

MACIEJ OGRODNICZUK*

INSTYTUT PODSTAW INFORMATYKI POLSKIEJ AKADEMII NAUK, WARSZAWA

Lingwistyka komputerowa dla języka polskiego: dziś i jutro

Słowa kluczowe: lingwistyka komputerowa, lingwistyka informatyczna, przetwarzanie języka naturalnego, technologie językowe, narzędzia i zasoby językowe dla polszczyzny.

1. Wprowadzenie

W 2006 roku na łamach „LingVariów” ukazał się entuzjastyczny tekst Marka Świdzińskiego na temat zdobyczy polskiej lingwistyki korpusowej, w którym autor przedstawiał korzyści, jakie językoznawstwu przyniosła, także w Polsce, rewolucja informatyczna. Sposób, w jaki środowisko językoznawcze przyjęło te zmiany, dość dobrze opisują z kolei w publicystycznych artykułach Piotr Żmigrodzki (2015a) i Adam Przepiórkowski (2015). W niniejszym tekście chciałbym spojrzeć na sytuację lingwistyki informatycznej w Polsce oraz na dalsze kierunki prac nad komputerowym przetwarzaniem polszczyzny z perspektywy o dziesięć lat późniejszej oraz z punktu widzenia przedstawiciela tej dyscypliny, często postrzeganej wyłącznie jako usługowa dla środowiska językoznawczego.

Zacznijmy zatem od kilku słów na temat przedmiotu naszej pracy i trudności, przed jakimi stoimy. Lingwistyka komputerowa (lingwistyka informatyczna, inżynieria lingwistyczna) to dyscyplina z pogranicza obu dziedzin – informatyki i językoznawstwa, której celem jest badanie języka naturalnego z punktu widzenia potrzeb jego przetwarzania i modelowania metodami komputerowymi oraz tworzenie elektronicznych zasobów i narzędzi przetwarzania języka na potrzeby usługowe. Powyższa definicja, łącząca istniejące dotąd definicje (w szczególności dwie: Janusza S. Bienia i Bonnie Webber¹), zwraca uwagę na dwa aspekty problemu: jest to jednocześnie dziedzina teoretyczna i stosowana. O pierwszym zdają się zapominać językoznawcy, o drugim – sami informatycy. Nie znam szczegółów tego rozstrzygnięcia, ale o stosunku polskiego środowiska językoznawczego do lingwistyki komputerowej jako dziedziny badań w obszarze nauk humanistycznych najlepiej świadczy niedawne usunięcie jej z panelu HS 2.6 w konkursach Narodowego Centrum Nauki, chociaż w projektach prowadzonych w moim rodzimym Zespole Inżynierii Lingwistycznej Instytutu Podstaw Informatyki PAN często więcej jest pracy ściśle językoznawczej (wykonywanej przez zaprzyjaźnionych lingwistów) niż informatycznej. Taka jest zresztą specyfika większości polskich zespołów zajmujących się przetwarzaniem języka: naszą domeną jest właśnie lingwistyka informatyczna, a nie

* maciej.ogrodniczuk@ipipan.waw.pl

¹ Bień 1988; Webber 2001, zob. też Piasecki 2008.

informatyka lingwistyczna, z konsekwencjami w postaci rezygnacji z wyścigu w rozwoju technologii anglojęzycznej na rzecz świadomej lingwistycznie technologii dla polszczyzny; stąd też nasza częsta obecność w najważniejszych polskich pismach językoznawczych.

Po stronie informatycznej grzechem głównym jest właśnie trudność z przyznawaniem się do częściowo usługowego wymiaru lingwistyki komputerowej i kłopoty komunikacyjne obu światów. Być może wynikają one z faktu postrzegania informatyków jako wykonawców, a nie partnerów w pracy badawczej. Usługowość rozumiem zatem nie jako gotowość nauczania lingwistów języka zapytań do wyszukiwarki czy pomoc w stworzeniu własnego korpusu, ale jako rzeczywistą współpracę w badaniach.

Mimo wspomnianych trudności mój punkt widzenia na stan lingwistyki informatycznej w Polsce jest optymistyczny: jesteśmy świadkami bardzo intensywnego rozwoju komputerowych zasobów i narzędzi dla języka polskiego², polskie środowisko lingwistyczno-informatyczne bierze udział w licznych grantach polskich i europejskich, uczestniczy w dużych infrastrukturach badawczych CLARIN i DARIAH, łączących nauki techniczne i humanistyczne, a sami użytkownicy humaniści coraz częściej sięgają po metody komputerowe naszego autorstwa. Trwa owocna współpraca między polskimi ośrodkami badawczymi, a komputerowy opis polszczyzny obejmuje coraz głębsze poziomy analizy. W dalszej części tekstu postaram się pokazać, w którym miejscu jesteśmy, i zaproponować kilka kierunków badań na najbliższe lata, jeśli nie dla całego środowiska, to przynajmniej dla naszego warszawskiego zespołu.

2. Stan technologii językowej dla polszczyzny okiem lingwisty komputerowego

2.1. Raport językowy

W roku 2011 warszawski Zespół Inżynierii Lingwistycznej Instytutu Podstaw Informatyki PAN brał udział w międzynarodowym projekcie CESAR, którego celem (wraz z innymi inicjatywami podobnego typu prowadzonymi pod wspólnym szyldem META-NET-u) było zebranie, scalenie oraz zrównoleglenie elektronicznych zasobów i narzędzi językowych dla trzydziestu języków europejskich, a następnie udostępnienie ich we wspólnym repozytorium. Miało to stworzyć nowe możliwości badawcze oraz technologiczne w dziedzinie przetwarzania języków naszego kontynentu. Projekt zakończył się, moim zdaniem, sporym sukcesem: doprowadził do powstania wielu ważnych zasobów, jak np. *Słownik gramatyczny języka polskiego* (Saloni i in. 2007/2015), doprowadził do unormowania kwestii licencyjnych oraz odświeżył stronę internetową CLIP (Ogrodniczuk, Przepiórkowski 2013)³ z rejestrem dostępnych publicznie narzędzi dla polszczyzny oraz polskich zespołów badawczych zajmujących się przetwarzaniem języka i prowadzonych przez nie projektów.

Innym ważnym wynikiem prac był tzw. raport językowy (Miłkowski 2012) upowszechniający wiedzę nt. technologii językowych dla polszczyzny i ich potencjalnych zastosowań. Jego ważnym zadaniem było też przedstawienie obrazu bieżącej sytuacji w dziedzinie przetwarzania

2 Por. <http://clip.ipipan.waw.pl/LRT> – strona zawierająca listę dostępnych w sieci narzędzi i zasobów dla polszczyzny.

3 Computational Linguistics in Poland ('Lingwistyka komputerowa w Polsce') – por. <http://clip.ipipan.waw.pl> (strona w języku angielskim).

języka w Polsce. Na podstawie oceny eksperckiej w skali od 0 (ocena bardzo niska) do 6 (bardzo wysoka) przedstawiał dane nt. dostępności, jakości i elastyczności technologii oraz zasobów w orientacyjnie zarysowanych dziedzinach. Wyniki tej oceny przytaczam w tabeli 1.

Tabela 1. Stan dostępnych technologii językowych dla języka polskiego
(przedruk, s. 31 raportu)

	Liczba	Dostępność	Jakość	Zakres	Dojrzałość	Trwałość	Elastyczność
Technologie językowe (narzędzia, technologie, aplikacje)							
Rozpoznawanie mowy	1	2	3	4	3	2	4
Synteza mowy	4	3	6	5	4	4	3
Analiza gramatyczna	4	4,5	4,5	4,5	4	4	3
Analiza semantyczna	1	1	3	1	1	2	2
Generowanie tekstu	1	1	1	1	1	1	2
Tłumaczenie maszynowe	3	4	3	3	3	4	3
Zasoby językowe (zasoby, dane, bazy wiedzy)							
Korpusy tekstów	3	2	4	4	5	5	3
Korpusy równoległe	3	1	4	4	5	5	5
Korpusy języka mówionego	1	0	3	3	2	2	2
Zasoby leksykalne	3	3	4	4	4	4	3
Gramatyki	3	2	4	4	3	2	2

Wykaz ten dobrze obrazuje stan prac sprzed pięciu lat: już wtedy doskonale radziliśmy sobie z syntezą mowy⁴, nieźle z analizą morfoskładniową (tu oznaczoną jako gramatyczna), jako dobrą oceniliśmy dostępność jednojęzycznych korpusów tekstów⁵, jako gorszą – korpusów równoległych; z dostępnością korpusów mówionych było gorzej niż źle. Raport zawiera więcej podobnych wniosków, głównie pesymistycznych: wskazuje na brak standaryzacji zasobów oraz na konieczność intensyfikacji prac nad głęboką analizą tekstu, zasobami ontologicznymi i elektronicznymi słownikami walencyjnymi.

4 Był to moment triumfu polskiego systemu Ivona w wielu międzynarodowych konkursach syntezy mowy, po którym firma będąca jego autorem została przejęta przez globalną firmę Amazon.

5 Ze względu na moment zakończenia prac nad Narodowym Korpusem Języka Polskiego (Przepiórkowski i in. (red.) 2012).

Ważnych informacji dostarcza także porównanie stanu technologii dla języka polskiego i innych języków europejskich; tabela 2 przedstawia jego wyniki dla zadania analizy tekstu.

Tabela 2. Stan technologii językowych dostępnych dla 30 języków europejskich dla zadania analizy tekstu (przedruk, s. 33 raportu)

Doskonała jakość	Bardzo dobra jakość	Dobra jakość	Średnia jakość	Słaba/zerowa jakość
	angielski	francuski hiszpański niderlandzki niemiecki włoski	baskijski bułgarski czeski duński fiński galisyjski grecki kataloński norweski polski portugalski rumuński słowacki słoweński szwedzki węgierski	chorwacki estoński irlandzki islandzki litewski łotewski maltański serbski

2.2. Sonda środowiskowa

Z perspektywy czasu widać, że zarówno usytuowanie polszczyzny w rodzinie języków europejskich, jak i wnioski na temat stanu rozwoju technologii były właściwe; kierunki najintensywniejszych prac prowadzonych w ciągu ostatnich lat pokrywają się z powyższymi postulatami. Powstaje duży semantyczny słownik walencyjny Walenty (Przepiórkowski i in. 2014), z którego zaczynają korzystać parsery, czyli narzędzia do głębokiej analizy tekstu (Patejuk, Przepiórkowski 2012; Jaworski, Przepiórkowski 2014; Woliński 2015), od wielu lat dynamicznie rozwija się Słowosieć (Piasecki i in. 2014), trwają prace nad poprawą jakości i dostępności korpusów równoległych i mówionych (Pęzik 2015, 2016).

Wiarygodną ocenę stanu obecnego najlepiej jednak przeprowadzić metodą konsultacji ze środowiskiem, postanowiłem zatem powtórzyć ocenę metodą prostej sondy. Na etapie tworzenia pytań na podstawie tabeli z raportu językowego okazało się, że określenia rodzajów narzędzi i zasobów są zbyt ogólne (stąd zapewne dobra ocena jakości „analizy gramatycznej”, który to termin byłby dziś raczej kojarzony z analizą składniową, a nie morfoskładniową), postanowiłem zatem przejąć listę bardziej szczegółowych kategorii z istniejących badań zagranicznych (np. Strik i in. 2002) i dostosować ją do specyfiki polskiego środowiska. Sonda w części dotyczącej obecnego stanu prac składała się z dwóch pytań:

1. Jak ocenia Pan/Pani dostępność narzędzi i zasobów dla polszczyzny?
Proszę o podanie Państwa opinii nt. stopnia możliwości wykorzystania istniejących narzędzi i zasobów dla języka polskiego.
 - a. Niedostępne
 - b. Dostępne w ograniczonym stopniu (np. za opłatą, w ograniczonym zakresie)
 - c. Dostępne w wystarczającym zakresie
 - d. Nie wiem
2. Jak ocenia Pan/Pani jakość technologii językowych dla polszczyzny?
Proszę o podanie Państwa opinii nt. ogólnej jakości dostępnej obecnie technologii (darmowej lub komercyjnej) dla języka polskiego, biorąc pod uwagę liczbę i wielkość zasobów, różnorodność i zakres dziedzinowy itp.
 - a. Brak lub słaba jakość
 - b. Średnia jakość
 - c. Dobra jakość
 - d. Doskonała jakość
 - e. Nie wiem

Lista ocenianych narzędzi i zasobów liczyła 29 pozycji:

- | | |
|--|------------------------------|
| 1. Rozpoznawanie mowy | 15. Analiza referencji |
| 2. Synteza mowy | 16. Analiza pragmatyczna |
| 3. Segmentacja tekstu | 17. Analiza dyskursu |
| 4. Analiza morfologiczna | 18. Korekta pisowni |
| 5. Synteza morfologiczna | 19. Generowanie tekstu |
| 6. Lematyzacja | 20. Tłumaczenie maszynowe |
| 7. Tagowanie (ujednoznacznianie morfoskładniowe) | 21. Streszczanie tekstu |
| 8. Parsowanie zależnościowe | 22. Ekstrakcja informacji |
| 9. Parsowanie powierzchniowe | 23. Korpusy języka pisanego |
| 10. Parsowanie głębokie | 24. Korpusy języka mówionego |
| 11. Rozpoznawanie nazw własnych | 25. Korpusy równoległe |
| 12. Ujednoznacznianie sensu słów | 26. Słowniki elektroniczne |
| 13. Analiza sentymentu/opinii | 27. Gramatyki formalne |
| 14. Głęboka analiza semantyczna | 28. Tezaurusy |
| | 29. Wordnety |

Ankieta została rozesłana do dwustu osób związanych ze środowiskiem lingwistyczno-informatycznym w Polsce, zebrano 23 odpowiedzi z 14 jednostek (uczelnia, instytutów badawczych i firm). Tabela 3 przedstawia wyniki sondy ograniczone do kategorii użytych w raporcie z 2011 roku i przeskalowane w sposób umożliwiający porównanie wyników obecnych z tymi sprzed pięciu lat.

Tabela 3. Stan dostępnych technologii językowych dla języka polskiego na podstawie sondy przeprowadzonej wśród przedstawicieli środowiska lingwistyczno-informatycznego

	Liczba	Dostępność	Jakość	Zakres	Dojrzałość	Trwałość	Elastyczność
Technologie językowe (narzędzia, technologie, aplikacje)							
Rozpoznawanie mowy	1	3	3	4	3	2	4
Synteza mowy	4	4	5	5	4	4	3
Analiza gramatyczna	4	4,5	5	4,5	4	4	3
Analiza semantyczna	1	1	2	1	1	2	2
Generowanie tekstu	1	3	2	1	1	1	2
Tłumaczenie maszynowe	3	3	2	3	3	4	3
Zasoby językowe (zasoby, dane, bazy wiedzy)							
Korpusy tekstów	3	4	4	4	5	5	3
Korpusy równoległe	3	3	3	4	5	5	5
Korpusy języka mówionego	1	3	3	3	2	2	2
Zasoby leksykalne	3	3	4	4	4	4	3
Gramatyki	3	3	3	4	3	2	2

Oznaczenia jaśniejszym kolorem odpowiadają poprawie oceny, ciemniejszym – jej pogorszeniu. Wynik negatywny należy raczej interpretować jako zmianę w sposobie postrzegania danej dziedziny (w porównaniu z technologiami podobnego typu albo dostępnymi dla innych języków) niż faktyczny spadek jakości danego narzędzia czy zasobu. Ze względu na obszerniejszą listę kategorii warto przyjrzeć się także wynikom dla technologii, które nie były uwzględnione na liście z 2011 r.: wśród nich za najbardziej dostępne (o dostępności w wystarczającym zakresie) uznano oprócz wskazanych w tabeli także mechanizmy segmentacji tekstu, tagowania i wordnety. Za niedostępne uznano z kolei narzędzia do głębokiej analizy semantycznej i pragmatycznej oraz analizy dyskursu. Najlepszą jakość przypisano natomiast technologii syntezy mowy, segmentacji tekstu, analizy i syntezy morfologicznej oraz korekty pisowni (doskonała jakość). Spośród technologii dostępnych najniżej oceniono jakość technologii do streszczania tekstu, a niewiele lepiej do parsowania głębokiego, ujednoznaczniania sensu słów, analizy sentymentu/opinii, generowania tekstu i, co już zostało pokazane w tabeli, tłumaczenia maszynowego. W mojej opinii oceny te odpowiadają obecnemu stanowi rozwoju wymienionych technologii.

3. Zadania dla lingwistyki komputerowej w Polsce

Postulaty, które teraz przedstawię, należy odczytać jako w dużej mierze subiektywną prezentację kierunków rozwoju lingwistyki komputerowej w Polsce. Mam nadzieję, że moje propozycje będą atrakcyjne dla obu środowisk naukowych – humanistycznego i informatycznego, pomogą odpowiedzieć na ich potrzeby i pozwolą im się jeszcze bardziej zbliżyć. W pewnym stopniu poruszają one problemy zasygnalizowane w odpowiedziach ankietowych, ale przedstawiam je jednak z perspektywy planów (i aspiracji) warszawskiego zespołu.

3.1. Korpus narodowy

Pierwsza wersja Narodowego Korpusu Języka Polskiego była organizacyjnym i badawczym przełomem, który jeszcze długo będzie wpływał na środowisko humanistyczne i lingwistyczno-informatyczne w Polsce. Prace prowadzone w ramach projektu rozwojowego MNiSW doprowadziły do powstania nie tylko samego zasobu korpusowego, ale także wielu narzędzi informatycznych do przetwarzania języka polskiego. Warto jednak zwrócić uwagę, że tryb projektowy doprowadził do powstania zbioru, który od zakończenia prac trwa w postaci zamrożonej – ostatnie teksty pochodzą z roku 2011, a automatyczne analizy lingwistyczne zostały wytworzone narzędziami o kilka generacji wcześniejszymi niż używane obecnie. Zapewnienie aktualności i stałego monitoringu NKJP, który stanowi przecież jeden z podstawowych zasobów referencyjnych dla rzeszy użytkowników, wydaje się koniecznością. A potrzeb jest dużo więcej: równie ważne wydaje się wielokierunkowe rozszerzanie bazy materiałowej, opracowanie metod reprezentacji i wizualizacji danych lingwistycznych, rozwój narzędzi analitycznych i eksploracyjnych. Próbę dokładniejszego zdefiniowania potrzeb w tym zakresie podjęliśmy już w ramach grupy roboczej „Korpusowa infrastruktura badawcza dla polszczyzny” powołanej w ramach konsorcjum DARIAH-PL, warto zatem powtórzyć kierunki, w których widzimy rozwój korpusu narodowego:

- 1) Rozszerzenie bazy materiałowej Narodowego Korpusu Języka Polskiego (o podkorpusy diachroniczne, podkorpusy równoległe, specjalistyczne, w tym gwarowy, dane mówione, podkorpus synchroniczny tworzony na bazie bibliotek cyfrowych i Internetu).
- 2) Opracowanie korpusowych standardów znakowania dla języka polskiego:
 - a) opracowanie formatu znakowania wszystkich tekstów (także historycznych i gwarowych) rozszerzonymi metadanymi i znacznikami morfosyntaktycznymi,
 - b) opracowanie formatu znakowania tekstów współczesnych znacznikami składniowymi, semantycznymi i dyskursywnymi.
- 3) Opracowanie metod reprezentacji danych multimedialnych:
 - a) opracowanie metod zrównoleglania tekstów ze źródłami (dane mówione, wideo, skany),
 - b) opracowanie metod wyszukiwania w plikach źródłowych.
- 4) Opracowanie interfejsu użytkownika korpusu uwzględniającego proste formułowanie zapytań oraz atrakcyjną i czytelną wizualizację wyników:
 - a) opracowanie intuicyjnych nakładek graficznych i interakcyjnych,

- b) opracowanie metod wizualizacji zaawansowanych struktur lingwistycznych (drzewa składniowe, koreferencje itp.).
- 5) Rozwój narzędzi do samodzielnego tworzenia korpusów.
- 6) Rozwój narzędzi do eksploracji i analizy danych korpusowych.
- 7) Monitoring źródeł internetowych na potrzeby korpusowe i leksykograficzne.

O istotności prac w każdym z tych kierunków świadczy fakt uruchomienia wielu projektów finansowanych ze środków Narodowego Programu Rozwoju Humanistyki oraz Narodowego Centrum Nauki, w których ramach powstają obecnie korpusy diachroniczne: XVI, XVII–XVIII i XIX wieku, korpus polskich tekstów prasowych (aktualnie z lat 1945–1954) czy korpus gwarowy. Jednocześnie powstają również reprezentacje zaawansowanych struktur lingwistycznych w rodzaju relacji referencyjnych czy metafor wraz z narzędziami analitycznymi. Projekty te, nawet jeśli wykorzystują dane korpusu narodowego, są jednak prowadzone niezależnie, bez intencji jego wzbogacenia, co wydaje mi się dużą stratą. A przecież podobnie do *Wielkiego słownika języka polskiego PAN* (Żmigrodzki 2015b) wielki korpus polszczyzny jest podstawowym zasobem językowym, który powinien istnieć w sposób stały i stanowić rodzaj „instytucji narodowej”. Prace nad nim nie mogą się ograniczać do jednego projektu – wymagają nadzoru środowiska humanistycznego, doradztwa naukowego na najwyższym poziomie i, co też trzeba powiedzieć, stałego finansowania. Szczególnie ten ostatni postulat wydaje się trudny do spełnienia, jak może się wydawać po lekturze długiej listy potrzeb, jednak najpilniejsza kwestia rozpoczęcia prac nad utrzymaniem korpusu nie wydaje się aż tak mało realistyczna i w moim przekonaniu wystarczyłoby do niej jednoczesne osiągnięcie trzech celów:

- 1) zapewnienie ciągłości życia NKJP, co stanowiłoby spełnienie części postulatów autorów monografii korpusu (Przepiórkowski i in. (red.) 2012) jego ustawicznego monitorowania oraz uzupełniania i ulepszania, w fazie wstępnej ograniczonego do dostosowania opisów lingwistycznych do najnowszych narzędzi analitycznych; tego rodzaju działania wymagałyby pracy zaledwie kilku osób,
- 2) zapewnienie nadzoru nad długofalowymi kierunkami rozwoju NKJP przez radę programową złożoną z wybitnych językoznawców,
- 3) opracowanie metody konstruowania grantów w sposób nagradzający włączanie ich wyników do zasobu korpusu narodowego i zapewniający finansowe wsparcie tego procesu.

Powyższa formuła jest oczywiście luźną propozycją, której celem jest rozpoczęcie dyskusji nad stworzeniem zasad funkcjonowania i rozwoju korpusu narodowego. Na jej bazie chciałbym niniejszym rozpocząć lobbing na rzecz tego ważnego zasobu i poprosić czytelników o poparcie tej inicjatywy.

3.2. Nowe tematy – i mozołnie naprzód

Drugim istotnym kierunkiem prac jest, moim zdaniem, formalny opis polszczyzny w zakresie od dawna podejmowanym dla języków zachodnich, a u nas wciąż traktowanym jako pieśń

przyszłości – chodzi przede wszystkim o kwestie związane z analizą dyskursu rozumianą jako generowanie struktur wypowiedzi wykraczających poza poziom zdania. Temat ten wpisuje się w prowadzone obecnie prace nad komputerowym opisem takich zjawisk, jak referencja, metafora oraz argumentacja w języku.

Oczywiście mój subiektywny punkt widzenia może nie odzwierciedlać istotności tematów dla pozostałej części środowiska lingwistyczno-informatycznego, sięgnę zatem raz jeszcze do ankiety, w której zadałem też trzecie pytanie, prosząc o podanie informacji o planach rozwoju danej technologii w zespole respondenta, ze wspomnianą wyżej szczegółową listą 29 technologii, i otrzymałem następujące odpowiedzi:

1. Jakość technologii jest wystarczająca, nie ma potrzeby prowadzenia dalszych prac.
2. Obecnie prowadzimy prace w tej dziedzinie.
3. Zamierzamy rozwijać narzędzia/zasoby z tej dziedziny w przyszłości.
4. Temat jest ważny, ale nie planujemy prowadzić prac w tej dziedzinie.
5. Nie uważamy tego tematu za istotny.
6. Trudno powiedzieć.

Intencją pytania było sprawdzenie, które problemy uważamy jako środowisko za rozwiązane, które są dla nas ważne (bo prowadzimy już prace w danej dziedzinie, zamierzamy je prowadzić lub uważamy temat za istotny), a które nie są warte badawczej uwagi. Ciekawym wynikiem sondy okazało się stwierdzenie, że praktycznie nie ma dziedzin, które zostałyby uznane za nieważne; za kwestie rozwiązane najczęściej ankietowanych uznało korektę pisowni (5 odpowiedzi z 23 udzielonych), dostępność korpusów języka pisanego (4), narzędzi do segmentacji tekstu, lematyzacji, analizy morfologicznej oraz takich zasobów jak korpusy równoległe i wordnety (3). Za tematy ważne uznano tłumaczenie maszynowe (19), rozpoznawanie mowy (17), syntezę mowy, tagowanie, analizę sentymentu/opinii i ekstrakcję informacji (16), a także streszczanie tekstu, rozwój korpusów języka pisanego, słowników elektronicznych i gramatyk formalnych (15). Wskazania te nie zmieniły się znacząco przy ograniczeniu odpowiedzi do technologii „przyszłościowych”, czyli odjęciu odpowiedzi uwzględniających prowadzone prace. W kontekście poprzedniego rozdziału komentarza wymaga obecność na obu listach korpusów języka pisanego; ocena ich dobrej dostępności może wynikać z zaufania do NKJP z jednoczesnym wskazaniem na potrzebę rozwoju korpusów dla polszczyzny – w tym również korpusu narodowego.

Oczywiście przy okazji kreślenia planów na przyszłość nie można pominąć wątku stałego poprawiania jakości dostępnych narzędzi i zasobów. Niedostigły w niektórych dziedzinach poziom 90-procentowej dokładności jest dla wielu celów niewystarczający, gdyż nawet jedno błędne wskazanie na każde dziesięć to za dużo, zwłaszcza gdy wynik jednego narzędzia (np. ujednoznacznienie części mowy w przetwarzanym tekście) jest źródłem danych dla kolejnego (np. analizatora składniowego wykorzystującego informację o częściach mowy). Poprawa narzędzi działających na różnych poziomach wiedzy lingwistycznej – analizatorów składniowych i semantycznych, mechanizmów do wykrywania nazw własnych, koreferencji – choć mało spektakularna, jest niezwykle ważna z perspektywy wykorzystania ich wyników w praktyce.

Oprócz dziedziny *stricte* badawczej nie można też zapomnieć o nie mniej ważnej części „usługowej”. Szlak jest już przetarty; środowisko informatyczne bierze udział w dwóch ważnych inicjatywach dotyczących humanistyki cyfrowej – infrastrukturach badawczych CLARIN-PL i DARIAH-PL. CLARIN-PL jest inicjatywą środowiska informatycznego wspierającą projekty badawcze w naukach humanistycznych i społecznych poprzez umożliwienie korzystania z dostępnych zasobów i narzędzi w zadaniach związanych z przetwarzaniem języka. Inicjatorem DARIAH-PL jest natomiast środowisko humanistyczne działające na rzecz tworzenia, rozwoju i utrzymania infrastruktury humanistyki cyfrowej w Polsce. Co ważne, instytucje reprezentujące główne polskie centra badawcze zajmujące się lingwistyką informatyczną⁶ są członkami obu inicjatyw, włączając się w badania humanistyczne w zakresie własnych zainteresowań i dostępnych środków.

Mimo wielu kłopotów właściwych polskiej nauce polskie środowisko lingwistyki informatycznej działa bardzo prężnie, a trzy główne tematy: składnia, semantyka i dyskurs, wydają się wyznaczać kierunki prac na najbliższe lata. Muszę jednak po raz kolejny zastrzec, że powyższy tekst przedstawia przede wszystkim moje osobiste poglądy, a w drugiej kolejności (na podstawie odpowiedzi z sondy) punkt widzenia informatyków zajmujących się przetwarzaniem języka, może zatem nie odzwierciedlać stanowiska użytkowników, których priorytety mogą być zgoła inne. W związku z tym zachęcam humanistów do współpracy, dzielenia się swoimi problemami, a także wynikami własnych prac. To najlepsza metoda na poprawianie jakości narzędzi dla polszczyzny, z których oba środowiska będą mogły korzystać.

Bibliografia

- Bień J.S. 1988: *Komputery a język naturalny*, „Biuletyn Polskiego Towarzystwa Informatycznego” VI (2–3), s. 6–7.
- Jaworski W., Przepiórkowski A. 2014: *Syntactic approximation of semantic roles*, [w:] A. Przepiórkowski, M. Ogrodniczuk (red.), *Advances in Natural Language Processing*, Lecture Notes in Artificial Intelligence, vol. 8686, Springer, Heidelberg, s. 193–201.
- Milkowski M. 2012: *Język polski w erze cyfrowej*, Springer, Berlin–Heidelberg (online: <http://doi.org/10.1007/978-3-642-30811-6>).
- Ogrodniczuk M., Przepiórkowski A. 2013: *CLIP – portal internetowy łączący projekty, narzędzia, zespoły związane z komputerowym przetwarzaniem języka polskiego*, „Język Polski” XCIII, s. 242.
- Patejuk A., Przepiórkowski A. 2012: *Towards an LFG parser for Polish. An exercise in parasitic grammar development*, [w:] N. Calzolari, K. Choukri, T. Declerck, M.U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk, S. Piperidis (red.), *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC 2012)*, ELRA, Stambuł, s. 3849–3852.
- Pęzik P. 2015: *Spokes – a search and exploration service for conversational corpus data*, [w:] *Selected Papers from CLARIN 2014*, Linköping University Electronic Press, Linköpings universitet, Linköping, s. 99–109.
- Pęzik P. 2016: *Exploring phraseological equivalence with paralela*, [w:] E. Gruszczyńska, A. Leńko-Szymańska (red.), *Polskojęzyczne korpusy równoległe. Polish-language Parallel Corpora*, Instytut Lingwistyki Stosowanej Uniwersytetu Warszawskiego, Warszawa, s. 67–81.
- Piasecki M. 2008: *Cele i zadania lingwistyki informatycznej*, [w:] P. Stalmaszczyk (red.), *Metodologie językoznawstwa. Współczesne tendencje i kontrowersje*, Lexis, Kraków.

6 Zob. <http://clip.ipipan.waw.pl/Centers>.

- Piasecki M., Maziarz M., Szpakowicz S., Rudnicka E. 2014: *PIWordNet as the Cornerstone of a Toolkit of Lexico-semantic Resources*, [w:] H. Orav, C. Fellbaum, P. Vossen (red.), *Proceedings of the Seventh International Global Wordnet Conference (GWC 2014)*, University of Tartu Press, Tartu, s. 304–312.
- Przepiórkowski A. 2015: *Inżynieria lingwistyczna a obecna sytuacja językoznawstwa polskiego*, „LingVaria” X, nr 2, s. 135–145.
- Przepiórkowski A., Bańko M., Górski R.L., Lewandowska-Tomaszczyk B. (red.) 2012: *Narodowy Korpus Języka Polskiego*, Wydawnictwo Naukowe PWN, Warszawa.
- Przepiórkowski A., Skwarski F., Hajnicz E., Patejuk A., Świdziński M., Woliński M. 2014: *Modelowanie własności składniowych czasowników w nowym słowniku walencyjnym języka polskiego*, „Polonica” XXXIII, s. 159–178.
- Saloni Z., Woliński M., Wołosz R., Gruszczyński W., Skowrońska D. 2007/2015: *Słownik gramatyczny języka polskiego*, wyd. 1: Warszawa 2007 (dokument elektroniczny), wyd. 2: Warszawa 2012 (dokument elektroniczny), wyd. 3: 2015 (online: <http://sgjp.pl>).
- Strik H., Daelemans W., Binnenpoorte D., Sturm J., de Vriend F., Cucchiarini C. 2002: *Dutch HLT resources. From BLARK to priority lists*, [w:] J. Hansen, B. Pellom (red.), *Proceedings of the Seventh International Conference on Spoken Language Processing (ICSLP 2002)*, Denver, USA, s. 1549–1552.
- Świdziński M. 2006: *Lingwistyka korpusowa w Polsce – źródła, stan, perspektywy*, „LingVaria” I, nr 1, s. 23–32.
- Webber B.L. 2001: *Computational perspectives on discourse and dialogue*, [w:] D. Schirin, D. Tannen, H. Hamilton (red.), *The Handbook of Discourse Analysis*, Blackwell Publishers Ltd.
- Woliński M. 2015: *Deploying the new valency dictionary Walenty in a DCG parser of Polish*, [w:] M. Dickinson, E. Hinrichs, A. Patejuk, A. Przepiórkowski (red.), *Proceedings of the Fourteenth International Workshop on Treebanks and Linguistic Theories (TLT14)*, Instytut Podstaw Informatyki PAN, Warszawa, s. 221–229.
- Żmigrodzki P. 2015a: *Wielki słownik języka polskiego PAN a obecna sytuacja językoznawstwa polskiego*, „LingVaria” X, nr 1, s. 109–119.
- Żmigrodzki P. 2015b: *Wielki słownik języka polskiego PAN – historia, stan obecny i perspektywy rozwoju po 2018 roku*, „Biuletyn Polskiego Towarzystwa Językoznawczego” LXXI, s. 177–187.

Summary

Natural language processing for Polish: today and tomorrow

Keywords: computational linguistics, natural language processing, language technology, language resources and tools for Polish.

The article attempts at framing directions for future work on computational processing of Polish in the face of recent intensive development of electronic tools and resources and close co-operation between Polish research centres involved in computational linguistics. The author regards renewing the work on the National Corpus of Polish as the most important topic, naming it the basic resource for Polish linguistics and listing the most urgent objectives: extension of the sources and linguistic representation as well as inclusion of diachronic, dialectical and parallel corpora. With respect to language technology, the author calls for enrichment of formal description of Polish with syntactic, semantic and discourse-feature representation and constant improvement of quality of tools and resources by means of co-operation between linguists and computer scientists.