

Marta Grabowska*

Zbiory *Big Data* oraz technologia chmury obliczeniowej dla humanistyki cyfrowej

Wprowadzenie

Po prawie dwóch wiekach epoki przemysłowej, w której maszyny zastępowały głównie siłę mięśni ludzkich, osiągnięcia wybitnych matematyków i inżynierów końca XIX i początku XX w., takich jak: David Hilbert (1862–1943), Kurt Goedel (1906–1978), John von Neumann (1903–1957), Charles Babbage (1871–1871) czy Lee De Forest (1873–1961)¹, skierowały wynalazczość w dziedzinie techniki na inne tory, tj. objęły sferę informacji. Pojawienie się komputerów stacjonarnych, a potem mobilnych, coraz większe pojemności pamięci dyskowych, miniaturyzacja, a także zaawansowane oprogramowania komputerowe i rozwój sieci teleinformatycznych przywiódł nas do punktu, w którym posługujemy się ogromnymi ilościami danych, zwanymi *Big Data*, oraz różnorodnymi narzędziami operowania nimi, otwierającymi nowe przestrzenie różnorodnych działań.

Terminem *Big Data* określa się ogromną ilość danych liczonych w terabajtach (1000 GB), petabajtach (1000 TB), exabajtach (1000 PB) i zettabajtach (1000 EX), niemieszczących się już w komputerach o pojemnościach kilkunastu gigabajtów, gromadzonych i przetwarzanych w środowisku nowej technologii, którą jest chmura obliczeniowa (ang. *cloud computing*)². Technologia ta, za pomocą zaawansowanych narzędzi operowania danymi, umożliwia lepszą organizację i zarządzanie wielkimi zbiorami danych i dzięki temu uzyskiwanie nowej wiedzy. Obejmuje ona nowe rozwiązania zarówno w zakresie hardware’u, jak i software’u, w tym też nowoczesnych narzędzi analizy i prezentacji danych.

* Dr hab. **Marta Grabowska**, prof. UW – Centrum Europejskie Uniwersytetu Warszawskiego, e-mail: mgrabowska@uw.edu.pl.

¹ A. Hodges, *Alan Turing: Enigma*, Wydawnictwo Albatros Andrzej Kuryłowicz S.C., Warszawa 2014.

² V. Mayer-Sconberger, K. Cukier, *Big data. Rewolucja, która zmieni nasze myślenie, pracę i życie*, MT Biznes sp. z o.o., Warszawa 2014.

Początki tworzenia wielkich zbiorów danych

Pojęcie *Big Data* wywodzi się z obszaru biznesu, gdzie gromadzone były i są nadal codziennie ogromne ilości danych w takich sektorach, jak np.: bankowość, handel, transport itd. Dane te mają zwykle bardzo różnorodny charakter. Mogą to być tzw. dane ustrukturalizowane (inaczej: uporządkowane czy strukturalne), np. dane o klientach: ich nazwiska, konta bankowe, informacje o transakcjach handlowych, nazwy czy kody kupowanych towarów, ale także dane nieustrukturyzowane, np. pochodzące z korespondencji e-mailowej z klientami czy z wypowiedzi w języku naturalnym umieszczanych np. na portalach społecznościowych, blogach lub w innych miejscach w sieci w postaci tekstu, obrazu czy dźwięku. Wszystkie te dane, napływające jednorazowo czy w sposób ciągły, stanowią wartościowy materiał analityczny. Dodatkowo, po okresach różnego rodzaju ograniczeń dostępu do danych – promowana obecnie polityka tzw. danych otwartych (ang. *open data*)³ zwiększa jeszcze zasoby, które mogą być przetwarzane w systemach informatycznych.

Początkowo w tradycyjnych bazach biznesowych, tj. w pojedynczych bazach relacyjnych, gromadzono dane ustrukturalizowane, posługując się językiem zapytań SQL (*Structured Query Language*)⁴, stworzonym przez firmę IBM w 1960 r. Jednak w przypadku dużych, międzynarodowych korporacji rozwiązania takie nie były wystarczające. W końcu lat 80. XX w., w związku z potrzebą wsparcia procesu zarządzania dużymi przedsiębiorstwami, powstały tzw. hurtownie danych (ang. *data warehouse*)⁵, w których gromadzono dane napływające z wielu ośrodków. Zastosowanie w nich lepszych narzędzi wyszukiwawczych wzbogaconych o elementy analityczne, jak np. Online Analytical Processing (OLAP)⁶, umożliwiło szybsze przeszukiwanie zagregowanych danych oraz ich analizę według określonych kryteriów. Operacje wykonywane na tych danych, zwane

³ P. Morin, *Open data structures. An introduction* [e-book] in OPEL (Open Paths to Enriched Learning). Edmonton: AU Press. 2013, <http://web.b.ebscohost.com/ehost/ebookviewer/ebook/ZTAwMHh3d19fNjM4OTU2X19BTg2?sid=d8b94be4-7c8e-4f69-9727-d3428c90cdfc@sessionmgr115&vid=3&format=EB&rid=1> (dostęp 7.05.2015).

⁴ M.M. David, L. Fesperman, *Advanced standard SQL dynamic structured data modeling and hierarchical processing* [e-book]. Boston, Artech House, 2013, <http://web.a.ebscohost.com/ehost/ebookviewer/ebook/ZTAwMHh3d19fNzUzNTk2X19BTg2?sid=96ce560b-b590-47f4-9f8c-18783f9988e6@sessionmgr4002&vid=3&format=EB&rid=1> (dostęp 6.05.2015).

⁵ A. Chodkowska-Gyuriks, *Hurtownie danych: teoria i praktyka*, Wydawnictwo Naukowe PWN S.A., Warszawa 2014.

⁶ A. Berson, S.J. Smith, *Data warehousing, data mining, and OLAP*, Mc Graw-Hill, New York 2001.

eksploatacją danych (ang. *data mining*)⁷, dawały możliwość wychwytywania pewnych trendów, które miały znaczenie w procesie podejmowania decyzji. Operacje mogły być wykonywane na danych historycznych, pokazujących, co się wydarzyło, na danych wpływających na bieżąco, pokazujących zmiany trendów w czasie rzeczywistym, oraz mogły być też dokonywane porównania czy generowane prognozy. Hurtownie danych miały charakter scentralizowany (dane gromadzone w jednym zbiorze) lub składały się z wielu zbiorów, tj. z tzw. minihurtowni, z których każda reprezentowała pewną kategorię danych, tworząc łącznie tzw. ang. *data mart*, co po polsku można określić jako „targowisko danych”. Hurtownie danych i stosowane tam narzędzia analityczne dały początek tzw. analityce biznesowej (ang. *business intelligence*)⁸. Ważną rolę odgrywały w nich wbudowane elementy automatycznej reakcji na różne zjawiska, tj. tzw. *feedback* (sprzężenie zwrotne), np. automatyczne generowanie zamówień adresowanych do producentów na brakujące towary za pomocą systemu EDI (*Electronic Data Interchange*)⁹, oraz rozbudowane metody wizualizacji danych. Do takich narzędzi analitycznych można też zaliczyć np. systemy ekspertowe¹⁰ z wbudowanymi procedurami wnioskowania czy różne rodzaje tzw. sieci neuronowych¹¹ (m.in. sieci rekurencyjne¹² stosowane do rozwiązywania problemów optymalizacyjnych czy skojarzeniowych, także samoorganizujące się mapy używane do wybierania i analizowania danych, ich klasyfikacji i wizualizacji), sieci radialne¹³ stosowane do rozwiązywania problemów klasyfikacyjnych i prognostycznych. Rodowód hurtowni danych wywodzi się również z firmy IBM. Ich głównym celem było lepsze zarządzanie przedsiębiorstwami i usługami oraz lepsza, bardziej trafna i odpowiadająca indywidualnym zapotrzebowaniom klientów alokacja towarów i usług, co w efekcie dawało obopólne korzyści: zarówno producentom, jak i klientom.

⁷ Ibidem.

⁸ *Business intelligence*, red. J. Suma, Wydawnictwo Naukowe PWN, Warszawa 2013.

⁹ M.C. Pparfatt, *What is EDI?: a guide to electronic data interchange*, PNCC Blackwell, Oxford 1992.

¹⁰ P. Jackson, *Introduction to experts systems*, Accisson Wesley Pub.co, Workingham (England) 1986.

¹¹ S. Osowski, *Sieci neuronowe do przetwarzania informacji*, wyd. 2 popr. i uzup., Oficyna Wydawnicza Politechniki Warszawskiej, Warszawa 2006.

¹² A. Chinae, *Understanding the principles of recursive neural networks. A generative approach to tackle model cpmplexity*, <http://arxiv.org/ftp/arxiv/papers/0911/0911.3298.pdf> (dostęp 2.05.2015).

¹³ L. Rutkowski, *Metody i techniki sztucznej inteligencji*, wyd. 2 zmien., Wydawnictwo Naukowe PWN, Warszawa 2012.

Technologia chmury obliczeniowej

Przełomem w organizowaniu i zarządzaniu tego typu danymi, których ilość z biegiem czasu zaczęła wzrastać lawinowo, stała się wspomniana wyżej nowa technologia chmury obliczeniowej (ang. *cloud computing*)¹⁴, której istota polega przede wszystkim na odciążeniu instytucji, firm i organizacji od całego skomplikowanego problemu gromadzenia i zarządzania wielkimi zbiorami danych i oddania go w ręce nowej kategorii wysoce wyspecjalizowanych usług zewnętrznych (ang. *outsourcingowych*) świadczonych przez firmy dysponujące zaawansowanymi technologiami zarówno w zakresie infrastruktury fizycznej, jak i oprogramowania narzędziowego oraz aplikacyjnego. Na usługi te składają się: serwery dysponujące wielkimi pojemnościami pamięci, których odległa lokalizacja, wobec możliwości współczesnych sieci telekomunikacyjnych, praktycznie nie odgrywa roli, rozbudowane oprogramowanie operacyjne i narzędziowe, którego zadaniem jest wszechstronne organizowanie napływających danych i m.in. zapewnienie ich bezpieczeństwa, wreszcie oprogramowania aplikacyjne umożliwiające różnego rodzaju operacje na tych danych, w szczególności dokonywanie ich wszechstronnej analizy oraz prezentacji¹⁵.

W literaturze przedmiotu można znaleźć stanowisko, że technologia chmury obliczeniowej z jednej strony wywodzi się z takich wcześniejszych osiągnięć techniki, jak videotex¹⁶, lecz z drugiej strony, co jest zaskakujące i co podaje V. Mosco, ma swój rodowód w koncepcji centralnie sterowanej gospodarki Związku Radzieckiego. Na początku lat 60. XX w. bowiem Nikita Chruszczow, chcąc zaradzić ciągłym brakom różnych towarów na rynku lub złej ich alokacji, korzystając z ówczesnych osiągnięć naukowych Norberta Wienera, twórcy cybernetyki, tj. nauki o sprawnym zarządzaniu¹⁷, i sformułowanej przez niego zasadzie tzw. sprzężenia zwrotnego, planował stworzenie centrum zarządzania gospodarką wspomaganego przez maszyny liczące (w tamtej epoce były to jeszcze duże komputery lampowe, tzw. *mainframe*). Decyzje miały być podejmowane centralnie

¹⁴ P. Mell, T. Grance, *The NIST definition of cloud computing. Recommendation of the National Institute of Standards and Technology*, Gaithersburg (MD), Computer Security Division, Information Technology Laboratory, National Institute of Standards and Technology, 2011, s. 2, <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf> (dostęp 2.05.2015).

¹⁵ V. Mosco, *To the cloud. Big data in a turbulent world*, Paradigm Publishers, Boulder (Colorado), London 2014.

¹⁶ A.F. Alber, *Videotex/teletext: principles and practices*, McGraw Hill Book Company, New York 1986.

¹⁷ N. Wiener, *Cybernetics: or control and communication in the animal and the machine*, M.I.T. Press, New York 1961.

na podstawie z jednej strony napływających drogą telexową informacji z sektora produkcji, a z drugiej – sygnalizowanych braków na rynku. Tak sterowana gospodarka miała być bardziej efektywna. Jak podaje V. Mosco, podobny model – pod nazwą Cybersyn – zaprojektowany został na początku lat 70. w Chile, gdzie upadającą gospodarkę socjalistyczną w okresie prezydentury Salvadora Allende’a próbowano ratować w ten sposób przed katastrofą. Oba kraje, ze względu na perturbacje polityczne, nie zrealizowały jednak do końca swoich projektów¹⁸.

Obecna technologia chmury obliczeniowej rozwinięta w Stanach Zjednoczonych Ameryki i w innych krajach charakteryzujących się wysokim poziomem rozwojem gospodarczego, jak np. Japonia, opiera się na wielu nowych osiągnięciach techniki oraz służy szerszym celom. Przede wszystkim takie rozwiązania, jak: powszechnie stosowane karty magnetyczne i technologia identyfikacji radiowej (*Radio Frequency Identification* – RFID)¹⁹, urządzenia mobilne i technologia zbliżeniowa (*Near Field Communication* – NFC)²⁰, telekomunikacja bezprzewodowa, miniaturyzacja urządzeń elektronicznych, a także digitalizacja ogromnych zbiorów różnego rodzaju dokumentów, wreszcie sam Internet, otworzyły nowe możliwości gromadzenia i przetwarzania niewyobrażalnych ilości danych już nie tylko w biznesie, ale także w administracji, sektorze zdrowia, transporcie, nauce itd. Za pomocą nowych narzędzi analitycznych stosowanych w tej technologii z danych wyłania się nowa wiedza.

Usługi świadczone przez firmy w ramach chmury obliczeniowej dzieli się zwykle na cztery rodzaje: Kolokacja, IaaS – Infrastructure as a Service, PaaS – Platform as a Service i SaaS – Software as a Service. W pierwszym przypadku użytkownikom udostępnia się jedynie pomieszczenia na serwery, w drugim – dodaje się infrastrukturę informatyczną (np. Elastic Compute Cloud EC2 firmy Amazon)²¹, w trzecim – dodaje się system operacyjny z elementami oprogramowania narzędziowego (np. Windows Azure)²², a w czwartym – udostępnia się również oprogramowanie aplikacyjne, np.

¹⁸ V. Mosco, op.cit., s. 22.

¹⁹ H. Lephamer, *RFID design principles*, Artech House, Inc., Boston, London 2008.

²⁰ NFC Forum: *NFC Data Exchange Format (NDEF). Technical specification. 2006-07-24*, NFC Forum, Inc., <http://www.eet-china.com/ARTICLES/2006AUG/PDF/NFCForum-TS-NDEF.pdf?SOURCE=DOWNLOAD> (dostęp 12.03.2014).

²¹ Amazon Web Services, *Amazon Elastic Computer Cloud. User Guide for Linux. API version 2014-10-01*, http://www.google.pl/url?sa=t&rct=j&q=&esrc=s&source=web&cd=12&ved=0CG8QFjAL&url=http%3A%2F%2Fawsdocs.s3.amazonaws.com%2FEC2%2Flatest%2Fec2-ug.pdf&ei=jZdLVEi2LuTgyQP_4CoBA&usq=AFQjCNFIg7f8P673xBxButxQueralRAHkg&bvm=bv.92885102,d.bGQ (dostęp 2.05.2015).

²² Microsoft Azure, <http://azure.microsoft.com/pl-pl/> (dostęp 3.05.2015).

systemy poczty elektronicznej (np. Gmail czy Hotmail), a nawet oprogramowanie aplikacyjne może być specjalnie tworzone dla indywidualnego klienta. Ponadto wymienia się również usługi Software+Service, tj. użytkowanie konkretnego oprogramowania firmy, która udostępnia usługi w chmurze (np. Microsoft i jego webowy Microsoft Business Productivity Online Service), CaaS – Communications as a Service – usługodawca zapewnia także platformę telekomunikacyjną oraz IPaaS – świadczenie integracji usług chmurowych²³. Usługi te mogą mieć charakter publiczny, prywatny, wspólny, tj. społecznościowy albo hybrydowy. Systemy muszą wykazywać się tzw. skalowalnością, tj. możliwością rozbudowy pod każdym względem, szczególnie tak, by sprostać gromadzeniu, analizowaniu i prezentacji napływających danych (*volume of data*), ich różnorodności (*variety of data*) oraz tempu ich napływu (*velocity of data*), a ponadto elastycznością i indywidualizacją usług oraz możliwością przetwarzania danych w zbiorach scentralizowanych, jak i w środowisku rozproszonym. Zagwarantowane muszą być: bezpieczeństwo i ochrona danych oraz ustalony sposób i zakres ich udostępniania²⁴.

W technologii chmury obliczeniowej upatruje się obecnie źródło przyspieszenia rozwoju gospodarczego i wzrostu nowych miejsc pracy oraz przyrostu PKB w granicach 2%. Unia Europejska wspiera rozwój usług przetwarzania w chmurze obliczeniowej, co wyrażone zostało zarówno w Europejskiej agendzie cyfrowej²⁵, jak i w Komunikacie Komisji do Parlamentu Europejskiego, Rady, Europejskiego Komitetu Ekonomiczno-Społecznego i Komitetu Regionów w sprawie wykorzystania potencjału chmury obliczeniowej w Europie²⁶. Usługodawcy mogą korzystać ze wzorów umów przygotowanych przez Unię Europejską na usługi świadczone w tym zakresie na jej terytorium, jak również dokonuje się porządkowania unijnego prawa w tym obszarze.

²³ P. Szmit, *Cloud computing historia, technologia, perspektywy*, Polska Agencja Rozwoju Przedsiębiorczości, Warszawa 2012, https://www.web.gov.pl/g2/big/2012_06/ebfa211-fla9f174c7517738f68df2d8b.pdf (dostęp 3.05.2015).

²⁴ J. Hurwitz, A. Nugent, F. Halper, M. Kaufman, *Big data for dummies*, John Wiley & Sons, Hoboken, NY 2013.

²⁵ Komisja Europejska. Komunikat Komisji do Parlamentu Europejskiego, Rady, Europejskiego Komitetu Ekonomiczno-Społecznego i Komitetu Regionów *Europejska agenda cyfrowa*. KOM(2010) 245 Bruksela 19.05.2010, wersja ostateczna.

²⁶ Komisja Europejska. Komunikat Komisji do Parlamentu Europejskiego, Rady, Europejskiego Komitetu Ekonomiczno-Społecznego i Komitetu Regionów w sprawie wykorzystania potencjału chmury obliczeniowej w Europie, Bruksela (2012). KOM(2012) 529, wersja ostateczna.

Zbiory *Big Data* w środowisku chmury obliczeniowej

Pomijając kwestie poziomu infrastruktury fizycznej, na działania w zakresie *Big Data* w środowisku chmury obliczeniowej składają się w zasadzie trzy etapy: 1) pozyskiwanie, gromadzenie i organizowanie danych, 2) analiza danych, 3) raportowanie i wizualizacja rezultatów²⁷.

Zasadniczą nowością w technologii chmury obliczeniowej w odniesieniu do pierwszego etapu, tj. gromadzenia i porządkowania danych, stało się wprowadzenie nowych, potężnych silników pozyskiwania, gromadzenia i porządkowania danych na przepastnych, wirtualnych obszarach Internetu i różnego rodzaju baz danych. Jednym z najbardziej znanych i publicznie dostępnych (otwartych) tego typu silników jest Hadoop firmy Apache²⁸, napisany w języku Java, zaprojektowany na podstawie zamkniętego oprogramowania o nazwie BigTable firmy Google²⁹. Najistotniejszą cechą Hadoop jest przejęta z BigTable procedura MapReduce³⁰, umożliwiająca z jednej strony skuteczne rozproszanie i zidentyfikowanie problemu (pytania) w rozproszonych systemach informacyjnych (funkcja *map*), z drugiej zaś sprowadzenie z powrotem już wstępnie i odpowiednio uporządkowanych i zagregowanych danych (funkcja *reduce*) jako rezultatu. Procedura ta, paralelnie wykonywana na wielkiej liczbie rozproszonych serwerów zawierających ogromne ilości danych oraz równoległe ich przetwarzanie przez ich rozbięcie na mniejsze grupy w celu przyspieszenia całego procesu (Hadoop Distributed File System – HDFS)³¹, daje możliwość błyskawicznego przeszukiwania, gromadzenia i porządkowania niewyobrażalnej ilości danych. Oba systemy, tj. Hadoop i BigTable, zapewniają pełną skalowalność. Hadoop został początkowo zaprojektowany dla biznesu prowadzonego przez Yahoo, lecz stał się wkrótce najpotężniejszym narzędziem pozyskiwania danych w technologii chmury obliczeniowej³².

²⁷ K. Bakshi, *Technologies of big data*, w: *Big data management. Technologies, and applications*, red. Wen-Chen Hu, N. Kaabouch, IGI Global, Hershey (USA) 2014, s. 1–22 (*Information Science References*).

²⁸ I.K. Saavas, G.N Sofianidon, M-Tachar Kechadi, *Applying the K-Means Algorithm in Big Raw Data Sets with Hadoop and MapReduce*, w: *Big data management. Technologies, and applications*, red. Wen-Chen Hu, N. Kaabouch, IGI Global, Hershey (USA) 2014, s. 23–46.

²⁹ E. Ciurana, *Developing with Google App Engine*, Apress, Distributed by New York Springer Verlag, Berkeley (CA) 2008.

³⁰ D. Miner, A. Schock, *MapReduce design patterns*, O'Reilly, Beijing, Koln 2013.

³¹ Hadoop, *HDFS Architecture guide*, http://hadoop.apache.org/docs/r1.2.1/hdfs_design.html (dostęp 3.05.2015).

³² E. Baldeschfieler, *Hadoop at Yahoo!*, <http://www.google.pl/url?sa=t&rct=j&q=&esrc=s&source=web&cd=6&ved=0CFMQFjAF&url=http%3A%2F%2Fwww.hadooper.cn%2Fdcf%2Fattach%2FY2xiOmNsYjpwZGY6ODE%3D&ei=xX5LVdLpGYupygOGiY>

Wspomniany BigTable, zaprojektowany przez firmę Google na podstawie GFS (Google File Systems)³³ – systemu do indeksowania i przetwarzania ogromnych ilości danych tej firmy, umożliwia organizowanie, indeksowanie i przechowywanie wielkich ilości danych „na stołach” składających się z „rzędów” i „kolumn”, lecz w porównaniu z możliwościami tradycyjnych baz SQL jest bardziej pojemny, gwarantuje danym większą odporność na zagrożenia, może przechowywać je w środowisku rozproszonym, tj. na wielu różnych serwerach oraz w środowisku wielowymiarowym. Jest jednak oprogramowaniem zamkniętym (prywatnym), choć elementy takiej właśnie organizacji danych są też zawarte w Hadoop).

Dopiero wobec tak zgromadzonych i zorganizowanych danych można zastosować narzędzia drugiego etapu, tj. narzędzia analityczne, w tym także prognostyczne, których istnieje duży wybór w zależności od charakteru danych i oczekiwanych rezultatów. Począwszy od prostych analiz statystycznych, tj. monitoringu danych napływających np. w czasie rzeczywistym w określonych odstępach czasu (np. co kilka sekund), dalej możliwe jest zastosowanie do nich bardziej zaawansowanych narzędzi analitycznych. Jak: drzewa klasyfikacyjne, wspomniane wcześniej sieci neuronowe (rekurencyjne), radialne, techniki klastrowe ustalające najbliższe sąsiedztwo określonych danych umożliwiające identyfikację zbliżonych danych według ich stopnia podobieństwa itd. Bardzo istotną rolę w całej konstrukcji projektu odgrywa oprogramowanie pośredniczące³⁴ między systemem operacyjnym a aplikacjami, tj. tzw. szyny integracyjne (ang. *middleware*) oraz tzw. API's – interfejsy programistyczne aplikacji (*Application Programming Interfaces*), które muszą zagwarantować swobodny przepływ danych między wszystkimi elementami platformy w środowisku rozproszonym, ale także zapewnić ich bezpieczeństwo zgodnie z zadanymi warunkami. Oprogramowania te zorientowane są na świadczone przez system usługi dla klienta (*Service Oriented Architecture – SOA*)³⁵. Projektują je zwykle osobne, wyspecjalizowane firmy.

Trzeci etap obejmuje prezentację danych za pomocą określonych narzędzi raportowania i wizualizacji w sposób statyczny lub dynamiczny.

HQDw&usg=AFQjCNGw_wSaaqGr32beq8POCxIne13KpA&bvm=bv.92885102,d.bGQ (dostęp 7.05.2015).

³³ J. Strickland, *How the Google File System works*, <http://computer.howstuffworks.com/internet/basics/google-file-system.htm> (dostęp 2.05.2015).

³⁴ D. Serain, *Middleware and enterprise application integration: the architecture of e-business solutions*, Springer Verlag, London 2002.

³⁵ T. Eri, *Service-Oriented Architecture: concepts, technology, and design* [e-book], NY, Prentice Hall Professional Technical Reference, Upper Saddle River 2005.

Aby cały zamierzony cel został osiągnięty, na poszczególnych poziomach, począwszy od infrastruktury fizycznej aż po prezentację rezultatów, architektura projektu musi być starannie przemyślana i zaprojektowana z uwzględnieniem indywidualnych wymagań klienta i zapewnienia skalowalności całego systemu.

Zbiory *Big Data* i technologia chmury obliczeniowej w humanistyce

Opisane wyżej przetwarzanie danych zapoczątkowane w sferze biznesu i obejmujące dane ustrukturalizowane stopniowo zaczęto uzupełniać przetwarzaniem danych nieustrukturalizowanych, tj. tekstów pochodzących z portali społecznościowych, e-maili klientów, nagrań dźwiękowych, blogów czy innych źródeł internetowych, a później także z zawartości dokumentów, danych graficznych, przestrzennych, bibliotek cyfrowych itd. Do przetwarzania tego typu danych stworzono z kolei narzędzia ogólnie zwane NoSQL (*Not only SQL*)³⁶.

Systemy NoSQL działają na bazie HBase, rozszerzonej wersji HDFS. HBase w stosunku do HDFS jest rozbudowana horyzontalnie i ma możliwość obsługiwanie znacznie większych ilości danych niż HDFS. Baza składa się z jednostki głównej (*Master Node*) i serwerów regionalnych (*Regional Servers*), na których posadowione są kolejne HFiles.

Dla humanistów niewątpliwie najbardziej interesujące jest przetwarzanie zawartości dokumentów znajdujących się w zasobach bibliotek cyfrowych oraz Internetu. Trwający już od lat proces digitalizacji zasobów bibliotecznych osiągający w niektórych krajach wysoki procent całości zbiorów, tworzenie bibliotek cyfrowych, cyfrowych zasobów czasopism naukowych i repozytoriów literatury naukowej, a także rosące lawinowo zasoby Internetu umożliwiają stosowanie narzędzi z obszaru ICT do ich analizy. Niewyobrażalny przyrost dokumentów podwajający się w coraz krótszym czasie skutkuje tym, że właściwie już żaden człowiek nie jest w stanie przeczytać wszystkich dokumentów, np. literatury naukowej nawet tylko ze swojej dziedziny, jej liczebność bowiem staje się tak duża i różnorodna (choćby ze względu na różnorodność języków), że jest to zadanie niewykonalne. Mówiąc obrazowo, człowiek stara się teraz zaprząć do czytania i analizowania dokumentów maszyny i narzędzia informatyczne wyższej generacji niż to było dotychczas. Za pomocą ekosystemu *Big Data* i *cloud computing* człowiek przenosi się w sferę, w której maszyna wyręczy go w dotychczasowych czynnościach czytania i analizowania

³⁶ A. Fowler, *NoSQL for dummies*, Jon Wiley & Sons, Inc. Hoboken (NY) 2015.

dokumentów w celu wytworzenia nowej wiedzy. Warunkiem jest jednak elektroniczna wersja dokumentów i zawartych w nich danych, toteż digitalizacja jak największej liczby dokumentów jest etapem wstępnym tych działań, jak i warunkiem *sine qua non* osiągnięcia sukcesu.

Do przetwarzania tekstów w tym ekosystemie stosuje się różne nowe narzędzia zaliczane do tzw. *Natural Language Processing* (NLP)³⁷. Nie chodzi tu o znane nam dotychczas systemy wyszukiwania informacji, jak np. systemy słów kluczowych, haseł przedmiotowych, itd., choć posłużenie się nimi może stanowić pierwszy etap naszego działania. Zgodnie z dotychczasową praktyką, po zrealizowaniu wyszukiwania za pomocą np. znanego nam katalogu przedmiotowego czy słów kluczowych i po uzyskaniu dostępu do wyszukanych dokumentów, sami siadaliśmy do pracy i analizowaliśmy te materiały. To na tym etapie wytwarzaliśmy nową wiedzę. W nowym ekosystemie chodzi o takie narzędzia, które zastąpią nas na tym drugim etapie pracy, tj. wybiorą i dokonają analizy dokumentów zgodnie z naszym zapotrzebowaniem i w wyniku tego działania same wytworzą nową wiedzę.

Jeśli wziąć pod uwagę, że narzędzia, takie jak Hadoop, umożliwiają zlokalizowanie, gromadzenie i organizowanie niewyobrażalnie wielkich zbiorów danych, głównie ustrukturalizowanych, których prawdopodobnie nigdy nie zdołalibyśmy zgromadzić i zanalizować sami, to do systemów NoSQL dla sfery tekstów można zaliczyć m.in. Cassandrę³⁸, platformę rozwiniętą początkowo dla Google, Amazon i Facebooka, w celu wzmocnienia mechanizmów wyszukiwania, a od 2009 r. będącą w gestii Apache i używaną przez IBM, Thompson-Reutersa i Twittera.

Apache Cassandra wzmocniona mechanizmami MapReduce z Hadoop oraz rozszerzonymi mechanizmami BigTable umożliwia lokalizowanie i porządkowanie ogromnych ilości danych, operując Cassandra Query Language³⁹, pracującym w środowisku rozproszonym (klastrow serwerów). Jest ona wysoce skalowalna oraz zapewnia redundancję uniemożliwiającą unieruchomienie systemu, gdy jeden z elementów (klastrow) okaże się niewydolny. Ponieważ posiada mechanizmy zarówno SQL jak i NoSQL, jest wykorzystywana przez takie firmy, jak np. DataStax⁴⁰, gdzie ustrukturalizowaną informację biznesową uzupełnia się informacjami pochodzącymi z analizy zawartości portali społecznościowych.

³⁷ G.G. Chowchury, *Natural language processing* "Annual Review of Information Science and Technology", vol. 37, no. 1/2003, s. 51–89.

³⁸ M. Brown, *Learning Apache Cassandra*, Packt Publishing Ltd., Birmingham 2015.

³⁹ Ibidem.

⁴⁰ Datastax, <http://www.datastax.com/> (dostęp 2.05.2015).

W ramach NPL stosuje się różne techniki, takie jak: syntaktyczną analizę zdań, morfologiczną i semantyczną analizę wyrazów, ekstrakcję terminów, faktów, wydarzeń, nazw geograficznych, a nawet wyławia się postawy uczuciowe. Stosuje się metody statystyczne, włącza elementy tłumaczenia maszynowego i uczenia się systemów. Można też w efekcie tych działań ustrukturalizować wyniki i na ich podstawie dokonywać porównań, analiz i prognoz. Zestawem różnego rodzaju technik NLP dysponuje np. IBM Watson⁴¹, superkomputer, który, posługując się różnymi algorytmami, jest już w stanie komunikować się z człowiekiem w języku naturalnym. Jest on w stanie zanalizować zawartość 1 mln książek na sekundę, a pierwsze jego praktyczne zastosowanie to analiza ogromnej ilości naukowej literatury medycznej w celu wspomaganie i formułowanie diagnostyki medycznej raka, szczególnie w sytuacjach natychmiastowej potrzeby informacji w warunkach klinicznych. System operuje ponad 100 różnymi technikami analizy języka naturalnego i w tym języku porozumiewa się także z użytkownikiem. Watson posiada architekturę zarządzania danymi nieustrukturalizowanymi (*Unstructured Information Management Architecture – UIMA*)⁴² firmy Apache i elementy silnika wyszukiwawczego Apache Hadoop służącego do działania w wirtualnym środowisku rozproszonym. Często, w celu zrealizowania indywidualnego zapotrzebowania użytkownika/klienta, firmy, a nawet sami użytkownicy dopisują odpowiednie fragmenty oprogramowania w celu wykonania przez system konkretnego zadania.

Jeśli chodzi o zawartość bibliotek cyfrowych i w ogóle różnego typu dokumentów w wersji cyfrowej użytych ewentualnie do tego typu analiz, których liczba nie jest już żadną przeszkodą dla tych systemów, to zasadniczym problemem jest obecnie sprawa praw autorskich. Wobec tego, że legislacja nie zawsze nadąża za rozwojem techniki, dotychczasowe prawo autorskie jeszcze nic nie mówi o tym, czy dokonywanie takich analiz za pomocą tego rodzaju silników na zbiorach *Big Data* w ekosystemie *cloud computing* jest czy nie jest dopuszczalne. Z drugiej strony można jednak przyjąć, że zamiast człowieka korzystającego z dokumentów na prawach bibliotecznych i tworzącego w wyniku tego procesu nową wiedzę, czyni to za niego maszyna. Trudno powiedzieć, czy jest tu jakaś różnica, czy nie. Aby nie złamać prawa autorskiego działanie silnika *Big Data*, na dokumentach musiałyby odbywać się na zasadach bibliotecznych lub być wykonywane na dokumentach ze sfery publicznej. Stąd też niektórzy eksperymenciści w świecie bibliotek na razie, wobec braku uregulowań w sferze prawa

⁴¹ IBM Watson, <http://www.ibm.com/smarterplanet/us/en/ibmwatson/> (dostęp 2.05.2015).

⁴² The Apache Software Foundation. *Getting started: Why Apache UIMA*, <https://uima.apache.org/doc-uima-why.html> (dostęp 2.05.2015).

autorskiego, ostrożnie przeprowadzają lub umożliwiają przeprowadzanie takich działań tylko na wybranych i wydzielonych zbiorach zaliczających się do sfery publicznej. Tak na przykład postąpiła the British Library, która przy współpracy z University College of London przystąpiła do działań eksperymentalnych, używając do tego celu platformy Microsoft Azure⁴³. Platforma ta, udostępniająca przede wszystkim usługi PaaS (Platform as a Service), jest w stanie zestawić dowolną architekturę narzędzi w celu wykonania określonego zadania. Inne biblioteki, jak np. amerykańskie, wprowadzają blokady dla każdego indywidualnego dokumentu objętego prawem autorskim, jeśli działanie silnika miało odbywać się na innych niż biblioteczne zasadach. W tym miejscu musimy sobie uświadomić, że to właśnie biblioteki będą zapewne tworzyć tego typu centra badawcze, tj. pracownie, gdzie dostępne będą omawiane wyżej narzędzia. Ponieważ jednak najczęściej dostęp do platform i narzędzi *Big Data* w ekosystemie *cloud computing* nie jest wolny od opłat – biblioteki będą pewnie musiały wykupywać licencje na korzystanie z tych narzędzi, które naukowcom/użytkownikom umożliwią prowadzenie tego typu badań. Centra takie, zawiązywane często jako konsorcja bibliotek, będą musiały też być w stałym kontakcie z firmami macierzystymi, które mogą gwarantować w ramach umów licencyjnych dostosowywanie narzędzi do indywidualnych potrzeb badaczy/użytkowników.

Humanistyka, która formalnie definiowana jest jako sfera dotycząca człowieka funkcjonującego w kontekście społecznym, włączając w to jego twórczości⁴⁴, jest tu raczej desygnatem szeroko pojętej produkcji treści prezentowanych za pomocą obszernych dokumentów, tj. książek, artykułów w czasopiśmie itd., w których mowa jest o myślach, ideach, pojęciach, koncepcjach, słowach, terminach, języku, wydarzeniach, osobach, gramatykach, tożsamościach, charakterach, lokalizacjach i miejscach – występujących w całkowicie nieustrukturalizowanej formie. Teksty te dają się podzielić na tomy, rozdziały, sekcje, a nawet strony, co mechanizmowi, takim jak BigTable, umożliwia wyodrębnienie mniejszych jednostek (np. pojedynczych stron) i błyskawiczne, paralelne, hierarchiczne ich przetwarzanie. Innymi możliwymi wymiarami są np.: język, organizacja myśli w ramach poszczególnych języków, okresy czasu. Wszystko to ułatwia działanie silnikom *Big Data*, które są w stanie błyskawicznie wydobyć wszystkie elementy ich kontekstu, zanalizować je i zaprezentować w do-

⁴³ J. Baker, *The British Library big data experiment. The British Library. Digital scholarship blog*, <http://britishlibrary.typepad.co.uk/digital-scholarship/2014/06/the-british-library-big-data-experiment.html> (dostęp 2.05.2015).

⁴⁴ A.R. Rogers, *The Humanities*, 2nd ed., Libraries unlimited, Inc., Littleton (Colorado) 1979.

godny dla badacza sposób. Warunkiem jednak ich skutecznego działania jest podział analizowanych tekstów na małe porcje, co wymaga dużych pojemności pamięci na procesy indeksowania. Operowanie całymi dokumentami nie zapewni oczekiwanych rezultatów.

Silniki te, jeśli nawet nie zadowolą nas w stu procentach, to są jednak już dziś w stanie zidentyfikować i zanalizować całość interesującego nas oraz dostępnego materiału na dany temat i przynajmniej wskazać kierunek dalszych działań. Najprawdopodobniej w najbliższym czasie żaden naukowiec nie będzie już w stanie prowadzić badań bez zaprzęgnięcia silników tego typu do pracy. W innym przypadku zabraknie mu życia, znajomości języków i czasu na analizę relewantnych źródeł, nie mówiąc o tym, że badania w tym kontekście prowadzone „na piechotę” mogą się okazać niepełne, dublowane i zbaczające z głównego nurtu wobec niemożności uwzględnienia przez badacza ważnych materiałów, których sam nigdy by nie zdobył i nie był w stanie zanalizować. Biblioteki muszą przygotować się do świadczenia tego rodzaju nowych usług, a warunkiem jest całkowita digitalizacja zbiorów. Na marginesie należy dodać, że w literaturze przedmiotu jako *Big Data* w sferze tekstów określa się liczebność zbiorów na minimum 1 mln zdigitalizowanych dokumentów. Do tego ekosystemu nie można zatem zaliczyć np. projektu Gutenberg, który zawiera tylko nieco ponad 45 tys. zdigitalizowanych dokumentów⁴⁵.

Przykład

Dobrym przykładem takiej właśnie cyberinfrastruktury jest The HathiTrust Research Centre (HTRC)⁴⁶, które to centrum jest partnerstwem dwóch uniwersytetów: Indiana University oraz University of Illinois posiadających łącznie ok. 10 mln zdigitalizowanych dokumentów pełnotekstowych na 2013 r. (w tym ponad 5 mln książek i ponad 274 tys. tytułów wydawnictw ciągłych), co daje łącznie ponad 3,5 mld stron tekstu. Zbiór ten w wersji papierowej zajmuje 125 mil półek, waży ponad 8,5 tys. ton, a w wersji elektronicznej odpowiada temu 473 TB powierzchni dyskowej. Dominują tam dokumenty w zakresie: języka, literatury, historii, socjologii, biznesu i ekonomii. Jest to The HathiTrust Digital Library – biblioteka cyfrowa o rozproszonej architekturze posadowiona w chmu-

⁴⁵ Project Gutenberg, *Free e-books – Projekt Gutenberg*, <https://www.gutenberg.org/> (dostęp 2.05.2015).

⁴⁶ S.T. Kowalczyk [i in.], *Big data at scale for digital humanities; An architecture for The HathiTrust Research Center*, w: *Big data management. Technologies, and applications*, red. W. Ch. Hu, N. Kaabouch, IGI Global, Hershey (USA) 2014, s. 270–294 (*Information Science Reference*).

rze. W ramach tego dziesięciomilionowego zbioru nieco ponad 3 mln to dokumenty ze sfery publicznej⁴⁷.

System pracuje na platformie Apache Cassandra oraz, w celu ułatwienia współpracy systemu operacyjnego z aplikacjami, posługuje się oprogramowaniem *middleware* WSO2, również firmy Apache (jest to tzw. szyna integracyjna)⁴⁸, gwarantującym właściwe usługi SOA. Do wyłaniania przez badacza wstępnej kolekcji dokumentów, które znajdują się w polu jego zainteresowania, służy otwarte oprogramowanie katalogowe Blacklight⁴⁹, umożliwiające przeszukiwanie zbioru zarówno za pomocą metadanych, jak i elementów pochodzących z pełnych tekstów dokumentów. Jest ono oparte na mechanizmach webowej, skalowalnej aplikacji „Ruby and Rails”⁵⁰ i współpracuje z oprogramowaniem indeksującym SOLR. Apache SOLR⁵¹, oparty na produkcie firmy Lucene, jest otwartym oprogramowaniem umożliwiającym indeksowanie danych pochodzących zarówno z metadanych o charakterze bibliograficznym, jak i ze zbiorów nieustrukturalizowanych. Za pomocą tych narzędzi tworzona jest wstępna kolekcja przygotowana do analizy. W zakresie metadanych system ma do dyspozycji nie tylko rekordy w formacie MARC (*Machine Readable Cataloguing*), ale także w formacie METS (*Metadata Encoding and Transmission Standard*)⁵². Z kolei do analizy danych i wizualizacji wyników stosuje się różne narzędzia aplikacyjne. Dla humanistyki polecane jest np. stworzone specjalnie w ramach projektu The Software Environment for the Advancement Scholarly Research (SEASR)⁵³ narzędzie pod nazwą Meandre⁵⁴, umożliwiające m.in. ilościową analizę słów w dokumentach oraz inne operacje, w tym także prezentowanie wyników np. w postaci chmury tagów, Simile Timeline⁵⁵ do wyławiania i prezentowania dat czy Google map⁵⁶ do wyławiania i prezentacji lokalizacji. Oprogramowanie to zostało przygotowane przy pomocy amerykańskiej fundacji wspierającej działalność

⁴⁷ Ibidem.

⁴⁸ BlueSoft Sp. z o.o., *WSO2 Enterprise Service Bus*, 2015, <http://www.bluesoft.net.pl/pl/technologie/wso2-esb/> (dostęp 8.05.2015); S.T. Kowalczyk [i in.], op.cit.

⁴⁹ Blacklight, <http://projectblacklight.org/> (dostęp 2.05.2015).

⁵⁰ B. Burd, *Ruby on Rails for dummies*, Wiley Publishing Inc., Hoboken, NY 2007.

⁵¹ SOLR, <http://lucene.apache.org/solr/> (dostęp 3.05.2015).

⁵² The Library of Congress, *Metadata Encoding and Transmission Standard*, ostatnio aktualizowane 17.04.2015, <http://www.loc.gov/standards/mets/> (dostęp 1.05.2015).

⁵³ SEASR, <http://www.seasr.org/documentation/> (dostęp 3.05.2015).

⁵⁴ Meandre, <http://www.seasr.org/meandre/> (dostęp 1.05.2014).

⁵⁵ *Reference Documentation for Timeline*, http://simile-widgets.org/wiki/Reference_Documentation_for_Timeline (dostęp 4.05.2015).

⁵⁶ P. Derbyshire, A. Derbyshire, *Getting startED with Google Apps*. [New York], Apress, 2010 (s. 489–515: Google Maps).

bibliotek, tj. The Andrew W. Mellon Foundation⁵⁷. Stosowane są też inne, zaawansowane algorytmy NLP oraz algorytmy do porównywania wyników uzyskanych z różnych przebadanych kolekcji. W projekcie HathiTrust Research Centre odpowiednie algorytmy badacz może wybrać z zestawu dostępnego na ekranie interfejsu systemu. Sam w ten sposób może kierować swoim badaniem, dobierając odpowiednie narzędzia i ukierunkowując system tak, aby uzyskać potrzebne wyniki.

Wyłania się tu wiele problemów związanych z dostępnością, jak i bezpieczeństwem danych, takich jak: wstępna identyfikacja użytkownika, dostęp do dokumentów objętych prawem autorskim, zabezpieczenie dostępu do rezultatów badań wykonanych przez poszczególnych badaczy (stworzonej w ten sposób nowej wiedzy) itd. Niemniej pierwsze próby zrealizowane dotychczas okazują się udane, a w ich wyniku otrzymujemy nową, wartościową wiedzę⁵⁸.

Zapewne z czasem systemy tego typu będą ulepszone i optymalizowane. Może doczekamy też chwili, gdy maszyna, po wykonaniu pracy badawczej, napisze za nas książkę! Wtedy to, co nazywany „pracą naukową” w obszarze humanistyki, osiągnie etap całkowitej automatyzacji, a rola człowieka, tak jak w innych zautomatyzowanych dziedzinach, ograniczy się do sterowania tym procesem.

Rozwój humanistyki cyfrowej wspiera też Unia Europejska. Zagadnienie to będzie przedmiotem osobnego opracowania.

Bibliografia

- Alber A.F., *Videotex/teletext: principles and practices*, McGraw Hill Book Company, New York 1986.
- Amazon Web Services, *Amazon Elastic Computer Cloud. User Guide for Linux. API version 2014-10-01*, http://www.google.pl/url?sa=t&rct=j&q=&esrc=s&source=web&cd=12&ved=0CG8QFjAL&url=http%3A%2F%2Fawsdocs.s3.amazonaws.com%2FEC2%2Flatest%2Fec2-ug.pdf&ei=jZdLVEi2LuTgyQPH_4CoBA&usq=AFQjCNFIg7f8P673xBxButxQueralRAHkg&bvm=bv.92885102,d.bGQ.
- Baker J., *The British Library big data experiment. The British Library. Digital scholarship blog*, <http://britishlibrary.typepad.co.uk/digital-scholarship/2014/06/the-british-library-big-data-experiment.html>.
- Bakshi K., *Technologies of big data*, w: *Big data management. Technologies, and applications*, red. Wen-Chen Hu, N. Kaabouch, IGI Global, Hershey (USA) 2014 (*Information Science References*).

⁵⁷ S.T. Kowalczyk [i in.], op.cit.

⁵⁸ Ibidem.

- Baldeschfieler E., *Hadoop at Yahoo!* http://www.google.pl/url?sa=t&rct=j&q=&esrc=s&source=web&cd=6&ved=0CFMQFjAF&url=http%3A%2F%2Fwww.hadooper.cn%2Fdct%2Fattach%2FY2xiOmNsYjpwZGY6ODE%3D&ei=xX5LVdLpGYupygOGiYHQDw&usg=AFQjCNGw_wSaggGr32beq8POCxlne13KpA&bvm=bv.92885102,d.bGQ.
- Berson A., Smith S.J., *Data warehousing, data mining, and OLAP*, Mc Graw-Hill New York 2001.
- Big data management, technologies, and Applications*, red. Wen-Chen Hu, N. Kaabouch, Information Science Reference (IGI Global), Hershey 2014 (Seria: Advances in Data Mining and Database Management).
- Blacklight*, <http://projectblacklight.org/>.
- BlueSoft Sp. z o.o. *WSO2 Enterprise Service Bus*, <http://www.bluesoft.net.pl/pl/technologie/wso2-esb/>.
- Brown M., *Learning Apache Cassandra*, Packt Publishing Ltd., Birmingham 2015.
- Burd B., *Ruby on Rails for dummies*, Wiley Publishing Inc., Hoboken, NY 2007.
- Business intelligence*, red. J. Suma, Wydawnictwo Naukowe PWN, Warszawa 2013.
- China A., *Understanding the principles of recursive neural networks. A generative approach to tackle model complexity*, <http://arxiv.org/ftp/arxiv/papers/0911/0911.3298.pdf>.
- Chodkowska-Gyuriks A., *Hurtownie danych: teoria i praktyka*, Wydawnictwo Naukowe PWN S.A., Warszawa 2014.
- Chowchury G.G., *Natural language processing*, Annual Review of Information Science and Technology”, vol. 37, no. 1/2003, s. 51–89.
- Ciurana E., *Developing with Google App Engine*, Apress, Distributed by New York Springer Verlag, Berkeley (CA) 2008.
- Datastax*, <http://www.datastax.com/>.
- David M.M., Feserman L., *Advanced standard SQL dynamic structured data modeling and hierarchical processing* [e-book]. Boston, Artech House, 2013, <http://web.a.ebscohost.com/ehost/ebookviewer/ebook/ZTAWMHh3d19fNzUzNTk2X19BTg2?sid=96ce560b-b590-47f4-9f8c-18783f9988e6@sessionmgr4002&vid=3&format=EB&rid=1>.
- Derbyshrine P., Derbyshrine A., *Getting start ED with Google Apps* [New York], Apress, 2010 (s. 489–515: Google Maps).
- Eri T., *Service-Oriented Architecture: concepts, technology, and design* [e-book], Prentice Hall Professional Technical Reference, Upper Saddle River, NY 2005.
- Fowler A., *NoSQL for dummies*, Jon Wiley & Sons, Inc. Hoboken (NY) 2015.

- Hadoop, *HDFS Architecture guide*, http://hadoop.apache.org/docs/r1.2.1/hdfs_design.html.
- Hodges A., *Alan Turing: Enigma*, Wydawnictwo Albatros Andrzej Kuryłowicz S.C., Warszawa 2014.
- Hurwitz J., Nugent A., Halper, F. Kaufman M., *Big data for dummies*, John Wiley & Sons. Inc., Hoboken, NY 2013.
- IBM Watson*, <http://www.ibm.com/smarterplanet/us/en/ibmwatson/>.
- Jackson P., *Introduction to experts systems*, Accisson Wesley Pub.co, Workingham (England) 1986.
- Komisja Europejska. Komunikat Komisji do Parlamentu Europejskiego, Rady, Europejskiego Komitetu Ekonomiczno-Społecznego i Komitetu Regionów *Europejska agenda cyfrowa*, KOM(2010) 245 Bruksela 19.05.2010 wersja ostateczna.
- Komisja Europejska. Komunikat Komisji do Parlamentu Europejskiego, Rady, Europejskiego Komitetu Ekonomiczno-Społecznego i Komitetu Regionów w sprawie wykorzystania potencjału chmury obliczeniowej w Europie, Bruksela (2012), KOM(2012) 529 wersja ostateczna.
- Kowalczyk S.T. [i in.], *Big data at scale for digital humanities; An architecture for The HathiTrust Research Center*, w: *Big data management. Technologies, and applications*, red. Wen-Chen Hu, N. Kaabouch, IGI Global, Hershey (USA) 2014 (*Information Science Reference*), s. 270–294.
- Lephamer H., *RFID design principles*, Artech House, Inc. Boston, London 2008.
- Mayer-Sconberger V., Cukier K., *Big data. Rewolucja, która zmieni nasze myślenie, pracę i życie*, MT Biznes sp. z o.o., Warszawa 2014.
- Meandre*, <http://www.seasr.org/meandre/>.
- Mell P., Grance T., *The NIST definition of cloud computing. Recommendation of the National Institute of Standards and Technology*, Computer Security Division, Information Technology Laboratory, National Institute of Standards and Technology, Gaithersburg, (MD) 2011, <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>.
- Microsoft Azure*, <http://azure.microsoft.com/pl-pl/>.
- Miner D., Schock A., *MapReduce design patterns*, O'Reilly, Beijing; Koln 2013.
- Morin P., *Open data structures. An introduction* [e-book] in OPEL (Open Paths to Enriched Learning). Edmonton: AU Press. 2013, <http://web.b.ebscohost.com/ehost/ebookviewer/ebook/ZTAwMHh3d19fNjM4OTU2X19BTg2?sid=d8b94be4-7c8e-4f69-9727-d3428c90cdfc@sessionmgr115&vid=3&format=EB&rid=1>.
- Mosco V., *To the cloud. Big data in a turbulent world*, Paradigm Publishers, Boulder (Colorado), London 2014.

- NFC Forum: *NFC Data Exchange Format (NDEF). Technical specification. 2006-07-24*, NFC Forum, Inc., <http://www.eet-china.com/ARTICLES/2006AUG/PDF/NFCForum-TS-NDEF.pdf?SOURCES=DOWNLOAD>.
- Osowski S., *Sieci neuronowe do przetwarzania informacji*, wyd. 2 popr. i uzup., Oficyna Wydawnicza Politechniki Warszawskiej, Warszawa 2006.
- Parfatt M.C., *What is EDI?: a guide to electronic data interchange*, PNCC Blackwell, Oxford 1992.
- Project Gutenberg, *Free e-books – Projekt Gutenberg*, <https://www.gutenberg.org/>.
- Reference Documentation for Timeline*, http://simile-widgets.org/wiki/Reference_Documentation_for_Timeline.
- Rogers A.R. *The Humanities*, 2nd ed., Libraries unlimited, Inc., Littleton (Colorado) 1979.
- Rutkowski L., *Metody i techniki sztucznej inteligencji*, wyd. 2 zmien., Wydawnictwo Naukowe PWN, Warszawa 2012.
- Saavas I.K., Sofianidon G.N, M-Tachar Kechadi, *Applying the K-Means Algorithm in Big Raw Data Sets with Hadoop and MapReduce*, w: *Big data management. Technologies, and applications*, red. Wen-Chen Hu, N. Ka-abouch, IGI Global, Hershey (USA) 2014, s. 1–22 (*Information Science References*).
- SEASR, <http://www.seasr.org/documentation/>.
- Serain D., *Middleware and enterprise application integration: the architecture of e-business solutions*, Springer Verlag, London 2002.
- SOLR, <http://lucene.apache.org/solr/>.
- Strickland J., *How the Google File System works*, <http://computer.howstuff-works.com/internet/basics/google-file-system.htm>.
- Szmit P., *Cloud computing historia, technologia, perspektywy*, Polska Agencja Rozwoju Przedsiębiorczości, Warszawa 2012, https://www.web.gov.pl/g2/big/2012_06/ebfa211f1a9f174c7517738f68df2d8b.pdf.
- The Apache Software Foundation, *Getting started: Why Apache UIMA*, <https://uima.apache.org/doc-uima-why.html>.
- The Library of Congress, *Metadata Encoding and Transmission Standard*, ostatnio aktualizowane 17.04.2015, <http://www.loc.gov/standards/mets/>.
- Wiener N., *Cybernetics: or control and communication in the animal and the machine*, M.I.T. Press, New York 1961.

Słowa kluczowe: humanistyka cyfrowa, technologia chmury obliczeniowej, zbiory *Big Data*

Key words: Digital Humanities, the Cloud Computing Technology, *Big Data Sets*

Big Data Sets and the Cloud Computing Technology for Digital Humanities

Abstract

The main aim of the paper is to draw a concept of sets of *Big Data* which occur in the sphere of digital humanities, to describe new IT tools applied for its analysis and the conditions of its usage in the environment of the *cloud computing* technology. The technological development, which took place in the recent years in the field of information and communications technologies (ICT), especially the change in the way of governing of data in relation to the development of the *cloud computing* technology enables nowadays gathering unimaginable amount of data (*Big Data sets*) being created in various spheres of human and machine world. The definition and the origin of development of such sets of data are presented in the paper as well as the construction of the ecosystem of the *cloud computing* technology in which *Big Data sets* are gathered and processed is described. *Big data sets* are reach sources of new information. Advanced IT analytical tools are applied to investigate its content in order to obtain a new knowledge. In the paper the emphasize is on sets of *Big Data* gathered within the digital humanities area and chosen examples of IT analytical tools currently applied for these purposes are described. Also, the conditions related to the usage of such new information infrastructure in relation to digital humanities are highlighted.