

Teksty Drugie 2014, 2, s. 106-128



Pierwszy rzut oka na stylometryczną mapę literatury polskiej

Jan Rybicki

 Jan Rybicki

Pierwszy rzut oka na stylometryczną mapę literatury polskiej

Wstęp

Nie ma lepszego sposobu na przypisanie tekstów ich autorom – oczywiście jeżeli z jakiegoś powodu nie możemy lub nie chcemy korzystać ze stron tytułowych – niż policzenie częstości kilkudziesięciu (lepiej: kilkuset) najczęstszych słów i zastosowanie do niej tej czy innej metody statystycznej (np. analizy skupień), najlepiej bez sprowadzania tych wyrazów do ich form podstawowych (bo w tekście literackim różnica między np. „mówi” a „powiedział” jest stylistycznie bardzo znacząca). Tradycja ilościowych metod stylometrycznej atrybucji autorskiej i stylometrycznej chronologii tekstów sięga jeszcze czasów przedkomputerowych¹, choć oczywiście pojawienie

Jan Rybicki – dr, pracownik IFA UJ, członek Komitetu Wykonawczego European Association for Digital Humanities. Prowadzi badania tekstów literackich stosując metody stylometrii komputerowej. Najnowsze publikacje to *Stylometryczna niewidzialność tłumacza* (2013) i *Collaborative Authorship: Conrad, Ford and Rolling Delta* (2014). Kontakt: jkrybicki@gmail.com

1 Do klasyki ilościowych badań tekstowych należy zaliczyć przede wszystkim Wincentego Lutostawskiego i jego *The origin and growth of Plato's logic: with an account of Plato's style and of the chronology of his writings* (Longmans, Green, London 1897); już wcześniej chronologią Platona zajmowali się w podobny sposób Lewis Campbell (*The Sophistes and Politicus of Plato*, Clarendon Press, London 1867) i Constantin Ritter (*Untersuchungen über Plato. Die echtheit und chronologie der platonischen schriften, nebst anhang: Gedankengang und grundanschauungen von Platos Theätet*, W. Kohlhammer, Stuttgart 1888). Z kolei pomysł na ilościowe ustalenie autentyczności listów św. Pawła miał w 1851 roku sam Augustus de Morgan.

się „mózgów elektronowych”, elektronicznych wersji tekstów literackich (i nieliterackich) oraz nowoczesnych metod statystyki komputerowej znacznie ułatwiło takie badania² i sprawiło, że podejmuje je coraz więcej badaczy na całym świecie³.

Nie da się nie zauważyć, że określanie autorstwa na podstawie słów najczęstszych – a więc, co wynika z praw Zipfa⁴, nie tylko tych najkrótszych, ale przede wszystkim tych najmniej „znaczących” – przenosi *focus* odczytywania tekstu ze słów takich jak *miłość*, *ojczyzna*, *Polska* czy *duch* albo *dobro*, i przede wszystkim z obszerniejszych związków między słowami (długość i szyk zdania, organizacja tekstu na jeszcze wyższych poziomach hierarchii, relacje intertekstualne), na nieefektywne i częstokroć niezauważalne (a przecież stanowiące często ponad połowę tekstu) słowa „funkcyjne”, „synsemantyczne”, „gramatyczne”, takie jak choćby królujące na szczytach list rangowych *i*, *się*, *w*, *nie*, *na* itd.

Tymczasem okazuje się, że choć analiza stylometryczna bywa często utożsamiana z atrybucją autorską, te same metody często zdają się mówić rzeczy ciekawe (albo przynajmniej tylko: sensowne) o związkach między tekstami

- 2 Pierwszym zastosowaniem maszyny liczącej (jeszcze nieelektronicznej) w badaniach stylometrycznych było studium Thomasa C. Mendenhalla *A mechanical solution of a literary problem*, „Popular Science Monthly” 1901 nr 60. Przełomowa publikacja to jednak już jak najbardziej komputerowa – i oparta na zliczaniu częstości najczęstszych słów – analiza autorstwa druków ulotnych z początków państwowości amerykańskiej, którą przeprowadzili statystycy Frederick Mosteller i David Wallace, *Inference and disputed authorship: the Federalist*, Addison-Wesley, Reading, 1964. Klasycznym zastosowaniem tych metod w badaniach literackich jest praca Johna Burrowsa *Computation into criticism: a study of Jane Austen's novels and an experiment in method*, Clarendon Press, Oxford 1987.
- 3 Wspomagane komputerowo badania literaturoznawcze prowadzą m.in. Hugh Craig z Uniwersytetu w Newcastle (Australia): H. Craig, A.F. Kinney *Shakespeare, computers, and the mystery of authorship*, Cambridge University Press, Cambridge 2009; Karina van Dalen-Oskam z Instytutu Huyghensa w Hadze (Holandia): *Names in novels: an experiment in computational stylistics*, „Literary and Linguistic Computing” 2013 nr 28, s. 359-370; David Hoover z New York University (USA): *Corpus stylistics, stylometry, and the styles of Henry James*, „Style” 2007 nr 41(2), s. 174-203; *Stylistics: prospect & retrospect*, Rodopi, Amsterdam 2007; Fotis Jannidis z Uniwersytetu w Würzburgu (Niemcy) i Matthew Jockers z Uniwersytetu w Lincoln (USA): *Macroanalysis. Digital methods and literary history*, University of Illinois Press, Champaign 2013; ten ostatni jest uczniem Franka Morettiiego (Stanford University, USA): *Distant reading*, Verso, New York 2013; *Graphs, maps, trees: abstract models for a literary history*, Verso, New York 2005.
- 4 Najbardziej znanemu prawu Zipfa („frekwencja danego słowa w danym korpusie znaków języka naturalnego jest odwrotnie proporcjonalna do jego miejsca w liście rangowej słownictwa tego tekstu”) towarzyszą dwa inne, równie ważne z perspektywy badań takich jak te prezentowane w niniejszym artykule: „im częstsze słowo, tym mniej konkretne znaczenie” i „im krótsze słowo, tym częściej występuje” (por. G.K. Zipf *Human Behaviour and the Principles of Least Effort*, Addison-Wesley, Boston 1949).

różnych autorów. Jeżeli zastosowana metoda jest skuteczna, pojedyncze teksty autora A grupują się razem; osobną grupę stanowią teksty autora B... Jednak kiedy do badanego korpusu dodamy jeszcze równie pięknie łączące się teksty autora C, będą one zwykle bliższe (na wykresie) tekstom – na przykład – A. Czy to oznacza, że autorzy A i C piszą „podobnie” do siebie niż A i B? Czy też jest to zwykły artefakt metody?

Najlepiej zacząć od eksperymentu. Poniżej prezentuję wyniki analizy wzajemnych podobieństw stylometrycznych między tekstami polskich pisarek i pisarzy, bazującej na częstościach najczęstszych słów.

Materiał

Badania przeprowadziłem na korpusie około 500 polskich tekstów literackich ze znaczną przewagą wszelkich podgatunków powieściowych z lat od 1775 (*Mikołaja Doświadczyńskiego przypadki*) do chwili bieżącej (koniec 2013 roku). W korpusie znaleźli się ponadto nieliczni reprezentanci poezji epickiej (*Pan Tadeusz*, *Beniowski* i *Zamek kaniowski*) oraz innych gatunków prozatorskich – przede wszystkim zbiory opowiadań (Schulz, Odojewski); te zostały włączone do korpusu z powodu niedostępności (lub wręcz nieistnienia) dłuższych tekstów literackich autorów ważnych z punktu widzenia historii literatury polskiej. Bazując na pracy Edera⁵, w celu zminimalizowania wpływu rozmiarów utworów na wyniki analiz unikałem tekstów poniżej 10 000 słów; najkrótszym tekstem w korpusie jest więc powieść Anny Nakwaskiej, *Aniela, czyli ślubna obrączka* z roku 1831 (13 571 słów), najdłuższym – *Lód* Jacka Dukaja z roku 2007 (391 104).

Wiek XVIII reprezentowany jest – niestety! – jedynie przez wspomniane już dzieło Ignacego Krasickiego; wiek XIX – już przez 111: od *Malwiny* po *Krzyżaków*, *Argonautów* i *Komorników*. Z wieku XX pochodzi 310 tekstów od *Popiołów*, *Próchna* i *Na srebrnym globie* po *Panią jeziora* czy *Dziwięć*. Lista 82 utworów z wieku XXI rozpoczyna się od *Pod mocnym aniołem* i kończy m.in. *Ostatnim rozdaniem*. Przy sporządzeniu korpusu korzystałem przede wszystkim z darmowych źródeł internetowych i z płatnych e-booków, konwertowanych następnie do zwykłego formatu tekstowego; tylko nieliczne książki zostały zeskanowane i poddane optycznemu rozpoznaniu znaków. To pozwoliło na stosunkowo szybkie sporządzenie korpusu, ale równocześnie wykluczyło możliwość użycia wielu tekstów, które zapewne powinny znaleźć się w takiej „reprezentatywnej” czy *horribile dictu*, „kanonicznej” próbie literatury polskiej.

5 M. Eder *Does size matter? Authorship attribution, small samples, big problem*, „Literary and Linguistic Computing”, publikacja online, 14.11.2013. <http://llc.oxfordjournals.org/content/early/2013/11/14/llc.fqt066.full?sid=5d729c48-0b31-480f-8d69-focde7e195de> (dostęp: 1.01.2014).

Jeżeli więc ośmielam się nazywać ten wybór reprezentatywnym, to z zastrzeżeniem, że jest to reprezentatywna próbka materiałów dostępnych obecnie w wersji elektronicznej – ale równocześnie skorzystałem z tego pretekstu, by włączyć do korpusu pozycje z repertuaru popularnego (by nie rzec – brukowego) i/lub młodzieżowego. Natomiast w sposób bardzo nieproporcjonalny potraktowałem literaturę *science fiction* i *fantasy*, która bezapelacyjnie króluje w polskich repozytoriach internetowych; z powodów chyba nie do końca zbadanych to ulubione podgatunki osób nałogowo/zawodowo związanych z komputerami – tych wszystkich informatyków, matematyków i fizyków cząstek elementarnych.

Oczywiście teksty zostały poddane obróbce korektorskiej; natomiast zgodnie z założeniami Craiga i Whippa⁶ nie zastosowałem ujednolicenia pisowni tekstów pochodzących z różnych epok; ten problem został zminimalizowany przez zastosowanie wysokich wartości *cullingu* (o czym poniżej).

Metoda

Analizę ilościową przeprowadziłem na ciągach częstości najczęściej występujących słów – w tym przypadku tych, które znalazły się na pierwszych 100-1000 pozycjach list rangowych. W celu wyeliminowania bezpośredniego wpływu tematyki utworów na wyniki analizy zastosowałem *culling* na poziomie 90 i 100% – oznacza to, że w analizie wykorzystano tylko te słowa, które wystąpiły równocześnie albo w 90% wszystkich tekstów, albo we wszystkich tekstach korpusu. W praktyce okazało się, że wszystkie teksty miały nie więcej niż 107 wspólnych słów i tylko 827 słów, które występowały w 9 na 10 tekstów – i właśnie takie były graniczne parametry przedstawionej analizy. Obie listy słów umieściłem w Dodatkach do niniejszego tekstu; warto zauważyć, że przeważają tam wyrazy odpowiedzialne za sposób prowadzenia narracji i kształt dialogów; brak natomiast słów bezpośrednio związanych z treścią utworów – nie wspominając już o imionach i nazwiskach bohaterów czy o nazwach własnych (choćby polskich miast). Co najważniejsze, takie ograniczenie – całkowicie niearbitralne, bo regulowane wyłącznie i automatycznie przez statystykę – likwiduje (jak już wspomniano powyżej) prosty wpływ różnic w pisowni „między dawnymi a młodszymi laty” na wyniki eksperymentu⁷.

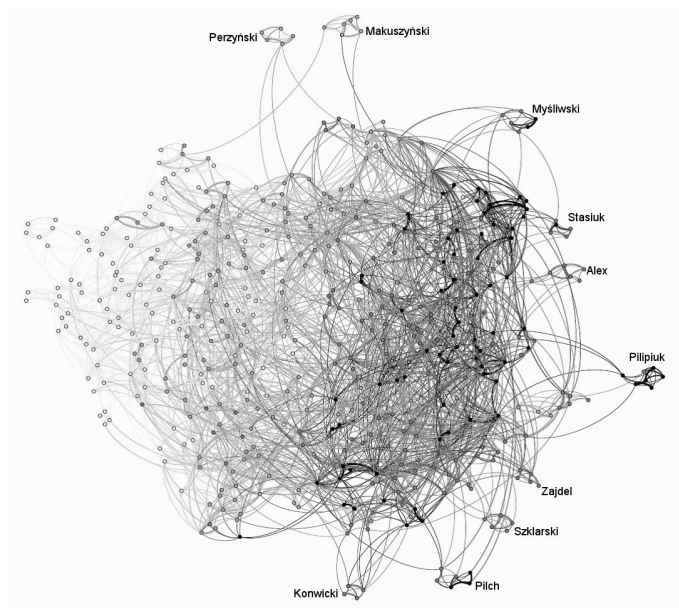
6 H. Craig, R. Whipp *Old spellings, new methods: automated procedures for indeterminate linguistic data*, „Literary and Linguistic Computing” 2010 nr 25(1), s. 37-52.

7 Więcej szczegółów na temat samej metody i jej teoretyczno-statystycznych podstaw przedstawia zamieszczony w tym samym tomie tekst Macieja Edera *Metody ścisłe w literaturoznawstwie i pułapki pozornego obiektywizmu – przykład stylometrii*.

Wszystkie te procedury – od wczytania tekstów elektronicznych po analizę statystyczną – zostały przeprowadzone za pomocą jednego pakietu „stylo” (wersja 0.5.3) stworzonego przez Macieja Edera przy niewielkim współudziale dwóch pozostałych autorów⁸; pakiet ten jest przeznaczony do środowiska programowania statystycznego R⁹. Natomiast do wizualizacji sieciowej wyników uzyskanych przez „stylo” posłużył open-source’owy program GEPHI¹⁰.

Wyniki

Już pierwszy rzut oka na wyniki analizy sieciowej słownictwa tekstów zawartych w korpusie pozwala zauważyć dość wyraźny sygnał chronologiczny (Wykres 1). Po lewej stronie wykresu przeważają jasnoszare punkty reprezentujące

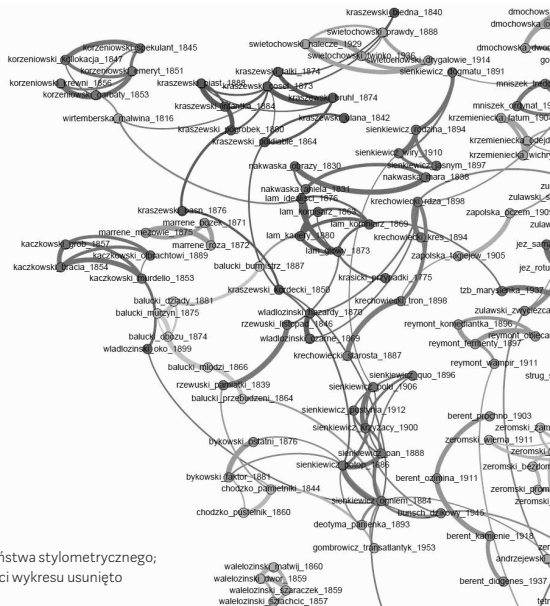


Wykres 1. Sieć podobieństwa stylometrycznego 503 powieści polskich; sygnał chronologiczny i zjawiska marginalne.

- 8 M. Eder, M. Kestemont, J. Rybicki *Stylometry with R: a suite of tools*, w: *Digital Humanities 2013: Conference Abstracts*, University of Nebraska-Lincoln, Lincoln 2013, s. 487-489.
- 9 R Core Team *A language and environment for statistical computing*. R Foundation for Statistical Computing, Wien 2013, <http://www.R-project.org/> (dostęp: 1.1.2014).
- 10 M. Bastian, S. Heymann, M. Jacomy *Gephi: an open source software for exploring and manipulating networks*. International AAAI Conference on Weblogs and Social Media 2009.

literaturę polską do końca XIX wieku. Czarne punkty to dzieła wydane po roku 2000; między nimi rozciąga się szara strefa literatury XX-wiecznej; w niej jednak widać z kolei wyraźny podział na pierwszą (nieco jaśniejszą) i drugą (nieco ciemniejszą) połowę wieku XX. Oczywiście strefy sąsiadujące ze sobą chronologicznie dość silnie na siebie zachodzą, ale po pierwsze, zjawisko to nie dotyczy literatury najwcześniejszej i najnowszej; po drugie, ewolucja od tekstów wczesnych do późnych jest wyraźna i zorientowana wzdłuż poziomej osi wykresu.

O ile makrostruktura tej swoistej literackiej mapy jest zdeterminowana chronologicznie, o tyle jej mikrostruktura ujawnia siłę sygnału autorskiego. Na przykładowym wycinku sieci (Wykres 2) najsilniejsze połączenia (najgrubsze linie) występują między tekstami tych samych autorów, często przy zachowaniu różnicowania siły tych połączeń. I tak w twórczości Kaczkowskiego największe podobieństwo stylometryczne występuje między *Murdelio* i *Braćmi ślubnymi* – może dlatego, że ukazały się w odstępie zapewne roku, a mimo że tylko *Murdelio* należy do cyklu *Ostatni z Nieczujów*. Z kolei dość wczesna, bo należąca jeszcze do wołyńskiego okresu twórczości swego autora powieść obyczajowa *Cate życie biedna* nieco odbiega od późniejszej, przeważnie historycznej twórczości Kraszewskiego. Trudno jednak mówić o podziałach gatunkowych czy podgatunkowych, skoro równie daleko jest do *Króla Piasta* czy *Hrabiny Cosel* również historycznym *Starej Baśni* i *Kordeckiemu*.



Wykres 2. Wycinek sieci podobieństwa stylometrycznego; sygnał autorski. Dla przejrzystości wykresu usunieto najsilniejsze połączenia.

Inni autorzy w przedstawionym wycinku nie pozostają już jednak w takiej samej *splendid isolation*: *Malwinę* Wirtemberskiej łączy stylometryczne podobieństwo o tej samej sile z *Garbatym* Józefa Korzeniowskiego, *Hrabinią Cosel* Kraszewskiego i *Idealistami* Jana Lama. Jeszcze silniej przyciągają się *Murdelio* z *Okiem proroka* Władysława Łozińskiego – do tego stopnia, że opisy perypetii Hanusza Bystrego mają stylometrycznie znacznie mniej wspólnego z innymi dziełami młodszego z braci-pisarzy; tym różnią się więc dwaj Łozińscy od siostr Brontë, których podobieństwo stylometryczne jest znaczne i zachowuje się nawet w przekładach na inne języki¹¹. Zresztą ofiara pojedynku z Karolem Cieszewskim i jego *Zaklęty dwór* czy *Czarny Matwij* też izolują się od stylometrii innych autorów.

Równocześnie jednak sygnał chronologiczny nie jest szczególnie zauważalny w ramach twórczości pojedynczych autorów. Prostą ewolucję cech stylometrycznych u jednego pisarza – taką, w której teksty układałyby się w łańcuch chronologiczny – zaburza fakt, że w większości przypadków utwory jednego autorstwa powiązane są między sobą wielokrotnie. Jednym z nielicznych wyjątków od tej reguły jest stylometria Wacława Berenta, którego *Próchno* (1903) rozpoczyna ewolucyjny szereg przebiegający przez *Oziminę* (1911) i *Żywe kamienie* (1918), a kończący się *Diogenesem w kontuszu* (1937). Być może to nie przypadek, bo właśnie o Berencie pisze Miłosz, że „zdołał w paru powieściach zawrzeć kolejne etapy swojego zmieniającego się poglądu na świat”¹². Skoro ewolucja poglądów, to może i cech stylometrycznych?

Zupełnie inaczej przedstawia się stylometryczny obraz twórczości Sienkiewicza – ta bowiem rozpada się na dwie niemal odrębne części. W górnej części omawianego wycinka wykresu znalazły się obyczajowe (i zwykle niżej oceniane, i na ogół późniejsze) utwory pierwszego polskiego noblisty; na dole za to gromadzą się największe osiągnięcia amatora licznych Marii – jego powieści historyczno-przygodowe, i to i te pierwsze, i ostatnie (łącznie z mniej udanym *Na polu chwały*). Owe dwa osobne światy łączy tylko jedna cienka linia prowadząca właśnie od najsłabiej przyjętej (i przez samego autora niepoważanej) opowieści sprzed wiktorii wiedeńskiej, *Na polu chwały*, do kompletnego fiaska pisarza, *Wirów*. Trudno powiedzieć, czy stylometria odzwierciedla w ten sposób sąsiedztwo chronologiczne obu powieści (odpowiednio, 1910 i 1906), czy właśnie ich niski poziom literacki, czy wreszcie podobieństwo tematyczne – w końcu z „górných,” obyczajowych powieści właśnie *Wiry* są równocześnie najbardziej polityczne. Oczywiście podobieństwo stylometryczne sienkiewiczowskich romansów przygodowych byłoby najłatwiej przypisać daleko idącej

11 J. Rybicki *Stylometryczna niewidzialność tłumacza*, „Przekładaniec” 2013 nr 27, s. 61-87.

12 Cz. Miłosz *Historia literatury polskiej do roku 1939*, Znak, Kraków 1998, s. 427.

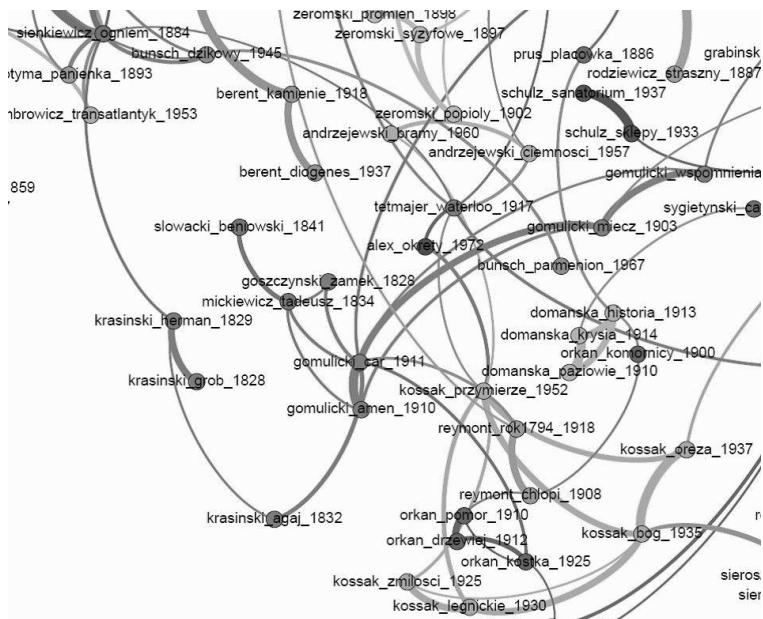
stylizacji Trylogii, *Krzyżaków* i *Quo vadis*, gdyby nie dwie sprawy. Po pierwsze, *Ogniem i mieczem*, *Potop* i *Pan Wołodyjowski* to stylizacja na staropolszczyznę XVII-wieczną; *Krzyżacy* to góralszczyzna udająca średniowiecze; *Quo vadis* przemawia językiem polskich przekładów literatury łacińskiej – więc nie jedna stylizacja, a trzy. Po drugie, trudno mówić o stylizacji języka *W pustyni i w puszczy*, nawet jeżeli jest to polska powieść stylizowana na angielską literaturę kolonialną – i nawet jeżeli „prawdziwym” językiem tej powieści jest... angielski¹³.

Ciekawie przedstawiają się związki dwóch sienkiewiczowskich skupisk z dziełami innych pisarzy. *Bez dogmatu* łączy się (poza opisywanym wycinkiem) z debiutem powieściowym Nałkowskiej (*Kobiety*, 1906), która potem oddala się znacznie od sienkiewiczowskich wzorów. Więcej interesujących rzeczy dzieje się jednak wokół „dolnego” Sienkiewicza. Silne podobieństwo stylo-metryczne między *Potopem* i *Kordeckim* musi stać się okazją do przypomnienia, że analiza została przeprowadzona wyłącznie na słowach najczęstszych (i w dodatku występujących we wszystkich lub prawie wszystkich tekstach w korpusie) – nie zaś na wyrazach, które miałyby jakikolwiek związek znaczeniowy z dwiema wizjami obrony Jasnej Góry przed Szwedami. *Kordecki* jest zresztą też podobny (tylko trochę mniej) do *Ogniem i mieczem* – obie te części Trylogii (wraz z *Krzyżakami*) przeciągają na stronę wieku XIX jak najbardziej XX-wieczny (ale przecież też historyczny i też stylizowany językowo) *Dzikowy skarb* Bunscha; bardzo silnie grawitują do siebie *Ogniem i mieczem* z *Panienką z okienka* Deotymy. Ze stylo-metrycznego punktu widzenia można upatrywać genezy języka Trylogii nie tylko (o czym już była mowa) w twórczości Kra-szewskiego, lecz również w romansach historycznych Piotra Jaksy Bykowski-go; nie da się też przejść do porządku dziennego nad linią łączącą *Ogniem i mieczem* z *Władysławem Hermanem* Zygmunta Krasińskiego. Ten obraz interak-cji Trylogii z szerszym repertuarem polskiej powieści historycznej dopełnia (poza wycinkiem) podobieństwo do *Waterloo* Przerwy-Tetmajera (1917).

Skoro wspomnieliśmy już o jednym z Trzech Wieszców... Wielka to szko-da, że tylko Krasiński popełniał powieści; z pewnym niepokojem o stabilność rodzajową korpusu włączyłem doń jeszcze *Pana Tadeusza* (którego najlepiej znany angielski przekład – pióra wybitnego amerykańskiego sławisty George’a R. Noyesa – pisany jest prozą i czyta się świetnie, zupełnie jak powieść) i *Beniowskiego* – *Maria* okazała się zbyt krótka; kryterium obszerności tekstu spełnił za to *Zamek kaniowski* (ponad 18 tysięcy słów). Polscy romantycy znaleźli się w tej samej okolicy naszej sieci (Wykres 3), ale bezpośredni związek łączy tylko Mickiewicza, Słowackiego i Goszczyńskiego – za mało danych,

13 J. Rybicki *Angielskie przekłady „W pustyni i w puszczy”*, w: *Wokół „W pustyni i w puszczy”*. W stulecie pierwodruku powieści, red. J. Axer, T. Bujnicki, Universitas, Kraków 2012, s. 555-570.

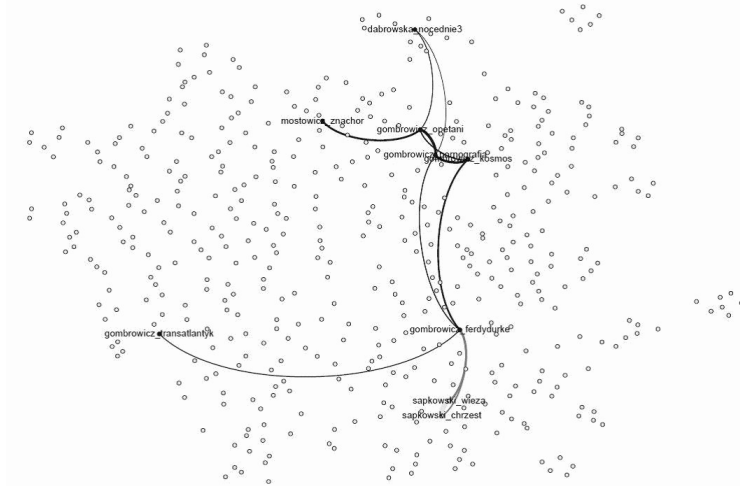
by sądzić, czy podobieństwo stylometryczne wynika ze wspólnoty rodzaju literackiego, czy z epoki literackiej. Nastoletnie powieści późniejszego autora *Nieboskiej* łączą się z romantyczną konkurencją dopiero za pośrednictwem pozytywisty – ale i poety – Gomułckiego. Bliskie sąsiedztwo do nich nie tylko *Ogniem i mieczem*, lecz (od drugiego końca) również XX-wiecznych powieści historycznych Orkana (*Kostka Napierski*), Reymonta (*Rok 1794*) i Kossak-Szczuckiej (*Przymierze*) może sugerować źródło ich archaizacyjnej stylizacji w twórczości trzech autorów wspomnianych poematów epickich. Warto zauważyć, że w lewym dolnym rogu wykresu – tego, w którym znalazły się teksty najstarsze i/lub powieści historyczne – dochodzi do lokalnego zakłócenia porządku chronologicznego.



Wykres 3. Wycinek sieci podobieństwa stylometrycznego; okolice polskich romantyków. Dla przejrzystości wykresu usunięto najsłabsze połączenia.

Ale najciekawszym zakłóceniem tego porządku – również powiązaniem z Sienkiewiczem – jest jednak bliskie sąsiedztwo i silne połączenie między *Ogniem i mieczem* i *Trans-Atlantykiem*. Rzecz niby oczywista: i Sienkiewicz, i Gombrowicz mówią w tych okolicach „Paskiem”, nawet jeżeli ten drugi czyni to w tonacji groteskowo-parodystycznej; a jednak nie należy zapominać, że na taki wynik nie miały wpływu ani szyk zdania, ani archaizowane elementy leksykalne (patrz Dodatki). Gombrowiczowski miecz parodii uderza zresztą

nie tylko w Sienkiewicza, lecz również w Ignacego Chodźkę i jego *Pamiętniki kwestarza*; efekt ten dziwić nie powinien, bo stylometryczne podobieństwo parodii literackiej i jej przedmiotu jest zjawiskiem znanym i dobrze opisanym¹⁴. W każdym razie *Trans-Atlantyk* tak skutecznie przyprawia sobie staropolską gębę, że dopływa najdalej na lewo (można powiedzieć, że na zachód!) wykresu ze wszystkich tekstów z drugiej połowy XX wieku. W ten sposób Gombrowicz jawi się jako autor stylometrycznie najbardziej proteuszowy z tych, których teksty znalazły się w omawianym korpusie. Jak widać na Wykresie 4, jego teksty – poza opisanym już spektakularnym rejsem – pojawiają się w trzech z czterech ćwiartek naszej sieci: *Ferdydurke* w prawej dolnej, *Opętani*, *Pornografia* i *Kosmos* po prawej u góry. Co ciekawe, sensacyjno-gotycki romans z 1939 roku wykazuje znaczące podobieństwo stylometryczne do *Znachora*. Czyżby dlatego, że jak Gombrowicz sam wyznaje, „zła literatura polska była dla mnie i ciekawa i pouczająca. Studiując [...] powieści Germana, Mniszkówny, Zarzyckiej, Mostowicza, odkrywałem rzeczywistość...”¹⁵. Już znacznie trudniej wytłumaczyć, dlaczego *Ferdydurke* łączy się z dwoma tekstami Sapkowskiego (*Chrzest ognia* i *Wieża jaskółki*). Albo skąd u tych samych *Opętanych* stylometryczne pokrewieństwo (prawda, że bardzo słabe) do *Nocy i dni*?

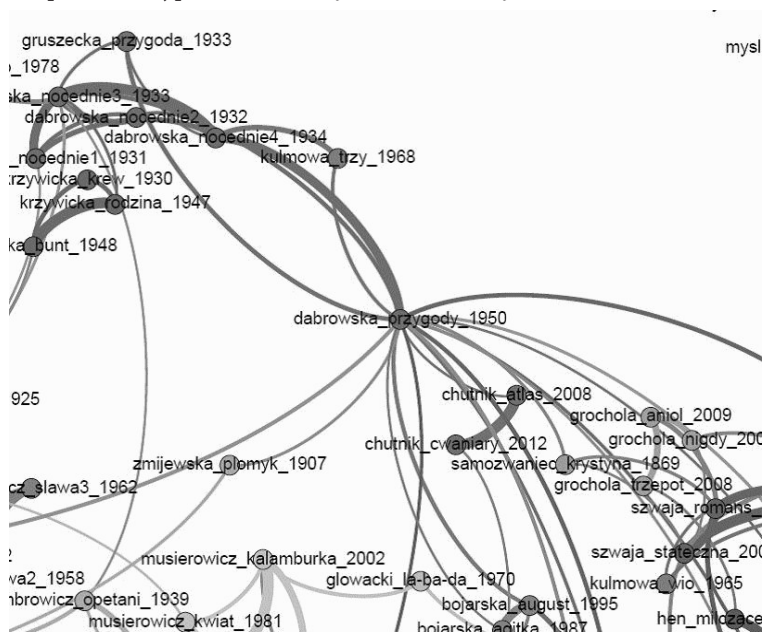


Wykres 4. Wycinek sieci podobieństwa stylometrycznego; autorski sygnał Gombrowicza. Dla przejrzystości wykresu usunięto wszystkie inne połączenia.

14 J. Burrows *Who wrote Shamela? Verifying the authorship of a parodic text*, „Literary and Linguistic Computing” 2005 nr 20(4), s. 437-450.

15 W. Gombrowicz *Dziennik 1953-1956*, Wydawnictwo Literackie, Kraków 1997, s. 108.

Skoro mowa o Marii Dąbrowskiej, to jeden z jej tekstów bije wszelkie rekordy pod względem liczby powiązań na wykresie. *Przygody człowieka myślącego* łączą się bowiem (Wykres 5) z aż czternastoma tekstami innego autorstwa. Czy dlatego, że *Przygody...* to rekonstrukcja nieukończonej powieści, połączenie gotowych rozdziałów z fragmentami brulionu autorki, pozbawione tak ważnej w twórczości Dąbrowskiej ostatniej redakcji autorskiej? Anna Kowalska pisze w swej przedmowie wręcz, że „to nie książka” ...¹⁶

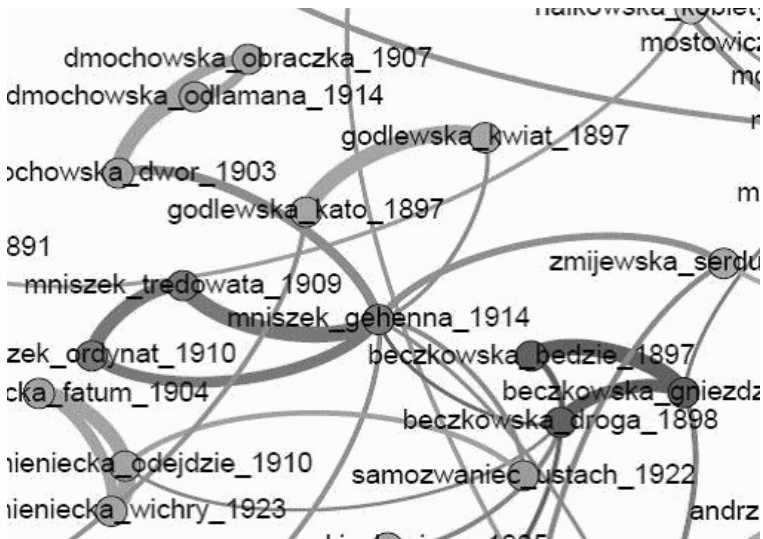


Wykres 5. Wycinek sieci podobieństwa stylometrycznego; *Przygody człowieka myślącego*. Dla przejrzystości wykresu usunięto najszabsze połączenia

Jest jeszcze jeden fragment opisywanej sieci, w którym szczególnie silne są powiązania między tekstami różnych autorów – a raczej, w tym przypadku, autorek. Chodzi o kilkanaście romansów z przełomu XIX i XX wieku autorstwa takich pisarek jak Mniszkówna, Grot-Bęczkowska, Dmochowska, Żmijewska czy Krzemieniecka. Węzłem łączącym stylometryczne cechy tej twórczości zdaje się (przynajmniej na podstawie Wykresu 6) *Gehenna*; ale równie mocno wiąże się z nimi parodia podgatunku, *Na ustach grzechu* Magdaleny Samozwaniec. I choć za główny przedmiot prześmiewczego dzieła córki

16 A. Kowalska, przedmowa do M. Dąbrowska *Przygody człowieka myślącego*, red. E. Korzeniewska, Czytelnik, Warszawa 1970, s. 9.

Wojciecha Kossaka przyjmuje się zazwyczaj *Trędowatą* (ta bowiem otwiera w dedykacji listę utworów sparodiowanych w *Na ustach grzechu*), to zapewne nie przypadek, że np. Grot-Bęczkowską wspomina Samozwaniec jako wzór własnych, nieudolnych prób literackich z dzieciństwa¹⁷.



Wykres 6. Wycinek sieci podobieństwa stylometrycznego; *Gehenna* i *Na ustach grzechu*. Dla przejrzystości wykresu usunięto najsłabsze połączenia.

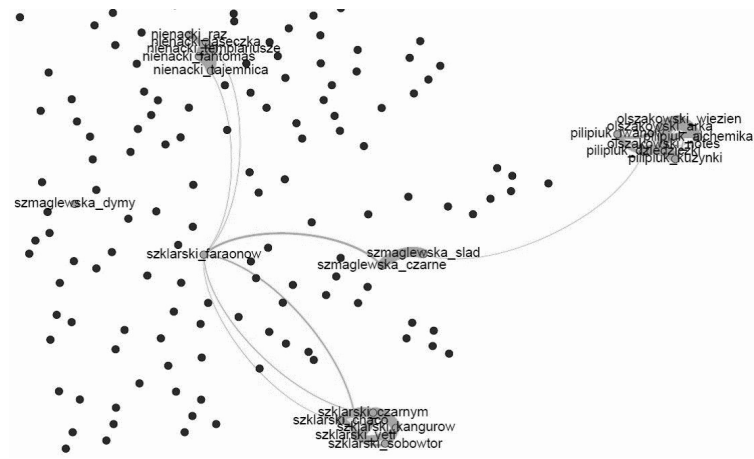
Innego rodzaju związek między autorami widać na Wykresie 7, gdzie silne podobieństwo łączy stylometrię trzech pisarzy: Kisielewskiego łączyła przecież bliska zażyłość z Tyrmandem („sporo myślałem sobie o Tyrmandzie. Lubiłem tego chłopca”¹⁸) i Iredyńskim („Iredyńskiego pijanego zgoła wczoraj wyrzuciłem na mordę”¹⁹). W tym samym wycinku sieci znalazły się też teksty Odojewskiego i Jana Józefa Szczepańskiego, ale, niestety, i Jan Józef Szczepański, jeszcze jeden z filarów pierwszego pokolenia *Tygodnika Powszechnego*, i Odojewski, o którym bardzo ciepło wyraża się w swych *Dziennikach* autor *Sprzysiężenia* („robi doskonałe wrażenie, podobno jest bardzo wybitny”²⁰), są tutaj tylko „przelotem” – ich teksty (odpowiednio, *Portki Odysa*

17 M. Samozwaniec *Maria i Magdalena*, Wydawnictwo Literackie, Kraków 1960, s. 10.

18 S. Kisielewski *Dzienniki*, Iskry, Warszawa 2001, s. 197.

19 Tamże, s. 98.

20 Tamże, s. 568.



Wykres 8. Wycinek sieci podobieństwa stylometrycznego; powieści dla młodzieży, powieści dla dorosłych. Dla przejrzystości wykresu usunięto najsłabsze połączenia.

Pilipiuk i Szklarski nie są jedynymi pisarzami wypchniętymi na zewnętrzną orbitę naszej sieci i zresztą znaleźli się tam w całkiem niezłym towarzystwie. Mamy tu bowiem (Wykres 1) nazwiska takie jak Myśliwski, Stasiuk, Pilch i *last but not least*, Konwicki. Na pewno błędem byłoby konkludować, że w związku z tym wszyscy ci pisarze mają ze sobą coś wspólnego (bo przecież stylometrycznie nie łączy ich nic i właśnie dlatego zostali wypchnięci na peryferia wykresu!); trudno jednak oprzeć się wrażeniu, że pewna przewaga wśród tych outsiderów pisarzy powszechnie i współcześnie uznanych, może sugerować jakiś związek między oryginalnością (czy choćby tylko osobliwością) stylometryczną i popularnością czy uznaniem – bo przecież konkludowanie na tych podstawach o jakiejś „jakości” czy choćby tylko „sukcesie” literackim byłoby już grubym nadużyciem. Tak czy inaczej zapewne czas najwyższy na jakieś – może „spokojniejsze” –

Wnioski

Powstała więc mapa powiązań stylometrycznych między (dokładnie) 503 polskimi tekstami literackimi wybranymi z tego, co było dostępne, w sposób, który autorowi tego skromnego studium (zresztą angliście-komparatyście) wydawał się najbardziej reprezentatywny. Jestem przekonany, że kto inny sporządzający taką samą listę wymieniłby na inne nawet i połowę tych tekstów. Nie to jednak jest w tym studium najważniejsze, bo wielkie dyskusje o istnieniu i ewentualnym kształcie kanonu, jakie przetoczyły się przez literatury zachodnie, na razie trochę jakby ominęły literaturę

polską²¹. Ważne jest to, że niezależnie od tego, jacy autorzy i z jakich epok literackich znaleźli się w korpusie, metoda analizy i wizualizacji oparta na częstościach najczęstszych słów potrafi ułożyć teksty tych autorów i tych epok w pewnym porządku, a znaczna część odstępstw od tego porządku daje się wytłumaczyć w sposób nieklózący się znacznie z tym, co wiemy z historii i komparatystyki literackiej. Podobną prawidłowość zauważyłem zresztą w analogicznym korpusie 500 tekstów literatury anglojęzycznej – i podobnie jak w przypadku omawianego tu korpusu polskiego nie jest ona prostą funkcją historycznych zmian leksykalnych w języku (jeszcze raz przydał się ten jakże prosty a jakże skuteczny zabieg cullingu)²².

Oczywiście wyniki te byłyby zapewne potraktowane znacznie poważniej, gdyby zostały uzyskane na cechach tekstu literackiego, do których jesteśmy bardziej przyzwyczajeni w praktyce interpretacji – gdyby nie opierały się na owych nieszczęsnych, pozbawionych kontekstu i w dodatku niezbyt „znaczących” słowach z górnych warstw listy rangowej. I – co należy dodać gwoli uczciwości – gdyby wiadomo było, dlaczego ciągi frekwencji słów najczęstszych tak skutecznie zdradzają autora tekstu. Jak pisze jeden z bardziej wpływowych stylometrów naszej doby, metody te „niebezpiecznie zakładają niezależności częstości jednych słów od drugich”, choć przyznaje, że „sprawdzają się mimo tego podejrzanego założenia”²³. Niewątpliwą zaletą stylometrii opartej na pojedynczych słowach jest przede wszystkim stosunkowa prostota metody i obliczeń; wszelkie próby automatycznego ustalania np. struktur zdaniowych mnożą problemy metodologiczne (szczególnie w języku tak silnie fleksyjnym, jakim jest polszczyzna); liczne próby dokonywane na n-gramach słownych wykazują znacznie mniejszą skuteczność niż na pojedynczych słowach²⁴; skuteczny *parsing* semantyczny materiału literackiego to pieśń niedalekiej już, ale jednak przyszłości²⁵. Na razie więc mapa literatury polskiej musi być taka, jaka jest; pierwsze mapy geograficzne też nie były bardzo dokładne i wykonane od razu najlepszymi metodami...

21 P. Wilczek *Czy istnieje kanon literatury polskiej?*, w: *Literatura polska w świecie. Tom I. Zagadnienia recepcji i odbioru*, red. R. Cudak, Gnome, Katowice 2006.

22 J. Rybicki *Visualizing literature: artistic statistics*, w: *The art of literature, art in literature*, red. B. Kucala, I. Curyło-Klag, M. Bleinert, Wydawnictwo UJ, Kraków 2014 (w druku).

23 S. Argamon *Interpreting Burrows's Delta: geometric and probabilistic foundations*, „Literary and Linguistic Computing” 2008 nr 23(2), s. 140.

24 M. Eder *Style-markers in authorship attribution: a cross-language study of the authorial fingerprint*, „Studies in Polish Linguistics” 2011 nr 6, s. 99-114.

25 Największe nadzieje wiąże piszący te słowa z działalnością Grupy Technologii Językowych (Politechnika Wrocławska), kierowanej przez dr. Macieja Piaseckiego.

Choć nie musimy wiedzieć, dlaczego papierek lakmusowy zmienia kolor pod wpływem kwasu czy zasady, a mimo to możemy stosować go na własnej skórze (przynajmniej w reklamach telewizyjnych) – trzeba przyznać, że rzeczywiście najdotkliwszą bolączką tego typu badań stylometrycznych jest brak teorii tłumaczącej w zadowalający sposób tak silną skuteczność stylometrycznego ustalania autorstwa (a więc podobieństwa między tekstami) na podstawie częstości najczęstszych słów. W fizyce „teoretycy” wskazują „doświadczalnikiem”, że powinni szukać bozonu Higgsa; komputerowy stylometr, czyli właśnie literaturoznawca „doświadczalny”, takiego komfortu nie ma.

Na szczęście nie jest też całkowicie bezbronny względem materiału swych badań. Istnieje przecież potężna, tradycyjna wiedza literaturoznawcza, agregat interpretacji i historii literatury, wypracowana przez pokolenia badaczy literatury, która aż prosi się o zestawienie, zderzenie, pogodzenie lub zanegowanie właśnie w swoistych eksperymentach literaturoznawczych, ilustrowanych seriami stylometrycznych wykresów. Warto powtórzyć: konfrontacja takiego doświadczenia z istniejącą wiedzą o literaturze może nie wyjaśni, dlaczego najczęstsze słowa tak chętnie zdradzają autorstwo – ale jest nadzieja, że bazując na wielkich, nieobejmowalnych gołym okiem przez czytelnika czy interpretatora zbiorach tekstów, pozwoli wykryć nowe grupy tekstów, nowe rodzaje pokrewieństw, nowe interpretacje porównawcze. A to już nie byłoby takie obojętne dla „mainstreamowego” badacza literatury, na przykład (choć nie tylko) jej historyka i specjalisty od stylistyki²⁶. Co więcej: współpraca takiego badacza ze stylometrami mogłaby pokazać – właśnie w obecnych warunkach uprawiania humanistyki – „że nie masz takowych terminów, z których by się *viribus unitis* przy boskich *auxiliach* podnieść nie można”.

26 Wyzwania i nadzieje związane z taką współpracą zostały zasygnalizowane w pracy J. Rybicki i M. Heydel, *The stylistics and stylometry of collaborative translation: Woolf's Night and Day in Polish*, „Literary and Linguistic Computing” 2013 nr 28(4), s. 708-717.

Dodatek 1:

Lista 827 najczęstszych słów występujących w 90% tekstów użytych w analizie (w malejącym porządku frekwencji).

1. i	38. tego	75. jakby	112. właśnie
2. się	39. może	76. ani	113. powiedział
3. w	40. ze	77. potem	114. mam
4. nie	41. dla	78. które	115. gdyby
5. na	42. gdy	79. siebie	116. niech
6. z	43. tu	80. jednak	117. im
7. to	44. ten	81. gdzie	118. która
8. do	45. nawet	82. niego	119. trzeba
9. że	46. pod	83. sam	120. wiem
10. a	47. ją	84. oczy	121. można
11. jak	48. pani	85. ty	122. były
12. o	49. ich	86. lecz	123. jako
13. ale	50. ma	87. te	124. wtedy
14. co	51. przed	88. żeby	125. mógł
15. po	52. teraz	89. niej	126. nagle
16. tak	53. on	90. nas	127. jestem
17. za	54. tam	91. aby	128. znowu
18. już	55. nic	92. ona	129. niż
19. od	56. który	93. no	130. panie
20. jest	57. przy	94. raz	131. więcej
21. go	58. była	95. tych	132. jeśli
22. jej	59. wszystko	96. są	133. ludzi
23. tylko	60. nim	97. chwili	134. wszyscy
24. mnie	61. tej	98. dobrze	135. wszystkie
25. mu	62. więc	99. zawsze	136. zaraz
26. było	63. by	100. aż	137. cię
27. ja	64. bardzo	101. ta	138. jeden
28. jego	65. kiedy	102. nich	139. zaś
29. pan	66. nad	103. kto	140. we
30. bo	67. bez	104. pana	141. sobą
31. tym	68. będzie	105. je	142. dalej
32. jeszcze	69. coś	106. albo	143. choć
33. czy	70. też	107. domu	144. której
34. był	71. u	108. ku	145. których
35. mi	72. być	109. nigdy	146. tę
36. przez	73. miał	110. nią	147. mój
37. sobie	74. ci	111. przecież	148. mówił

149. lub	190. rzeczy	231. sposób	272. wszystkim
150. nikt	191. dlaczego	232. panu	273. swoim
151. drzwi	192. czego	233. chce	274. jesteś
152. człowiek	193. rękę	234. dzieci	275. drugi
153. ręce	194. zupełnie	235. swoją	276. temu
154. między	195. znów	236. masz	277. głosem
155. nam	196. ciebie	237. wśród	278. całe
156. twarz	197. swoje	238. głowy	279. bóg
157. miała	198. dziś	239. całą	280. szybko
158. ludzie	199. zresztą	240. stał	281. jaki
159. trochę	200. ktoś	241. wiele	282. oto
160. taki	201. proszę	242. ile	283. swoich
161. dopiero	202. którą	243. was	284. stronę
162. którego	203. ziemi	244. czasem	285. później
163. głowę	204. pokoju	245. powiedziała	286. tymczasem
164. cóż	205. bardziej	246. prawda	287. razie
165. chciał	206. głos	247. wiedział	288. mówić
166. także	207. ojciec	248. stary	289. parę
167. wszystkich	208. my	249. powiedzić	290. chociaż
168. chwilę	209. długo	250. dni	291. mimo
169. czas	210. twarzy	251. czowieka	292. tutaj
170. chyba	211. sama	252. samo	293. razy
171. którym	212. moje	253. stało	294. jedno
172. czasu	213. zaczął	254. trzy	295. razu
173. lat	214. głową	255. obok	296. rzecz
174. słowa	215. jakieś	256. wiesz	297. taka
175. prawie	216. nimi	257. strony	298. zapytał
176. razem	217. musi	258. lepiej	299. gdzieś
177. dwa	218. oni	259. byli	300. widział
178. tyle	219. jakie	260. mną	301. kiedyś
179. takie	220. swego	261. dość	302. naprawdę
180. coraz	221. którzy	262. ojca	303. świecie
181. kilka	222. mówi	263. świat	304. końcu
182. cały	223. dnia	264. serce	305. innych
183. myśli	224. dlatego	265. matka	306. dwie
184. jakiś	225. moja	266. widać	307. mają
185. życie	226. każdy	267. myśl	308. cicho
186. dzień	227. mogę	268. nocy	309. musiał
187. wreszcie	228. będę	269. niby	310. wam
188. wie	229. wcale	270. mieć	311. takiego
189. życia	230. pierwszy	271. mogła	312. takim

313.	dwóch	354.	miejscu	395.	kobieta	436.	samego
314.	góry	355.	miejsca	396.	jedną	437.	chcesz
315.	nieco	356.	głowie	397.	wielki	438.	mówię
316.	ręką	357.	muszę	398.	miasta	439.	jakże
317.	chcę	358.	taką	399.	często	440.	całej
318.	kobiety	359.	jednej	400.	drugiej	441.	nikogo
319.	jednego	360.	inne	401.	daleko	442.	został
320.	swój	361.	będą	402.	widzę	443.	właściwie
321.	świata	362.	zbyt	403.	jaką	444.	mówiąc
322.	noc	363.	mamy	404.	wody	445.	poszedł
323.	jednym	364.	innego	405.	chciała	446.	patrząc
324.	skąd	365.	życiu	406.	wielkie	447.	młody
325.	swej	366.	czasie	407.	drodze	448.	szedł
326.	miejsce	367.	raczej	408.	zapewne	449.	zdawało
327.	spojrzał	368.	koło	409.	powoli	450.	dać
328.	poza	369.	jedna	410.	czuł	451.	włosy
329.	oczach	370.	byłem	411.	ziemię	452.	podczas
330.	dziecko	371.	wy	412.	dom	453.	wziął
331.	swojej	372.	koniec	413.	wiadomo	454.	pieniądze
332.	odpowiedział	373.	piersi	414.	górze	455.	da
333.	pewno	374.	matki	415.	słowo	456.	woli
334.	serca	375.	cała	416.	szczęście	457.	nasze
335.	mało	376.	ciągle	417.	wiedzieć	458.	przynajmniej
336.	roku	377.	duszy	418.	spokojnie	459.	dobry
337.	każdym	378.	zrobić	419.	trudno	460.	stąd
338.	inaczej	379.	tą	420.	wyszedł	461.	zdaje
339.	natychmiast	380.	drogę	421.	iść	462.	widok
340.	nogi	381.	wolno	422.	miłość	463.	niczego
341.	mojej	382.	zaczęła	423.	siedział	464.	słów
342.	śmierci	383.	idzie	424.	moim	465.	jakaś
343.	takich	384.	jutro	425.	stanie	466.	nasz
344.	kilku	385.	pracy	426.	choćby	467.	głośno
345.	drogi	386.	samym	427.	należy	468.	kogo
346.	pewnie	387.	mniej	428.	moją	469.	dawno
347.	dużo	388.	znaczny	429.	stała	470.	słońce
348.	usta	389.	oczu	430.	myślał	471.	okna
349.	mówiła	390.	miałem	431.	tobie	472.	rzucił
350.	pół	391.	wobec	432.	nami	473.	zwykle
351.	jakoś	392.	śmierć	433.	robić	474.	podniósł
352.	boże	393.	dał	434.	jakąś	475.	one
353.	mieli	394.	chodzi	435.	sprawy	476.	będziesz

477. siły	518. żadnych	559. światło	600. łatwo
478. mocno	519. początku	560. pierwsze	601. starego
479. prawo	520. wrażenie	561. stronie	602. wieczór
480. ust	521. wprost	562. ramionami	603. mieście
481. samej	522. moich	563. wszędzie	604. musiała
482. lata	523. godziny	564. żyć	605. ostatni
483. inny	524. byłoby	565. wielkim	606. pokój
484. czemu	525. wzrok	566. przede	607. jesteśmy
485. źle	526. nowe	567. miało	608. żyje
486. wielu	527. ręki	568. jakiegoś	609. drugiego
487. wielką	528. rozumiem	569. rąk	610. przyjdzie
488. tuż	529. sami	570. łyzy	611. długie
489. dobre	530. dzieje	571. okno	612. milczeniu
490. takiej	531. rano	572. pięć	613. stoi
491. widzi	532. przyszedł	573. spokój	614. prędej
492. zamiast	533. cztery	574. czarne	615. powietrze
493. krew	534. krwi	575. żaden	616. stołu
494. mąż	535. wczoraj	576. drugą	617. kim
495. końca	536. innym	577. rozmowy	618. widzieć
496. mogą	537. byle	578. drzewa	619. wielkiej
497. chwila	538. możesz	579. ciała	620. kogoś
498. stole	539. jaka	580. głębi	621. ramiona
499. powiem	540. trzech	581. panem	622. jakiejś
500. żadnego	541. ciało	582. lekko	623. prawa
501. inni	542. sprawę	583. mogli	624. powinien
502. wszystkiego	543. powodu	584. ciężko	625. żona
503. uśmiechem	544. wrócił	585. najbardziej	626. słychać
504. ręki	545. jakimś	586. pierwszej	627. kieszeni
505. wszedł	546. mały	587. ściany	628. twoje
506. myślę	547. same	588. czoło	629. ramię
507. mogło	548. tobą	589. sprawa	630. słyszał
508. wzrokiem	549. znam	590. naszego	631. słońca
509. zrobił	550. wieczorem	591. wyraźnie	632. postać
510. jakim	551. niebo	592. widząc	633. wieku
511. żadnej	552. przykład	593. drugim	634. miały
512. widzisz	553. dół	594. droga	635. sto
513. stanął	554. naszych	595. światła	636. równie
514. całym	555. dobra	596. żal	637. rok
515. robi	556. czegoś	597. twój	638. cisza
516. znalazł	557. imię	598. zwrócił	639. uwagi
517. bądź	558. nasze	599. będziemy	640. różne

641. pełne	682. rękami	723. naszym	764. krok
642. każdej	683. znaleźć	724. twarzą	765. słowem
643. białe	684. prosto	725. los	766. boku
644. oko	685. sił	726. własne	767. rękach
645. samą	686. celu	727. wysoko	768. nadzieję
646. znał	687. drugie	728. daleka	769. czekać
647. daj	688. pamięci	729. nasza	770. godzinę
648. pytanie	689. małe	730. stół	771. porządku
649. każdego	690. środka	731. własnej	772. komu
650. uwagę	691. udało	732. najlepiej	773. druga
651. znak	692. okiem	733. odpowiedzi	774. złe
652. zęby	693. przyszła	734. zostać	775. bok
653. wielkiego	694. dwadzieścia	735. ustach	776. chwilą
654. dłoni	695. powrotem	736. powie	777. własną
655. zna	696. nieba	737. którymi	778. myślą
656. jemu	697. obraz	738. brak	779. podobne
657. stały	698. całego	739. każde	780. nadziei
658. roboty	699. oka	740. wygląda	781. daje
659. powiedz	700. została	741. ludźmi	782. nogami
660. przyszło	701. mógłby	742. rana	783. złego
661. wielka	702. głowa	743. bliżej	784. pamięć
662. prawdę	703. powietrzu	744. pomocy	785. wiedzą
663. nowy	704. sen	745. długi	786. wielkich
664. dziesięć	705. pierwszym	746. wziąć	787. pierwszej
665. pewnego	706. znać	747. nóg	788. jakiej
666. stara	707. ostatnie	748. palce	789. późno
667. część	708. wyszła	749. któremu	790. zostało
668. dzięki	709. dawna	750. możemy	791. leży
669. otworzył	710. rady	751. pierwszego	792. poznać
670. człowiekiem	711. szkoda	752. robił	793. drzwiami
671. głęboko	712. nikomu	753. patrzeć	794. stron
672. czekał	713. dłużej	754. pewnym	795. kolei
673. kroków	714. drogą	755. niewiele	796. wodą
674. zatrzymał	715. czeka	756. dziwne	797. ruchu
675. mówią	716. drzwiach	757. godzin	798. trzeci
676. tymi	717. szukać	758. każda	799. położył
677. ruch	718. krótko	759. piękne	800. zostanie
678. stać	719. leżał	760. siedzi	801. tacy
679. obu	720. nowego	761. trzymał	802. snu
680. myśleć	721. uśmiech	762. szeroko	803. stali
681. byłby	722. spać	763. stare	804. każdą

805. wyjść	811. dobrego	817. własnym	823. znajdzie
806. jakich	812. głosu	818. niedawno	824. wysoki
807. muszą	813. ostatniej	819. żadna	825. otwarte
808. odpowiedź	814. innej	820. żadne	826. nową
809. mogły	815. gorzej	821. ostatnich	827. ostatnim
810. nigdzie	816. nogach	822. nowych	

Dodatek 2:

Lista 107 najczęstszych słów występujących we wszystkich tekstach użytych w analizie (w malejącym porządku frekwencji).

1. i	28. był	55. coś	82. ludzi
2. się	29. przez	56. też	83. wszyscy
3. w	30. sobie	57. u	84. wszystkie
4. nie	31. tego	58. miał	85. zaraz
5. na	32. może	59. ani	86. we
6. z	33. ze	60. potem	87. sobą
7. to	34. dla	61. siebie	88. tę
8. do	35. tu	62. gdzie	89. nikt
9. że	36. ten	63. niego	90. rękę
10. a	37. nawet	64. sam	91. taki
11. jak	38. pod	65. oczy	92. dopiero
12. o	39. ich	66. te	93. chciał
13. ale	40. ma	67. ona	94. czas
14. co	41. przed	68. raz	95. czasu
15. po	42. teraz	69. tych	96. lat
16. tak	43. on	70. są	97. razem
17. za	44. tam	71. dobrze	98. dwa
18. już	45. nic	72. zawsze	99. takie
19. od	46. który	73. aż	100. coraz
20. jest	47. przy	74. ta	101. cały
21. go	48. była	75. nich	102. dzień
22. tylko	49. wszystko	76. kto	103. długo
23. było	50. nim	77. nigdy	104. musi
24. jego	51. kiedy	78. właśnie	105. dni
25. bo	52. nad	79. która	106. trzy
26. tym	53. bez	80. trzeba	107. drugi
27. czy	54. będzie	81. można	

Abstract

Jan Rybicki

JAGIELLONIAN UNIVERSITY (KRAKÓW)

First glimpse at a stylometric map of Polish literature

The author presents the analysis of circa 500 Polish literary texts, based on the analysis of frequency of usage of particular words and the visualisation of the outcome through network analysis. The result is discussed as a "map" of stylometric relationships between particular texts. The presence of two major "signals" was detected: auctorial and chronological. The most interesting exceptions were discussed in detail.