DEPREL ペア。
10. MISC: その他のアノテーション。

表 2 UD 品詞タグ (UPOS)

Open class words	Closed class words	Other
ADJ 形容詞	ADP 側置詞	PUNCT 句読点
ADV 副詞	AUX 助動詞	SYM 記号
INTJ 感嘆詞	CCONJ 並列接続詞	X その他
NOUN 名詞	DET 限定詞	
PROPN 固有名詞	NUM 数詞	
VERB 動詞	PART 接辞	
	PRON 代名詞	
	SCONJ 従属接続詞	

表 3 UD 係り受けタグ (DEPREL)

	Nominals	Clauses	Modifier Words	Function Words
Core arguments	nsubj 主語 obj 目的語 iobj 間接目的語	csubj 節主語 ccomp 節目的語 xcomp 節補語		
Non-core dependents	obl 斜格補語 vocative 呼称語 expl 形式語 dislocated 外置語	advcl 連用修飾節	advmod 連用修飾語 discourse 談話要素	aux 動詞補助成分 cop 繫辞 mark 標識
Nominal dependents	nmod 体言による連体修飾語 appos 同格 nummod 数量による修飾語	acl 連体修飾節	amod 用言による連体修飾語	det 決定語 clf 類別語 case 格表示
Coordination	MWE	Loose	Special	Other
conj 接続 cc 接続語	fixed 固着 flat 並列 compound 複合	list 細目 parataxis 隣接表現	orphan 親なし goeswith 泣き別れ reparandum 言い損じ	punct 句読点 root 親 dep 未定義

# text = kamuy tura okay=an									
1	kamuy	kamuy	NOUN	名詞	-	3	obl	-	-
2	tura	tura	ADP	後置副詞	-	1	case	-	-
3	okay	okay	VERB	自動詞	-	0	root	-	SpaceAfter=No
4	=an	=an	PART	人称接辞	-	3	nsubj	-	-

# text = カムイ トウラ オカヤン									
1	カムイ	kamuy	NOUN	名詞	-	3	obl	-	-
2	トウラ	tura	ADP	後置副詞	-	1	case	-	-
3-4	オカヤン	-	-	-	-	-	-	-	-
3	オカイ	okay	VERB	自動詞	-	0	root	-	-
4	アン	=an	PART	人称接辞	-	3	nsubj	-	-

# text = камуй тура okayн									
1	камуй	kamuy	NOUN	名詞	-	3	obl	-	-
2	тура	tura	ADP	後置副詞	-	1	case	-	-
3-4	okayн	-	-	-	-	-	-	-	-
3	okay	okay	VERB	自動詞	-	0	root	-	-
4	ан	=an	PART	人称接辞	-	3	nsubj	-	-

図1 アイヌ語 UD の CoNLL-U データ

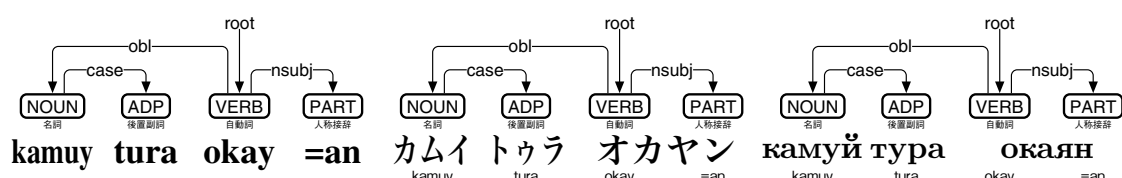


図2 deplacy によるアイヌ語 UD の可視化

心主義によって言語横断的な記述が可能だという点にあり、Mel'čuk (1988) 依存文法をコンピュータ向けに洗練した UD においても、言語に関わらない記述、という特長が前面に押し出されている。UD における文法構造記述は、句構造を考慮せず、全てを単語間のリンクとして表現する。これにより、言語横断的な文法構造記述を可能としている。

UD 係り受けコーパスの交換用フォーマットとして、CoNLL-U と呼ばれるタブ区切りテキスト (文字コードは UTF-8) が規定されている。CoNLL-U の各行は各単語に対応しており、表 1 に示す 10 個のタブ区切りフィールドで構成される。ID・FORM・LEMMA は、単語そのものに関するフィールドである。UPOS・XPOS・FEATS は、単語の品詞と形態素属性に関するフィールドである。HEAD・DEPREL・DEPS は、単語の係り受けに関するフィールドである。

UDにおける係り受け関係は、単語間の有向グラフを HEAD と DEPREL で記述する。HEAD は、その単語に入る有向枝のリンク元 ID を示しており、DEPREL は、その有向枝における係り受けタグである。ただし、HEAD が 0 の場合、その枝に入るリンク元は存在しない。リンクの本数は単語の個数に等しく、各リンクのリンク先は、全て互いに異なっている。すなわち、各単語から出るリンクは複数の可能性があるが、各単語に入るリンクは 1 つだけである。なお、リンクはループしない。

UD の係り受けリンクは、Mel'čuk (1988) 依存文法の後裔にあたり、いわゆる動詞中心主義である。動詞をリンク元として、主語や目的語へとリンクする。修飾関係においては、被修飾語から修飾語へとリンクする。ただし、側置詞 (前置詞や後置詞) を体言の修飾語だとみなす [Nivre (2015)] 点が、Mel'čuk (1988) とは異なっている。また、コピュラ文においては動詞中心主義を採らず、補語をリンク元として、主語や繫辞へとリンクする。

UD は単語長を規定しておらず、各言語ごとに、自由に単語長を決めることができる。アイヌ語 UD では、田村 (1996) 『アイヌ語沙流方言辞典』を、作業上の単語認定に用いている。なお、接尾辞・接頭辞については、人称接辞と動名詞接尾辞 (-i と -p) だけを語とみなし、それ以外は前後の語にくっ付けている。

アイヌ語 UD の例として、「kamuy tura okay=an」「カムイ トウラ オカヤン」「камуй тура окаян」の CoNLL-U データを図 1 に示す。LEMMA と XPOS は『アイヌ語沙流方言辞典』に従っている⁽³⁾。また、これらの CoNLL-U を比較すべく、deplacy [安岡 (2020)] で可視化した (図 2)。UD 依存構造は全く同一だが、「オカヤン」や「окаян」は、文字の途中に単語境界がある点に注意されたい。

2. アイヌ語 Universal Dependencies コーパスの作成

ローマ字・カタカナ・キリル文字で書かれたアイヌ語文書に対し、係り受け解析エンジン esupar のアイヌ語 DeBERTa モードで仮コーパスを作成し、その結果をアイヌ語 UD エディターで編集する、という手順で、アイヌ語 UD コーパスを作成した。以下、それぞれのアイヌ語 UD コーパスについて、概要を述べる。

2.1 アイヌ神謡集

知里 (1923) 『アイヌ神謡集』は、本文 124 ページに 13 編のアイヌ神謡を収録しており、見開き左ページ (偶数ページ) にローマ字で書かれたアイヌ語を、見開き右ページ (奇数ページ) に日本語訳を配置している (図 3)。各編の構成は以下のとおり。

1. 「銀の滴降る降るまはりに」 (11 ページ 230 行) × 2
2. 「トワトワト」 (7 ページ 136 行) × 2
3. 「ハイクンテレケ ハイコシテムトリ」 (6 ページ 121 行) × 2
4. 「サンパヤ テレケ」 (5 ページ 104 行) × 2
5. 「ハリツクンナ」 (4 ページ 83 行) × 2

⁽³⁾ アイヌ語 UD の XPOS では、固有名詞を名詞から分離し、数詞を連体詞から分離した上で、複他動詞を他動詞に統合し、さらに記号を加えた [安岡 (2021b)]。

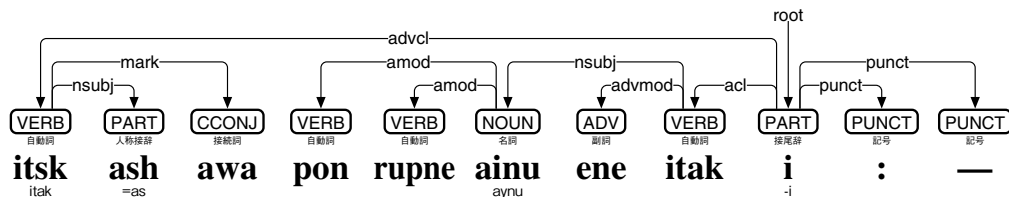
tapan petpo teeta rehe tane rehe
ukaepita eki kushnena.”
hawash chiki itakash hawe ene okai :—
“Nennamora tapan petpo teeta rehe
tane rehe erampeuteka!
teeta kane shinnupurita tapan yetpo
‘Kanchiwetunash’ ari ayea korka
tane shirpan kushu ‘Kanchiwemoire’ ari
aye ruwe tashi anne!”
itskash awa ponrupneainu ene itaki :—
“Pii tuntun, pii tun tun!”
connobetapne ehawan chiki,
ushinritpita aki kushnena!”
hawash chiki itakash hawe eneokai :—
“Nennamora eshinrichihi erampeuteka!
otteeta Okikirmui kimta oman wa,
kucha karita keneinunpe kar aike
ne inunpe apekar wa sattek okere,
Okikirmui oararkehe oterke ko oararkehe
hotari. Nawaanpe Okikirmui rushka kushu
ne inunpe pet otta kor wa san wa,
oshura wa isam ruwe ne.
Orowano ne inunbe petesoro mom aineno,
atuioero oslma, tu atuipenrur re atuipenrur
chieshirkik shiri kamuiutar nukar wa,

此の川の前の名と今の名を
言つて見ろ。」
聞くさ、私の言ふことには
「誰が此の川の前の名
今の名を知らないものか!
昔、えらかつた時代には此の川を
流れの早い川と言つてゐたのだが
今は世が衰へてゐるので流れの遅い川と
言つてゐるのさ。」
云ふさ小男の云ふことには
「ビイトントンビイトントン
本當にお前そんな事を云ふなら
お互の素性の解合ひをやらう。」
聞いて私の云ふことには
「誰がお前の性素を知らないものか!
大昔、オキ、リムイが山へ行つて
狩獵小屋を建てた時様の木の爐縁を作つたら
その爐縁が火に當つてうからに乾いてしまつた。
オキ、リムイが片方を踏むと片一方が
上る、それをオキ、リムイが怒つて
其の爐縁を川へ持つて下り
捨てしまつたのだ。
それから其の爐縁は流れに沿ふて流れていつて
海へ出で、彼方の海此方の海波
に打つけられる様を神様たちが御覧になつて、

図3 『アイヌ神謡集』「ホテナオ」70～71 ページ

text = itskash awa ponrupneainu ene itaki : —

1	itsk	itak	VERB	自動詞	—	9	advcl	—	SpaceAfter=No
2	ash	=as	PART	人称接辞	—	1	nsubj	—	—
3	awa	awa	CCONJ	接続詞	—	1	mark	—	—
4	pon	pon	VERB	自動詞	—	6	amod	—	SpaceAfter=No
5	rupne	rupne	VERB	自動詞	—	6	amod	—	SpaceAfter=No
6	ainu	aynu	NOUN	名詞	—	8	nsubj	—	—
7	ene	ene	ADV	副詞	—	8	advmod	—	—
8	itak	itak	VERB	自動詞	—	9	acl	—	SpaceAfter=No
9	i	-i	PART	接尾辞	—	0	root	—	—
10	:	:	PUNCT	記号	—	9	punct	—	—
11	—	—	PUNCT	記号	—	9	punct	—	—



6. 「ホテナオ」(3 ページ 66 行) × 2
7. 「コンクワ」(6 ページ 125 行) × 2
8. 「アトイカトマトマキ、クントテアシフム、フム!」(8 ページ 193 行) × 2
9. 「トーロロハンロクハンロク!」(2 ページ 43 行) × 2
10. 「クツニサクトンクトン」(2 ページ 30 行) × 2
11. 「此の砂赤い赤い」(3 ページ 70 行) × 2
12. 「カツパレウレウカツパ」(3 ページ 53 行) × 2
13. 「トヌペカランラン」(2 ページ 43 行) × 2

『アイヌ神謡集』のローマ字表記は『アイヌ語沙流方言辞典』と異なっており、また、誤植も散見される⁽⁴⁾。われわれのアイヌ語 UD コーパスでは、『アイヌ神謡集』の表記をそのまま FORM に入れ、LEMMA と単語長を『アイヌ語沙流方言辞典』に合わせた。単語長の差は、MISC の SpaceAfter=No で吸収した。たとえば、図 3 左 70 ページ 10 行目「itskash awa ponrupneainu ene itaki : ー」に対しては、FORM は誤植も含めてそのままとし、LEMMA は「itak =as awa pon rupne aynu ene itak -i : ー」としている(図 4)。なお、『アイヌ語沙流方言辞典』に見当たらない単語については、基本的に片山(2003)の単語認定に依っている。

2.2 アイヌ語會話字典

神保・金澤(1898)『アイヌ語會話字典』は、本文 2 段組 278 ページの段組左側に日本語を、右側にローマ字のアイヌ語訳を配置している(図 5)。Bugueva(2011)は『アイヌ語會話字典』を拡張する形で、トピック別アイヌ語會話辞典(全 3847 見出し)を公開している。

『アイヌ會話字典』のローマ字表記は『アイヌ語沙流方言辞典』と異なっており、特に単語長の認定が全く違う。われわれのアイヌ語 UD コーパスでは、『アイヌ會話字典』の表記をそのまま FORM に入れ、LEMMA と単語長を『アイヌ語沙流方言辞典』に合わせた。単語長の差は、FORM 中の空白や、MISC の SpaceAfter=No で吸収した。たとえば、図 5 右側 18~19 行目「Tambeta ne shomo k'eiwange.」に対しては、FORM は「ta ne」に空白を含みつつ、LEMMA は「tan pe tane somo k= eywanke .」としている(図 6)。

2.3 Аинско-русский словарь

Добротворский(1875)アイヌ語・ロシア語辞典の補遺第 12 章(図 7)には、キリル文字で書かれた樺太アイヌ語の対話文が収録されている[寺田・安田(2019)]。この対話文については、阪口(2021)によるローマナイゼーションと日本語訳、および詳細な解説があり、これを参照しつつアイヌ語 UD コーパスの作成をおこなった。図 7 左ページ本文 3 行目「Танъ котанъ охтá утáса—анъ кусý áреги анъ.」に対するアイヌ語 UD を、図 8 に示す。なお、阪口(2021)は「охтá」のローマナイゼーションを「ohtà」としているが、われわれは『アイヌ語沙流方言辞典』に合わせて「or ta」とした。樺太アイヌ語を沙流アイヌ語に合わせるかどうかについては、もちろん議論の余地があると考えられる。

⁽⁴⁾ 図 3 の左 70 ページには、6 行目「yetpo」→「petpo」、10 行目「itskash」→「itakash」、12 行目「eonnohetapne」→「sonnohetapne」の誤植がある[佐藤(2004)]。右 71 ページには、15 行目「性素」→「素性」、18 行目「らうから」→「からから」の誤植がある。

(21)	
イノル(祈)	Inonno-itak.
イボ(疣)	Erum-tambu.
イバラ(棘)	Ai-ush-ni; Ai-o-ni.
イビキ(鼾)	Etoro.
イマ(今)	Tane; Tanepo.
今馬で来たところだ	Tanepo ku unima o wa k'ek na.
今参りました	Tane ariki an ruwe ne.
イモ(芋)	Emó; Chiurip.
イモート(妹)	Mataki; Matapa; Tureshpo.
イヤ	Kopan; Kochan.
こんな物は己は厭だ	Tambe ne no ambe onakne ku kopan.
あの人は厭ひだ	Nei aiuu ku etunne.
イリケチ(入口)	Soigeta., Apa-ushta.
イル(入用)	Eiwange.
これはもーいらない	Tambeta ne shomo k'ei- wange.
イル(射)	Tukan; Ak.

図5 『アイヌ語會話字典』21 ページ

text = Tambeta ne shomo k'eiwange.

1	Tam	tan	DET	連体詞	-	2	det	-	SpaceAfter=No
2	be	pe	NOUN	形式名詞	-	6	obj	-	SpaceAfter=No
3	ta ne	tane	ADV	副詞	-	6	advmod	-	-
4	shomo	somo	ADV	副詞	-	6	advmod	-	-
5	k'	k=	PART	人称接辞	-	6	nsubj	-	SpaceAfter=No
6	eiwange	eywanke	VERB	他動詞	-	0	root	-	SpaceAfter=No
7	.	.	PUNCT	記号	-	6	punct	-	-

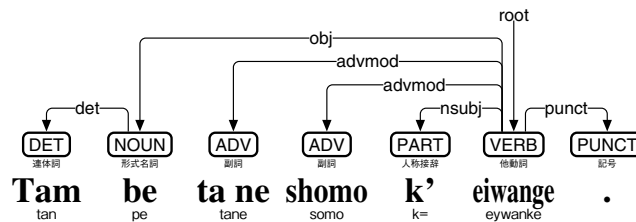


図6 「Tambeta ne shomo k'eiwange.」のアイヌ語 UD

12. СЛОНЕНИЯ И СПРАЖЕНИЯ. ПЕРЕСТАНОВКА СЛОВЪ. ЧАСТИЦЫ И ИХЪ УПОТРЕБЛЕНИЕ.

1. Рѣчь Чивоканке къ о. Симеону, приѣхавшему исповѣдывать Кусунайскую команду и уѣзжавшему.

Танъ котанъ охъ утаса—анъ кусу ареги анъ. Охъ-роно охъ анъ тренькайнъ, тѣе танъ котанъ та охъ. Танъ котанъ та охъ тренькайнъ, анъ пукара. Трѣкоро иту—акскара, трѣкоро анской прѣйки. Тамъ ипхонъ, трѣкоро ий прѣйкере анъ. Тамъ кусу пирка на охъичи—пирка. Ипхонъ охъ кусу—кара.

«И прѣйти въ это село погостить. Прожить я долго благополучно, теперь прѣхать сюда ты. Прибыть ты сюда благополучно; я тебя увидалъ. Мы очень знакомы, и я очень тебя благодарю. Ты далъ мнѣ бисеру. Весьма благодарю тебя. За это говорю тебѣ „счастливый путь“, говорю отъ души. Самъ я отправлюсь послѣ».

2. Душъ ссорящихся (утуа—айну) и желающаго мириться.

— Анокай анъ тренькайнъ, ипхонъ ѣмъ ка хѣнъ. Наванъ кусу поно—нишка уранкара кара э и ки карачики, ороано укораму пирка анъ ки нанго. Иртасиано укораму укораму (или укораму укораму, или укораму—укораму) анъ ки нанго. Нанго ороано укораму, хаманъ ки нанго, нанго—короне—но, поно—поно—нишка уранкара кара э и ки карачики укораму пирка анъ ки. Нанго ороано укораму, анъ уо хемакаре, анъ ки нанго. Ке!

Отвѣтъ обиженаго нѣсколько, но также желающаго мириться.

— Сопнока а ийурамукосма. Энъ эъ ка э ки кусу по ай, тамъ соно айю ахтунъ, анъ ки кусу. Манинъ ханъ перъ—анъ итакъ, хаманъ кунъ ги кусу. Хамеуъ по-

но—нишка уранкара кара анъ оваракара, наванъ кусу та-не ороано укораму, хаманъ ки кунинъ, ороно—такъ ги, укораму хайта, хаманъ ки нанго—короне по. Сне утахта по кусу, порово хамеиакка пирка нанго—короне по. Кван ороано поно—поно ку ки киту эми. Нанъ ийи харахъ кано ка уранго—пононо хаме—нишка пирка нанго—короне—но. Нанго ороано иртасиано укораму пирка анъ кичики пирка нанго—короне—но. Нанго хемака.

Чивоканке на эти рѣчи замѣчаетъ: «анъ—уко—хеманаре—кики, пирка нанго». То есть: «если помирится, то вѣроятно будетъ ладно». Значить, одному нужно „помириться“. — Мы согласны не ссориться между собою. Поэтому если немного взаимно облизаться, то вѣроятно наступитъ взаимное согласіе, вѣроятно мы продолжимъ другъ къ другу дорогу. Отмытъ конечно не будетъ разладина, а будетъ, когда мы немного уступимъ другъ другу, взаимный миръ. Взаимное озлобленіе исчезнетъ. Иду! (?)».

(1) Прим. ред. Отвѣтъ обиженаго не переводитъ авторъ словаря, смѣтъ котораго не далъ ему возможности ослѣпить извѣстныя въ приложеніи къ Словарю статьи, кромѣ первой, и даже назвать ихъ, какъ видно изъ печатаемаго на слѣдующей страницѣ «Оглавленіи» того, что онъ предполагалъ вставить въ 415—420-й рукописи, озаглавленной такъ «Материалы для изученія Айну и ихъ языка». Разборъ сочиненія Пичеловъ, замечательный жемчужина № 4 приложенія къ Словарю, составляетъ отдѣльную статью, помещенную не въ «материалахъ». Въ «материалахъ» же между остальными слѣдующими за «Разборомъ» статьи большое количество невошедшихъ листовъ, не успѣвшихъ принести на себя другіе «материалы» равно утѣшеннаго трудомъ науки.

図7 『Аинско-русский словарь』補遺第12章

text = Танъ котанъ охъ утаса—анъ кусу ареги анъ.

1	Танъ	tan	DET	連体詞	—	2	det	—	—
2	котанъ	kotan	NOUN	名詞	—	3	nmod	—	—
3	ох	or	NOUN	位置名詞	—	5	obl	—	SpaceAfter=No
4	та	ta	ADP	格助詞	—	3	case	—	—
5	утаса	u-tasa	VERB	自動詞	—	9	advcl	—	SpaceAfter=No
6	—	-	PUNCT	記号	—	5	punct	—	SpaceAfter=No
7	анъ	=an	PART	人称接辞	—	5	nsubj	—	—
8	кусу	kusu	CONJ	接続助詞	—	5	mark	—	—
9	ареги	ar-iki	VERB	自動詞	—	0	root	—	—
10	анъ	=an	PART	人称接辞	—	9	nsubj	—	SpaceAfter=No
11	.	.	PUNCT	記号	—	9	punct	—	—

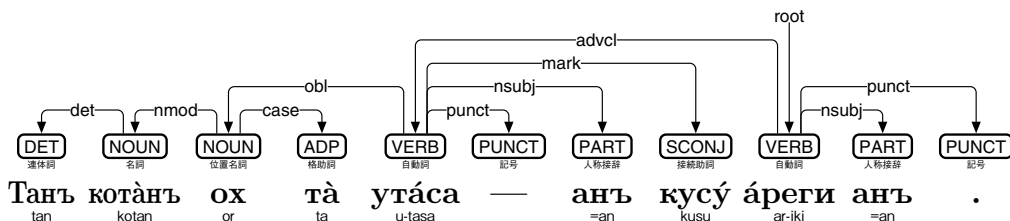


図8 「Танъ котанъ охъ утаса—анъ кусу ареги анъ。」のアイヌ語UD

2.4 アイヌ語訳『五倫名義解』

加賀家文書館(別海町)所蔵のアイヌ語訳『五倫名義解』(整理番号 K3-21)は、室・空谷(1855, 1858)『五倫名義解』に加賀伝蔵がアイヌ語訳を施したもので、文久～慶応年間に書かれたものである[深澤(2014a)]。以下に示す5章と刊記で構成される。

1. 「父子有親」6 ページ 12 文×2
2. 「君臣有義」6 ページ 14 文×2
3. 「夫婦有別」6 ページ 13 文×2
4. 「長幼有序」7 ページ 13 文×2
5. 「朋友有信」8 ページ半 13 文×2
6. 刊記 5 ページ半 11 文×2

各ページには、日本語が3行ずつ書かれており、その横にカタカナでアイヌ語訳が記されているが、各章の表題はアイヌ語訳されていない(図9)。アイヌ語訳に小書きカタカナは使われておらず、拗音も促音も小書きにしない上、末子音が母音を伴って書かれている。深澤(2014b)が指摘するとおり、母音の混同(イとエ、ウとヲ)も散見される。しかも書き直しが多く、非常に読みにくい。図9のアイヌ語訳に対するアイヌ語 UD を、図10に示す。ただし、「アルシヤナ」に「earsayne」を当てていいのか、「ウバカシ」は「uwepakasnu」なのか、など多くの疑問点が残っており[安岡・安岡(2023)]、現在も引き続き作業中である。

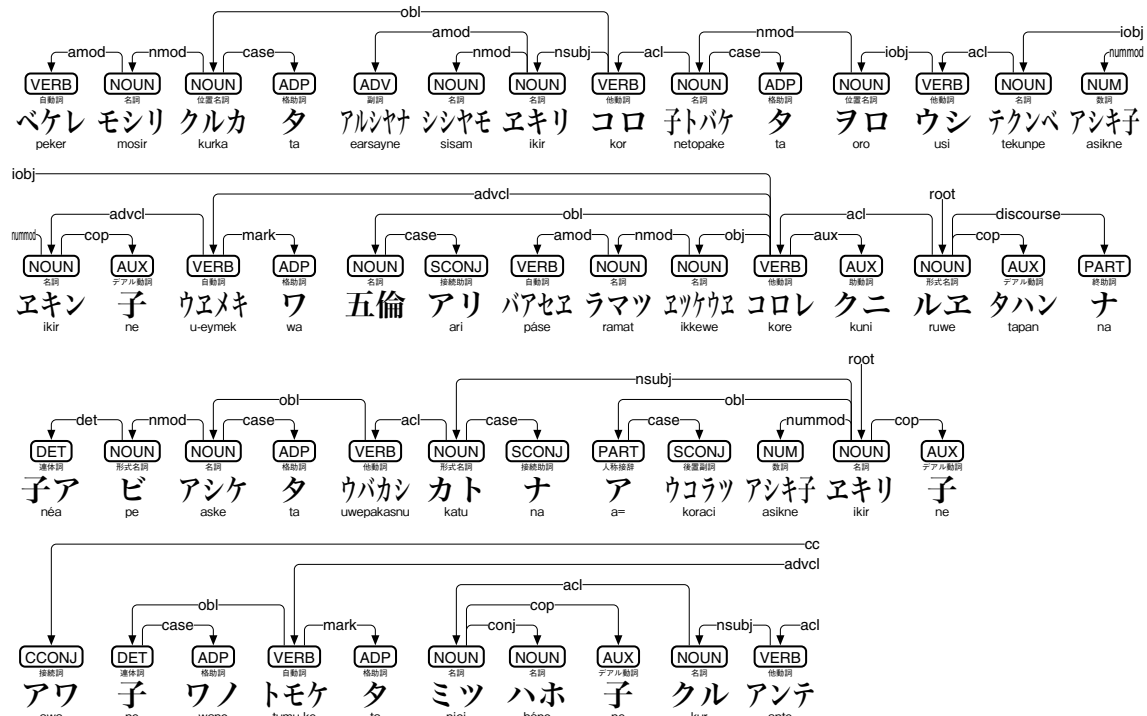


図10 アイヌ語訳『五倫名義解』冒頭部のアイヌ語 UD

2.5 国立アイヌ民族博物館ガイドブック

ウポポイ (民族共生象徴空間) はアイヌ語を第一言語としており [小林 (2023), 深澤 (2023)], その方針は、国立アイヌ民族博物館 (2020a,b, 2021) にも踏襲されている。



図 11 『国立アイヌ民族博物館ガイドブック』 21～22 ページ

国立アイヌ民族博物館 (2020a) 『国立アイヌ民族博物館ガイドブック』は、見開き左ページ (奇数ページ) にアイヌ語で解説を書き、見開き右ページ (偶数ページ) にその日本語訳と英語訳を載せる、という方針で編集されている (図 11)。アイヌ語はカタカナで書かれており、表題に限ってローマ字が添えられている。図 11 の「カムイ トゥラ オカヤン」「kamuy tura okay=an」に対するアイヌ語 UD を、図 1・2 に示す。ただ、各解説には執筆者が記されており、アイヌ語 UD コーパスを作成した場合、著作権処理をどうおこなうべきか悩ましい。

3. アイヌ語 Universal Dependencies の可能性

われわれが作成したアイヌ語 UD コーパスは、係り受け解析エンジン esupar の訓練に用いている。esupar のアイヌ語モジュールを訓練して解析精度を上げることで、さらなるアイヌ語 UD コーパスの作成が楽におこなえる。いわば循環システムだと考えてよい。このようなシステムがうまくいっているのは、UD の言語横断性に加え、単語長と LEMMA を田村 (1996) 『アイヌ語沙流方言辞典』に押し込んだ点が、功を奏したと言える。

ただ、樺太アイヌ語や釧路アイヌ語など、多種多様なアイヌ語を、全て沙流アイヌ語に押し込んでいいものだろうか。この点は、われわれにとっても非常に悩ましい。多種多様なアイヌ語をそのまま言語処理しようとする、それぞれの分量が少なくなってしまうため、解析精度が下がってしまう。多種多様なアイヌ語を保持したまま解析精度を維持するには、FORMに原文を入れた上で、LEMMAを『アイヌ語沙流方言辞典』に接地する、という両天秤な手法しか、うまくいくやり方を見つけきれていない。

アイヌ語 UD は万能ではない。実際、いくつかの文を記述する際に「綻び」が出てきているのも、また事実である。たとえば「shichorpok chikushte shienka chikushte」⁽⁵⁾は、佐藤 (2004) の指摘どおり「shi」を分離する方が適切 (図 12) なのだが、これはアイヌ語 UD としては、かなり特異な事例である。このような特異な事例を踏まえつつ、より適用範囲の広いアイヌ語コーパスを作成していくには、どうすべきか。われわれの今後の研究に期待されたい。

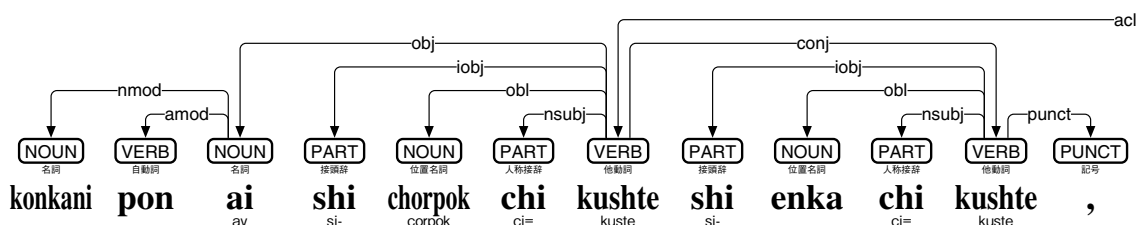


図 12 「konkani ponai shichorpok chikushte shienka chikushte,」のアイヌ語 UD 改良案

謝 辞

本発表に用いた係り受け解析エンジン esupar は、学際大規模情報基盤共同利用・共同研究拠点公募型共同研究『単語間に区切りのない書写言語における係り受け解析エンジンの開発』の成果である。また、アイヌ語 UD エディターとコーパス管理システムの開発、およびそれらを用いたコーパス作成作業は、文部科学省『AI 等の活用を推進する研究データエコシステム構築事業』の支援を受けている。

文 献

- Hajime Senuma, and Akiko Aizawa (2017). “Toward Universal Dependencies for Ainu.” *NoDa-LiDa 2017 Workshop on Universal Dependencies*, pp. 133–139.
- Hajime Senuma, and Akiko Aizawa (2018). “Universal Dependencies for Ainu.” *LREC 2018: Eleventh International Conference on Language Resources and Evaluation*, pp. 2354–2358.
- 知里幸恵 (1923). 『アイヌ神謡集』 郷土研究社, 東京.
- 安岡孝一 (2021a). 「アイヌ語 Universal Dependencies 再考」 東洋学へのコンピュータ利用, 第 34 回研究セミナー, pp. 25–53.

⁽⁵⁾ 知里 (1923) 『アイヌ神謡集』「銀の滴降る降るまはりに」4 ページ 2 行目。

- 安岡孝一 (2021b). 「Universal Dependencies によるアイヌ語テキストコーパス」 情報処理学会研究報告, 2021-CH-127:5, pp. 1–8.
- 佐藤知己 (1996). 「アイヌ語を記述するのに必要な文字セットについて」 JIS 符号化文字集合調査研究委員会第 2 分科会 (WG2) 資料, JCS-2-8-02.
- Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman (2021). “Universal Dependencies.” *Computational Linguistics*, 47:2, pp. 255–308.
- Lucien Tesnière (1959). *Éléments de Syntaxe Structurale*. Paris: C. Klincksieck.
- Igor A. Mel’čuk (1988). *Dependency Syntax: Theory and Practice*. New York: State University of New York Press.
- Joakim Nivre (2015). “Towards a Universal Grammar for Natural Language Processing.” *CICLing 2015: 16th International Conference on Intelligent Text Processing and Computational Linguistics*, pp. 3–16.
- 田村すず子 (1996). 『アイヌ語沙流方言辞典』 草風館, 東京.
- 安岡孝一 (2020). 「Universal Dependencies にもとづく多言語係り受け可視化ツール deplacy」 人文科学とコンピュータシンポジウム「じんもんこん 2020」 論文集, pp. 95–100.
- 佐藤知己 (2004). 「知里幸恵『アイヌ神謡集』の難読箇所と特異な言語事例をめぐって」 北海道立アイヌ民族文化研究センター研究紀要, 10, pp. 1–32.
- 片山龍峯 (2003). 『「アイヌ神謡集」を読みとく』 片山言語文化研究所, 武蔵野.
- 神保小虎・金澤庄三郎 (1898). 『アイヌ語會話字典』 金港堂書籍, 東京.
- Anna Bugaeva (2011). “Internet Applications for Endangered Languages: A Talking Dictionary of Ainu.” 早稲田大学高等研究所紀要, 3, pp. 73–81.
- M. M. Добротворский (1875). *Аинско-русский словарь*. Казань: Университетская типография.
- 寺田吉孝・安田節彦 (2019). 「M. M. ドブロトウヴォールスキーのアイヌ語・ロシア語辞典 (26)」 北海学園大学学園論集, 178, pp. 121–149.
- 阪口諒 (2021). 「『アイヌ語ロシア語辞典』中のアイヌ語樺太方言テキスト」 千葉大学大学院人文社会科学部研究プロジェクト報告書, 第 358 集, pp. 43–55.
- 室鳩巢・空谷茂潤 (1855). 『五倫名義解』 此君園, 江戸.
- 室鳩巢・空谷茂潤 (1858). 『五倫名義解』 宗谷御用所, 宗谷.
- 深澤美香 (2014a). 「加賀家文書のアイヌ語資料と加賀伝蔵」 千葉大学大学院人文社会科学部研究プロジェクト報告書, 第 274 集, pp. 21–48.
- 深澤美香 (2014b). 「加賀家文書における表記の特徴と傾向」 千葉大学大学院人文社会科学部研究プロジェクト報告書, 第 274 集, pp. 49–72.
- 安岡孝一・安岡素子 (2023). 「アイヌ語訳『五倫名義解』 Universal Dependencies への挑戦」 東洋学へのコンピュータ利用, 第 36 回研究セミナー, pp. 3–37.
- 小林美紀 (2023). 「アイヌ語を第一言語に」 国立アイヌ民族博物館 (編) 『ウアイヌコロ コタン アカラ ウポポイのことばと歴史』 国書刊行会, 東京 pp. 97–111.
- 深澤美香 (2023). 「国立アイヌ民族博物館のアイヌ語による展示解説文と「私たち」」 国立ア

アイヌ民族博物館 (編) 『ウアイヌコロ コタン アカラ ウポポイのことばと歴史』 国書刊行会, 東京 pp. 112–153.

国立アイヌ民族博物館 (2020a). 『国立アイヌ民族博物館ガイドブック』 国立アイヌ民族博物館, 白老.

国立アイヌ民族博物館 (2020b). 『アヌココロ アイヌ イコロマケンル an=ukokor aynu ikor oma kenru』, 国立アイヌ民族博物館パンフレット (日本語), 白老.

国立アイヌ民族博物館 (2021). 『ゴールデンカムイ トウラノ アプカシアン』, 国立アイヌ民族博物館第 2 回特別展示, 白老.

関連 URL

Universal Dependencies for Ainu	https://github.com/KoichiYasuoka/UD-Ainu
係り受け解析エンジン esupar	https://github.com/KoichiYasuoka/esupar
アイヌ語 UD エディター	https://koichiyasuoka.github.io/UD-Ainu/editor/
トピック別アイヌ語会話辞典	https://ainu.ninjal.ac.jp/topic/
国立国会図書館デジタルコレクション	
『アイヌ神謡集』	https://dl.ndl.go.jp/pid/1909336
『アイヌ語會話字典』	https://dl.ndl.go.jp/pid/993685