

Morphological Disambiguation for Tocharian

Interdisciplinary Project Report

Overview

- Introduction
- What is Tocharian?
- The Problem
- Approach
- Evaluation
- Results
- Conclusion

Introduction

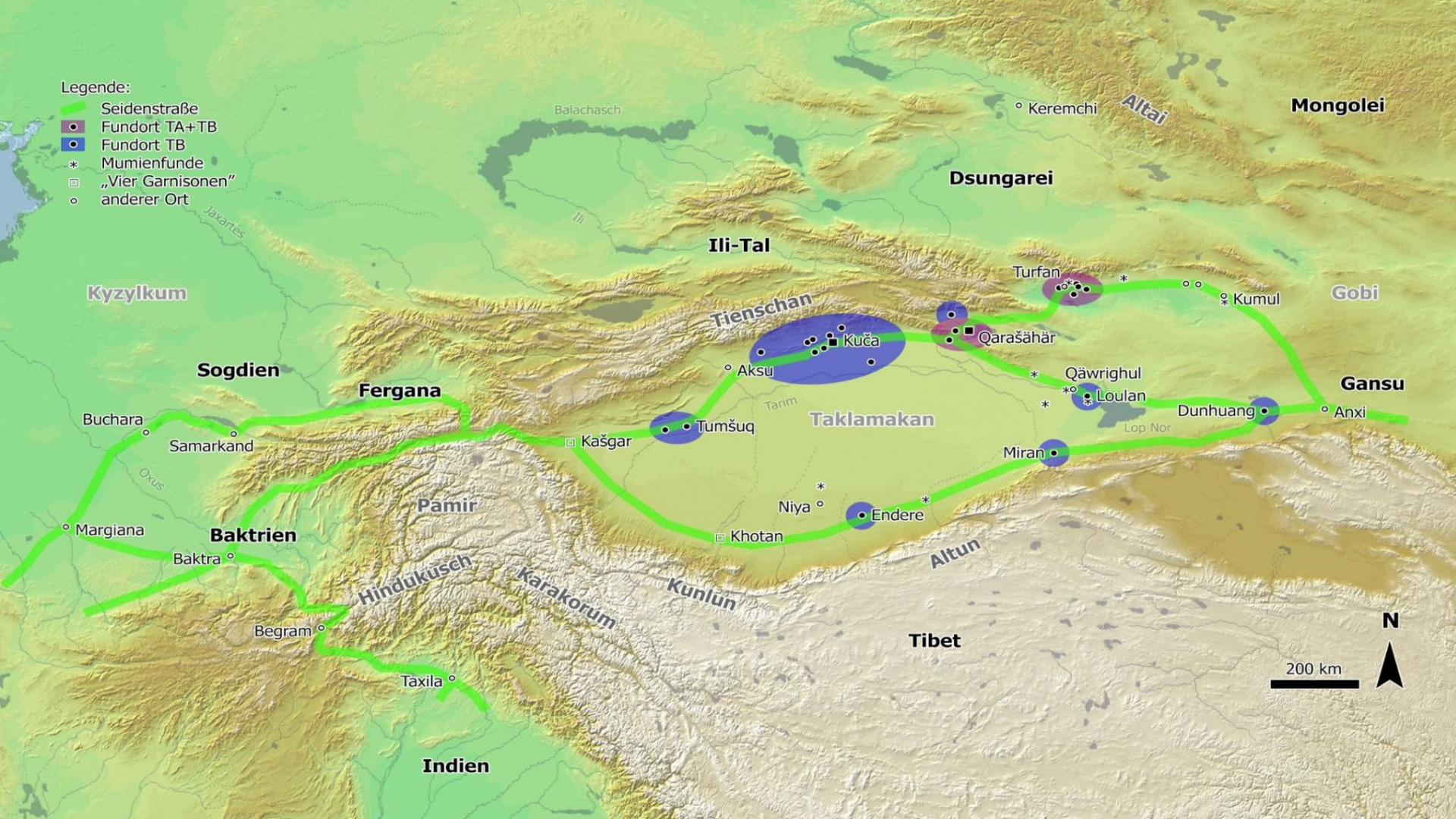
- Master's Student in Data Science
- Interdisciplinary Project
 - “Solve a practical problem in interdisciplinary project work”
 - Problem should not be primarily based in IT or Mathematics.
- Got in contact with Project “Tarim Brahmi”
 - “allow the comprehensive paleographic investigation ... ” [of Tocharian]
 - by linking relevant data about manuscripts in one place

The Tocharian Language(s)

- Was discovered around 1900
 - along the routes of the Silk Road
 - in the Tarim Basin
- Actually two very similar languages
 - Tocharian A (TA) - found to the east
 - Tocharian B (TB) - found to the west
- 4th to 10th century CE
- Manuscripts are often translations of Sanskrit
 - Buddhist, scientific, administrative documents.

Legende:

- Seidenstraße
- Fundort TA+TB
- Fundort TB
- Mumienfunde
- „Vier Garnisonen“
- anderer Ort



200 km





Finland

Russia

Estonia
Latvia
Lithuania

Belarus

Ukraine
Moldova
Romania

Kazakhstan

Mongolia

Greece
Bulgaria
Turkey

Georgia

Azerbaijan

Uzbekistan

Tajikistan

IT'S HERE

North Korea

South Korea

Japan

Syria
Lebanon
Israel
Jordan

Iraq

Iran

Afghanistan

Pakistan

China

India

Nepal

Bhutan

Bangladesh

Taiwan

Egypt

Saudi Arabia

United Arab Emirates

Oman

Sudan

Eritrea

Yemen

South Sudan

Ethiopia

DRC

Uganda

Kenya

Somalia

Sri Lanka

Myanmar (Burma)

Laos

Thailand

Cambodia

Vietnam

South China Sea

Philippines

Palawan

Negros

Mindanao

Basilan Island

Celebes Sea

Bay of Bengal

Andaman Sea

Gulf of Thailand

Malaysia

Singapore

Java Sea

Indonesia

Banda Sea

Central African Republic

Rwanda

Burundi

Tanzania

Papua New Guinea

Indo-European Languages

| TA | TB | english | Latin | Sanskrit | Persian |
|--------|--------|---------|--------|----------|---------|
| känt | kante | hundred | centum | śatām | sad |
| pracar | procer | brother | frāter | bhrāṭṛ | barâdar |

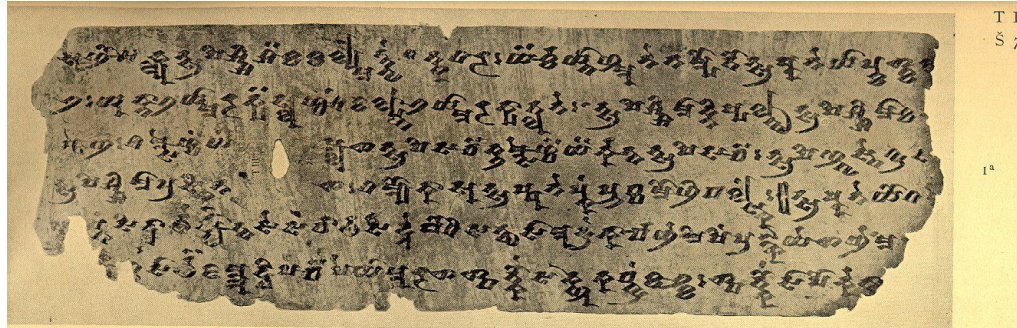


Handwritten text in an ancient script, likely Tamil, on a heavily damaged and torn piece of aged paper. The text is arranged in approximately 10 horizontal lines, though many characters are obscured by large holes and irregular tears in the paper. The ink is dark, and the paper is yellowed with age. The script consists of small, rounded characters with some distinct vertical strokes. The overall appearance is that of a fragment from an ancient manuscript or record.

Project “Tarim Brahmi”

- Cooperation between University of Vienna & Austrian Center for Digital Humanities and Cultural Heritage
- “... *link the text witnesses to their digital facsimiles on the character level and to publish this material together with a **TEI-encoded dictionary** in an online database*” [7]
- “... *all quantifiable features of all characters, ligatures and words will be extracted and compared using software tools*” [7]
- Currently working on the dictionary

Current Progress



T II
Š 7:

1^a

sātkatār

| | |
|------------------|---|
| Lemma: | sātkā- |
| Word class: | Finite verb form |
| Language: | TA |
| Lexeme variants: | sātkatār |
| Lexeme family: | <ul style="list-style-type: none"> sātkā- <i>Root</i> <ul style="list-style-type: none"> Kausativum 1 "to spread" Grundverb "to spread out" <ul style="list-style-type: none"> Present 3 Subjunctive 5 <ul style="list-style-type: none"> sātkālune verbal abstract Preterite 1 <ul style="list-style-type: none"> sātko preterite participle |
| Person: | 3 |
| Number: | Singular |
| Tense/Mood: | Present |
| Stem: | Present |
| Stem class: | 3 |
| Voice: | Middle |

Inflectional paradigm

Present

| | Singular | | Plural | | Dual | |
|--------|----------|----------|--------|--------------------------|--------|--------|
| | Active | Middle | Active | Middle | Active | Middle |
| First | | | | | | |
| Second | | | | | | |
| Third | | sātkatār | | sātkantār sātkantr-ām | | |

Preterite

| | Singular | | Plural | | Dual | |
|--------|--------------|--------|--------------|---------------|--------|--------|
| | Active | Middle | Active | Middle | Active | Middle |
| First | | | | | | |
| Second | | | | | | |
| Third | stāk sātkā-ṃ | | satkar ām | satkar- ām | | |

Occurrences: 4

sātkatār

- 1 A 1 a1 ^{<at>} ṅom klyu tsraṣiṣi śāk kālymentwaṃ **sātkatār** : yārḱ yṅāṃmune nam poto tsraṣṣuneyā
- 2 A 2 b2 kāsu ṅom klyu amoktsāp kālyme kālyme **sātkatār** : yārḱā yāmāl māskatār potal (k)r(o)-_{b3}-pal
- 3 A 37 b5 _{<b5>} -trā : kāsu ṅom klyu **sā(tkatār)** _{<b6>} -ntrā kutsmātāṅ nāmeṅc **///**
- 4 A 79 b5 : kāswoneyo yā(mu) nu pāl ākntsāṣi **sātkatār** tri āpāytwāṃ : 1 (yā)sluntāṃ pe

Transcription

- ... (kā)-
- a1 -su ṅom klyu tsraṣiṣi śāk kālymentwaṃ sātkatār : yārḱ yṅāṃmune nam poto tsraṣṣuneyā p_ukāṣ kāl(pnā)-
- a2 -l : yuknāl ymārāk yāsluṅcās kālpnāl ymārāk yātluṅe : 1 tsraṣiṣi māḱ niṣpalntu tsraṣiṣi māḱ (śkaṃ) (ṣṅā)-
- a3 -ṣṣeṅ : nāmseṅc yāsluṣ tsraṣisac kumseṅc yārḱant tsraṣisac : tsraṣiṅ waste wrasa(śsi)
- a4 tsraṣiṣi mā praski naṣ : tām̄yo kāsū tsraṣṣune p_ukaṃ pruccamo ṅi pālskaṃ : || tsraṣṣuneyo tām(n)e (ne)-
- a5 -ṣ (pra)staṃ siddhārthes lānt se sarvārthasiddhe bodhisattu sāmudraṃ kārp ṅemiṣiṃ praṅkā yeṣ ṅemi - -
- a6 - l-i - sārth jambudvipac pe yāmurāṣ ṣpāt koṃsā kṅukac wraṃ kālk : ṣpāt koṃsā pokenā - - - -
- b1 - (kā)lk ṣpāt koṃsā lyomaṃ kālk ṣpāt koṃsā wāltṣ pāltwāyo oplāsyo wraṃ oplāṣ oplā kārnm(ām) (kālkorā)-
- b2 -ṣ pāṅ kursār wā ār(p/ṣ)lāsyo rarkusāṃ tkaṅā kālk : tmāṣ rākṣtsāṣi dvipaṃ yeṣ tmāṣ yakṣāṣi - -
- b3 baladvipaṃ yeṣ tmāṣ śtwar-wāknā ārṣlās(l/y)o rarkuṅcās iṣāṅās kcāk śtwar-wāknā spe(ṣinā)-
- b4 -s«ṭā» klumtsāsyo sopis sāgares lānt lāṅci waṣt pāsāntās sāwes empeles (n)ā(kā)-
- b5 -s āsuk kātkorāṣ sāgareṃ lāntāṣ cindāmaṅi wṃar toriṃ kālpāt poṅcām jambudvipis e(kro)-
- b6 -rṅe wawik ślak śkaṃ || ṣāṃnernaṃ || māskī kātkālām klāṅkeṅc tsraṣiṅ sāmuddrā : traidhātuk sams(ār) (tsra)-
- ... -(ṣṣuneyo)

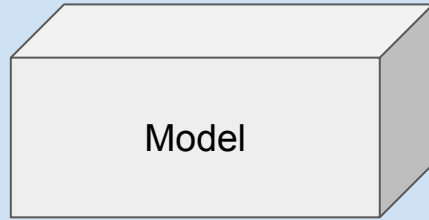
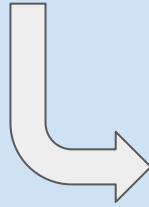
The Problem

The Problem

- Create a dictionary of all tokens + morphological & grammatical information
 - Lemma
 - Grammatical Tags in multiple categories (case, pos, gender, ...)
- Manually annotating each token/type is very time consuming
- Solution: Train predictors for lemma and grammatical information
 - Iterate over the tokens of all documents
- This project was **not**:
 - beating a benchmark
 - comparing different architecture
 - exploring the best parameter setting

Goal

Token



Lemma + gramm.Tags

`"Worten" => "Wort" + [plural,dativ,...]`

Morphological Disambiguation

- Finding the correct grammatical parse for the morphemes of an inflected word
 - commonly includes lemmatization
- Definition & Approach depends on language
 - english: inflective (barely), but very ambiguous
 - turkish: agglutinative, vast vocabulary, modular

dogs → dog(N+pl)

are → be(aux+pres+3+pl)

are → be(pres+2+sg)

saglamlastirmek → saglam/(adj) las(verb+become) tir(verb+caus) mak(noun+nom)

= “the thing that causes something to become strong”

Tocharian Declension

- 10 Cases
 - 4 Primary Cases: Nominative, Genitive, Accusative, Vocative (only TB)
 - 6 Secondary Cases: Perlative, Comitative, Allative, Ablative, Locative, Instrumental (only TA)
- Secondary cases attach to [stem] + [acc.]
- Seems inflective

| "King" | |
|--------|--------|
| nom.sg | walo |
| acc.sg | lānt |
| all.sg | lāntäs |

Derivational Morphology

- Productive morphemes can stack together
 - Seems agglutinative

| | | |
|-----------------|-------------|-----------|
| snai | preposition | “without” |
| snaitstse | adjective | “poor” |
| snaitstsãññe | noun | “poverty” |
| snaitstsãññeşşe | adjective | “pitiful” |

Ultimately “snaitstsãññeşşe” is an adjective.

Mode of Disambiguation

$\text{disambiguate}(x) = \{\text{lemma}(x), p \in \text{pos}, c \in \text{case}, \dots, v \in \text{voice}, \dots\}$

pos {noun, adj, verb, uninfl, unkn}

case {nom, gen, acc, voc*, per, com, all, abl, loc, ins*}

gender {m, f, n}

number {sg, pl, du}

person {1, 2, 3}

tense {prs, sbj, impf, opt, pret, imp}

voice {act, mid}

*language specific

Data Exploration

```
<pos-n type="entry" id="B_B_prahati" lemma="prahati">
<meaning>Indian nightshade, Solanum indicum [a medical ingredient]</meaning>
<pos val="pos-n"/>
<stem_gender/>
<nominative val="infl-nom-i"/>
<accusative val="infl-acc-i"/>
<plural val="infl-pl-nta"/>
<etym>Skt. &lt;i>prahati&lt;i>; Khot. &lt;i>prahatti&lt;i>;
<word id="B_B_prahati">
<case val="case-nom"/>
<case val="case-acc"/>
<number val="num-sg"/>
<gender/>
<phon id="B_prahati">
<form>prahati</form>
<accent-rule val="y-p-Aə-a"/>
<accent-rule val="y-p-Ua-a"/>
<accent-rule val="y-p-Ui-i"/>
</phon>
<phon id="B_brāhati">
<form>brāhati</form>
<accent-rule val="y-p-Ua-ā"/>
<accent-rule val="y-p-Aə-a"/>
<accent-rule val="y-p-Ui-i"/>
</phon>
</word>
<word id="B_B_prahatinta">
<case val="case-nom"/>
<case val="case-acc"/>
<number val="num-pl"/>
<gender/>
<phon id="B_prahatinta">
<form>prahatinta</form>
<accent-rule val="y-p-Aə-a"/>
<accent-rule val="y-p-Ua-a"/>
<accent-rule val="y-p-Ui-i"/>
</phon>
</word>
</pos-n>
<pos-n type="entry" id="B_B_yirmakka" lemma="yirmakka">
<meaning>~ treasurer (lit: measurer)</meaning>
<pos val="pos-n"/>
<stem_gender/>
<note>
< bibl xmlns="http://www.tei-c.org/ns/1.0">
< title corresp="ref:adams1999"/>
</ bibl>; s.v. < hi xmlns="http://www.tei-c.org/ns/1.0">yirmakka* </ hi>. </ note>
<word id="B_B_yirmakka">
<case val="case-nom"/>
<number val="num-sg"/>
<gender/>
<note>
< bibl xmlns="http://www.tei-c.org/ns/1.0">
< title corresp="ref:adams1999"/>
</ bibl>; s.v. < hi xmlns="http://www.tei-c.org/ns/1.0">yirmakka* </ hi>. </ note>
<phon id="B_yirmakka">
<form>yirmakka</form>
THT 3206| - n- tkālñe - nti dharma ka - ʒa ka ste - ñe nañ mā -ce - sañ khya - ā - k-ā -
THT 1450.d| - k- r-ñ- škau sa i - kUce tne sā rñi kwā lype l-e - ntsa a nte ke ñā ssa - ʒu ma šāk maitya -
pelaikne -s- ʒ- e -n- rt-i tāñ ersna we šai - kauwa pyappayantsa - piš yāknēsā ploryaimē - aṃ aṃaik - e -
tre n-a -
THT 3216| - fñantse -rā - pāskeṃ : na - ñ- sk- ntā -y- kau rocā -e -
PK AS 16.2|ke ktsen bram-nākte mant weña || pañdurāñkāhēna || wi-ppewāñne kṣattaryi śpālmeṃ : ñākteṃ śāmnāṃts
yārkwento- aṃaṃ cai : pelaikneṣe wāñtre ceṃ saimtsa : centsak saimtsa kantār se šaiṣṣe : 1 || tusaṃ warñai kṣatyi
poša te nauṃ yārkwento aṃaṃ takāre || tumēṃ mant cai ksa alyaiṃ alannēṃ śāma eñkālñe nāki kīrsormē ālyaucaś
weñare || katarosine || śukentane trenkāltsa perne peṃyo muskitār : eñkālñente ʒarntsa ywārc yārtoṃ lkāntār
wertsyaṃne : calle ʒ wesāṃ miṣenta lauke tarkam eñkālñe : wāto wṣeññai saimtsa wes šayem omte pintwātsa : 1 ||
tumēṃ cai eñkālñente nāki kīrsormē koṭanmasa warrttone lateṃ • tū no kUce yāknēsā || śawāññe-kwanane || pw
eñkālñenta rerinoṃ šaul śāwāñte aīrpāce : bram-nākti ra yayaṭaṃ wāto wṣeññai saim yaucaś : ompalskoññeś
spelkkesoñc kīrsōṃ nāki kleśāmnāṃts : śle-maīyyā ywārc ersante abhiññēnta piš ʒṣp no : 1 || se tane teri ste
ente pañākti šaiṣṣene mā tsōmoṃ takāṃ • twak māka krātayuk preściyaṃne kUce kai oroteste cāmpamñeci bodhisatvi
takāṃ cai ot tāmpak-yāknēsā rṣāki māskentār cenāṃts omte aīrpāce šaul šaitsi skeyessontāṃts cāmpamñe ʒai
tarāltse šaiṣṣe tāntsi ālyine amalākampa tasemane po wāñtarwa lkātsi raddhisa yatiñ nāciyai klautsaīsa
klyaucci kātkor ekamātte karsatsoi • eṃṣke nemce ylai-ñākte bram-ñākteśa warñai ñākteṃts yarkesa yamaṣṣālyi takāṃ
tu yāknēsā aurtšana aiśāmnānta pārkāñ-me • eṃṣke tot naivasamññānsamññāyataṃ tāntsi
THT 1389.n| - ca ke lk- cakene pa -te ye kle - tuwak -
Dd 6.1| - krent ʒpane lka -
PK Bois A27| - kunacāṃttre ākṣa • -
THT 147.5| - repacyeṃ wikastsi -t-ʒc rāmttār no rāse -
THT 1439.f| - cchati gacch- mā - seṃ -
PK Bois A21| - parra tārka se -
THT 2192| - l-a - ʒpā yi -
THT 592.b| - -ntsa snai kārstau snai ʒotri - wātkałtsa śāłtaññe ramtā - aikne cpi aksāṣṣāṃ - ksa cpi nervvanne
ʒāñ o -
THT 1857| -r- tw- ntse - ly- s-
THT 3891.a| - pā - sna - : dvā -
THT 2386.k| - nt-r-e -e -i -I ke -e l-ai -
THT 274|kUce toṃ wñāwa āyorntā maittreyeṃcā śāmtsīsā papāṣṣorñe tonts pontants ścōññā nesāṃ šaiṣṣentse kentsa
pārñna mā stāmoṃ ʒāli stāna onolmi papāṣṣorññesa pārna kUc- āyornta yāmorntasa 30-5 kUce toṃn yāmāṃ wāñtarwa
eyñāke rā ksa śāmnā ñāktents ñākteṃ maittreyeṃ lkāṃ su āksau ñāś centsā tonneṃ kā ks- ālām pālko takāṃ
akñātsaññēsā ñākti lañco wāñtreśi mā cai lkāñ-ne klyomñēsā 30-6 papāṣṣorñe takāṃñ kā ʒeṃ ra ksa cok tāñksā
twāṣṣāṃ ñi śāriputra - maittreyeṃ kUce ysāṣṣāna pyappayin rā kātāṃ ñā -ā - kauc krUī -i - śāmnantso cet- sū
lkātsi 30-7 mā yataṃ sū yāñmātsi meyyāsa epreñēsā ś-t- māktewne lantūññēsā mā r- āmokāñt- āklora ma etreUññai
meyyāsā krent yāmorsa āyorsa yataṃ śāmtai maittreyeṃc 30-8 papāṣṣorññeṣṣe krentāṃ yayaṭaṃñ kołkentsā ai - sā
āyorṣṣe aiśāññeṣṣe pantaintsā śāmnā śāñmeṃ maittreyeṃc po šaiṣṣente - ce saṃtk- ewāñ-m- onwāññe lām saṃsāṣṣe
pelemēṃ 30-9 śleṃ te yatka puññākte śāriputri prāsāñne śkas yāknēsā maīwa keṃ tary yāltse po šaiṣṣenne kodyāññā
sumerñtā naittāre po waskāte -
THT 3506| - ma kā śy- n- m-a mai -
THT 1334.m| - meṃ postāṃ - ʒ- ā ya -
THT 2819.v| - iś -
PK Bois B104|ñi yaitkorsa maṃt pyām - yaṃ cewāmpa klaina - oksaiñ trai • keroccapāṃ -
THT 3830.c| - ma -
THT 2795.y| - wa wā -
THT 3291.g| - śtā - k- ne rā -
THT 2379.u| -p- l-e twe - stha k-a - ne kekmoṃā šaiṣṣe - wpsa ntai : a l-e - lyelkorme tāñwā - l- ʒm- ly- ka
-ts- śc tāññe ʒārmntsa - ceṃnts no tā kUśa ñ- me śpā ke ʒe kā o -
THT 3916.c| - tu vi ru - y- hā -
THT 2762.c| -s- wi - ś- ści -
THT 3204| - tu lmi - l-r pu -
IOL Toch 257| - r- yamaṣṣa - y-ntār n- w- ln- -meṃ ʒai kUce ʒpā pone ñorameṃ ʒai toy - ña stānampa tasemane •
se ʒuktante pañākte -e kārśanalle 10-7 || tvamoghe • dvipa • preñke • trāṃṃ • waste • kṣat- -ṣuwa ra
yamaṣṣeñcañca pelaikñenta aksāṣṣeñca - ñe leleku ʒai cwī no kUśalamūlānta a - laklesā ñeñusku ʒai ñakti aiśai
yamaṣṣa - o -
THT 2377.t| - pā ʒṣā -l- ś- ne - ñāsā lkā - lna ra ki t kātṭwoy yontwe māñce -n-e - ñ-e ne me -t- ñā kte wa -y-
```

Labeled Data

| | TB | TA |
|-----------------|-------------|--------------|
| # types | 11268 | 3955 |
| # tokens | 29298 | 13220 |
| # documents | 7057 | 1635 |
| mean doc length | 20.4 tokens | 33.24 tokens |

41.9%

of samples are annotated (TB)

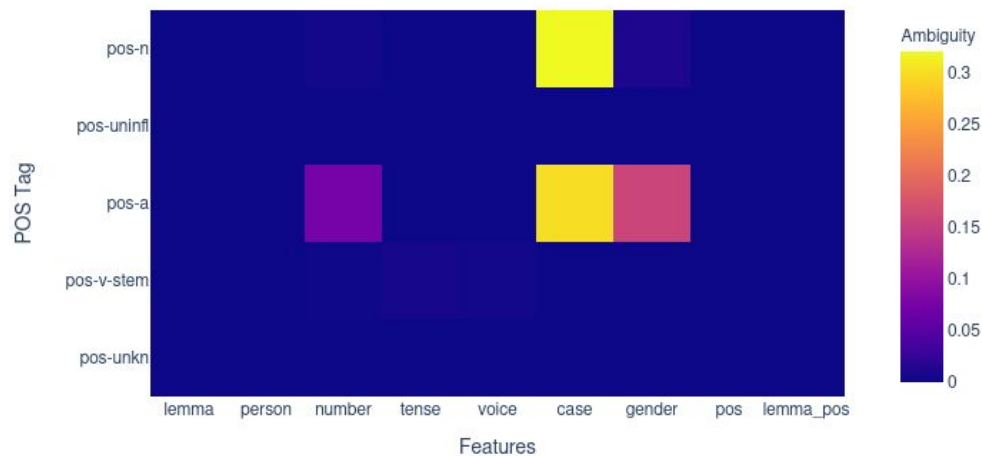
44.4% for TA

| labeled data (N=11268) | unlabeled data (N=15560) |
|------------------------|--------------------------|
| 0.7 | 0.2 |
| .1 | |

Ambiguity: How many words have multiple parses?

- Influences choice of model
 - low ambiguity: character-level
 - high ambiguity: word-level & context
- answer: **30.04%** (3386)
- What kind of ambiguities?

Ambiguity of Tags



şamāññeşşe → şamāne + **nom** + m + adj
şamāññeşşe → şamāne + **acc** + m + adj
şamāññeşşe → şamāne + [**nom-acc**] + m + adj

Approach

CoNLL-SIGMORPHON [4]

- Challenge in 2018
- Task 1: Inflection
 - given: lemma + tags
 - generate: inflected form
 - with low (10^2), medium (10^3) and high (10^4) numbers of training samples
- On a large number of different languages
 - inflective and agglutinative Ls represented

- Adequate reference for this project
 - *this is the inverse problem*
 - *many languages in the challenge (inflective and agglutinative)*
 - *according to this we have a medium to high number of training samples*

A string generation task

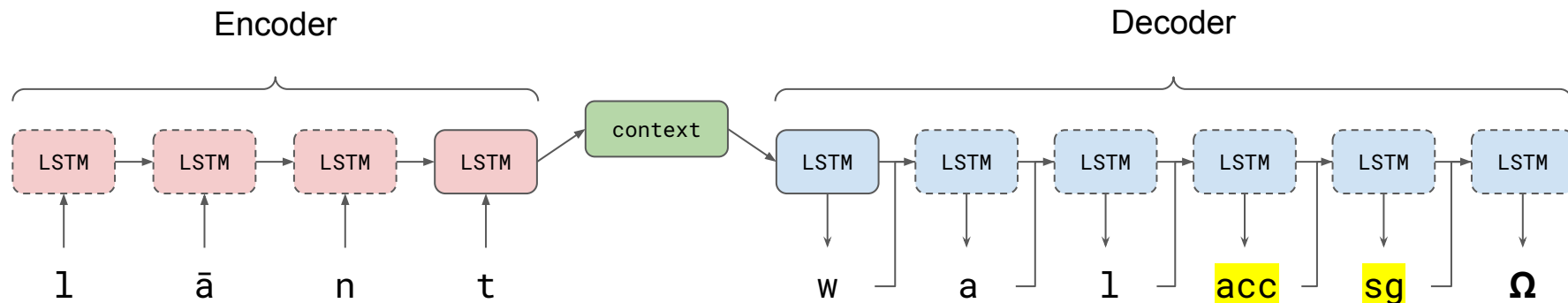
- Common approach to SIGMORPHON
 - and well performing (eg. BME-HAS Acs [5])
- Include the grammatical tags in the target string
 - treat tags as characters
- Train (neural) predictor to generate lemma + tags
 - character by character
 - solves both lemmatization and grammatical tagging

l+ā+n+t+ä+ś → w+a+l+all+sg

| "King" | |
|--------|--------|
| lemma | wal |
| nom.sg | walo |
| acc.sg | lānt |
| all.sg | lāntäs |

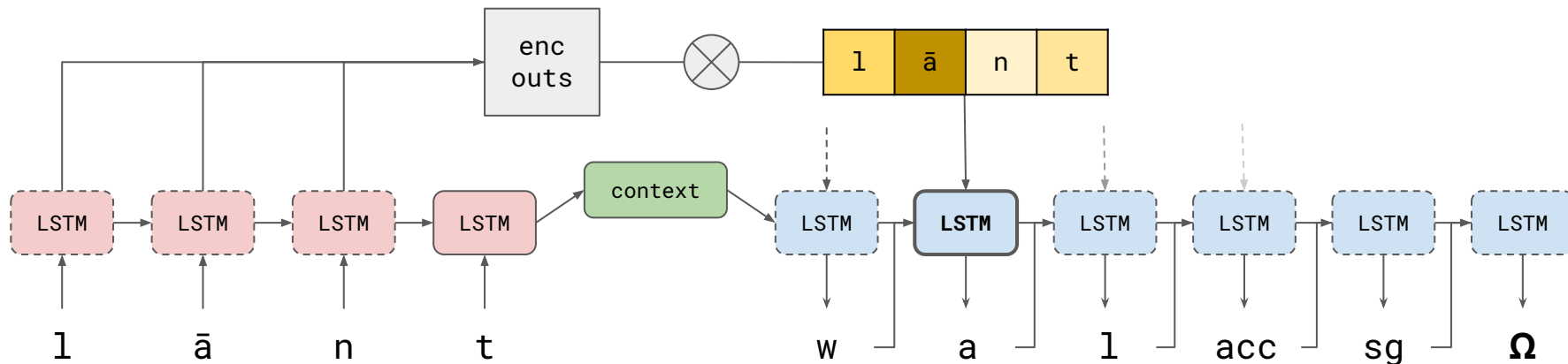
Sequence-to-Sequence Model

- predicts output sequence given input sequence
 - uses some form of RNN
 - enables different input & output lengths
- in this case: encoder-decoder seq2seq model [2]
- common for translation tasks
 - language generation tasks in general
 - often on word-level
 - in our case: character level



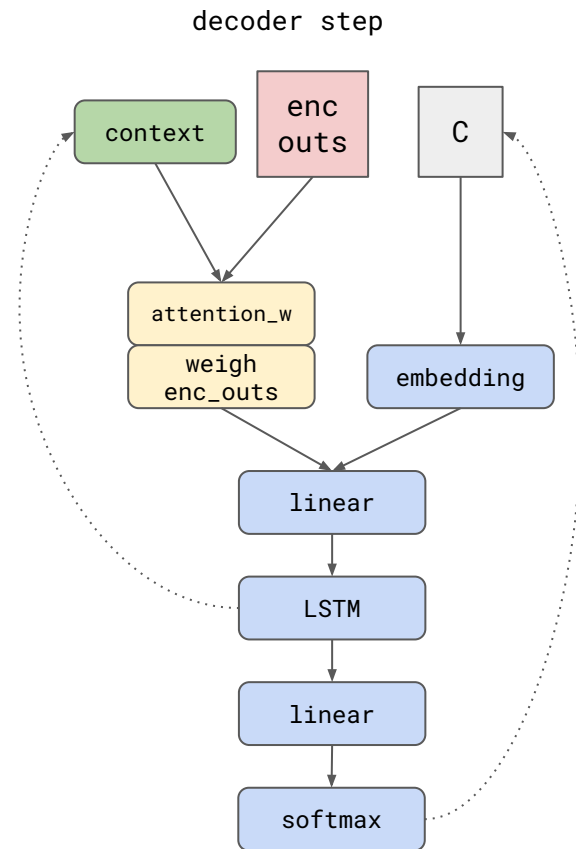
Attention

- Extension of encoder-decoder architecture
 - decoder has access to all encoder outputs
 - weighs encoder outputs (=drawing attention to certain elements)
- Advantages
 - counters long sequence bottleneck
 - provides feedback of the model
 - increase performance



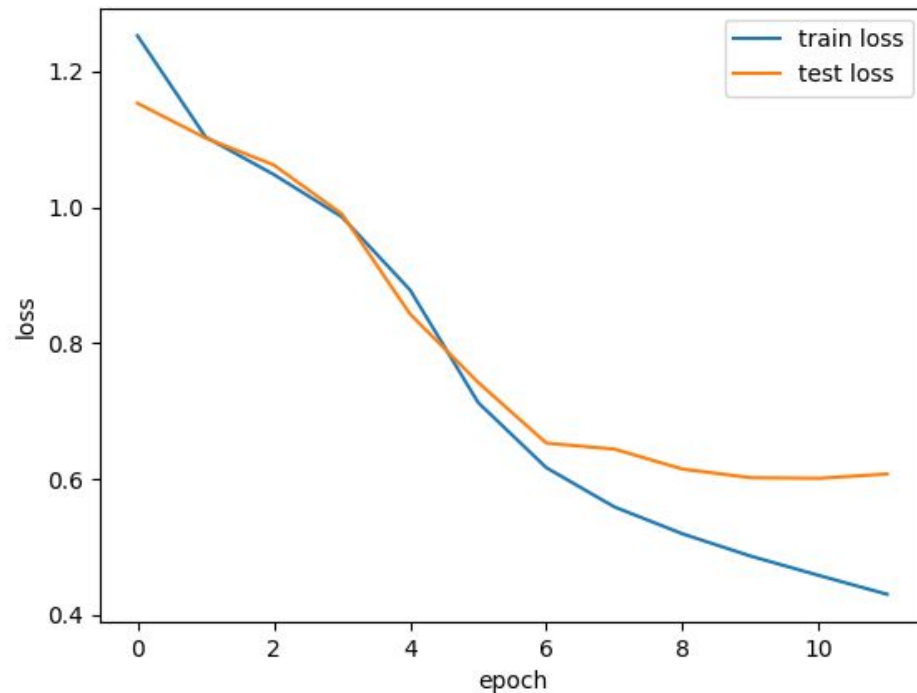
Bahdanau Attention [1]

- “Additive attention”
 - as opposed to multiplicative attention (Luong [3])
- computed in decoder step
- Steps
 - calculate alignment from enc_outs and context
 - softmax alignment (= attention_weights)
 - enc_outs * attention_weights
 - concat attended & embedded



Training Process

| | |
|----------------|---------------|
| loss | cross entropy |
| optimizer | adam |
| learning rate | 1e-4 |
| batch_size | 1 |
| embedding_size | 50 |
| context_size | 100 |



Evaluation & Results

Evaluation

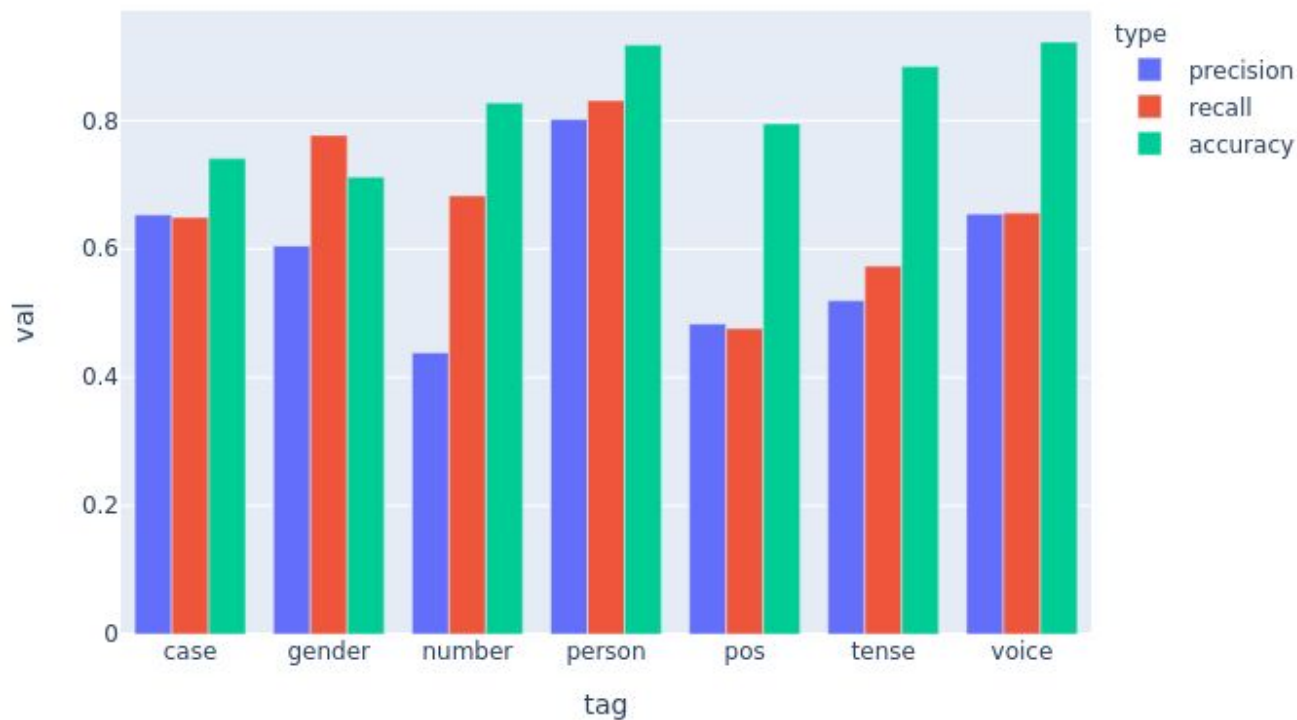
- Split the per-sample evaluation in to two parts
 - Lemma: Levenshtein Distance
 - Gramm.Classes: Precision, Recall, Accuracy
- Levenshtein Distance
 - count of operations to transform given string into target string
 - operations: add, delete, substitute (characters)
- Precision, Recall, Accuracy
 - calculated per gram.category
 - tuning for recall would be ideal

1.21 / 1.87

Levenshtein Distance

test / val

Classification Metrics (TB val set)



| average | |
|-----------|-------|
| precision | 57.2% |
| recall | 60.8% |
| accuracy | 82.2% |

Conclusion

- Scores are nowhere near perfect
 - they do not need to be, since they will be checked
 - team was actually positively surprised with the results
- Benefits of a black box model
 - Team “Tarim Brahmi” did not need to know much about ML
 - I did not have to learn Tocharian
 - = less interdisciplinary work?
- Neural networks are hard to debug
 - have a proper experimental setup
 - do not get carried away
- Interdisciplinary Work
 - Communication is very important (vocabulary!)
 - Ideally: Linguists in control up until preprocessing step
 - How much do they need to know? How much do I need to know?
 - I was possibly very lucky considering the data situation

References

- [1] **Neural Machine Translation by Jointly Learning to Align and Translate** Bahdanau et al. 2015
- [2] **Sequence to sequence learning with neural networks** Sutskever et al. 2014
- [3] **Effective Approaches to Attention-based Neural Machine Translation** Luong et al. 2015
- [4] **The CoNLL–SIGMORPHON 2018 Shared Task: Universal Morphological Reinflection** Cotterell et al. 2018
- [5] **BME-HAS System for CoNLL–SIGMORPHON 2018 Shared Task: Universal Morphological Reinflection** Ács 2018
- [6] **CeTOM Website** <https://www.univie.ac.at/tocharian/?home> (25.05.2021)
- [7] **Project Tarim Brahmi** <https://www.oeaw.ac.at/acdh/projects/tarim-brahmi/> (25.05.2021)