

# Pre pokročilých používateľov

15th February 2024



# Vaši sprievodcovia pre dnešnú cestu



**Mirjam**  
User Success  
Team



**Sara**  
User Success  
Team

(Dušan Katuščák, doplnky a preklad)

[s.mansutti@readcoop.eu](mailto:s.mansutti@readcoop.eu)

[m.elattal@readcoop.eu](mailto:m.elattal@readcoop.eu)

# Digitálna knižnica - Texty

## **Sprístupnenie z digitálnych repozitárov:**

- 1. Digitálna knižnica obrázková (len nasnímané obrázky)**
- 2. Digitálna knižnica plnotextová (full texts) (obrázky +  
OCR/HTR)**
- 3. Digitálna knižnica hybridná (čiastočne obrázky + čiastočne  
OCR/HTR)**

# Obsah

- 1. Úvod
- 2. Trénovanie & Značkovanie/Tagovanie
- 3. Analýza rozloženia & Základné čiary
- 4. Polia modelov (Beta) & Modely tabuliek
- 5. Zverejňovanie – Stránky Transkribus
- 

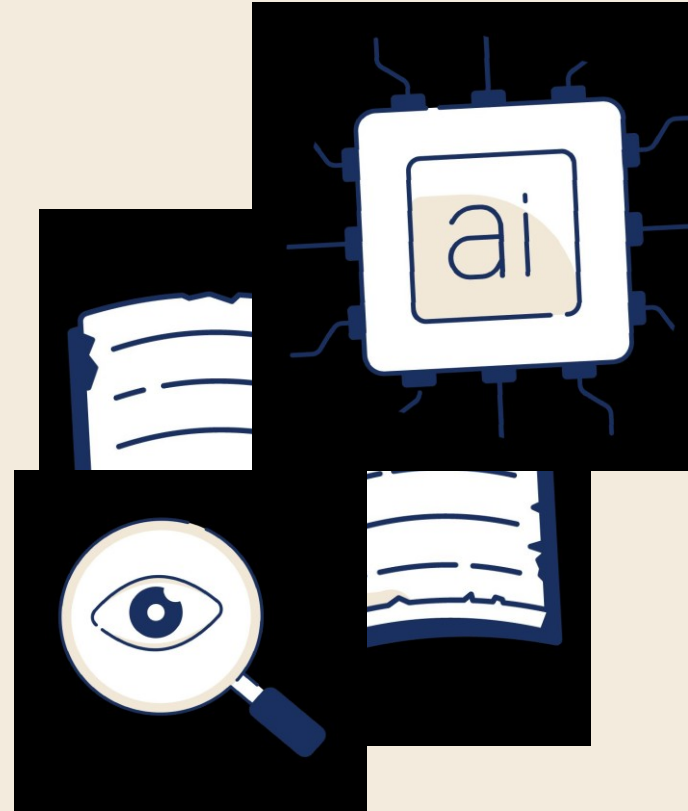




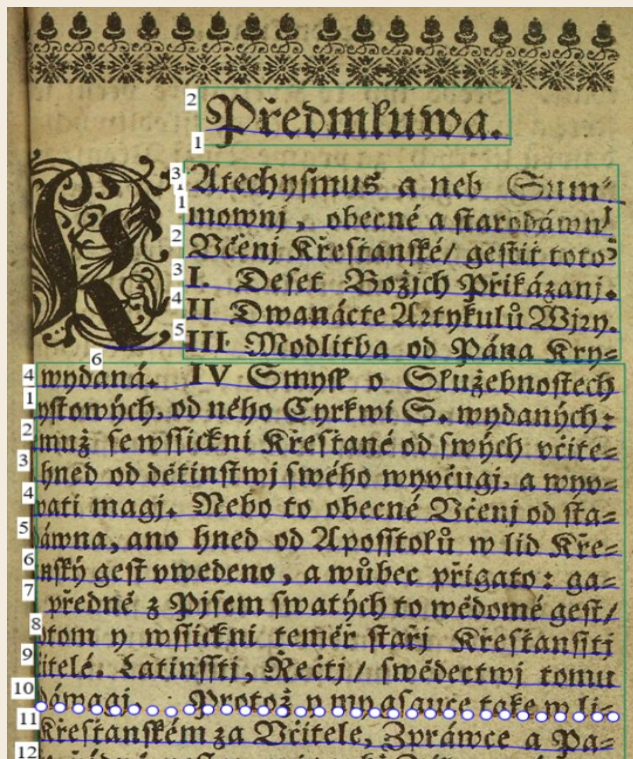
Úvod

# Čo je Transkribus?

Transkribus je váš partner, ktorý pomocou umelej inteligencie (AI) zjednodušuje časovo náročnú a namáhavú prácu s historickými dokumentmi.



## České dokumenty – práca študentov na Opavskej univerzite (Halfarová)



### 1 Předmluva.

- 1 Atechysmus a neb Sum
- 2 mowni, obecné a starodawn | |
- 3 Vcenj Křestanské / gestit toto
- 4 I. Deset Božich Prikázanj.
- 5 II Dwanácte Artykulu Wjry.
- 6 III Modlitba od Pána Kry=

- 1 wydanä. IV. Smysk o. Skuzebnostech
- 2 ystowých. od něho Cyrkwí S. wydaných:
- 3 muzi se wssickni Křestane od swych: učite= | |
- 4 hned od dětinstwí swého wyučuj, a wyo=
- 5 wati magj. Nebo to obecné Vcenj od sta= |
- 6 dawna, ano hned od Aposstolů w lid Kre= |
- 7 nský gest wvedeno, a wübec prigato: ga= | |
- 8 předné z Písem swatých to wědomě gest /
- 9 ptom y wssickni teměr staří Křestansiti
- 10 litelé: Latinssti, Recti / swědectwi tomu
- 11 dáwaj. Protož y myjsauče take wli-

## České dokumenty – práca študentov na Opavskej univerzite (Taufrová)



- 1 Čert a Prawda.
- 2 To gest:
- 3 welmi pěkně smyšlené, utěšené
- 4 Hystorie a Rozprávky
- 5 w několika stech,
- 6 kterěz se pro wyraženi mysli a pro
- 7 zasmáni při dlouhé chvíli, y w každé
- 8 weselé společnosti, časem také y drobet
- 9 pro wybrauseni rozumu dobře užije-
- 10 wati mohau.
- 11 Wydané podlé rukopisu
- 12 Hylarya, Jokosa, Astucya.
- 13 Gakozto II. Díl
- 14 k Zrcadlu Possetilostj.

Region 2

Region 3

- 1 Kraméryusowým nákladem.

Region 4

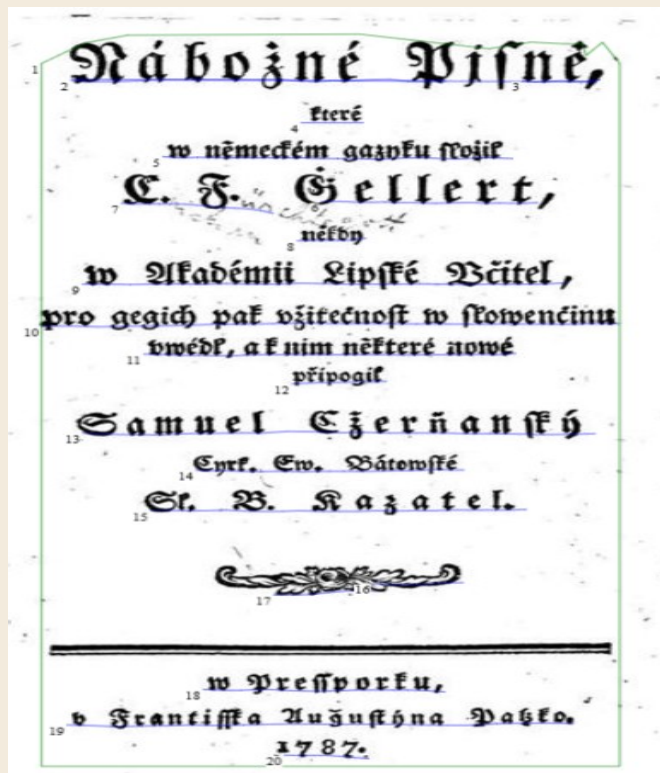
Region 5

- 1 W Praze, 1796.
- 2 kdostání w České Expedycy w Dominy-
- 3 kánské ulicy, u Hrabů w Nie. 373.





# České dokumenty – práce studentov na Opavskéj univerzite (Gajdošová)

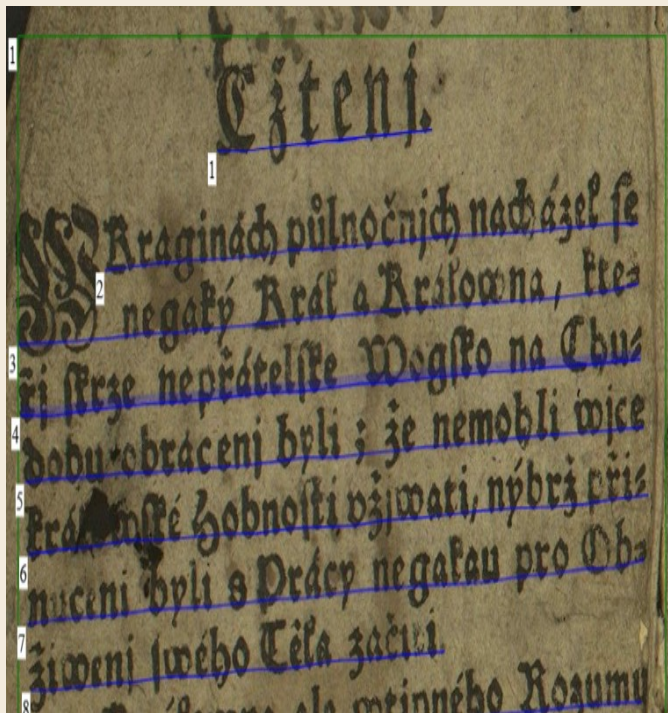


Region 1

## Nábožné Písně

které  
w německém gazyku složil  
Gellert  
C. J.  
někdy  
w Akademii Lipské Učitel,  
pro gegich pak užitečnost w slowenčinu  
uwědl, a k nim některé nové  
připogil  
Samuel Cžerňanský  
Cyrk. Ew. Bátowské  
Sl. B. Kazatel.

W Pressporku,  
U Frantířka Augustýna Pazko.  
1787

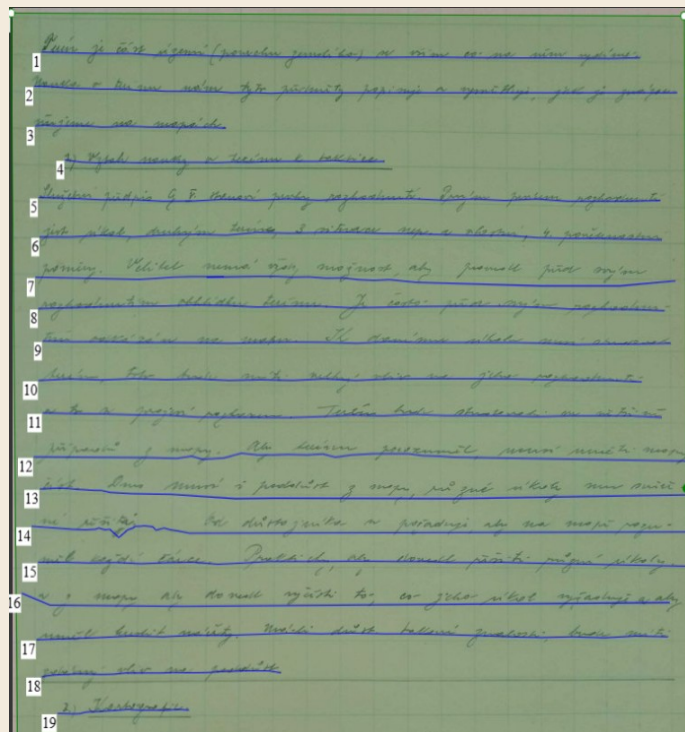


Region 1

Čtení.

W kraginách půlnočnjich nacházel se  
negaký Král a Králowna, kte-  
řj skrze nepřátelske Wogsko na Chu-  
dobu obráčení byli ; že nemohli wjce  
králowské hodnosti užjwati, nýbrž při-  
nuceni byli s Prácy negakau pro Ob-  
živeni swého Těla začiti.

## České dokumenty – práce studentů na Opavské univerzitě (Němec)



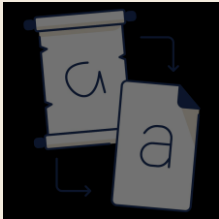
Region 1

- 1 Terén je část území (povrchu zemského) se vším co na něm vydíme
- 2 Nauka o terénu nám tyto předměty popisuje a vysvětluje, jak je znázor-
- 3 ňujeme na mapách
- 4 1) Vztah nauky o terénu k taktice.
- 5 Služební předpis G V. stanoví prvky rozhodnutí. Prvým prvkem rozhodnutí
- 6 jest úkol, druhým terén, 3 situace nep. a vlastní, 4. povětrnostní
- 7 poměry. Velitel nemá vždy možnost, aby provedl před svým
- 8 rozhodnutím obhlídku terénu. Je často před svým rozhodnu-
- 9 tím odkázán na mapu. K danému úkolu musí stanovit
- 10 terén, toto bude mít velký vliv na jeho rozhodnutí.
- 11 a to se projeví rozkazem. Terén bude sledovati ve většině
- 12 případů z mapy. Aby terénu porozuměl, musí uměti mapy
- 13 číst. Dnes musí i provést z mapy různé úkoly mu svěře-
- 14 ně řešiti. Od důstojníka se požaduje, aby na mapě rozu-
- 15 měl každé čárce. Prakticky, aby dovedl řešiti různé úkoly.
- 16 a z mapy aby dovedl vyčísti to, co jeho úkol vyžaduje a aby
- 17 uměl kreslit náčrty. Má-li důstojník takové znalosti, bude mít
- 18 zdárný vliv na poddůst.
- 19 2, Kartografie

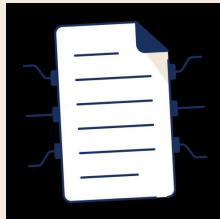
# Začínáme: Prvé kroky v Transkribuse

- 1. Registrácia a prehľad používateľského rozhrania
- 2. Vytvorenie zbierky
- 3. Nahrávanie súborov
- 4. Použitie kreditu

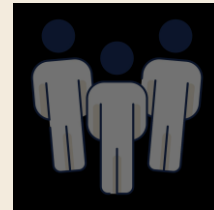
# Čo Transkribus umožňuje?



**Manuálny a automatický prepis  
ručne písaných a tlačených  
dokumentov**



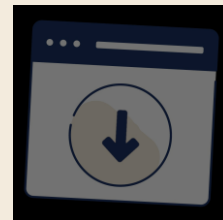
**Trénovanie modelov  
umelej inteligencie**



**Spolupráca**



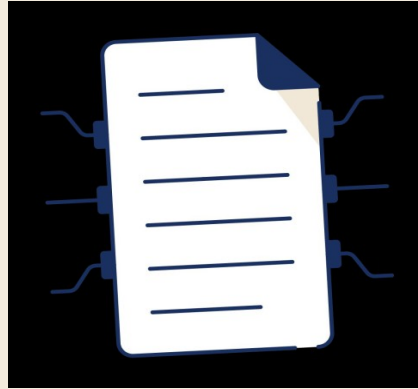
**Tagovanie štruktúry  
a obsahu  
dokumentov**



**Export dokumentov v  
rôznych formátoch**

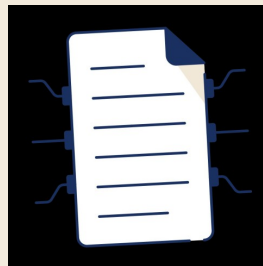
# Trénovanie modelov umelej inteligencie

**Strojové učenie:**



**Umožňuje strojom učiť sa z (označených alebo neoznačených) údajov, identifikovať vzorce a robiť predpovede s minimálnym zásahom človeka.**

# Trénovanie modelov AI



## Modely umelej inteligencie:

algoritmy vytvorené počas tréningového procesu systému strojového učenia

predstavujú výstup tréningu/školenia      získané vedomosti.

<https://help.transkribus.org/text-recognition>



# Trénovanie modelov AI

- **Ground Truth** (Training Data, Základná pravda):  
Označené údaje pre tréning, ktoré umožňujú modelu identifikovať vzory a robiť predpovede pre tieto označenia na základe nových údajov.  
= všetky strany, ktoré boli prepísané ručne

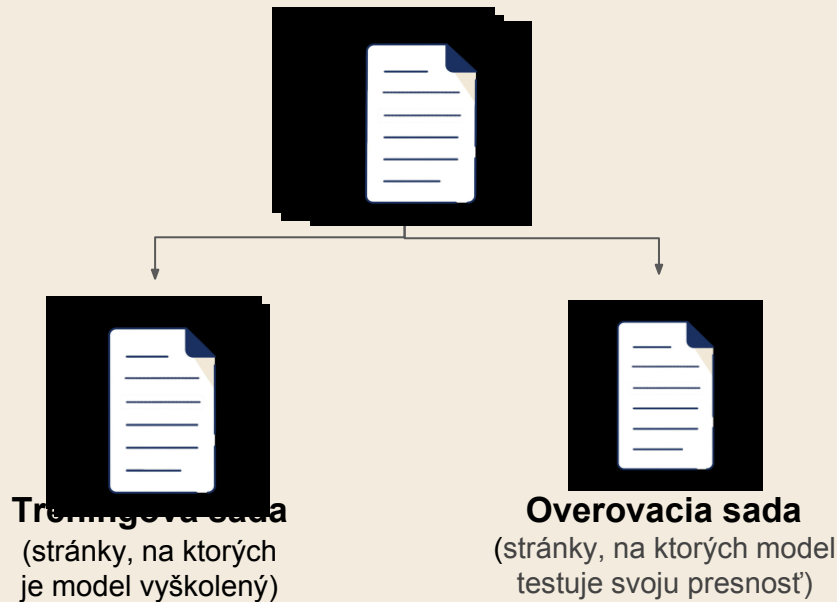
## Tréningová sada (Training set)

Súbor príkladov, ktoré sa používajú na úpravu parametrov modelu  
= dáta, na ktorých sú postavené poznatky v neurónovej sieti

- **Overovacia sada (Validation Set)**

Súbor príkladov, ktoré sa používajú na objektívne posúdenie výkonnosti modelu  
= údaje použité na doladenie parametrov modelu počas jeho tréningu

## Ground Truth (Základná pravda)



Dobrá overovacia sada: to je 10% tréningovej sady + obsahuje všetky príklady (znaky, glyfy)

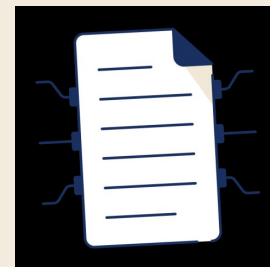


Tréning modelov

<https://help.transkribus.org/text-recognition>



# Trénovanie modelov AI



## Modely trénovateľné s Transkribusom:

Text

Riadky

Bloky textu

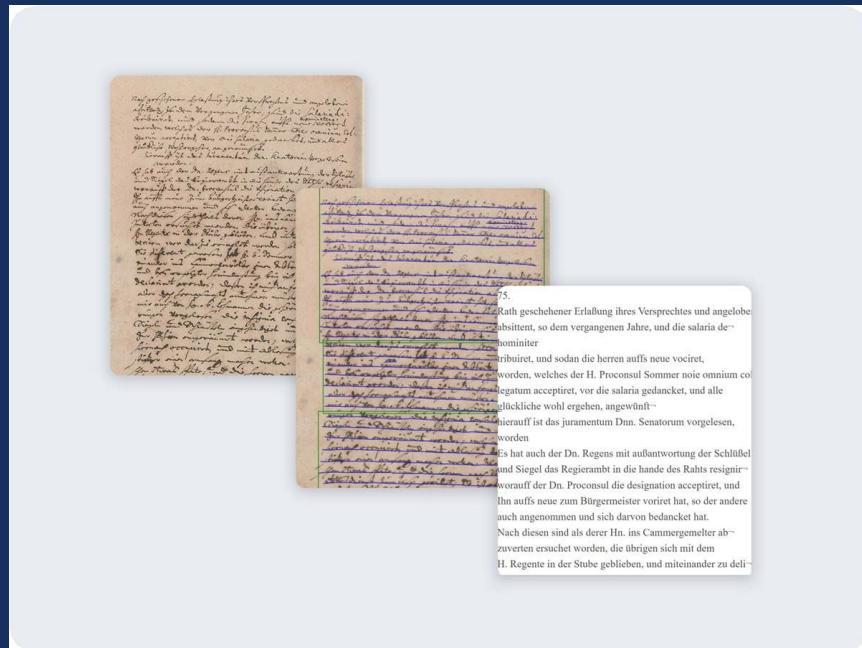
The screenshot shows the Transkribus model training interface. On the left, there is a list of model types: Text Recognition Model, Baselines Model, Field Model, and Table Model. On the right, there is a green button labeled '+ Train New Model'. Arrows point from the labels 'Text', 'Riadky', and 'Bloky textu' to the corresponding model types in the list.

- Text Recognition Model
- Baselines Model
- Field Model
- Table Model

+ Train New Model

# Analýza rozloženia (segmentácia)

- 1. Automatická analýza rozloženia (segmentácie)
- 2. Rozšírené nastavenia konfigurácie rozloženia (segmentácie)
- 3. Manuálna úprava rozloženia (segmentácie)
- 4. Základné modely
- 5. Modely polí
- 6. Tabuľky Modely
- 7. Noviny



# Trénovanie textových modelov

# Modely textu

Pred tréningom modelu:

potrebujete **25 až 75 strán (5000-15000 slov)** prepísaného materiálu (**GT\_Základná pravda**), v závislosti od typu dokumentu (tlačený alebo písaný rukou)

**2 možnosti:**

**1. Ručný prepis stránky**

<https://help.transkribus.org/transcribing-manually>

**2. Použitie hotového modelu**, ktorý bol trénovaný na podobnom skripte (ak je k dispozícii) a manuálna oprava prepisu



# Textové modely

## 1. možnosť: manuálny prepis dokumentov

1. Vyberte stránky, ktoré chcete zahrnúť do **GT\_Základnej pravdy**
2. Spustíte rozpoznávanie rozloženia textu – segmentácia (Layout Recognition)
3. Prepísať od začiatku:

**Označte** slová, ktoré nemôžete prečítať ako nejasné alebo "medzera"

**Riadky**, ktoré zostali prázdne: sa v tréningu neberú do úvahy

**Skratky**: udržiavané/riešené/označené: záleží na tom, čo očakávate ako konečný výstup

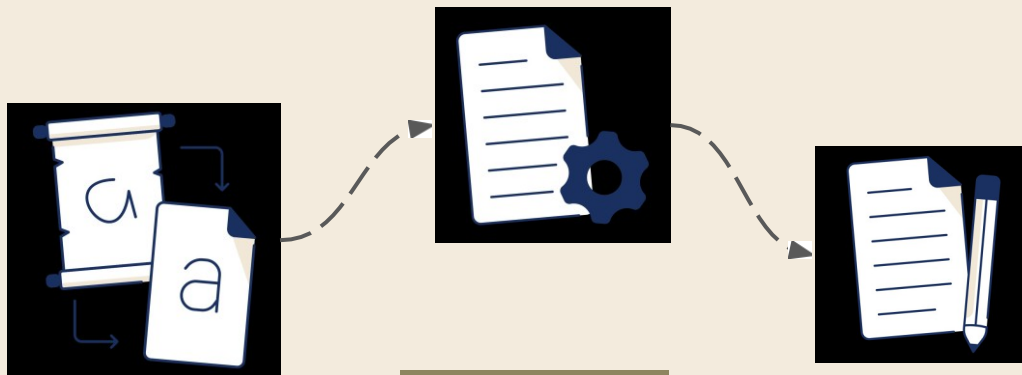
**Uložte stránku ako GT "Základnú pravdu"!**



# Textové modely

**2. možnosť:** použitie modelu/supermodelu a následná oprava automatických prepisov

1. Vyberte stránky, ktoré chcete zahrnúť do Základnej pravdy (GT)
2. Spustenie rozpoznávania textu
3. Oprava automatických prepisov
4. Uložte stránku ako "Základnú pravdu,, (GT)

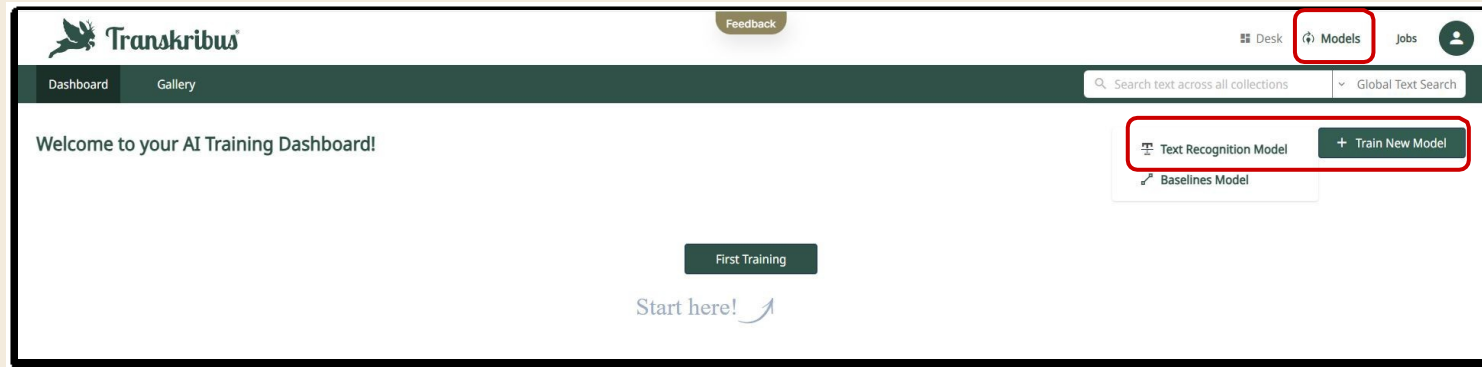




# Textové modely

Po vytvorení prepisov (Základná pravda):

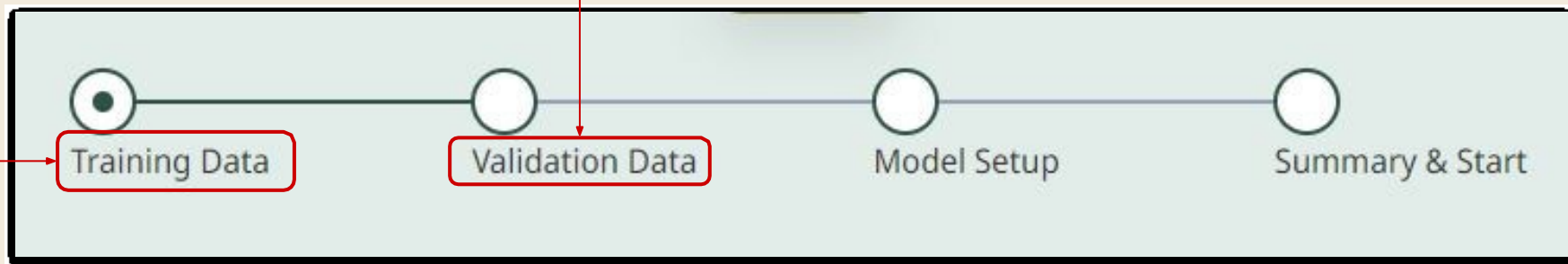
- prejdite do sekcie "Modely"
- kliknite na "Train New Model - Text Recognition Model"
- vyberte zbierku s prepismi (Základná pravda) Ground Truth



# Textové modely

○ Vyberte stránky na:

1. Tréning/školenie (stránky, na ktorých je model školený)
2. Validáciu (strany, na ktorých model testuje svoju presnosť). Dobrá validačná sada: 10% tréningovej sady + obsahuje všetky príklady




# Textové modely

Rozšířené možnosti:

Base Model Recommended

Select a pre-existing model to use as the base for your own model.

 [Select Model](#)

Advanced Settings (optional) ^

Training Cycles optional

Training Cycles

Enter the number of times you want the model to go through the entire training dataset.

Early stopping optional

Early stopping

Enter when you want to use early stopping to prevent overfitting.

Reverse Text (RTL) Optional

Select if you want the text to be written in a right-to-left direction.

# Textové modely

## Rozšírené možnosti:

- **Základný model (Base model):** pomocou základného modelu (Base model) tréning nezačína od nuly, ale od toho, čo sa už naučilo v tréningovom procese tohto modelu



# Textové modely

Rozšírené možnosti:

- **Tréningové cykly (Training cycles (epochs)):** Maximálny počet prechodov modelu cez celú množinu tréningových údajov. Pri prvom tréningu ponechajte predvolený počet 100 tréningových cyklov
- **Predčasné zastavenie (Early stopping):** Minimálny počet cyklov tréningu. Predvolená hodnota je 20: ak po 20 epochách CER validačnej sady neklesne, tréning sa zastaví

# Textové modely

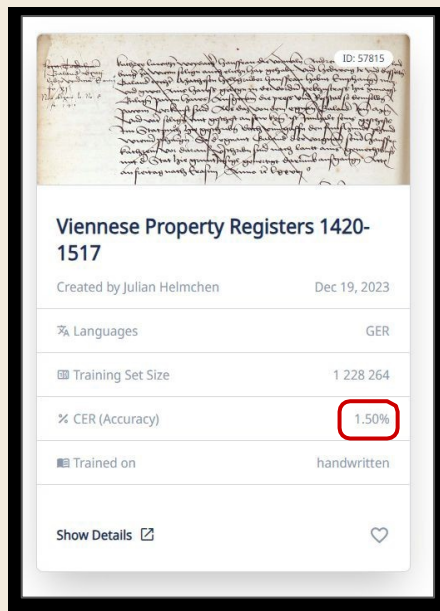
## Rozšírené možnosti:

- **Obrátený text (Reverse text (RTL)):** Ak bol text na obrázku napísaný sprava doľava, ale v textovom editore bol prepísaný zľava doprava
- **Použitie existujúcich polygónov (Use existing line polygons):** Pozn.: používať iba v prípade, že ste upravili mnohoúhelníky v *Transkribus Expert*
- **Tréning s rozpisom skratiek (Train Abbrevs with expansion):** Trénuje model tak, aby automaticky označoval skratky a pridal ich rozpis
- **Vynechať riadky s tagmi nejasné/medzera (Omit lines by tag unclear/gap):** Táto možnosť vynecháva riadky obsahujúce slová označené ako gap/unclear.

# Textové modely

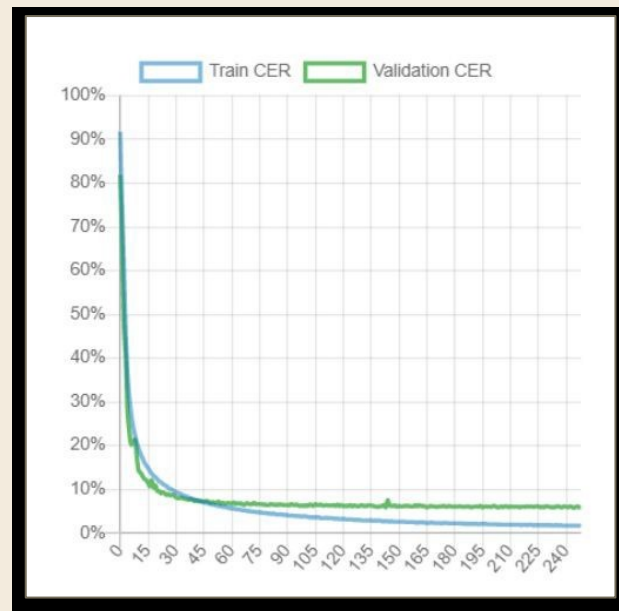
Po dokončení tréningu sa môžete pozrieť na podrobnosti modelu:

1. CER (Chybovosť znakov = Character Error Rate)
2. Krivka učenia



The screenshot shows a model card for 'Viennese Property Registers 1420-1517'. The card includes a thumbnail of a handwritten document, the model name, creator (Julian Helmchen), creation date (Dec 19, 2023), languages (GER), training set size (1 228 264), and CER (Accuracy) of 1.50%. The CER value is highlighted with a red circle. The model is trained on handwritten data.

Property	Value
Created by	Julian Helmchen
Created	Dec 19, 2023
Languages	GER
Training Set Size	1 228 264
CER (Accuracy)	1.50%
Trained on	handwritten



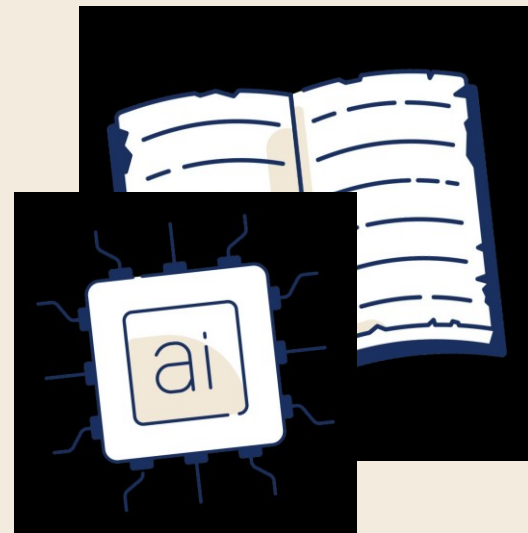
# Textové modely

	CER (chybovosť znakov)	Tréningová sada
<b>Tlačný text</b>	0,5-2%	~ 5.000 words / 25 pages
<b>Jedna ruka</b> - jednoduché písanie	2-4%	10.000+ words / 50+ pages
<b>Niekoľko rúk</b> - zistené	4-6%	10.000+ words per hand / 150+ pages
<b>Veľa rúk</b> - z toho istého obdobia a regiónu – nie všetky zistené počas tréningu	6-8%	100.000+ words / 500+ pages



# Textové modely

- Ruky, ktoré nie sú nijako zistené, alebo načmárané poznámky oveľa horšie výsledky, tak potom:
- Zdvojnásobte počet tréningových dát 20-25% zníženie chybovosti
- **Existujúce modely** sa môžu použiť ako východiskový krok (Base model - základný model) na zníženie požadovaného množstva nových údajov



# Textové modely

Verejný holandský rukopisný vzor: [Dutch Margaretha Turnor 17th Century](#)

Trained by The Utrecht Archives; Training set: 178 pages, Validation set: 20 pages

## Dutch Margaretha Turnor 17th



by The Utrecht Archives

Nov 28, 2022

🌐 Languages

DUT

📄 Training Set Size

36 289

📊 % CER (Accuracy)

3.10%

📅 Centuries

17

📖 Trained on

handwritten

# Model ID

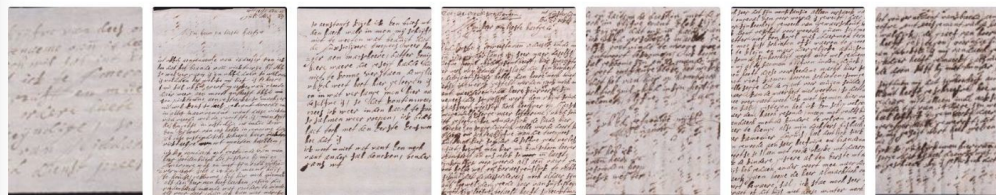
48329

### Model description

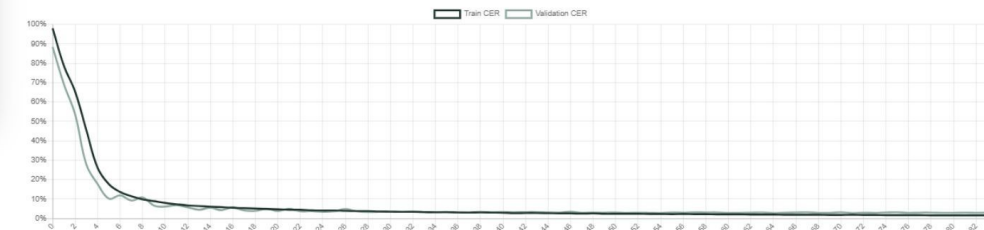
This is the first model created by the Utrecht Archives. It is based on a thousand letters Margaretha Turnor wrote to her husband during the late 17th century. She managed the castle of Amerongen, while her husband worked abroad as a diplomat for the Dutch Republic. Her letters provide an insight into family life in the Dutch Republic as well as the political situation in the country.

### Training data

[View all >](#)



### Training stats



# Textové modely

Verejný model írskej gaelčiny: [Irish, Gaelic and Roman type \(Seanchló agus Cló Rómhánach\)](#)

Trained by Gerard Farrell; Training set: 243 pages, Validation set: 3 pages

Public Model

**Irish, Gaelic and Roman type (Seanchló agus Cló Rómhánach) v.3**

by farrelgn@tcd.ie Nov 4, 2023

🌐 Languages IRI

📄 Training Set Size 70 965

📊 CER (Accuracy) 1.20%

📖 Trained on print

[Edit](#) [Show Description](#)

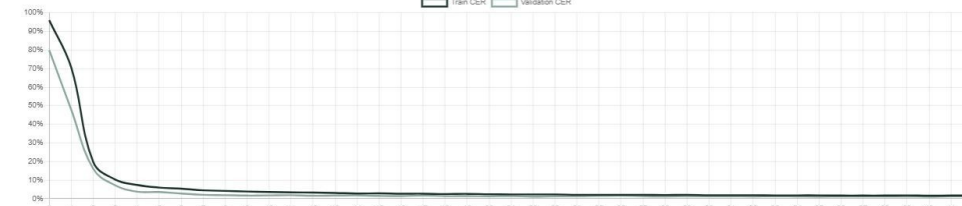
### Model description

Model for reading Irish Gaelic (Gaelige) type or seanchló (common pre-mid-20th century). Can also read Irish in the standard Roman typeface used today. This model was trained on over 70,000 words of material in various typefaces from the 17th century to the early 20th, leaning more heavily towards books published from the mid-19th century in Cló Newman. The model can, however, handle text printed in earlier fonts, such as Cló Petrie, which was used in O'Donovan's edition of the Annals of the Four Masters, and the earlier Cló Moxon used in Bedell's Irish version of the Old Testament (1685). Dotted consonants are transcribed as the consonant followed by a 'h', following modern Irish convention, and the Tironian 'y' is transcribed as 'agus'. Around 30% of the training material also consisted of modern printed Irish texts.

### Training data



### Training stats



Iteration	Train CER	Validation CER
0	95%	85%
1	25%	20%
2	15%	12%
3	12%	10%
4	11%	10%
5	10%	10%
6	10%	10%
7	10%	10%
8	10%	10%
9	10%	10%
10	10%	10%
11	10%	10%
12	10%	10%
13	10%	10%
14	10%	10%
15	10%	10%
16	10%	10%
17	10%	10%
18	10%	10%
19	10%	10%
20	10%	10%
21	10%	10%
22	10%	10%
23	10%	10%
24	10%	10%
25	10%	10%
26	10%	10%
27	10%	10%
28	10%	10%
29	10%	10%
30	10%	10%
31	10%	10%
32	10%	10%
33	10%	10%
34	10%	10%
35	10%	10%
36	10%	10%
37	10%	10%
38	10%	10%
39	10%	10%
40	10%	10%
41	10%	10%
42	10%	10%



## Tagovanie/Značkovanie

<https://help.transkribus.org/tagging>

# Tagging

## a. Štrukturálne tagy (Structural Tags):

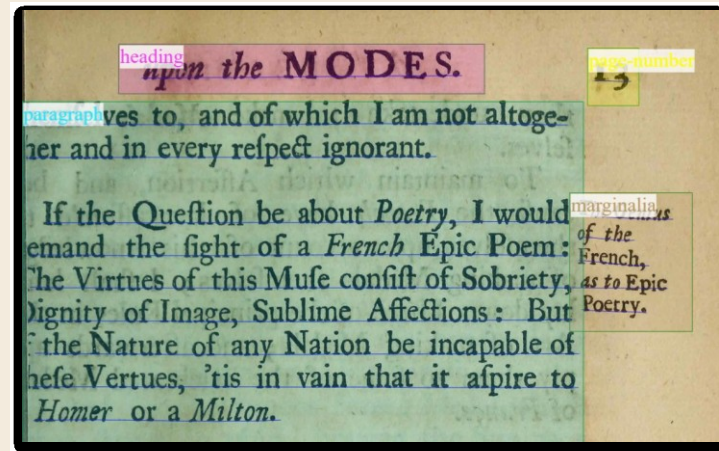
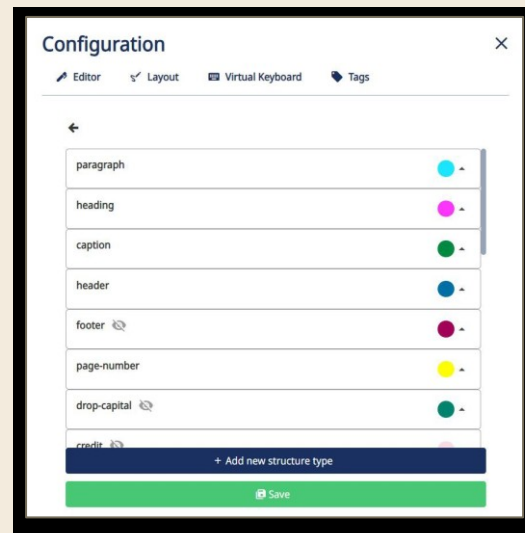
Slúžia na označenie prvkov štruktúry dokumentu

Editor dokumentov: prejdite na **Konfigurácia**  
Rozloženie (Layout)

Riadenie typov štruktúry (Manage Structure Types)

Povoľte viditeľnosť značiek, ktoré chcete použiť/pridajte ďalšie značky

Vyberte tvar , kliknite pravým tlačidlom myši a pridajte štrukturálnu značku



# Tagovanie/značkovanie

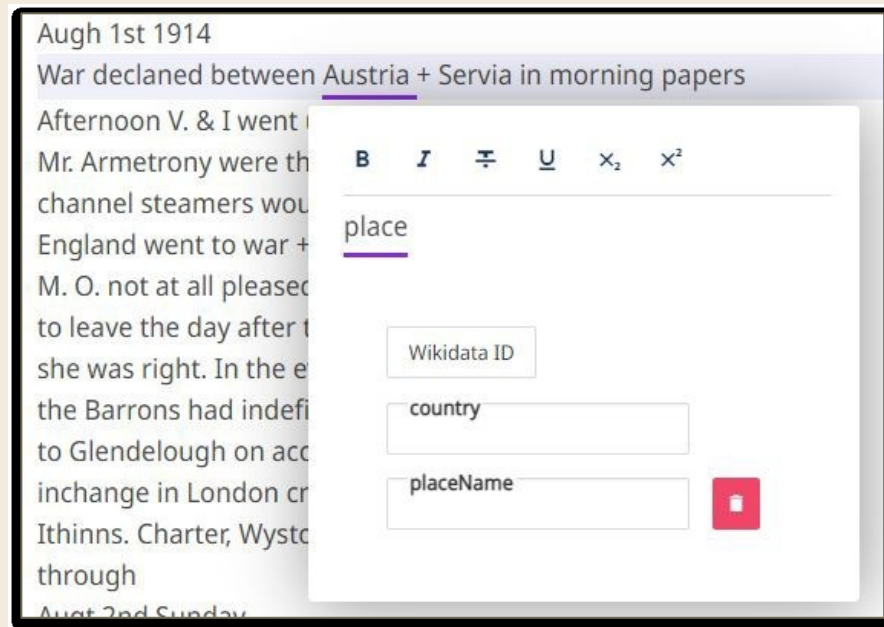
**b. Textové tagy/značky:** slúžia na označenie prepisu a pridanie atribútov vo vnútri textov

Textový editor: v editore vyberte kurzorom slovo, kliknite na príslušnú značku a pridajte vlastnosť

Správa textových značiek:

**Konfigurácia** Upravujte značky v nastaveniach kolekcie: pridávajte / odstraňujte značky a upravujte atribúty

[Example](#)



The screenshot shows a text editor interface. The text being edited is: "Augh 1st 1914 War declaned between Austria + Servia in morning papers Afternoon V. & I went Mr. Armetrony were th channel steamers wou England went to war + M. O. not at all pleased to leave the day after t she was right. In the e the Barrons had indefi to Glendelough on acc inchange in London cr Ithinns. Charter, Wyste through Augt 2nd Sunday". The word "Austria" is highlighted with a purple underline. A context menu is open over the word, showing a toolbar with icons for Bold (B), Italic (I), Strikethrough (ABC), Underline (U), Subscript (x₂), and Superscript (x²). Below the toolbar, the word "place" is entered and underlined. There are three input fields: "Wikidata ID", "country", and "placeName". A red square icon with a white document symbol is also visible.

# Skratky

According to your needs, you can decide to train the model to:

1. **Ponechajte skrátenú formu v prepise: jednoducho prepíšte skratky ako sú v dokumente**

Nerozpisujeme

output: Skratka v texte

2. **Rozpisovanie skratiek:** Neurónové siete sú často schopné naučiť sa rozpoznávať a používať rozšírenia, najmä ak sa objavujú často napíšte rozšírenie skratky do prepisu, venujte dôslednú pozornosť

Rozpisujeme skratky (pozorne, rovnako)

output: Skratky. + rozšírenia v texte

3. **Tagujeme a trénujeme skratky vrátane rozpisu :** označte skratku a pridajte zodpovedajúci rozpis do vlastnosti "Rozšírenie" Pri tréňovaní modelu vyberte možnosť tréňovať skratky

Tagy vrátane rozšírení

output: možnosť získať iba skratky, skratky. po ktorých nasledujú ich rozpis alebo náhrada

# Skratky

V konfigurácii tréningu  
začiarknite políčko

**Train Abbrevs with  
expansion  
(Trénovať model s  
rozpisom skratiek)**

**Text Recognition Model**

Training Data ✓ Validation Data ✓ Model Setup ○ Start ○

Remove Title

X Diary of John Henry Fisher - Copy

< Back

English ⓘ

Search

Centuries

Base Model Recommended

Select a pre-existing model to use as the base for your own model.

Select Model

**Advanced Settings (optional)**

Training Cycles optional

100

Enter the number of times you want the model to go through the entire training dataset.

Early stopping optional

20

Enter when you want to use early stopping to prevent overfitting.

Reverse Text (RTL) optional

Select if you want the text to be written in a right-to-left direction.

Use existing line polygons for training optional

**Train Abbrevs with expansion** optional

Omit lines by tag optional

unclear


gap



# Skratky

- Verejný model [UCL–University of Toronto #7](#) trénovaný na riešenie skratiek v stredovekých rukopisoch
  - Training set: 330 pages, Validation set: 30

**UCL–University of Toronto #7**




by Bentham Project (University College London), DEEDS-project (University of Toronto) Dec 13, 2022

🌐 Languages	LAT
📄 Training Set Size	140 158
📊 % CER (Accuracy)	1.70%
📅 Centuries	13-15
📖 Trained on	handwritten
# Model ID	48734


**Model description**

Seventh iteration of the collaborative UCL–University of Toronto model for processing medieval Latin manuscripts, particularly those containing a large quantity of abbreviated words. E-mail: [criley@ucl.ac.uk](mailto:criley@ucl.ac.uk).

**Training data** View all >



**Training stats**




Epoch	Train CER	Validation CER
0	100%	100%
1	~50%	~50%
2	~10%	~10%
3	~5%	~5%
4	~3%	~3%
5	~2%	~2%
6	~1.7%	~1.7%
7	~1.7%	~1.7%
8	~1.7%	~1.7%
9	~1.7%	~1.7%
10	~1.7%	~1.7%
11	~1.7%	~1.7%
12	~1.7%	~1.7%
13	~1.7%	~1.7%
14	~1.7%	~1.7%
15	~1.7%	~1.7%
16	~1.7%	~1.7%
17	~1.7%	~1.7%
18	~1.7%	~1.7%
19	~1.7%	~1.7%
20	~1.7%	~1.7%
21	~1.7%	~1.7%
22	~1.7%	~1.7%
23	~1.7%	~1.7%
24	~1.7%	~1.7%
25	~1.7%	~1.7%
26	~1.7%	~1.7%
27	~1.7%	~1.7%
28	~1.7%	~1.7%
29	~1.7%	~1.7%
30	~1.7%	~1.7%
31	~1.7%	~1.7%
32	~1.7%	~1.7%
33	~1.7%	~1.7%
34	~1.7%	~1.7%
35	~1.7%	~1.7%
36	~1.7%	~1.7%
37	~1.7%	~1.7%
38	~1.7%	~1.7%
39	~1.7%	~1.7%
40	~1.7%	~1.7%
41	~1.7%	~1.7%
42	~1.7%	~1.7%
43	~1.7%	~1.7%
44	~1.7%	~1.7%
45	~1.7%	~1.7%
46	~1.7%	~1.7%
47	~1.7%	~1.7%
48	~1.7%	~1.7%
49	~1.7%	~1.7%
50	~1.7%	~1.7%
51	~1.7%	~1.7%
52	~1.7%	~1.7%
53	~1.7%	~1.7%
54	~1.7%	~1.7%
55	~1.7%	~1.7%
56	~1.7%	~1.7%
57	~1.7%	~1.7%
58	~1.7%	~1.7%
59	~1.7%	~1.7%
60	~1.7%	~1.7%
61	~1.7%	~1.7%
62	~1.7%	~1.7%
63	~1.7%	~1.7%
64	~1.7%	~1.7%
65	~1.7%	~1.7%
66	~1.7%	~1.7%
67	~1.7%	~1.7%
68	~1.7%	~1.7%
69	~1.7%	~1.7%
70	~1.7%	~1.7%
71	~1.7%	~1.7%
72	~1.7%	~1.7%
73	~1.7%	~1.7%
74	~1.7%	~1.7%
75	~1.7%	~1.7%
76	~1.7%	~1.7%
77	~1.7%	~1.7%
78	~1.7%	~1.7%
79	~1.7%	~1.7%
80	~1.7%	~1.7%
81	~1.7%	~1.7%
82	~1.7%	~1.7%
83	~1.7%	~1.7%
84	~1.7%	~1.7%
85	~1.7%	~1.7%
86	~1.7%	~1.7%
87	~1.7%	~1.7%
88	~1.7%	~1.7%
89	~1.7%	~1.7%
90	~1.7%	~1.7%
91	~1.7%	~1.7%
92	~1.7%	~1.7%
93	~1.7%	~1.7%
94	~1.7%	~1.7%
95	~1.7%	~1.7%
96	~1.7%	~1.7%
97	~1.7%	~1.7%
98	~1.7%	~1.7%
99	~1.7%	~1.7%
100	~1.7%	~1.7%

[Example](#)

# Skratky

- Model trévaný na stredovekých latinských dokumentoch (1520) na rozpoznávanie značky "skratka" vrátane vlastníctva "rozpisu skratiek" Training set: 177 pages, Validation set: 30 pages

Hertziana_1520_abbrevs	
	
🌐 Languages	VAR
📄 Training Set Size	59 225
📊 CER (Accuracy)	19.80%
📖 Trained on	handwritten
# Model ID	38873

[Example](#)

مکتب جمعی

دینی ، اجتماعی ، تربیوی ، ادبی ، علمی و فنی در .

سنه : ۱ — ۱۵ اگستوس ۱۳۳۶ — نومرو : ۱

## مقصد (۱)

خلفك افكارینی تئور مقصدیله چقاردیمیز بو مجموعه قارشینده دهرین برهمنولیت دویوز و بولی  
قیمتلی بروظیفه تلقی ایدیورز . اکر ، بووظیفه ایلهده ، مملکتك پك محتاج اولدیهی معاری دویمولرینی  
فعالیت دولرینی . . . اویاندیرمه خدمت ایده بیلیرمهك بختیارز .  
« مکتب جموعه سی » ، ساغاد ، مطبعه وسائره اجر تئرنیک پك بهالی اولدیهی بوزمانده عرفان تولید  
ایتمك ، هرکس ایچون فاندولی اولمق، اوزره ساحه انتشاره آتمق .

# Tréningové modely pre RTL písmo

# RTL skripty

**5 verejných modelov** ([public models](#)) pre rôzne RTL skripty v Transkribus

2 verzie osmansko-tureckého tlačového modelu

Vaybertaytsh typ písma (jidiš)

Rukopis jidiš (model Dybbuk)

Zmes historických hebrejských písiem a jazykov (DiJeSt 2.0)



The screenshot shows the 'Text Recognition' section of the Transkribus interface. On the left, there are filters for 'Favorite Models' (0), 'Public Models' (5), and 'Private Models' (20871). Below these is a search bar and a 'Languages' filter. The main area displays a table of models with columns for Name, Words, Language, and CER.

Name	Words	Language	CER
OttomanTurkish_Print_v2	248 083	TUR	7.60%
Vaybertaytsh.YidTakNL	66 497	YID, HEB	0.90%
OttomanTurkish_Print_1	180 854	TUR	7.20%
The Dybbuk for Yiddish Handwriting	144 985	YID	4.40%
DiJeSt 2.0	773 726	HEB, YID, LAD, JUD	2.00%

# RTL skripty

Ako v súčasnosti prepisovať a trénovať údaje RTL v Transkribuse:

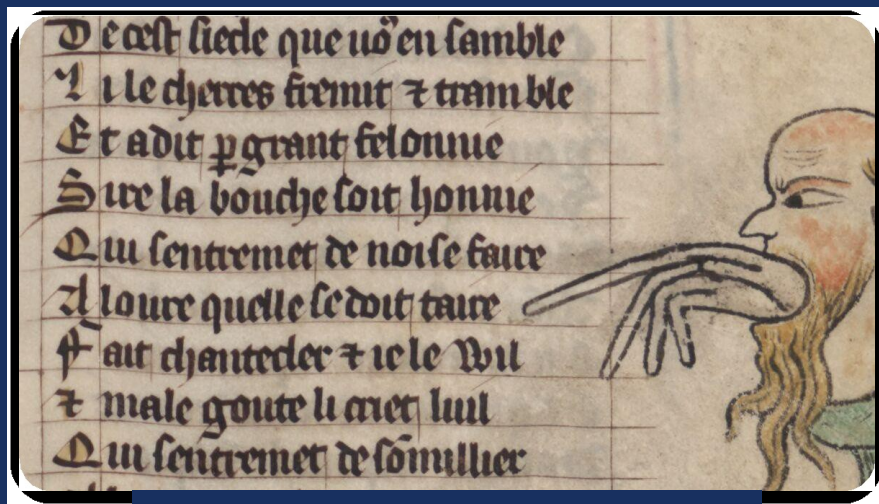
Manuálne spustenie segmentácie (rozpoznávania rozloženia) alebo označovanie rozloženia (oblasti textu + základné čiary) manuálne

- Prepis textu z **left-to-right** v textovom editore (zľava – doprava)
- V konfigurácii tréningu Rozšírené nastavenia vyberte **Reverse Text (RTL)** tak, aby bol výstupný text napísaný v smere sprava doľava

[Example](#) DiJeSt 2.0 model

**Vízia:**

- Podpora RTL pre webovú aplikáciu
- Prispôsobovanie konfigurácie tréningu

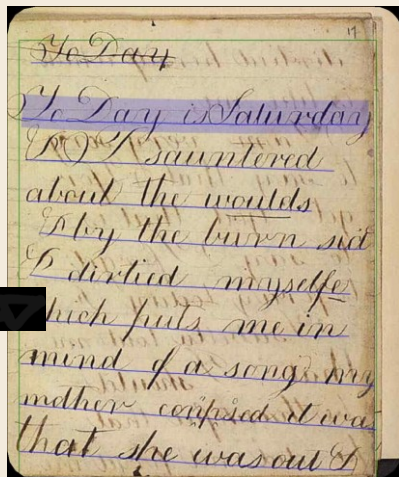


*Paris, BnF, Fr. MS 12584 (13th century)*

Rozpoznávanie rozloženia (Segmentácia)

# Čo sa stane, keď sa stránka rozpozná?

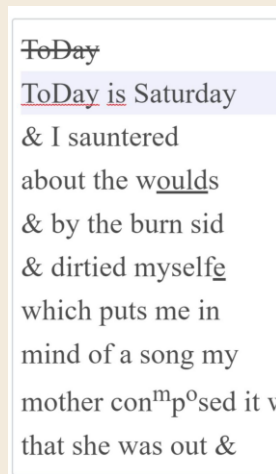
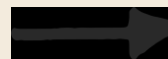
 Recognize



1. krok

**Rozpoznávanie rozloženia**

(Základné čiary (Baselines) & Bloky textu (Text regions))



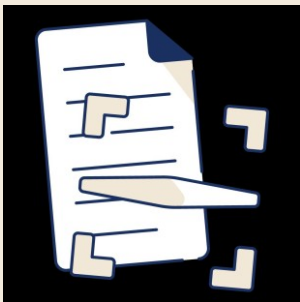
2. krok

**Rozpoznávanie textu**

# Rozpoznávanie rozloženia (segmentácia)

## 1. krok

### Rozpoznávanie rozloženia



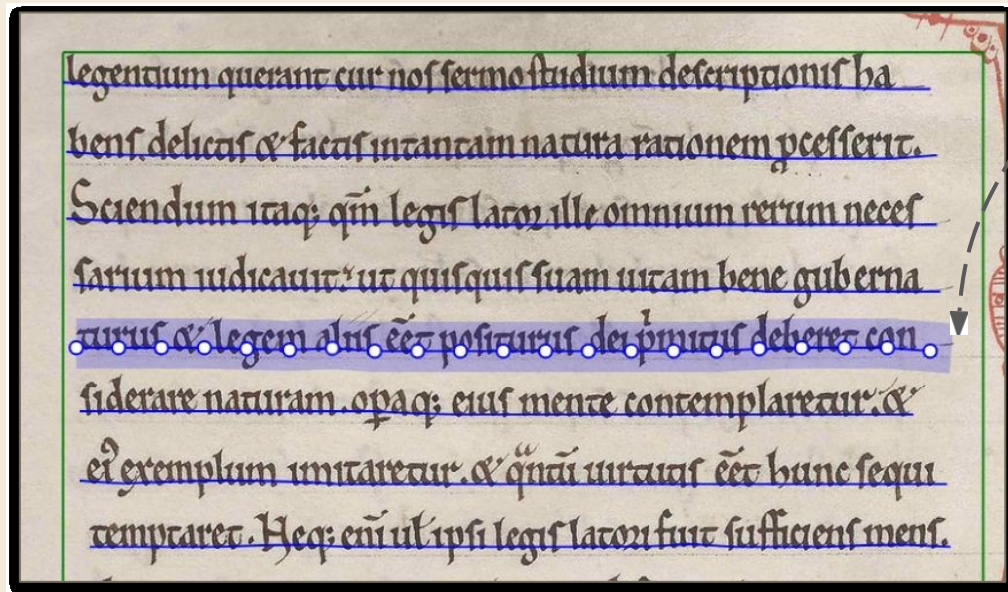
- Analýza rozloženia obrazu dokumentu
- Obrázok je potrebné rozdeliť na textové oblasti a základné čiary
- Základ pre rozpoznávanie a pre transkripciu (prepis)



# Tri piliere rozloženia (segmentácie)

## 1) **Základná čiara** (Baseline):

Členená čiara prebiehajúca pozdĺž spodnej časti riadka rukou písaného textu



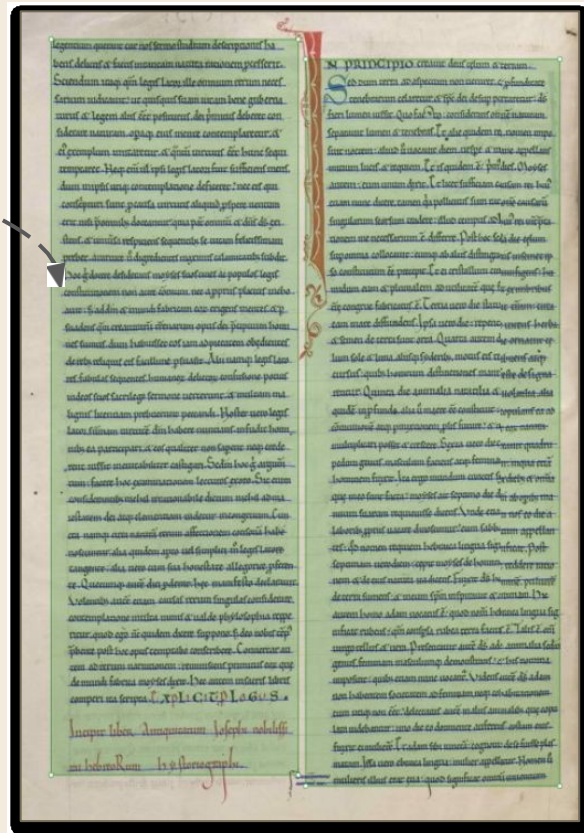
# Tri piliere rozloženia (segmentácie)

1) Základné čiary (Baselines)

2) **Bloky textu (Text region):**  
obdĺžnikový tvar obklopujúci text

Pri predvolenej analýze rozloženia sú základné čiary zoskupené do blokov textu (textových oblastí na základe ich súradníc (prístup zdola nahor))

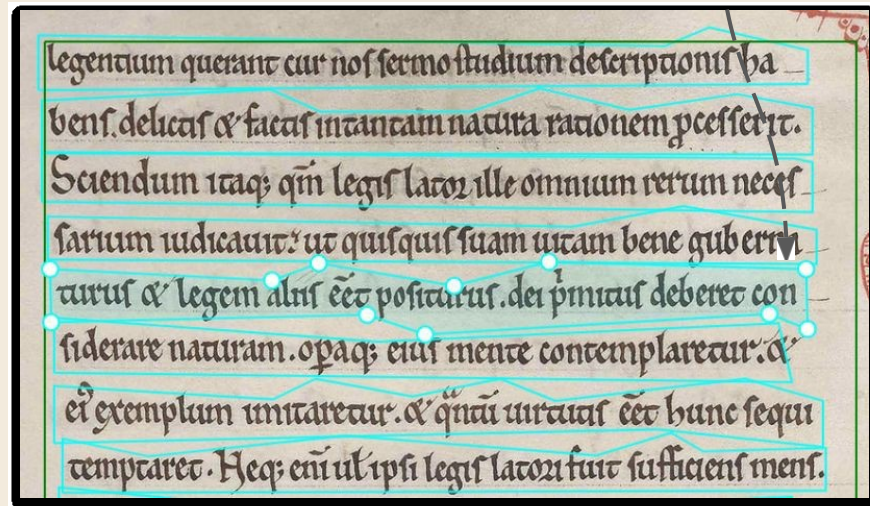
Bloky textu (Text region)



# Tri piliere rozloženia (segmentácie)

- 1) Baseline
- 2) Text region
- 3) Polygóny riadku (Line Polygons:** mnohoúhelníky, obklopujúce všetok rukou písaný text v riadku

Pri spustení tréningu textu alebo rozpoznávania textu sa mnohoúhelníky čiar vypočítajú algoritmom, počnúc **základnými čiarami**



Tréning a rozpoznávanie textu prebiehajú na úrovni **základných čiar !!!**

# Rozpoznávanie rozloženia (segmentácia)

Kvalitu konečného rozpoznania (segmentácie) môže ovplyvniť:

## 1) Nepresné základné čiary (baselines):

- Zistí sa príliš málo základných čiar (východiskových hodnôt) alebo príliš veľa základných čiar (východiskových hodnôt)

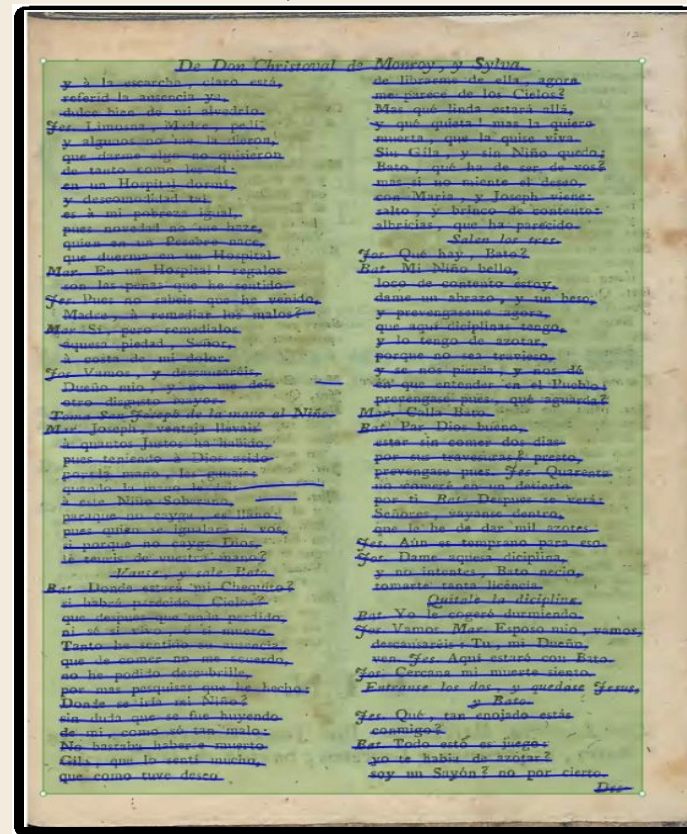


# Rozpoznávanie rozloženia (segmentácia)

Kvalitu konečného rozpoznania (segmentácie) môže ovplyvniť:

## 2) Nepresné bloky textu:

- Nesprávne poradie čítania riadkov;
- Príliš málo blokov textu/príliš veľa blokov textu (text regións)



# Rozpoznávanie rozloženia (segmentácia)

Kvalitu konečného rozpoznania (segmentácie) môže ovplyvniť:

### 3) Nepresné polygóny (Inaccurate polygons):

- Aj keď sú základné čiary správne, modely nedokážu správne prepísať text.
- Riadkové mnohoúhelníky nepokrývajú väčšinu tela písmen/  
Polygóny čiar zahŕňajú aj ďalšie (neželané) prvky na strane



Sept 1 / the 1836  
~~of the ...~~  
I have been thinking of  
writing to you for some  
time but have not had  
time to do so. I am  
well & hope you are  
the same. I have not  
heard from you for some  
time & I am sorry to  
hear that you are not  
well. I hope you will  
soon be better. I have  
not much news to write  
at present. I am  
Dear Sir,  
I have the honor to  
acknowledge the receipt  
of your letter of the  
10th inst. & am glad to  
hear that you are  
well. I have not much  
news to write at present.  
I am Dear Sir,  
I have the honor to  
acknowledge the receipt  
of your letter of the  
10th inst. & am glad to  
hear that you are  
well. I have not much  
news to write at present.  
I am Dear Sir,  
I have the honor to  
acknowledge the receipt  
of your letter of the  
10th inst. & am glad to  
hear that you are  
well. I have not much  
news to write at present.  
I am

## Nepresné základné čiary

# Nepresné základné čiary

Example





# Nepresné základné čiary – čo robiť?

Riešenia:

- 1) Použitie iného verejného modelu základných čiar (Baseline model)**
- 2) Zmeňte pokročilé nastavenia (advanced settings)**
- 3) Vytrénujte model základnej čiary (Train a baseline model)**

# Nepresné základné čiary

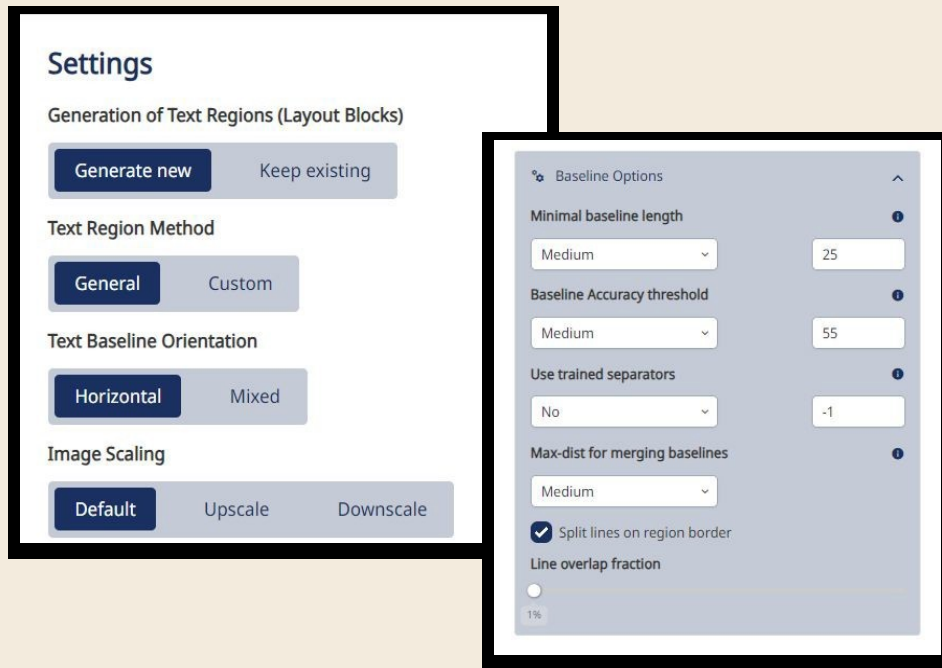
## 1) Použitie iného verejného modelu základných čiar (Baseline model)

:

- Zmiešaná orientácia riadkov (Mixed Line Orientation)
- Horizontálna orientácia riadkov (Horizontal Line Orientation)
- Univerzálne riadky (Universal Lines)

# Nepresné základné čiary

## 2) Zmeňte pokročilé nastavenia (advanced settings)



The image shows two overlapping screenshots of a settings interface. The left screenshot shows the main 'Settings' panel with options for text region generation and scaling. The right screenshot is a zoomed-in view of the 'Baseline Options' section, which includes several adjustable parameters.

**Settings**

Generation of Text Regions (Layout Blocks)

**Generate new** Keep existing

Text Region Method

**General** Custom

Text Baseline Orientation

**Horizontal** Mixed

Image Scaling

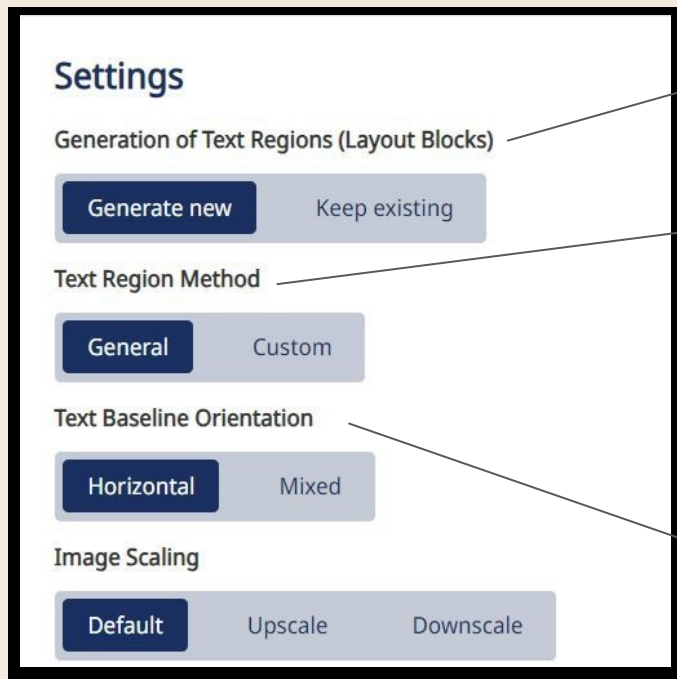
**Default** Upscale Downscale

**Baseline Options**

- Minimal baseline length: Medium (25)
- Baseline Accuracy threshold: Medium (55)
- Use trained separators: No (-1)
- Max-dist for merging baselines: Medium
- Split lines on region border
- Line overlap fraction: 1%

# Nepresné základné čiary

## 2) Zmeňte pokročilé nastavenia (advanced settings)



**Generate new:** Generovať ďalšie textové oblasti /

**Keep existing:** Zachovať existujúce oblasti textu (použite to s poľami a tabuľkami)

Po zistení sú riadky zoskupené do textových oblastí. K dispozícii sú dve metódy zoskupovania:

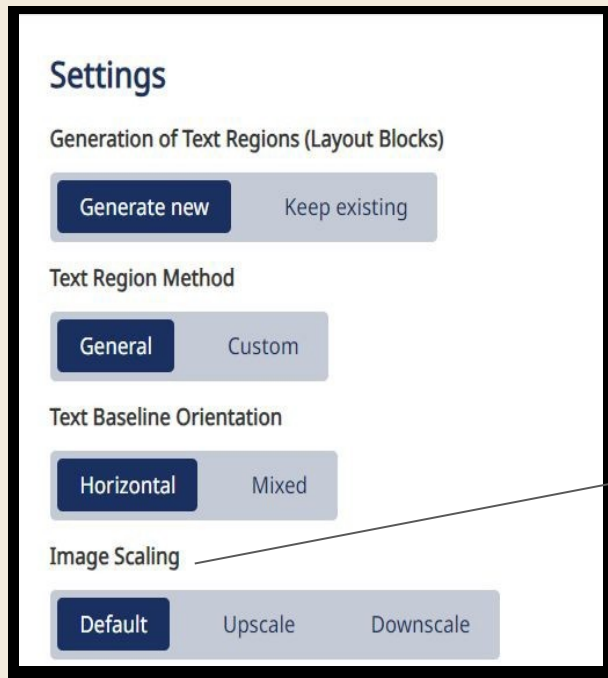
**General** (Všeobecné): zoskupí čiary zľava doprava

**Custom** (Vlastné): aglomeračné zoskupovanie založené na bode úplne vľavo každej čiary

Voľba **General**: Výber orientácie riadka textu na zlepšenie klastrovania (zoskupovania)

# Nepresné základné čiary

## 2) Zmeňte pokročilé nastavenia (advanced settings)

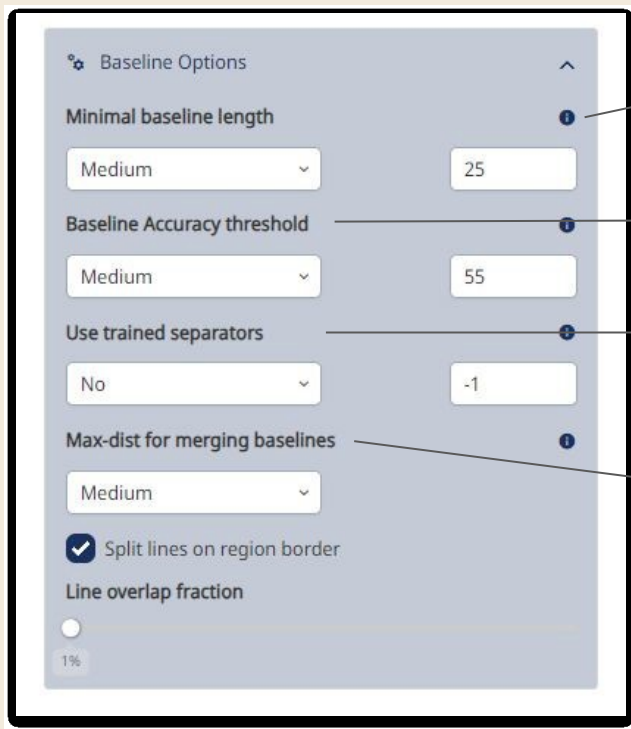


Škálovanie obrázka:

**Upscale** obrázky s nízkym rozlíšením alebo  
**Downscale** obrázky s vysokým rozlíšením  
(túto funkciu použite len v prípade, že  
rozpoznávanie rozloženia nezistí žiadne alebo  
len niekoľko riadkov)

# Nepresné základné čiary

## 2) Zmeňte pokročilé nastavenia (advanced settings)



### Minimálna dĺžka základnej čiary

(Minimal baseline length): Minimálna dĺžka riadkov **v pixeloch** (pre tabuľky je lepšie nastaviť ho na hodnotu Nízka)

### Prah presnosti základnej čiary (Baseline

Accuracy threshold): Stredné a nízke poskytujú lepšie výsledky

**Použitie tréovaných separátorov (Use trained separators)** Ak zvýšite túto hodnotu, okolité čiary sa zvyčajne zlučujú

### Max vzdialenosť pre spojenie základných čiar

(Distance for merging baselines):

**Low:** Zlúčia sa iba najbližšie čiary

**Medium**

**High:** vzdialené základné čiary sa zlúčia

# Nepresné základné čiary

## 2) Zmeňte pokročilé nastavenia (advanced settings)

Baseline Options

Minimal baseline length

Medium 25

Baseline Accuracy threshold

Medium 55

Use trained separators

No -1

Max-dist for merging baselines

Medium

Split lines on region border

Line overlap fraction

1%

**Rozdeliť čiary v rámci bloku** (Split lines on region border)

Iba ak zachováte existujúce bloky textu:

Delené čiary na hranici regiónu:

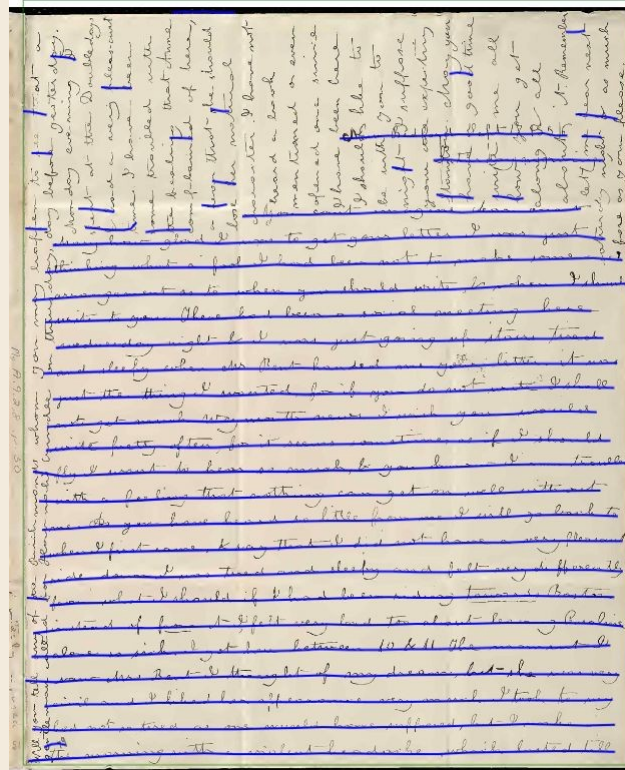
Aby čiary striktno dodržiavali hranicu regiónu.

Dôležité pre tabuľky!

# Nepresné základné čiary

Example 1

Example 2





# Nepresné základné čiary

3) Ak vám verejné modely a rozšírené nastavenia neposkytnú dobrý výsledok, tak:

**Trénujte Model pre základné čiary (Baselines model) vášho špecifického dokumentu**

Všetky stránky musia mať podobné rozloženie!

MURRAY, MARGARET D.		
10/20/08	Marks Received on Examination.	Jacket
10/20/08	Recommendation of Exam. Board.	"
10/28/08	Authority of Sec. to appoint.	B 14
10/28/08	Appointed. Reported 11/8/08	Jacket.
11/18/09	Req. Trans. to Mare Island, Cal.	E 10
2/3/10	Req. trans. to Wash. & Req. for Mare Island, Cal. withdrawn.	Jacket
3/21/13	Telegrams re- resignation. Miss Taylors	jacket
Bureau M. & S., Navy Department, Incc. 1 Jan. '11		

(2) MURRAY, MARGARET D.		
3/20/13	Tenders resignation.	Jacket.
5/10/13	Authority of Dept. to accept.	"
5/18/13	Resigned. (M.I.)	"
4/14/15	Miss Delano req. infor. (ans. 4/17/15. K 9	
3/24/14	3/R to Ruff	
2826 Calvert St., Baltimore, Md.		
4/15/34 - 2101 Sh. Paul St. Balti. Md.		

# Tréning modelu základných čiar (Baseline Model)

MURRAY, MARGARET D.		
10/20/08	Marks Received on Examination.	Jacket
10/20/08	Recommendation of Exam. Board.	"
10/28/08	Authority of Sec. to appoint.	B 14
10/28/08	Appointed. Reported 11/2/08	Jacket.
11/18/09	Req. Trans. to Mare Island, Cal.	E 10
2/3/10	Req. trans. to Wash. & Req. for Mare Island, Cal. withdrawn.	Jacket
3/21/13	Telegrams re- resignation.	Miss Taylors jacket
Bureau U. & S. Navy Department, 16,000. 1 Jan. '11		

(2) MURRAY, MARGARET D.		
3/20/13	Tenders resignation.	Jacket.
5/10/13	Authority of Dept. to accept.	"
5/16/13	Resigned. (M.I.)	"
4/14/15	Miss Delano req. infor. (ans. 4/17/15. K 9	
3/24/19	<i>3/A to R-FF</i>	
2826 Calvert St., Baltimore, Md.		
<i>1/15/34 - 2101 St. Paul St. Balti. Md.</i>		

MURRAY, MARGARET D.		
10/20/08	Marks Received on Examination.	Jacket
10/20/08	Recommendation of Exam. Board.	"
10/28/08	Authority of Sec. to appoint.	B 14
10/28/08	Appointed. Reported 11/2/08	Jacket.
11/18/09	Req. Trans. to Mare Island, Cal.	E 10
2/3/10	Req. trans. to Wash. & Req. for Mare Island, Cal. withdrawn.	Jacket
3/21/13	Telegrams re- resignation.	Miss Taylors jacket
Bureau U. & S. Navy Department, 16,000. 1 Jan. '11		

(2) MURRAY, MARGARET D.		
3/20/13	Tenders resignation.	Jacket.
5/10/13	Authority of Dept. to accept.	"
5/16/13	Resigned. (M.I.)	"
4/14/15	Miss Delano req. infor. (ans. 4/17/15. K 9	
3/24/19	<i>3/A to R-FF</i>	
2826 Calvert St., Baltimore, Md.		
<i>1/15/34 - 2101 St. Paul St. Balti. Md.</i>		

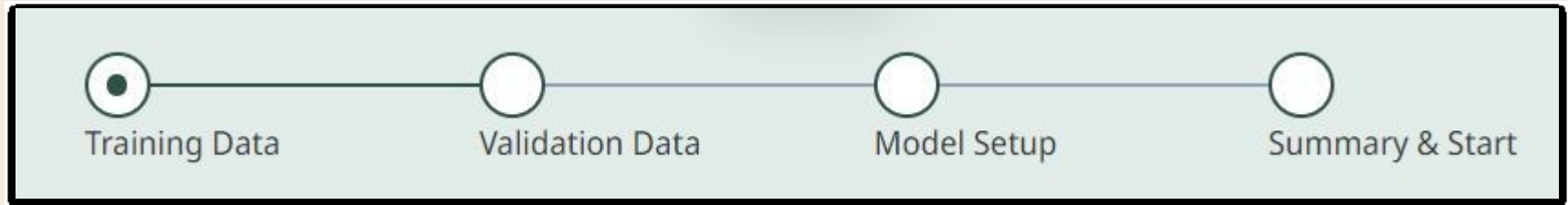
# Tréning modelu základných čiar (Baseline Model)

**Pripravte si aspoň 50 strán GT so správnymi základnými čiarami:**

- Nakreslite všetky základné čiary manuálne alebo opravte automatické rozpoznávanie rozloženia
- Nakreslite základné čiary iba pre časti, ktoré chcete prepísať

# Tréning modelu základných čiar (Baseline Model)

- Vyberte tréningové údaje (Training Data)
- Vyberte overovacie údaje (Validation Data)
- Nastavenie modelu (Model setup)
- Rozšírené nastavenia



# Modely pre základné čiary (Baselines Models)

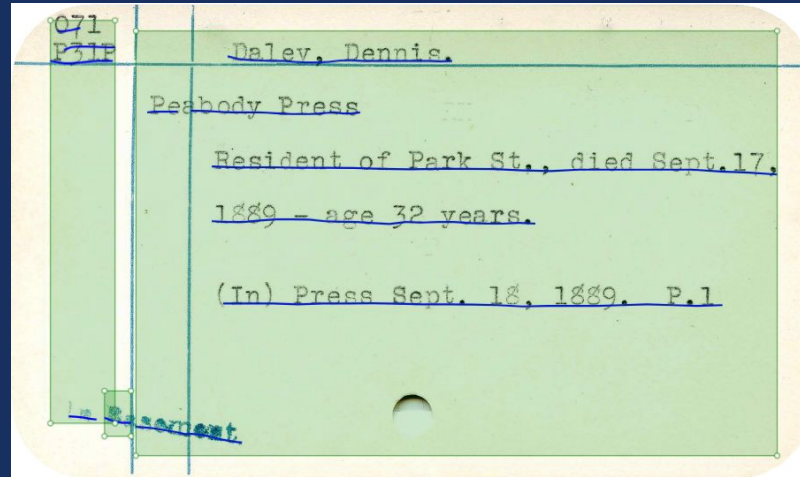
Po zaškolení môžete použiť svoj prispôsobený *Model pre základné čiary* (Baselines model) pre váš dokument! Zobrazí sa v zozname vašich súkromných Modelov rozloženia (Layout Models)

The screenshot displays the Text Recognition interface. At the top, there is a 'Text Recognition' header with a 'Layout' tab selected. Below this, a search bar contains 'NL-RISA\_199\_226'. To the right, there are buttons for 'Start Recognition' and 'Advanced Settings', along with credit information: 'Credits needed: 0.00' and 'Available: 0.00'.

The main area shows a list of models under the 'Private Models' tab, which is highlighted with a red box. The list has columns for 'NAME' and 'WORDS'. The models listed are:

NAME	WORDS
[Redacted]	v3
[Redacted]	v2
[Redacted]	v1
Medieval manuscript with glosses	2 366
Notes and miscellaneous materiel	8 468

On the right side, a detailed view of a 'Private Model' is shown. The model name is '[Redacted]\_v2', created by 's.mansutti@readcoop.eu' on '19/12/2022'. The 'CER (Accuracy)' is highlighted with a red circle and shows a value of '5.19%'. Other details include 'Trained on: handwritten' and 'Languages'.

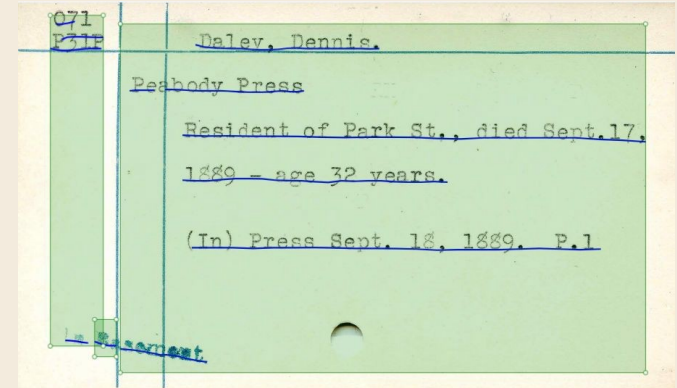
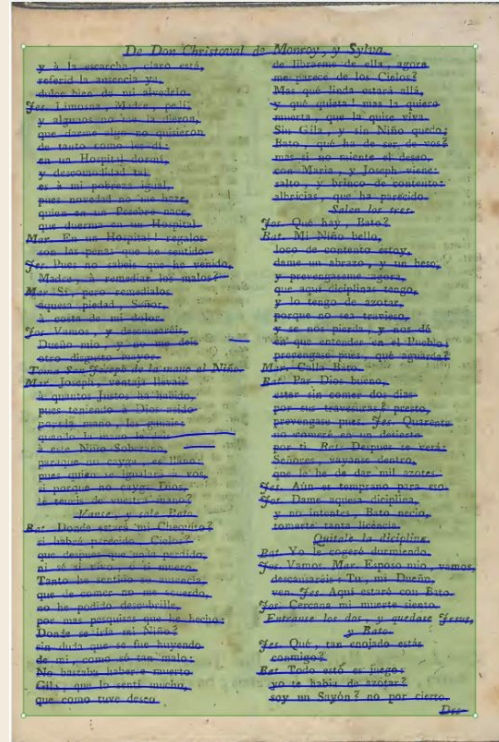


Nepresné bloky textu

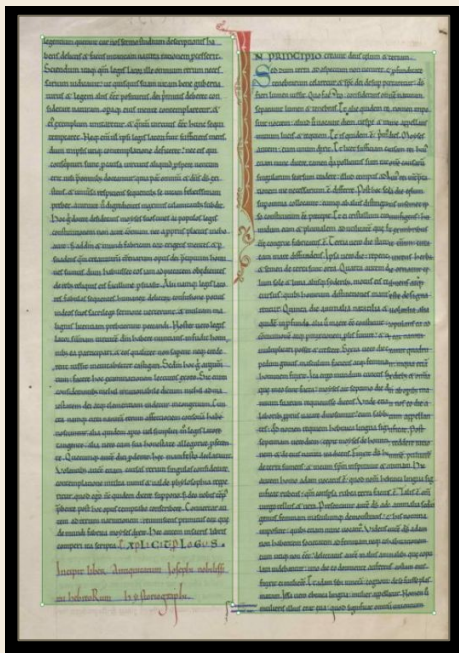
# Nepresné bloky textu

Example 1

Example 2



# Analýzy rozloženie/segmentácia (rozpoznanie textu)



## Bloky textu:

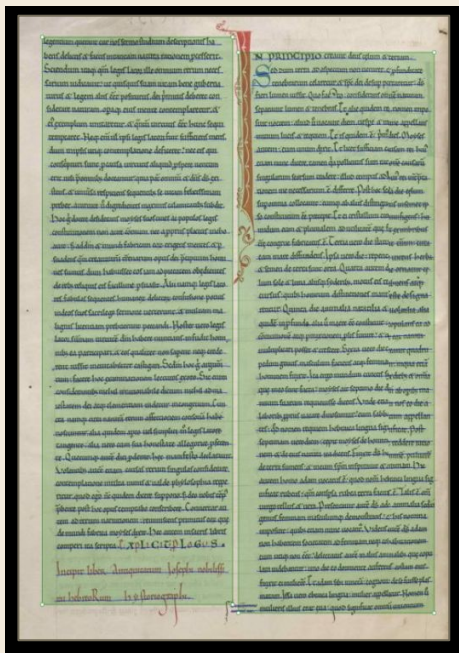
### Prístup zdola nahor

(s predvoleným rozpoznávaním textu a rozloženia):

1. Rozpoznanie základných čiar
2. Agregácia východiskových hodnôt v textových oblastiach na základe ich súradníc
3. Základné čiary a polygóny sa tvoria v momente rozpoznávania textu (Text Recognition)



# Analýzy rozloženie/segmentácia (rozpoznanie textu)



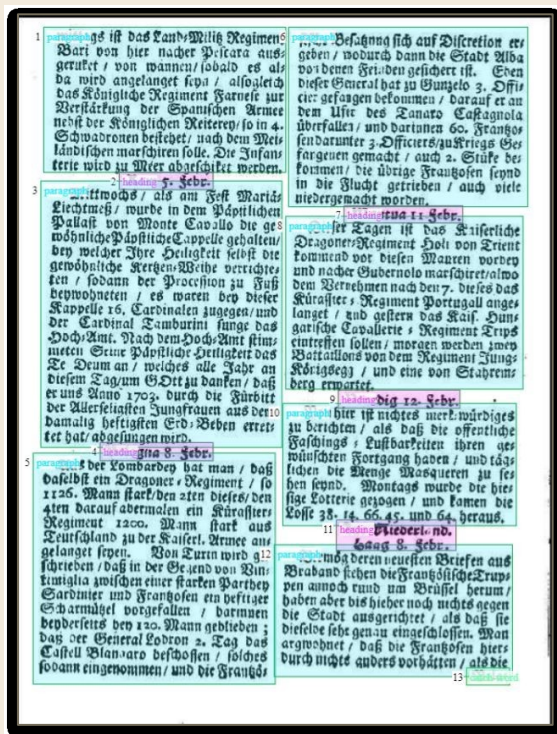
## Bloky textu:

## Prístup zdola nahor

V tomto prístupe môžete upraviť iba nastavenia:

1. Metóda oblasti textu
2. Orientácia základnej čiary textu

# Analýzy rozloženie/segmentácia (rozpoznanie textu)



## Bloky textu:

## Prístup zhora nadol

1. Rozpoznávanie blokov textu pomocou **Modelu pol'a (Field Model):** *polia sú v blokoch stránky*
2. Rozpoznanie základných čiar (Layout Recognition)
1. *Základné čiary a polygóny sa tvoria v momente rozpoznávania textu (Text Recognition)*

# Vormerkblatt

Name: **Name:** H u r t h, Oberstlttn.

Geburtsjahr u. Ort: **Year:** 1884 **Place:** au

Heimatzuständigkeitsort **Place:** au, Sudetenland

$\left. \begin{array}{l} \text{vor} \\ \text{nach} \end{array} \right\}$  dem Umsturz 1918:

Assentjahr: 1903

## Modely poľa

# Modely poľa (Beta)

Haupt-Grundbuchheft (Offentjahrgang)		1904...	Blatt-Nr.	625	
Vor- und Zuname		Johann Klüpfel <i>Klupfau</i>			
Geburts-	Ort	<i>Innsbrück</i>	Heimatsberechtigt in	Orts- gemeinde	<i>Innsbrück</i>
	Bezirk	<i>Innsbrück</i>		Bezirk	<i>Innsbrück</i>
	Comitat	<i>⁄</i>		Comitat	<i>⁄</i>
	Land	<i>Tirol</i>		Land	<i>Tirol</i>
		Geburts- jahr	18.83		
		Religion	<i>kathol.</i>		
		Kunst, Gewerbe, sonstiger Lebensberuf	<i>Lindler</i>		
April 1904 nach der Losreihe auf drei Jahre in der en Jahre in der Reserve und zwei Jahre in der Landwehr, zum 3. April 1. Tirol. Karl Lager					

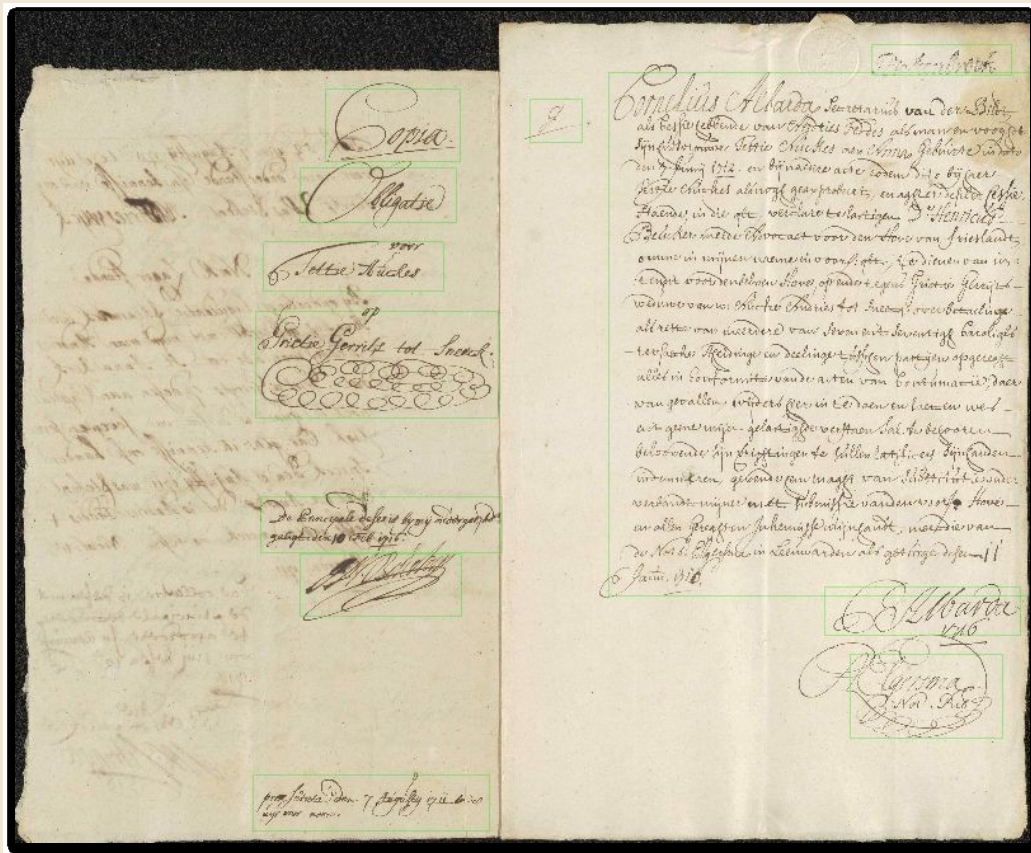
# Modely poľa (Beta)

Modely poľa je možné trénovať na:  
**automatické rozpoznávanie a  
označovanie určitých prvkov (dát)  
rozloženia dokumentu.**

- Bloky textu - Textové oblasti (polia)
- Priradenie značiek štruktúry pre tieto oblasti

Haupt-Grundbuchheft (Offentjahrgang)		1904...	Blatt-Nr.		625			
Vor- und Zuname		Name: Johann Hüpfauer Hüpfauer						
Geburts-	Ort	Innsbrück	Heimatsberechtigt in	Orts-gemeinde	Innsbrück	Geburts-jahr	Jahrgang	1883
	Bezirk	Innsbrück		Bezirk	Innsbrück		Religion	kathol.
	Comitat	✓		Comitat	✓	Kunst, Gewerbe, sonstiger Lebensberuf		Lieferant
	Land	Tirol		Land	Tirol			
Juli 1904 nach der Losreihe auf drei Jahre in der en Jahre in der Reserve und zwei Jahre in der Landwehr, zum 3. Aug. d. Tirol. Milit. 3. Quart.								

# Blogy textu (Text regions)



# Noviny: Segmentácia rozloženia



# Segmentácia formulára

*all*

Vater: vater separiert 100% 1951

Mutter: mutter separiert 100% name 100% 3 D

Staatsangehörigkeit: Staatsangehörigkeit 99%

Personalakt.: Geburtsname 100% Geburtsort 99% Geburtsort 99%

Familienname: name 100%

Vornamen: vorname 99%

Geburtsname: datum 100% Geburtsort: ort 99%

Glaubensbek.: Religion 100% Kreis: Beruf 100% rov.

Beruf: 1. Beruf 100% 3.  
4. 5. 6.

Mitglied und Seiltug  
i. d. NSDAP  
oder einer ihrer  
Gliederungen

Familienangehörige	Geburts- tag	mo- nat	jahr	Geburtsort (Kreis, Provinz) Standesamt	Glaubens- bek.	Aus- zugs- verm.	Mitglied und Seiltug i. d. NSDAP oder einer ihrer Gliederungen	Vermerke
<b>Kinder:</b>								

Verheiratet seit verheiratet vu0020seit 99% Standesamt Standesamt 97%

Standesamtsnummer 99% dem verheiratet mit 95%

geb. verheiratet mit geboren am 97%

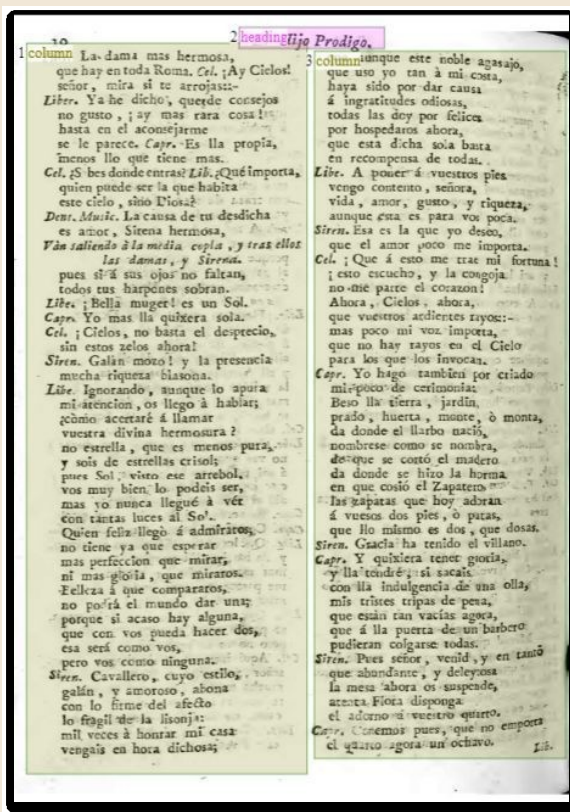
in verheiratet mit geboren in 100%

Wohnung: Wohnung 98%

Verdruck Nr. 304b (Wahlbild unvollst.) 9.48 380 000 In 150 Nr. 945 Staatsdruckerei Berlin 2706



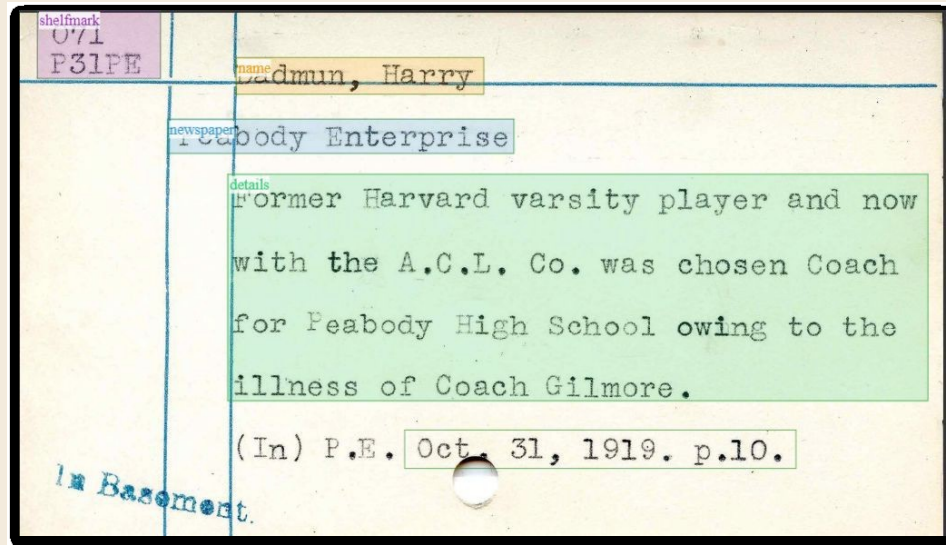
# Viac stĺpcové rozloženie textu



# Modely poľa (Beta)

## Pripravte si cca 50 strán tréningových dát:

- Nakreslite textovú oblasť okolo relevantných informácií, ktoré chcete extrahovať
- Prirad'ujte štrukturálne značky (voliteľné)



# Modely poľa (Beta)



Desk

Models

Sites

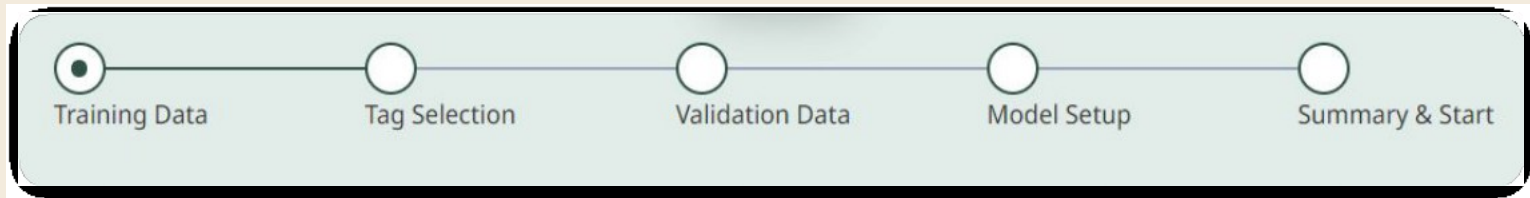
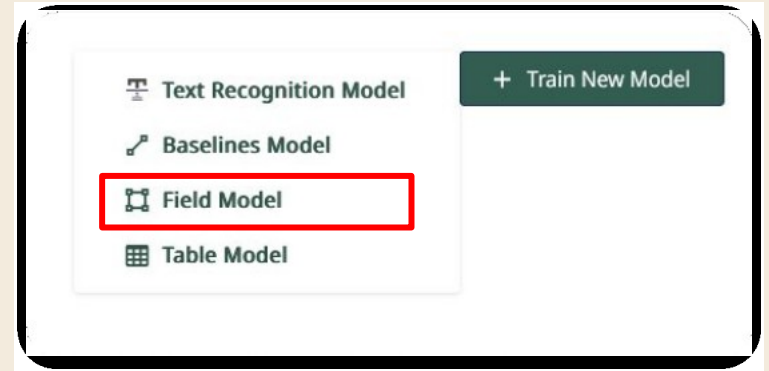
Jobs



**Modely (Models)** Transkribus je miesto, kde môžete trénovať a spravovať svoje modely.

# Modely poľa (Beta)

- Tréningové údaje (Training Data)
- Výber značky (tagov) (Tag Selection)
- Overovacie údaje (Validation Data)
- Nastavenie modelu (Model Setup)
- Rozšírené nastavenia (Cykly tréningu a miera učenia)



# Spracovanie dokumentov s poľami

- 1) **Vytvorenie Ground Truth pre rozpoznávanie polí:**
  - minimálne 50 strán
  - Viac strán so zložitým rozložením

# Spracovanie dokumentov s poľami

- 1) Vytvorenie Ground Truth pre rozpoznávanie polí
- 2) **Trénovanie modelu rozpoznávania polí**

# Spracovanie dokumentov s poľami

- 1) Vytvorenie Ground Truth pre rozpoznávanie polí
- 2) Trénovanie modelu rozpoznávania polí
- 3) **Použitie modelu rozpoznávania polí na zostávajúce strany**

# Spracovanie dokumentov s poľami

- 1) Vytvorenie Ground Truth pre rozpoznávanie polí
- 2) Trénovanie modelu rozpoznávania polí
- 3) Použitie modelu rozpoznávania polí na zostávajúce strany
- 4) **Spustenie rozpoznávania rozloženia na detekciu čiar:**

## Nastavenia:

- **Model základnej čiary (Baseline model):** Horizontal/Mixed Text Line Orientation/Model trained by you
- **Zachovanie existujúcich blokov - oblastí textu** (môže pomôcť) Minimálna dĺžka základnej čiary: (low) nízka
- **Rozdelené čiary na hranici regiónu**



# Spracovanie dokumentov s poľami

- 1) Vytvorenie Ground Truth pre rozpoznávanie polí
- 2) Trénovanie modelu rozpoznávania polí
- 3) Použitie modelu rozpoznávania polí na zostávajúce strany
- 4) Spustenie rozpoznávania rozloženia na detekciu čiar
- 5) Rozpoznávanie textu**
- 6) Verejný model / Privátny model, ktorý ste vyškoli, → možnosť aplikovať rôzne modely v rôznych oblastiach

# Spracovanie dokumentov s poľami

- 1) Vytvorenie Ground Truth pre rozpoznávanie polí
- 2) Trénovanie modelu rozpoznávania polí
- 3) Použitie modelu rozpoznávania polí na zostávajúce strany
- 4) Spustenie rozpoznávania rozloženia na detekciu čiar
- 5) Rozpoznávanie textu  
Verejný model / Vami vytrénovaný model
- 6) Korekcie (optional)
- 7) **Export**

# Spracovanie dokumentov s poľami

	A	B	C	D	E	F
1	TranskribusFilename	shelfmark	name	newspaper	details	reference
2	00729.jpg	071 D218	Dynan, Mary E.	Salem Evening News	Received diploma in October, 1910 from N.E. Institute of Anatomy, Sanitary Science and Embalming. Was first Peabody girl to graduate as an embalmer, also the youngest in the state.	Oct. 10, 1910. P.5
3	00730.jpg	071	Dynan, Mary E.	Salem Evening News	Of 17 Franklin St. was granted an Undertaker's license from the Board of Health. She passed a successful examination in embalming before the State Board and was the first woman in town to be granted such a license.	June 10, 1911. P.5
4	00731.jpg		Dynan, Timothy I.	Salem Evening News	Who died at his home, 30 Chestnut St. was a baker by trade and an active member of organized labor. He was employed by Jackson and Tortat until he met with an accident 5 years ago.	July 22, 1920. P.7
5	00732.jpg		Dynamite	Salem Evening News	Left unguarded, boys wander into magazine of the Essex Trap rock where enough is stored to blow the city to pieces. They take two sticks with them.	June 21, 1918. P.2

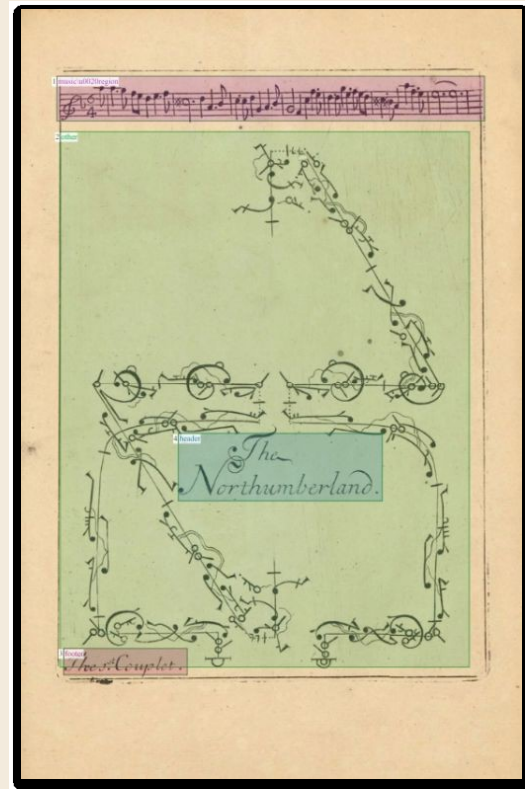
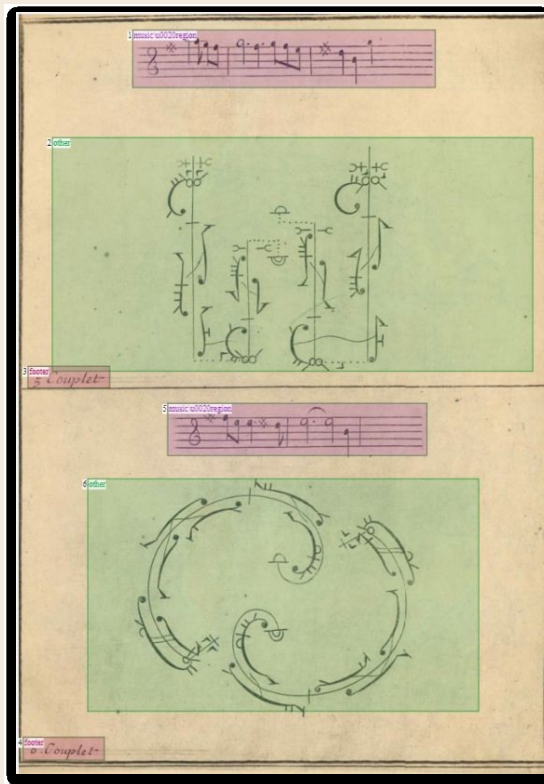
# Príklady Modelov polí

**Ground Truth:**

30 strán

5 tagov

[Example](#)



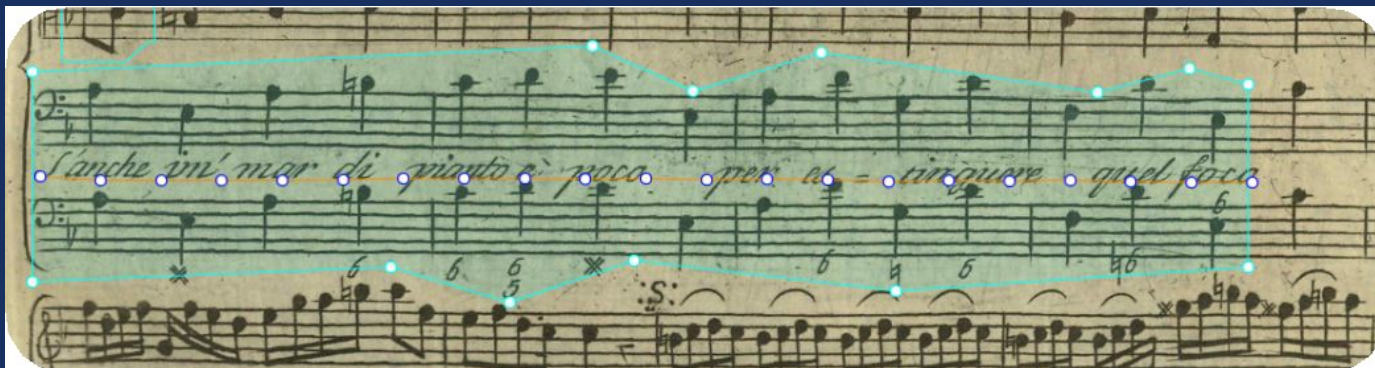
# Príklady Modelov polí

Example

1 **Header**  
**Anno 1746.** (Num. 17.) 26. Februarii.  
**Wienerisches**  
**DIARIUM.**  
 10 **Imprimt** Ihrer K. k. Kaiserl. auch zu Singarn, und Wáheim K. k. Maj. Streybeten  
**In dem neuen Michaeler-Haus/ bey Joh. Peter v. Spelen.**

4 **Imperator-singl**  
 5 **beobacht** **Italien.**  
**Genna 29. Jan.**  
 6 **ergangen**  
 7 **ergangen**  
 8 **ergangen**  
 9 **ergangen**  
 10 **ergangen**  
 11 **ergangen**  
 12 **ergangen**  
 13 **ergangen**  
 14 **ergangen**  
 15 **ergangen**  
 16 **ergangen**  
 17 **ergangen**  
 18 **ergangen**  
 19 **ergangen**  
 20 **ergangen**  
 21 **ergangen**  
 22 **ergangen**  
 23 **ergangen**  
 24 **ergangen**  
 25 **ergangen**  
 26 **ergangen**  
 27 **ergangen**  
 28 **ergangen**  
 29 **ergangen**  
 30 **ergangen**

1 **ergangen**  
 2 **ergangen**  
 3 **ergangen**  
 4 **ergangen**  
 5 **ergangen**  
 6 **ergangen**  
 7 **ergangen**  
 8 **ergangen**  
 9 **ergangen**  
 10 **ergangen**  
 11 **ergangen**  
 12 **ergangen**  
 13 **ergangen**  
 14 **ergangen**  
 15 **ergangen**  
 16 **ergangen**  
 17 **ergangen**  
 18 **ergangen**  
 19 **ergangen**  
 20 **ergangen**  
 21 **ergangen**  
 22 **ergangen**  
 23 **ergangen**  
 24 **ergangen**  
 25 **ergangen**  
 26 **ergangen**  
 27 **ergangen**  
 28 **ergangen**  
 29 **ergangen**  
 30 **ergangen**

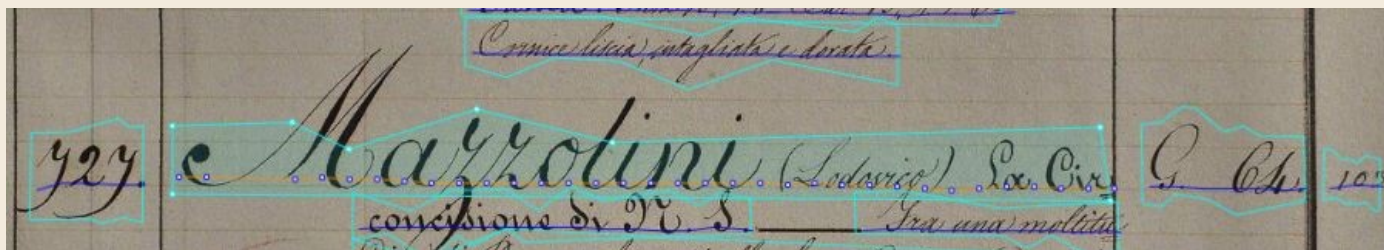


Nepresné polygóny (mnohouholníky)

# Inaccurate Polygons

[Example 1](#)

[Example 2](#)



# Field Model trained on Line Polygons

Prepare about 50 pages of training data:

- Adjust the line polygons manually





# Field Model trained on Line Polygons

- Training data
- Tag selection: TRAIN ON LINE POLYGONS
- Validation data
- Model setting
- Advanced settings (Training Cycles and Learning Rate)

Field Recognition Model

Training Data Tag Selection Validation Data Model Setup Start

Remove	Title	Example polygons
X		Example polygons

< Back

Next >

1 documents selected

Recognise untagged regions  
Select if you want include untagged regions in your training.

Train on line polygons  
Instead of training on tags, your model will be trained on line polygons.

# Field Model trained on Line Polygons



	Region 1
1	44
2	-
3	1
4	-
5	error
6	-
7	2
8	---
	Region 2
1	27



	Region 1
1	44
	Region 2
1	21
	Region 3
1	non sperare di smorzare col tuo pianto l'ira mi-
	Region 4
1	-a
	Region 5
1	s'anche in' mar di pianto poco per es-tinguere quel foco
	Region 6
1	ch'arde gel di gelo-si-a per estinguere quel fo-co
	Region 7
1	ch'arde al gel di gelo-si-a
	Region 8
1	Da Capo

APPLICATIONS FOR EDUCATIONAL INSTITUTIONS 1

NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION
Acadia University (1)	F Wolfville, N.S.		General	1/9/11	Has done share, 4/23/12 (ans. 2)
Americus Institute	Americus, Ga.	10000	Building	1/17/11	D + Low
Albemarle Normal & Ind. Inst.	Albemarle, N. C.		General	2/4/11	Low
Asbury College	Wilmore, Ky.		Buildings & Industrial Plant	2/13/11	D
Adrian College	P Adrian, Mich.		Endowment	3/13/11	Denominational, 3/6/11
Alabama State Normal School	Florence, Ala.	25000	Building	3/10/11	State institution, 2/15/11
Antioch College	Yellow Springs, O.	100000	Endowment	3/23/11	Not sufficiently developed, 2/17/11
American Church, Inst. for Negroes	New York, N. Y.		General	4/10/11	D
Amherst College	P Amherst, Mass.	50000	Increase Salaries	F 5/10/11	Has done share, 5/19/11
Alma College	P Alma, Mich.		Library Building	5/18/11	Denominational, 3/24/11
American International College	Springfield, Mass.		General	1/10/11	Not sufficiently developed, 2/10/11
Alberta Ladies' College (1)	Red Deer, Alta.		General	12/15/11	Not sufficiently developed, 1/14/11
Allen University	Columbia, S.C.		Library Building	4/26/12	Low
Acadia University (2)	F Wolfville, N. S.	25000	Library Building	5/10/12	Has done share, 4/23/12; D
Abingdon Presbytery	Orange, Va.	500	Building	5/14/12	D + Low
Amity College	College Springs, Ia.		Endowment	5/13/12	Not sufficiently developed, 10/15/11
Albert Lea College	P Albert Lea, Miss.		Endowment	6/18/12	Denominational, 1/12/11
Anderson College (1)	Anderson, S. C.		General	8/10/12	D + Low
Alabama University of	University, Ala.		Memorial Building	10/10/12	State institution, 1/15/11

## Table Models

# Modely pre tabuľky (Beta)

97

NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION
Lutheran Ladies' Seminary	Reading, Minn.		Sanctory	1/10/10	Not sufficiently developed, 2/10/10
Lombard College	P. Yaleburg, Ill.	50 000	Science Building	2/1/10	D
Loyan Female College	Russellville, Ky.	15 000	Building and Equipment	2/1/10	Low
Livingston College (1)	P. Salisbury, N.C.		Sanctory	2/1/10	D
Linden Hall Seminary	Leitch, Pa.	50 000	Library and Science Building	2/1/10	Seminary, 12/12/10
Livingston College (2)	P. Salisbury, N.C.		Land	2/1/10	D
Laurinburg Norm. & Ind. Inst.	Laurinburg, N.C.		General	2/1/10	Normal
Lenoir College (1)	Hickory, N.C.	70 000	Science Bldg., Compt. Bldg., and Education	2/1/10	Not sufficiently developed, 2/10/10
La Grange College (1)	La Grange, Mo.		Endowment and Buildings	2/1/10	Not sufficiently developed, 2/10/10
Lexington College	Lexington, Mo.		Endowment	2/1/10	Not sufficiently developed, 2/10/10
Lafayette College (1)	P. Easton, Pa.		Engineering Bldg. & Endowment	2/1/10	As done share, 2/1/10
Lenox College	P. Hopkinton, Iowa	12 500	Dept. of Agriculture	1/1/10	As done share, 1/10/10
Lincoln Memorial University (1)	P. Cumberland Gap, Tenn.		Building	1/1/10	As done share, 1/10/10
Lincoln Institute	Jefferson, Ky., Mo.		Library Building	2/1/10	Low
Lutheran College (projected)	Seguin, Tex.		Building	2/1/10	D
Leeds Industrial School	Lewisburg, Pa.		Building	2/1/10	Low
La Grange College (2)	La Grange, Mo.		Library Building	2/1/10	Not sufficiently developed, 2/10/10
Lindenwood College for Women	P. St. Charles, Mo.	10 000	Building	2/1/10	D

97

NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION
Lutheran Ladies' Seminary	Reading, Minn.		Sanctory	1/10/10	Not sufficiently developed, 2/10/10
Lombard College	P. Yaleburg, Ill.	50 000	Science Building	2/1/10	D
Loyan Female College	Russellville, Ky.	15 000	Building and Equipment	2/1/10	Low
Livingston College (1)	P. Salisbury, N.C.		Sanctory	2/1/10	D
Linden Hall Seminary	Leitch, Pa.	50 000	Library and Science Building	2/1/10	Seminary, 12/12/10
Livingston College (2)	P. Salisbury, N.C.		Land	2/1/10	D
Laurinburg Norm. & Ind. Inst.	Laurinburg, N.C.		General	2/1/10	Normal
Lenoir College (1)	Hickory, N.C.	70 000	Science Bldg., Compt. Bldg., and Education	2/1/10	Not sufficiently developed, 2/10/10
La Grange College (1)	La Grange, Mo.		Endowment and Buildings	2/1/10	Not sufficiently developed, 2/10/10
Lexington College	Lexington, Mo.		Endowment	2/1/10	Not sufficiently developed, 2/10/10
Lafayette College (1)	P. Easton, Pa.		Engineering Bldg. & Endowment	2/1/10	As done share, 2/1/10
Lenox College	P. Hopkinton, Iowa	12 500	Dept. of Agriculture	1/1/10	As done share, 1/10/10
Lincoln Memorial University (1)	P. Cumberland Gap, Tenn.		Building	1/1/10	As done share, 1/10/10
Lincoln Institute	Jefferson, Ky., Mo.		Library Building	2/1/10	Low
Lutheran College (projected)	Seguin, Tex.		Building	2/1/10	D
Leeds Industrial School	Lewisburg, Pa.		Building	2/1/10	Low
La Grange College (2)	La Grange, Mo.		Library Building	2/1/10	Not sufficiently developed, 2/10/10
Lindenwood College for Women	P. St. Charles, Mo.	10 000	Building	2/1/10	D

Modely tabuliek automaticky rozpoznávajú riadky a stĺpce a tým zlepšujú extrakciu a analýzu tabuľkových údajov.

# Modely pre tabuľky(Beta)

- Modely sa učia rozpoznávať riadky, stĺpce alebo obe
- Zatiaľ žiadne všeobecné modely, ale školenia pre konkrétne zbierky/dokumenty
- Nie sú potrebné oddeľovače (separátory)
- S dostatkom tréningových údajov dokáže model spracovať viacero typov tabuliek









# Viacriadkové bunky

9

NUMM welke inde STUK Omtrent	DATUM der Registratie	NUMM der ontvangens STUKKEN van IEDER STUK.	AANTWOORD van het OBJEKT van IEDER STUK.	Afschrift of Particulier de de ontvangens Stukken heeft gehooren of aan wien men schrijft.	ZAKELIJKE INHOUD VAN IEDER STUK.	Aan wien de zaak is Gedemandeerd.	DISPOSITIE		Aanmerkingen.
							180000		
11	1800	111	111		De inhoud van de stukken is als volgt: ...	111	111		
12	1800	112	112		De inhoud van de stukken is als volgt: ...	112	112		
13	1800	113	113		De inhoud van de stukken is als volgt: ...	113	113		
14	1800	114	114		De inhoud van de stukken is als volgt: ...	114	114		
15	1800	115	115		De inhoud van de stukken is als volgt: ...	115	115		

# Modely pre tabuľky

Ground Truth tvorba v editore:

Tabuľka

- Stĺpce
- Riadky

APPLICATIONS FOR EDUCATIONAL INSTITUTIONS						1
NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION	
Acadia University (1)	Woolfville, N.S.		General	1/9/11	Has done share, 11/23/12 (ans 2)	
Americus Institute	Americus, Ga.	10000	Building	11/7/11	D + Low	
Albemarle Normal + Ind. Inst.	Albemarle, N.C.		General	2/1/11	Low	
Asbury College	Wilmore, Ky.		Buildings + Industrial Plant	2/13/11	D	
Adrian College	Adrian, Mich.		Endowment	3/13/11	Denominational, 3/1/11	
Alabama State Normal School	Montgomery, Ala.	25000	Building	3/10/11	State institution, 2/15/11	
Antioch College	Yellow Springs, O.	100000	Endowment	3/2/11	Not sufficiently developed, 2/15/11	
American Church, Inst. for Negroes	New York, N.Y.		General	4/1/11	D	
Amherst College	Amherst, Mass.	50000	Increase Salaries	F 5/1/11	Has done share, 5/14/11	
Alma College	Alma, Mich.		Library Building	5/18/11	Denominational, 3/24/11	
American International College	Springfield, Mass.		General	1/10/11	Not sufficiently developed, 3/2/109	
Alberta Ladies' College (1)	Red Deer, Alberta		General	12/15/11	Not sufficiently developed, 1/14/11	
Allen University	Columbia, S.C.		Library Building	4/26/12	Low	
Acadia University (2)	Woolfville, N.S.	25000	Library Building	5/15/12	Has done share, 4/23/12; D	
Kingdon Presbytery	Stafford, Va.	500	Building	5/10/12	D + Low	
Amity College	College Springs, Ia.		Endowment	5/13/12	Not sufficiently developed, 10/15/11	
Albert Lea College	Pell City, Miss.		Endowment	6/18/12	Denominational, 11/2/11	
Anderson College (1)	Anderson, S.C.		General	8/10/12	D + Low	
Alabama University of	University, Ala.		Memorial Building	10/20/12	State institution, 1/15/109	

# Modely pre tabuľky

Stránky GT:

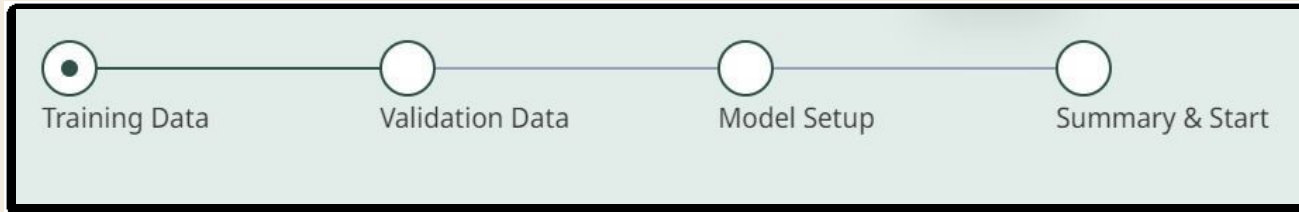
- **Jednoduché dabuľky:** 20 strán GT
- **Ťažké tabuľky:** 50 strán GT
- **mix rôznych tabuliek:** 50 až 100 strán GT v závislosti od počtu tabuliek

APPLICATIONS FOR EDUCATIONAL INSTITUTIONS						1
NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION	
Acadia University (1)	Woolfville, N.S.		General	1/9/11	Has done share, 11/23/12 (ans. 2)	
Americus Institute	Americus, Ga.	10000	Building	11/7/11	D + Low	
Albemarle Normal & Ind. Inst.	Albemarle, N.C.		General	2/1/11	Low	
Asbury College	Wilmore, Ky.		Buildings & Industrial Plant	2/13/11	D	
Adrian College	Adrian, Mich.		Endowment	3/13/11	Denominational, 3/1/11	
Alabama State Normal School	Montgomery, Ala.	25000	Building	3/10/11	State institution, 2/15/11	
Antioch College	Yellow Springs, O.	100000	Endowment	3/2/11	Not sufficiently developed, 2/12/11	
American Church, Inst. for Negroes	New York, N.Y.		General	4/1/11	D	
Amherst College	Amherst, Mass.	50000	Increase Salaries	F 5/1/11	Has done share, 5/14/11	
Alma College	Alma, Mich.		Library Building	5/18/11	Denominational, 3/24/11	
American International College	Springfield, Mass.		General	1/10/11	Not sufficiently developed, 3/2/109	
Alberta Ladies' College (1)	Red Deer, Alberta.		General	12/12/10	Not sufficiently developed, 1/1/109	
Allen University	Columbia, S.C.		Library Building	4/26/12	Low	
Acadia University (2)	Woolfville, N.S.	25000	Library Building	5/1/12	Has done share, 4/23/12; D	
Kingdon Presbytery	Stafford, Va.	500	Building	5/14/12	D + Low	
Amity College	College Springs, Ia.		Endowment	5/13/12	Not sufficiently developed, 10/15/10	
Albert Lea College	Pellott, La., Miss.		Endowment	4/18/12	Denominational, 11/2/11	
Anderson College (1)	Anderson, S.C.		General	8/15/12	D + Low	
Alabama University of	University, Ala.		Memorial Building	10/2/12	State institution, 1/15/109	

# Modely pre tabuľky

Tréning ([beta.transkribus.eu](https://beta.transkribus.eu)):

- Training data
- Validation data
- Model setting
- Advanced settings: Training Cycles and Learning Rate



# Modely pre tabuľky

Ground Truth: 20 strán

2 APPLICATIONS FOR EDUCATIONAL INSTITUTIONS					
NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION
Allegheny College (1)	P Meadville, Pa.		Chemistry Building	12/12/12	Has done share, 12/17/12
Allegheny College	Allegheny, Ore.		Endowment	12/16/12	Discontinuation, 12/17/12
Assiut College	Assiut, Egypt		Building	12/17/12	Outside field of work (Egypt), 12/17/12
Allegheny College (2)	P Meadville, Pa.		Library Building	12/21/12	Has done share, 12/17/12 (ans. 1)
Atlanta Normal & Ind. Inst. (1) C	Atlanta, Ga.		General and Land	12/21/12	Low; Gen 1914, 4/15/14
Alabama School of Trade & Ind.	Rayland, Ala.	25 000	Land	3/10/13	Planning stage, 3/18/13
American University (projected)	Washington, D.C.	140 000	Building	3/12/13	D
Adelphi College	Brooklyn, N. Y.		Endowment	5/2/13	Not disposed, 5/5/13
Urbington Lit. & Ind. School (1) C	Annapolis, Md.	7 000	Building	4/1/13	Low; Gen 1914, 2/18/14
Allegheny College (1)	Meadville, Pa.		Building	12/10/13	Discontinuation & not developed, 4/1/14
Austin College (1)	Sherman, Tex.		Library Endowment	1/10/14	Discontinuation, 1/15/14
Atlanta University	C P Atlanta, Ga.		Endowment	1/15/14	Gen 1914, 2/20/14
Allegheny County Academy (1)	Cumberland, Md.		Endowment	1/17/14	Academy, 1/5/14
Austin College (2)	Sherman, Tex.		Organ	1/17/14	No organs for institutions, 1/20/14

# Modely pre tabuľky

## Rozpoznávanie s tabuľkovými modelmi

### Processed pages

46					
APPLICATIONS FOR EDUCATIONAL INSTITUTIONS					
NAME OF INSTITUTION	TOWN	AMOUNT	OBJECT	DATE	DISPOSITION
Emporia, College of (2) Elgin Academy (2)	P Emporia, Kan. Elgin, Ill.		Organ Endowment	7/24/14 10/16/14	No organs for institutions, 10/6/14 Gen 1914, 10/30/14
"Ewing School"	Ewing, Va.		Building	6/3/15	
Emory and Henry College	P Emory, Va.	25 000	Endowment	12/1/15	
Elk Creek Training School	Elk Creek, Va.		Dormitory	12/10/15	
Elisee High School	Hemp, N.C.		Piano	1/3/16	
Emporia, College of (2)	P Emporia, Kan.		Rebuilding of Carnegie Library		"Gen 1914" 2/2/16
Elmira College (1) 1917 ↓	P Elmira, N.Y.	20 000 or 25 000	Library Building	1/4/16	"Gen 1914" 2/7/16
Ellsworth College (2)	P Iowa Falls, Iowa		Buildings and endowment	3/3/17	"Gen 1914," 4/6/17
Elmira College (2)	P Elmira, N.Y.	50 000	Buildings and endowment	3/10/17	"Gen 1914," 4/6/17
Edenton Ind. & Norm. College	Edenton, N.C.		General	10/8/17	
Emory and Henry College	Emory, Va.		Enlargement, and equipment	7/6/18	
Elizabethtown College	Elizabethtown Pa.	50 000	Library	1/28/19	Gen.
Esqfield Preparatory School	Esqfield, Pa.		Building	3/20/19	Gen. 3/26/19
Ellsworth College	Iowa Falls, Va.		Buildings and endowment	5/1/19	Gen 6/2/19

# Spracovanie dokumentov s tabuľkami

- 1) **Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek:**

# Spracovanie dokumentov s tabuľkami

- 1) **Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek**
- 2) **Trénovanie modelu rozpoznávania tabuliek**



# Spracovanie dokumentov s tabuľkami

- 1) **Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek**
- 2) **Trénovanie modelu rozpoznávania tabuliek**
- 3) **Použitie modelu rozpoznávania tabuľky na zostávajúce strany**

# Processing documents with tables

- 1) Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek
- 2) Trénovanie modelu rozpoznávania tabuliek
- 3) Použitie modelu rozpoznávania tabuľky na zostávajúce strany
- 4) Spustenie rozpoznávania rozloženia na detekciu riadkov:
- 5) **Nastavenia:**
- 6) **Model Základnej čiary (Baseline model):** Horizontal/Mixed Text Line Orientation/Model trained by you
  - Zachovanie existujúcich oblastí textu
  - Zmena mierky obrázka
  - Minimálna dĺžka základnej čiary:Low
  - Rozdelené čiary na hranici regiónu

# Processing documents with tables

- 1) **Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek**
- 2) **Trénovanie modelu rozpoznávania tabuliek**
- 3) **Použitie modelu rozpoznávania tabuľky na zostávajúce strany**
- 4) **Spustenie rozpoznávania rozloženia na detekciu riadkov:**
  - **Rozpoznávanie textu (Text Recognition)**  
Verejný model / Súkromný model, ktorý ste trénovali

# Processing documents with tables

1. **Vytvorenie GT „základnej pravdy“ pre rozpoznávanie tabuliek**
2. **Trénovanie modelu rozpoznávania tabuliek**
3. **Použitie modelu rozpoznávania tabuľky na zostávajúce strany**
4. **Spustenie rozpoznávania rozloženia na detekciu riadkov**
5. **Rozpoznávanie textu (Text Recognition)**
6. **Korekcie (Correction (voliteľné))**
7. **Export (Excel)**

# Modely polí a tabuliek: Súhrn



začnite s približne 40-60 stranami GT

50 strán pre Modely polí

- jednoduché tabuľky: 10/20 strán
- Zložité tabuľky: 30-50 strán
- Mix rôznych tabuliek: minimálne 50 strán

Príprava tréningových údajov pomocou editora rozloženia

- Oblasti kreslenia a tagovania pre modely polí (= priradiť tagy štruktúry)
- Kreslenie tabuliek pre tabuľkové modely



Pracovný postup pre prácu s tabuľkami a poľami:

1. rozpoznať oblasti alebo tabuľky
2. potom základné čiary
3. potom text



# Výpočty presnosti transkripcie

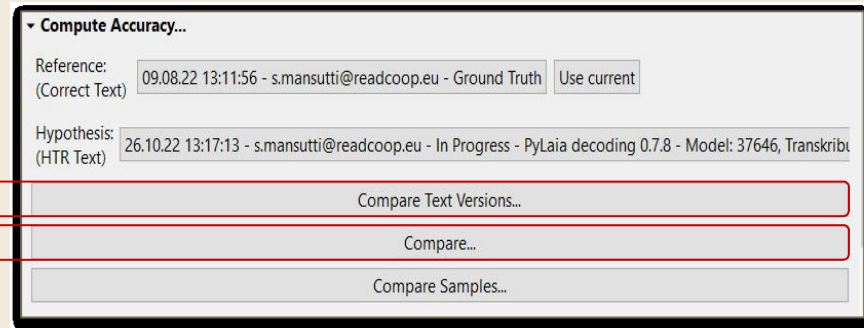


# Výpočty presnosti transkripcie

Dve verzie tej istej stránky:

1. **Reference** (Ground Truth)
2. **Hypothesis** (HTR Automatic Transcription)

- **Porovnajte textové verzie (pozrite si rozdiely medzi dvoma vybratými verziami)**
- **Porovnať...(Compare)**
- **(porovnáva tieto dva prepisy a vypočítava chybovosť slov a chybovosť znakov)**



# Výpočty presnosti transkripcie

## Porovnať textové verzie

Ground Truth - model "Transkribus  
English handwriting M3b" bez jazykového  
modelu:



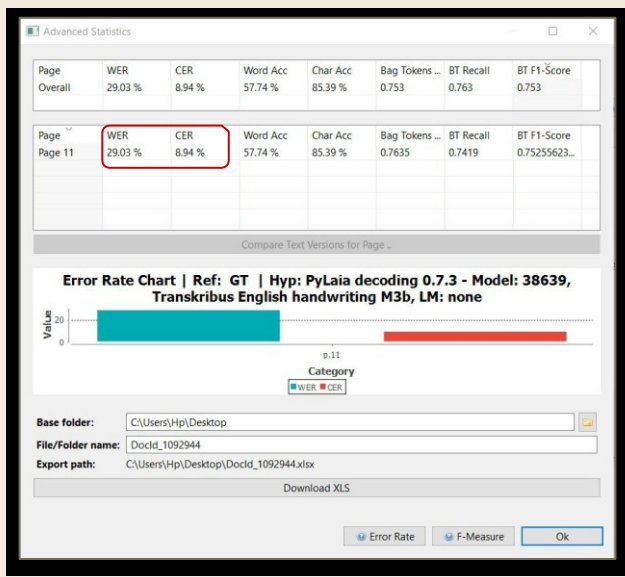
Ground Truth - "Transkribus anglický  
rukopis M3b" model s jazykovým  
modelom:



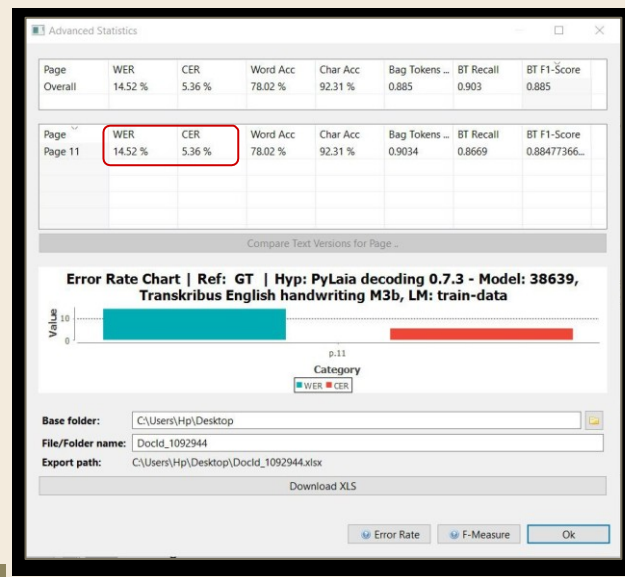


# Výpočty presnosti transkripce

**Porovnat'**...Ground Truth - Model  
"Transkribus English handwriting M3b"  
bez jazykového modelu:



Ground Truth - "Transkribus anglický rukopis M3b" model s jazykovým modelom:



# Výpočty presnosti transkripcie

Compare → Advanced Compare → Baselines

Compare: Advanced Compare

Type: Baselines

Pages (2): 1

Options: default (case sensitive)

Reference: GT Select hypothesis by toolname: TrHtr recognition 2.3.0 - Model: 51170, The Text Titan I

Compare

Previous Advanced Compare Results

Created	Status	Queries	Duration	Scope	Type	Results
26.09.23 10:35:06	Completed	Page(s) : 1   Ref: GT   Hyp : TrHtr recognition 2.3.0 - Model: 51...	0.52 sec.	Document ...	Baselin...	P/R/F1: 0.74/0.98/0.84 (p1: 0.74/0.98/0.84)
26.09.23 10:34:46	Completed	Page(s) : 1   Ref: GT   Hyp : Transkribus LA 0.0.5, Model: 49272...	0.60 sec.	Document ...	Baselin...	P/R/F1: 0.87/0.99/0.93 (p1: 0.87/0.99/0.93)
26.09.23 10:34:36	Completed	Page(s) : 1   Ref: GT   Hyp : Transkribus LA 0.0.5, Model: 51962...	0.59 sec.	Document ...	Baselin...	P/R/F1: 0.93/0.97/0.95 (p1: 0.93/0.97/0.95)

Options Cancel

Predvolené  
rozloženie s  
rozpoznávaním  
textu

Základný model  
orientácie  
zmiešanej čiary

Základný model  
univerzálnych  
línii

# Kontrola kvality

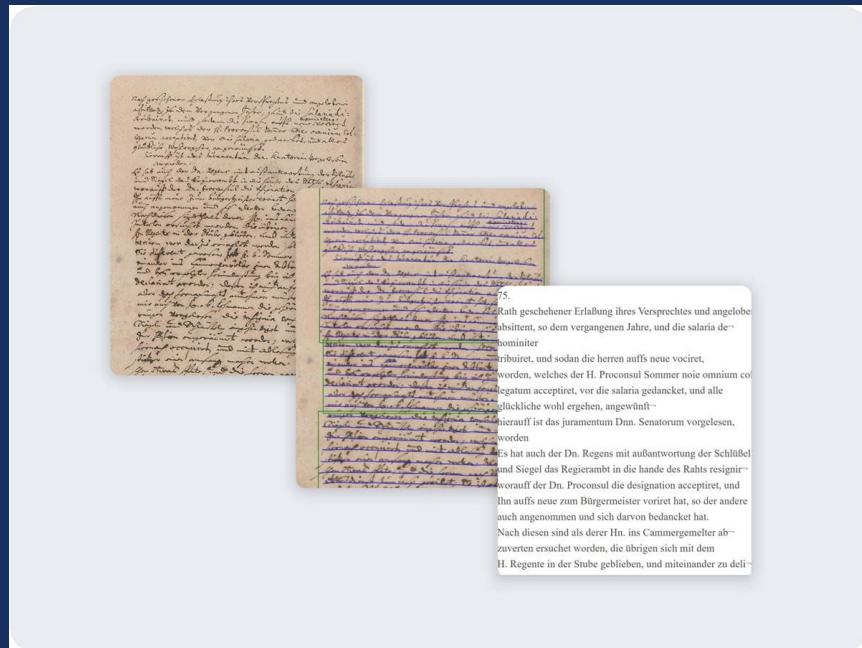
The screenshot displays the Transkribus Quality Control interface. At the top, the Transkribus logo is on the left, and navigation links for Desk, Models, Sites, Connect, and Jobs are on the right. Below the navigation bar, the breadcrumb path is "Quality Control > Sample 01 > Task 01".

Five evaluation metrics are shown in a row:

- Size: 100 Pages
- Layout Evaluation: 100%
- Transcript Evaluation: 98%
- Tags Evaluation: 97%
- Attributes Evaluation: 98%

Below the metrics is a table listing errors:

PAGE	STATUS	ERRORS	
Page #33	Error	Transcript	<a href="#">See Page</a>
Page #78	Error	Transcript, Tags	<a href="#">See Page</a>
Page #84	Error	Tags	<a href="#">See Page</a>
Page #98	Error	Tags	<a href="#">See Page</a>



# Publikačné modely v Transkribus

# Publikačné modely



Používatelia sa rozhodnú publikovať svoje vlastné modely, pretože

Sú hrdí na svoju prácu, a preto ju chcú sprístupniť aj ostatným používateľom, ktorí pracujú s podobnými skriptami a jazykmi

Musia publikovať čo najviac

Majú záujem o spoluprácu s inými vedcami na súvisiacich projektoch

Môžu vedieť o iných kolegoch alebo výskumných projektoch, ktoré by chceli použiť model, ale nemôžu zdieľať tréningové údaje

[Zenodo](#) Komunita pre publikovanie súborov údajov GT  
plánuje zahrnúť priame rozšírenie od spoločnosti Transkribus

# Publikačné modely

## Ako publikovať model:

Kontaktujte nás prostredníctvom [info@readcoop.eu](mailto:info@readcoop.eu) alebo prostredníctvom [contact form/help center](#) aby ste nás informovali, že chcete zverejniť svoj model v rámci spoločnosti Transkribus

- Požiadavky: veľkosť tréningovej sady ~ 50 000 slov, CER 7%-5% alebo nižšia . Ak ide o model vyškolený na skript alebo jazyk, ktorý zatiaľ nemôžeme ponúknuť, tieto kritériá neplatia
- Poskytnúť stručný opis modelu, ktorý pomôže ostatným používateľom pochopiť použitý obsah školenia; Užitočné je aj prídanie reprezentatívneho obrázka alebo úryvku
- Povedzte nám, kto by mal byť uvedený ako tvorca modelu - môže to byť jedna alebo viac osôb alebo celý výskumný projekt
- Viditeľnosť tréningových údajov: môžu byť zachované v súkromí (z dôvodov ochrany údajov) alebo zdieľané, aby boli aj údaje o školeniach verejné



 Desk

 Models

 Sites

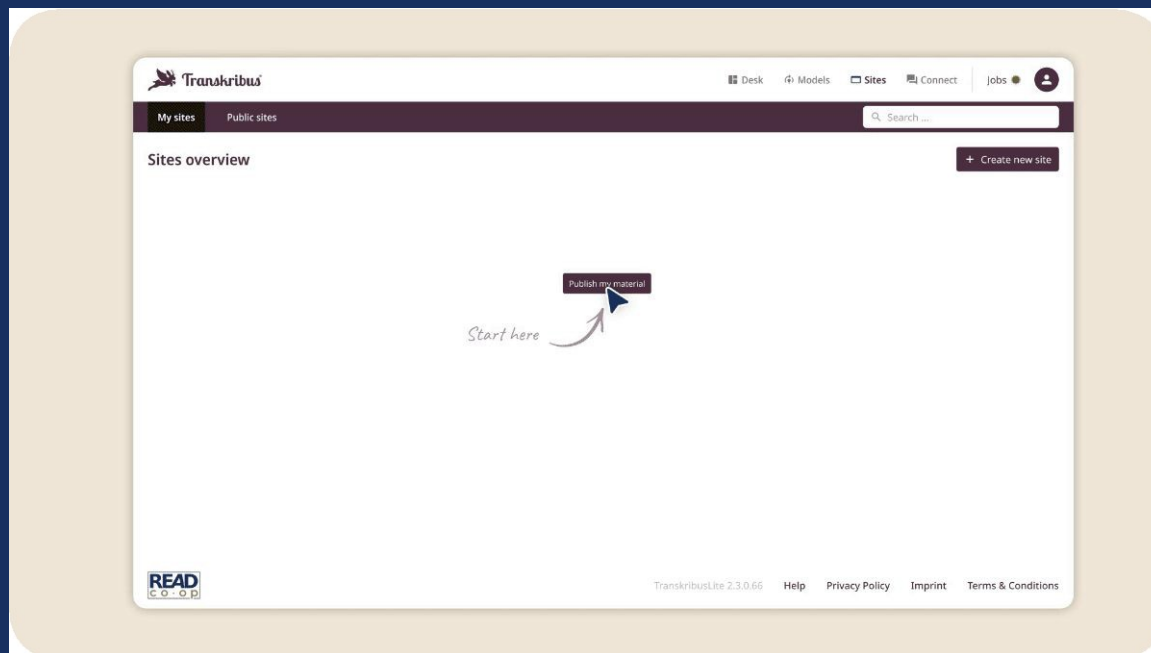
 **Connect**

Jobs 



Plánované na rok 2024

Transkribus **Connect** je miesto, kde sa **exchange** stane.



# Transkribus stránky





 Desk

 Models

 Sites

Jobs



Transkribus **Connect** je miesto, kde sa **exchange** stane.

# Plány predplatného

## Individual

0 €

Ideal for Genealogists & Students  
/month incl. 20% VAT\*

- ✓ AI Text Recognition
- ✓ Custom AI Training
- ✓ DOCX & PDF Export

Choose plan

## Scholar

14.9 €

Tailored for Individual Researchers  
/month incl. 20% VAT\*

- ✓ Collaboration Tools
- ✓ Advanced AI Tools
- ✓ Transkribus Sites

Choose plan

## Organisation

—

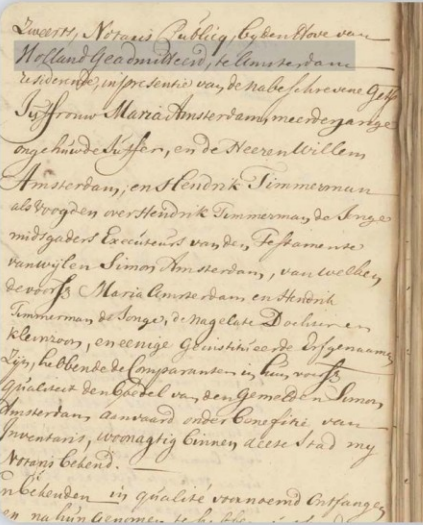
For Research & Cultural Institutions

- ✓ User Management
- ✓ Dedicated Success Manager
- ✓ API Access

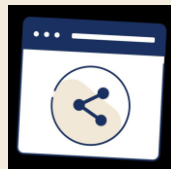
Get in Touch

100 Free Credits / Month

# Transkribus stránky - vlastnosti



Op Heden den 27sten Maart Ao. 1790  
Compareerden voor my Philip  
Zweerts, Notaris Publicq, by den Hove van  
**Holland geadmitteerd, te Amsterdam**  
residerende, in presentie van de nabeschrevene Get.  
Juffrouw Maria **Amsterdam** meerderjange  
ongehuwde Suffer, en de Heeren Willem  
**Amsterdam**, en Hendrik Timmerman  
als Voogden over Hendrik Timmerman de Jonge  
midsgaders Executeurs van den Testamente  
van wijlen Simon **Amsterdam**, van welken  
devoorsz Maria **Amsterdam** en Hendrik  
Timmerman de Jonge, de nagelate Dochter en  
kleinzoon, en eenige geinstitueerde erfgenaamen  
Zijn, hebbende de Comparanten en hun voorsz  
qualiteit den boedel van, den gemelden Simon  
**Amsterdam**, aanvaard onder Conscriptie van  
Inventaris, woonagtig binnen dese Stad my  
Notaris bekend.  
En bekenden in qualiteit voornoemd Ontfangen  
en na hungenomen te hebben, uithanden  
van Juff. Anna westerveen weduwe van  
simon **Amsterdam** voornoemd alle de meubilen



Jednoduché zdieľanie materiálu

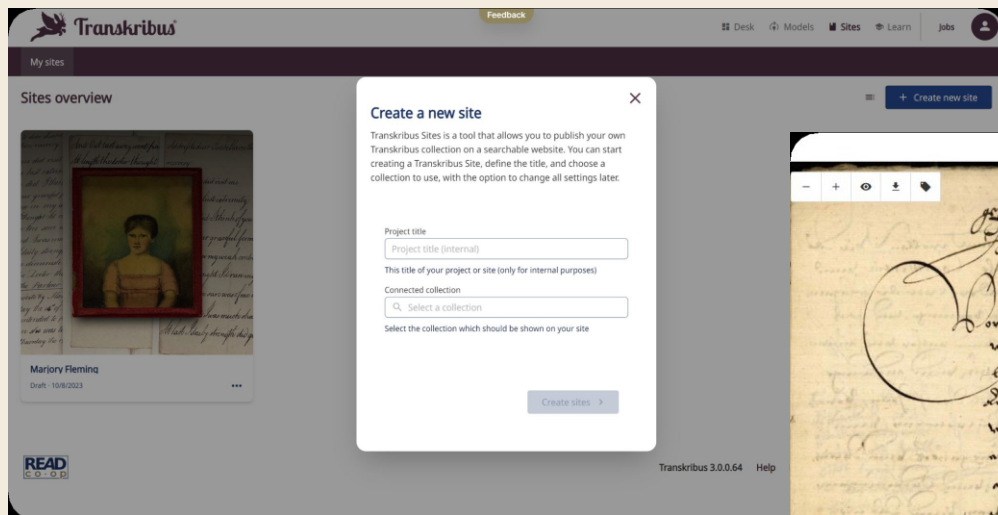


Pohľad strana vedľa strany  
(obrázok-prepis)



Vylepšené možnosti vyhľadávania

# Transkribus stránky

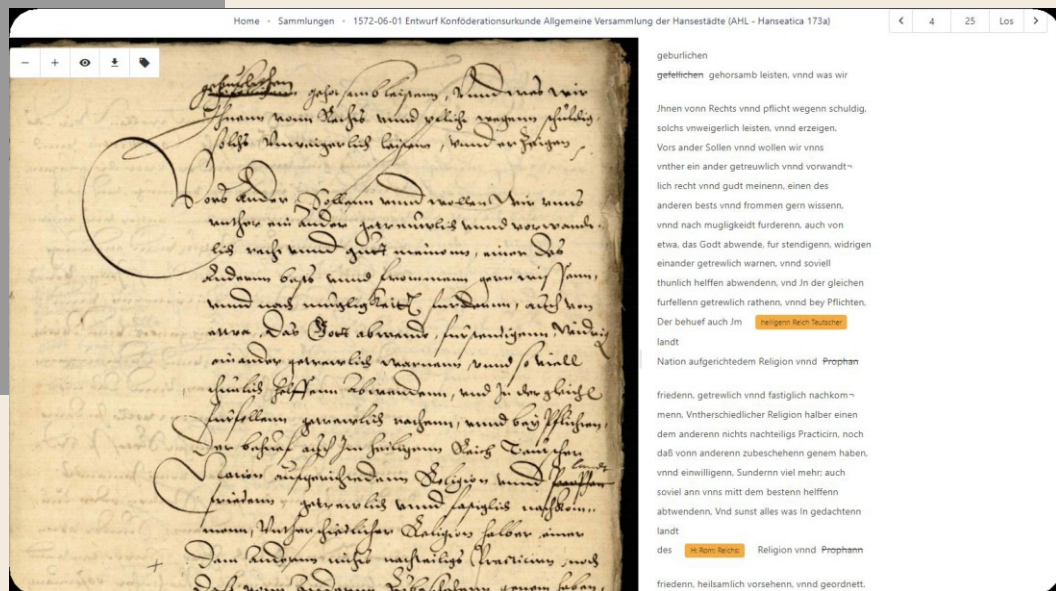


The screenshot shows the Transkribus interface with a 'Create a new site' modal dialog open. The dialog has a title bar with a close button (X) and a '+ Create new site' button. The main text reads: 'Transkribus Sites is a tool that allows you to publish your own Transkribus collection on a searchable website. You can start creating a Transkribus Site, define the title, and choose a collection to use, with the option to change all settings later.'

The form contains the following fields:

- Project title:** A text input field with the placeholder 'Project title (internal)'.
- This title of your project or site (only for internal purposes):** A text input field.
- Connected collection:** A dropdown menu with the placeholder 'Select a collection'.
- Select the collection which should be shown on your site:** A text input field.

At the bottom of the dialog is a 'Create sites >' button. In the background, the 'Sites overview' page is visible, showing a card for 'Marjory Fleming' with a portrait and the text 'Draft: 10/8/2023'. The bottom of the interface shows the 'READ' logo and the version 'Transkribus 3.0.0.64 Help'.



The screenshot displays a transcription page for a historical document. The top navigation bar includes 'Home', 'Sammlungen', and the document title '1572-06-01 Entwurf Konföderationsurkunde Allgemeine Versammlung der Hansestädte (AHL - Hanseatica 173a)'. On the right, there are navigation icons for back, page 4, page 25, and search (Los).

The main content area is split into two columns:

- Left column:** A high-resolution image of a handwritten manuscript page in Gothic script. The text is dense and written in dark ink on aged paper.
- Right column:** A transcription of the text from the manuscript. The text is in a modern, clean font. It begins with 'gebürlichen gehorsamb leisten, vñnd was wir' and continues with several lines of text. There are some highlighted words in orange, such as 'heiligen Reich' and 'Nation'.

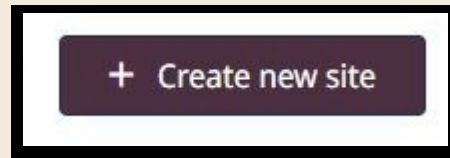
The transcription text is as follows:

gebürlichen gehorsamb leisten, vñnd was wir  
gebürlichen gehorsamb leisten, vñnd was wir  
/hnen vom Rechts vñnd pflicht wegn schuldig,  
solchs vñnweigerlich leisten, vñnd erzeigen.  
Vors ander Sollen vñnd wollen wir vnns  
vñnther ein ander getrewlich vñnd vorwandt-  
lich recht vñnd gndt meinnen, einen des  
anderen bests vñnd frommen gern wissen,  
vñnd nach muglichkeit fuderern, auch von  
etwa, das Godt abwende, fur stendigen, widrigen  
einander getrewlich warnen, vñnd soviel  
thunlich helfen abwendenn, vñnd In der gleichen  
furfellenn getrewlich rathenn, vñnd bey Pflichten.  
Der behuf auch im **heiligen Reich** / **Teutsche**  
landt  
Nation aufgerichtetem Religion vñnd Prophan  
friedenn, getrewlich vñnd fastiglich nachkom-  
menn, Vñnterschiedlicher Religion halber einen  
dem anderen nichts nachtheiligs Practicirn, noch  
daß vonn anderen zubeschehenn genem haben,  
vñnd einwilligen, Sundern viel mehr: auch  
soviel ann vnns mitt dem bestenn helfenn  
abwendenn, Vñnd stund alles was In gedachtem  
landt  
des **heiligen Reich** Religion vñnd Prophan  
friedenn, heilsamlich vorsehenn, vñnd geordnet.

# Vaša prvá stránka Transkribus

## Vytvorenie novej stránky

- Názov projektu
- Vlastná webová adresa([app.transkribus.eu/sites/yourchosenname](https://app.transkribus.eu/sites/yourchosenname))
- Prepojené zbierky



# Vaša prvá stránka Transkribus

## **3 editovateľné stránky:**

- **Domov**
- **O**
- **Preskúmať**

upravovať stránky a zobrazovať aktualizácie súčasne, vedľa seba

# Vaša prvá stránka Transkribus

## Domov: ( Home - Domovská stránka)

- Titul
- Stručný opis obsahu/stránky
- Obrázok pozadia domovskej stránky

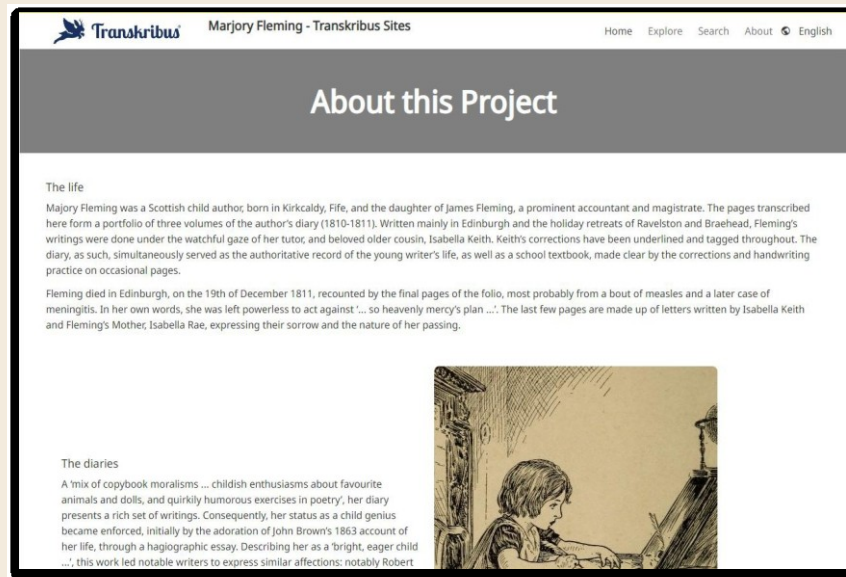


# Vaša prvá stránka Transkribus

## O (About)

(Vysvetlenie projektu,  
obsah, tím...):

- Toľko sekcií, koľko chcete
- Každá časť: nadpis - text - obrázok (voliteľné)



The screenshot displays the 'About this Project' page for Marjory Fleming on the Transkribus website. The page features a dark header with the Transkribus logo and navigation links (Home, Explore, Search, About, English). The main content area is titled 'About this Project' and contains two sections: 'The life' and 'The diaries'. The 'The life' section provides a biographical overview of Marjory Fleming, mentioning her birth in Kirkcaldy, her family, and her education. The 'The diaries' section describes the content of her diaries, including moralisms, childish enthusiasms, and humorous exercises. An illustration of a young girl writing at a desk is positioned to the right of the 'The diaries' text.

Transkribus Marjory Fleming - Transkribus Sites Home Explore Search About English

### About this Project


**The life**

Majory Fleming was a Scottish child author, born in Kirkcaldy, Fife, and the daughter of James Fleming, a prominent accountant and magistrate. The pages transcribed here form a portfolio of three volumes of the author's diary (1810-1811). Written mainly in Edinburgh and the holiday retreats of Ravelston and Braehead, Fleming's writings were done under the watchful gaze of her tutor, and beloved older cousin, Isabella Keith. Keith's corrections have been underlined and tagged throughout. The diary, as such, simultaneously served as the authoritative record of the young writer's life, as well as a school textbook, made clear by the corrections and handwriting practice on occasional pages.

Fleming died in Edinburgh, on the 19th of December 1811, recounted by the final pages of the folio, most probably from a bout of measles and a later case of meningitis. In her own words, she was left powerless to act against '... so heavenly mercy's plan ...'. The last few pages are made up of letters written by Isabella Keith and Fleming's Mother, Isabella Rae, expressing their sorrow and the nature of her passing.

**The diaries**

A 'mix of copybook moralisms ... childish enthusiasms about favourite animals and dolls, and quirkily humorous exercises in poetry', her diary presents a rich set of writings. Consequently, her status as a child genius became enforced, initially by the adoration of John Brown's 1863 account of her life, through a hagiographic essay. Describing her as a 'bright, eager child ...', this work led notable writers to express similar affections: notably Robert



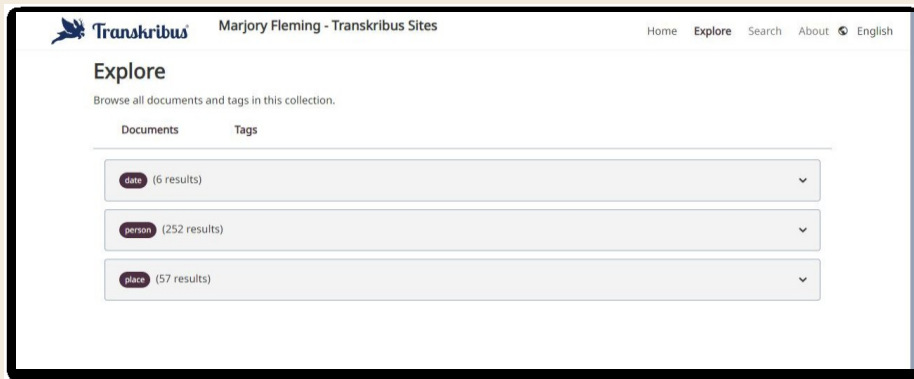


# Vaša prvá stránka Transkribus

## Preskúmať (Explore)

(Ako chcete nakonfigurovať stránku vyhľadávania):

- Povolenie značiek prehľadávania
- Povolené značky (ak ste použili značky vo vašich dokumentoch Transkribus)
- Povoľiť filtre a filter rokov (na základe metadát dokumentov Transkribus)

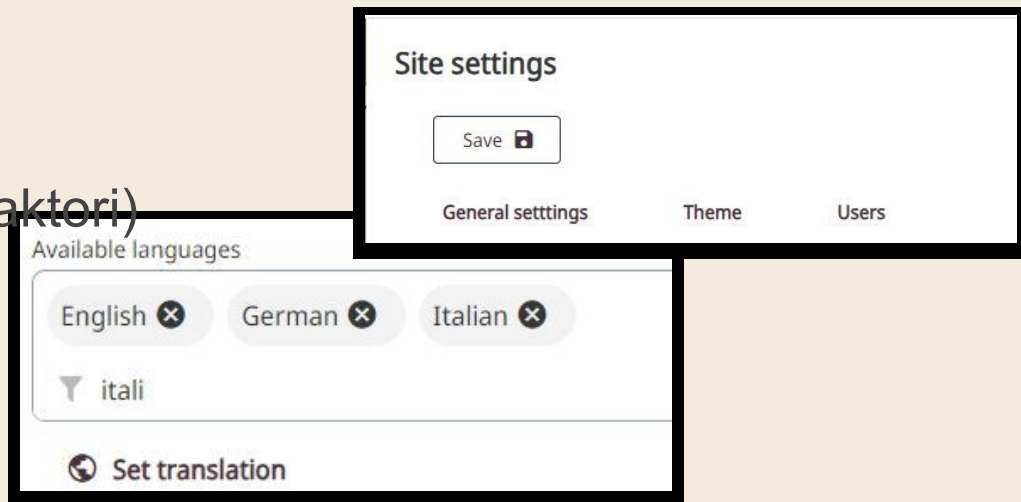


# Vaša prvá stránka Transkribus

[Read&Search - Demo \(transkribus.eu\)](https://transkribus.eu)

## Ďalšie nastavenia (Other settings):

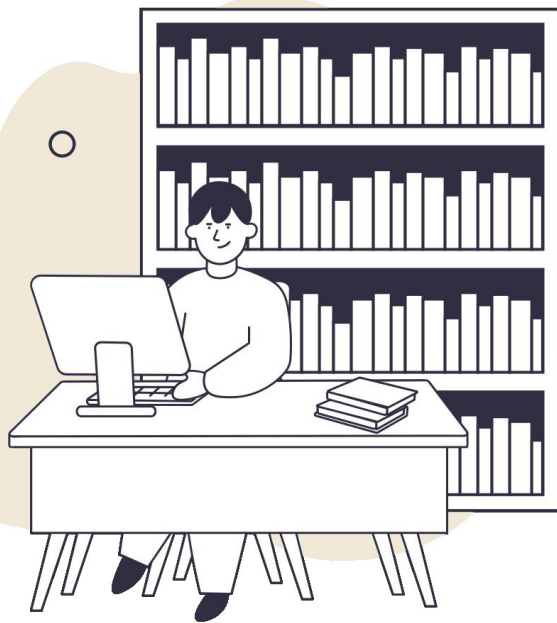
- Jazyky + možnosť úpravy prekladov
- Súkromie
- Motív (logo a farba)
- Používatelia (vlastník, redaktori)





Čas na otázky





Hands-on  
session  
Praktické  
sedenie



# Help Center

<https://help.transkribus.org/>



# Thank you!

Website: <https://transkribus.org/>

Email addresses:

[s.mansutti@readcoop.eu](mailto:s.mansutti@readcoop.eu)

[m.elattal@readcoop.eu](mailto:m.elattal@readcoop.eu)

[info@readcoop.eu](mailto:info@readcoop.eu)



Unlocking the past, together

