

MTA SZTAKI DSD

Department
of Distributed
Systems

KOPI
A fordítási plágiumok keresője

Pataki Máté
Kovács László

- n MTA SZTAKI Elosztott Rendszerek Osztály
- n 1995. óta létezik
- n 12 teljes állású munkatárs, és diákok
- n Kutatás, fejlesztés, (online) szolgáltatások
- n Munkák eloszlása:
 - n 80% EU-s k+f pályázatok
 - n 15% Hazai pályázatok
 - n 5% Belső projektek és szolgáltatások
- n Három fő terület:
 - n Digitális könyvtárak és archívumok
 - n Csoportmunkát támogató technológiák
 - n Webes rendszerek

World Wide Web

Government Portals
Infrawebs

Brein E-VOTING
E-ADMINISTRATION
Web4Groups

Workflow Promóció
Forum
Collaborative Filtering

Csoportmunka

SZTAKI Szótár
KOPI

GeneSyS
StreamOnTheFly
EUTIST-AMI
Abilities

CORES

Select Rating

Digitális Könyvtárak

AQUA

HEKTÁR

DELOS NoE 1

DELOS

ORG

DELOS NoE 2

További információk

<http://dsd.sztaki.hu>

- n A plágium probléma a
 - n Felsőoktatás területén
 - n Középiskolában is egyre inkább
 - n Tudományos életben
 - n Digitális könyvtárak számára
 - n Könyvkiadóknak
 - n Cégek esetében is (pl. honlapok tartalma)
 - n Wikipedia

1. Sok a diák
2. Hasznos anyagok az interneten
3. Digitális szakdolgozatok
4. Jó nyelvtudás

n 1-3 → könnyű plagizálás

n Plágiumkeresők

n KOPI

n +4 → fordítási plágiumok

n ???

A KOPI Plágiumkereső ismertetése

- n KOPI Online Plágiumkereső és Információs Portál - internetes **hasonlóság** és **plágiumkereső** szolgáltatás
- n MTA SZTAKI Elosztott Rendszerek Osztály
- n <http://kopi.sztaki.hu/>

sztaki kopi Kopi Online Plágiumkereső és Információs Portál

Szótár KOPI NDA Kereső

Plágiumkeresés és dokumentumkezelés

Feltöltés Dokumentumaim Plágiumkereső Futó keresések

Válassza ki azokat a dokumentumokat amelyekkel plágiumkeresést szeretne végezni.

Cím	Szerző	Feltöltés dátuma	Szűrő bekapcsolása
<input checked="" type="checkbox"/> Cikk 3 PZsP 01b 2% (30 szó) egyezés	-	2005.09.30.	Szerkeszt Részletes
<input checked="" type="checkbox"/> Szabványok a kórházi informatikában - Absztrakt 30% (30 szó) egyezés	-	2005.09.30.	Szerkeszt Részletes
<input checked="" type="checkbox"/> A kutatók kötelesek	IP	2005.09.13.	Szerkeszt Részletes

magyar | english

Betűméret - +
Nagy kontraszt
Sugó

KOPI

Kezdőlap
Fórum

Felhasználó: a

Beállításaim
Üzenetek
Plágiumkeresés

Kilépés
Admin

A KOPI Plágiumkereső története

- n **2001:** elkezdődtek az **alapkutatások**, hogy miként lehetne egy nyelvfüggetlen, magyar nyelven is jól működő plágiumkeresőt elkészíteni
- n **2003:** állami támogatással elkezdődött a KOPI Portál **fejlesztése** (ITEM pályázat, IHM-OM)
- n **2004:** elindult a **publikus plágiumkereső szolgáltatás** magyar és angol nyelven
- n **2006:** számos tanár elkezdi használni a szolgáltatást, **több mint 1000 felhasználó**

A KOPI Plágiumkereső története

- n **2007: fejlesztések az első három év tapasztalatai alapján**
- n **2009: felhasználóink száma már több mint 10 000, egyes egyetemi karok használják már a KOPI Plágiumkeresőt rendszeresen**
- n **2010: új kutatásba kezdtünk, hogy miként lehetne fordítási plágiumokat felismerni és megtalálni**
- n **2011: a világon elsőként beépítettük a KOPI Plágiumkeresőbe fordítási plágiumok megtalálására képes algoritmust, amely a teljes angol Wikipédiában keres**

- n Sok esetben nem szándékos a plagizálás
- n Nem oktatják az egyetemeken a helyes idézés módját
- n Mekkora hasonlóságot várunk el
 - n 0% - nincs irodalomkutatás
 - n 10%
 - n 50%
 - n 100% - egyértelműen plágium
- n Diákok és tanárok egyaránt használják a KOPI Plágiumkeresőt

- n BME, 400-500 diák, 5 feladat, 6 év
 - n 2007: 9 pár, 2009: 4 pár, 2010: 2 pár
- n Statisztika
 - n Közel 20 000 felhasználó
 - n 30 000 dokumentum
 - n Körülbelül 25 000 000 dokumentumrészlet
 - n Ebből 20 000 000 magyar

- n Feladat
 - n Működő szolgáltatás magyaroknak
 - n Az angol eredeti szöveg megtalálása a magyar fordítás ismeretében
- n Egyéb felhasználási területek
 - n Párhuzamos korpusz építése
 - n Létező fordítások keresése
 - n Hírek, cikkek, anyagok terjedésének a vizsgálata
 - n Idézetkereső

- n Test cases for plagiarism detection software, Debora Weber-Wulff, HTW Berlin, 2010
- n 48 különböző plágiumkereső, 42 teszt
- n *The biggest gap in all the plagiarism checkers was the **inability to locate translated plagiarism**. While this is widely expected as the technology to make such detections simply is not there.*

- n CLEF 2010
- n Potthast: Overview of the 2nd International Competition on Plagiarism Detection
 - n *After analyzing all 17 reports, certain algorithmic patterns became apparent to which many participants followed independently. ... In order to simplify the detection of cross-language plagiarism, non-English documents in D are **translated to English using machine translation** (services).*

Irodalom – fordítási plágiumok

- n Európában fontos téma
- n Az algoritmusok nyelvpár-függők
- n Magyar nyelvben három fő akadály
 - n nem kötött szórend
 - n ragozás
 - n jelentős nyelvtani különbség az angol nyelvtől
- n rosszak az automatikus fordítók (erre)

Az új algoritmus

- n Mondatalapú
- n szó, n-szó, tagmondat,
bekezdés, dokumentum



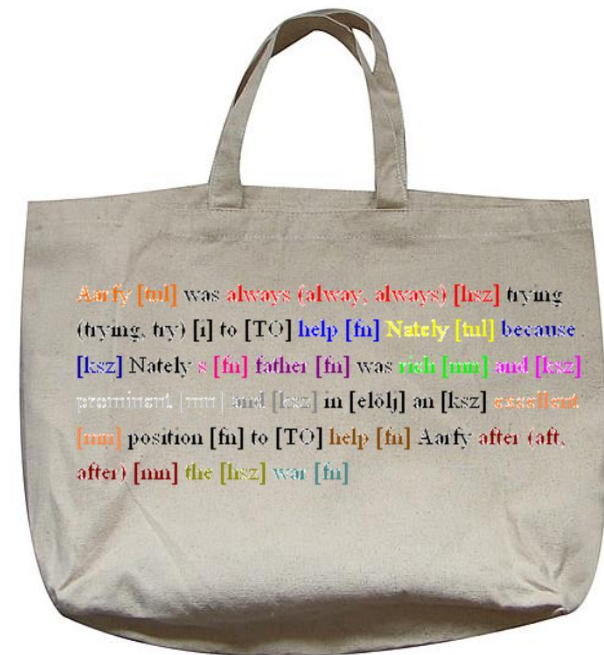
- n Hasonlósági metrika

$$\text{Sim}(x,y) = \min (\alpha \cdot | S_x \cap S_y | - \beta \cdot | S_x \setminus S_y | , \alpha \cdot | S_y \cap S_x | - \beta \cdot | S_y \setminus S_x |)$$

- n Lapos szótár, szószedet

Az új algoritmus

- n Bag of words jellegű algoritmus
- n előnyök
 - n nem kell szóegyértelműsítést alkalmazni
 - n nem kell szinonimaegyértelműsítést /
-szűrést alkalmazni
 - n nem érzékeny a szavak
sorrendjére
- n hátrányok
 - n keresési tér nagy
 - n lineáris keresési idő



- n Angol Wikipedia
 - n 31GB XML
 - n 3 800 000 szócikk
 - n SZTAKI Desktop GRID
 - n Letölthető szöveges változat több nyelven:
<http://kopiwiki.dsd.sztaki.hu/>

- n Google Translate
 - n Csak teszteléshez
 - n Találati arány egyezik a kézi fordításával



WIKIPEDIA
The Free Encyclopedia

- n <http://www.wikipedia.org>
- n <http://translate.google.com>
- n <http://kopi.sztaki.hu>

Statisztikák

		találatok száma →				
		1	2	3	4	5
mondatok száma →	1	0,555709				
	2	0,802606	0,308813			
	3	0,9123	0,583218	0,17161		
	4	0,961035	0,766092	0,400344	0,095365	
	5	0,982688	0,874424	0,603594	0,264845	0,052995
	6	0,992309	0,934587	0,754096	0,453091	0,170722
	7	0,996583	0,966664	0,854397	0,620362	0,327637
	8	0,998482	0,98329	0,916785	0,750418	0,490307
	9	0,999325	0,991732	0,953742	0,842869	0,634853
	10	0,9997	0,995952	0,974854	0,904482	0,750449


Találati arány

sztaki kopi

Szótár	KOPI	NDA	Kereső
--------	------	-----	--------

Teszt

A plágiumkeresőt úgy tesztelheti, hogy Petőfi Sándor verseiből szűr be egy-két versszakot ide. A rendszer már párszor tíz szavas egyezést is képes kijelezni. Petőfi verseket az alábbi oldalon találhat: <http://mek.oszk.hu>



Kezdőlap

Tartalom: **Mit tud**, **Kinek szánjuk**, **Hol keres**, **A KOPI használatáról röviden**, **Történet**, **Kapcsolat**

Üdvözljük a KOPI plágiumkereső portálon!

KOPI - A fordítási plágiumok keresője

"plágium: szellemi tolvajlás, más művének közlése saját név alatt, a mű alap gondolatának vagy részleteinek felhasználása a szerzőre való hivatkozás nélkül" (Magyar Értelmező Szótár)

Napjainkban egyre gyakrabban találkozunk szó szerint lemásolt, plagizált tartalommal. Ennek felkutatására számos megoldás született már, ezek közül magyar nyelven a SZTAKI Elosztott Rendszerek Osztálya által üzemeltetett **KOPI Plágiumkereső a legismertebb.**

magyar | **english**

Betűméret - +
 Nagy kontraszt
 Súgó

KOPI

Kezdőlap

Plágiumkeresés
 Feltöltés
 Dokumentumaim
 Plágiumkereső
 Futó keresések
 Üzenetek
 Fórum

Felhasználó:

Beállításaim
 Kilépés

Dokumentumok

Jogszabályok

1.

Cím:

Szerző:

Dokumentum:

2.

<input checked="" type="checkbox"/>	Miért kell akadálymentesíteni?	Pataki Máté	2010.11.14.	<input type="button" value="Szerkeszt"/>
				<input type="button" value="Részletes"/>

3.

- Egynyelvű keresés - dokumentumok összehasonlítása:
 - egymással
 - minden felhasználó dokumentumaival
- Többnyelvű keresés (**tesztüzem**) - dokumentumok összehasonlítása:
 - az angol Wikipédiával
 - a magyar Wikipédiával

From: KOPI
Date: 2012.01.24.
Subject: 1 dokumentum összehasonlítása az angol Wikipédiával.

[\[Üzenet törlése\]](#)

2 hasonló mondatot talált a rendszer 3 Wikipédia cikkben:

1. **Rövidítés** (3)

Rövidítésnek (latinul abbreviatura) nevezünk közszavak és tulajdonnevek rövidített formáit, melyek szinte kizárólag írott formában élnek, azaz amelyeket kiejtve teljes alakjukban használunk.

- rövidítés és mozaikszó egy szó, kifejezés vagy név rövidített formája
- megjegyzés 1: rövidítésnek n
- formáit, melyek szinte kizáróli
- teljes alakjukban használunk.

(utca), km (kilométer), É (észak

- (utca), km (kilométer), É (ész

1. **Pete Seeger** (7)

Seeger was born in French Hospital, Midtown Manhattan, the youngest of three sons.

- Pete Seeger Manhattan közepén, a Midtown-nak is hívott városrész francia kórházában született.

His father, Charles Louis Seeger Jr. was a prominent musicologist, composer, and music professor.

- Apja, ifj. Charles Louis Seeger, zeneszerző és zenetudós volt, aki az elsők között vizsgálta mind az amerikai népzene, mind a nem-európai gyökerekből fakadó zenét.

His stepmother, Ruth Crawford Seeger, was one of the most significant female composers of the twentieth century.

- Nevelőanyja, Ruth Crawford Seeger egyike volt a huszadik század legkiemelkedőbb női zeneszerzőinek.

<input type="checkbox"/>	12 cikk magyar  8% (8 mondat) egyezés	kopiwiki	2011.11.17.	Szerkeszt Részletes
<input type="checkbox"/>	Prince Sisouk GER  28% (2 mondat) egyezés	Wiki	2011.11.30.	Szerkeszt Tömör
<p>Kiadó: -</p> <p>dokumentum: Eredeti: PrinceSisoukGer.docx Szöveges: PrinceSisoukGer.docx.txt Dokumentum hossza: 80 szó Csak én láthatom</p> <p>Nyelv: Német</p> <p>Megjegyzés:</p> <p>Kulcsszavak:</p> <p>Kivonat:</p> <p>Hasonló dokumentumok: 2 mondat egyezés az angol Wikipédiával: Sisouk na Champassak</p>				
<input type="checkbox"/>	Prince Sisouk HUN  80% (4 mondat) egyezés	Wikipedia	2011.11.30.	Szerkeszt Részletes

Demó – Nem talált mondatok

HUN: A Carnatic ez? - robbant ki.

ENG: Am I on the Carnatic?" -1

HUN: A detektívnek minden oka megvolt, hogy így okoskodjon.

ENG: The detective was not far wrong in making this conjecture. -2

HUN: A detektív is hasztalanul fáradozott, hogy ő legyen a nézeteltérésben a főszereplő.

ENG: As vainly did the detective endeavor to make the quarrel his. -3

HUN: - A hídon.

ENG: "On the bridge." 2

HUN: - Addig óvadék ellenében szabadlábra helyezem mindkettőjüket.

ENG: "Meanwhile, you are liberated on bail." -3

HUN: A gentlemannek különben is kész volt a terve a továbbiakra.

ENG: Mr. Fogg's course, however, was fully decided upon. -6

HUN: A gépész azonban történetesen épp e napon felment a fedélzetre, megkereste Mr. Foggot, és meglehetősen élénk vitát folytatott vele.

ENG: On this day the engineer came on deck, went up to Mr. Fogg, and began to speak earnestly with him. -4

<http://kopi.sztaki.hu>

Köszönöm a figyelmet!

Web: <http://dsd.sztaki.hu>

Email: Mate.Pataki@sztaki.hu