

# Transkribus & Pero OCR

Miniúvod k HTR a OCR technologiím

Anna Michalcová  
ÚČJTK, FF UK

12. 10. 2022



# OCR vs. HTR

- **OCR (Optical Character Recognition)**
  - technologie, která umí rozpoznávat jednotlivé znaky (může být problematické v případě starého textu)
  - pokud OCR umožňuje vytvořit trénovací data, zvýší se i úspěšnost jeho použití na konkrétní dokument
- **HTR (Handwritten Text Recognition)**
  - pracuje s celými řádky a dekoduje řádek jako celek
- Hlavní rozdíl z pohledu uživatele spočívá v tom, že fáze analýzy rozložení/segmentace je integrována do mechanismu OCR, zatímco v případě HTR se jedná o samostatný krok v pracovním postupu

Pero	Transkribus
s hiewem mluwim, a ni <sup>z</sup> to wiedieti <sup>z</sup> e	shniewem mluwím. ani <sup>z</sup> to wiedieti <sup>z</sup> e-
Patřitiť slusie snevmieš vi	Patřitiť flulie Gnevmieš. vii.
přeznameť su knihy genezis	přeznaňeť fu knihý geneziš.
gesto mluwie oftwořeňie swie	gešto mluwie oftwořeňie fwie-
ta, opočeti lidske° rodu, orozdielení ze	ta. opočeti lidkeo rodu. orozdieleni ze-
mie, ofmiesfeni iazvkow, ogiti židow	nie. ofmieffení iazyków. osgiti židow
skem až dogipta. Zgewny su knihy s	fkém aždoegipta. Zgewny fu knihy k-
rodus toczyž wystie sdesieti pohroma	rodus toczyž wyťtie. sdefieti pohroma-
mi, sdesaterem přikazanim sducho	mi. sdefaterem přikazaňim sducho
wnimi a božskymi přikazanimi Q	wními a božkyymi přikazaňimi. W-
hotowie su knihy lewitikus, w nichžto	hotowie fu knihý lewitikus. w nichžto
wfeli <sup>a</sup> swiecena obiet Nebrž wseli	wfeli <sup>a</sup> fwiecena obiet Nebrž wfeli
ka temierz serčenje a rucha aaronowa	ka temierz ferčenje a rucha aaronowa-
a weslken řad kniežsky udychawa ne	a weslken řad kniežky wdychawa ne-
beska tagemstwie. Ale knihy počtu. y	belka tagemftwie. Ale knihý počtu. y
zdali neosahuge wfe° početne° vmie	zdali neofáhuge wf eo vmie-
nie a proročtwie Balaamowa a rlti	ňie a proroctwie Balaamowa. a. xlti.
a dwu stanowiffti napusti tavnosti. De	a dwu ftanowiffti napúfti tayofti. De-
vtronomivm knihy druhe° zakona.	vtronomiyňi knihy druheo zakona-
y nowe° zakona podobenstwie. yzda	y noweo zakona podobenftwie. Yzda-
li netak ma tiech genž prwnie su ia	li netak ma tiech genž prwnie fu. ia
koby wsak prwnievsie nowa byla	koby wlfak prwnieyľie nowa byla
wsecka zwetchych až dofawad. penta	wľecka zwetchych až dofawad. Penta-
tenkus totižto pacery knihy ge° gimíž	teukus toczyto patery knihy geo gimí ž
to pieti slowy chlubi sie a postol že chce	to píeti flowy chlubi fie a poľtol že chce-
mluwiti wtierkwi božie Job prziklad	mluwiti wćierkwi božie Job prziklad
trpieliwosti, kterychž tavnosti swu řeči	trpieliwofti. kterýchž tayofti fwú řeči

skniewem mluwim. ani<sup>z</sup>to wiedieti <sup>z</sup>e  
**S** Patřitiť slusie snevmieš. vii.  
 přeznameť su knihy genezis.  
 gesto mluwie oftwořeňie swie  
 ta. opočeti lidske° rodu. orozdieleni ze-  
 mie. ofmiesfeni iazyków. osgiti židow  
 skem až doegipta. Zgewny su knihy s  
 rodus toczyž wystie sdesieti pohroma  
 mi. sdesaterem přikazanim sducho  
 wnimi a božskymi přikazanimi Q  
 hotowie su knihy lewitikus, w nichžto  
 wfeli<sup>a</sup> swiecena obiet Nebrž wseli  
 ka temierz serčenje a rucha aaronowa  
 a weslken řad kniežsky udychawa ne  
 beska tagemstwie. Ale knihy počtu. y  
 zdali neosahuge wfe° početne° vmie  
 nie a proročtwie Balaamowa. a. xlti.  
 a dwu stanowiffti napusti tavnosti. De  
 vtronomivm knihy druhe° zakona.  
 y nowe° zakona podobenstwie. yzda  
 li netak ma tiech genž prwnie su. ia  
 koby wsak prwnievsie nowa byla  
 vsecka zwetchych až dofawad. penta  
 tenkus totižto pacery knihy ge° gimíž  
 to pieti slowy chlubi sie a postol že chce  
 mluwiti wtierkwi božie Job prziklad  
 trpieliwosti. kterychž tavnosti swu řeči

# Transliterace vs. transkripce

- PERO OCR i Transkribus využívají transliteraci
- oproti paleografickému přepisu nejde o nápodobu rukopisu (např. rozlišování různých grafických podob písmen s výjimkou rozlišování dlouhého a kulatého s), ale o přepis znak za znak
- „**Transliterace** je pohodlné řešení otázek, vlastně pohodlné přesunutí všech sporných otázek z vydavatele na činitele. **Transkripce** je odnětí těch svízelných zhodnocení, co je (z hlediska vyjadřovacího) obsah a co forma, uživatelům edic a podání již jistého názoru na text a jeho obsah. Tedy jistá interpretace.“
- DAŇHELKA, Jiří (2013). Dopis Jiřího Daňhelky Marii Skalické, in: *Textologie a starší česká literatura*. Ed. J. Sichálek. Praha: Ústav pro českou literaturu Akademie věd České republiky v. v. i., s. 213–214.



# PERO OCR

- „Pokročilá extrakce a rozpoznávání obsahu tištěných a rukou psaných digitalizátů pro zvýšení jejich přístupnosti a využitelnosti má za cíl **vytvořit nástroje a technologie pro zpřístupnění obsahu digitalizovaných historických dokumentů s využitím nejnovějších poznatků v oblasti počítačového vidění, strojového učení a jazykového modelování.** Hlavním řešitelem projektu je Fakulta informačních technologií Vysokého učení technického v Brně, Moravská zemská knihovna má roli spoluřešitele (NAKI II, 2018–2022).“

- <https://pero-ocr.fit.vutbr.cz/document/documents>

výsledek OCR



# Transkribus

- Transkribus je platforma pro **rozpoznávání textu, analýzu obrazu a rozpoznávání struktury historických dokumentů**. Platforma byla vytvořena v rámci dvou projektů EU tranScriptorium a READ. Byl vyvinut univerzitou v Innsbrucku. Od 1. července 2019 platformu řídí a dále rozvíjí **READ-COOP**.
- nahrání dokumentů, transkripční nástroje zdarma × automatický přepis placený (kreditový systém)

1-1 The Dean stated, that this meeting was called in  
1-2 consequence of a letter he had received from several  
1-3 members which he desired might be read. It was  
1-4 accordingly read by the clerk and is as follows.  
1-5 "Edinburg, 25 Nov. 180  
1-6 11  
1-7 Sir, - We request that you will call a meeting  
1-8 of the Members of the Faculty, for the purpose of  
1-9 considering a bill lately brought into the House of  
1-10  
1-11  
1-12  
1-13 We do not presume to suggest any particular  
1-14  
1-15 with regard to it, than that it should take place  
1-16  
1-17 and at such distance of time as will give the Mem  
1-18 bers of the Faculty full opportunity to bestow that  
1-19 attention on the subject which its importance calls  
1-20  
1-21 /Signed  
1-22 5 Jeffrey  
1-23 William Erokine, Henry Cockburn.  
1-24 George Jos. Bell. I Henry Mackenzie. John A Murray  
1-25 Tho. Thomson  
1-26 The following resolutions were then  
1-27 moved by Mr Jeffrey, and seconded by Mr John  
1-28 Murray:  
1-29 That it is the opinion of this meeting

# Transkribus

- **kdo všechno Transkribus využívá?**

- archivy: Berlínský státní archiv, Zeit.punkt (digitalizace novin ze Severního Porýní-Vestfálska z let 1801–1945), Archiv města Amsterdam, Národní archiv Finska, Arolsen Archives (International Center on Nazi Persecution) – archivy nyní prohledatelné

- univerzity, akademie: k vytvoření korpusu textů, umožňuje pracovat na transkripci z různých míst

- Kati: University of Innsbruck / Federal State of Tirol

- EnrichEuropean+ (podpora hromadné práce pro evropské kulturní dědictví)

- Österreichischer Bibelübersetzer (Berlin-Brandenburgischen Akademie der Wissenschaften)

- Rakouská akademie věd + READ-COOP (optické rozpoznávání hudby)

1-1 The Dean stated, that this meeting was called in  
1-2 consequence of a letter he had received from several  
1-3 members which he desired might be read. It was  
1-4 accordingly read by the clerk and is as follows.  
1-5 "Edinburg, 25 Nov. 180  
1-6 11  
1-7 Sir, - We request that you will call a meeting  
1-8 of the Members of the Faculty, for the purpose of  
1-9 considering a bill lately brought into the House of  
1-10 Lords, entitled, An Act touching the Administration  
1-11 of Justice in Scotland, and touching Appeals to the  
1-12 any particular  
1-13 we have no farther wish  
1-14 with regard to it, than that it should take place  
1-15 upon whatever day will be most convenient for you,  
1-16 the Mem  
1-17 bers of the Faculty full opportunity to bestow that  
1-18 attention on the subject which its importance calls  
1-19 for.  
1-20 for.  
1-21 /Signed  
1-22  
1-23 William Euskins, Henry Cockburn,  
1-24 George Jos. Bell, Henry Mackenzie, John A Murray  
1-25 Tho. Thomson  
1-26 The following resolutions were then  
1-27 moved by Mr Jeffrey, and seconded by Mr John  
1-28 Murray:  
1-29 That it is the opinion of this meeting



# Transkribus

## APPLICATION

Please, send a **short CV** and **cover letter** to Jan Odstrčilik | [jan.odstrcilik@oeaw.ac.at](mailto:jan.odstrcilik@oeaw.ac.at) before 11<sup>th</sup> September 2022. As a subject, use „HTR School 2022“.

Please don't forget to indicate the team/module you would like to take (i.e. Carolingian Latin, late medieval Latin, medieval German, medieval Czech). If you have any questions, don't hesitate to ask.

**Michael Berger**, University of Vienna | [michael.berger@univie.ac.at](mailto:michael.berger@univie.ac.at)

**Tim Geelhaar**, Bielefeld University | [tim.geelhaar@uni-bielefeld.de](mailto:tim.geelhaar@uni-bielefeld.de)

**Tobias Hodel**, University of Bern | [tobias.hodel@unibe.ch](mailto:tobias.hodel@unibe.ch)

**Sarah Hutterer**, Universität Wien | [sarah.hutterer@univie.ac.at](mailto:sarah.hutterer@univie.ac.at)

**Anna Michalcová**, Czech Academy of Sciences | [a.michalцова@ujc.cas.cz](mailto:a.michalцова@ujc.cas.cz)

**Stephan Müller**, University of Vienna | [stephan.mueller@univie.ac.at](mailto:stephan.mueller@univie.ac.at)

**Jan Odstrčilik**, Austrian Academy of Sciences | [jan.odstrcilik@oeaw.ac.at](mailto:jan.odstrcilik@oeaw.ac.at)

**Steffen Patzold**, University of Tübingen | [steffen.patzold@uni-tuebingen.de](mailto:steffen.patzold@uni-tuebingen.de)

**Leon Pürstinger**, Austrian Academy of Sciences | [leon.puerstinger@oeaw.ac.at](mailto:leon.puerstinger@oeaw.ac.at)

**Helmut Reimitz**, Princeton University | [hreimitz@princeton.edu](mailto:hreimitz@princeton.edu)

**Dennis Wegener**, University of Vienna | [dennis.wegener@univie.ac.at](mailto:dennis.wegener@univie.ac.at)

**Vicent Bosch**, Transkriptorium, Universitat Politècnica de València | [vbosch@transkriptorium.com](mailto:vbosch@transkriptorium.com)

**Gerda Heydemann**, Freie Universität Berlin | [gerda.rummel-heydemann@fu-berlin.de](mailto:gerda.rummel-heydemann@fu-berlin.de)

**Daniela Mairhofer**, Princeton University | [daniela.mairhofer@princeton.edu](mailto:daniela.mairhofer@princeton.edu)

**Joan Andreu Sanchez**, Transkriptorium, Universitat Politècnica de València | [jandreu@prhlt.upv.es](mailto:jandreu@prhlt.upv.es)

**Alejandro Toselli**, Transkriptorium, Universitat Politècnica de València | [ahector@prhlt.upv.es](mailto:ahector@prhlt.upv.es)

**Enrique Vidal**, Transkriptorium, Universitat Politècnica de València | [evidal@prhlt.upv.es](mailto:evidal@prhlt.upv.es)



CALL FOR PARTICIPATION  
WINTER SCHOOL | OCTOBER - DECEMBER 2022

INTRODUCTION INTO HTR  
HANDWRITTEN TEXT RECOGNITION  
TECHNOLOGIES OF MEDIEVAL MANUSCRIPTS  
LATIN|GERMAN|CZECH



Organisation:  
Institute for Medieval Research of  
the Austrian Academy of Sciences  
Manuscript, Rare Books and  
Archival Studies Initiative  
(MARBAS), Princeton University

In cooperation with:  
CRC 1288 Practices of  
Comparing, University Bielefeld  
Department of German Studies,  
University of Vienna  
Digital Humanities,  
Walter Benjamin Kolleg,  
Universität Bern  
History Department, University  
of Tübingen  
tranSkriptorium, Valencia  
[tranSkriptorium.com](http://tranSkriptorium.com)

Contact:  
Dr. Jan Odstrčilik  
Austrian Academy of Sciences  
[Jan.Odstrcilik@oeaw.ac.at](mailto:Jan.Odstrcilik@oeaw.ac.at)

4 Zoom online sessions | Oct 21, Nov 4 and 25, Dec 9  
3-day-workshop in person at Vienna | December 19-21



# Transkribus

## Polish – General Model

Free Public AI Model for Handwritten Text Recognition with [Transkribus](#)

[Transkribus](#) > [Public Models](#) > [Polish – General Model](#)

[« Back to all public models](#)

A general model for historical and modern Polish handwriting.

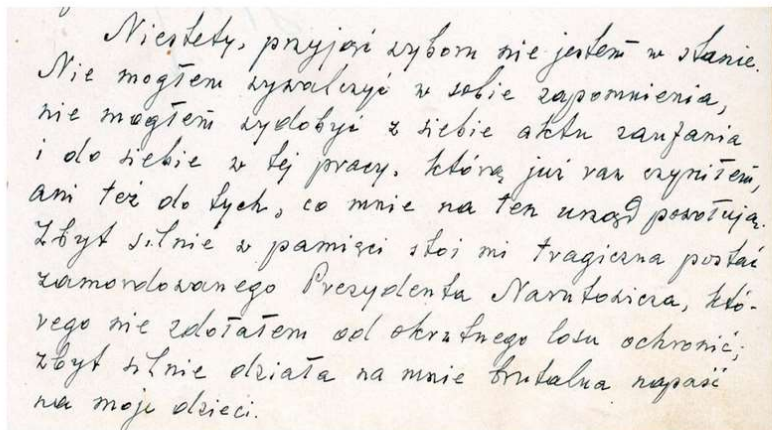


Image Source: [Wikipedia](#)

### Model Overview

**Name:**  
Transkribus Polish M2

**Creator:**  
Transkribus Team

**Model ID:**  
44976

**Century:**  
n/a

**Languages:**  
Polish

**Script:**  
Latin alphabet

**Engine:**  
PyLaia

**Material:**  
Handwritten



**Štátna vedecká knižnica v Banskej Bystrici**  
a  
**Filozofická fakulta Univerzity Mateja Bela v Banskej Bystrici**

**Vás pozývajú**  
na vedeckú konferenciu s medzinárodnou účasťou

**Digital humanities**  
**Nástroje sprístupňovania historického dedičstva**

usporiadanú v dňoch 12. – 13. októbra 2022

**Štátna vedecká knižnica v Banskej Bystrici**  
Lazovná 9



AGENTÚRA  
NA PODPORU  
VÝSKUMU A VÝVOJA



# Transkribus

- nyní hlavně podpora Transkribus LITE, ale nemá zatím všechny funkce
- <https://transkribus.eu/lite/home>

1-1 The Dean stated, that this meeting was called in  
1-2 consequence of a letter he had received from several  
1-3 members which he desired might be read. It was  
1-4 accordingly read by the clerk and is as follows.  
1-5 "Edinburg, 25 Nov. 180  
1-6 11  
1-7 Sir, - We request that you will call a meeting  
1-8 of the Members of the Faculty, for the purpose of  
1-9 considering a bill lately brought into the House of  
1-10 Lords, entitled, "An Act for the Administration  
1-11 of Justice in Scotland, and touching Appeals to the  
1-12 House of Lords.  
1-13 We do not presume to suggest any particular  
1-14 day for that purpose, as we have no farther wish  
1-15 with regard to it, than that it should take place  
1-16 upon whatever day will be most convenient for you,  
1-17 and at such distance of time as will give the Mem  
1-18 bers of the Faculty full opportunity to bestow that  
1-19 attention on the subject which its importance calls  
1-20 for.  
1-21 /Signed  
1-22 5 Jeffrey  
1-23 William Erokine, Henry Cockburn.  
1-24 George Jos. Bell. I Henry Mackenzie. John A Murray  
1-25 Tho. Thomson  
1-26 The following resolutions were then  
1-27 moved by Mr Jeffrey, and seconded by Mr John  
1-28 Murray:  
1-29 That it is the opinion of this meeting



Server Overview Layout Metadata Tools

Layout Analysis

Method: CITIab Advanced Configure...

Current page Pages (1): 1-1

Find Text Regions Only use un

Find Lines in Text Regions

Run

Text Recognition

Method: HTR (CITIab)

Models... Train...

Run...

Compute Accuracy

Other Tools

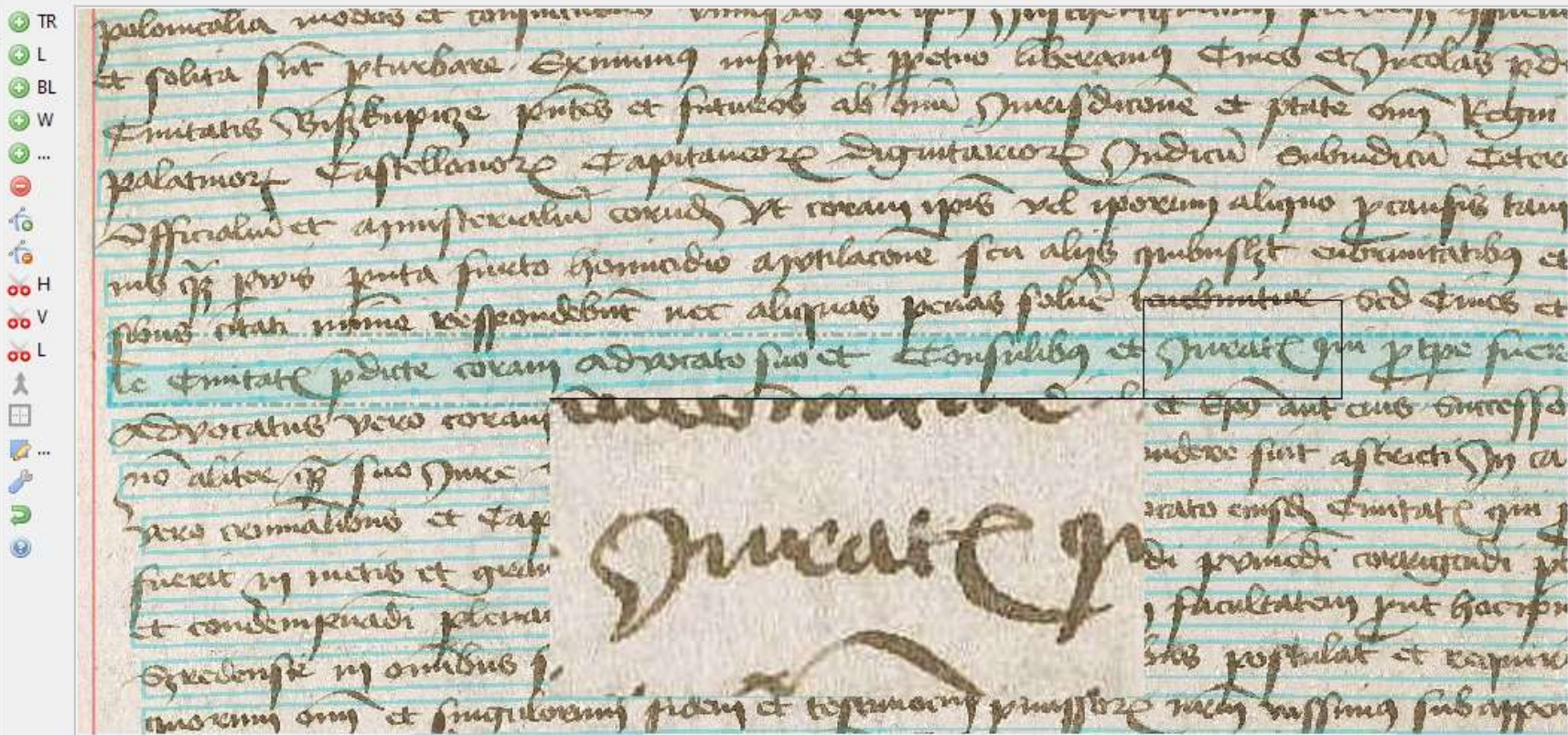
P2PaLA...

Text2Image...

Current page Pages (1): 1-1

Add Baselines to Polygons

Add Polygons to Baselines



1-10 sibus· citati· minime· respondebunt· nec· aliquas· penas· solvere· tenebuntur· sed·  
cives· et· Inco↵

1-11 le· civitatis· predicte· coram· advocato· suo· et· consulibus· et· ↵



Děkuji za pozornost!

