
Sztuczna inteligencja w rozwoju bibliograficznych baz danych

dr hab. inż. Przemysław Korytkowski, prof. ZUT, prof. OPI-PIB

Zachodniopomorski Uniwersytet Technologiczny w Szczecinie
Ośrodek Przetwarzania Informacji – Państwowy Instytut Badawczy

Lublin, 27 września 2024



1

Ekstrakcja cytowań

Projekt "System informatyczny do analizy bibliografii polskich tekstów naukowych"

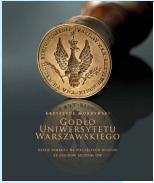
Finansowanie z program Nauka dla społeczeństwa II, MNiSW

Okres realizacji: luty 2025 – luty 2028

Wykonawcy: ZUT i OPI-PIB

2

Problem badawczy



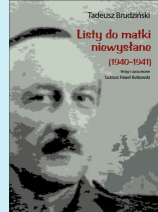
USTANOWIENIE GODŁA

Potrzeba badań

Godło Uniwersytetu Warszawskiego, mimo swojej ważnej, symbolicznej roli w historii uczelni, nie stało się tematem zbyt wielu opracowań naukowych. Informacje o ustanowieniu, umieszczeniu i przygotowaniu tego znaku, wyjaśnienie symboliki jego elementów można nieprzeanalizować w pracach o dziejach uniwersytetu lub ułożeniu w okresie Królestwa Polskiego, jednak jedynie obcasie studiów monograficznych skłaniając ten symbol uczelni w ujęciu historycznym opublikował Stefan K. Kuczyński. Ze względu na to, że praca ta obejmowała dość wąski zakres czasowy – od założenia uczelni do końca XX w. – nie było w niej miejsca na pogłębione badanie samego obłocznostwo ustanowienia godła.

Znak uczelni – orzeł „dławił” polski pod koroną zamkniętą, w oceniana pięciu psian, utrzymując w pierwszym szeregu gałkę ławną, a w lewym pałkowcu – zasał przycięty w kształcie 1871 r. na wniosek Komisji Królewskiej Wymiaru Religijnych i Oświecenia Publicznego? Dlaczego zdecydowano się na zawieszenie znaku obrotującego się do Ota Baluga, historycznego godła dawnej Korony Królestwa Polskiego? Wybór takiego jego postaci nie był wcale oczywisty, gdyż nie nastąpiło żadne – jak uważali w Komisji – reprezentowane były przez przedstawienia różnych symboli heraldyki charakterystycznych dla nowoczesności uniwersyteckich i akademickich: ksiąg symbolizujących Biblię lub prawa i ustawy uczelni¹ albo obłocznostwo z ornamentacją². Utworzył w godle Ota Baluga miało czytelnie konotacje patriotyczne, ale mogło budzić powątpiewanie w odnośniku do wyrażenia postępowości i nowoczesności. Po interwencji wielkiego księcia Konstantego uniwersytet został przywrócić do stanu sprzed 1871 r. W styczniu 1873 r. wprowadzono orszak królewskie godła Królestwa Polskiego. Okoliczności przyjęcia pierwotnego znaku oraz fakty jego zastąpienia innym symbolem składają do podstawy badań historycznych, których celem jest głębsze zrozumienie znaczenia godła uczelnianego w pierwszych latach konstytucyjnego okresu Królestwa Polskiego. Jedną z podstawowych tez wyrażonych przez Kuczyńskiego podlega, że symboliwość utworzona przez uniwersytet znaku Ota Baluga była specjalnym przypadkiem, którym mogły się cieszyć tylko dwie elity: uczelnia i Węskio Polskie³. Analiza historii tego jedynego niu tu, z której korzystano we wcześniejszych badaniach, pozwała zaktualizować tę tezę.

- ¹ S.K. Kuczyński, *Godła Uniwersytetu Warszawskiego*, Warszawa: Wydawnictwo Uniwersytetu Warszawskiego, 2005, s. 132-134.
- ² Wskazywać należy, że nie ma wątpliwości, iż wybrany przez komisję orzeł, mimo swojej symboliki, nie był wcale oczywistym wyborem, gdyż nie nastąpiło żadne – jak uważali w Komisji – reprezentowane były przez przedstawienia różnych symboli heraldyki charakterystycznych dla nowoczesności uniwersyteckich i akademickich: ksiąg symbolizujących Biblię lub prawa i ustawy uczelni¹ albo obłocznostwo z ornamentacją².
- ³ S.K. Kuczyński, *Godła Uniwersytetu Warszawskiego*, Warszawa: Wydawnictwo Uniwersytetu Warszawskiego, 2005, s. 132.



czas nie było niczym nadzwyczajnym dla dzieci z zamożnych rodzin, a do takich można niewątpliwie zaliczyć Brudzińskich⁴. W 1910 r. rodzina przeniosła się do Warszawy, gdzie Józef Brudziński został organizatorem i lekarzem naczelnym powstającego od podstaw Szpitala Karola i Marii dla dzieci, fundacji Zofii Szlenkierówny⁵.

Po przeniesieniu się rodziny do Warszawy Tadeusz rozpoczął w 1910 r. naukę w klasie wstępnej prywatnego 8-klasowego Gimnazjum Emiliana Konopczyńskiego w Warszawie, która przetrwał w 1911 r. z powodów zdrowotnych, a wrócił do tejże szkoły dwa lata później, już do trzeciej klasy⁶. W tym samym roku wstąpił do harcerstwa⁷. O niepodległościowej i patriotycznej atmosferze, w której wzrastał, świadczy zachowany fragment dzienniczka jego siostry Janiny. Pisała w nim o wydarzeniach jesieni roku 1915, gdy Warszawa została już zajęta przez wojska niemieckie, a konspiracyjna dotąd Polska Organizacja Wojskowa zaczęła propagować wstępowanie do Legionów Polskich:

Gdy wyszłam z Tadiem⁸ z domu deszcz zaczął padać, a ja pomimo słotnego dnia jakos uroczyście i radośnie byłam nastroszona. O, bo dziś pierwszy raz w moim życiu miałam zaśpiewać „Boże coś Polskę” w kościele, w katedrze! Umówiliśmy się wszyscy to jest 5⁹, 10¹⁰ drużyna i zawieszki (dajmyż) Zawiszy¹¹ żeby zaśpiewać „Boże coś Polskę” i „Z dymem pożarów¹². Mieli

- ⁴ T.P. Rutkowski, *op. cit.*, s. 44-66.
- ⁵ *Ibidem*, s. 66-75.
- ⁶ Wojskowe Biuro Historyczne im. gen. broni Kazimierza Sosnkowskiego (dalej: WBH), Centralne Archiwum Wojskowe (dalej: CAW), Kolekcja Akt Personalnych (dalej: KAP), Akta personalne (dalej: AP) Tadeusza Brudzińskiego, sygn. 1.481.B.13841, Tadeusz Brudziński, Curriculum vitae, b.d., b.pg.
- ⁷ *Ibidem*.
- ⁸ Tadeuszem Brudzińskim.
- ⁹ Chodzi o Drużynę Skautów im. Zawiszy Czarnego, od 1916 – 16 Warszawską Drużynę Harcerzy im. Zawiszy Czarnego.
- ¹⁰ *Z dymem pożarów* – pieśń skomponowana przez Józefa Nikurczica, do której tekst napisał w 1847 Kornel Ujejski. Popularna od Wiesławy Ludów, w czasie powstania styczniowego pełniła rolę hymnu narodowego.

Lokalizacja przypisu

Table 15. The frequency of oblique cases in the text

Case	Frequency	Percentage
Nominative	17	33.33%
Accusative	17	33.33%
Dative	17	33.33%
Genitive	17	33.33%
Locative	17	33.33%
Instrumental	17	33.33%
Prepositional	17	33.33%
Conjunctive	17	33.33%
Other	17	33.33%

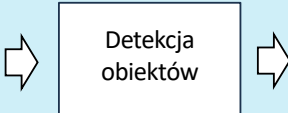
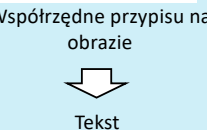


Table 16. The frequency of oblique cases in the text

Case	Frequency	Percentage
Nominative	17	33.33%
Accusative	17	33.33%
Dative	17	33.33%
Genitive	17	33.33%
Locative	17	33.33%
Instrumental	17	33.33%
Prepositional	17	33.33%
Conjunctive	17	33.33%
Other	17	33.33%



Passage na 40 z 3¹

Informacje ogólne i metadane zrzutu²

Passage z przypisem na 40 z 3³

Passage z przypisem na 40 z 3⁴

Passage z przypisem na 40 z 3⁵

Passage z przypisem na 40 z 3⁶

Passage z przypisem na 40 z 3⁷

Passage z przypisem na 40 z 3⁸

Passage z przypisem na 40 z 3⁹

Passage z przypisem na 40 z 3¹⁰

Passage z przypisem na 40 z 3¹¹

Passage z przypisem na 40 z 3¹²

Passage z przypisem na 40 z 3¹³

Passage z przypisem na 40 z 3¹⁴

Passage z przypisem na 40 z 3¹⁵

Passage z przypisem na 40 z 3¹⁶

Passage z przypisem na 40 z 3¹⁷

Passage z przypisem na 40 z 3¹⁸

Passage z przypisem na 40 z 3¹⁹

Passage z przypisem na 40 z 3²⁰

Passage z przypisem na 40 z 3²¹

Passage z przypisem na 40 z 3²²

Passage z przypisem na 40 z 3²³

Passage z przypisem na 40 z 3²⁴

Passage z przypisem na 40 z 3²⁵

Passage z przypisem na 40 z 3²⁶

Passage z przypisem na 40 z 3²⁷

Passage z przypisem na 40 z 3²⁸

Passage z przypisem na 40 z 3²⁹

Passage z przypisem na 40 z 3³⁰

Passage z przypisem na 40 z 3³¹

Passage z przypisem na 40 z 3³²

Passage z przypisem na 40 z 3³³

Passage z przypisem na 40 z 3³⁴

Passage z przypisem na 40 z 3³⁵

Passage z przypisem na 40 z 3³⁶

Passage z przypisem na 40 z 3³⁷

Passage z przypisem na 40 z 3³⁸

Passage z przypisem na 40 z 3³⁹

Passage z przypisem na 40 z 3⁴⁰

Passage z przypisem na 40 z 3⁴¹

Passage z przypisem na 40 z 3⁴²

Passage z przypisem na 40 z 3⁴³

Passage z przypisem na 40 z 3⁴⁴

Passage z przypisem na 40 z 3⁴⁵

Passage z przypisem na 40 z 3⁴⁶

Passage z przypisem na 40 z 3⁴⁷

Passage z przypisem na 40 z 3⁴⁸

Passage z przypisem na 40 z 3⁴⁹

Passage z przypisem na 40 z 3⁵⁰

Passage z przypisem na 40 z 3⁵¹

Passage z przypisem na 40 z 3⁵²

Passage z przypisem na 40 z 3⁵³

Passage z przypisem na 40 z 3⁵⁴

Passage z przypisem na 40 z 3⁵⁵

Passage z przypisem na 40 z 3⁵⁶

Passage z przypisem na 40 z 3⁵⁷

Passage z przypisem na 40 z 3⁵⁸

Passage z przypisem na 40 z 3⁵⁹

Passage z przypisem na 40 z 3⁶⁰

Passage z przypisem na 40 z 3⁶¹

Passage z przypisem na 40 z 3⁶²

Passage z przypisem na 40 z 3⁶³

Passage z przypisem na 40 z 3⁶⁴

Passage z przypisem na 40 z 3⁶⁵

Passage z przypisem na 40 z 3⁶⁶

Passage z przypisem na 40 z 3⁶⁷

Passage z przypisem na 40 z 3⁶⁸

Passage z przypisem na 40 z 3⁶⁹

Passage z przypisem na 40 z 3⁷⁰

Passage z przypisem na 40 z 3⁷¹

Passage z przypisem na 40 z 3⁷²

Passage z przypisem na 40 z 3⁷³

Passage z przypisem na 40 z 3⁷⁴

Passage z przypisem na 40 z 3⁷⁵

Passage z przypisem na 40 z 3⁷⁶

Passage z przypisem na 40 z 3⁷⁷

Passage z przypisem na 40 z 3⁷⁸

Passage z przypisem na 40 z 3⁷⁹

Passage z przypisem na 40 z 3⁸⁰

Passage z przypisem na 40 z 3⁸¹

Passage z przypisem na 40 z 3⁸²

Passage z przypisem na 40 z 3⁸³

Passage z przypisem na 40 z 3⁸⁴

Passage z przypisem na 40 z 3⁸⁵

Passage z przypisem na 40 z 3⁸⁶

Passage z przypisem na 40 z 3⁸⁷

Passage z przypisem na 40 z 3⁸⁸

Passage z przypisem na 40 z 3⁸⁹

Passage z przypisem na 40 z 3⁹⁰

Passage z przypisem na 40 z 3⁹¹

Passage z przypisem na 40 z 3⁹²

Passage z przypisem na 40 z 3⁹³

Passage z przypisem na 40 z 3⁹⁴

Passage z przypisem na 40 z 3⁹⁵

Passage z przypisem na 40 z 3⁹⁶

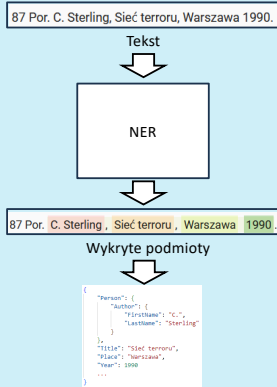
Passage z przypisem na 40 z 3⁹⁷

Passage z przypisem na 40 z 3⁹⁸

Passage z przypisem na 40 z 3⁹⁹

Passage z przypisem na 40 z 3¹⁰⁰

Przetworzenie przypisu

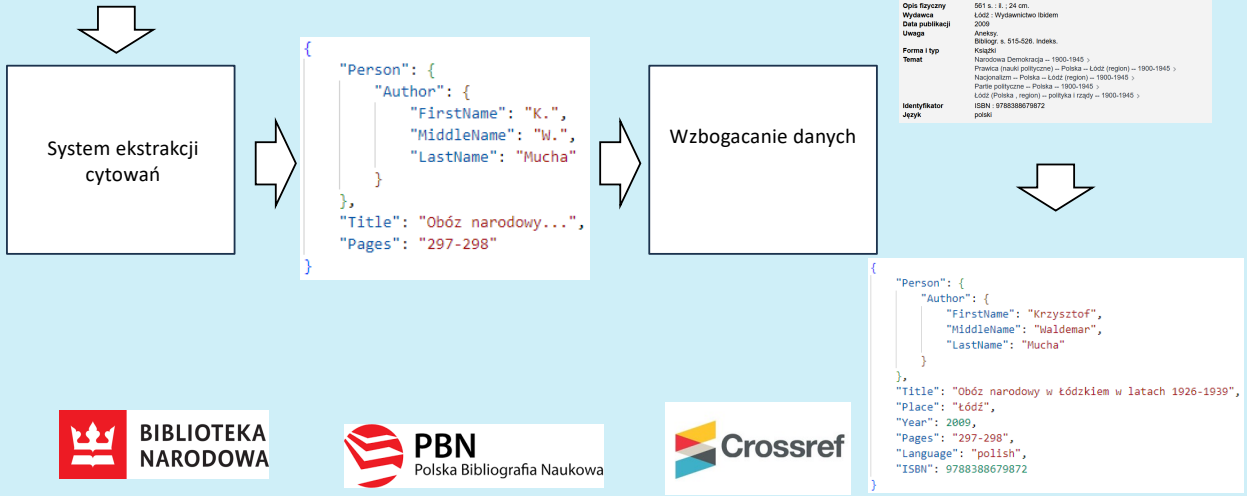


98 M. Grzechnik, 'Ad Maiorem Poloniae Gloriam' Polish Inter-colonial Encounters in Africa in the Interwar Period, „The Journal of Imperial and Commonwealth History” 2020, nr 48 (5), s. 826–845; B. Balogun, Polish Lebensraum: the colonial ambition to expand on racial terms, „Ethnic and Racial Studies” 2018, nr 41 (14), s. 2561–2579.

5

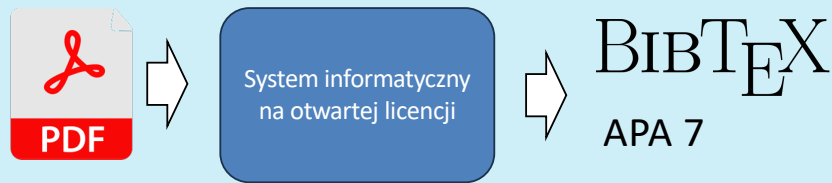
Identyfikacja i wzbogacanie

64 K.W. Mucha, Obóz narodowy..., s. 297–298.



6

Wynik projektu

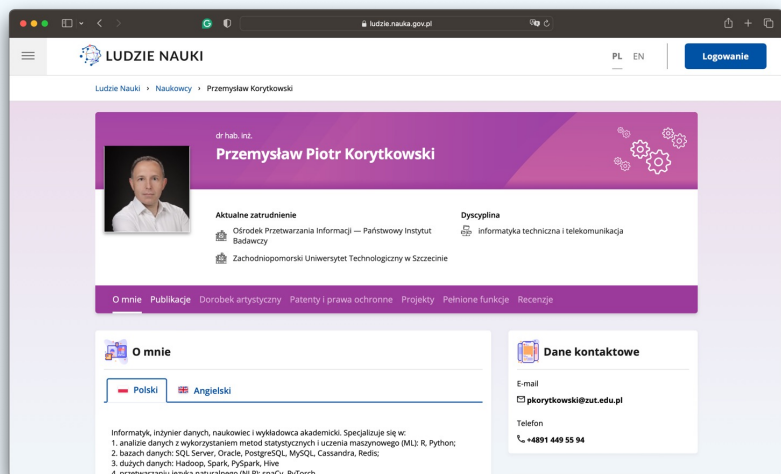


7

Ludzie Nauki

ludzie.nauka.gov.pl

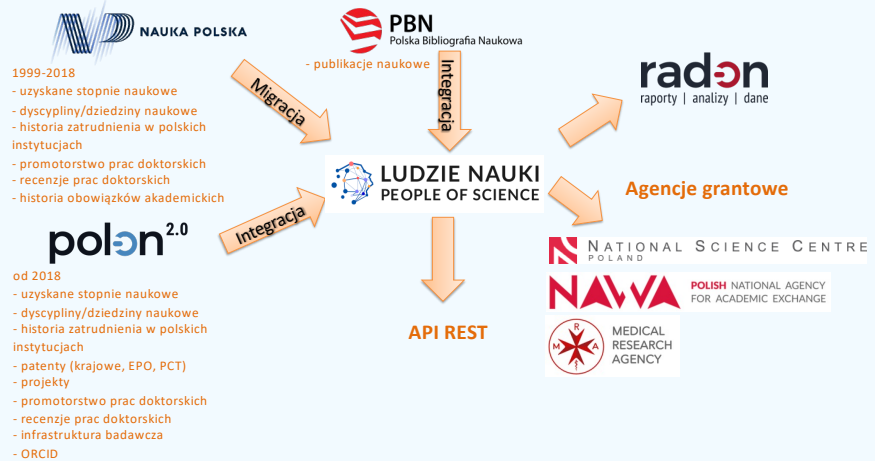
200 000+ profili naukowców
 2 000 000+ metadanych publikacji
 130 000 osiągnięć artystycznych
 12 000+ patentów
 36 000+ projektów



Państwowe dane = wiarygodne dane

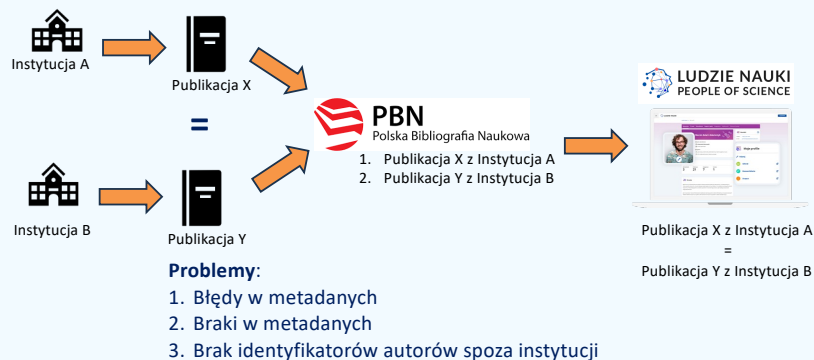
8

Docelowy ecosystem Ludzi Nauki



9

Jakość danych bibliograficznych



10

Złoty Rekord Publikacji

Złoty rekord to termin używany w dziedzinie zarządzania danymi i odnosi się do pojedynczej i wiarygodnej reprezentacji danych, która pochodzi z różnych źródeł i jest przechowywana w centralnej bazie danych.

Złoty Rekord Publikacji to pojedyncza, najbogatsza i wiarygodna reprezentacja metadanych o publikacji naukowej, której przynajmniej jednym z autorów jest pracownik polskiej instytucji naukowej.

11

Złoty Rekord Publikacji

1. **Gromadzenie danych** z różnych źródeł w celu stworzenia kompleksowego obrazu podmiotu danych.
2. **Czyszczenie i standaryzacja** danych, poprawianie błędów i standaryzacja danych w celu zapewnienia ich spójności i dokładności.
3. **Usuwanie duplikatów**, identyfikacja i łączenie duplikatów i sprzecznych informacji w celu stworzenia pojedynczej, wiarygodnej reprezentacji metadanych.
4. **Wzbogacanie danych**, dodawanie brakujących lub uzupełniających informacji do złotego rekordu, aby uczynić go bardziej kompletnym i użytecznym.
5. **Walidacja danych**, sprawdzanie dokładności złotego rekordu i upewnianie się, że jest on zgodny z regułami biznesowymi i standardami jakości danych.

12

Deduplikacja metadanych publikacji



13

Złoty Rekord Publikacji

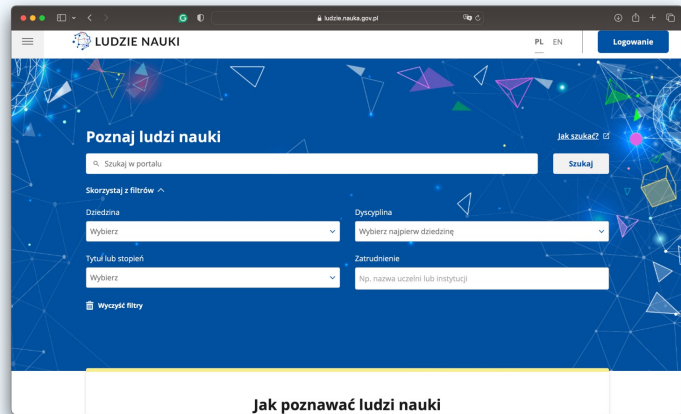
1. Wstępne wyszukiwanie par tytułów z wykorzystaniem algorytmów TF oraz Hierarchical Navigable Small World
2. Zgodność musi być na poziomie 0,9 odległości cosinusowej
3. Wartości następujących atrybutów muszą się zgadzać (te same wartości lub NULL): DOI, ISBN, ISSN, ISMN, rok, MNISW_ID
4. Jeśli w tytule cyfry się nie zgadzają, to para publikacji nie jest łączona (np. tom I i tom II)
5. Tytuł musi być dłuższy niż 16 znaków
6. Nazwa czasopisma i nazwa wydawnictwa musi być zgodne na poziomie 0,8 odległości cosinusowej
7. Tytuły książki w przypadku rozdziałów muszą być zgodne na poziomie 0,9 odległości cosinusowej

14

Plany – wyszukiwarka semantyczna

Wyszukiwanie semantyczne to podejście, które stara się uzyskać jak najdokładniejsze wyniki poprzez zrozumienie intencji użytkownika oraz powiązań między słowami i kontekstem samego zapytania.

Planujemy wykorzystać duże modele językowe (LLM).



15

Konkluzje

Połączone Otwarte Dane

Dla publikacji naukowych:

1. DOI – Digital Object Identifier
2. ORCID – Open Researcher and Contributor ID
3. ROR – Research Organization Registry



By Michael Hausenblas, James G. Kim, five-star Linked Open Data rating system developed by Tim Berners-Lee. - <http://5stardata.info/en/>, CCO, <https://commons.wikimedia.org/w/index.php?curid=64408651>

16

Dziękuję za uwagę!

dr hab. inż. Przemysław Korytkowski, prof. ZUT, prof. OPI-PIB

Zachodniopomorski Uniwersytet Technologiczny w Szczecinie
Ośrodek Przetwarzania Informacji – Państwowy Instytut Badawczy

Lublin, 27 września 2024