

MTA SZTAKI DSD

Department
of Distributed
Systems

Fordítási plágiumok keresése

MSZNY 2011

Pataki Máté

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 1-4 → fordítási plágiumok
 - n ???

- 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 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)

- 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 Lehetséges megoldások
 - n Jelenlegi egynyelvű algoritmus továbbfejlesztése
 - n n-gram
 - n 3-6 szó: 2-4 nagyságrend
 - n + szórend
 - n Teljesen új algoritmus
 - n ???

Az új algoritmus

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



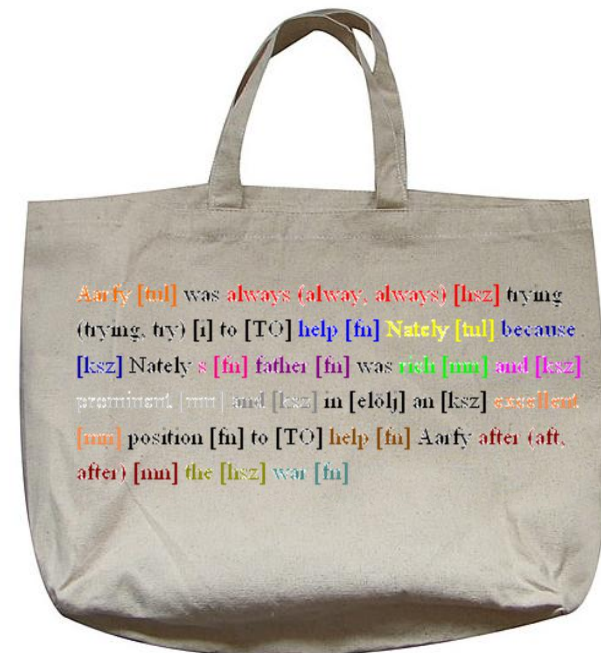
- n Hasonlósági metrika

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

- 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



Az új algoritmus

- n Fontosabb változók
- n Stopszavak
- n Tulajdonnevek / ismeretlen szavak
- n Hasonlóság mértéke / hol vágunk (f+ / f-)
- n Szótár mérete
- n Talált/nem talált szavak értéke (szófajonként)

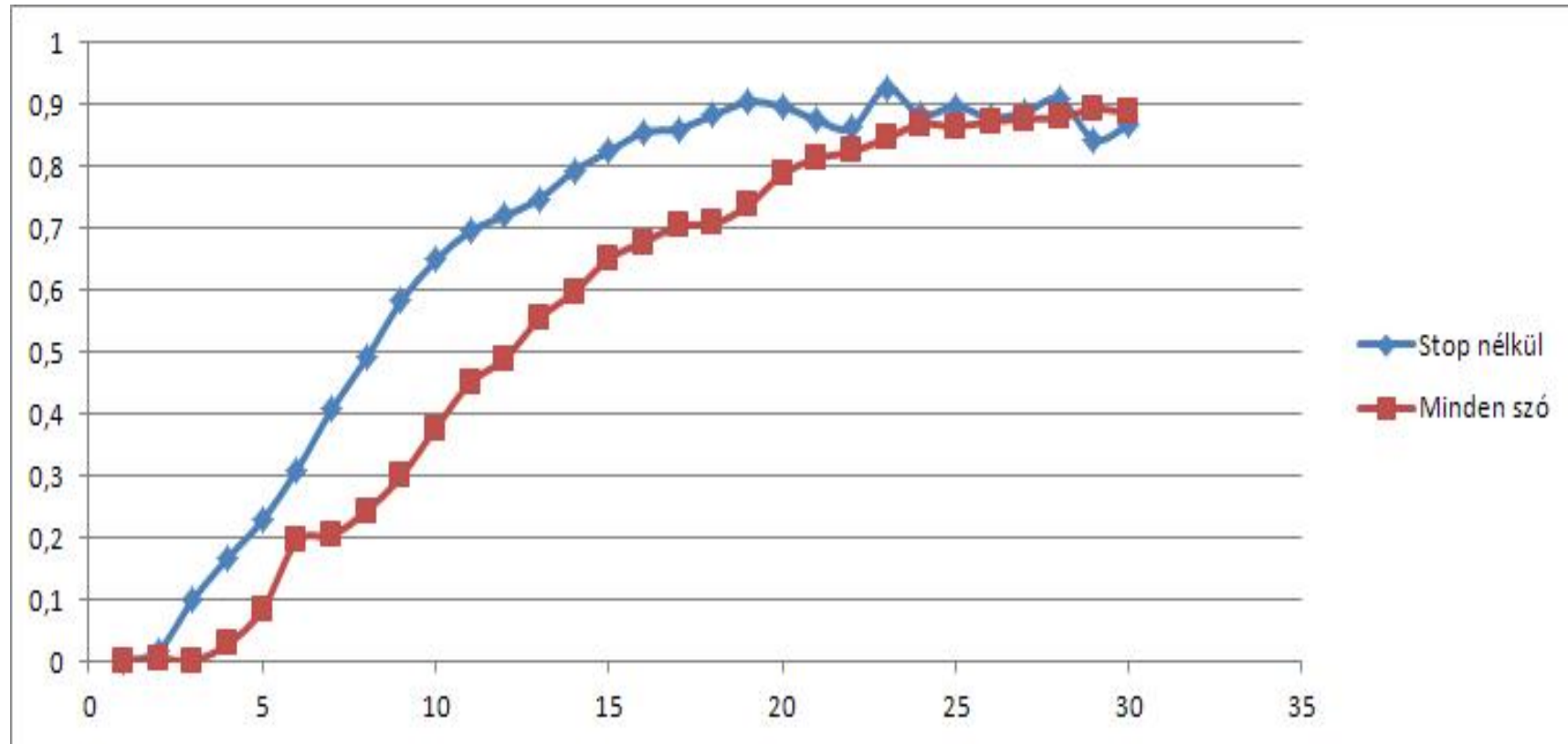
Aarfy [tul] was always (alway, always) [hsz] trying (trying, try) [i] to [TO] help [fn] Nately [tul] because [ksz] Nately s [fn] father [fn] was rich [mn] and [ksz] prominent [mn] and [ksz] in [előlj] an [ksz] excellent [mn] position [fn] to [TO] help [fn] Aarfy after (aft, after) [mn] the [hsz] war [fn]

Aarfy mindig (mindig, mind) igyekezett (igyekezik, igyekezett) Natelyn segíteni (segít) mert (mert, mer) Nately apja (apa) gazdag és befolyásos ember volt (volt, van) aki kitűnő állása (állás) révén (révén, rév) segíthetett (segít) volna (van) Aarfyn a háború után

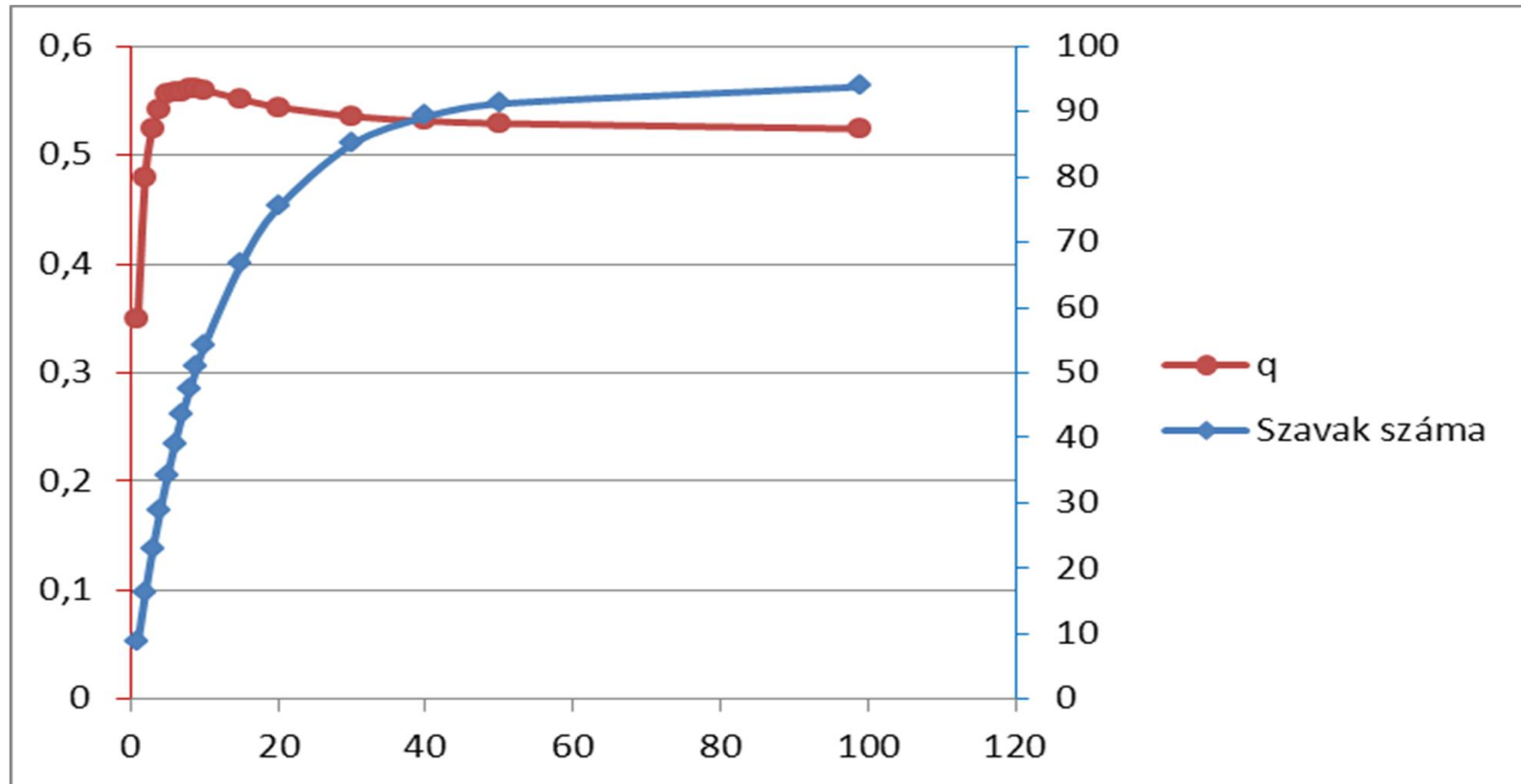
- n Angol Wikipedia
 - n 31GB XML
 - n 3 800 000 szócikk
 - n SZTAKI Desktop GRID
 - n Lesz letölthető szöveges változat több nyelven
- n Hunspell
- n Google Translate
 - n Csak teszteléshez
 - n Találati arány egyezik a kézi fordításéval



WIKIPEDIA
The Free Encyclopedia



Mondat hossza és a megtalált fordítások aránya



Szósák mérete és találati arány viszonya

Demó – Wikipedia random article



Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia

Interaction
Help
About Wikipedia
Community portal
Recent changes
Contact Wikipedia

Toolbox
Print/export

Languages
Afrikaans
العربية
Deutsch
Español
Français
עברית
Nederlands
日本語
Polski
Português
Русский
Suomi

Log in / create account

Article [Discussion](#)

[Read](#) [Edit](#) [View history](#)

Search

Minitel

From Wikipedia, the free encyclopedia

The **Minitel** is a Videotex online service accessible through the telephone lines, and is considered one of the world's most successful pre-World Wide Web online services. It was launched in France in 1982 by the PTT (Poste, Téléphone et Télécommunications; divided since 1991 between France Télécom and La Poste). From its early days, users could make online purchases, make train reservations, check stock prices, search the telephone directory, have a mail box, and chat in a similar way to that now made possible by the Internet.

In February 2009, France Telecom indicates the Minitel network still has 10 million monthly connections, among which 1 million on the 3611 (directory). France Telecom is planning to retire the service on the 30th of June 2012.^[1]

Contents [hide]

- Business model
 - Finances
 - Phonebook
- Technical
- Minitel and the Internet
- Minitel in other countries
- References
- External links

Business model [\[edit\]](#)

Millions of terminals were handed out free to telephone subscribers, resulting in a high penetration rate among businesses and the public. In exchange for the terminal, the possessors of Minitel would not be given free "white page" printed directories (alphabetical list of residents and firms), but only the yellow pages (classified commercial listings, with advertisements); the white pages were accessible for free on Minitel, and they could be searched by a reasonably intelligent search engine; much faster than flipping through a paper directory.

France Télécom estimates that almost 9 million terminals — including web-enabled personal computers (Windows, Mac OS, and Linux) — had access to the network at the end of 1999, and that it was used by 25 million people (of a total population of 60 million).

The Minitel allowed access to various categories of services:

- phone directory (free)
- mail-order retail companies
- airline or train ticket purchases
- information services
- databases
- message boards



Minitel 1. Built 1982



Demó – Google translate

+You Web Images Videos Maps News Gmail More ▾

SZTAKI DSD ▾ ⚙



Translate

From: English ▾



To: Hungarian ▾

Translate

English Hungarian German

The Minitel is a Videotex online service accessible through the telephone lines, and is considered one of the world's most successful pre-World Wide Web online services. It was launched in France in 1982 by the PTT (Poste, Téléphone et Télécommunications; divided since 1991 between France Télécom and La Poste). From its early days, users could make online purchases, make train reservations, check stock prices, search the telephone directory, have a mail box, and chat in a similar way to that

Hungarian German English

A Minitel egy Videotex online szolgáltatás elérhető a telefonvonalak, és ezt tartják a világ egyik legsikeresebb pre-World Wide Web online szolgáltatások. Ben indult Franciaországban 1982-ben a PTT (Poste, telefon et télécommunications; osztva 1991 óta között a France Télécom és a La Poste). A kezdetektől, a felhasználók tudta, hogy az online vásárlások, hogy a vonat foglalás, ellenőrizze tőzsdei árfolyamok, keresés a telefonkönyvben, van egy postaláda, és a chat-ben hasonló módon, hogy most már lehetővé tette az interneten.

2009 februárjában, a France Telecom jelzi Minitel hálózat még van 10 millió havi kapcsolatok, melyek közül 1 millió a 3611 (könyvtár). France Telecom azt tervezi, hogy nyugdíjba vonul a szolgáltatás 30-án 2012. [1]



New! Click the words above to view alternate translations. [Dismiss](#)

Keresés az angol Wikipédiában

Magyar/angol/német szöveg:

A Minitel egy Videotex online szolgáltatás elérhető a telefonvonalak, és ezt tartják a világ egyik legsikeresebb pre-World Wide Web online szolgáltatások. Ben indult Franciaországban 1982-ben a PTT (Poste, telefon et télécommunications; osztva 1991 óta között a France Télécom és a La Poste). A kezdetektől, a felhasználók tudta, hogy az online vásárlások, hogy a vonat foglalás, ellenőrizze tőzsdei árfolyamok, keresés a telefonkönyvben, van egy postaláda, és a chat-ben hasonló

Keres

A bemenet nyelve **magyar**.

6 mondatot találtam.

0. A Minitel egy Videotex online szolgáltatás elérhető a telefonvonalak, és ezt tartják a világ egyik legsikeresebb pre-World Wide Web online szolgáltatások.

The Minitel is a Videotex online service accessible through the telephone lines, and is considered one of the world's most successful pre-World Wide Web online services. (15, [Minitel](#)).

Pre-World Wide Web online services are online service providers that predate the creation of the world wide web in the early 1990s and did not (initially) use TCP/IP. (0, [Category:Pre-World Wide Web online services](#)).

*Where to find Web Services on the Web: Investigating Web Services on the World Wide Web (2008) (0, [Web service](#)).

* Where to find Web Services on the Web: Investigating Web Services on the World Wide Web. (0, [WS-Discovery](#)).

Online Matrimony Services in India This World Wide Web based service has millions of users Datta, Damayanti. (0, [Matrimonial websites](#)).

A list of the holdings is available on the World Wide Web. (0, [Wells Coates](#))

Missing translations: France, online, Poste, Telecom, Minitel, ben, chat, pre, án, Télécom, World, web, PTT, Wide, Ben

Találatok:

1. [Minitel](#) (59)

The Minitel is a Videotex online service accessible through the telephone lines, and is considered one of the world's most successful pre-World Wide Web online services.

- A Minitel egy Videotex online szolgáltatás elérhető a telefonvonalak, és ezt tartják a világ egyik legsikeresebb pre-World Wide Web online szolgáltatások. (15)

It was launched in France in 1982 by the PTT (Poste, Téléphone et Télécommunications; divided since 1991 between France Télécom and La Poste).

- Ben indult Franciaországban 1982-ben a PTT (Poste, telefon et távközlés; osztva 1991 óta között a France Télécom és a La Poste). (13)

From its early days, users could make online purchases, make train reservations, check stock prices, search the telephone directory, and chat in a similar way to that now made possible by the Internet.

- A kezdetekből, a felhasználók tudta, hogy az online vásárlások, hogy a vonat foglalás, ellenőrizze tőzsdei árfolyamok, keresés a telefonkönyvben, van egy postaláda, és a chat-ben hasonló módon, hogy most már lehetővé tette az interneten. (28)


- 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

Találathi arány

sztaki kopi

Szótár	KOPI	NDA	Kereső
<input type="text"/>		<p>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</p>	
<input type="button" value="Teszt"/>			
<h2>Kezdőlap</h2> <hr/> <p>Tartalom: Mit tud, Kinek szánjuk, Hol keres, A KOPI használatáról röviden, Történet, Kapcsolat</p> <p>Üdvözöljük a KOPI plágiumkereső portálon!</p> <p>KOPI - A fordítási plágiumok keresője <i>"plágium: szellemi tolvajlás, más művének közlése saját név alatt, a mű alap gondolatának vagy részleteinek felhasználása a szerzőre való hivatkozás nélkül" (Magyar Értelmező Szótár)</i></p> <p>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.</p>		<p>magyar english</p> <p>Betűméret - +</p> <p>Nagy kontraszt</p> <p>Súgó</p> <p>KOPI</p> <p>Kezdőlap</p> <p>Plágiumkeresés</p> <ul style="list-style-type: none"> Feltöltés Dokumentumaim Plágiumkereső Futó keresések <p>Üzenetek</p> <p>Fórum</p> <p>Felhasználó:</p> <p>Beállításaim</p> <p>Kilépés</p> <p>Dokumentumok</p> <p>Jogszabályok</p>	

Plágiumkeresés és dokumentumkezelés

Feltöltés

Dokumentumaim

Plágiumkereső

Futó keresések

Üzenetek

Kérem válasszon plágiumkeresési formát:

- Egynyelvű keresés - dokumentumok összehasonlítása:
 - minden felhasználó dokumentumaival
- Többnyelvű keresés (**tesztüzem**) - dokumentumok összehasonlítása:
A dokumentum túl hosszú (2980 szó), a tesztüzemben maximum 2500 szavas dokumentumokkal lehet többnyelvű keresést indítani.

Plágiumkeresés indítása >>

Cím	Szerző	Feltöltés dátuma	
12 cikk	kopiwiki	2011.11.17.	Törlés

Amennyiben valamelyik dokumentumot mégse szeretné, hogy résztvegyen a keresésben, eltávolíthatja a "törlés" gombbal a listából.

From: KOPI
Date: October 7, 2011
Subject: 1 dokumentum összehasonlítása az angol Wikipédiával.

[\[Delete message\]](#)

6 hasonló mondatot talált a rendszer 2 Wikipédia cikkben:

1. **Pete Seeger** (6)

Seeger was born in French Hospital, Midtown Manhattan.

- Pete Seeger Manhattan közepén, a Midtown-nak is hívott városrész francia kórházában született. (4)

His parents were living with his grandparents in Patterson, New York, from 1918 to 1920.

- Szülei 1918 és 1920 között a nagyszülőkkel együtt a New York állambeli Pattersonban éltek. (11)

His father, Charles Louis Seeger Jr., was a composer and pioneering ethnomusicologist investigating both American folk and non-Western music.

- Apja, ifj. Charles Louis Seeger, zeneszerző és zenetudós volt, aki az elsők között vizsgálta mind az amerikai népzénet, mind a nem-európai gyökerekből fakadó zenét. (3)

His mother, Constance de Clyver Edson, was a classical violinist and teacher.

- Édesanyja, Constance de Clyver Edson klasszikus hegedűművész és tanár volt. (9)

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. (1)

His half-sister, Peggy Seeger, also a well-known folk performer, was married for many years to British folk singer Ewan MacColl.

- Féltestvére, Peggy Seeger, aki szintén ismert népzenei előadó volt, hosszú évekig Ewan MacColl brit folkénekesrel élt házasságban. (7)


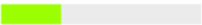

2. **Mike Seeger** (2)

His father, Charles Louis Seeger Jr., was a composer and pioneering ethnomusicologist, investigating both American folk and non-Western music.

- Apja, ifj. Charles Louis Seeger, zeneszerző és zenetudós volt, aki az elsők között vizsgálta mind az amerikai népzénet, mind a nem-európai gyökerekből fakadó zenét. (3)

His sister Peggy Seeger, also a well-known folk performer, was married for many years to British folk singer Ewan MacColl.

- Féltestvére, Peggy Seeger, aki szintén ismert népzenei előadó volt, hosszú évekig Ewan MacColl brit folkénekesrel élt házasságban. (9)

<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