

# DIGITÁLNY PRÍBEH KOMENSKÉHO UČEBNICE *ORBIS PICTUS* (1798)

## THE DIGITAL STORY OF COMENIUS' TEXTBOOK *ORBIS PICTUS* (1798)

Lucia Nižníková – Michaela Mikušková

Univerzitná knižnica Univerzity Mateja Bela v Banskej Bystrici

### Abstrakt

**Účel:** Príspevok patrí do oblasti digitálnych humanitných vied a dokumentuje využitie umelej inteligencie pri spracovaní kultúrneho dedičstva. Zameriava sa na praktické uplatnenie nástrojov HTR+ (Handwritten Text Recognition) a platformy Transkribus v kontexte automatického rozpoznávania a transkripcie historických tlačí a skúmania ich využiteľnosti pri vytváraní univerzálnych modelov transkripcie. Cieľom príspevku nie je analýza, ani porovnávanie rôznych nástrojov transkripcie rukopisných textov a historických tlačí, ale zhodnotenie niekoľkoročnej práce so softvérom Transkribus na demonštrácii výsledkov, ktoré autorky dosiahli ako súčasť aplikovaného výskumu v projekte SKRIPTOR (2020–2024).

**Metodológia / prístup:** Výskum bol realizovaný v prostredí *Transkribus Expert Client* a *Transkribus App* s využitím technológií HTR+ a PyLaia. Postup zahŕňal digitalizáciu tlače *Orbis Pictus* (1798) pomocou nástrojov *ScanTent* a *DocScan*, manuálnu segmentáciu textu, tvorbu vzorky *Ground Truth* a tréning modelov automatickej transkripcie pre štyri typy písma (antikva, kurzíva, fraktúra, švabach) a štyri jazyky (latinčina, maďarčina, nemčina, čeština). Vytrénované modely boli overované na historických tlačiach *Orbis Pictus* (1820) a *Adparatus ad Historiam Hungariae* (1735). Na základe ukazovateľov chybovosti CER a WER sa vyhodnocovala ich univerzálnosť.

**Výsledky:** Najlepšie výsledky pri transkripcii tlače *Orbis Pictus* (1798) boli dosiahnuté v Modeli 11 s chybovosťou CER 1,00 %. Následne vytrénovaný nový model *Bel\_Adparatus\_1735\_model 3* s chybovosťou CER 0,93 % sa ukázal ako stabilný a spoľahlivý pri transkripcii tlače *Adparatus ad Historiam Hungariae*, avšak nevykazoval univerzálnu použiteľnosť pre iné tlače s odlišnými fontmi. Porovnanie technológií potvrdilo vyššiu presnosť PyLaia oproti HTR+ a efektívnosť využitia *base modelov*. Výskum poukázal aj na limitujúce faktory ako sú kvalita digitalizátu, rozloženie textu a výskyt špecifických grafém. Testovanie supermodelu *Transkribus Print M1* preukázalo jeho vysokú presnosť a použiteľnosť pre viacjazyčné historické tlače.

**Originalita / hodnota:** Štúdia prináša komplexný pohľad na aplikáciu umelej inteligencie v oblasti spracovania a sprístupňovania historických tlačí v slovenskom prostredí. Originálnym prínosom je vytrénovanie modelu pre štvorjazyčný dokument so štyrmi typmi písma a overenie možností jeho univerzálného využitia. Výsledky a závery prezentované v štúdiu potvrdzujú, že Transkribus predstavuje efektívny nástroj na automatickú transkripciu písomného kultúrneho dedičstva a otvára nové možnosti pre interdisciplinárny výskum v oblasti digitálnych humanitných vied.

**Kľúčové slová:** automatická transkripcia, historické tlače, Transkribus, digitálne humanitné vedy

-----  
Táto práca bola podporená Agentúrou na podporu výskumu a vývoja na základe zmluvy č. APVV-19-0456 *Inovatívne sprístupnenie písomného dedičstva Slovenska prostredníctvom systému automatickej transkripcie historických rukopisov (SKRIPTOR)*.

This paper was supported by the Slovak Research and Development Agency under contract no. APVV-19-0456 *Innovative disclosure of written heritage of Slovakia through the automatic transcription of historical manuscripts (SKRIPTOR)*.

<https://doi.org/10.24040/123456789-1266>



Táto publikácia je šírená pod licenciou Creative Commons Attribution 4.0 International Licence CC BY.

## Abstract

*Purpose:* This paper contributes to digital humanities research and demonstrates the use of Artificial Intelligence in the processing of written cultural heritage. It focuses on the practical application of HTR+ (Handwritten Text Recognition) tools and the Transkribus platform in the context of automatic recognition and transcription of historical prints and examines their usability in creating universal transcription models. The aim of the paper is not to analyze or compare various tools for the transcription of handwritten texts and historical prints, but to make the evaluation of several years work with Transkribus software based on the results achieved by the authors as part of the applied research project SKRIPTOR (2020–2024).

*Methodology / Approach:* The research was carried out with the Transkribus Expert Client and the Transkribus App software using HTR+ and PyLaia technologies. The procedure included digitizing the *Orbis Pictus* (1798) textbook using the *ScanTent* and *DocScan* tools, manual text segmentation, creation of the Ground Truth sample, and training automatic transcription models for four font types (antiqua, cursive, fraktur, Schwabacher) and four languages (Latin, Hungarian, German, Czech). The best trained models were verified on historical prints of *Orbis Pictus* (1820) and *Adparatus ad Historiam Hungariae* (1735). Their universality was evaluated using the CER and WER indicators.

*Results:* The best results for the transcription of the *Orbis Pictus* (1798) print were achieved in Model 11 with a character error rate CER 1.00%. The subsequently trained model *Bel\_Adparatus\_1735\_model 3* with a CER error rate of 0.93% proved to be stable and reliable in transcribing the *Adparatus ad Historiam Hungariae* print book, but did not show universal applicability for other prints with various fonts. A comparison of technologies confirmed the higher accuracy of PyLaia compared to HTR+ and the effectiveness of using base models. The research also pointed to limiting factors in the quality of the digitized material, the layout of the text, and the appearance of specific graphemes. Testing the Transkribus Print M1 supermodel demonstrated its high accuracy and usability for historical multilingual prints.

*Originality / Value:* The study provides a comprehensive overview of the application of Artificial Intelligence in processing and making historical prints accessible in the Slovak environment. An original contribution is the training of a model for a quadrilingual document with four types of font and the verification of its universal applicability. The results and conclusions presented in the study confirm that Transkribus is an effective tool for automatic transcription of written cultural heritage and opens up new possibilities for interdisciplinary research in the field of digital humanities.

*Keywords:* automatic transcription, historical prints, Transkribus, digital humanities

Umelá inteligencia dnes zasahuje prakticky do každej oblasti nášho pracovného i osobného života. Mení spôsob, akým žijeme, mení naše návyky. Ovplyvňuje priemysel, výrobu, zdravotníctvo, školstvo, vedu nevynímajúc, a prináša so sebou výzvy i riziká. Aktuálne v spoločnosti prebieha debata o tom, ako umelú inteligenciu uchopiť a čo najefektívnejšie aplikovať do praxe. Možnosťami jej využívania vo výskume sa zaoberajú aj vedci z oblasti humanitných vied. Ide o spájanie rôznych vedných odborov (nielen v rámci humanitných vied) s informačnými technológiami a počítačovou vedou a použitie digitálnych nástrojov ako geografické informačné systémy, text mining, 3D modelovanie, databázová analýza a ďalších na analýzu a vizualizáciu humanitných dát. Prvé náznaky prepájania humanitných vied a výpočtovej techniky možno spätne vystopovať až k priekopníckej práci z konca 40. rokov 20. storočia a k modelom, ktoré inšpirovali archívne projekty v Oxforde na začiatku 70. rokov 20. storočia. Od konca 80. rokov až do začiatku 21. storočia sa rozvinula prvá vlna digitálnych humanitných vied, ktorá sa zaoberala spôsobmi štruktúrovania humanitných dát s cieľom efektívne komunikovať s výpočtovou technikou (Burdick et al., 2012, s. 8). *Digital humanities* dnes prinášajú nové formy výskumu, výučby a publikovania, vrátane kolaboratívnych a transdisciplinárnych prístupov. Typickým príkladom aktivít v tejto oblasti je použitie softvéru na analýzu veľkých textových korpusov (napr. literárnych diel, historických dokumentov), vytváranie digitálnych naratívov, rekonštrukcia historických miest v 3D forme, tvorba interaktívnych digitálnych sprievodcov, ale aj budovanie digitálnych archívov a databáz s cieľom sprístupniť historické a kultúrne pramene a písomníctvo. Digitalizácia prameňov uľahčuje ich sprístupnenie vedcom aj širšej verejnosti a otvára možnosti pre ich ďalší výskum z celkom nových perspektív.

Umelá inteligencia dokáže napodobňovať ľudské myslenie, učiť sa na základe veľkého množstva dát a čo je najdôležitejšie, zlepšovať svoje výkony v závislosti od získaných skúseností. Práve na tejto vlastnosti sú založené

európske projekty *tranScriptorium*<sup>1</sup>, riešený v rokoch 2013 – 2015 a financovaný zo 7. rámcového programu, a projekt *READ – Recognition and Enrichment of Archival Documents*<sup>2</sup>, realizovaný v rokoch 2016 – 2019 a financovaný zo zdrojov programu Horizont 2020 (Muehlberger, 2021). Ich výsledkom je inovatívny nástroj *Transkribus*, ktorý umožňuje automatické rozpoznávanie textu a transkripciu historických dokumentov. Predpokladom jeho využitia v praxi sú za pomoci umelej inteligencie vytrénované modely na prepis rôznych rukopisných štýlov a typov písma (tlačeného a strojopisného) vrátane špecifických grafém a rukopisných skratiek, a to bez ohľadu na jazyk dokumentu.



Obr. 1 Ponuka platformy Transkribus. Zdroj: <https://readcoop.eu/transkribus>

Metódy rozpoznávania rukopisných textov (HTR – Handwritten Text Recognition) využívajú umelé neurónové siete (ANN – Artificial Neural Networks) ako jeden z algoritmov strojového učenia. Ich využitie v praxi funguje na základe vytvorenia základného modelu, ktorý sa cvičí na vybranej vzorke dát. Na základe postupného zväčšovania objemu dát sa trénujúci model zdokonaľuje a zvyšuje sa jeho presnosť. Transkribus ponúka stovky vytrénovaných verejných modelov, ktoré boli vyškolené na určitý typ písma, rukopisu a typ dokumentu. Hoci bol primárne vyvíjaný ako nástroj na automatickú transkripciu rukopisných textov, veľmi dobré výsledky vykazuje aj na tlačiach a strojopisných zbierkach (Smida, 2023). Po ukončení pilotného projektu READ sa na jeho udržateľnosti a na ďalšom vývoji podieľala organizácia READ-COOP SCE združujúca rôzne vedecké, výskumné, vzdelávacie a archívne inštitúcie. Do roku 2024 boli dostupné dve verzie Transkribusu:

- *Transkribus Expert Client* – pôvodná softvérová verzia, ktorá sa inštaluje priamo do osobného počítača; od roku 2024 sa nepočíta s jej ďalším vývojom, ale s postupným presunom jej funkcií do webovej verzie<sup>3</sup>,
- *Transkribus Lite* (neskôr premenovaná na *Transkribus app*) – webová verzia, na ktorú je zameraný súčasný a budúci vývoj.

Obe platformy ponúkajú aj nástroje na editáciu, spracovanie a tagovanie transkribovaných textov. Výsledky transkripcie bolo spočiatku možné sprístupniť verejnosti na portáli *Read&Search*<sup>4</sup>, platenej platforme Transkribusu, ktorý sprístupňoval dokumenty zo zbierky vytvorenej na platforme Transkribus Expert Client online formou. Po prechode klienta na webovú verziu na tento účel slúži prostredie *Transkribus Sites*<sup>5</sup>. Rozhranie bohaté na funkcie je ideálne na sprístupnenie historických dokumentov a vyhľadávanie na webe. Alternatívne je, po exporte v podporovaných formátoch, možné výstupy transkripcie a zdrojové obrázky sprístupňovať aj na iných webových sídlach alebo v digitálnych repozitároch. Pridanou hodnotou je možnosť exportovať dokument aj podľa tagov.

Na tieto projekty nadväzuje slovenský projekt SKRIPTOR z rokov 2020 – 2024. Predmetom výskumu boli historické dokumenty, ich automatická transkripcia a uplatnenie digitálnych a informačných technológií v oblasti humanitných vied. Výskumným problémom projektu je tvorba čo najlepších modelov automatického rozpoznávania textov historických slovacikálnych dokumentov, ktoré predstavujú súčasť európskeho písomného dedičstva a prezentujú našu kultúru v kontexte európskej vedy, kultúry, politiky, hospodárstva a vzdelávania. Špecifickým výskumným problémom projektu je digitalizácia, atraktívna prezentácia a inovatívne sprístupnenie týchto dokumentov širokej odbornej a laickej verejnosti.

Na účely aplikovaného výskumu projektu Skriptor bola použitá platforma *Transkribus Expert Client*. V prvej polovici obdobia trvania projektu boli v expert klientovi na trénuvanie modelov a automatickú transkripciu dokumentov k dispozícii dve technológie:

<sup>1</sup> TranScriptorium, <https://cordis.europa.eu/project/id/600707>

<sup>2</sup> Recognition and Enrichment of Archival Documents, <https://cordis.europa.eu/project/id/674943>

<sup>3</sup> V lete 2023 vývojový tím Transkribus avizoval ukončenie vývoja funkcií v prostredí expert klienta a presun vývoja na online verziu Transkribus Lite /Beta. Na rozvoj online platformy a presun funkcií dostupných zatiaľ len v expert klientovi boli nasmerované všetky vývojové aktivity. Napriek tomu ešte v januári 2024 mali používatelia expert klienta dostupný upgrade na novšiu verziu 1.27.0. Definitívne ukončenie platformy Transkribus Expert Client bolo napokon naplánované na jún 2024. Pre nových používateľov už nie je k dispozícii.

<sup>4</sup> <https://readcoop.eu/readsearch/>

<sup>5</sup> <https://www.transkribus.org/sites#features>

- HTR+ vyvinutá University of Rostock,
- PyLaia vyvinutá Universitat Politècnica de València.

PyLaia je zdokonalením pôvodnej HTR+ technológie a používateľom umožňuje ľubovoľné nastavenie rôznych parametrov sieťovej štruktúry pri tréningu modelu. Podpora HTR+ bola zo strany vývojového tímu Transkribus ukončená v novembri 2022. Táto zmena sa výrazne dotkla dovedy vytrénovaných modelov riešiteľského kolektívu. V súčasnosti existuje už celý rad nástrojov, ktoré slúžia na automatické rozpoznávanie zdigitalizovaných, pôvodne tlačovaných dokumentov, medzi nimi OCR4all, eScript, Rescribe, Pero.cz, ABBYY Cloud OCR SDK, Online OCR, Kofax, Omnipage, a iné. Zámerom tohto článku nie je analýza, ani porovnávanie týchto nástrojov, ale zhodnotenie využiteľnosti softvéru Transkribus s možnosťami a funkcionalitami, ktoré aktuálne ponúka. V rámci aplikovaného výskumu projektu Skriptor sa riešiteľky zamerali na automatickú transkripciu unikátnych tlačovaných dokumentov v jazykoch nášho kultúrneho regiónu, medzi ktoré patrí slovenčina, čeština, latinčina, či maďarčina, a na možnosti ich sprístupnenia. Na tento účel boli vybrané dve tlač: učebnica J. A. Komenského *Orbis Pictus* (1798) a dielo Mateja Bela *Adparatus ad Historiam Hungariae* (1735).

Prvé štvorjazyčné vydanie Komenského učebnice *Orbis Pictus*<sup>6</sup> z dielne prešporského tlačiara Šimona Petra Webera z roku 1798 je súčasťou historického knižničného fondu Univerzitnej knižnice Univerzity Mateja Bela v Banskej Bystrici. Na účely projektu bola učebnica zaujímavá práve tým, že ide o štvorjazyčné vydanie, pričom každý jazyk je vytlačený iným typom písma – latinčina antikvou, maďarčina kurzívou, nemčina fraktúrou a čeština švabachom. Podobný výskum, avšak na rukopisnom prameni realizovali Capurro et al. (2023). Cieľom aplikovaného výskumu bolo vytrénovať použiteľný univerzálny model automatickej transkripcie pre štyri jazyky a štyri typy písma. Stoosemdesiatštyristranová publikácia s rozmermi 18,9 x 11,5 cm je vytlačená na ručne odlietavanom papieri. V textovej časti sú početne zastúpené obdĺžnikové drevoryty (Nižníková & Mikušková, 2022). Text je formálne usporiadaný do štyroch blokov. Latinský, maďarský a nemecký v stĺpcoch so zachovaním medzier tak, aby boli jednotlivé jazykové verzie súvisiace s rovnakými pojmami na jednej úrovni. Česká mutácia textu má klasickú knižnú formu cez celú stranu. Súčasťou dokumentu sú abecedné registre základných pojmov v latinčine, maďarčine, nemčine a češtine.

Práca s dokumentom na platforme Transkribus pozostáva z niekoľkých fáz. Vytvorenie konta v *Transkribus app* je bezplatné. Po prihlásení získava používateľ, ktorý chce platformu využívať zdarma, každý mesiac voľných 50 kreditov. Táto forma prihlásenia však umožňuje prístup len k obmedzenej ponuke funkcionalít. Prístup k ďalším nástrojom sa v ďalších typoch predplatného (*Scholar, Team, Organisation*)<sup>7</sup> zvyšuje priamo úmerne s poplatkami. Kredity sa aktuálne používajú na segmentáciu dokumentu a na tréning modelov, ostatné práce v aplikácii sú zatiaľ bezplatné.

### Zhotovenie digitalizátu a nahranie na platformu

Základným predpokladom na prácu s Transkribusom je zhotovenie digitalizátu dokumentu, ktorý bude slúžiť na tréning modelov a následnú transkripciu, prípadne transkripciu s použitím už vytrénovaného vlastného alebo verejného modelu. Transkribus momentálne pracuje s obrázkami vo formátoch JPEG/JPG a PNG a so súborami formátu PDF. Snímky dokumentu by mali byť vytvorené vo vývojármi odporúčanej kvalite priemerne 300 DPI. Najvyššiu kvalitu snímania umožňujú profesionálne skenery (400 – 600 DPI), ale takéto vysoké rozlíšenie je zbytočné, pretože nijako neprispieva k zlepšeniu rozpoznania a automatickej transkripcie textu. Transkribus vykazuje vynikajúce výsledky aj s obrázkami nafotenými vlastným zariadením, napr. fotoaparát, mobilný telefón, tablet a i., pokiaľ disponujú kvalitnou optikou. Na uľahčenie snímania dokumentu vo vyhovujúcej kvalite bola pre potreby projektu READ vyvinutá prenosná pracovná pomôcka *ScanTent* s vlastným podsvietením zabezpečujúcim rovnomerné osvetlenie dokumentu (nevýhodou je obmedzený rozmer plochy určenej na snímanie do veľkosti A3) a aplikácia DocScan voľne stiahnuteľná do smartfónu, ktorá umožňuje snímanie dokumentu aj bez neustáleho stláčania spúšte fotoaparátu. Na digitalizáciu tlače *Orbis Pictus* (1798) sme využili pomôcku *ScanTent* a aplikáciu *DocScan*, výsledkom bolo deväťdesiatštyri snímok dvojstrán s rozlíšením 72 DPI<sup>8</sup>.

<sup>6</sup> KOMENSKÝ, J. A. (1798). Joann. Amos Comenii Orbis pictus, in hungaricum, germanicum et slavicum translatus et hic ibive emendatus. Sumtibus & Typis Simonis Petri Weber.

<sup>7</sup> Cenová a kreditová politika na platforme sa v posledných rokoch mení, preto si treba odsledovať aktuálne platné podmienky. Kredity sa v procese automatickej transkripcie textu priebežne odpočítavajú v závislosti od parametrov nastavených pred spustením automatickej transkripcie, aj od rozsahu a typu dokumentu, ktorý sa prepisuje (historická tlač, strojopisný text, rukopisný dokument). Viac informácií o poplatkoch pre rôzne typy predplatného na <https://www.transkribus.org/plans>

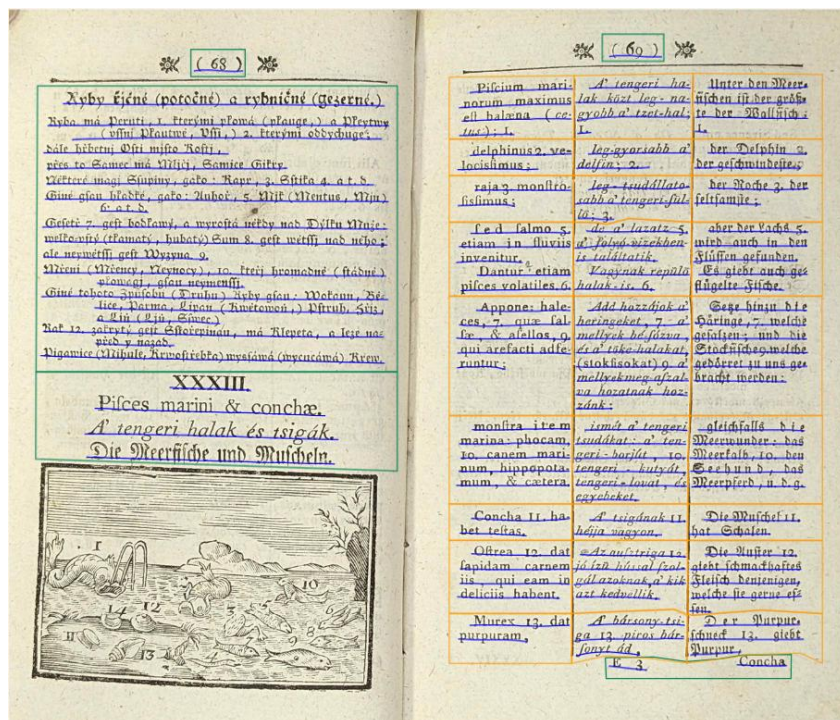
<sup>8</sup> Pri ďalších prácach na výskume sme použili vydanie z iného roku a v dvoch rôznych rozlíšeniach: *Orbis Pictus* z roku 1820 z fondu Univerzitnej knižnice UMB zdigitalizovaný na profesionálnom skeneri v Centre vedecko-technických informácií SR v Bratislave v rozlíšení 600 DPI a *Orbis Pictus* z roku 1820 stiahnutý z portálu Digitálna knižnica a digitálny archív Slovenskej národnej knižnice (DIKDA) v rozlíšení 96 DPI.

Nevyhnutnou podmienkou na vkladanie digitalizátov na platformu je vytvorenie vlastnej zbierky. Do nej sa následne pomocou funkcie **+Upload** nahrávajú zdigitalizované snímky dokumentu v niektorom z predpísaných formátov. Ak sa snímky nahrávajú jednotlivito vo formátoch JPEG a PNG, nemali by presahovať veľkosť 10 MB. V prípade nahrávania dokumentu v PDF formáte by jeho veľkosť nemala byť väčšia ako 200 MB a maximálny počet strán v súbore by nemal presiahnuť hodnotu 3000 (Bôbová et al., 2024, s. 41–45). Na tomto mieste sa žiada zdôrazniť, že všetky snímky a modely, ktoré sa nachádzajú v zbierke používateľa, sú – hoci uložené na serveroch platformy – v jeho výlučnom vlastníctve. Každý používateľ má právo rozhodnúť sa, či, do akej miery a komu umožní prístup do svojej zbierky.

### Segmentácia textu

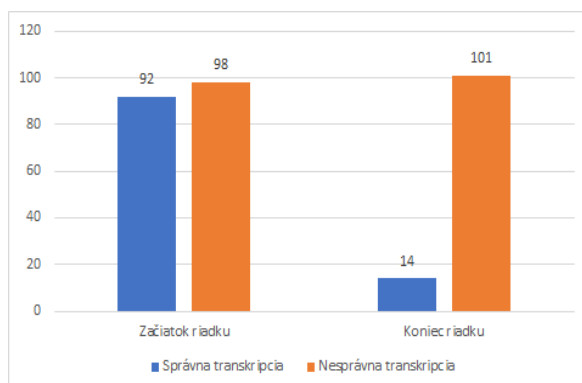
Prvým krokom práce s naimportovanými digitalizátmi je ich segmentácia, ktorou sa určuje orientácia a poradie čítania textu na snímkach. Zahŕňa vymedzenie textových blokov (Text Regions) a riadkov v nich (Baselines) a má vplyv na kvalitu trénovaného modelu a následný automatický prepis dokumentu. Segmentáciu je možné vykonať automaticky alebo manuálne. Vzhľadom na štruktúru dokumentu Orbis Pictus (kombinácia klasického knižného textu a stĺpcového členenia) sme pristúpili k manuálnemu označeniu textových rámcov zohľadňujúcich členenie textu, jeho jazykové mutácie a typy písma. V zadefinovaných textových rámcoch sme spustili automatickú detekciu riadkov. Ako vo svojom článku upozorňuje Ost (2024), automatická segmentácia nie je bezchybný proces. Najmä pri štruktúrovaných dokumentoch zvyčajne treba urobiť kontrolu a vykonať väčšie, alebo menšie korekcie v delení riadkov a správnom čítaní ich poradia (Bôbová et al., 2024, s. 72–76). Technologické zmeny na platforme Transkribus (november 2022) priniesli nové funkcionality a zefektívnenie niektorých jestvujúcich. Tie sa prejavili najmä na úrovni segmentácie dokumentu, keď pôvodnú metódu analýzy textu CITLab Advanced LA nahradila metóda Transkribus LA. Tá ponúka možnosť nastavenia viacerých parametrov delenia riadkov a ich automatického rozdeľovania, ak sa text nachádza v tesnej blízkosti hranice textového rámca, čo bol aj prípad dokumentu Orbis Pictus.

Na trénovanie modelu a automatický prepis dokumentu vplyva správne zadefinovanie základných čiar (Baselines). Tie predstavujú najdôležitejší referenčný bod na rozpoznávanie textu. Opisujú líniu (polyčiaru), ktorá sa tiahne pozdĺž celej spodnej časti riadka. Pri automatickej segmentácii riadkov je potrebné hranice začiatku a konca základnej čiary manuálne korigovať. Podľa definovanej základnej čiary sa vytvorí oblasť riadka, ktorá predstavuje polygón zahŕňajúci text dokumentu na príslušnom riadku. Korekcia základných čiar v rozsiahlejších dokumentoch môže byť časovo náročná. Menší rozsah tlače Orbis Pictus nám však umožnil detailné porovnanie správnosti prepisu chybné určených základných čiar. Chyby v určení základnej čiary boli rôzneho rozsahu (nezahŕňali celý znak, jeho časť, diakritické znamienka, rozdeľovníky a i.), pričom sa vyskytovali na začiatku i na konci riadka.



Obr. 2 Ukážka segmentácie tlače Orbis Pictus. Zdroj: Transkribus

Z nižšie uvedeného grafu vidieť, že aký vplyv má nesprávne definovaná línia základnej čiary na prepis. Softvér pri transkripcii pomocou použitia vybraného modelu dokáže v porovnaní so stanovenou základnou čiarou zachytiť širšiu oblasť riadka. V takom prípade môže mať oblasť riadka pozitívny vplyv na správny prepis grafém na hranici začiatku a konca základnej čiary.



Graf 1 Chybovosť transkripcie pri nesprávne určenej línii základnej čiary na začiatku a konci riadka

Potvrdilo sa tiež, že nesprávne nedefinovaný polygón označujúci text riadka nemá pri správne zadefinovanej základnej čiare vplyv na prepis textu. Chybovosť prepisu pri nesprávne určenej línii základnej čiary bola vyššia na konci riadka ako na jeho začiatku.

### Príprava vzorky *Ground Truth*

Po segmentácii textu sa na automatickú transkripciu dokumentu použije starostlivo vybraný model iného používateľa<sup>9</sup> alebo sa vytrénuje nový model špecifický pre dokument, ktorý je predmetom skúmania. Pred trénovaním nového modelu automatickej transkripcie je potrebné pripraviť vzorku dát, ktorá sa použije na trénovanie. V strojovom učení sa na označenie takýchto dát používa termín *Ground Truth (GT)*<sup>10</sup>. Strany možno prepísať manuálne alebo použiť vhodný model iného používateľa. Následne sa tento prepis dôsledne skontroluje, skoriguje a strany sa označia ako *Ground Truth*. Druhý spôsob významne šetrí ľudskú prácu i čas a bol použitý pri viacerých výskumoch (Capurro et al., 2023; Marsili et al., 2025). Odporúčaný počet transkribovaných slov pre tlačeneé dokumenty<sup>11</sup> je minimálne 5 000.

V prípade historických textov metodika na prácu so softvérom Transkribus odporúča ísť cestou diplomatického prepisu, ktorý presne odráža to, čo sa nachádza v pôvodnom dokumente<sup>12</sup>. V závislosti od vôle zachovať čo najvyššiu mieru autenticity dokumentu je preto dôležité vopred sa rozhodnúť, či k prepisu dokumentu budete pristupovať metódou:

- *transkripcie* – písomného vyjadrenia (vyslovovaných alebo cudzím grafickým systémom napísaných) slov a textov z hľadiska ich výslovnosti prostriedkami určitého grafického systému (Mistrík, 1993, s. 343), tzn. počíta s nahradením starších a nepoužívaných znakov modernými, alebo
- *transliterácie* – prevodu z jednej grafickej sústavy do druhej, pri ktorom každému písmenu jedného grafického systému zodpovedá vždy písmeno druhého systému (rovnaké písmeno alebo spojenie písmen), takže je možný aj jednoduchý spätný prevod do jazyka originálu (Mistrík, 1993, s. 343), tzn. prepisu znakov v ich pôvodnej podobe.

Štvorjazyčnosť dokumentu, s ktorým sme pracovali, a použité znakové sady (fonty) poskytujú zaujímavý materiál na ďalšie analýzy (grafologický a jazykový výskum), preto sme uplatnili metódu transliterácie (aj s chybami, ktoré vznikli pri tlači dokumentu). Na prepis sa v Transkribuse používajú bežné sady UNICODE, pričom softvér obsahuje aj funkciu na import nových znakov. Pri manuálnom prepise špecifických znakov, ktoré softvér neobsahoval, sme

<sup>9</sup> Použiť sa dá voľne dostupný model, ktorý najlepšie zodpovedá typu a jazyku dokumentu, použitému fontu, prípadne obdobiu jeho vzniku. Aktuálne sú na platforme dostupné aj supermodely *The Text Titan I a II*, ktoré však platforma odporúča používať na automatickú transkripciu rukopisných dokumentov.

<sup>10</sup> *Ground Truth* = overená, skutočná informácia použitá na trénovanie a testovanie modelov umelej inteligencie a na porovnanie s výsledkami týchto modelov. Je to základný referenčný bod, ktorý model porovnáva s vlastnými predikciami, aby overil svoju presnosť a naučil sa identifikovať vzory.

<sup>11</sup> Viac o príprave dát na trénovanie modelu na <https://help.transkribus.org/data-preparation>

<sup>12</sup> Viac ku konvenciám manuálneho prepisu <https://help.transkribus.org/data-preparation>

spočiatku používali náhradné znaky, ktoré boli významovo a foneticky podobné (Bôbová et al., 2023, s. 170-192). Na vyladenie transliterácie sme nakoniec využili nástroj na tvorbu kombinácií UNICODE znakov<sup>13</sup>. Každý jazyk a font obsahoval grafémy, z ktorých mnohé sa v dnešnej dobe už nepoužívajú. Ich najvyšší výskyt bol v českom jazyku:

- latinčina (antikva) → æ, œ, ě
- maďarčina (kurzíva) → ő, ű
- nemčina (fraktúra) → ß, ț, ă, Ț, ũ
- čeština (švabach) → č, ě, ě, ě, ě, ě, ě, ě, ě, ě

Spoluhlásky v češtine sa vyskytovali aj vo forme kapitálok. Všetky fonty obsahovali grafému dlhé s (ř), ktorá nahrádza okrúhle s na začiatku a uprostred slova.

Na to, aby model dosiahol čo najlepšie výsledky, je nevyhnutné, aby vzorka GT obsahovala dostatočné zastúpenie všetkých znakov. Potvrdilo sa to aj v prípade tlačených dokumentov. Pri tréovaní modelu na tlači *Orbis Pictus* (1798) boli v cvičnom súbore zastúpené znaky ř, ě, č, ě, ě, ě, ě, œ. Správnosť prepisu týchto znakov nebolo možné dostatočne overiť na overovacom súbore tréovaného modelu, pretože sa v ňom nemuseli vyskytovať. Preto sme pristúpili k analýze na celom transkribovanom texte. Ich prepis bol prevažne nesprávny, zvyčajne išlo o zámenu œ → æ, ě → e/č, ě → g, ř → r/i/ff/k, č → c/č/e.

### Nastavenie parametrov tréovania

Na tréovanie sa neodporúča vyberať presvietené snímky a strany, ktoré obsahujú ťažko čitateľné miesta, škrty, rukopisné poznámky a podobne. Vybrané reprezentatívne strany sa následne rozdelia do cvičného (*training set*) a overovacieho súboru (*validation set*), ideálne v pomere 10 : 1. Softvér ponúka možnosť manuálneho aj automatického výberu strán do tréovania. Ďalším krokom je nastavenie parametrov, ktoré môžu zlepšiť ukazovatele vytréovaného modelu. Prednastavený maximálny počet cyklov (opakovaní tréovania) je 250 a odporúčame ho ponechať. Zvyšovaním počtu cyklov sa proces tréovania predlžuje a zvyšovanie nemusí mať pozitívny vplyv na výslednú úspešnosť modelu. Voľba predčasného zastavenia tréovania (*Early Stopping*) je predvolená na 20 cyklov, čo znamená, že ak sa hodnoty modelu zlepšujú, tréovanie bude aj po dosiahnutí 20 cyklov pokračovať. Ak však už hodnoty nevykazujú zlepšenie, tréovanie sa automaticky zastaví a vyhodnotí. Funkcia obráteného textu (*Reverse Text*) sa používa v prípadoch, keď je smer písania na snímke opačný ako pri prepise (napr. originál je napísaný sprava doľava a prepísaný text smeruje zľava doprava). Pri tréovaní je možné predvolene použiť polygóny, ktoré boli vygenerované pri segmentácii dokumentu (*Use existing line polygons for training*). V prípade zrušenia tejto voľby sa pri tréovaní modelu vytvoria nové. Voľba tréovania skratiek (*Train Abbrevs with expansion*) pomáha dosiahnuť lepšie výsledky pri rozpoznávaní skratiek. Ak sa na stranách vybraných do vzorky *Ground Truth* nachádzajú slová označené tagmi Nejasný (*unclear*) alebo Medzera (*gap*), z procesu tréovania ich možno vylúčiť voľbou *Omit lines by tag*. Vynechá sa tak nielen označené slovo, ale aj celý riadok, keďže tréovanie prebieha na úrovni riadkov.

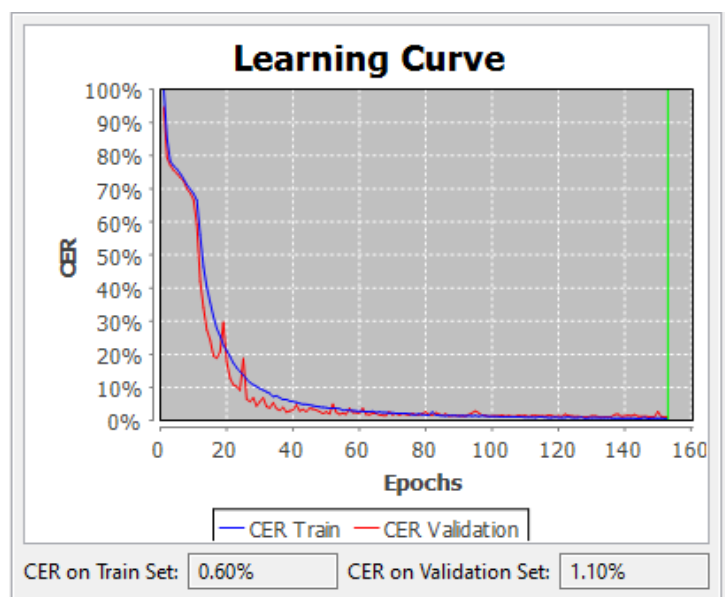
	Character Error Rate	Training Data
<b>Printed text</b>	0,5-2%	~ 5.000 words / 25 pages
<b>Single hand</b> - simple writing	2-4%	10.000+ words / 50+ pages
<b>Several hands</b> – all seen during training	4-6%	10.000+ words per hand / 150+ pages
<b>Many hands</b> from same period and region – not all seen during training	6-8%	100.000+ words / 500+ pages

Obr. 3 Odporúčaný počet prepísaných slov vo vzorke GT podľa typu dokumentu a akceptovateľná miera chybovosti znakov po vytréovaní modelu

<sup>13</sup> <https://onlineunicodetools.com/add-combining-characters>

### Hodnotenie úspešnosti modelu

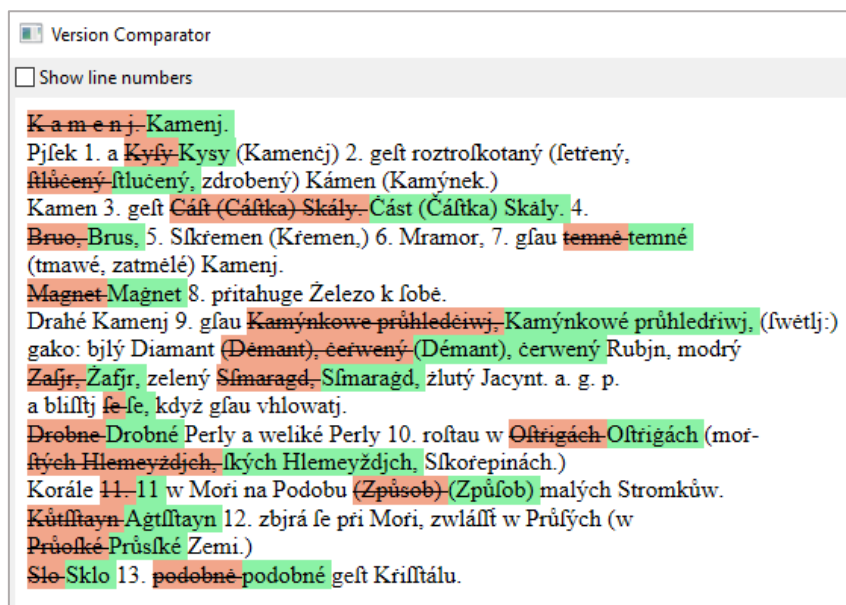
Po vytrénovaní modelu expert klient ponúkne vyhodnotenie v podobe grafu a percentuálneho vyjadrenia chybovosti znakov (os y)<sup>14</sup> v automaticky prepísanom texte. Výsledok tréningu CER on Train Set 0,60 % a CER on Validation Set 1,10 % znamená, že v cvičnom súbore bolo bezchybne určených až 99,40 % znakov a v overovacom súbore 98,90 %. Predchádzajúce výskumy v rámci projektu READ ukazujú, že hodnoty CER automaticky prepísaných rukopisných strán do 10 % sa považujú za uspokojivé, v najlepších prípadoch sa pohybujú okolo 5 %. Vynikajúce výsledky modelov tréňovaných na tlačенých stranách dosahujú chybovosť približne 1 – 2 % (Muehlberger et al., 2019), preto môžeme výsledky modelu na obr. 4 považovať za veľmi uspokojivé. Druhá časť vyhodnotenia tréningu modelu obsahuje grafické zobrazenie procesu tréningu, tzv. krivku učenia sa (*Learning Curve*). Po spustení tréningového procesu sa model opakovane trénuje a overuje. V prípade modelu na obr. 4 sme ponechali nastavenie na 250 cyklov (*Epochs*), pričom model sa prestal zlepšovať pri 153. cykle (os x) a tréning sa automaticky zastavilo.



Obr. 4 Výsledok tréningu modelu – grafické znázornenie. Zdroj: Transkribus

<sup>14</sup> Miera chybovosti znakov CER (*Character Error Rate*) sa udáva v percentách. Krivka sa vždy začína na 100 % a tým, ako sa model trénuje a zlepšuje, postupne klesá. Miera chybovosti znakov porovnáva celkový počet znakov (n) vrátane medzier s minimálnym počtom vložením (i), nahradením (s) a vymazaním (d) znakov potrebných na dosiahnutie rovnakého výsledku ako vo vzorke *Ground Truth*. Vzorec na výpočet miery chybovosti znakov:  $CER = [(i + s + d) / n] * 100$

Úspešnosť modelu možno vyhodnotiť aj na úrovni jednotlivých strán. Komparáciou textových verzií pomocou nástroja *Compare Text Versions* získame podrobný prehľad toho, čo model prepísal správne a kde v porovnaní s verziou *Ground Truth* urobil chybu.



Obr. 5 Nesprávne prepísané slová v overovacom súbore (červené) porovnané so správnym prepisom vo vzorke *Ground Truth* (zelené). Zdroj: *Transkribus*

Nástroj *Compare* umožňuje v percentách vyjadriť mieru chybovosti na úrovni slov *WER* (*Word Error Rate*). Tá sa môže na jednotlivých stranách výrazne líšiť, treba si však uvedomiť, že zväčša ide o chyby spojené s interpunkciou (chýbajúca alebo nadbytočná bodka, čiarka, dvojbodka a pod.) a s diakritikou (krátka samohláska namiesto dlhej, spoluhláska namiesto spoluhlásky s mäkčeňom a pod.), ktoré nemajú žiadny alebo minimálny vplyv na zrozumiteľnosť textu.

V texte získanom po automatickej transkripcii a nevyhnutných korektúrach je možné zvýrazniť významné údaje a označiť jednotlivé časti štruktúry dokumentu. Platforma *Transkribus* ponúka nástroje na prácu s textom prostredníctvom tagovania (značkovania). K dispozícii sú dva druhy tagov. Textové tagy zdôrazňujú významné informácie v texte, napr. osobné mená, geografické názvy, dátumy, skratky a pod. Štrukturálne tagy definujú jednotlivé časti štruktúry dokumentu ako nadpisy, odstavce, ilustrácie, čísla strán, marginálne poznámky a i. (Bôbová et al., 2024, s. 113-122). Okrem exportu snímok a prepísaného textu je z platformy možné exportovať aj textové tagy využiteľné na tvorbu registrov.

### Tréovanie vlastných modelov

Tréovanie samostatných modelov pre jednotlivé fonty skúmanej tlače pomocou technológie *HTR+* v prvej fáze projektu neprinieslo výsledky, ktoré by zodpovedali vývojárami deklarovanej nízkej chybovosti pri tlačených dokumentoch. V ukazovateľoch *CER* a *WER* na overovacom súbore a aj vizuálne (porovnávaním textových verzií) v spoločných modeloch pre všetky fonty modely vykazovali vysokú chybovosť vo fonte švabach. Už spomínané technologické zmeny v roku 2022 priniesli odstavenie dosiaľ používanej technológie *HTR+* a nasadenie nového stroja *PyLaia*. Z tohto dôvodu sme s časovým odstupom 1,5 roka samostatné modely pretréovali novou technológiou. Tá mala spočiatku výrazné nedostatky, ktoré sa prejavovali najmä vyššou chybovosťou. Ako naprotektickejší font sa na základe ukazovateľov *CER* a *WER* na overovacom súbore ukázala kurzíva a s fontom švabach si technológia neporadila ani po opakovanom spustení tréovania.

Zdokonaľovanie technológie na automatický prepis historických dokumentov, nové funkcionality a nástroje stále priamejšie smerujú k snahe vytvárať univerzálne modely, ktoré budú schopné automaticky prepisovať typovo analogické rukopisné, tlačené či strojové dokumenty. Súbežne so samostatnými modelmi sme sa preto rozhodli tréovať aj spoločné modely, ktoré nezohľadňovali typ fontu, ani jazyk. Dokument, ktorý bol predmetom nášho výskumu, sme nakoniec automaticky transkribovali dovedy najlepším modelom *Model 10* a podrobili ho redakčným úpravám. Relatívne malý rozsah – 184 strán textu formátu *A5* – nám umožnil detailne analyzovať

chybovosť až na úroveň diakritiky a interpunkcie. V zmysle metodiky pre automatizované vyhodnocovanie úspešnosti modelu pomocou ukazovateľa CER platí, že aj malá chyba pri prepise je štatisticky plnohodnotnou chybou. To znamená, že chýbajúcu čiarku, „u“ namiesto „v“, „á“ namiesto „ä“, medzeru navyše alebo veľké písmeno namiesto malého softvér vyhodnotí ako rovnocenné chyby. Rovnaký princíp sme použili v procese redakčnej úpravy pri typológii týchto chýb a ich následnom štatistickom vyhodnotení (tab. 1) Položka *Nesprávne prepísaná graféma/interpunkcia* znamená nesprávny prepis samotnej grafémy, najčastejšie „v“ namiesto „u“, „l“ alebo „k“ namiesto „f“, ale aj nesprávny prepis interpunkčného znamienka (dĺžeň, dvojbodka, trojbodka, apostrof) na inak správne prepísanej graféme, najčastejšie „ö“ namiesto „ó“ alebo „ô“ či „ú“ namiesto „ü“ alebo „ů“.

Najvyššia chybovosť na úrovni strany sa vyskytla vo vecnom registri v maďarčine tlačenej kurzívou. Viacerí riešitelia projektu Skriptor konštatovali problematický prepis číslic v rukopisných textoch (Nagy, 2021; Tomeček, 2025), chybovosť v prepise číslic sa však v tlači *Orbis Pictus /1798* vyskytovala ojedinele s výnimkou už spomínaných registrov. Tie naopak vykazovali vysokú chybovosť vo viacerých aspektoch. Dôvodom bol spôsob rozloženia textu, malé medzery medzi oddelenými stĺpcami a vysoký výskyt interpunkčných znamienok (bodky), ktoré boli použité na oddelenie názvov kapitol a stránkovania. Ako najproblematickejšie písmo sa dovedy dlhodobo ukazoval švabach, ktorý zvyšoval chybovosť čítania znakov aj v spoločných modeloch. Dokazujú to hodnoty namerané pri samostatnom tréningu fontov pomocou technológie HTR+. Ukazovatele získané po pretréningu tých istých súborov technológiou PyLaia však v miere chybovosti „favorizovali“ kurzívu. Podrobná analýza chybovosti v tabuľke toto zistenie potvrdzuje.

Typ chyby	antikva	kurzíva	fraktúra	švabach
nesprávne prepísaná graféma / interpunkcia	430	541	386	548
chýbajúca graféma / interpunkcia / medzera	50	66	61	33
graféma / interpunkcia / medzera, kde nemá byť	39	48	43	48
nesprávne prepísaná diakritika (bodka, čiarka, bodkočiarka, dvojbodka, pomlčka, otáznik)	128	48	51	24
chýbajúca diakritika (bodka, čiarka, bodkočiarka, dvojbodka, pomlčka, zátvorka)	65	27	94	20
diakritika, kde nemá byť (pomlčka, bodka, dvojbodka)	70	50	45	43
<i>Správnosť čítania špeciálnych znakov s nízkym výskytom v tréningovom súbore</i>				
ligatúra ae	1			
ligatúra oe	15			
ligatúra tz			4	
<b>Celkový počet chýb</b>	<b>798</b>	<b>780</b>	<b>684</b>	<b>716</b>

Tab. 1 Typy a počet chýb po automatickej transkripcii textu Modelom 10

Domnievame sa, že najčastejším dôvodom nesprávneho prepisu bola nižšia kvalita samotnej tlače (najmä slabo vytlačené písmo), presvit z opačnej strany listu, otláčenie tlačiarskej farby z predchádzajúcej strany a nečistota na papieri.

Na zvyšovanie kvality modelov sme v rôznych fázach projektu používali tri základné postupy a ich kombinácie:

- oprava manuálnej transkripcie – po vytréningu modelu a analýze chýb na úrovni znakov a slov zistíme, že softvér upozornil aj na naše vlastné chyby. V tejto fáze je teda priestor na opätovnú kontrolu prepisu vzorky GT a opravu nezrovnalostí.
- zvýšenie počtu slov v cvičnom súbore – v niektorých prípadoch je odporúčaný počet slov zahrnutých do tréningu modelu automatickej transkripcie nedostatočný a tu je na mieste zvýšiť počet strán. Vo väčšine prípadov platí pravidlo, že čím vyšší je výskyt znaku v cvičnom súbore, tým lepšie sa ho stroj naučí rozpoznávať. Touto metódou sa dá znížiť miera chybovosti znakov v overovacom súbore aj o niekoľko percentuálnych bodov.
- použitie základného modelu (*base model*) – táto metóda vychádza z predpokladu, že základný model si pamätá, čo sa naučil, a teda každé ďalšie tréningovanie ho (teoreticky) zlepšuje. Pri tomto postupe možno ako *base model* použiť najlepší vlastný model alebo zodpovedajúci model iného používateľa.

Z hľadiska ďalšej využiteľnosti ktoréhokolvek modelu je dôležité, aby bol proces jeho vzniku transparentne popísaný. Na tento účel platforma Transkribus ponúka možnosť vytvoriť redakčné vyhlásenie (*Editorial Declaration*). Táto funkcia obsahuje súbor preddefinovaných položiek s možnosťou tvorby vlastných popisov. V redakčnom vyhlásení sa uvádza aj zoznam znakov UNICODE použitých na prepis špeciálnych znakov, ktoré dokument obsahuje. Tento krok je nevyhnutný najmä v prípade transliterácie historických rukopisov a historických

tlačí, pretože ďalším používateľom uľahčuje interpretáciu vytrénovaného modelu a rozhodovanie o jeho použití na prepis iných obsahovo, či formálne príbuzných dokumentov, alebo ako *base modelu*. Je vhodné, ak je redakčné vyhlásenie dostupné aj v anglickom jazyku, keďže komunita používateľov platformy Transkribus je medzinárodná a model môže mať potenciál byť aplikovateľný aj na iné tlače nielen územne, jazykovo a provenienčne slovacikálneho charakteru.

Jedným z najočakávanejších okamihov našej niekoľkoročnej práce bolo použitie najlepšieho vytrénovaného modelu na automatickú transkripciu formálne, obsahovo a typovo porovnateľného dokumentu s cieľom overiť použiteľnosť, resp. posúdiť mieru univerzálnosti tohto modelu. Na tento účel sme zvolili Model 11 s ukazovateľmi CER 0,60 % na cvičnom a CER 1,00 % na overovacom súbore. Modelom sme automaticky prepísali vybrané strany učebnice J. A. Komenského *Orbis Pictus* vydanom v tej istej tlačiarni v roku 1820. K dispozícii sme mali dve vzorky digitalizátov v rôznej kvalite. Prvá vzorka zdigitalizovaná v kvalite 600 DPI v Centre vedecko-technických informácií SR v Bratislave bola v minulosti reštaurovaná, preto v dôsledku pevnosti väzby nebolo možné nasnímať všetky strany bez viditeľného ohybu textu na vnútornom okraji strán. Druhá vzorka zdigitalizovaná v kvalite 96 DPI (dostupná na portáli DIKDA) obsahuje snímku každej strany jednotlivo a bez ohybu v dôsledku väzby. Dve rôzne vzorky nám tak okrem overenia funkčnosti vytrénovaného modelu na podobnej historickej tlači umožnili preveriť aj vplyv kvality digitalizátu na automatickú transkripciu. Na overenie oboch hypotéz sme si vybrali dve strany obsahujúce ten istý text skladajúci sa zo všetkých štyroch fontov – antikva, kurzíva, fraktúra a švabach. Funkčnosť modelu sa nám na oboch vzorkách potvrdila. Model 11 použitý na pôvodnej tlači *Orbis Pictus* (1798) vykazoval chybovosť na úrovni CER 2,55 % a WER 7,50 %. Na stranách *Orbis Pictus* z roku 1820 Model 11 vykazoval vyššiu chybovosť, ale na tú zrejme mala vplyv aj viditeľne horšia kvalita tlače originálneho diela. Vzorka zdigitalizovaná v kvalite 96 DPI vykazovala chybovosť CER 3,80 % a WER 12,61 %, vzorka zdigitalizovaná v kvalite 600 DPI vykazovala chybovosť CER 4,98 % a WER 16,33 %. Na základe výsledkov možno konštatovať, že kvalita digitalizátu (počet DPI) nemá dramatický vplyv na presnosť prepisu. Nepotvrdili sa tak závery, ktoré vo svojom článku publikovali Harish & Raghavendra (2024). Väčší počet chýb v prípade kvalitnejšieho digitalizátu má za dôsledok mierny ohyb strán v mieste väzby. Vyššia hodnota WER nespochybňuje efektívnosť a účelnosť automatickej transkripcie a ide o nezávažné chyby v diakritike (chýbajúci/nadbytočný dĺžeň a i.) a interpunkcii (rozdeľovník na konci slova, čiarka namiesto bodkočiarky a pod.), ktoré nemajú žiadny vplyv na zrozumiteľnosť textu.

Model 11 sme v ďalšom kroku použili na prepis vybraných strán historického diela Mateja Bela *Adparatus ad Historiam Hungariae*<sup>15</sup> z roku 1735. Ide o sedemstodvadsaťosemstranové dielo polyhistora Mateja Bela napísané v latinskom jazyku, ktoré je tiež súčasťou historického knižničného fondu Univerzitetnej knižnice Univerzity Mateja Bela v Banskej Bystrici. Textová časť obsahuje len dva fonty – antikvu a kurzívu. Formálne usporiadanie textu pozostáva z nadpisov, ozdobných iniciálok, marginálií a poznámok pod čiarou. Niektoré časti diela sú usporiadané do stĺpcov a tabuliek. Ozdobnú časť tvoria okrem iniciálok aj vinety a vlasy. V texte sa nachádzajú špecifické ligatúry (ct, st, Æ, &, œ), grafémy (ê, è, ù, â, à, á, ù, û, ÿ, W) a znaky (§, =, ¿, „). Historický dokument bol zdigitalizovaný na profesionálnom skeneri v Centre vedecko-technických informácií SR v Bratislave v rozlíšení 600 DPI.

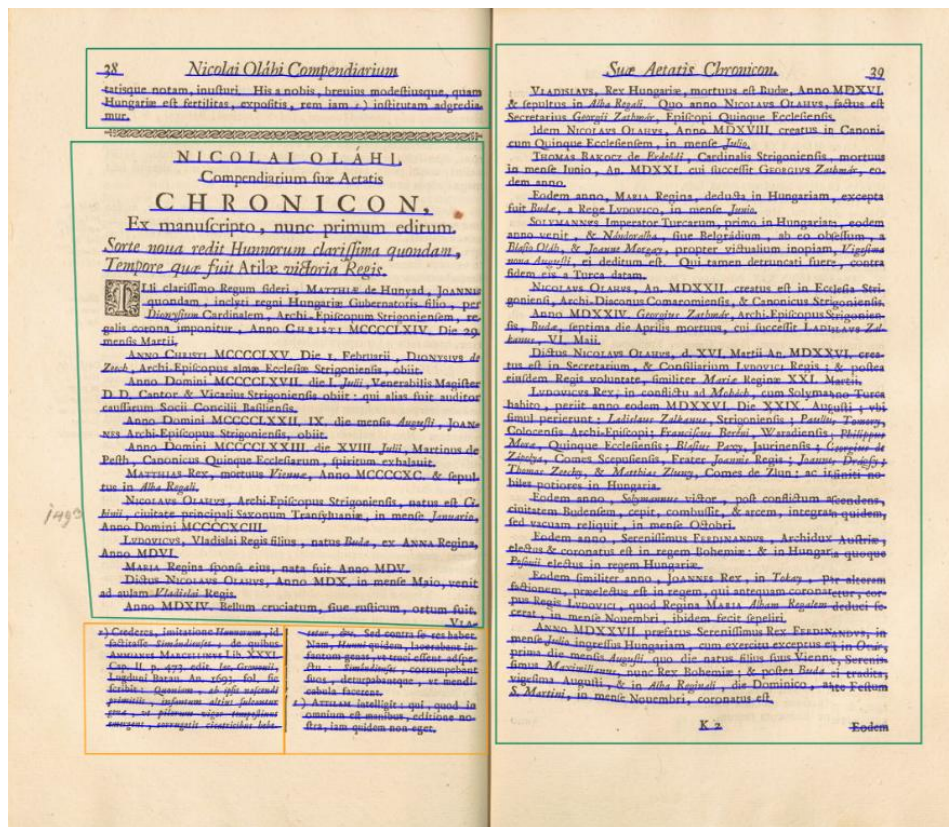
Vzhľadom na príbuzné fonty písma, geografickú blízkosť tlačiarní a obdobie vydania diel *Orbis Pictus* (1798) a *Adparatus ad Historiam Hungariae* (1735) autorky predpokladali dobré výsledky automatickej transkripcie s použitím Modelu 11. Výsledný prepis však tieto očakávania nepotvrdil. Za dôvody vysokej chybovosti na *Adparate* možno považovať:

- nové typy špecifických grafém, na ktorých nebol Model 11 trébovaný,
- vyšší výskyt kapitálok, ktoré dielo *Orbis Pictus* neobsahovalo v dostatočnej miere.

Významným faktorom zlyhania Modelu 11 sa javí aj kvalita papiera a použitie viacerých druhov antikvy. Na zlepšenie prepisu sme sa preto rozhodli vytrébovať nový model. Príprava vzorky GT vychádzala z opravy chýb na stranách prepísaných Modelom 11<sup>16</sup>. Následne autorky vykonali kontrolu a dodatočné opravy prepisu, čím bola vzorka GT pripravená na trébovanie špecifického modelu Bel\_Adparatus\_1735\_model 1. Model 11 bol pri trébovaní použitý ako *base model*. Nejasné znaky označené ako *unclear* boli z trébovania modelu vylúčené.

<sup>15</sup> Bel, M. (1735). *Adparatus ad Historiam Hungariae, sive collectio miscella, Monumentorum ineditorum partim, partim editorum, sed fugientium. Conquisit, in Decades partitus est, & Praefationibus, atque Notis illustravit, Matthias Bel. Typis Joannis Paulli Royer.*

<sup>16</sup> Opravu strán zabezpečili študenti Mária Kmecová a Dominik Polička z Katedry filozofie Filozofickej fakulty Univerzity Mateja Bela v Banskej Bystrici v rámci povinnej praxe v Univerzitetnej knižnici UMB, ktorým týmto ďakujeme za spoluprácu.



Obr. 6 Ukážka segmentácie tlačie Adparatus. Zdroj: Transkribus

Výsledky vytrénovaného modelu boli uspokojivé, poukázali však na niekoľko chýb v príprave vzorky GT v podobe nesprávneho prepisu<sup>17</sup>. Po ich odstránení bol vytrénovaný Bel\_Adparatus\_1735\_model 2. V modeli Bel\_Adparatus\_1735\_model 3 sme opätovne otestovali teóriu o vylúčení existujúcich polygónov. Tretí model na tlačí *Adparatus ad Historiam Hungariae* nakoniec dosiahol chybovosť len na úrovni 0,93 %.

Univerzálnosť tohto modelu sme recipročne overili na tlačí *Orbis Pictus* z roku 1820. Podobne ako v prípade overovania funkčnosti Modelu 11 sme zvolili tie isté strany obsahujúce všetky typy fontov. Výsledky boli neuspokojivé, dosahovali chybovosť na úrovni CER 23,99 % a WER 53,55 % v prípade digitalizátov v kvalite 96 DPI a CER 26,99 % a WER 57,5 % na digitalizátoch v kvalite 600 DPI. Vyššia chybovosť bola spôsobená zlou čitateľnosťou textu v blízkosti väzby dokumentu. Model Bel\_Adparatus\_1735\_model 3 mal problém s čítaním špecifických grafém s dnes už nepoužívanou interpunkciou (ö, ü, ç, è, ñ, ÿ), ktorá sa v tlačí *Adparatus ad Historiam Hungariae* nevyskytovala. Keďže Bel\_Adparatus\_1735\_model 3 bol vytrénovaný na fontoch antikva a kurzíva, predpokladali sme dobrý výsledok transkripcie v týchto fontoch. V kurzíve však bola chybovosť na zvlášť vysokej úrovni a to nielen pri prepise špecifických grafém. Chybovosť vo fontoch fraktúra a švabach nebola až taká prekvapivá. Použitím Modelu 11 ako *base modelu* sa teda nevytvoril univerzálny model vhodný na prepis oboch tlačí. Prehľad a popis všetkých vytrénovaných modelov sa nachádza v prílohe tohto článku.

V roku 2023 platforma Transkribus už obsahovala také množstvo dát, že vývojárom poslúžili ako základ na vytrénovanie prvých takzvaných supermodelov. Tie sú založené na výkonnejšej technológii a trénované na veľkom objeme rôznorodých dokumentov. Supermodely úľahčujú využívanie platformy, pretože (teoreticky) umožňujú transkripciu veľkých zbierok obsahujúcich rôzne jazyky a viac typov písma bez nevyhnutnosti vytvárania vlastného modelu (Park, 2025). Začiatkom roka 2024 na výročnej konferencii Transkribus komunity TUC24 avizoval vývojový tím aplikovanie veľkých jazykových modelov (Large Language Models – LLM) do platformy. Postupné zdokonaľovanie nástrojov umelej inteligencie využívajúcich veľké jazykové modely umožnilo dosahovať veľmi dobré výsledky prepisu rukopisných dokumentov aj bez špecificky vytrénovaných modelov (Humphries et al., 2024). Výsledky testovania supermodelu *Text Titan 1 ter* vykazujú jeho vysokú univerzálnosť a zároveň vyššiu

<sup>17</sup> Skúsenosti s poukázaním na vlastné chyby pri príprave súboru GT sa potvrdili pri trénovaní modelov pre tlačie *Orbis Pictus* aj *Adparatus ad Historiam Hungariae*. Rovnaké skúsenosti potvrdili aj ďalší autori, napr. Griffiths (2024), Raghallaigh et al. (2023).

presnosť v porovnaní s poprednými LLM nástrojmi ako ChatGPT, Gemini, Claude a Mistral (Noble, 2025). Aj výsledky práce na úlohách projektu Skriptor a výsledné modely viedli k vytvoreniu dvoch agregovaných modelov – *Slovak Supermodel M1 (SSM1)* pre rukopisné texty a *Slovak Supermodel print&typewriter1 (SSPT1)* pre historické tlače a strojopisný text (Bôbová et al., 2024, s. 100-103). Dokumenty a dáta spracované riešiteľmi boli s veľkou pravdepodobnosťou použité aj na tréningovanie jedného z prvých supermodelov *Transkribus Print M1*, ktorý obsahuje aj tlače v slovenskom a českom jazyku.

*Supermodel Transkribus Print M1* sme otestovali na automatickom prepise tlače *Orbis Pictus* (1820). Výsledky transkripcie boli na vysokej úrovni. Prepis bol čitateľný vo všetkých štyroch fontoch, chyby sa vyskytovali len v prípade diakritiky a interpunkcie s minimom chýb v prepise grafém a číslic. Transkripcia však nezohľadňovala špecifické, dnes už nepoužívané grafémy použité predovšetkým vo fontoch kurzíva a švabach (napr. ō, ů, č, è, ñ, ÿ). Okrem toho je zaujímavé, že supermodel bol pravdepodobne vytrénovaný na vzorke, ktorá zohľadňovala fonetické čítanie jednotlivých slov a grafému g automaticky prepisovala na grafému j (gsou → jsou, gako → jako, gezerné → jezerné, plauge → plauje a i.). Pre používateľov platformy Transkribus, ktorí nemajú záujem o prepis formou transliterácie, je teda použitie tohto modelu jasnou voľbou. Pre potreby nášho projektu by prepis dokumentu *Orbis Pictus* v kvalite 96 DPI za použitia tohto modelu a následnej opravy špecifickej interpunkcie vykazoval chybovosť CER 7,78 % a WER 21,97 %. Ak sa používateľ platformy uspokojí s transkripciou, chybovosť modelu je CER 3,06 % a WER 10,56 %. Do tejto chybovosti je započítaný aj prepis grafémy g, ktorú model prepisoval ako j.

## Závery

Štvorročná práca s platformou Transkribus nám dáva priestor na formulovanie týchto záverov:

- manuálna segmentácia celého dokumentu a manuálny prepis (alebo poloautomatický prepis s následnou korektúrou) vybraných strán, ktoré sa použijú ako vzorka *Ground Truth* na vytrénovanie nového modelu, sú najzdĺhavejšou časťou prípravných prác predchádzajúcich automatickej transkripcii celého dokumentu. S pribúdajúcimi skúsenosťami s prácou na platforme sa však časová dotácia na tieto činnosti znižuje v prospech výskumníka.
- tréningu modelu výrazne pomáha vynechanie problematických strán (preexponované snímky, nečistoty na papieri, machule) a ťažko čitateľných miest v texte (škrťance, zhustené písmo na konci riadkov, zahnuté písmo v strede pri väzbe, slabá tlač). Vidieť to najmä na CER overovacích súboru, kde rozdiely na jednotlivých stranách dosiahli viac ako 1 %, a na ukazovateli WER, a rozdiel chybovosti na dvoch overovacích stranách toho istého modelu dosiahol aj viac ako 10 %. K rovnakému záveru dospeli aj riešitelia, ktorí skúmali rukopisné dokumenty (Nagy, 2021; Bôbová, 2023).
- navýšenie strán vo vzorke *Ground Truth* oproti odporúčanému rozsahu nemá výrazný efekt na zlepšenie modelu. Platí to pre rukopisné aj tlačené dokumenty. Dôkazom sú modely 10 a 11, ktoré sme trénovali s odstupom takmer jedného roka, pričom Model 11 bol trénovaný na celej tlači s rozsahom 184 strán – na ukazovateľoch vidieť len minimálne zlepšenie. Navyšovanie však má zmysel v prípade, že sa v texte (a teda aj v cvičnom súbore) v menšom počte nachádzajú špeciálne znaky a ligatúry. Znak by sa mal v texte vyskytovať v dostatočnom počte (aspoň 50x pri tlačiach a 500x pri rukopisoch), aby sa ho stroj naučil (spoľahlivo) čítať.
- vo všetkých prípadoch tréningovania modelu sa dosiahli lepšie výsledky s predvoľbou tréningovania bez existujúcich polygónov. Tento záver platí aj pre tréningovanie modelov na tlači *Adparatus ad Historiam Hungariae*. Na rovnakých cvičných a overovacích stranách dosahoval *Bel\_Adparatus\_1735\_model 3* trénovaný bez použitia existujúcich polygónov lepšie výsledky ako *Bel\_Adparatus\_1798\_model 2* (na úrovni chybovosti znakov CER ide o rozdiel 0,30 %).
- pri tréningu modelov pomocou PyLaia je automaticky prednastavená voľba *Early stopping* na 20 cyklov. Odporúčame toto nastavenie ponechať, pretože softvér pokračuje v tréningu modelu, ak sa výsledky nezlepšujú. Dokazuje to nastavenie parametrov pri tréningu vybraných modelov. Model *Bel\_Adparatus\_1735\_model 2* nechal Transkribus trénovať 250 cyklov aj napriek tomu, že krivka výsledkov na úrovni chybovosti znakov (CER) sa približne od 65 cyklu nezlepšovala. V prípade modelu *Bel\_Adparatus\_1735\_model 3* sa krivka v priebehu 250 cyklov tréningovania menila. Model *Bel\_Adparatus\_1735\_model 3* dosahoval o 0,30 % lepšie výsledky v CER. Pri modeloch trénovaných na tlači *Orbis Pictus* (1798) sa v jednom prípade zastavilo tréningovanie po 153 cykloch (Model 10), v ďalšom po 92 cykloch (Model 11). Rozdiel na úrovni CER bol iba 0,10 % a to aj napriek tomu, že Model 11 obsahoval podstatne rozsiahlejšiu vzorku cvičných a overovacích dát.
- jednoznačne sa ukazuje, že je prínosom aplikovať *base model* pri tréningu modelov na tom istom dokumente: zodpovedajúci model iného používateľa alebo vlastný model. Vidieť to na grafoch modelov

Bel\_Adparatus\_1735\_model 1 a 2. Transkribus s base modelom vykázal na krivke učenia už pri prvom cykle nižšiu neúspešnosť (v rozmedzí 87 – 95 % CER na overovacom súbore), pričom bez použitého base modelu je neúspešnosť pri prvom cykle takmer vždy 100 %.

- zdokonalenie modelu jeho obohatením o ďalší typ písma prostredníctvom *base modelu* trénovanom na inom dokumente sa nepotvrdilo. Model Bel\_Adparatus\_1735\_model 3 vytrénovaný na tlačí *Adparatus ad Historiam Hungariae* sme v rámci overovania jeho funkčnosti použili na vzorke jednej dvojstrany tlače *Orbis Pictus* (1820). Výsledky boli neuspokojivé. Napriek tomu, že pri jeho trénovaní bol ako *base model* použitý Model 11, nový model nevykazoval dobrú čitateľnosť špecifických znakov a fontov fraktúra a švabach.
- v súvislosti s AI sa veľa hovorí a píše o jej nestabilite, náchylnosti halucinovať (vytvárať nesprávne alebo zavádzajúce výsledky) a podobne. Faktom však je, že moderné modely umelej inteligencie, najmä veľké jazykové modely (LLM), sú navrhnuté tak, aby rozpoznávali vzory v rozsiahlych súboroch údajov a generovali vierohodné odpovede na základe týchto vzorov. Z tohto dôvodu sme sa rozhodli najúspešnejší model Bel\_Adparatus\_1735\_model 3 preveriť jeho opakovaným trénovaním s tými istými parametrami. Výsledok svedčí o vysokej konzistentnosti Transkribusu – model Bel\_Adparatus\_1735\_model 3a sa v ukazovateli CER mierne zlepšil (0,27 % oproti 0,31 %), ukazovateľ WER zostal rovnaký.
- textové verzie (DOCx, PDF, Excel) získané „z obrázkov“ automatickou transkripciou sú fultextovo prehľadateľné a vďaka rozvoju ďalších nástrojov umelej inteligencie na automatizovaný preklad je možné využiť prekladače z rôznych, aj menej používaných jazykov. Takto transkribovaný a preložený dokument môže slúžiť na ďalší vedecký výskum a jeho spracovanie v podobe pramennej edície (Tomeček & Nagy, 2024).
- využitie platformy má zmysel predovšetkým pri práci s rozsiahlymi dokumentmi alebo zbierkami. V prípade dokumentov s menším počtom strán nemusí pomer času a práce na vytrénovaní nového modelu na jednej strane a výsledný automaticky transkribovaný text s väčšou či menšou mierou chybovosti na druhej strane vychádzať v prospech úsilia výskumníka. V tomto subjektívnom hodnotení vychádzame z toho, že s expert klientom sme sa od začiatku projektu učili pracovať metódou pokus – omyl, samoštúdiom vtedy relatívne obmedzeného počtu metodík dostupných na webových stránkach projektu READ a občasnou výmenou skúseností medzi členmi riešiteľského kolektívu. Motivačne nepôsobilo ani to, že uprostred projektu vývojový tím odstavil dosiaľ používanú technológiu a dovtedajšie výsledky trénovania modelov v podstate zostali bezcenné, resp. nepoužiteľné, ak sme chceli porovnávať výsledky s modelmi, ktoré sme v nasledujúcom období trénovali novou technológiou. Napriek tomu rýchly vývoj platformy postupne prinášal zlepšenia, ktoré by nám počiatočné práce na segmentácii a príprave GT vzorky výrazne pomohli a ušetrili množstvo času. Zároveň však treba zdôrazniť, že získané skúsenosti priniesli benefit vo forme vynikajúceho modelu na prepis takmer osemstostranového diela Mateja Bela a jeho príprava a trénovanie si vyžiadali odhadom desatinu času a úsilia v porovnaní s prácami na Komenského učebnici *Orbis Pictus*.

Po skúsenostiach s platformou Transkribus a na základe vyvedených záverov odporúčame jej využívanie na automatickú transkripciu historických dokumentov. Kontinuálny vývoj, zlepšovanie a pridávanie funkcionalít, ako aj rastúci záujem používateľov sú dôkazom, že medzi nástrojmi na automatický prepis starých rukopisov a tlačí má svoje významné miesto. V súčasnosti je Transkribus najčastejšie používaným nástrojom HTR v oblasti kultúrneho dedičstva (Nockels et al., 2022). Široké využitie platformy sa ukazuje nielen v oblasti trénovania modelov, ale aj rozpoznávania typov rukopisov (Cuéllar & Boadas, 2025), či trénovania modelov na automatické rozpoznávanie entít nachádzajúcich sa v dobových textoch (Sánchez-Salido et al., 2023). Platforma dnes poskytuje používateľom aj širokú podporu: organizuje webináre, semináre a workshopy na rôzne témy nielen v angličtine, ale aj v iných jazykoch<sup>18</sup>, k dispozícii sú manuály a videá na prácu so softvérom, každoročne sa koná výročná konferencia *Transkribus User Conference*. V roku 2020 získal Transkribus cenu *Horizon Impact Award*, ktorá sa udeľuje projektom so spoločenským dopadom na celú Európu i za jej hranicami. V januári 2023<sup>19</sup> komunitu tvorilo viac ako 100 000 individuálnych a inštitucionálnych používateľov (vrátane národných knižníc, akademických knižníc, archívov, múzeí a výskumných inštitúcií<sup>20</sup>) a na serveroch bolo uložených 43 miliónov spracovaných digitalizovaných strán rukopisných a tlačených dokumentov. V októbri 2025 bol už počet používateľov trojnásobný, počet strán presiahol 50 miliónov a počet vytrénovaných HTR AI modelov prekročil hranicu 20 000.

<sup>18</sup> <https://www.transkribus.org/events>

<sup>19</sup> <https://edinburgh-innovations.ed.ac.uk/case-studies/transforming-scholarship-in-the-archives>

<sup>20</sup> <https://readcoop.org/members>

## Príloha

Model	Dátum vytvorenia	ID modelu	Technológia	Cvičný súbor			Overovací súbor					Počet cyklov
				počet riadkov	počet slov	CER	počet riadkov	počet slov	CER	WER 1	WER 2	
Print_antikva_1	19.8.2022	44152	CITlab HTR+	383	917	0,16 %	77	198	2,64 %	8,04 %	5,20 %	50
Print_kurziva_1	19.8.2022	44148	CITlab HTR+	373	1 018	0,13 %	87	236	2,11 %	6,49 %	9,38 %	50
Print_fraktura_1	21.8.2022	44173	CITlab HTR+	370	1 025	0,08 %	88	243	1,23 %	1,94 %	7,18 %	50
Print_schwabach_1	19.8.2022	44150	CITlab HTR+	151	908	0,04 %	31	170	4,78 %	8,74 %	18,71 %	50
Print_antikva_2	12.3.2024	60660	PyLaia	383	917	1,10 %	77	198	2,20 %	8,09 %	6,25 %	201
Print_kurziva_2	12.3.2024	60662	PyLaia	373	1 018	1,20 %	87	236	4,30 %	9,73 %	14,80 %	201
Print_fraktura_2	12.3.2024	60664	PyLaia	370	1 025	0,10 %	88	243	3,00 %	6,45 %	13,88 %	201
Print_schwabach_2	13.3.2024	60669	PyLaia	151	908	99,80 %	31	170	100 %	28,49 %	33,18 %	201

Tab. 2 Prehľad samostatných modelov vytrénovaných v rokoch 2022 – 2024

Model	Dátum vytvorenia	ID modelu	Technológia	Cvičný súbor			Overovací súbor					Počet cyklov
				počet riadkov	počet slov	CER	počet riadkov	počet slov	CER	WER 1	WER 2	
Model 1	8.8.2022	43995	CITlab HTR+	653	2 047	0,29 %	202	605	2,33 %	6,42 %	8,56 %	50
Model 2	8.8.2022	44001	CITlab HTR+	601	1 878	0,25 %	254	773	3,09 %	6,65 %	9,74 %	50
Model 3	18.8.2022	44136	CITlab HTR+	653	2 047	0,28 %	202	605	2,23 %	7,74 %	7,31 %	50
Model 4	18.8.2022	44117	CITlab HTR+	601	1 878	0,25 %	254	773	2,69 %	6,65 %	9,00 %	50
Model 5	23.8.2022	44220	CITlab HTR+	1 318	4 107	0,36 %	202	603	1,02 %	2,43 %	3,34 %	50
Model 6	23.8.2022	44218	CITlab HTR+	1 266	3 938	0,44 %	254	772	1,53 %	5,08 %	3,72 %	50
Model 7	24.8.2022	44247	CITlab HTR+	653	2 049	0,17 %	202	603	2,26 %	4,87 %	8,35 %	50
Model 8	24.8.2022	44251	CITlab HTR+	601	1 880	0,17 %	254	772	2,64 %	6,85 %	9,14 %	50
Model 9a	5.11.2022	45904	PyLaia	1 881	5 818	1,00 %	243	751	1,20 %	3,69 %	5,12 %	101
Model 9	11.11.2022	46060	PyLaia	1 826	5 791	2,80 %	248	778	2,80 %	12,63 %	7,30 %	101
Model 10	1.4.2023	51087	PyLaia	1 877	5 811	0,60 %	247	760	1,10 %	4,56 %	3,47 %	153
Model 11	1.3.2024	60244	PyLaia	9 854	30 892	0,60 %	1 170	3 653	1,00 %	1,84 %*	5,05 %*	92
Bel_Adparatus_1735_model 1	14.8.2025	385017	PyLaia	1 477	10 337	0,32 %	203	1 626	1,50 %	3,39 %**	8,49 %**	250
Bel_Adparatus_1735_model 2	16.8.2025	386277	PyLaia	1 347	9 411	0,33 %	203	1 623	1,23 %	1,92 %**	8,49 %**	250
Bel_Adparatus_1735_model 3	6.9.2025	395837	PyLaia	1 347	9 411	0,31 %	203	1 623	0,93 %	2,33 %**	6,56 %**	250
Bel_Adparatus_1735_model 3a	23.9.2025	405517	PyLaia	1 347	9 411	0,27 %	203	1 623	0,96 %	2,33 %**	6,95 %**	250

Tab. 3 Prehľad spoločných modelov vytrénovaných v rokoch 2022 – 2025<sup>21</sup>

**Model 1:** ID43995 MIX print\_antikva\_kurziva\_fraktura\_svbach\_1 (počet slov: 2047, training set: pages 1 – 8, validation set: pages 9 – 10) – vytvorený 8.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,29 %, bez base modelu, bez existujúcich polygónov

**Model 2:** ID44001 MIX print\_antikva\_kurziva\_fraktura\_svbach\_2 (počet slov: 1878, training set: pages 1 – 6, 9 – 10, validation set: pages 7 – 8) – vytvorený 8.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,25 %, bez base modelu, bez existujúcich polygónov

**Model 3:** ID44136 MIX print\_antikva\_kurziva\_fraktura\_svbach\_3 (počet slov: 2047, training set: pages 1 – 8, validation set: pages 9 – 10) – vytvorený 18.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,28 %, bez base modelu, bez existujúcich polygónov, po oprave chýb v manuálnom prepise vzorky *Ground Truth*

**Model 4:** ID44117 MIX print\_antikva\_kurziva\_fraktura\_svbach\_4 (počet slov: 1878, training set: pages 1 – 6, 9 – 10, validation set: pages 7 – 8) – vytvorený 18.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,25 %, bez base modelu, bez existujúcich polygónov, po oprave chýb v manuálnom prepise vzorky *Ground Truth*

**Model 5:** ID44220 MIX print\_antikva\_kurziva\_fraktura\_svbach\_5 (počet slov: 4107 training set: pages 1 – 8, 11 – 15, validation set: pages 9 – 10) – vytvorený 23.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,36 %, bez base modelu, bez existujúcich polygónov, zvýšenie počtu slov v cvičnom súbore

**Model 6:** ID44218 MIX print\_antikva\_kurziva\_fraktura\_svbach\_6 (počet slov: 3938, training set: pages 1 – 6, 9 – 15, validation set: pages 7 – 8) – vytvorený 23.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,44 %, bez base modelu, bez existujúcich polygónov, zvýšenie počtu slov v cvičnom súbore

<sup>21</sup> \* Do overovacieho súboru modelu 11 bolo zaradených 9 strán (v predchádzajúcich modeloch vždy po 2 strany), WER1 a WER2 Modelu 11 preto uvádzame v tabuľke nasledovne: WER1 = najnižšia hodnota z deviatich, WER2 = najvyššia hodnota z deviatich. \*\* Do overovacích súborov Bel\_Adparatus\_1735\_model 1, Bel\_Adparatus\_1735\_model 2, Bel\_Adparatus\_1735\_model 3 a Bel\_Adparatus\_1735\_model 3a boli zaradené 3 strany, WER1 a WER2 modelov preto v tabuľke uvádzame nasledovne: WER1 = najnižšia hodnota z troch, WER2 = najvyššia hodnota z troch.

**Model 7:** ID44247 MIX print\_antikva\_kurziva\_fraktura\_svbach\_7 (počet slov: 2049 training set: pages 1 – 8, validation set: pages 9 – 10) – vytvorený 24.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,17 %, s base modelom ID43995 MIX print\_antikva\_kurziva\_fraktura\_svbach\_1 (CERtrain:0,29 %, CERvalidation: 2,33 %), bez existujúcich polygónov

**Model 8:** ID44251 MIX print\_antikva\_kurziva\_fraktura\_svbach\_8 (počet slov: 1880 training set: pages 1 – 6, 9 – 10, validation set: pages 7 – 8) – vytvorený 24.8.2022, technológia HTR+, Nr. of Epochs 50, CER (train) = 0,17 %, s base modelom ID44001 MIX print\_antikva\_kurziva\_fraktura\_svbach\_2 (CERtrain:0,25 %, CERvalidation: 3,09 %), bez existujúcich polygónov

**Model 9a:** ID45904 MIX print\_antikva\_kurziva\_fraktura\_svbach\_9a (počet slov: 5818, training set: pages 1 – 7, 9, 11 – 20, validation set: pages 8, 10) – vytvorený 5.11.2022, pretrénovanie Modelu 8 technológiou **PyLaia**, Nr. of Epochs 101, CER (train) = 1,00 %, bez base modelu, bez existujúcich polygónov

**Model 9:** ID46060 MIX print\_antikva\_kurziva\_fraktura\_svbach\_9 (počet slov: 5791, training set: pages 1 – 7, 9, 11 – 20, validation set: pages 8, 10) – vytvorený 11.11.2022, vytrénovaný nanovo technológiou PyLaia, Nr. of Epochs 101, CER (train) = 2,80 %, bez base modelu, bez existujúcich polygónov, viac slov v cvičnom súbore v porovnaní s Modelom 8

**Model 10:** ID51087 MIX print\_antikva\_kurziva\_fraktura\_svbach\_10 (počet slov: 5811, training set: pages 1 – 15, 17, 19 – 20, validation set: pages 16, 18) – vytvorený 1.4.2023, technológia PyLaia, Nr. of Epochs 153, CER (train) = 0,60 %, bez base modelu, bez existujúcich polygónov, bez použitia predvolenej funkcie *Deslant*

**Model 11:** ID60244 MIX print\_antikva\_kurziva\_fraktura\_svbach\_11 (počet slov: 30892, training set: pages 1 – 5, 7, 9, 10, 12 – 19, 21, 22 – 32, 35 – 50, 52 – 60, 62 – 93, validation set: pages 6, 8, 11, 20, 23, 33, 34, 51, 61) – vytvorený 1.3.2024, technológia PyLaia, Nr. of Epochs 92, CERtrain = 0,60 %, bez base modelu, bez existujúcich polygónov (pôvodne s existujúcimi polygónmi, ale výsledok bol nepoužiteľný), tréning na všetkých stranách prepísaného a editovaného dokumentu

**Bel\_Adparatus\_1735\_model 1:** ID385017 (počet slov: 10141, training set: pages 4 – 7, 9 – 12, 14 – 21, 23, 24, 42, 44, 46, 72, validation set: pages 8, 13, 22) – vytvorený 14.8.2025, technológia PyLaia, Nr. of Epochs 250, CER (train) = 0,32 %, s base modelom ID60244 Model 11 (CERtrain:0,60 %, CERvalidation: 1,00 %), s využitím existujúcich polygónov, s vylúčením slov označených ako unclear

**Bel\_Adparatus\_1735\_model 2:** ID386277 (počet slov: 9215, training set: pages 4 – 7, 9 – 12, 14 – 21, 23, 24, 42, 44, 46, 72, validation set: pages 8, 13, 22) – vytvorený 16.8.2025, technológia PyLaia, Nr. of Epochs 250, CER (train) = 0,33 %, s base modelom ID60244 Model 11 (CERtrain:0,60 %, CERvalidation: 1,00 %), s využitím existujúcich polygónov, s vylúčením slov označených ako unclear, po oprave chýb v manuálnom prepise vzorky *Ground Truth*

**Bel\_Adparatus\_1735\_model 3:** ID395837 (počet slov 9215, training set: pages: 4 – 7, 9 – 12, 14 – 21, 23, 24, 42, 44, 46, 72), validation set: pages: 8, 13, 22) – vytvorený 6.9.2025, technológia PyLaia, Nr. of Epochs 250, CER (train) = 0,31 %, s base modelom ID60244 Model 11 (CERtrain:0,60 %, CERvalidation: 1,00 %), bez existujúcich polygónov, s vylúčením slov označených ako unclear

**Bel\_Adparatus\_1735\_model 3a:** ID405517 (počet slov 9215, training set: pages: 4 – 7, 9 – 12, 14 – 21, 23, 24, 42, 44, 46, 72), validation set: pages: 8, 13, 22) – vytvorený 22.9.2025, technológia PyLaia, Nr. of Epochs 250, CER (train) = 0,27 %, s base modelom ID60244 Model 11 (CERtrain:0,60 %, CERvalidation: 1,00 %), bez existujúcich polygónov, s vylúčením slov označených ako unclear; tréning s totožnými parametrami ako model 3 s časovým odstupom 15 ní na overenie (ne)stability AI

## Použitá literatúra

- Bel, M. (1735). *Adparatvs ad Historiam Hvngrariae, sive collectio miscella, Monumentorum ineditorum partim, partim editorum, sed fugientium. Conquisiut, in Decades partitus est, & Praefationibus, atque Notis illustrauit, Matthias Bel*. Typis Joannis Paulli Royer.
- Bôbová, M., Katuščák, D., Kurhajcová, A., Maliniak, P., Mikušková, M., Nagy, I., Nižníková, L., Kunec, P., & Tomeček, O. (2023). *Automatická transkripcia slovacikálnych historických dokumentov*. Belianum. Vydavateľstvo Univerzity Mateja Bela v Banskej Bystrici. <https://doi.org/10.24040/2022.9788055720203>
- Bôbová, M., Katuščák, D., Kurhajcová, A., Maliniak, P., Mikušková, M., Nagy, I., Nižníková, L., & Tomeček, O. (2024). *Automatická transkripcia historických dokumentov v prostredí webovej aplikácie Transkribus: Metodická príručka pre účastníkov workshopu*. Belianum. Vydavateľstvo Univerzity Mateja Bela v Banskej Bystrici. <https://doi.org/10.24040/2024.9788055721439>
- Burdick, A., Drucker, J., Lunenfeld, P., Presner, T., & Snapp, J. (2012) *Digital Humanities*. The MIT Press. <https://doi.org/10.7551/mitpress/9248.001.0001>
- Capurro, C., Provatorova, V., & Kanoulas, E. (2023). Experimenting with training a neural network in transkribus to recognise text in a multilingual and multi-authored manuscript collection. *Heritage*, 6(12), 7482–7494. <https://doi.org/10.3390/heritage6120392>
- Cuellar, Á., & Boadas, S. (2025). Artificial Intelligence for Calligraphic Writer Identification: The Case of Lope de Vega's Autographs. *Hipogrifo. Revista de literatura y cultura del Siglo de Oro*, 12(1), 517-532. <https://doi.org/10.13035/H.2025.13.01.36>
- Griffiths, R. (2024). Handwritten Text Recognition (HTR) for Tibetan Manuscripts in Cursive Script. *Revue d'Etudes Tibétaines*, (72), 43-51. [https://d1i1jdw69xsqx0.cloudfront.net/digitalhimalaya/collections/journals/ret/pdf/ret\\_72\\_03.pdf](https://d1i1jdw69xsqx0.cloudfront.net/digitalhimalaya/collections/journals/ret/pdf/ret_72_03.pdf)
- Harish, R., & Raghavendra Rao, G. N. (2024). Transcription of Ancient Indian Manuscripts Through Artificial Intelligence – Current Status of Technology and the Way Forward. In: *Artificial Intelligence: Theory and Applications* (vol. 844, s. 339-352). Springer Nature, Singapore. [https://doi.org/10.1007/978-981-99-8479-4\\_25](https://doi.org/10.1007/978-981-99-8479-4_25)
- Humphries, M., Leddy, L. C., Downton, Q., Legace, M., McConell, J., Murray, I. & Spence, E. (2024). Unlocking the Archives: Large Language Models Achieve State-of-the-Art Performance on the Transcription of Handwritten Historical Documents. In *ArXiv*. <https://arxiv.org/html/2411.03340v1>
- Komensky, J. A. (1798). *Joann. Amos Comenii Orbis pictus, in hungaricum, germanicum et slavicum translatus et hic ibive emendatus*. Sumtibus & Typis Simonis Petri Weber.
- Marsili, G., & Hassam, S. N. (2025). From archives to museum and back: Transcribing, digitizing, and enriching cultural heritage and manuscript legacy data of the Villa del Casale of Piazza Armerina. *Digital Applications in Archaeology and Cultural Heritage*, 38, e00441. <https://doi.org/10.1016/j.daach.2025.e00441>
- Mistrík, J. (1993). Prepis z iného písma. In *Encyklopédia jazykovedy*. Obzor.
- Muehlberger, G., Seaward, L., Terras, M., Ares Oliveira, S., Bosch, V., Bryan, M., Colutto, S., Déjean, H., Diem, M., Fiel, S., Gatos, B., Greinoecker, A., Grüning, T., Hackl, G., Haukkovaara, V., Heyer, G., Hirvonen, L., Hodel, T., Jokinen, M., ... Zagoris, K. (2019). Transforming scholarship in the archives through handwritten text recognition: Transkribus as a case study. *Journal of Documentation*, 75(5), 954–976. <https://doi.org/10.1108/JD-07-2018-0114>
- Muehlberger, G. (2021). H2020 Project READ (Recognition and Enrichment of Archival Documents) – 2016 – 2019. [https://www.academia.edu/22653102/H2020\\_Project\\_READ\\_Recognition\\_and\\_Enrichment\\_of\\_Archival\\_Documents\\_-\\_2016-2019](https://www.academia.edu/22653102/H2020_Project_READ_Recognition_and_Enrichment_of_Archival_Documents_-_2016-2019)

- Nagy, I. (2021). Možnosti aplikácie metódy digitálnej transkripcie historických rukopisných textov pri sprístupňovaní archívnych fondov. *Slovenská archivistika*, 51 (2), 53-67.
- Noble, A. (2025, June 20). The new Text Titan I ter and how it compares to ChatGPT, Gemini, and other LLMs. *Blog.transkribus.org*.  
<https://blog.transkribus.org/en/new-text-titan-i-ter-and-how-it-compares-to-chatgpt-gemini-and-other-llms>
- Nockels, J., Gooding, P., Ames, S., & Terras, M. (2022). Understanding the application of handwritten text recognition technology in heritage contexts: A systematic review of Transkribus in published research. *Archival Science*, 22(3), 367-392. <https://doi.org/10.1007/s10502-022-09397-0>
- Ost, K. (2024). Möglichkeit und Grenzen einer groß angelegten Volltexterschließung von Inkunabeln mit Transkribus. *Zeitschrift Für Bibliothekswesen Und Bibliographie*, 71(3), 169–181.  
<https://doi.org/10.3196/186429502471338>
- Sánchez-Salido, E., Menta, A., & García-Serrano, A. (2023). Seeking information in Spanish Historical Newspapers: The Case of Dario de Madrid (18th and 19th Centuries). *Digital Humanities Quarterly*, 17(4).  
<https://dhq.digitalhumanities.org/vol/17/4/000735/000735.html>
- Park, F. (2025, April 22). What are Super Models and how do they work? *Blog.transkribus.org*.  
<https://blog.transkribus.org/en/what-are-super-models-and-how-do-they-work>
- Raghallaigh, B. O., Palandri, A., & Cárthaigh, C. M. (2022). Handwritten Text Recognition (HTR) for Irish-Language Folklore. *Proceedings of the 4th Celtic Language Technology Workshop within LREC2022*, 121–126.  
<https://doi.org/> <https://aclanthology.org/2022.cltw-1.17>
- Smida, M. (2023). Možnosti automatickej transkripcie v platforme Transkribus na príklade správy o vybavovaní sťažností občanov v období komunistickej diktatúry. *Acta historica Neosoliensia*, 26(1), 125-148.  
<https://doi.org/10.24040/ahn.2023.26.01.125-148>
- Tomeček, O., & Nagy I. (2024). *Reambulačný protokol mesta Banská Bystrica z roku 1820: edícia prameňa pomocou nástroja na automatickú transkripciu historických dokumentov*. Belianum. Vydavateľstvo Univerzity Mateja Bela v Banskej Bystrici. <https://doi.org/10.24040/2024.9788055721903>
- Tomeček, O. (2025). O možnostiach automatického prepisu historických rukopisných dokumentov pomocou platformy Transkribus. *Historica. Revue pro historii a příbuzné vědy*, 16(1), 87-108.  
<https://doi.org/10.15452/Historica.2025.16.0005>

## O autorkách

*Lucia Nižníková* pracuje ako vedúca Oddelenia podpory vedy v Univerzitetnej knižnici Univerzity Mateja Bela v Banskej Bystrici. Okrem metodologickej činnosti v oblasti evidencie publikačnej činnosti vedeckých a pedagogických zamestnancov na univerzite sa venuje najmä témam súvisiacim s otvoreným prístupom a otvorenou vedou, a ich aplikáciou v praxi. Po absolvovaní troch akreditovaných vzdelávacích programov v Centre vedecko-technických informácií na tieto témy poskytuje akademickej obci univerzity konzultačné služby, realizuje vzdelávacie podujatia a publikuje články. Bola spoluorganizátorkou medzinárodnej konferencie ILIDE (Innovative Library in Digital Era) a knihovníckej konferencie Bibliosféry. V rokoch 2020 – 2024 bola členkou riešiteľského kolektívu projektu aplikovaného výskumu SKRIPTOR podporeného Agentúrou na podporu výskumu a vývoja, na ktorý od roku 2025 nadväzuje projekt DIGINARCH-AI zameraný na modernizáciu služieb bádateľom v štátnych archívoch prostredníctvom digitalizácie vybraných archívnych fondov a sprístupnenia ich obsahu pomocou nástrojov umelej inteligencie.

E-mail: [lucia.niznikova@umb.sk](mailto:lucia.niznikova@umb.sk)

ORCID: <https://orcid.org/0000-0002-7450-0193>

*Michaela Mikušková* je absolventkou Katedry knižničnej a informačnej vedy Filozofickej fakulty Univerzity Komenského v Bratislave. V súčasnosti zastáva pozíciu riaditeľky Univerzitetnej knižnice Univerzity Mateja Bela v Banskej Bystrici. Venuje sa vzdelávaniu a publikovaniu v oblasti otvorenej vedy, stála pri vzniku prvých politík a začiatkoch budovania infraštruktúry otvoreného prístupu na univerzite. Niekoľko rokov bola členkou organizačného výboru konferencií ILIDE (Innovative Library in Digital Era) a Bibliosféry, vystúpila na viacerých odborných podujatiach. V rokoch 2020 – 2024 bola súčasťou riešiteľského kolektívu projektu SKRIPTOR, aktuálne participuje na projekte DIGINARCH-AI (oba podporené Agentúrou na podporu výskumu a vývoja).

E-mail: [michaela.mikusкова@umb.sk](mailto:michaela.mikusкова@umb.sk)

ORCID: <https://orcid.org/0000-0003-4465-2812>