

# VALÓS TÉRBEN – AZ ONLINE TÉRÉRT

**Networkshop 31: országos konferencia**

2022. április 20–22.  
Debreceni Egyetem

Szerkesztette: Tick József, Kokas Károly, Holl András

HUNGARNET Egyesület  
Budapest, 2022



A kötet megjelenését támogatta az  
Energiaügyi Minisztérium

Szerkesztette: Tick József, Kokas Károly, Holl András

Tipográfia és tördelés: Vas Viktória

Workshop

2022. április 20–22. Debreceni Egyetem, konferencia előadásainak közleményei

ISBN 978-615-82243-0-7

DOI: [10.31915/NWS.2022](https://doi.org/10.31915/NWS.2022)

Kiadja a HUNGARNET Egyesület  
az MTA Könyvtár és Információs Központ közreműködésével  
Budapest  
2022

Borítókép: [freepik.com](https://www.freepik.com)

## TARTALOMJEGYZÉK

Előszó .....	5
Lencsés Ákos: A nyílt tudomány pénzügyi vonatkozásai .....	7
Farkas Katalin: Centenáriumi média-adattár és virtuális kiállítás létrehozásának tanulságai az SZTE Klebelsberg Könyvtárban .....	13
Bódog András: A nyílt archívumi információs rendszer (OAIS) szabványának honosítása.....	20
Perlaki Attila: Oktatást segítő gamifikációs alkalmazások, mint szakdolgozati témák .....	27
Csapó Noémi – Dani Erzsébet: APPropó fejlődés – A Bács-Kiskun Megyei Katona József Könyvtár mobilapplikációja.....	32
Simon András: Integrált könyvtári rendszerek tranzakciós rekordjainak vizsgálata, a könyvtári állomány digitalizálásának tervezésekor.....	41
Németh Márton: Az OSZK Webarchívum nemzetközi kapcsolatai.....	58
Antal Péter: A mesterséges intelligencia kihívásai a XXI. század társadalmára .....	70
Hajdu Csaba – Szilágyi Zoltán: Modern robotikai technológiai ismeretek oktatása „Teljes spektrumú” oktatási módszerrel .....	77
T. Nagy László – Boda István Károly – Tóth Erzsébet: E-tananyagfejlesztés virtuális 3D környezetben.....	84
Palencsárné Kasza Marianna: Digitális átállás – Minőség – lehetőségek az EQAVET terén.....	92
Nagy Gyula: Nemzetközi kitekintés a felsőoktatási könyvtárak világára: a EUGLOH könyvtári workshopja .....	99
Babocsay Gergely: Az európai természettudományi gyűjtemények digitális integrációja: határ a csillagos ég.....	108
Somorjai Noémi: Egyenlőtlenségek a tudományos kutatás területén. Az amatőr kutatók szerepe .....	114
Molnár Dániel – Dani Erzsébet: Robotok a könyvtárban: Hogyan válhat a robotika a könyvtári mindennapok részévé? .....	122
Horváthné Felföldi Helga: Digitalizáció a szakképzésben. A Szakmajegyzékben szereplő szakmák digitáliskompetencia jártassági szintjeinek felülvizsgálata .....	130
Kalcsó Gyula: Ne csak útra csomagoljunk! Miért fontos a csomagolás a digitális megőrzésben? .....	138
Karsa Zoltán István – Szeberényi Imre: A CIRCLE felhő elmúlt évtizede .....	146
Bobák Barbara – Kasza Péter: Az MI lehetőségei a kora újkori filológiában: Johannes Michael Brutus <i>Rerum Ungaricarum</i> libri kéziratának digitális kiadása (esettanulmány) .....	154
Egyed-Gergely Júlia – Vajda Róza, Gárdos Judit – Horváth Anna – Meiszterics Enikő – Micsik András – Martin Dániel – Marx Attila – Pataki Balázs – Siket Melinda: Szociológia, kutatási adatok, mesterséges intelligencia: lehetőségek és tapasztalatok .....	161
Szemes Botond – Bajzát Tímea – Fellegi Zsófia – Kundráth Péter – Horváth Péter – Indig Balázs – Dióssy Anna – Hegedüs Fanni – Pantyelejev Natali – Sziráki Sarolta – Vida Bence – Kalmár Balázs – Palkó Gábor: Az ELTE Drámakorpuszának létrehozása és lehetőségei.....	170



Sebestyén Ádám: Az ELTEdata szemantikus adatbázis legújabb fejlesztései.....	179
Szlamka Erzsébet: Új trendek a tanulási eredmények tanúsításában .....	185
Tóth Máté – Héjja Balázs: Webshop indítása közkönyvtári környezetben.....	192
Etlinger Mihály – Hernády Judit: A kiadás hagyatéka / a hagyatéka kiadása: A Régi Magyar Költők Tárának hálózati kiadásáról.....	199
Varga Emese – Makkai T. Csilla: „Ki a fenének kell collstok?” A digitális szöveg rejtett mértékegységei .....	204
Dobás Kata – Fazekas Júlia: ITIdata – Egy irodalmi adatbázis fejlesztése Wikibase alapon és ennek hasznosítása Kosztolányi Dezső forrásjegyzékénél .....	211
Sörény Edina: Kézai Simon Program – digitális családi fotóarchívum.....	219
Fülöp Tiffany – Molnár Tamás – Hoczopán Szabolcs: Open Monograph Press e-könyvplatform a Szegedi Tudományegyetemen .....	227
Palkó Gábor: Mesterséges intelligencia, digitális bölcsészet, kulturális örökség: trendek és eredmények.....	235
Pergéné Szabó Enikő – Bátfai Mária Erika: A tudományos publikálás támogatása a Debreceni Egyetemi és Nemzeti Könyvtárban .....	241
Csirmazné Rezi Éva: Nemzetközi kiadványazonosítók és kötelezpéldányok kezelése az OSZK OKP (Országos Könyvtári Platform) rendszerében .....	250
Alföldi István – Dióssy Anna Laura: Digitálisan született kutatási anyagok megőrzése: a relációs adatbázis mint born-digital objektum .....	262
Fekete Norbert: HTR-modellépítés és kézírásfelismerés nagyméretű, többszerzős szövegtörzsen. A Transkribus alkalmazása az Arany János hivatali iratokon.....	271
Horváth Péter – Kundráth Péter – Palkó Gábor: ELTE Népdalkorpusz – magyar népdalok gépileg annotált adatbázisa .....	276
Nagy György: IKT eszközök alkalmazása az alsó tagozatos környezetismeret órákon.....	284
Köpösdí Zsuzsa – Molnár Tamás: Multimédiás, interaktív és adaptív tananyagok létrehozásának lehetőségei H5P keretrendszerrel .....	289
Jankó Tamás: Munka 4.0 – Ipar 4.0 – Szakképzés 4.0 – : A digitális kompetencia jövőbeni fejlesztési útjai .....	296
Békésiné Bognár Noémi Erika – Nagy Andor: Megújuló könyvtári statisztika: az egységes adatstruktúra és a korszerű megjelenítés kialakításának útján .....	304
Bolya Máttyás: Kézírtos dallamlejegyzések feldolgozása MI-vel támogatott digitális környezetben .....	310
Maróthy Szilvia – Seláf Levente – Vigyikán Villó: Régi magyar verskorpusz összeállítása stilometriai és számítógépes metrikai kutatásokhoz .....	324
Szúcs Kata Ágnes: Kézírtos források transzformációinak lehetőségei a közgyűjteményekben.....	330
Fellegi Zsófia: A digitális filológia infrastruktúrái. A DigiPhil megújulásáról. ....	338
Mihály Eszter: Mi az a dHUpla? A Digitális Bölcsészeti Platform bemutatása.....	345
Nemeskey Dávid Márk – Palkó Gábor: Szemantikus névelém-azonosítás magyar nyelvű szövegeken (a HuWikifier bemutatása) .....	359

## Kéziratok forrásainak transzformációinak lehetőségei a közgyűjteményekben Transformations of manuscripts in public collections

Szűcs Kata Ágnes  
Digitális Bölcsészeti Központ  
Országos Széchényi Könyvtár  
[szucs.kata@oszk.hu](mailto:szucs.kata@oszk.hu)

### Abstract

Providing digital accessibility for manuscript resources in public collections is feasible at various levels. Though what does digital processing of manuscripts mean? Filling in record fields of an online catalogue, uploading scanned images into Transkribus, and digital publishing a TEI XML file can all line up behind the concept. In my paper, I discuss the changes a manuscript undergoes and the multiple media formats it becomes interpretable and provided for the public. From a practical point of view, I demonstrate how the use of Transkribus is transforming a workflow of a public collection. The paper also investigates methods for building HTR models by comparing two basic theoretical strategies in practice. The first approach is to add the previously generated Ground Truth (GT) as a Base Model (BM). The other is to train an entirely new model by merging the old and new GTs without a base model.

**Keywords:** htr, handwritten text recognition, Transkribus, digital processing, digital edition, manuscript, public collection

### 1. Bevezetés

A Digitális Bölcsészeti Központ Kiss József-projektje egyrészt mint digitális forráskiadás, másrészt mint kéziratok feldolgozását és közzétételét kidolgozó közgyűjteményi feladat és végül mint automatikus kézírásfelismerést szorgalmazó vállalkozás is fontos célokat tűzött ki maga elé és valósított meg.

AdHUpla<sup>1</sup> – Digital Humanities Platform – digitális szövegkiadásokat és a szövegfeldolgozáson alapuló kreatív tartalmakat tesz közzé. Tehát egy olyan publikálási környezet, amely a közgyűjtemények szöveges forrásainak digitális megjelenítésére szolgál. Az elmúlt két évben, a Móricz-kutatócsoport<sup>2</sup> munkája mellett, ez a közgyűjteményi pilot projekt segítette feltérképezni és kialakítani a megjelenítéstől kezdve az xml publikáció tagkészletén át, a kép és szöveg összekapcsolásának lehetőségeit.

<https://dhupla.hu/collection/kiss-jozsef-levelazes>

A Kiss József-projektben továbbá kialakult egy olyan workflow, amely bármely kézirat esetében alkalmazható lépéseket ír elő a feldolgozáshoz és nem utolsósorban a publikálás elkészítéséhez, másrészt pedig a szoftverhasználatra is létrejött egy részletesen kidolgozott, mégis rugalmas ajánlás. Ez főképp a Transkribus platformot,<sup>3</sup> amely történelmi dokumentumok

1 Digital Humanities Platform, [www.dhupla.hu](http://www.dhupla.hu) hozzáférés: 2022. június 23.

2 A munkáról bővebben a Móricz-műhely jelenleg költöztetés alatt álló oldalán lehet olvasni, <https://pim.hu/hu/digitalis-bolcseszeti-kozpont/moricz>, hozzáférés: 2022. június 23.

3 READ-COOP, Transkribus, <https://readcoop.eu/transkribus> hozzáférés: 2022. június 23.

szövegfelismerésére, képelemzésére és szerkezetfelismerésére alkalmas, és az Oxygen tei-xml szerkesztő szoftvert<sup>4</sup> helyezi előtérbe.

A közgyűjteményekben fellelhető kéziratos források digitális szolgáltatása különböző szinteken lehetséges. A kéziratok mediális transzformációja ebben a projektben az alábbi, egymáshoz is szervesen kapcsolódó állomások köré szerveződött:

- Fizikai példány
- Digitalizált másolat(ok)
- Transkribus dokumentum
- Kétrétegű PDF
- XML-file
- Online publikáció (dHUplá)

A kéziratok minden formájának megvan a maga jelentősége és önmagában olyan információkat hordoz, amelyeket más változatok nem. A Transkribus dokumentum már egy átírással rendelkező, szegmentált képet, illetve xml dokumentumot tartalmazó objektum, amely ebben a formában lehetőséget teremtett nem csak a digitális forráskiadás létrehozására, hanem automatikus kézírásfelismerő modell építésére is. A továbbiakban esettanulmány jelleggel az e téren elért újabb eredményekről lesz szó.

## 2. Automatikus kézírásfelismerés (Handwritten Text Recognition – HTR)

A szoftver a gyakorlatban is alkalmasnak bizonyult a közgyűjteményekben fellelhető különböző kéziratok kezelésére és publikálására. A jelenlegi workflowban egyrészt az átírások elkészítésére és ezek különböző formátumokba történő exportálásra használjuk, illetve folytattuk az automatikus kézírásfelismertetést is a Kiss József-levelezésen. A Transkribus szervezet<sup>5</sup> és a körülötte kialakult felhasználói közösség<sup>6</sup> pedig több téren is segítette a kézírásfelismerésben végzett munkánkat.

Egy korábbi publikációmban<sup>7</sup> bemutattam az általánosságoktól a konkrét kimenetelig a Transkribus program használatát, a mesterséges intelligencián alapuló kézírásfelismerő-modell építésének módját, illetve a Kiss József kézírásán tanult első modellek eredményeit is. Az egy kéz által írt leveleken tanult eddigi legjobb modellünk 6,94-es hibaszázalékkal rendelkezik.

Hosszútávú terveink közé tartozik egy minél általánosabb, a 19–20. századi magyar írók kézírását könnyen felismerő modell létrehozása. Efelé tett következő lépés az volt, hogy a Kiss József-levelezés másik felét, azaz a 19. századi magyar költőnek, A Hét című hetilap szerkesztőjének címzett leveleket is bevontuk a modellépítésbe. Az így keletkezett vegyes kézírásos korpusz már elegendő mennyiségű adattal rendelkezett egy általánosabb magyar

4 Oxygen XML Editor, <https://www.oxygenxml.com/> hozzáférés: 2022. június 23.

5 A kezdeményezés 2016 óta létezik és az Európai Unió 2020-as Horizon programja által finanszírozott kutatási projektből nőtte ki magát. 2019 július 1-jétől Read COOP néven és immár vállalati alapokon folytatja megkezdett munkáját. European Commission CORDIS research results, *Recognition and Enrichment of Archival Documents*, hozzáférés: 2022. június 23., <https://cordis.europa.eu/project/id/674943>. Vö.: READ-COOP, *Our Story*, hozzáférés: 2022. június 23., <https://readcoop.eu/our-story/>.

6 „Transkribus users | Facebook”, hozzáférés: 2022. június 23., <https://www.facebook.com/groups/614090738935143>.

7 Szűcs Kata Ágnes, „Automatikus kézírás-felismertetés Kiss József levelezésén”, in *Online térben az online térért: Networkshop 30: országos online konferencia*. 2021. április 6-9. Eötvös Loránd Tudományegyetem (Networkshop, HUNGARNET Egyesület, 2021), 73–80, <https://doi.org/10.31915/NWS.2021.8>.



nyelvű kézírást felismerő modell létrehozásához, amely azóta nyilvánosan is hozzáférhető a Transkribus oldalán<sup>8</sup> és az asztali alkalmazáson keresztül is.

Neve: *Transkribus Hungarian Handwriting 19th–20th Century*

Epoch-ok száma: 250

Szavak száma: 74.862

Sorok száma: 16.630

## 2.1. Vegyes kézírásos tanult modellek

A jelenlegi modellhez felhasznált kéziratok a Petőfi Irodalmi Múzeumban találhatóak, melyek között szerepelnek borítékok, képeslapok, hagyományos és fejléces levelek, névjegykártyák. A levélírók Kiss József és családja, illetve a századforduló írói, újságírói és művészei. Ez nagyságrendileg 300 darab változó hosszúságú és minőségű levelet jelent. A levelezés további kéziratjai az Országos Széchényi Könyvtárban találhatóak, amelyekkel terveink szerint a modell tanítókorpusza a későbbiekben bővülni fog.

A vegyes kézírásmodell építéséhez alapvetően több adat szükséges, mivel nagyobb változatosság mutatkozik az egyes írásképek között, mint egyetlen ember esetében. Például, a Kiss József kézírásán tanult modellnél körülbelül 1/3 mennyiségű adat állt rendelkezésre, mégis 2,27%-kal jobb eredményt sikerült elérni.<sup>9</sup>

Az egyik legfontosabb mérőszám a HTR modell építésekor, a Character Error Rate (CER) on Validation Set,<sup>10</sup> azaz, az ellenőrző korpuszon ejtett karakterhibaérték aránya. A jelenlegi legjobb eredmény 9,19% – és bár akadnak olyan levelek a korpuszban, amelyeket szinte tökéletesen képes átírni, azt nem lehet elvárni ettől a modelltől, hogy ezt bármilyen magyar nyelvű kézirattal megtegye.

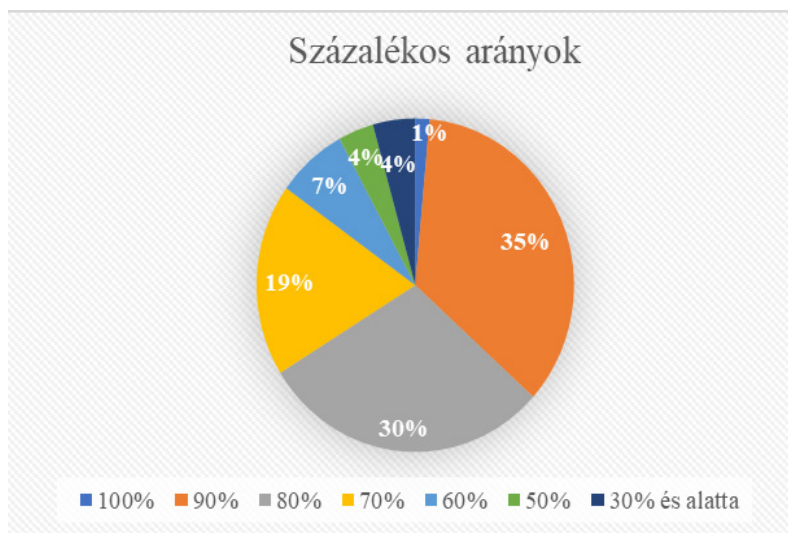
Egyrészt azért, mert a modell egy adott korpuszon, Kiss József-levelezésen tanult (1265 oldal), és így jelenleg ennek a felismerésére a legalkalmasabb. A 9,19-es érték pedig egy előjelzett átlag, amelyet a Validation Set nevű elkülönített ellenőrzőkorpuszon (138 oldal) számol ki az algoritmus.

Tehát a modell ezen a tanulási folyamatból kimaradó kézírathalmazon átlagosan közel 90%-os pontossággal meg tudja állapítani, hogy milyen karakterekből áll össze egy levél szövege, amely jó indikátora annak, hogy hogyan teljesítene egy új, ismeretlen kéziraton. Ez az eredmény egyes leveleknél valóban 90-, vagy akár 100%-os pontosságot jelent, de más, nehezebb kézírások/írásképek vagy oldaltükrök esetében 60-, 20- vagy akár 0%-os eredményt is produkálhat. Ezeket az olvasatokat átlagolja az algoritmus a CER meghatározásakor. Egy általános, a magyar nyelvű kéziratokat széleskörűen felismerő modellhez a jelenleginél sokkal több adatra van szükség.

8 READ-COOP, *Transkribus Hungarian Handwriting 19th–20th Cent.*, hozzáférés: 2022. június 23., <https://readcoop.eu/model/hungarian/>.

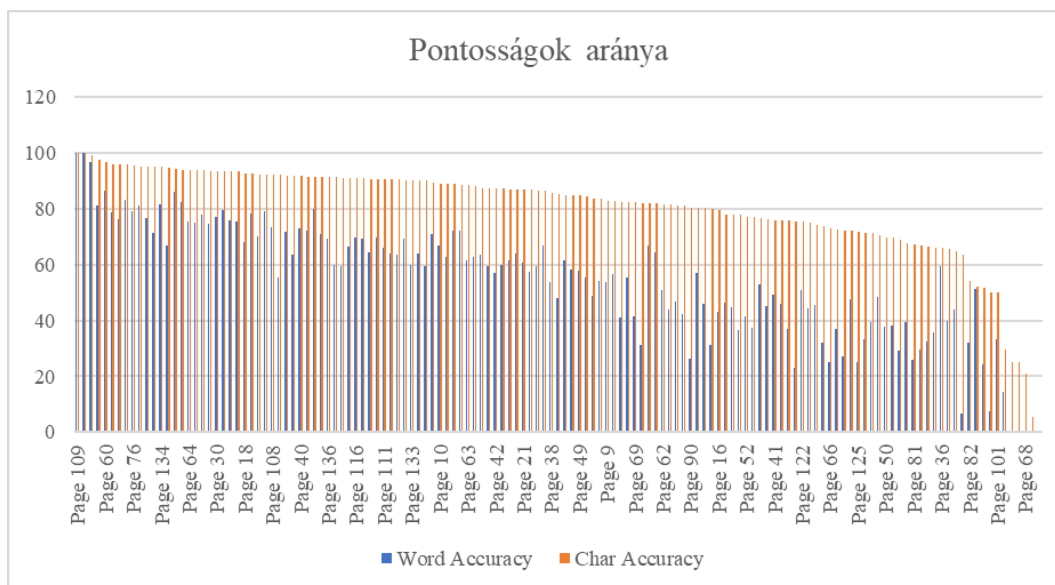
9 CER on VS: 6,94%  
Epoch-ok száma: 200  
Szavak száma: 74.862 25.348  
Sorok száma: 16.630 5.634

10 Vö.: Elisabeth Heigl, *CER? Don't Worry!*, hozzáférés: 2022. június 23., <https://rechtsprechung-im-ostseeraum.archiv.uni-greifswald.de/cer-dont-worry/>.



1. ábra A CER on Validation Set százalékos eloszlása: 100% tökéletes, átírás 30% és alatta, jelentős hibákat tartalmazó átírás.

Egy másik mérőszám, a Word Error Rate (WER) a szavak hibaarányát jelzi. A helyesen átírt szavak megoszlása a karakterekhez képest nagyobb kilengéseket mutat, mert egy karakter hibás felismerésekor az adott szót is, amelyben az szerepel automatikusan hibásnak fogja értékelni a program. Ez abban az esetben is igaz, ha csak egy nagybetű-kisbetű tévesztésről, esetleg a szó végén lévő írásjel elhagyásáról van szó.



2. ábra Szavak pontosságának az alakulása a karakterpontossághoz képest.

## 2.2. Modellek egymásba építése

A vegyes kéziratkorpusz folyamatos feldolgozása új kérdéseket vetett fel a modellek egymásba építésével kapcsolatban. Mivel a modellépítéshez szükséges átírt kéziratok folyamatosan készülnek, ezért egy-egy nagyobb mennyiségnél fontossá vált, hogy milyen módon tudjuk azokat beépíteni a már meglévő modellbe. Ezáltal be tudtunk kapcsolódni egy aktív nemzetközi diskurzushoz, melyből kiderült, hogy ezek a kérdések nemcsak a Transkribus használói, hanem a fejlesztői számára sem egyértelműek.





Két elméleti lehetőség van, de az, hogy melyik a jó megoldás az a kéziratok korpusztól is függ. Eddigi tapasztalataink szerint befolyásolja az eredményeket, hogy mennyire homogén a kézírás, időben közel vagy távolabb vannak-e egymástól az egyes dokumentumok, műfajuk megegyezik-e stb.

Az első opció, hogy az új anyagba alapmodellként építjük be az eddigi kéziratokon tanult modellt és így fejlesztjük tovább.<sup>11</sup> Ebben az esetben a modell minden egyes tanulási ciklusa egy már meglévő modellre, azaz alapmodellre (Base Modellre – BM) épül. Az alapmodellek „emlékeznek” arra, amit már megtanultak, ezért elméletileg minden egyes új tanulási ciklus javítja a modell minőségét. Az új modell tanul az elődjéből, és így egyre jobbá válik.<sup>12</sup>

A második, hogy az elkészült új anyagot Ground Truth-ként (GT) adjuk hozzá a régihez és ebből hozunk létre új modellt. A GT vagy alapigazság egy statisztikából ismert alapfogalom, amely egy adott kérdéssel kapcsolatos igazság ismeretére vonatkozik, ez az ideális elvárt eredmény. Kéziratok esetében ez egy dokumentum 100 %-ban helyes gépelt példányát jelenti, amely a Mesterséges Intelligencia betanítására szolgál.<sup>13</sup>

Mindegyik modell HTR+ technológiával készült és a modellépítésnél be volt kapcsolva az *omit line by tags* funkció, melynek hatására az algoritmus automatikusan figyelmen kívül hagyja azokat a sorokat, ahol bizonytalan olvasatot vagy olvashatatlan szövegrészt jelöltünk. Ez minimum 2–3% javulást eredményezett a modelleknél, annak ellenére, hogy ezáltal kb. 1000 sorral kevesebb lett a GT. Az első lépésben elkülönítettünk egy kisebb Tesztkorpuszt (TK). Ez az anyag egyáltalán nem vett részt a modellépítésben, így befolyásoltság (bias) nélküli tesztelésre adott lehetőséget a modellépítés során.

### 2.3. A Base Modellel tanított modellek

Az első vegyeskézírásból készült modellünk 484 oldalnyi szövegen tanult (6.716 sor), alapmodell nélkül 10,38%-os lett. Ebbe építettük be a Kiss József kézírásán tanult modellt (CER on VS: 6,94%), amellyel 0,27%-ot javult a teljesítménye.

A második vegyeskézírásból készült modell 298 új oldalnyi szövegen tanult és azt a modellt (ID: 34757 Vegyes\_kézírás\_13) használtuk Base Modellként, amelybe Kiss József kézírása is be volt építve alapmodellként. Amikor az első vegyeskézírás korpuszon épült modellt, (ID: 34600 Vegyes\_kézírás\_11) használtuk alapmodellként, 0,1%-kal kaptunk rosszabb eredményt. Ebből az látszik, hogy egy alapmodell alapmodelljeként beépített – jelen esetben a Kiss József kézírásán tanult – modell jelenléte vagy hiánya alig okozott teljesítmény béli különbséget. Harmadik lépésként az összevont vegyes kézíráskorpuszhoz építettük be a Kiss József-kézírása alapmodellét, amely a legjobb eredményt hozta az első fázisban (9,92%). Tehát nagyobb hatása van az alapmodellnek, ha nincs annyira „eltemetve.”

11 Vö.: Dirk Alvermann, *Use Case: “Model Booster”*, hozzáférés: 2022. június 23., <https://rechtsprechung-im-ostseeraum.archiv.uni-greifswald.de/use-case-model-booster/>.

12 READ-COOP, *Base Models*, hozzáférés: 2022. június 23., <https://readcoop.eu/glossary/base-models/>.

13 READ-COOP, *Ground Truth*, hozzáférés: 2022. június 23., <https://readcoop.eu/glossary/ground-truth/>.

	50 epoch	100 epoch	150 epoch	200 epoch
BM 34757 Vegyes kézírás_13	10,05%	10,28%	10,35%	0%
BM: 34600 Vegyecs kézírás_11	10,68%	10,15%	10,42%	10,20%
BM: 31617 Kiss József kézírása + Összevont GT	12,27%	10,28%	9,92%	10,30%

3. ábra. Második vegyes kézírásmodellek alapmodellel.

Az alapmodellel történő modellépítés egyik vitathatatlan előnye, hogy viszonylag kevés epoch alatt éri el a maximumát, ami idő- és energiatakarékos megoldás lehet bizonyos esetekben. Az epoch a Transzkibus terminológiájában azt az iterációs ciklust jelenti, ahányszor az algoritmus a tanulási folyamatot elvégzi az adott korpuszon.

## 2.4. Összevont Ground Truth alapú modellek

A második szakaszban az összevont GT-s modellekkel folytattuk a tesztelést.

Először kombináltuk az első és második vegyes kézírásmodell anyagát, és BM nélkül az algoritmus 10,34%-os CER értéket produkált a korpuszon. Ehhez hozzátettük a Kiss József kézírásával készült levelek korpuszát, ezen már 9,31%-os eredményt értünk el.

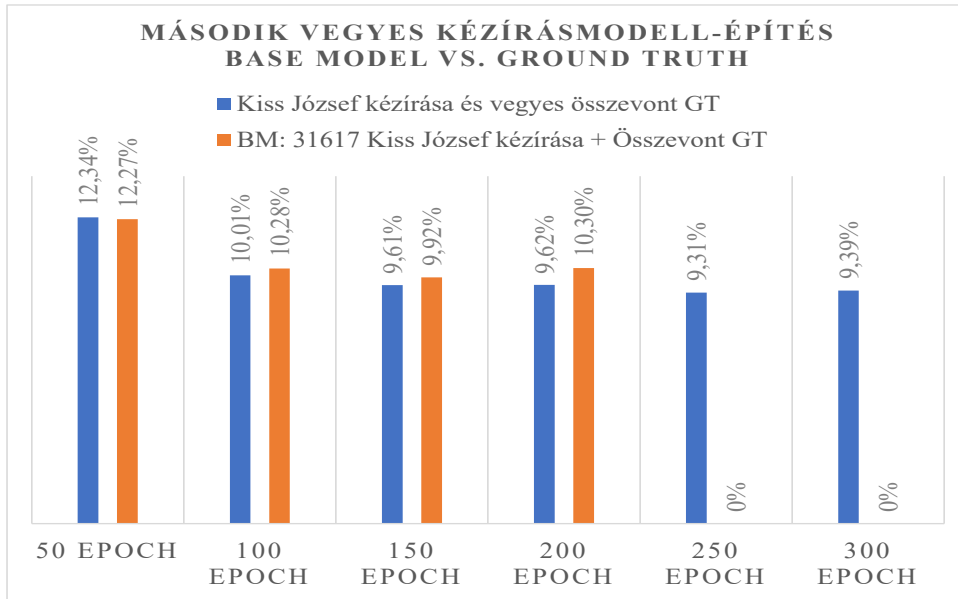
A különböző tesztek lefuttatása után a TK-t is beépítettük a modellbe és így értük a vegyes korpusz eddigi legjobb, 9,19%-os eredményét.

	50 epoch	100 epoch	150 epoch	200 epoch	250 epoch	300 epoch
Vegyes összevont GT	12,71%	10,96%	10,72%	10,55%	10,36%	10,34%
Kiss József kézírása és vegyes összevont GT	12,34%	10,01%	9,61%	9,62%	9,31%	9,39%
Kiss József kézírása és vegyes összevont GT + Teszt Korpusz	12,16%	10,14%	9,89%	9,19%	9,19%	0%

4. ábra. Második vegyes kézírásmodellek összevont alapigazsággal.

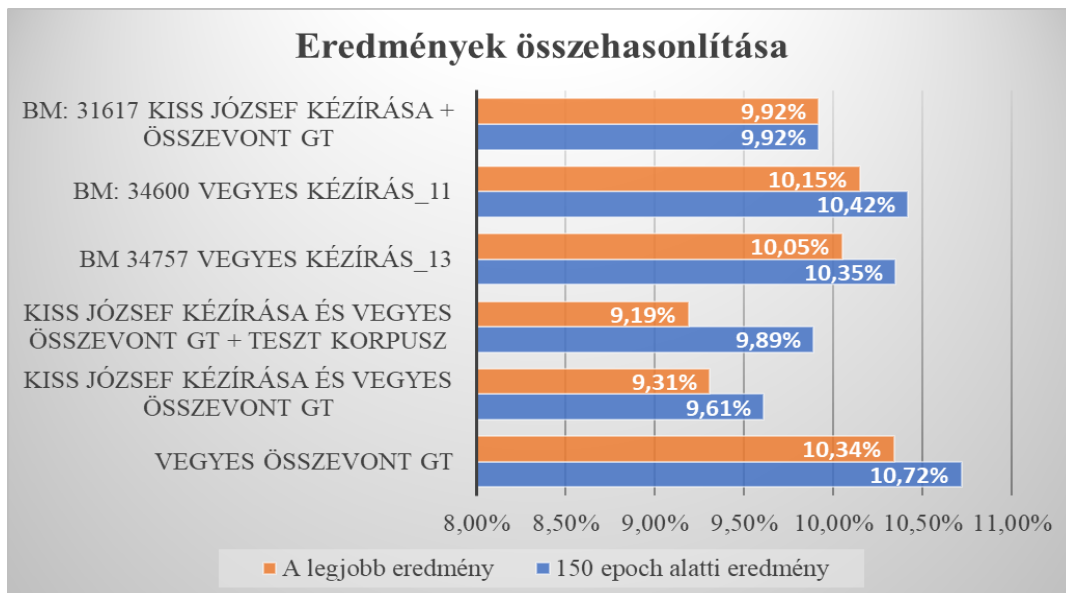
## 2.5. Az eljárások összehasonlítása

A Base Modelles és Ground Truth-os eljárásokat nehéz összehasonlítani, mert mechanizmusukból fakadóan más-más esetben hoznak jó eredményt. A diagrammon is látszik, hogy az alapmodellel rendelkező modell kevesebb epoch alatt hozott jó eredményt és egy bizonyos mennyiségű ismétlés (150 epoch) felett már túltanulta magát, így nőni kezdett a hibaszázalék aránya. A GT-s modellek épp ellenkezőleg, eleinte sokkal rosszabb eredményeket hoztak és csak 150 epoch után kezdett az érték 10% alá csökkenni.



5. ábra. Az eljárások összehasonlítása.<sup>14</sup>

150 epoch-nál még mindegyik modell értékelhető eredményeket hozott és összehasonlítva ezeket már körvonalazódnak a végső eredmények is. A Base Modellel rendelkező modell itt érte el a maximumát, a Ground Truth-os modellek közül a végén a második és az első helyet cserél, ugyanis 250 epochon 0,12%-kal jobban teljesít a TK beépítése utáni adattöbblet miatt.



6. ábra. Az hibaszázalékok összehasonlítása 150 epoch-nál.

<sup>14</sup> 250 epoch-nál 0%-os érték szerepel, mert az egyvel korábbi szakaszon is romló tendenciát mutatott a modell, így utána nem futtattuk le a betanítást.

### 3. Kitekintés

Terveink közé tartozik, hogy folytatjuk az automatikus kézírásfelismerés kutatását és tesztelését. A következő lépés az, hogy az OSZK-ban található Kiss József-levelezés elemeit is beépítsük a modellbe, hogy minél inkább rátanítsuk azt a vegyes kézírás felismerésére.

A legjobb eredmény érdekében a különböző magyar nyelvű projekteknek, amelyek automatikus kézírásfelismertetést használnak, a jövőben össze kell fogniuk, hogy a saját korpuszokon betanított modelleket egymásba építve egyre általánosabb érvényű automatikus kézírásfelismertető eszköz jöjjön létre a magyar nyelvű források kutatható közzétételéhez. Meg kell találni azokat az eszközöket, amelyek segítségével a közgyűjteményekben rejtőző kéziratos kincsek a digitális térben hozzáférhetővé, feldolgozhatóvá, kutathatóvá válnak. A most nyilvánosságra hozott első magyar kézírásfelismerő modell ennek a folyamatnak volt fontos mérföldköve.

### 4. Felhasznált irodalmak

Alvermann, Dirk. *Use Case: "Model Booster."* Hozzáférés: 2022. június 23.

<https://rechtsprechung-im-ostseeraum.archiv.uni-greifswald.de/use-case-model-booster/>.

European Commission CORDIS Research Results. *Recognition and Enrichment of Archival Documents | READ Project | Fact Sheet | H2020.* Hozzáférés: 2022. június 23.

<https://cordis.europa.eu/project/id/674943>.

Heigl, Elisabeth. *CER? Don't Worry!*, Hozzáférés: 2022. június 23.

<https://rechtsprechung-im-ostseeraum.archiv.uni-greifswald.de/cer-dont-worry/>.

READ-COOP. *Base Models.* Hozzáférés: 2022. június 23.

<https://readcoop.eu/glossary/base-models/>.

READ-COOP. *Ground Truth.* Hozzáférés: 2022. június 23.

<https://readcoop.eu/glossary/ground-truth/>.

READ-COOP. *Our Story.* Hozzáférés: 2022. június 23. <https://readcoop.eu/our-story/>.

READ-COOP. *Transkribus Hungarian Handwriting 19th–20th Cent.* Hozzáférés: 2022. június 23. <https://readcoop.eu/model/hungarian/>.

Szűcs, Kata Ágnes. „Automatikus kézírás-felismertetés Kiss József levelezésén”. In *Online térben az online térért: Networkshop 30: országos online konferencia.* 2021. április 6-9. Eötvös Loránd Tudományegyetem, 73–80. HUNGARNET Egyesület, 2021.

<https://doi.org/10.31915/NWS.2021.8>.

*Transkribus users | Facebook.* Hozzáférés: 2022. június 23.

<https://www.facebook.com/groups/614090738935143>.