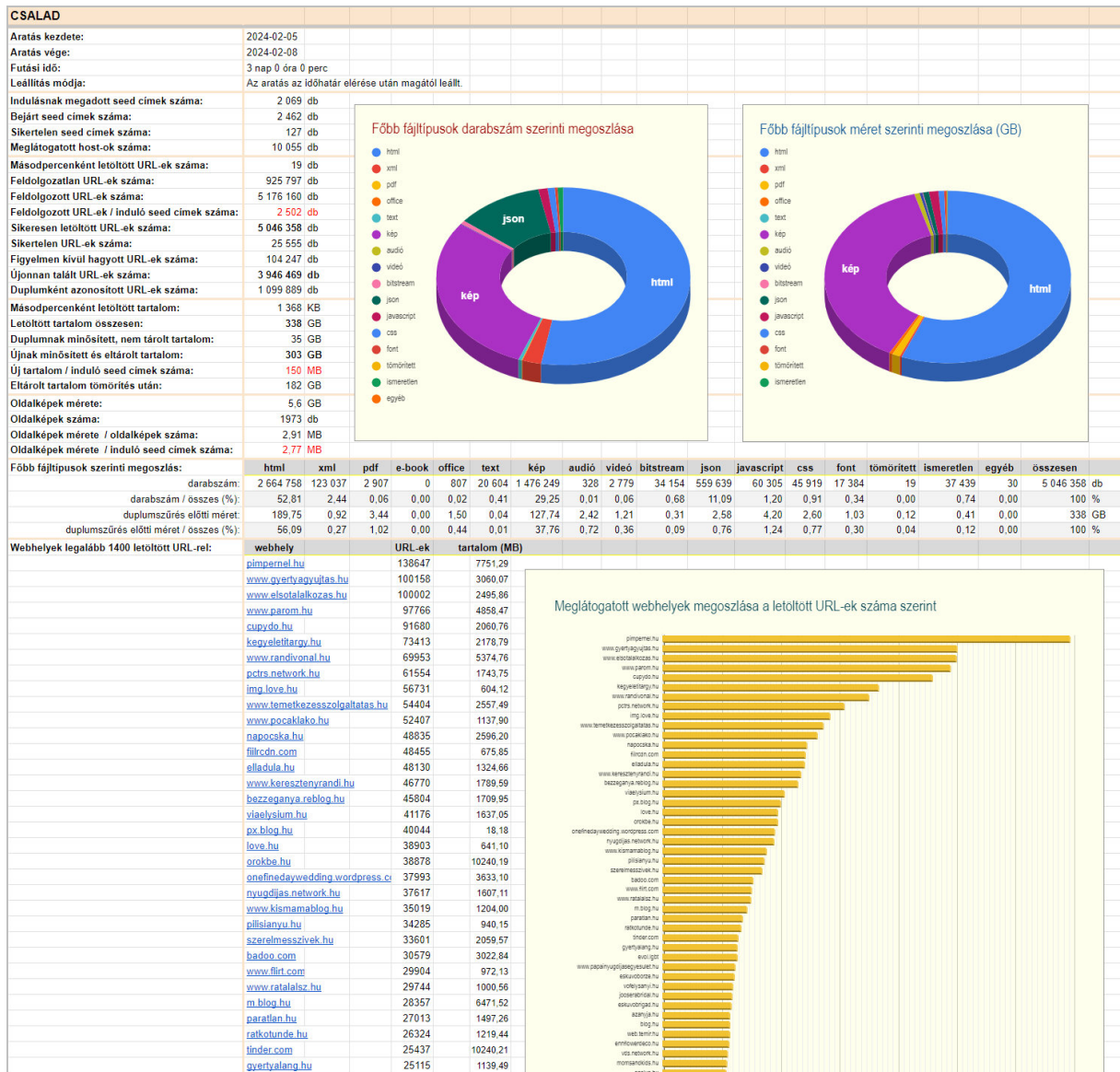


Az OSZK Webarchívum 2024. februári hírei

Archiválás és nyilvántartás

A januárban csaknem 24 napig futó webtér szintű jobok mellett nem akartuk elindítani az eredetileg az év első hónapjára beütemezett tematikus aratásokat is, hogy ne terheljük túl a szervert, ezért ezeket a februáriakkal együtt, a korábbi gyakorlathoz képest sűrűbben indítottuk el. A hónap végére sikerült behozni a tavalyi technikai problémák miatt kialakult lemaradást, így márciustól újra az eredeti ütemben folynak majd a tömeges mentések. A februári aratások felsorolása a beszámoló végén olvasható, részletes statisztikáik pedig a szintén ott található linken nézhetők meg.

Február 5. és 8. között első alkalommal mentettük önálló részgyűjtemény formájában a CSALAD kódnevű címlistában szereplő 2 ezer webhelyet. Ez a válogatásunk a következő főbb témakörökből áll: párkapcsolat, esküvő, házasság, szexualitás, erotika, születés, anyaság, család, gyerekevelés, idősök, halál. A három napnyi futásidő alatt 5 millió URL-t mentett le a robot 338 gigabájt össz méretben, amiből 303 volt az új tartalom. Elkészítettük a CSALAD és a tavaly összeállított ELETMOD részgyűjtemények XML formátumú metaadat rekordjait is, melyeket feltöltöttünk a honlapunkra.



A „Párkapcsolat, család” nevű részgyűjtemény első aratásának statisztikája

Létrehoztunk egy újabb válogatást MEZGAZ néven, amiben a mezőgazdasággal és az élelmiszeriparral kapcsolatos honlapokat és blogokat gyűjtjük. A címlista már több mint 2,5 ezer tételt tartalmaz, első aratására március elején kerül majd sor. Minden ilyen gyűjtési időszakban, így most is, több száz tétellel gyarapodott a nyilvántartásunk olyan címekkel, amelyeket eddig nem ismertünk és valamilyen más témakörbe vagy műfajba tartoznak. Az elektronikus időszaki kiadványok táblázatába például 114 URL került be február folyamán.

Az elmúlt hetekben 44 olyan webhely mentéseivel bővült a webarchívum nyilvános része, amelyek szolgáltatáshoz nem kell egyedi engedélyeket kérni. Főként települési és vármegyei önkormányzatok, könyvtárak és művelődési házak, illetve nemzeti parkok honlapjainak archivált verziói kerültek ki a publikus felületre, ahol három-négy megjelenítővel is visszanezhetők ezek. Az első tesztmentés után minden esetben megpróbáltuk optimálisan paraméterezni a Web Curator Tool keretrendszerben a Hertirix robotot és néha a PC-n futtatható HTTrack programmal is készítettünk egy alternatív mentést, aminek az eredményét WARC formátumra konvertáltuk. Ennek ellenére az archív változatokban vannak hibák és hiányok, amelyek részben az archiváló és/vagy a megjelenítő szoftverek képességeivel magyarázhatók, részben pedig az egyes fájlok és az egyes aratások méreteire, a futási időre és a linkek követési mélységére beállított korlátokkal. Tanulságos megfigyelni, hogy ugyanaz a mentés mennyire másként jelenik meg a különböző visszanező szoftverekben, illetve hogy mennyire függ az archiválás sikeressége az eredeti webhelyen használt technikai megoldásoktól.

Technikai ügyek

A tesztszerveren a Web Curator Tool beüzemelésével kapcsolatos problémák okainak felderítése folyt. Az ékezetkódolási hibát végül sikerült elhárítani, de a WCT és a Heritrix közötti kommunikáció még nem stabil. A következő feladat a WCT tesztadatokkal való feltöltése lesz.

Újabb megbeszélések voltak az érintett informatikusokkal arról, hogy a kozterkep.hu oldalról begyűjtött mintegy 265 ezer, Creative Commons licenc alatt közzétett fénykép metaadatait hogyan kellene betölteni a Digitális Képtár adatbázisába és hogy milyen további részfeladatok kellene még ahhoz, hogy ezek a köztéri műalkotásokat ábrázoló fotók megjelenhessenek a DKA-ban. Megállapodtunk abban is, hogy a konvertálást és az adatok betöltését először az új MEK fejlesztésére szolgáló, de a DKA rendszerének másolatát is tartalmazó tesztszerveren végezzük el, hogy ne akadályozzuk a hagyományos katalogizálást és az esetlegesen felmerülő hibákkal ne veszélyeztessük az éles szolgáltatást.

A Nyelvtudományi Kutatóközponttal való együttműködésünk keretében elkészültek és az átadás idejére egy ideiglenes depozit tárhelyre kerültek azok az index- és checksum-fájlok, melyek a 2018 óta végzett webtér aratások technikai metaadatait tartalmazzák. Az OutbackCDX programmal készült indexek összmérete 321 gigabájt lett, a duplumszűrésre használt ellenőrző összegeket tartalmazó JDB adatbázisok pedig 544 gigabájtot foglalnak el egybecsomagolás és tömörítés után.

Ismeretterjesztés

Február 12-én online „gyorstalpaló” tanfolyamot tartottunk a Vörösmarty Mihály Könyvtár informatikusának a WCT használatáról, mivel a székesfehérvári kollégák is ezt a keretrendszert fogják használni a Fejér vármegyei regionális webarchívumhoz.

A honlapunk „A projektről” menüpontja alatt található „Előadások, prezentációk, publikációk” oldalon egy külön szekciót alakítottunk ki az OSZK webarchiválási projektjéről (is) szóló blogbejegyzésekről. A lista jelenleg 16 posztot tartalmaz, közülük kettő angol nyelvű, melyek az International Internet Preservation Consortium, illetve a Society of American Archivists blogjaiban jelentek meg.

Összeállítottunk egy kérdőív tervezetet, amivel azt szeretnénk felmérni, hogy a megújítás alatt álló, „Az internet archiválása mint közgyűjteményi feladat” elnevezésű tanfolyamunkkal kapcsolatban milyen igényei vannak a leendő résztvevőknek, melyek azok a résztemák, amelyek iránt leginkább érdeklődnek.

Az OSZK podcast csatornáján megjelent egy 40 perces interjú a webarchívum munkatársával, Drótos Lászlóval „Az internet úttörője” címmel.

Az elmúlt hetekben lefutott tematikus aratások

Könyvtárak, levéltárak, múzeumok és galériák (2131 db seed URL)
Irodalom, irodalomtudomány és -történet (1431 db seed URL)
Kulturális intézmények, művelődési házak, rendezvényhelyszínek (966 db seed URL)
Természet- és műszaki tudományok, szakterületek (2253 db seed URL)
Képző-, előadó-, zene- és filmművészet (8733 db seed URL)
Közoktatás és egyéb képzések (7239 db seed URL)
Bölcsészet- és társadalomtudományok, szakterületek (6131 db seed URL)
Életmód, szabadidő, hobbi (8597 db seed URL)
Kormányzat, önkormányzatok, politikai és civil szervezetek (7020 db seed URL)
Történelem, hely- és családtörténet (1465 db seed URL)
Podkasztkok (4721 db seed URL)
Elektronikus periodikák (8737 db seed URL)
Média, sajtó, műsorszórás (937 db seed URL)
Könyv- és egyéb kiadók, kereskedők (1477 db seed URL)
Párkapcsolat, család (2069 db seed URL)

Az egyes aratások részletes statisztikai adatai a <https://webarchivum.oszk.hu/szelektiv-aratasok/> weblapon nézhetők meg. A projekt hírei a <https://webarchivum.oszk.hu/a-projektrol/hirek-esemenyek/> oldalon kísérhetők figyelemmel. Kapcsolati cím: webarchivum@oszk.hu