



TITLE:

Universal Dependenciesによるアイヌ語テキストコーパス

AUTHOR(S):

安岡, 孝一

CITATION:

安岡, 孝一. 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: 人文科学とコンピュータ研究会報告 2021, 2021-CH-127(5): 1-8

ISSUE DATE:

2021-08

URL:

<http://hdl.handle.net/2433/266186>

RIGHT:

ここに掲載した著作物の利用に関する注意 本著作物の著作権は情報処理学会に帰属します。本著作物は著作権者である情報処理学会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします。; © 2021 Information Processing Society of Japan

Universal Dependencies によるアイヌ語テキストコーパス

安岡孝一 (京都大学人文科学研究所附属東アジア人文情報学研究センター)

書写言語としてのアイヌ語は、カタカナ・キリル文字・ローマ字(ラテンアルファベット)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文法構造を持ち、これらの言語向けの文法記述手法は、アイヌ語に太刀打ちできない。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。本発表では、カタカナ・キリル文字・ローマ字で書かれたアイヌ語を、Universal Dependencies で記述する際の困難さについて、考察する。

1 はじめに

アイヌ語 Universal Dependencies は、東京大学の瀬沼甫が開発 [1, 2] し、2017 年 3 月 31 日に GitHub で公開(以下「瀬沼版」と呼ぶ)された。最初の公開は『アイヌ神謡集』[3]の「ホテナオ」のみであり、その後データが追加されることもなく、4 年後には GitHub リポジトリが閉鎖されてしまった。これに対し筆者は、手元に残された瀬沼版のコピーに、新たなデータを追加すべく検討を加えた [4]。しかし、残念ながら瀬沼版は、アノテーションの基準が不明瞭な(他者と共有されていない)点が多く、直接これに追加をおこなうのは、筆者には無理だった。

アイヌ語は、いわゆる「文字を持たない言語」であり、当然のことながら正書法が存在しない。書写言語としてのアイヌ語は、カタカナ・キリル文字・ローマ字(ラテンアルファベット)等を借りて記述されており、その「借り方」も一定していない。たとえば「ホテナオ」[3]の「itskash awa ponrupneainu ene itaki」に対し、[5]では「itak as a ūa pon rupne aĵnu ene itaki」と、[6]では「itak as a wa pon rupne aynu ene itak i」と、[7]では「itak=as awa pon rupne aynu ene itak i」「イタカシ アワ ポンルプネアイヌ エネ イタキ」と、それぞれ記されている。あるいは、キリル文字 [8]なら「итакас ава пон рубне айну эне итаки」と記することも可能だろう。

そのような書写言語としてのアイヌ語を、Universal Dependencies [9](以下「UD」と呼ぶ)は、本当に記述できるのだろうか。瀬沼版のような「特定の表記法に寄せる」やり方ではなく、複数の表記法を同じように扱えるアイヌ語 UD は、本当に実現可能なのだろうか。これらの点について、日本語 UD との比較を交えつつ考察する。

2 Universal Dependencies の概要

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。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 の構造的統語論 [10]に源を発し、Мельчук の有向グラフ記述 [11]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、Мельчук 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。これは、Мельчук の有向グラフ記述が、単語間のリンクという形態を取っていたからであり、そういう割り切りの結果として、言語横断的な文法構造記述が可能としているのである。

UD 依存構造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、以下に示す 10 個のタブ区切りフィールドで構成される。

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^{a)}。
5. XPOS: 言語固有の品詞タグ。

^{a)} ADJ・ADP・ADV・AUX・CCONJ・DET・INTJ・NOUN・NUM・PART・PRON・PROPN・PUNCT・SCONJ・SYM・VERB・X の 17 種類。

表 1: 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 決定詞 clf 類別詞 case 格表示
Coordination	MWE	Loose	Special	Other
conj 接続 cc 接続詞	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

- FEATS: UD で規定された言語普遍的形態素属性のリスト。言語固有の拡張も可。
- 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。
- DEPREL: UD で規定された言語普遍的係り受けタグ (表 1)。HEAD が 0 の場合は root とする。言語固有の拡張も可。
- DEPS: 複数の係り受け元を持つ場合、全ての HEAD:DEPREL ペア。
- MISC: その他のアノテーション。

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有り得るが、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、Мельчук 依存文法の後裔であり、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点 [12] が、Мельчук とは異なっている。ちなみに、コンピュータ文においては、補語をリンク元として、主語へとリンクする。

本稿執筆時点で最新の UD は、2021 年 5 月 16 日発表の UD 2.8.1 である。UD 2.8.1 は、114 の言語にまたがるツリーバンクだが、筆者の見る限り、抱合語はチュクチ語 UD とユピック語 UD (セントローレンス島) の 2 つだけのようである。

3 アイヌ語 UD の再設計

CoNLL-U のタブ区切りフィールドのうち、FEATS と DEPS を除く 8 つのフィールドを用いる形で、アイヌ語 UD の再設計をおこなった。以下に、概要を述べる。

3.1 単語長と語境界

書写言語としてのアイヌ語は、通常、分かち書きされており、空白で語境界を示している。ただし、語境界には合意が無い (図 1)。作業上の語境界を決めるしか無さそうである。

筆者の知る限り、オンライン検索可能なアイヌ語辞書で最大のものは、『アイヌ語沙流方言辞典』 [13] である。国立アイヌ民族博物館アイヌ語アーカイブ^{b)}

^{b)}https://ainugo.nam.go.jp

itskash awa ponrupneainu ene itaki
itak as a ũa pon rupne ajnu ene itaki
itak as a wa pon rupne aynu ene itak i
itak=as awa pon rupne aynu ene itak i
イタカシ アワ ポンルプネアイヌ エネ イタキ
итакас ава пон рубне айну эне итаки

図 1: 「itskash awa ponrupneainu ene itaki」の語境界

に収録されており、ローマ字とカタカナで検索可能である。この『アイヌ語沙流方言辞典』を、作業上の単語認定に用いることにする。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る。

FORM は対象の原文そのままの表記を、LEMMA は『アイヌ語沙流方言辞典』のローマ字表記を用いる。作業上の単語長と原文の語境界との間に生じる齟齬は、以下の2つの方法で吸収する。

- MISC に SpaceAfter=No を入れる
- DEPREL に goeswith を入れる

これらの方法で記述できない場合は、以下の方法を用いてもよい。

- ID に単語の範囲を入れる(縮約語)
- FORM に空白を含める

3.2 品詞付与と係り受け

XPOS は『アイヌ語沙流方言辞典』の品詞分類に従う。ただし、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。

UPOS は、XPOS から半自動で付与する(表2)。接続詞のうち CCONJ とすべき語(「awa」と「korka」)があるので、その点は注意されたい。

DEPREL は、表1のうち、depを除く36種類を用いる。人称接辞へのリンクは、nsubj・obj・detを基本とするが、二重主語などの場合には expl を使う[4]。デアル動詞へのリンクは cop とし、コピュラ文とみなす。自動詞や副詞が連体修飾に用いられている場合は、被修飾語から amod でリンクする。他動詞が連体修飾に用いられている場合は、被修飾語から acl でリンクする。

図1の6つの文に対し、単語長を定めた上で、品詞付与と係り受けを手作業でおこなった。結果を図2に示す。語境界の齟齬は、SpaceAfter=No・goeswith・

表 2: アイヌ語 UD 再設計のための UPOS・XPOS

UPOS	XPOS
NOUN	名詞・位置名詞・形式名詞
PROPN	固有名詞
PRON	代名詞
VERB	完全動詞・自動詞・他動詞
AUX	助動詞・デアル動詞
ADV	副詞
NUM	数詞
DET	連体詞
CCONJ	「awa」「korka」
SCONJ	接続詞・接続助詞・後置副詞
ADP	格助詞・副助詞
PART	終助詞・人称接辞・接尾辞・接頭辞
INTJ	間投詞
PUNCT	記号

縮約語の3つの方法で吸収した。単語数が微妙に異なっているものの、ローマ字・カタカナ・キリル文字のいずれの例に対しても、CoNLL-Uを同様に記述できると言える。比較のために、図2の各データを deplacy [14] で可視化した(図3)ので、合わせて見比べてほしい。

4 アイヌ語 UD と他言語 UD の比較

『アイヌ神謡集』[3]の各文には、知里幸恵による日本語訳が付されている。これを平行コーパスとみなして、アイヌ語 UD と日本語 UD の比較(図4)をおこなってみよう。なお、アイヌ語 UD は筆者の手作業[4]で、日本語 UD は SuPar-UniDic [15](旧仮名 UniDic + 青空文庫 BERT) で作成した。また、図4には参考として、Julie Kaizawa の英訳 [7] と Márta Müller の洪訳 [16] も、Trankit [17] で UD 化して示した。

結論から書くと、アイヌ語 UD の依存構造は、日本語 UD や英語 UD やハンガリー語 UD とは、かなり異なっている。図4(a)は、主語を表す人称接辞「ash」が、動詞の直後に位置している。図4(b)のコピュラ文「Pon Horkeusani ene」は、「Pon Horkeusani」が補語、「e」が主語、「ne」が繫辞だが、主語を「eani anak」で強調しているために、「e」の方を expl にせざるを得ない。図4(c)は、いわゆる連他動詞(「chorpok kushte」と「enka kushte」)がアイヌ語に特徴的であり、そもそも他の言語に訳しきれない。

情報処理学会研究報告
IPSJ SIG Technical Report

text = itskash awa ponrupneainu ene itaki

1	itsk	itak	VERB	自動詞	-	0	root	-	SpaceAfter=No
2	ash	=as	PART	人称接辞	-	1	nsubj	-	-
3	awa	awa	CCONJ	接続詞	-	1	cc	-	-
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	-	1	conj	-	-

text = itak as a ũa pon rupne ajnu ene itaki

1	itak	itak	VERB	自動詞	-	0	root	-	-
2	as	=as	PART	人称接辞	-	1	nsubj	-	-
3	a	awa	CCONJ	接続詞	-	1	cc	-	-
4	ũa	-	X	-	-	3	goeswith	-	-
5	pon	pon	VERB	自動詞	-	7	amod	-	-
6	rupne	rupne	VERB	自動詞	-	7	amod	-	-
7	ajnu	aynu	NOUN	名詞	-	9	nsubj	-	-
8	ene	ene	ADV	副詞	-	9	advmod	-	-
9	itak	itak	VERB	自動詞	-	10	acl	-	SpaceAfter=No
10	i	-i	PART	接尾辞	-	1	conj	-	-

text = itak as a wa pon rupne aynu ene itak i

1	itak	itak	VERB	自動詞	-	0	root	-	-
2	as	=as	PART	人称接辞	-	1	nsubj	-	-
3	a	awa	CCONJ	接続詞	-	1	cc	-	-
4	wa	-	X	-	-	3	goeswith	-	-
5	pon	pon	VERB	自動詞	-	7	amod	-	-
6	rupne	rupne	VERB	自動詞	-	7	amod	-	-
7	aynu	aynu	NOUN	名詞	-	9	nsubj	-	-
8	ene	ene	ADV	副詞	-	9	advmod	-	-
9	itak	itak	VERB	自動詞	-	10	acl	-	-
10	i	-i	PART	接尾辞	-	1	conj	-	-

text = itak=as awa pon rupne aynu ene itak i

1	itak	itak	VERB	自動詞	-	0	root	-	SpaceAfter=No
2	=as	=as	PART	人称接辞	-	1	nsubj	-	-
3	awa	awa	CCONJ	接続詞	-	1	cc	-	-
4	pon	pon	VERB	自動詞	-	6	amod	-	-
5	rupne	rupne	VERB	自動詞	-	6	amod	-	-
6	aynu	aynu	NOUN	名詞	-	8	nsubj	-	-
7	ene	ene	ADV	副詞	-	8	advmod	-	-
8	itak	itak	VERB	自動詞	-	9	acl	-	-
9	i	-i	PART	接尾辞	-	1	conj	-	-

```
# text = イタカシ アワ ポンルプネアイヌ エネ イタキ
1-2   イタカシ  -      -      -      -      -      -      -      -
1     イタク   itak  VERB  自動詞  -      0      root   -      -
2     アシ     =as  PART  人称接辞 -      1      nsubj  -      -
3     アワ     awa  CCONJ 接続詞  -      1      cc     -      -
4     ポン     pon  VERB  自動詞  -      6      amod   -      SpaceAfter=No
5     ルプネ   rupne VERB  自動詞  -      6      amod   -      SpaceAfter=No
6     アイヌ   aynu  NOUN  名詞    -      8      nsubj  -      -
7     エネ     ene  ADV   副詞    -      8      advmod -      -
8-9   イタキ   -      -      -      -      -      -      -      -
8     イタク   itak  VERB  自動詞  -      9      acl    -      -
9     イ       -i   PART  接尾辞  -      1      conj   -      -
```

```
# text = итакас ава пон рубне айну эне итаки
1     итак    itak  VERB  自動詞  -      0      root   -      SpaceAfter=No
2     ас      =as  PART  人称接辞 -      1      nsubj  -      -
3     ава     awa  CCONJ 接続詞  -      1      cc     -      -
4     пон     pon  VERB  自動詞  -      6      amod   -      -
5     рубне   rupne VERB  自動詞  -      6      amod   -      -
6     айну    aynu  NOUN  名詞    -      8      nsubj  -      -
7     эне     ene  ADV   副詞    -      8      advmod -      -
8     итак    itak  VERB  自動詞  -      9      acl    -      SpaceAfter=No
9     и       -i   PART  接尾辞  -      1      conj   -      -
```

図 2: 再設計版アイヌ語 UD の CoNLL-U データ

5 おわりに

前章での考察にもとづく限り、日本語 UD における「知識」は、残念ながら、アイヌ語 UD に援用できる可能性が低い。英語 UD や他の欧米諸語も同様である。これら以外の UD、特に抱合語との比較もおこなうべきではあるが、現時点のチュクチ語 UD は 6124 語、ユピック語 UD (セントローレンス島) は 2568 語と、係り受けコーパスとしては少量である。これらの抱合語 UD における「知識」は、人手でアイヌ語 UD に適用できる可能性はあるが、アイヌ語の機械学習に援用できる可能性は極めて低い。

本発表では、『アイヌ神謡集』[3]と『アイヌ語沙流方言辞典』[13]を下敷きに、ローマ字・カタカナ・キリル文字で書かれたアイヌ語 UD の作成に関して考察した。このうち「ホテナオ」と「クツニサクトクトン」については、アイヌ語 UD と日本語 UD の平行コーパスを筆者自身が作成中であり、適宜 GitHub で公開⁹⁾している。ただ、今後、このアイヌ語 UD を筆者自身が拡張し続けるべきかは、非常に悩ましい。

筆者のアイヌ語の能力は世間並であり、筆者一人では判断のつかない文法的事項も多々ある。ちゃんと開発をおこなうには、共同研究等のチームを組んで、継続的に開発をおこなえるよう環境を整える必要があるだろう。

参考文献

- [1] Hajime Senuma, Akiko Aizawa: Toward Universal Dependencies for Ainu, NoDaLiDa 2017 Workshop on Universal Dependencies (May 2017), pp.133-139.
- [2] Hajime Senuma, Akiko Aizawa: Universal Dependencies for Ainu, LREC 2018: Eleventh International Conference on Language Resources and Evaluation (May 2018), pp.2354-2358.
- [3] 知里幸恵: アイヌ神謡集, 東京: 郷土研究社 (1923 年 8 月).
- [4] 安岡孝一: 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 (2021 年 7 月), pp.25-53.
- [5] Ainaj jukaroj, Sapporo: Hokkajda Esperanto-Ligo (1979).
- [6] 切替英雄: 『アイヌ神謡集』辞典, 北大言語学研究報告, 第 2 号 (1989 年 6 月).

⁹⁾<https://github.com/KoichiYasuoka/ud-ainu>

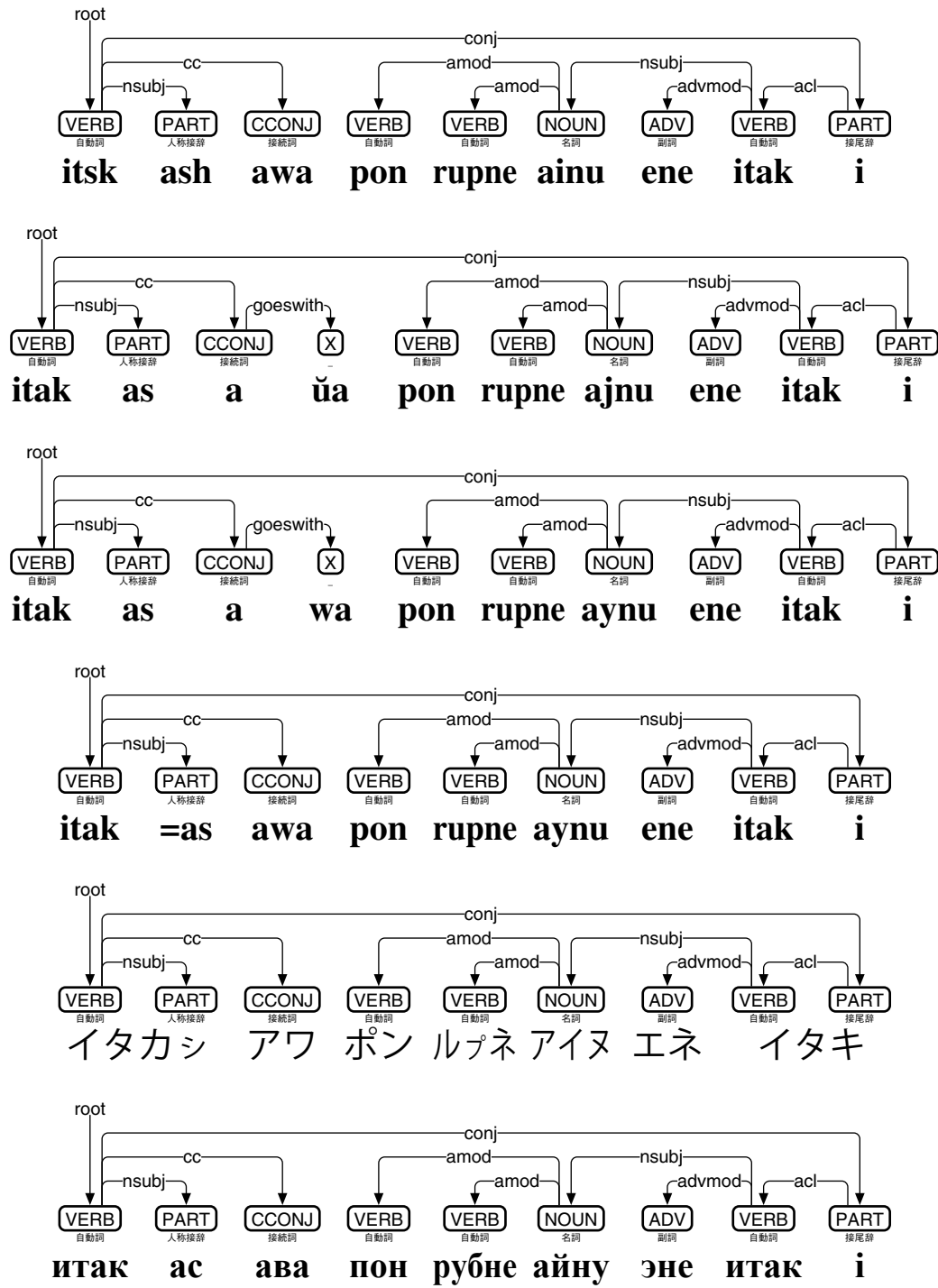
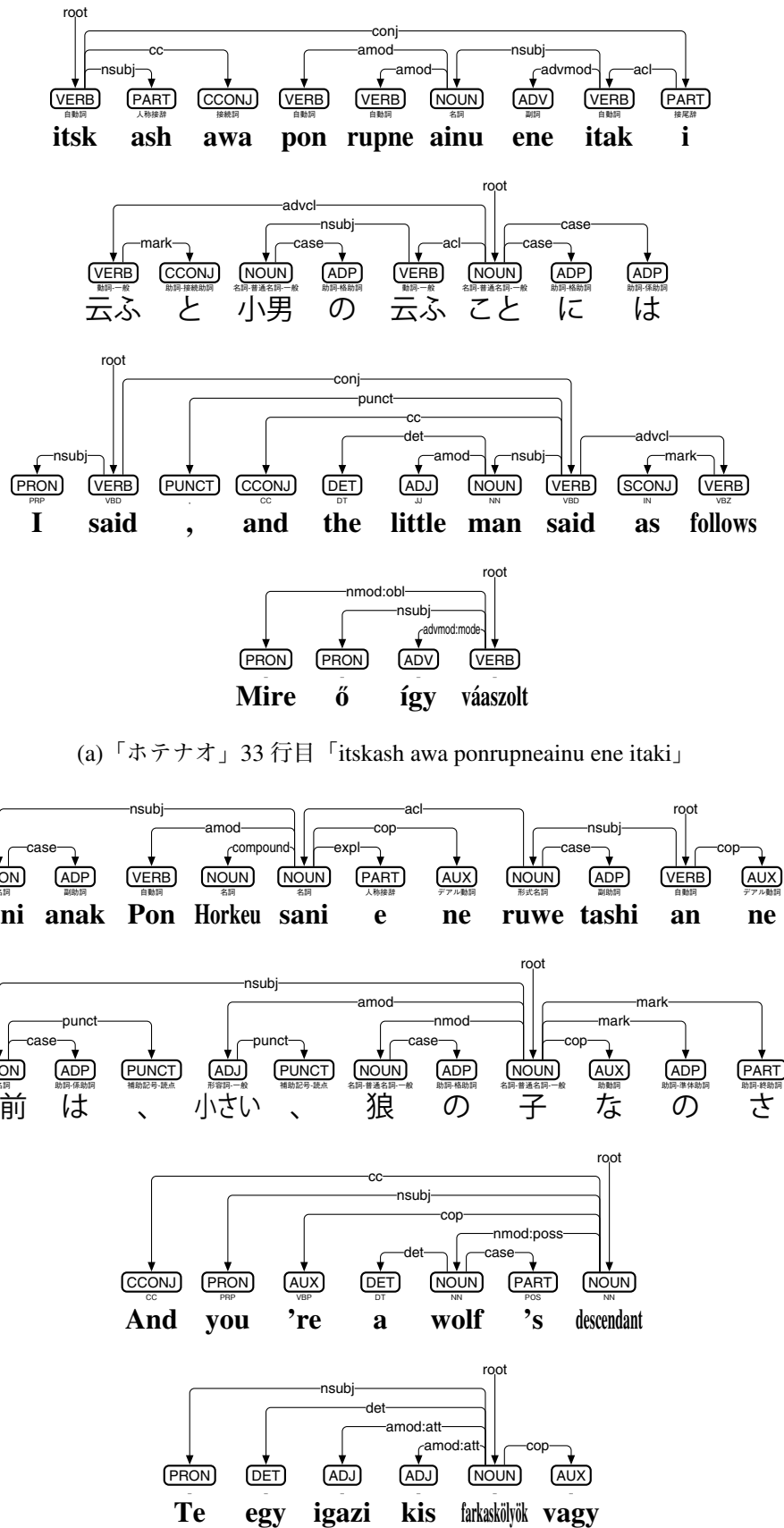
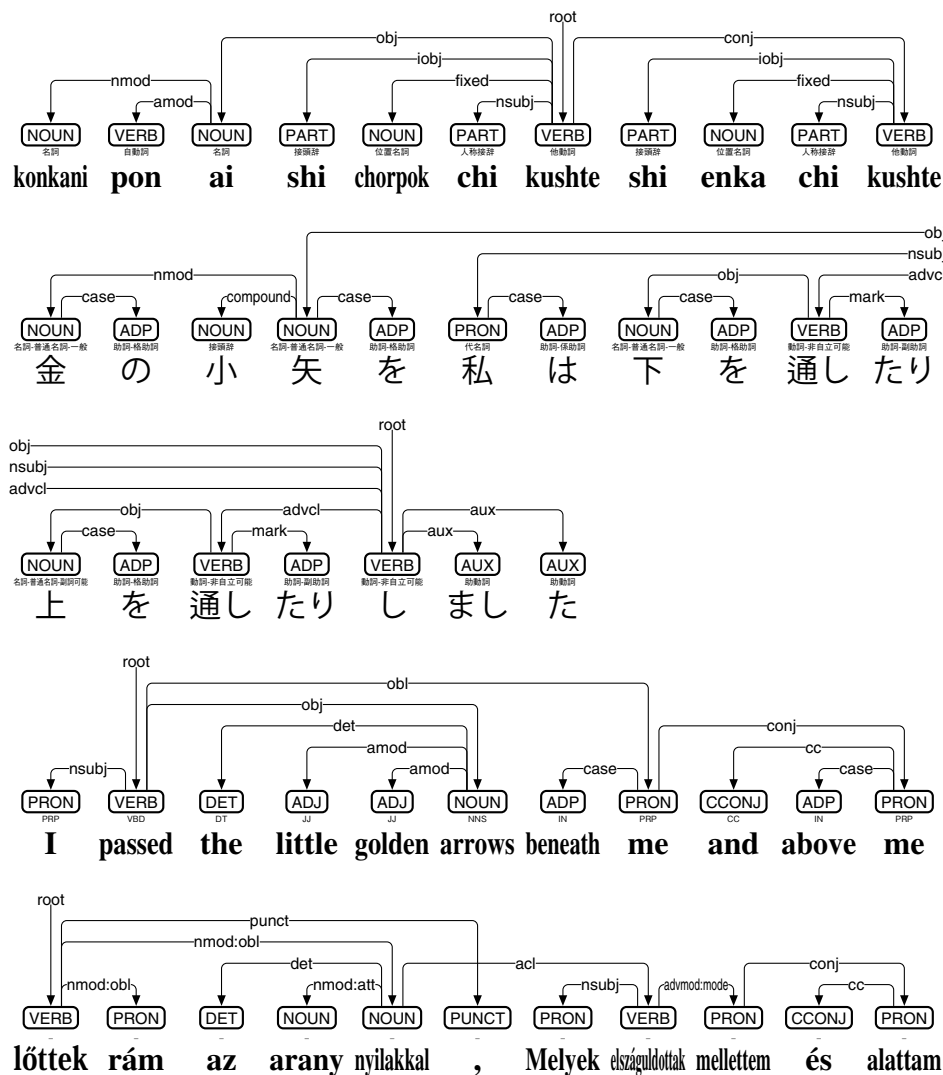


図 3: deplacy によるアイヌ語 UD の可視化





(c) 「銀の滴降る降るまはりに」 20～21 行目 「konkani ponai shichorpok chikushte shienka chikushte」

図 4: アイヌ語 UD と他言語 UD (日本語・英語・ハンガリー語) の比較

- [7] 片山龍峯: 「アイヌ神謡集」を読みとく, 武蔵野: 片山言語文化研究所 (2003 年 5 月).
- [8] 荻原眞子, 丹菊逸治: 千徳太郎治のピウスツキ宛書簡 — 「ニシパ」へのキリル文字の手紙 —, 千葉大学ユーラシア言語文化論集, 第 4 号 (2001 年 3 月), pp.187-226.
- [9] Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, Daniel Zeman: Universal Dependencies, Computational Linguistics, Vol.47, No.2 (June 2021), pp.255-308.
- [10] Lucien Tesnière: Éléments de Syntaxe Structurale, Paris: C. Klincksieck (1959).
- [11] Igor A. Mel'čuk: Dependency Syntax: Theory and Practice, New York: State University of New York Press (1988).
- [12] Joakim Nivre: Towards a Universal Grammar for Natural Language Processing, CICLing 2015: 16th International Conference on Intelligent Text Processing and Computational Linguistics (April 2015), pp.3-16.
- [13] 田村すず子: アイヌ語沙流方言辞典, 東京: 草風館 (1996 年 9 月).
- [14] 安岡孝一: 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 (2020 年 12 月), pp.95-100.
- [15] 安岡孝一: 世界の Universal Dependencies と係り受け解析ツール群, 第 3 回 Universal Dependencies 公開研究会 (2021 年 6 月).
- [16] Márta Müller: A bagolyisten dala, epubli (2015 szeptember).
- [17] Minh Van Nguyen, Viet Dac Lai, Amir Pouran Ben Veyseh and Thien Huu Nguyen: Trankit: A Light-Weight Transformer-based Toolkit for Multilingual Natural Language Processing, EACL 2021: 16th Conference of the European Chapter of the Association for Computational Linguistics (April 2021), System Demonstrations, pp.80-90.