

Recovering Text from Endangered Languages Corrupted PDF documents

Nicolas Stefanovitch

European Commission, Joint Research Center, Ispra, Italy



ComputEL-5 Workshop
May 26-27 2022

Universal Declaration of Human Rights in Nivkh (niv)

Қатыгун сик правоғун Декларация Генеральная Ассамблея резолюция 217А (III) 10 к`лолоң 1948 ань хавыл ғера провозглашайра жад

Нивң қ`атыгун жара, чуғун жара сик намадивңчоғр иввут т`охтта, адяй сик п`ңафқ-ңафқ ңазин правоғун ивңы-худ куғытнд жара салғанжунвд жара сикғун мир жара; адяй пиған сикак к`икр п`ид-жаымди қ`оғл ивкғун нивгун правоғун лақ қаврта жа-жаңы варвар актғундох вылңута, худғундох нивгун сик макыр қ`оғл легу осқита жагут урла қ`оғл ивң нивгунах п`ңалзайгута; адяй айф тымди мир айғай – ху мирух нивгун кулғытң туғс жара ахтнд жара ивңы адяй к`лунд жара жекинд жара ағзута – худғун нивгун язра ныйнынд; адяй нивгун жекирқир п`сиңрукғундох озн вопун ухмуна жа ғавргуйнырқир нивгун право властығун торкир ытңуна; адяй қ`атыгун п`ңафқ-ңафқрох ңафқңалагуйнын ихмына жагуйн; адяй объединенных наций қ`атыгун п`Уставух адяй қ`атыгун правоғун ахтхымлыта, нивгун п`и намад ивд яймта азымығундиң ңаңғундиң салғат жумта жатот социальный прогресс жагуйныңа адяй нивгун малғо куғытнд ивт жунвд ургута; адяй государствоғун – членгун Организация объединенных наций ройныт т`охтта,

:àòüãóí ñèè ìðàâî4óí Äãëëàðàöëý
Äãíàðàëüíàý Äññàíàëëäý ðãçîëpöëý 217A (III)
10 è`ëîëî2 1948 àíü 1àâüë 4ãðà
ìðîîîçãëèàèèðà 1àä

íèà2 6`àòüãóí 1àðà, ÷ó4óí 1àðà ñèè íàìàèà2÷í43 èàãóò
ò`í5òòà, àäýë ñèè ì`2àò6-2àò6 2àçèí ìðàâî4óí èà2ü-òóä èó4üòíà
1àðà ñàè8àí1óíàä 1àðà ñèè4óí ìèð 1àðà; àäýë ìè8àí ñèèàè è`èè3
ì`èà-1àüíàè 6`ì8è èàè4óí íèããóí ìðàâîãóí èà6 6àãðòà 1à-1à2ü
àððààð àèò4óíàí5 àüè2óòà, òóä4óíàí5 íèããóí ñèè ìàèüð 6`ì8è
èããó ññ6èòà 1àãóò óðèà 6`ì8è èà2 íèããóíà5 ì`2àèçàèãóòà; àäýë
àèò üüíàè ìèð àè8àè – òó ìèðóò íèããóí èóè4üò2 óó4ñ 1àðà à5òíà
1àðà èà2ü àäýë è`éóíà 1àðà 1àèèíà 1àðà à4çóòà – 1óä4óí íèããóí
ýçðà íüéíüíà; àäýë íèããóí 1àèè3èè3 ì`ñè2ðóè4óíàí5 íçí àñíóí
óòíóíà 1à 9àãððãóéíü3èè3 íèããóí ìðàâî àèàñòü4óí óìèèè3 üò2óíà;
àäýë 6`àòüãóí ì`2àò6-2àò6ðí5 2àò62àèàãóéíü èòíüíà 1àãóéí;
àäýë íáüããèíáíüò íàèè 6`àòüãóí ì`Õòòàãóò àäýë 6`àòüãóí
ìðàâî4óí à5òðüíèüòà, íèããóí ì`è íàìà èàà ýéíà àçíüòü4óíàè2
3à294óíàè2 ñàè8àò 1óíà 1àòíò ñíòèàèüíüé ìðíðãññ 1àãóéíü2à
àäýë íèããóí ìàè8í èó4üòíà èàò 1óíàà óðãóòà; àäýë àññòààðòòàí4óí
– ÷èáíãóí ìðãñèçàèè íáüããèíáíüò íàèè ðíéíüò ò`í5òòà,

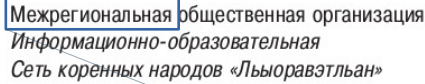
Universal Declaration of Human Rights in Nivkh (niv)

Қатыгун сик правоғун Декларация Генеральная Ассамблея резолюция 217А (III) 10 к`лолоң 1948 ань жавыл ғера провозглашайра жад

Нивң қ`атыгун жара, чуғун жара сик намадивңчоғр` иввут т`охтта, адяй сик п`ңафқ-ңафқ ңазин правоғун ивңы-худ куғытнд жара салғанжунвд жара сикғун мир жара; адяй пиған сикак к`икр` п`ид-жаымди қ`оғл ивкғун нивгун правоғун лақ қаврта жа-жаңы варвар актғундоҳ вылңута, худғундоҳ нивгун сик макыр қ`оғл легу осқита жагут урла қ`оғл ивң нивгунаҳ п`ңалзайгута; адяй айф тымди мир айғай – ху мирух нивгун кулғытң туғс жара ахтнд жара ивңы адяй к`лунд жара жекинд жара ағзута – худғун нивгун язра ныйнынд; адяй нивгун жекиркір` п`сиңрукғундоҳ озн вопун ухмуна жа ғавргуйныркір` нивгун право властьғун торкир` ытңуна; адяй қ`атыгун п`ңафқ-ңафқроҳ ңафқңалагуйнын ихмына жагуйн; адяй объединенных наций қ`атыгун п`Уставух адяй қ`атыгун правоғун ахтхымлыта, нивгун п`и намад ивд яймта азмытғундиң раңғундиң салғат жумта жатот социальный прогресс жагуйныңа адяй нивгун малғо куғытнд ивт жунвд ургута; адяй государствоғун – членгун Организация объединенных наций ройныт т`охтта,

Қ'атыгун сик правоғун Декларация Генеральная Ассамблея резолюция 217 А (III) 10 қ'лолоң 1948 ань жавыл ғера провозуглавшайра жад

Нивң қ'атыгун жара, чуғун жара сик намадивңчоғр` иввут т'охтта, адяй сик п'ңафқ-ңафқ ңазин правоғун ивңы-худ куғытнд жара салғанжунвд жара сикғун мир жара; адяй пиған сикак к'икр` п'ид-жаымди қ'оғл ивкғун нивгун правоғун лақ қаврта жа-жаңы варвар актғундож вылңута, худғундож нивгун сик макыр қ'оғл легу осқита жагут урла қ'оғл ивң нивгунаж п'ңалзайгута; адяй айф тымди мир айғай – ху мирух нивгун кулғытң туғс жара ахтнд жара ивңы адяй к'лунд жара жекинд жара ағзута – худғун нивгун язра ныйнынд; адяй нивгун жекиркір` п'сиңрукғундож озн вопун ухмуна жа ғавргуйныркір` нивгун право властьғун торкир` ытңуна; адяй қ'атыгун п'ңафқ-ңафқрож ңафқңалагуйнын ихмына жагуйн; адяй объединенных наций қ'атыгун п'Уставух адяй қ'атыгун правоғун ахтхымлыта, нивгун п'и намад ивд яймта азмытғундиң раңғундиң салғат жумта жатот социальный прогресс жагуйныңа адяй нивгун малғо куғытнд ивт жунвд ургута; адяй государствоғун – членгун Организация объединенных наций ройныт т'охтта,



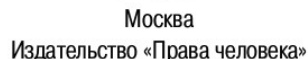
```
<pages>
<page id="1" bbox="0.000,0.000,822.047,566.929" rotate="0">
<textbox id="0" bbox="145.561,481.195,327.469,515.196">
<textline bbox="145.969,505.195,327.469,515.196">
<text font="FLIIF+PragmaticaCondC" bbox="145.969,505.195,152.819,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">M</text>
<text font="FLIIF+PragmaticaCondC" bbox="153.069,505.195,157.279,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">e</text>
<text font="FLIIF+PragmaticaCondC" bbox="157.529,505.195,162.849,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000"></text>
<text font="FLIIF+PragmaticaCondC" bbox="163.099,505.195,167.589,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">p</text>
<text font="FLIIF+PragmaticaCondC" bbox="167.839,505.195,172.849,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">e</text>
<text font="FLIIF+PragmaticaCondC" bbox="173.089,505.195,178.099,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">r</text>
<text font="FLIIF+PragmaticaCondC" bbox="178.339,505.195,183.349,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">i</text>
<text font="FLIIF+PragmaticaCondC" bbox="183.589,505.195,188.599,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">o</text>
<text font="FLIIF+PragmaticaCondC" bbox="188.839,505.195,193.849,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">h</text>
<text font="FLIIF+PragmaticaCondC" bbox="194.089,505.195,198.339,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">n</text>
<text font="FLIIF+PragmaticaCondC" bbox="198.589,505.195,202.539,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">y</text>
<text font="FLIIF+PragmaticaCondC" bbox="202.789,505.195,207.149,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">h</text>
<text font="FLIIF+PragmaticaCondC" bbox="207.399,505.195,211.439,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">a</text>
<text font="FLIIF+PragmaticaCondC" bbox="211.689,505.195,215.739,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000">j</text>
<text font="FLIIF+PragmaticaCondC" bbox="215.989,505.195,218.319,515.196" colourspace="ICBased" ncolour="[0.13725, 0.12157, 0.12549]" size="10.000"> </text>
```

Сив я' латдан илена
ненэция'' права
лэтрамбава езмня
сертавы декларация

Тиканэ резолюцияда 217 А (III)
Генеральной Ассамблея нарка пэвдя
Иры юдимтей яля 1948 похна нямвы

```
<textbox id="1" bbox="109.392,311.031,316.352,403.031">
<textline bbox="109.392,383.031,316.352,403.031">
<text font="FLIIDD+Nenezred" bbox="109.392,383.031,122.632,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:2)</text>
<text font="FLIIDD+Nenezred" bbox="122.632,383.031,133.572,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:3)</text>
<text font="FLIIDD+Nenezred" bbox="133.572,383.031,143.332,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:4)</text>
<text font="FLIIDD+Nenezred" bbox="143.332,383.031,150.692,403.031"
colourspace="Separation" ncolour="1" size="20.000"></text>
<text font="FLIIDD+Nenezred" bbox="150.692,383.031,161.052,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:5)</text>
<text font="FLIIDD+Nenezred" bbox="161.052,383.031,170.812,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:6)</text>
<text font="FLIIDD+Nenezred" bbox="170.812,383.031,178.172,403.031"
colourspace="Separation" ncolour="1" size="20.000"></text>
<text font="FLIIDD+Nenezred" bbox="178.172,383.031,189.112,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:7)</text>
<text font="FLIIDD+Nenezred" bbox="189.112,383.031,199.432,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:8)</text>
<text font="FLIIDD+Nenezred" bbox="199.432,383.031,207.972,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:9)</text>
<text font="FLIIDD+Nenezred" bbox="207.972,383.031,219.532,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:10)</text>
<text font="FLIIDD+Nenezred" bbox="219.532,383.031,229.852,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:8)</text>
<text font="FLIIDD+Nenezred" bbox="229.852,383.031,240.792,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:11)</text>
<text font="FLIIDD+Nenezred" bbox="240.792,383.031,248.152,403.031"
colourspace="Separation" ncolour="1" size="20.000"></text>
<text font="FLIIDD+Nenezred" bbox="248.152,383.031,259.092,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:3)</text>
<text font="FLIIDD+Nenezred" bbox="259.092,383.031,270.632,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:7)</text>
<text font="FLIIDD+Nenezred" bbox="270.632,383.031,280.372,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:12)</text>
<text font="FLIIDD+Nenezred" bbox="280.372,383.031,291.312,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:11)</text>
<text font="FLIIDD+Nenezred" bbox="291.312,383.031,301.632,403.031"
colourspace="Separation" ncolour="1" size="20.000">(cid:8)</text>
<text font="FLIIDD+Nenezred" bbox="301.632,383.031,308.992,403.031"
colourspace="Separation" ncolour="1" size="20.000"></text>
<text font="FLIIDD+Nenezred" bbox="308.992,383.031,316.352,403.031"
colourspace="Separation" ncolour="1" size="20.000"></text>
```

Extracted text: 000000 0000 0000 00 0000 00



1	2	l	—	?	?
Н	н	l	—	ë	А
?	?	?	?	?	?
Н	О	П	С	Т	У
?	?	?	?	?	?
д	е	ж	з	и	й
.	?	?	1	2	?
с	т	у	ф	х	ц
?	?				
ю	я				

Corrupted Nenets (yrk) font

d	;	t	ê	H
d	;	t	ê	H
N	ú	E	Q	O
N	ú	E	Q	O
A	c	G	l	-
A	c	G	l	-
.	ç	h	1	ã
.	ç	h	1	ã

Valid Portuguese (por) font

Recovery process

```
...  
line 90 => 1 10 6 7 4 3 17  
line 91 => 3 6 5 6 5 4  
line 92 => 1 10 3 4 8 2 3 11 2 4  
line 93 => 4 3 4 9 6 6 5 4  
line 94 => 6  
line 95 => 9 5 4  
line 96 => 1 11 3 3 11 9 5 6  
line 97 => 1 3 3 3 12 2 4 12 2 4  
line 98 => 11 12 5 10 11  
line 99 => 2 5 2 8 7 5 2 6  
line 100 => 1 10 14 5 4  
line 101 => 1 11 3 6 8 2 11 5  
line 102 => 5 4 2 11 11 5 3 5 5  
...
```

Unicode sequence (input from user):

Организация Объединенных Наций цельгундож
ёскиндфурнд

→ token lengths: [11, 12, 5, 10, 11]

CID sequence (extracted from the document) + automatic guess of space and dot characters:

:24:48:35:32:45:40:39:32:54:40:63 :24:33:58:37:36:40:45:37:45
:45:59:53 :23:32:54:40:41 :54:37:43:60:8:51:45:36:46:9 :16:49:42:40
:45:36:52:51:48:45:36

→ token lengths: [11, 12, 5, 10, 11]

Unique match!

```
new 24 -> О  
new 48 -> р  
new 35 -> г  
new 32 -> а  
new 45 -> н  
new 40 -> и  
new 39 -> з  
...
```

28 symbols recovered → 34% of the characters

I.0001: :23:40:34:6 :10:15:32:50:60:35:51:45 :5:32:48:32:2 :55:51:8:51:45 :5:32:48:32 :49:40:42 :45:32:44:32:36:40:34:6:55:46:8:7 :40:34:34:51:50
I.0002: :50:15:46:9:50:50:32:2 :32:36:63:41 :49:40:42 :47:15:6:32:52:10:3:6:32:52:10 :6:32:39:40:45 :47:48:32:34:46:8:51:45 :40:34:6:59:53:51:36 :42:51:8:59:50:45:36
I.0003: :5:32:48:32 :49:32:43:12:32:45:5:51:45:34:36 :5:32:48:32 :49:40:42:8:51:45 :44:40:48 :5:32:48:32 :32:36:63:41 :47:40:12:32:45 :49:40:42:32:42 :42:15:40:42:7
I.0004: :47:15:40:36:3:5:32:59:44:36:40 :10:15:46:12:43 :40:34:42:8:51:45 :45:40:34:35:51:45 :47:48:32:34:46:35:51:45 :43:32:10 :10:32:34:48:50:32 :5:32:3:5:32:6:59
I.0005: :34:32:48:34:32:48 :32:42:50:8:51:45:36:46:9 :34:59:43:6:51:50:32:2 :53:51:36:8:51:45:36:46:9 :45:40:34:35:51:45 :49:40:42 :44:32:42:59:48 :10:15:46:12:43



I.0001: Ни:34:6 :10:15:а:50:ьгун :5:ара:2 :55:уғун :5:ара сик на:44:ади:34:6:55:оғ:7 и:34:34:у:50
I.0002: :50:15:оҗ:50:50:а:2 адяй сик :47:15:6:аф:10:3:6:аф:10 :6:азин :47:ра:34:оғун и:34:6:ыхуд куғы:50:нд
I.0003: :5:ара сал:12:ан:5:ун:34:д :5:ара сикғун :44:ир :5:ара адяй :47:и:12:ан сикак к:15:ик:7
I.0004: :47:15:ид:3:5:аы:44:ди :10:15:о:12:л и:34:кғун ни:34:гун :47:ра:34:огун ла:10 :10:а:34:р:50:а :5:а:3:5:а:6:ы
I.0005: :34:ар:34:ар ак:50:ғундох :34:ыл:6:у:50:а:2 худғундох ни:34:гун сик :44:акыр :10:15:о:12:л



Нивң қ`атьгун жара, чуғун жара сик намадивңчоғр иввут
т`охтта, адяй сик п`ңафқ-ңафқ ңазин правоғун ивңы-худ куғытнд
жара салғанжунвд жара сикғун мир жара; адяй пиған сикак к`икр
п`ид-жаымди қ`оғл ивкғун нивгун правоғун лақ қаврта жа-жаңы
варвар актғундох ылңута, худғундох нивгун сик макыр қ`оғл

221=>жымди қ`оғл уйгид
Ниғвң дуфтож вылңуд Санги
Нивң қ`атьгун жара, чуғун жара сик намадивңчоғр
Организация Объединенных Наций цельғундох ёскиндфурнд
эна положенияғун ивғай напы п`жатыжать қаврна, п`ңафқ-ңафқ
самоуправляющаяся жа ғаврд лу,
Генеральная Ассамблея туң сик
Декларация задача жагун провозглашайдра
Чу п`жоғара, сикак маңра жаң общество
жекинд
удовлетворить
зақоўид.
Произвольно
Техническое
п`Уставух
Конституцияғир
Ин
Қ`атьңгун
Эғлгун
Ыткғун
Раңғымкмунд

```
...
line 175 => 7 6 5 8 4 2 6 3 4
line 176 => 8 11 2 4 7 5 5 6
line 177 => 10 6 6 5 4
...
line 220 => 12 7 8 4 10 4 7
line 221 => 5 5 6
line 222 => 7 5 12 4 1 15
...
```


Automatically guesses characters for space and dot

Sorts characters by number of instances + displays sample lines containing them:

```
* cid 25 count: 2
line 101 :25:роиз:34:ольно :5:а:50 на:7:50:и:6 нарух:50:и:6 йа гра:38:данс:50:34:о ас:10:а:44
line 209 :25:15:и :47:ра:34:оғун :5:ара куғы:50:ндғун :5:ара осу:57:ес:50:34:ляй:6:ы ни:34:гун
```

Sorts lines, by number of unrecovered characters + display sample:

```
line 172 ['6', '44', '5', '5'] роғуйны:6 :44:едицинское :5:ара социальное :5:ара иғрығрынд
line 34 ['44', '6', '6', '34'] леле:44:а:6:6 образо:34:ание
```

Checks consistency of user input

Automatically guesses upper case letters based on document statistics

```
GUESS CAPS
try: :1:аңы
MATCH! жаңы
new 1 -> Ж
```

2022
2021
2020
2019
2017
2016
2015
2014
2013
2012
2011

News

2022

February 12, 2022: The Gilyak [niv] translation is available for review, thanks to Nicolas Stefanovitch.

February 12, 2022: Corrections to the Asturian [ast] translation, thanks to David Mediavilla.

February 11, 2022: The Nenets [yrk] translation is available for review, thanks to Nicolas Stefanovitch.

n	1	2	3	4	5	6
niv	10	33	72	93	97	98
yrk	17	29	60	86	98	100

Table 1: Proportion in percent of unique sequences of token lengths for sequence length n in the UDHR for two different languages

language	niv	yrk
unique characters	81	68
text length (words)	1430	1530
input length (words)	57	76

Table 2: Unique character count and total word length for the UDHR declaration in two languages, and the total number of input words necessary for full recovery of these texts

Discussion, conclusion, questions...

See you at the poster presentation! :-)