

## Szerkesztői környezet TEI-alapú szövegkiadásokhoz

### An Editor Framework for Digital Scholarly Editions in TEI

Mihály Eszter  
Országos Széchényi Könyvtár, Digitális Bölcsészeti Központ (OSZK DBK)  
[mihaly.eszter@oszk.hu](mailto:mihaly.eszter@oszk.hu)

Micsik András  
Számítástechnikai és Automatizálási Kutatóintézet, Elosztott Rendszerek Osztály (SZTAKI DSD)  
[micsik@sztaki.hu](mailto:micsik@sztaki.hu)

#### Absztrakt

Az OSZK Digitális Bölcsészeti Központ által fejlesztett új platform, a dHUpla (Digital Humanities Platform - dhupla.hu) elsősorban digitális szövegkiadások publikálására jött létre. Ennek háttérében a szövegkorpuszok előállításához egy teljes szerkesztőségi rendszer kialakítására volt szükség, amely egyrészt felhasználóbarát módon teszi lehetővé a szerkesztést, másrészt funkcióival támogatja a szövegek jelölőnyelvi kódolását. Ezt a felületet a SZTAKI és az OSZK közösen fejleszti egy XML szerkesztő bővítményeként. A szövegek kódolás alapja a TEI (Text Encoding Initiative) szabványa, a keretrendszer e nemzetközi ajánlás bonyolult konstrukcióinak bevitelét, többek között kontextus-menüvel, beszúrható mintákkal, Schematron-validációval segíti. A fejlesztésben központi szerepet játszanak a szöveg feldolgozását támogató további eszközök is: névterekkel, adatbázisokkal való összekapcsolás, MI-alapú névelem-felismerés, valamint különböző automatizált műveletek, úgymint PDF- és képkonverzió, vagy adatvizualizáció támogatása. A manuálisan és a gép által végezhető részfolyamatok minden esetben kiegészítik egymást, megteremtve ezzel egy minőségi, a digitális közeg adta lehetőségeket kiaknázó szövegkiadási módszert.

**Kulcsszavak:** digitális szövegkiadás, TEI, XML szerkesztőfelület

#### Abstract

The Digital Humanities Centre (DBK) of the Hungarian National Library has developed a platform called dHUpla (Digital Humanities Platform - dhupla.hu) for publishing digitized text editions. For this, the creation of a complete editing environment was necessary to support editors and digital humanists and to help them in the input and validation of correct TEI XML encoding. This environment was implemented as an XML editor extension jointly by DBK and SZTAKI (Institute for Computer Science and Control). The extension contains custom toolbars, templates, Schematron validation and a set of scripts to automate steps of the conversion to TEI XML. Scripts support named entity recognition and linking, PDF generation, extraction of data for visualizations and other task automations. These operations combine automated execution with manual supervision to reach high quality TEI production.

**Keywords:** digital scholarly editions, TEI, XML editor

## Bevezetés

Az Országos Széchényi Könyvtár Digitális Bölcsészeti Központ (DBK) egyik elsődleges feladata egy korszerű online platform fejlesztése a közgyűjteményekben őrzött szöveges források kezelésére, amely egységes kutatói környezetet jelent az irodalomtudomány, a nyelvtudomány, és más humán tudományok számára. A dHUpla<sup>1</sup> (Digital Humanities Platform) 2021 óta üzemel (2021 decemberéig a Petőfi Irodalmi Múzeum, majd az Országos Széchényi Könyvtár szolgáltatásaként), és számos szövegkiadást tesz nyilvánosan elérhetővé, amelyhez entitástár, valamint kreatív tartalmak, vizualizációk is társulnak. E szolgáltatások háttérében a szövegkorpuszok előállításához egy teljes szerkesztőségi rendszer kialakítására volt szükség, amely egyrészt felhasználóbarát módon teszi lehetővé a szerkesztést, másrészt funkcióival támogatja a szövegek jelölőnyelvi kódolását. Ezt a felületet a SZTAKI és az OSZK közösen fejleszti egy XML-szerkesztő bővítményeként. Az alábbiakban röviden bemutatjuk a dHUpla funkcióit, majd rátérünk a szerkesztői felület részletes ismertetésére.

## A dHUpla infrastruktúra

A dHUpla elsődleges célja, hogy a kutatók, olvasók számára egységes felületen, filológiai igényességgel hozzáférhetővé tegye a magyar kulturális örökség különböző intézményekben őrzött, eddig ismeretlen vagy méltatlanul elfeledett szöveges tartalmú, elsősorban kéziratos forrásait. Emellett olyan szerkesztőségi keretrendszert is kínál a tartalmak digitalizálásához, amely biztosítja a források egységes és színvonalas feldolgozását.

A platform alapjait és felépítését korábban már részletesen bemutattuk [1]. A rendszer rugalmas és moduláris, a tartalom előállítása többféle úton történhet. Az átírt szövegek a nemzetközileg támogatott TEI (Text Encoding Initiative) szabvány szerint annotált XML formátumban készülnek, amely a nemzetközi integráció mellett lehetővé teszi többek között a dokumentumok gépi feldolgozását és értelmezését (linked open data), szemantikus hálók kiépítését (semantic web), adatgazdagítást (data enrichment), illetve a távoli olvasás (distant reading) különböző aspektusait. A szolgáltatás kapcsolatot létesít különböző névterekkel, bibliográfiai forrásadatbázisokkal, illetve nyelvi elemző szoftverekkel. Mindezek segítségével a legkülönbözőbb elemzések, korpuszlekérdezések, adatvizualizációk válnak megvalósíthatóvá.

Az infrastruktúra középpontjában a git verziókövető szoftver áll, amely végül teljes mértékben kiváltotta egy XML-adatbázis használatának szükségességét, nagyban leegyszerűsítve a rendszer használatát, karbantartását és fejlesztését. A dHUpla git-ben lévő források (szöveg, programkód) alapján publikál, így a git repository-k birtokában bárhol újraépíthető a teljes dHUpla honlap. A projektek mind önálló git repository-ban vannak, a publikáláshoz szükséges transzformációt docker containerek végzik, minden egyes projekthez meg lehet adni saját ún. buildert, amelyekben tetszőleges programnyelvet lehet használni. A HTML tartalmon túl Apache Solr indexfájl is előállítunk, amelynek segítségével a legkülönbözőbb facettás keresések is lehetővé válnak.

A publikációs felületen több módon lehetőség nyílik az átírt szöveg és az eredeti facsimile együttes vizsgálatára, a digitális objektumok különböző szempontú rendezésére, szűrésekre elvégzésére. Az egyes gyűjtemények megjelenésének, funkcióinak konfigurálása egyszerű szöveges fájlokban (yaml) történik.

---

1 [dhupla.hu](http://dhupla.hu)

A kéziratok átírása során ezenkívül olyan kézírásfelismerő-modell épül, amely folyamatosan bővülve alkalmassá válik a magyar nyelvű kézírások automatikus felismertetésére, azaz mesterséges intelligencián alapuló gépi feldolgozására (Handwritten Text Recognition).

## A TEI szerkesztői felület

A dHUpla oldalára szánt szövegeket a DBK által készített TEI XML specifikációk szabályai szerint kell létrehozni, mivel a TEI igencsak rugalmas XML formátum, és több különböző ábrázolási mód is választható ugyanarra a szövegelemre. A TEI XML-ek részletes szerkesztésére az Oxygen XML Editor<sup>2</sup> programot használjuk. Az Oxygen XML Editor (röviden továbbiakban Oxygen) beépített TEI támogatással rendelkezik, de ez a már említett teljesen általános, korlátozások nélküli TEI-n alapul. Az Oxygenen belül lehetőség van ún. frameworkök kifejlesztésére, amely egy személyre szabott szerkesztői felületet tud nyújtani, eszköztárakkal, menüvel, specializált megjelenítéssel és ellenőrzési lehetőségekkel.

A frameworkben alkalmazott megoldásokból elsőként magát a szövegszerkesztő ablakot vesszük sorra. Itt háromféle nézetből lehet választani, ebből számunkra kettő a lényeges: az első magát az XML kódolást mutatja, ahol a rideg valóság tárul elénk, és ellenőrizhetjük, illetve változtathatjuk az XML-t. A szerzői nézet ezzel szemben egy CSS-sel barátságossá alakított, olvasható nézetet kínál, ahol az elemek szerkesztése vezetett módon lehetséges (1. ábra).

» Kedves ▾ » **fiam**!

» ▾ » S o k á i g ▾ töprengtem azon, hogy ennek a ▾ » JENŐNEK▾ – ha már itthon lebzsel – miképpen lehetne valami kereset forrást szerezni? Végre is azon gondolat érlelődött meg lelkemben,

» **Betoldás:** **Felelős személy:** 324324 **Írószköz:** ceruza ▾ **Felelős egységesített név:** Móricz Zsigmond **Ok:** beszúrás ▾

**Típus:** -üres- ▾ **Hely:** fölé ▾ **Forrás:**  **Ugyanaz a művelet:**  **Művelet azonosító:**

**miszerint**

legjobb, legcélszerűbb lenne, ha oly munkát találhatnánk részére, mit it▾ » t▾hon is elvégezhet.

» Ily munka pedig – nézetem szerint – csak másolások, ▾ » címírások▾

»

stb. végzése lehetne. Ez annális inkább ▾ » **Bizonytalan olvasat:** **Ok:** sérült ▾ qudrálna▾ neki, mert mellettök gondolkozni nem kell és jó

» **foljó**

írása van

1. ábra: Az ún. szerzői nézetben segítséget kapunk a speciális elemek kitöltéséhez

Az 1. ábrán megfigyelhetjük, hogy a speciális szövegelemek előtt álló » jelre kattintással kinyitható és becsukható az elem tartalma. Az ábrán egy <add> és egy <unclear> elem attribútumait láthatjuk.

A metaadatok beírását részletesen kidolgozott űrlapok segítik (2. ábra). A piros aláhúzás hibás kitöltést jelez, a magyarázat a legelső sorban olvasható. Az ellenőrzést Schematron<sup>3</sup> szabályokkal végezzük.

2 [https://www.oxygenxml.com/xml\\_editor.html](https://www.oxygenxml.com/xml_editor.html)

3 <https://www.schematron.com/>

TEI teiHeader fileDesc sourceDesc msDesc msPart physDesc objectDesc supportDesc

**További azonosításra szolgáló információk a levél egyes részeihez**

**Fizikai leírás**

**Objektum leírás**

**Fizikai hordozó adatai**

**Anyag:**

**Terjedelem**

**Mennyiség:**

**Mértékegység:**

**Méret**

**Magasság:**

**Szélesség:**

**Mértékegység:**

**Állapotleírás**

Szöveges információ (bekezdés):

Fekete tintairás.

A terjedelem (extent) elemben szerepelnie kell egy mennyiség (measure) elemnek piece vagy folio mértékegység (unit) értékkel.

2. ábra: A TEI fejléc kitöltése és ellenőrzése

A 3. ábrán láthatóak a szerkesztők munkáját megkönnyítő eszköztárak. A felső sorban található menüvel a gyakori TEI elemeket lehet beilleszteni a szövegbe. A középső sor főleg a fejléc, vagyis a metaadatok kitöltéséhez nyújt segítséget. A legalsó sorból feldolgozási parancsokat lehet indítani. Ebben a sorban nem látszik az összes lehetséges művelet, van még keresés több más névtérben is, illetve PDF készítés mint menüpont.



3. ábra: Az általunk készített eszköztárak

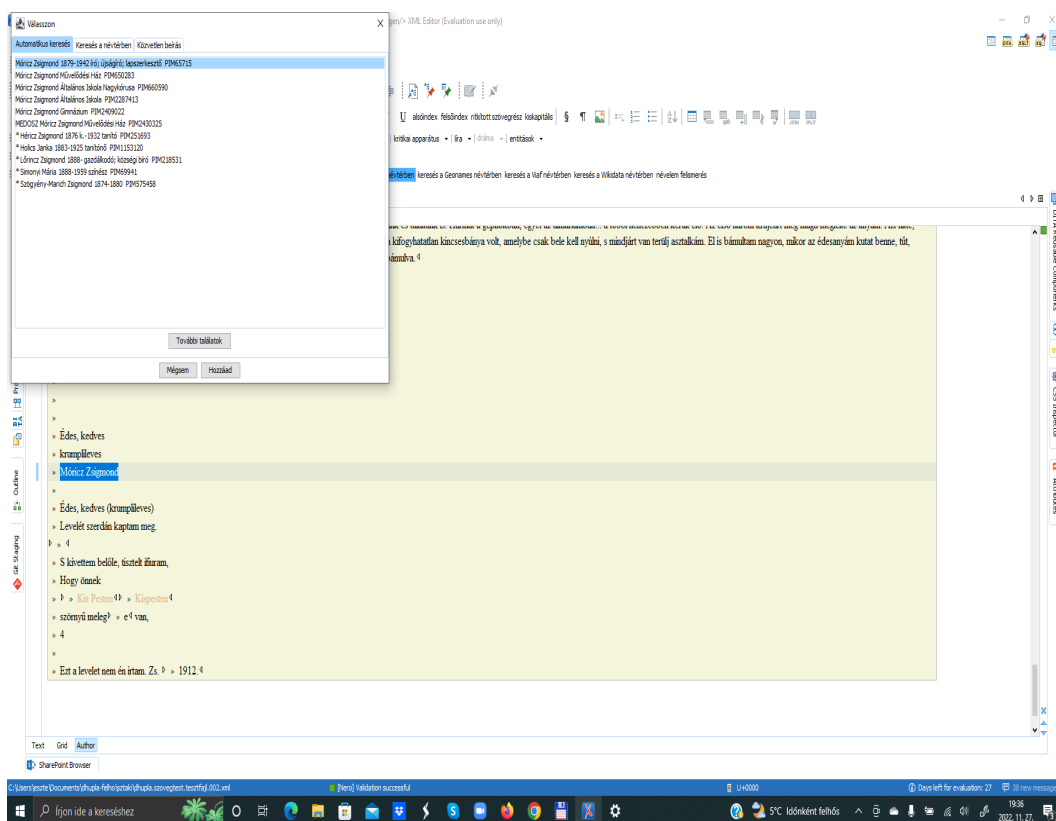
## Példa munkafolyamat

A kéziratok kézi és automatikus gépi átírásához, valamint az ún. text-image linking (kép és szöveg összekötése) elvégzéséhez a Transkribus programot használjuk, amely TEI XML formátumban is tud exportálni, azonban ez a formátum nem szabványos, és a DBK előírásainak sem felel meg. Az Oxygen eszköztárból indítható Transkribus export konverzió segítségével azonban megfelelő formára tudjuk hozni a TEI XML-t.

Ezután a szerkesztőben ellenőrizhetjük az XML-t, javíthatjuk a kapott Schematron hibajelzéseket, és egyéb finomításokat tudunk tenni a szövegjelölésekben. Ehhez a fázishoz kapcsolódik a TEI fejléc kitöltése, amelyhez egy teljes sablont tudunk egyetlen gombnyomással beilleszteni, majd a kapott űrlapokat kitölteni.

A szöveg feldolgozásában segít a megemlített névelemek (személyek, települések, szervezetek stb.) bejelölése. Erre saját fejlesztésű névelem-jelölőket lehet használni, amely a HuSpaCy [2] nyelvi feldolgozó eszköz segítségével beilleszti a megfelelő TEI XML elemeket (pl. <persName>) a talált nevekhez. Ezek a nevek még azonosítatlanok, vagyis nincsenek egy gondozott névtérben a megfelelő névelemhez kötve. Erre a célra az eszköztár keresőgombokat tartalmaz, amely a kijelölt szöveget megkeresi az adott névtérben (pl. PIM<sup>4</sup>, GeoNames, Wikidata), és a szerkesztő választásának megfelelően beírja az XML-be a névtér azonosítót (4. ábra).

4 <https://opac-nevter.pim.hu/search>



4. ábra: Keresés a szövegben kiválasztott névre az Oxygen szerkesztőfelületén

Ha elkészült a TEI XML feldolgozás, az elemzési és terjesztési feladatokat is támogatja a framework. Az XML-ből például PDF-et tudunk előállítani, amely nem tartalmazza ugyan az XML-ben kódolt összes részletet, viszont jól olvasható és megosztható.

Dokumentumok csomagjaiból (pl. egy író levelezése) a Metaadatok kinyerése funkció segítségével táblázatokat kapunk, amelyek aztán különböző összesítések és vizualizációk alapjául szolgálhatnak. Az egyik ilyen vizualizációs eszköz hamarosan elkészül: a levelek feladását és kézhezvételét mutatja térképen.

## Összefoglalás

A digitális szövegkiadások TEI alapú publikálása rengeteg új lehetőséget ad a digitális bölcsészeknek, amelyből párat jelen cikkben említettünk. Nehéz viszont odáig eljutni, hogy a szövegek jó minőségű TEI XML formátumba kerüljenek, mivel a folyamat sok szakértelmet és szakértői munkát igényel. Ennek a feldolgozási munkának a megkönnyítésére készítettük el a bemutatott Oxygen XML Editor framework megoldást, amely számos hasznos funkcióval segíti a szerkesztést, de ezen felül lehetőséget teremt a TEI-n belüli saját kódrendszerek megalkotására és az elkészített TEI dokumentumok e szerinti ellenőrzésére, és egy alaminőség biztosítására a konverziók végén. A jelenlegi framework a levelezések XML kódolására fókuszál, de látjuk már a lehetőséget többféle „szakosodott” leírási mód (pl. drámák, versek, naplók) egyidejű támogatására is. A keretrendszer fejlesztése jelenleg a házon belüli alkalmazás és tesztelés fázisában van.

## Irodalomjegyzék

- [1] Mihály, Eszter (2022) *Mi az a dHUpla?: A Digitális Bölcsészeti Platform bemutatása*. In: Valós térben - Az online térért, Networkshop 31: országos konferencia. 2022. április 20–22. Debreceni Egyetem. Kiadja a HUNGARNET Egyesület az MTA Könyvtár és Információs Központ közreműködésével, Budapest, pp. 345–358. ISBN 978-615-82243-0-7 DOI: [10.31915/NWS.2022.44](https://doi.org/10.31915/NWS.2022.44)
- [2] Orosz György, Szántó Zsolt, Berkecz Péter, Szabó, Gergő, Farkas Richárd (2022). HuSpaCy: an industrial-strength Hungarian natural language processing toolkit. In XVIII. Magyar Számítógépes Nyelvészeti Konferencia.



The background is a complex digital artwork. It features a grid of squares, each containing a different texture or color, ranging from warm oranges and yellows on the left to cool blues and teals on the right. A bright, glowing light source is positioned in the center, creating a lens flare effect that radiates across the grid. The overall composition is symmetrical and has a high-tech, futuristic feel.

# ÚJ TECHNOLÓGIÁKKAL, ÚJ TARTALMAKKAL A JÖVŐ DIGITÁLIS TRANSZFORMÁCIÓJA FELÉ

32. Networkshop: országos konferencia

2023. április 12–14.

Pannon Egyetem, Veszprém



# ÚJ TECHNOLÓGIÁKKAL, ÚJ TARTALMAKKAL A JÖVŐ DIGITÁLIS TRANSZFORMÁCIÓJA FELÉ

**32. Networkshop: országos konferencia**

2023. április 12–14.  
Pannon Egyetem, Veszprém

Szerkesztette: Tick József, Kokas Károly, Holl András

HUNGARNET Egyesület  
Budapest, 2023





Szerkesztette: Tick József, Kokas Károly, Holl András

Tipográfia és tördelés: Vas Viktória

Workshop

2023. április 12–14. Pannon Egyetem, Veszprém konferencia előadásainak közleményei

ISBN 978-615-82243-1-4

DOI: [10.31915/NWS.2023](https://doi.org/10.31915/NWS.2023)

Kiadja a HUNGARNET Egyesület  
az MTA Könyvtár és Információs Központ közreműködésével

Budapest

2023

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

## TARTALOMJEGYZÉK

Előszó.....	5
Király Sándor, Balla Tamás: Flipped classroom az sqlsuli.hu-ban.....	7
Wirágh András: Abaújszántótól Zombolyáig. Megjegyzések egy új sajtóadatbázishoz .....	14
Albert Ágota Katalin: Az EGT-tagállamok adatvédelmi felügyeleti hatóságainak szankcionálási gyakorlata az oktatási szektorban a GDPR alkalmazása óta .....	19
Simon András: Digitális dokumentumok gyűjteménykezelési gyakorlatának támogatása a digitális tartalmak számossága, mérete és féleségeik vizsgálatával .....	24
Bódog András: Az Annif gépi tárgyszavazó rendszer magyarországi adaptációjának feltételei és lehetőségei .....	31
Dezső Krisztina: A Pécsi Egyetemtörténeti Gyűjtemény online adatbázisai és digitális gyűjteményei .....	36
Ungváry Rudolf, Király Péter: Nemzeti könyvtárak és az OSZK MARC21 állományainak összehasonlító elemzése néhány adatmező alapján .....	42
Szemes-Révész Enikő Evelin: Kapocs a tudáshoz – A könyvtár szerepe a civilek és a tudomány kapcsolatában .....	50
Tóth Zoltán: Az RO-Crate alapú kutatási objektum csomagolás keretrendszere az ELKH ARP platformban .....	54
Király Roland, Király Sándor, Palotai Martin Marcell: Neurális hálózatok oktatási alkalmazását támogató keretrendszer Virtual (VR) és Augmented Reality (AR) eszközökkel .....	60
T. Nagy László: Mesterséges intelligencia, multimédia, tanulástámogatás .....	69
Horváth Péter: Egy automatikusan generált rímshótár fejlesztése és a magyar kanonikus költészet rímshavainak néhány jellemzője .....	77
Héjja Balázs, Tóth-Jávorka Brigitta, Tóth Máté: Digitális tartalomfejlesztés közkönyvtári környezetben .....	85
Koczka Ferenc: Szemelvények egy felsőoktatási rendszer informatikai védelmének tapasztalataiból .....	91
Bolya Mátyás: A digitális gyűjtésrekonstrukció lehetőségei: az Ethiofolk projekt .....	99
Dobás Kata, Sidó Zsuzsa, Szabó-Reznek Eszter: A Kolozsvári Állami Magyar Színház jelmezterveinek digitalizációja és felvitele az ITIdata adatbázisba .....	108
Köpösdí Zsuzsa: H5P-ben létrehozható interaktív és adaptív tananyagok .....	116
Fülöp Tiffany, Molnár Tamás, Hoczopán Szabolcs: Komplex kutatástámogató szolgáltatási portfólió az SZTE Klebelsberg Könyvtárban .....	122
Vass Johanna: Az Open Science könyvtári vonatkozásai .....	129
Antal Péter, Czeglédi László: A digitális oktatás módszertana a gyakorlatban .....	135
Máray Tamás: A szuperszámítástechnika mint európai stratégiai ágazat .....	143
Frankó Máté, Zeller Rozália: Szoftveres Cutter-keresés az SZTE Klebelsberg Könyvtárban .....	151
Zsiborács Judit, Dési Ádám Dániel, Nagy Attila Árpád, Urbán Katalin: Tudományometriai műhely könyvtári környezetben .....	157



Palkó Gábor, Szekrényes István, Bobák Barbara: A Digitális Örökség Nemzeti Laboratórium webszolgáltatásai automatikus kézírás-felismertetéshez .....	164
Szűcs Kata Ágnes: Adatvizualizációs lehetőségek a bölcsészettudományban .....	170
Leitgéb Mária: A BME Építészettörténeti és Műemléki Tanszék repozitóriuma .....	178
Mihály Eszter, Micsik András: Szerkesztői környezet TEI-alapú szövegkiadásokhoz .....	186
Dobás Kata, Fellegi Zsófia, Palkó Gábor: A kis gömböc meséje - az ITIdata irodalomtudományos adatbázis fejlesztése 2022–2023-ban .....	192
Alföldi István, Szemigán Dorottya Henrietta, Palkó Gábor, Fellegi Zsófia: Kutatói e-mail hagyaték archiválása és feldolgozása .....	199