



AKADEMIA GÓRNICZO-HUTNICZA  
IM. STANISŁAWA STASZICA W KRAKOWIE  
WYDZIAŁ ELEKTROTECHNIKI, AUTOMATYKI,  
INFORMATYKI I ELEKTRONIKI  
KATEDRA INFORMATYKI

ROZSTRZYGANIE WIELOZNACZNOŚCI WYRAZU  
ZA POMOCĄ RELACYJNEGO OPISU ZNACZENIA

EWA THŁON I TOMASZ PĘDZIMAŹ  
INFORMATYKA

PRACA MAGISTERSKA  
PRACA NAPISANA POD KIERUNKIEM:  
PROF. DR HAB. WIESŁAWA LUBASZEWSKIEGO

Kraków, 2008

## SPIS TREŚCI:

ROZDZIAŁ I – WSTĘP .....	5
ROZDZIAŁ II – SŁOWNIK SEMANTYCZNY .....	8
2.1. Słownik semantyczny języka polskiego .....	8
2.2. Definicje w ramach semantyki .....	9
2.3. Wieloznaczność słów .....	11
2.3.1. Podejście składnikowe .....	11
2.3.2. Podejście prototypowe .....	13
2.3.3. Podejście relacyjne .....	14
2.4. Realizacje słownika semantycznego .....	15
2.5. Słownik semantyczny Katedry Lingwistyki Komputerowej UJ .....	17
2.5.1. Relacje syntagmatyczne w słowniku semantycznym UJ .....	17
2.5.2. Próbkę słownika semantycznego .....	19
ROZDZIAŁ III – ROZSTRZYGANIE WIELOZNACZNOŚCI .....	21
3.1. Model .....	23
3.2. Obróbka zdania .....	25
3.3. Rozumowanie w modelu .....	28
3.4. Analiza wyników .....	29
ROZDZIAŁ IV – IMPLEMENTACJA .....	30
4.1. Generacja pliku modelu .....	30
4.2. Implementacja ujednoznaczniania .....	32
ROZDZIAŁ V – ANALIZA ZACHOWANIA ALGORYTMU .....	40
5.1. Testy autorskie .....	43
5.1.1. Zbiór testowy i otrzymane wyniki .....	43
5.1.2. Analiza .....	45
5.2. Teksty języka potocznego .....	46
5.2.1. Zbiór testowy i wyniki .....	46
5.2.2. Analiza .....	47
5.3. Teksty dzieł literackich .....	48
5.3.1. Zbiór testowy i wyniki .....	48
5.3.2. Analiza .....	52

5.4. Notatki PAP .....	53
5.4.1. Zbiór testowy i wyniki .....	53
5.4.2. Analiza .....	60
5.5. Wszystkie notatki zawierające wybrane słowo .....	61
5.5.1. Zbiór testowy i wyniki .....	61
5.5.2. Analiza .....	61
5.6. Ograniczony zbiór relacji .....	62
5.6.1. Zbiór testowy i wyniki .....	62
5.6.2. Analiza .....	63
5.7. Sztucznie wygenerowane ciągi słów .....	64
5.7.1. Zbiór testowy i wyniki .....	64
5.7.2. Analiza .....	65
ROZDZIAŁ VI – PODSUMOWANIE .....	67
BIBLIOGRAFIA .....	70

## SPIS RYSUNKÓW:

Rys.1 Trójkąt semiotyczny C.K. Ogdena i A. Richardsa a wieloznaczność .....	9
Rys.2 Rozkład zdania Chomsky'ego .....	13
Rys.3 Fragment modelu dla słowa wieloznacznego <i>zamek</i> .....	24
Rys.4 Kroki przetwarzania tekstu wejściowego przykładowego zdania.....	27
Rys.5 Moduły aplikacji i ich wzajemne zależności .....	33
Rys.6 Dzieła literackie – wyniki testów.....	52
Rys.7 Wybrane notatki PAP – wyniki testów .....	60
Rys.8 Wszystkie notatki PAP – wyniki testów .....	61
Rys.9 Ograniczony zbiór relacji – wyniki testów.....	62
Rys.10 Ograniczony zbiór relacji – porównanie wyników .....	63
Rys.11 Sztucznie wygenerowane ciągi słów – wyniki testów .....	65
Rys.12 Miejsce opracowanego modułu w algorytmach przetwarzania języka .....	67

## SPIS TABEL:

Tab.1 Relacje paradygmatyczne słownika semantycznego.....	16
Tab.2 Relacje syntagmatyczne słownika semantycznego .....	18
Tab.3 Format przykładowego testu.....	41

# ROZDZIAŁ I

## WSTĘP

Natura języka sprawia, że nie istnieje jeden, uniwersalny sposób wyrażania myśli. Wieloznaczność pojawia się na każdym z poziomów rozumienia – od wyrazów i ich związków, poprzez zdania, a na całych wypowiedziach kończąc. Z mnogością znaczeń najczęściej mamy do czynienia w słownictwie – poszczególne wyrazy, w zależności od intencji mówiącego, mogą być rozumiane na różne sposoby.

Warto zaznaczyć, że zjawisko to i wszystkie zagadnienia z nim związane, rzadko bywają przeszkodą w komunikacji międzyludzkiej. Wieloznaczność stanowi tak integralny element języka, że właściwa interpretacja sposobu rozumienia słowa dokonywana jest przez ludzki mózg automatycznie i nieświadomie – dla języka angielskiego na przykład, około 40% słów może być rozumiane na więcej niż 1 sposób.<sup>1</sup> Wieloznaczność staje się jednak problemem dla aplikacji komputerowych automatyzujących procesy przetwarzania tekstu. Sposobem na jego rozwiązanie jest konstrukcja słowników semantycznych i właściwe ich wykorzystanie w celu wskazywania poprawnych w danym kontekście znaczeń przetwarzanych wyrazów.

Wyróżnić można kilka ogólnych etapów pracy hipotetycznego systemu przetwarzającego język naturalny:

1. Rozpoznawanie mowy (opcjonalnie) – etap zamiany dźwięku na zapis słowny,
2. Tokenizacja i segmentacja – wydzielenie z tekstu podstawowych, nierozdzielnych jednostek ,
3. Analiza morfosyntaktyczna – formalny opis tokenów ze względu na ich właściwości składniowe, rozpoznanie znanych form wyrazowych,
4. Ujednoznacznianie sensu słów – rozstrzygnięcie niejednoznaczności w przypisaniu znaczenia do poszczególnych wyrazów,
5. Analiza składniowa – przypisanie do poszczególnych wyrazów językowych struktur składniowych ,
6. Analiza semantyczna – przypisanie do wyrazów językowych wyrazów języka formalnego,
7. Analiza dyskursu – analiza powiązań znaczeniowych i pełnego znaczenia wypowiedzi.<sup>2</sup>

---

<sup>1</sup> Y.Ravin, C.Leacock *Polysemy, Theoretical and Computational Approaches*, Oxford University Press 2000, s.1.

<sup>2</sup> D.Jurafsky, J.H.Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, Prentice Hall 2000, s. 3-4.

Opracowany algorytm realizuje czwarty z wymienionych powyżej punktów analizy tekstu, służącym za wejście do bardziej zaawansowanych algorytmów przetwarzania wypowiedzi języka naturalnego. Tak uzyskana wiedza o analizowanym tekście może zostać wykorzystana w systemach z wielu dziedzin, jak choćby:

- Automatyczne tłumaczenie tekstu – obecnie duża liczba rozwiązań problemu opiera się na podejściu statystycznym. Duży korpus tekstów pozwala na zmniejszenie problemu wieloznaczności wypowiedzi poprzez statystyczne wybranie odpowiedniego terminu. Zastosowanie naszego algorytmu, poprzez wybranie konkretnego terminu spośród wielu terminów, pozwoli na poprawę jakości generowanego tłumaczenia. Istnieją potwierdzone przypadki polepszenia wyników działania tłumaczenia opartego na statystyce przez dodanie modułu ujednoznacniającego tekst.<sup>3</sup>
- Wyszukiwanie pozycji w katalogach zasobów w poszukiwaniu zadanego terminu (przykładem może być elektroniczna biblioteka książek) – zastosowanie algorytmu pozwala na doprecyzowanie zapytania poprzez wskazanie bardziej zawężonego terminu, bądź powiększenie zapytania o słowa powiązane z wprowadzonymi. Pewną analogią do powyższego algorytmu jest podpowiadanie zapytania użyte w silniku wyszukiwarki Google. Jednak w tamtym przypadku metoda opiera się na podejściu statystycznym, stąd analogia jest jedynie na poziomie użytkowym.<sup>4</sup>
- Kolejną dziedziną, w której użycie naszego algorytmu może poprawić jakość usługi jest pozycjonowanie reklam dla konkretnego tekstu – dla istniejących rozwiązań, jako przykład można podać Google AdWords, istotne jest umieszczenie reklamy przy treści jak najbardziej się do niej odnoszącej. Rozwiązanie firmy Google opiera się na pozycjonowaniu reklam na stronach, które zawierają zestaw słów kluczowych z nią powiązanych.<sup>5</sup>

Wymienione wyżej przykłady stanowią najbardziej oczywiste zastosowania opisywanego algorytmu. Należy jednak zaznaczyć, że praktycznie każde zagadnienie które opiera się na przetwarzaniu tekstu, podczas którego występuje zjawisko wieloznaczności może zostać przy jego użyciu ulepszone. Dalszymi, możliwymi zastosowaniami może być automatyczne tłumaczenie tekstów, wyszukiwarki dla pojęć (a nie słów, jak jest obecnie), uczenie maszynowe i inne.

Niniejsza praca jest opisem próby zmierzenia się z problemem występowania wieloznaczności w procesie przetwarzania języka naturalnego i jej celem było opracowanie i implementacja algorytmu pozwalającego na wybór najbardziej prawdopodobnego ze znaczeń wyrazów wieloznacznych wypowiedzi, w zależności od kontekstu ich występowania.<sup>6</sup>

---

<sup>3</sup> Y.S.Chang, H.T.Ng, D.Chiang, *Word Sense Disambiguation Improves Statistical Machine Translation*, [http://www.isi.edu/~chiang/papers/chan\\_wsd\\_in\\_mt.pdf](http://www.isi.edu/~chiang/papers/chan_wsd_in_mt.pdf), s.8.

<sup>4</sup> *Frequency Asked Questions for Google Suggest*: <http://labs.google.com/suggestfaq.html>.

<sup>5</sup> *Pozycjonowanie w AdWords*: <http://adwords.google.pl/support/bin/answer.py?answer=65133&topic=9352>.

<sup>6</sup> Kontekst można zdefiniować jako zbiór słów współwystępujących w rozpatrywanej wypowiedzi i dostarczających informacji o jej znaczeniu, do których to ograniczamy przestrzeń poszukiwań.

Istnieje wiele technik zmierzania się z tym zagadnieniem, takich jak na przykład:

- metody probabilistyczne,
- sieci neuronowe,
- metody regułowe,
- drzewa decyzyjne,
- gramatyki formalne i inne.

W trakcie realizowania pracy wybrano podejście relacyjne, w którym to pojęcia i związki pomiędzy nimi zachodzące uporządkowane są w formie relacyjnej sieci semantycznej.

Podczas implementacji wykorzystano istniejące komponenty stworzone w ramach prac Grupy Lingwistyki Komputerowej Katedry Informatyki Akademii Górniczo-Hutniczej oraz Katedry Lingwistyki Komputerowej Uniwersytetu Jagiellońskiego: Słownik Fleksyjny Języka Polskiego w postaci biblioteki CLP<sup>7</sup> oraz Słownik Semantyczny Języka Polskiego<sup>8</sup> który to dostarczył wiedzy na temat wspomnianych wyżej bytów i ich zależności. Algorytm został zaimplementowany w postaci biblioteki języka C++, działającej zarówno na platformie Windows jak i Unix/Linux.

Poniżej omówiona została zawartość poszczególnych części pracy.

W rozdziale drugim przedstawiony został krótki przegląd zagadnień związanych z pojęciem wieloznaczności – potrzeba ujednoznaczniania treści na poziomie semantycznym. Omówione ponadto zostały różne podejścia dotyczące problemu postrzegania znaczeń przez ludzi i wywodzącą się z jednego z nich idea słownika semantycznego.

Rozdział trzeci opisuje sposób, w jaki słownik semantyczny może być wykorzystany w problemie rozstrzygania znaczeń wyrazu w zależności od kontekstu ich występowania, oraz szczegóły opracowanego algorytmu ujednoznaczniającego.

W rozdziale czwartym zaprezentowano szczegółowy sposób implementacji omówionego wcześniej algorytmu. Przedstawiony został podział na moduły, przybliżone zostały kolejne kroki algorytmu, opis najważniejszych struktur i cel ich zastosowania.

Rozdział piąty zawiera opis zbiorów testowych wraz z otrzymanymi dla nich wynikami. Dla każdego rodzaju testów dokonano analizy zwróconych rezultatów.

Szósty rozdział zawiera podsumowanie i omówienie wyników pracy, a także przegląd możliwości dalszego rozszerzenia zarówno algorytmu, jak i słownika semantycznego.

---

<sup>7</sup> Zob. [http://winnie.ics.agh.edu.pl/proj\\_uk/fleksbaz/](http://winnie.ics.agh.edu.pl/proj_uk/fleksbaz/).

<sup>8</sup> Zob. [http://winnie.ics.agh.edu.pl/proj\\_re/slse/index.html](http://winnie.ics.agh.edu.pl/proj_re/slse/index.html).

## ROZDZIAŁ II

### SŁOWNIK SEMANTYCZNY

#### 2.1. SŁOWNIK SEMANTYCZNY JĘZYKA POLSKIEGO

Zdolność rozstrzygnięcia wieloznaczności wyrazu na poziomie syntaktycznym znacząco wpływa na polepszenie skuteczności działania algorytmów z dziedziny przetwarzania języka naturalnego – wykorzystanie odpowiedniej konstrukcji baz odmian słów języka jest pierwszym etapem działania większości zaawansowanych algorytmów przetwarzania tekstu. Specyfika języka powoduje jednak obecność wieloznaczności także na poziomie semantycznym. Zjawisko to nazywane jest polisemią, gdy słowa wieloznaczne wywodzą się ze wspólnego źródła (przykładem czego jest słowo *język*, w rozumieniu *organu ciała* lub *mowy ludzkiej*) lub – w przeciwnym razie, czyli braku posiadania wspólnych korzeni – homonimią.<sup>9</sup>

Homonimia może występować w:

- morfologii fleksyjnej – gdy pokrywają się jedynie niektóre formy odmiany, np. słowo *dam*, jako forma rzeczownika *damy* i czasownika *dawać*,
- słowotwórczej – np. *ranny* jako *zraniony*, lub *ranny* w rozumieniu *poranka*,
- słownictwie – np. *bal* jako *przyjęcie* (z języka francuskiego), lub *bal* jako *ktoda drewna* (z języka niemieckiego),
- składni – np. wyrażenie *zdrada przyjaciela* które może być rozumiane jako *zdradzenie przyjaciela* lub *zostanie zdradzonym przez niego*.

Słownik semantyczny to zbiór słów wraz z informacjami o ich znaczeniach. Znaczenie to określane jest poprzez zestaw relacji łączących wyraz definiowany z wyrazami go definiującymi.

Wykorzystanie słownika semantycznego potrafiącego przekształcić pisemną reprezentację wyrazu w struktury symbolizujące jego znaczenie jest ważne dla zapewnienia skuteczności działania algorytmów przetwarzania języka naturalnego. Brak tego mechanizmu uniemożliwia przeniesienie

---

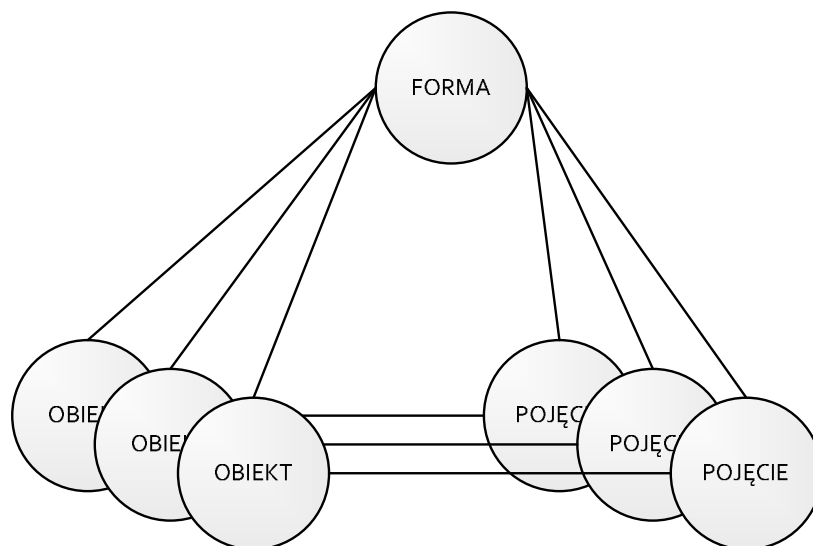
<sup>9</sup> E. Olinkiewicz, K. Radzyńska, H. Styś, *Słownik Encyklopedyczny – Język polski*, Wydawnictwo Europa 1999.



rozumowania z poziomu wyrazów, na poziom operacji na pojęciach połączonych ze sobą relacjami semantycznymi, co jest sposobem symulacji rozumienia tekstu przez maszynę.

## 2.2. DEFINICJE W RAMACH SEMANTYKI

*Semantyka* to dyscyplina badająca relacje pomiędzy nazwami, a przedmiotami których dotyczą. Semantyka zajmuje się badaniem znaczeń, interpretacją znaków, zdań i wyrażeń języka. Problem wieloznaczności w semantyce, pojawia się przy odwzorowaniu obiektów i ich abstrakcji, czyli pojęć w nazwy – często bowiem zdarza się, że różnym obiektom odpowiada ta sama nazwa (rysunek 1).



**Rys. 1.** Trójkąt semiotyczny C.K. Ogdena i A. Richardsa a wieloznaczność

Poniżej omówiono podstawowe pojęcia związane z semantyką i zagadnieniami dotyczącymi przetwarzania języka naturalnego.

Najmniejszą, nierozkładalną część znaczenia nazywamy *cechą semantyczną* (inaczej *sem*). Cechy semantyczne znaczenia danej nazwy nie posiadają jednakowej rangi – można je pogrupować w trzy *strefy konotacyjne*:

- I strefa konotacyjna – obejmuje tzw. *kryterialne cechy semantyczne*, cechy wspólne dla wszystkich przedmiotów oznaczanych przez dany wyraz, pojęcie czy też znak (*desygnaty* danej nazwy), np. *kot* ma *ogon*,
- II strefa konotacyjna – obejmuje tzw. *konotacje semantyczne*, cechy łączone z danym przedmiotem stereotypowo, na zasadzie doświadczenia kulturowego z daną klasą obiektów; dla słowa *kot* mogą to być takie cechy jak: *niechęć do psów, łapanie myszy, przynoszenie nieszczęścia przez koty o czarnej sierści* itp. – nie można cech tych uznać za kryterialne, gdyż można znaleźć koty niespełniające ich,
- III strefa konotacyjna – obejmuje cechy indywidualne łączone z danym obiektem przez danego użytkownika języka, wynikające z jego osobistych przekonań lub doświadczeń, np. *przynoszenie kapci przez kota*.<sup>10</sup>

Zestawienie wszystkich cech semantycznych znaczenia danego znaku językowego tworzy jego *definicję kognitywną* – najbardziej wyidealizowany model, stanowiący idealnie typowy przykład egzemplarza danej kategorii.

Sumę znaczeń danego znaku językowego nazywamy *polem semantycznym* (lub *znaczeniowym*). Obejmuje ono zarówno wszystkie konotacje (cechy współzależne łączone przez nazwę, tworzące jej sens i treść) jak i denotacje (zbiór desygnatów, czyli cech wspólnych obiektu) wyrazu. Do pola semantycznego danego znaku językowego należą te znaki, których kryterialne cechy semantyczne mieszczą się w jego konotacjach.

Na podstawie istnienia cech semantycznych jesteśmy w stanie przyporządkowywać obiekty do określonych *kategorii semantycznych*, czyli zbiorów obiektów wchodzących w zakres określonego znaczenia. Kategorie mogą być całkowicie rozłączne, zawierać się w sobie, przecinać czy też pokrywać – determinując właściwą dla danego języka i kultury taksonomię świata.

Uważa się, że podział na kategorie jest jedną z podstawowych zasad ludzkiego myślenia (E. Sapir, 1978). Za łączenie w grupy informacji o podobnym charakterze odpowiedzialna jest pamięć semantyczna. Jest ona określana jako rodzaj pamięci skojarzeniowej, odpowiedzialnej za etykietowanie zbiorów informacji o podobnym charakterze.

---

<sup>10</sup> M.Brzożowska *Etymologia a konotacja wybranych nazw kamieni w: Etnolingwistyka. Problemy języka i kultury* nr 12, Lublin 2000.

## 2.3. WIELOZNACZNOŚĆ SŁÓW

Wyróżnić można trzy główne podejścia do semantyki i zagadnienia znaczenia wyrazu:

- o składnikowe – opracowane m.in. przez P.E. Goddarda i Noama Chomsky'ego, wywodzące się z filozofii i logiki, a stanowiące rozwinięcie myśli Arystotelesa,
- o prototypowe – które opracowali: Marvin Lee Minski, Roger Schank, Charles J. Fillmore, a wywodzące się z psychologii,
- o relacyjne – opracowane przez George'a A. Millera.<sup>11</sup>

### 2.3.1. PODEJŚCIE SKŁADNIKOWE

Zjawisko mnogości znaczenia wyrazu jest od dawna obserwowane i badane. Już w IV wieku p.n.e. Arystoteles w dziele „Organon” porusza zagadnienie bytu i równocześnie wskazuje na wielość jego znaczeń wymieniając cztery podstawowe:

- o byt akcydentalny,
- o byt wedle kategorii,
- o byt jako prawda i fałsz,
- o byt potencjalny i aktualny.

Dzieli byty według dziesięciu kategorii:

- o substancja albo istota,
- o jakość,
- o ilość,
- o relacja,
- o działanie,
- o doznawanie,
- o miejsce,
- o czas,
- o posiadanie,
- o położenie.

Arystoteles zauważa, że tę samą myśl można wyrazić na wiele sposobów różniących się jedynie stylistyką języka, ale nie kategorią (jak na przykład wyrażenia: *człowiek powraca do zdrowia* i *człowiek jest powracający do zdrowia*).

---

<sup>11</sup> Y.Ravin, C.Leacock, *Polysemy: Theoretical and Computational Approaches*, Oxford University Press 2000, s.7 – 18.

Warto zaznaczyć, że pomimo zwrócenia uwagi na fakt wieloznaczności bytu, Arystoteles nie określa żadnej relacji pomiędzy nimi – istnienie tej złożonej relacji pomiędzy słowem a jego znaczeniem zostało po raz pierwszy dostrzeżone przez stoików. Zauważyli oni, że jedno pojęcie może być opisane przez wiele różnych słów (zjawisko synonimii), a także że pojedyncze słowo może oznaczać wiele pojęć (zjawiska polisemii i homonimii).

Zgodnie z tym podejściem, każdy istniejący obiekt należy do którejś z wymienionych wyżej kategorii. Możemy powiedzieć, że rozpatrywany obiekt należy do pewnej kategorii, jeśli posiada zestaw cech koniecznych (inaczej definiujących) oraz wystarczających (bazowych) dla niej.

Tradycyjnie, przyjmuje się kilka założeń co do tego podejścia:

- o kategorie składają się z list właściwości połączonych operatorami logicznymi takimi jak koniunkcja czy alternatywa,
- o kategorie uporządkowane są w strukturze hierarchicznej, gdzie pojęcia znajdujące się na tym samym poziomie współdzielą odziedziczone właściwości bazowe pojęć z poziomów wyższych, ale mają wzajemnie wykluczające się zbiory właściwości definiujących.

Nowoczesna wersja tego podejścia powstała z dokonanej przez Jerrolda Katza i Jerry'ego Fodora próby kompilacji klasycznego ujęcia teorii znaczenia, z pracami Noama Chomsky'ego opublikowanymi w 1957 roku.<sup>12</sup> Chomsky zauważył, że sama składnia nie tworzy podstawowej struktury języka, a słowa to znaki o określonych właściwościach, które mogą funkcjonować jedynie w odpowiednim kontekście semantycznym. Dowodem jest słynne zdanie jego autorstwa: "Bezbarwne zielone idee wściekle śpią" (ang. "Colorless green ideas sleep furiously") – jest ono poprawne gramatycznie, chociaż bezsensowne. Jego poprawność gramatyczną dowodzi rozkład logiczny, zgodnie z zaproponowanym przez Chomsky'ego fragmentem gramatyki języka angielskiego, w postaci:

$S \rightarrow NP VP$

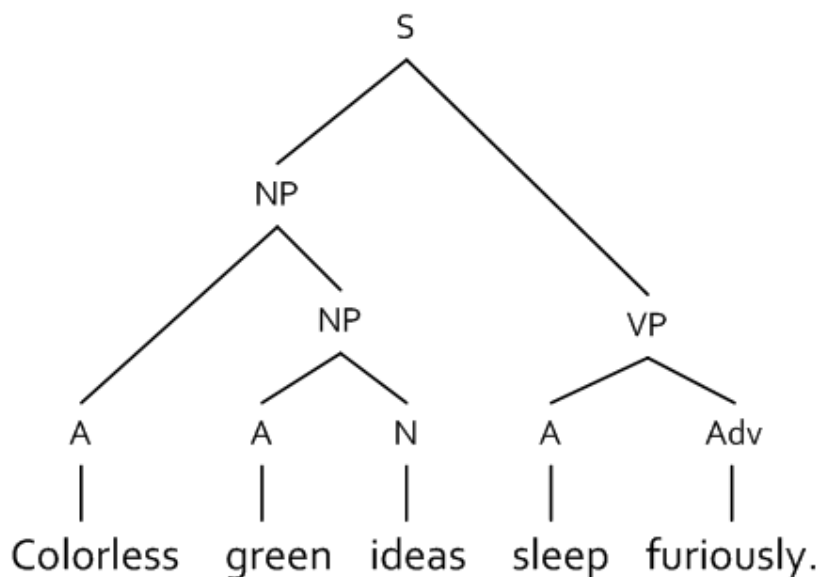
$NP \rightarrow Det N1$

$N1 \rightarrow (AP) N1 (PP)$

---

<sup>12</sup> N.Chomsky, *Syntactic Structures*, Mouton, The Hague 1957.

Rozkład powyższego zdania wygląda następująco:



Rys. 2. Rozkład zdania Chomsky'ego

Celem Katza i Fodora było wyartykułowanie zasad uniwersalnej teorii semantycznej, która tłumaczyłaby reguły rządzące relacjami pomiędzy słowami a ich znaczeniami. Konfrontacja podejścia klasycznego z rzeczywistością okazała się trudna ze względu na rozmycie granic kategorii semantycznych.

### 2.3.2. PODEJŚCIE PROTOTYPOWE

Problemem dla podejścia klasycznego, który wystąpił przy próbie podziału na kategorie był fakt, że różne znaczenia słowa wykazują podobieństwa poprzez zazębianie i przecinanie się właściwości, na takiej samej zasadzie jak cechy gatunków.

W psychologii, teoria kategoryzacji poprzez podobieństwo została przedstawiona w latach siedemdziesiątych przez Eleanor Rosch. Stwierdziła ona, że ludzie nie kategoryzują obiektów poprzez posiadanie przez nich cech koniecznych i wystarczających, ale na zasadzie podobieństwa ich do prototypów – typowych przedstawicieli danej kategorii.

Podejście to uznaje zatem istnienie hierarchicznych kategorii pojęć, ale zmienia sposób przypisywania do nich obiektów poprzez wprowadzenie reprezentującego ją prototypu. Rosch prezentuje dwa rodzaje prototypów:

- o pojedynczy obiekt posiadający możliwie największą liczbę charakterystycznych cech,
- o kilka obiektów, z których każdy posiada oddzielny zestaw charakterystycznych cech, niekoniecznie do siebie zbliżonych – to ten rodzaj grupowego prototypu został przyjęty przez językoznawców.

W swoich badaniach Rosch wielokrotnie wykazywała, że ludzie poznają nowe pojęcia poprzez przybliżanie ich znanymi sobie, istniejącymi obiektami. Prototypy są więc silnie zależne od cech indywidualnych i kulturowych, a stwierdzenie czy dany obiekt należy do kategorii jest tym prostsze i szybsze, im bardziej obiekt jest podobny do prototypu.

Podejście to zostało zaadaptowane przez wielu lingwistów, między innymi wymienionego powyżej Charlesa J. Fillmore'a.

### 2.3.3. PODEJŚCIE RELACYJNE

W podejściu relacyjnym, słowa zgodnie ze swoim znaczeniem zorganizowane są w formę relacyjnej sieci semantycznej. Tak jak w przypadku podejścia prototypowego, model semantyczny pracuje na domenach semantycznych. W idealnym przypadku, znalezienie znaczenia słowa sprowadza się do znalezienia jego lokalizacji w strukturze modelu.

Istnieje wiele typów relacji, które mogą być uwzględnione w strukturze sieci semantycznej. Możemy je podzielić na:

- o relacje paradygmatyczne – dotyczące słów reprezentujących podobne pojęcia i często pojawiających się w zbliżonym kontekście, ale rzadko razem ze sobą (np. słowa *gładki* i *jedwabisty*), a także pojęcia będące swoimi generalizacjami lub specjalizacjami,
- o relacje syntagmatyczne – dotyczące słów współwystępujących często w wypowiedziach języka naturalnego (np. *gładki* i *jedwab*).

Jako najbardziej podstawowe uznaje się następujące relacje paradygmatyczne:

- o synonimię – jedno znaczenie może być wyrażane przez kilka słów,
- o antonimię – przeciwieństwo znaczeń,
- o hiperonimię – relacja generalizacji,
- o hiponimia – relacja specjalizacji.

Przykładem realizacji podejścia relacyjnego jest słownik semantyczny WordNet omówiony w następnym rozdziale.

## 2.4. REALIZACJE SŁOWNIKA SEMANTYCZNEGO

Przykładami słowników semantycznych są: WordNet (język angielski), EuroWordNet (język czeski, duński, estoński, francuski, hiszpański, niemiecki i włoski), Multilingual Central Repository, Global WordNet i inne.

WordNet to słownik semantyczny języka angielskiego stworzony przez Georga A. Millera i jego współpracowników z Princeton University. Budowę słownika rozpoczęto w 1985 roku. Początkowo projekt traktowany był na zasadzie eksperymentu i miał zbadać granice możliwości relacyjnego opisu pojęć poprzez sprawdzenie, jak wiele wiedzy można w ten sposób zapisać. Zgodnie z informacjami podawanymi na stronie projektu, w tym momencie baza słów zawiera 206941 elementów.

Słownik WordNet grupuje rzeczowniki, czasowniki, przymiotniki i przysłowki w kognitywne zbiory „prawie synonimów” nazywane synsetami, z których każdy wyraża osobne pojęcie. Są one ze sobą połączone za pomocą relacji semantycznych i leksykalnych. Synset jest zapisywany w postaci zbioru jednostek leksykalnych reprezentujących poszczególne leksemy np.:

{ anger, choler, ire – (a strong emotion; a feeling that is oriented toward some real or supposed grievance) }

{ anger, angriness – (the state of being angry) }

Pomiędzy synsetami zachodzą relacje semantyczne. Założono, że pomiędzy dwoma synsetami relacja taka zachodzi wtedy i tylko wtedy, gdy występuje ona dla wszystkich par elementów z obu synsetów.

W podejściu relacyjnym wykorzystywane są oba typy relacji semantycznych, jednak istniejące słowniki w większości przypadków (jak np. WordNet) ograniczają się jedynie do relacji paradygmatycznych dzielących pojęcia na kategorie semantyczne:

<u>NAZWA RELACJI</u>	<u>OPIS</u>	<u>PRZYKŁAD</u>
<b>SYNONYMY</b>	Synonimia	<i>zamek(błyskawiczny) – synonimy – suwak</i>
<b>IS A PART OF</b>	meronimia, relacja bycia częścią	<i>zamek(budowla) – is a part of – krajobraz</i>
<b>CONSISTS OF</b>	holonimia, przeciwieństwo meronimii	<i>zamek(budowla) – consists of – budynki mieszkalne</i>
<b>IS A</b>	hiperonimia, przeciwieństwo hiponimii	<i>budowla – is a – zamek(budowla)</i>
<b>IS A KIND OF</b>	hiponimia, relacja bycia rodzajem	<i>zamek(budowla) – is a kind of – budowla</i>

**Tabela 1.** Relacje paradygmatyczne słownika semantycznego

Relacje te klasyfikują pojęcia, opisując hierarchiczne zależności pomiędzy nimi. Sama znajomość taksonomii nie jest jednak wystarczająca do wyłuskania kontekstu użycia wyrazu, ponieważ zależności tego typu rzadko występują w zdaniach w języku naturalnym. Dla zapewnienia skuteczności wyszukiwania kontekstu, należy uwzględnić drugi typ relacji, z którymi to mamy w zdaniach najczęściej do czynienia. Mowa tu o relacjach syntagmatycznych, opisujących akcje i stany, w jakich może dane pojęcie się znajdować.

Zarówno twórcy WordNet'u, jak i innych słowników, zapowiadają rozszerzenie ich o tego typu relacje, jednak jak na razie pozostaje to w sferze planów. Pozostałymi propozycjami uzupełnienia WordNetu są:

- Braki w połączeniach znaczeń z właściwym ich użyciem,
- Słabe odzwierciedlenie wag znaczeń w języku naturalnym,
- Zapewnienie wzajemnego wykluczania się znaczeń w ramach synsetów,
- Lepsze osadzenie w języku współczesnym.<sup>13</sup>

<sup>13</sup> P.Hans, <http://www.fi.muni.cz/gwc2004/pres/panel/Hanks/hanks-panel.pdf>.



## 2.5. SŁOWNIK SEMANTYCZNY KATEDRY LINGWISTYKI KOMPUTEROWEJ UJ

W trakcie implementacji wykorzystaliśmy słownik semantyczny stworzony poprzez Katedrę Lingwistyki Komputerowej Uniwersytetu Jagiellońskiego. Słownik ten znajduje się w fazie mocno rozwojowej i – jak na razie – zawiera niecałe 52000 pojęć.

### 2.5.1. RELACJE SYNTAGMATYCZNE W SŁOWNIKU SEMANTYCZNYM UJ

Słownik opracowany na Uniwersytecie Jagiellońskim tym różni się od innych dostępnych słowników semantycznych, że oprócz wymienionych powyżej relacji paradygmatycznych, został on rozszerzony o relacje syntagmatyczne opisujące akcje i stany w jakich może znajdować się rozpatrywany obiekt.

Pełny spis relacji syntagmatycznych wykorzystanych w słowniku przedstawiony został w tabeli 2:

<b><u>NAZWA RELACJI</u></b>	<b><u>OPIS</u></b>	<b><u>PRZYKŁAD</u></b>
<b>ACTION</b>	czynności wykonywane przez pojęcia, co do których nie można określić ich zgodności czy też niezgodności z przeznaczeniem	<i>mięta – action – pachnieć</i>
<b>ACTION NEGATIVE</b>	akcja negatywna (czyli niezgodna z przeznaczeniem) z punktu widzenia pojęcia	<i>admiral – action negative – poddanie się</i>
<b>ACTION NEGATIVE RT (related to)</b>	pojęcia, których dotyczyć może akcja negatywna	<i>admiral – action negative rt poddanie się – przeciwnik</i>
<b>ACTION PASSIVE</b>	akcja bierna, w której podmiot jest odbiorcą czynności	<i>admiral – action passive – zranienie</i>
<b>ACTION PASSIVE RT (related to)</b>	pojęcia powiązane z istniejącą akcją bierną	<i>admiral – action passive rt zranienie – przeciwnik</i>
<b>ACTION POSITIVE</b>	akcja pozytywna (zgodna z przeznaczeniem) z punktu widzenia pojęcia	<i>admiral – action positive – zwyciężanie</i>

<b>ACTION POSITIVE RT (related to)</b>	Pojęcia powiązane z istniejącą akcją pozytywną	<i>admirał – action positive rt zwyciężanie – przeciwnik</i>
<b>ACTOR</b>	aktor	<i>bicie – actor – gracz</i>
<b>DESTINATION</b>	przeznaczenie, określa do czego służy dany przedmiot	<i>admirał – destination – dowodzenie</i>
<b>DESTINATION RT (related to)</b>	pojęcia powiązane z istniejącą relacją przeznaczenia	<i>admirał – destination rt – ludzie</i>
<b>OBJECT</b>	przedmiot mogący być celem danej czynności	<i>bicie – object – pionek</i>
<b>ROLE</b>	rola w jakiej obiekt może występować, specjalizacja	<i>admirał – role – szef sztabu</i>
<b>ROLE RT (related to)</b>	pojęcia powiązane z relacją roli	<i>admirał – role rt – flota</i>
<b>SIMILAR TO</b>	relacja bycia podobnym do	<i>zamek(budowla) – similar to – twierdza</i>
<b>SOURCE</b>	materiał wykonania	<i>zamek(w drzwiach) – source – metal</i>
<b>STATE</b>	stan w jakim może znajdować się pojęcie	<i>admirał – state – czytanie</i>
<b>STATE NEGATIVE</b>	stan negatywny w jakim może znajdować się pojęcie	<i>admirał – state negative – tchòrzostwo</i>
<b>STATE NEGATIVE RT (related to)</b>	pojęcia powiązane ze stanem negatywnym	<i>admirał – state negative rt tchòrzostwo – przeciwnik</i>
<b>STATE POSITIVE</b>	stan pozytywny w jakim może znajdować się pojęcie	<i>admirał – state positive – lojalność</i>
<b>STATE POSITIVE RT (related to)</b>	pojęcia powiązane ze stanem pozytywnym	<i>admirał – state positive rt lojalność – przełożony</i>

**Tabela 2.** Relacje syntagmatyczne słownika semantycznego

Relacje typu RT (related to) umożliwiają wprowadzenie relacji trójstronnych, określając zachodzenie pewnej relacji warunkowo, w stosunku do określonego obiektu.

Powyższe relacje odzwierciedlają zależności występujące pomiędzy pojęciami, rozszerzając taksonomię. Dzięki ich istnieniu możliwe jest określenie w sposób bardziej precyzyjny związków występujących pomiędzy pojęciami.

## 2.5.2. PRÓBKA SŁOWNIKA SEMANTYCZNEGO

Poniżej przedstawiona została próbka słownika semantycznego dla hasła ADMIRAŁ:

Hasło słownika: admiral (oficer marynarki)

CATEGORY: human  
IS A PART OF: marynarka, flota,  
IS A:  
IS A KIND OF: marynarz  
DESTINATION: dowodzić/dowodzenie, walczyć/walka, żeglować/żeglowanie/żegluga  
ROLE:  
RELATED TO: flota, eskadra, flotylla IS: dowódca, szef sztabu  
ACTION:  
POSITIVE:  
RELATED TO: przeciwnik IS: pokonać, zwyciężyć/zwycięstwo, zatopić/zatopienie  
NEGATIVE:  
RELATED TO: morze, woda, ocean IS: utonąć  
RELATED TO: przeciwnik IS: przegrać/przegrana, poddać się  
RELATED TO: burta IS: wypaść (za)  
PASSIVE: choroba, niedyspozycja  
RELATED TO: przeciwnik IS: dostać, oberwać, rana, zraniony, ranny, obrażenia, zginąć  
STATE:  
POSITIVE: odwaga/odważny, służba/służyć  
RELATED TO: podwładny IS: autorytet, szacunek/szanowany, popularność/popularny, ceniony  
RELATED TO: przełożony IS: posłuszny/posłuszeństwo, lojalny/lojalność  
NEGATIVE: rezerwa, emerytura  
RELATED TO: przeciwnik, żywioł IS: tchòrz/tchòrzostwo/tchòrzliwość

Hasło słownika: admiral (motyl)

CATEGORY: animal  
IS A PART OF: natura  
IS A:  
IS A KIND OF: motyl  
DESTINATION: zapylić/zapylać/zapylenie  
ROLE:  
ACTION:  
    POSITIVE:  
        RELATED TO: kwiat IS: (zbierać) nektar  
    NEGATIVE: choroba, niedyspozycja  
    PASSIVE: zginąć  
STATE:  
    POSITIVE: żyć/żywy  
    NEGATIVE: martwy

Powyższy przykład prezentuje relacje paradygmatyczne i syntagmatyczne dla wyrazu *admiral*. Pierwsze z nich pozwalają na kategoryzację pojęć, drugie odzwierciedlają związki pomiędzy elementami świata rzeczywistego. Ich obecność daje nam możliwość wychwycenia dodatkowych zależności pomiędzy słowami obecnymi w tekście i może posłużyć do rozstrzygnięcia ewentualnych wieloznaczności wyrazów.

## ROZDZIAŁ III

### ROZSTRZYGANIE WIELOZNACZNOŚCI

Stworzony algorytm stanowi kolejny element, mający na celu polepszenie procesu maszynowego rozumienia tekstu. Jego użycie pozwala na poprawę jakości działania narzędzi z dziedziny przetwarzania języka naturalnego, obejmujących automatyczne tłumaczenie tekstu, kategoryzację tekstów językowych, pozycjonowanie reklam kontekstowych, wyszukiwanie informacji na podstawie wyrażeń języka naturalnego i inne.

Idea działania algorytmu opiera się na następujących obserwacjach:

[1] Dla każdego potrafiącego budować kontekst wypowiedzi słowa, można wskazać zbiór innych słów bliskich mu znaczeniowo,

[2] Dodawanie do wypowiedzi słów bliskich znaczeniowo słowom budującym kontekst skutkuje wzmocnieniem go,

Oraz:

[3] Właściwe znaczenie słowa wieloznacznego w wypowiedzi determinowane jest przez kontekst, w którym ono występuje, który tworzony jest przez istniejącą relacje syntagmatyczne.

Zadanie algorytmu, którym jest wskazanie najtrafniej pasującego do kontekstu znaczenia słowa wieloznacznego, sprowadza się zatem do znalezienia dla istotnych słów wypowiedzi, najbliższego im znaczeniowo homonimu. Odległość jest tu obliczana poprzez ilość i typ relacji semantycznych, za pomocą których możemy przemieszczać się między słowami. Wiedzę o zachodzących pomiędzy pojęciami relacjach dostarcza słownik semantyczny.

Algorytm w procesie ujednoznaczniania opiera się na trzech zmiennych:

- o grafie – w którym zapisane są informacje o powiązaniach semantycznych między słowami na podstawie informacji słownika semantycznego,
- o zmiennych liczbowych – określających horyzont przeszukiwania grafu i wagi poszczególnych relacji,
- o aktualnym wejściowym ciągu znaków do sprawdzenia.

W pierwszym kroku algorytmu, wyekstrahowane zostają wszystkie występujące w wejściu słowa. Następnie poszukiwane są odpowiadające im wierzchołki. W przypadku słów wielosegmentowych (np. *energia słoneczna*) znaczenie ma kolejność ich występowania w zdaniu (np. w zdaniu *panna była młoda*, algorytm nie wyszuka wielosegmentowego słowa *panna młoda*). Na tym etapie ma też miejsce podział słów na jedno- i wieloznaczne, którego kryterium jest ilość znaczeń słowa w grafie. W kolejnym kroku, dla każdego słowa, które zostało uznane za wieloznaczne uruchamiany jest algorytm przeszukiwania grafu wszerz. Poszukuje on skierowanych ścieżek łączących to słowo z pozostałymi znalezionymi. Wagi relacji na poszczególnych ścieżkach są sumowane. Otrzymana wartość określa odległość znaczeniową pomiędzy słowami. Algorytm wybiera to, którego suma jest minimalna, a ilość znalezionych słów największa (algorytm potraktuje tak samo fakt znalezienia jednego słowa w odległości 1, jak dwóch słów w odległości 2). Dodatkowo zwracany jest procentowy udział tej sumy w ramach sum wszystkich ścieżek, co jest interpretowane jako stopień pewności wyboru danej ścieżki.

Ze względu na trudność określenia faktycznych granic wypowiedzi, aspekt ten został pominięty – zakłada się, że na wejściu programu znajduje się wystarczająca ilość informacji do określenia kontekstu. Przy analizie tego zagadnienia, konieczne jest uwzględnienie składni języka polskiego, w celu określenia granic zdań wtrąconych i zależności pomiędzy poszczególnymi ich częściami. Ze względu na naturę języka, odseparowanie wpływu kontekstu wypowiedzi od jej składni jest niemożliwe. Stworzenie pełnego algorytmu obejmującego oba te zagadnienia wykracza poza ramy tej pracy magisterskiej i zostało pominięte w dalszych rozważaniach.<sup>14</sup>

Ujednoznacznienie tekstu przeprowadzane jest dla wejściowej wypowiedzi. Na żadnym z etapów działania algorytmu nie jest wykonywana operacja korekty pisowni, dlatego zakładamy poprawność zarówno słownika semantycznego jak i wejściowego ciągu znaków.

---

<sup>14</sup> Analiza składniowa odbywa się najczęściej w oparciu o pewną sformalizowaną gramatykę, omówienie istniejących rozwiązań można znaleźć w: A. Mykowiecka, *Inżynieria lingwistyczna, Komputerowe przetwarzanie tekstów w języku naturalnym*, Warszawa 2007.

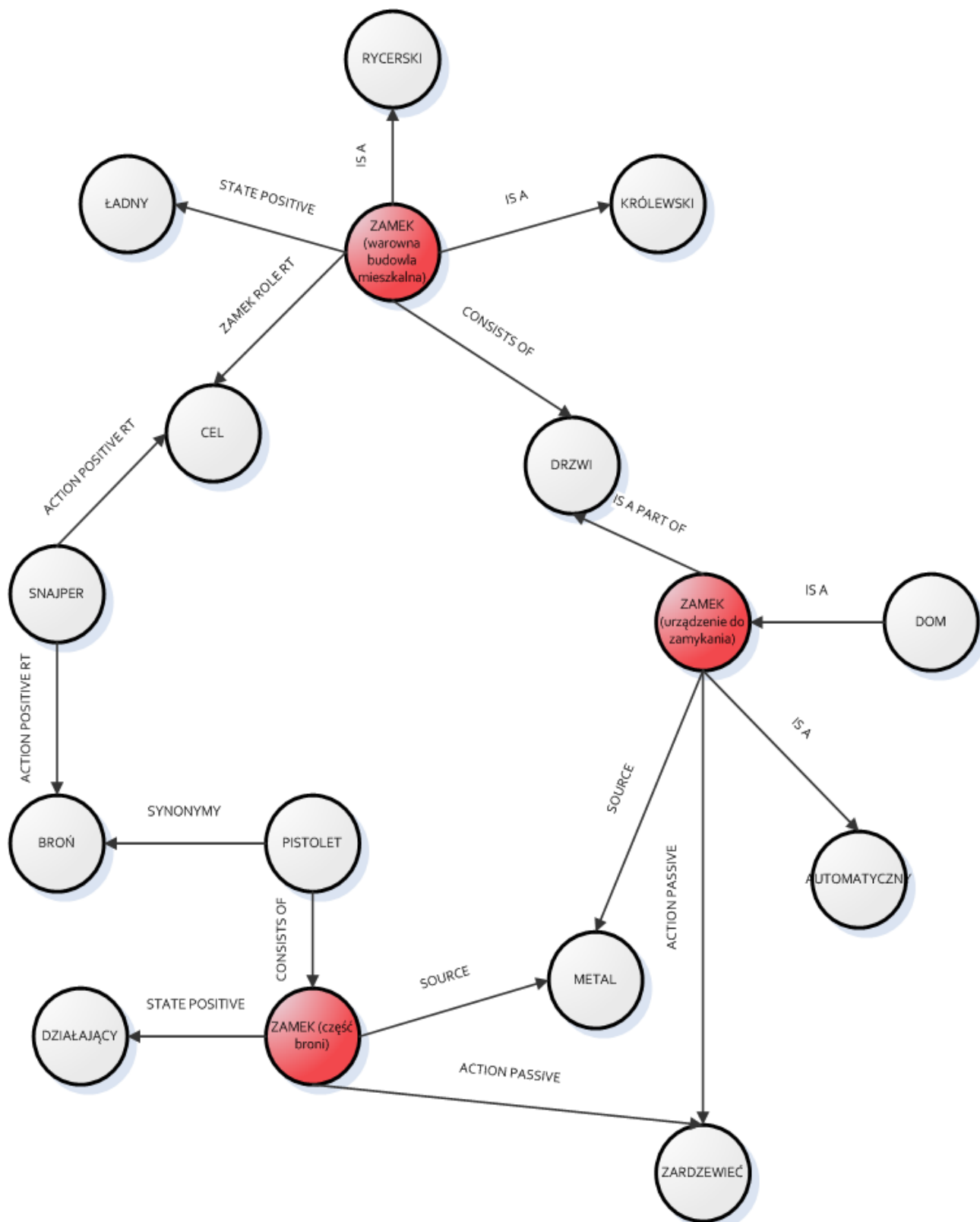
### 3.1. MODEL

W celu przeprowadzenia procesu ujednoznaczniania słów, konieczna jest wiedza o wzajemnych zależnościach pomiędzy wyrazami. W tym celu posługujemy się słownikiem semantycznym, który stanowi statyczną reprezentację relacji zachodzących pomiędzy pojęciami języka naturalnego. Algorytm, o który oparto działanie programu jest algorytmem grafowym. Aby możliwe było przeprowadzanie procesu rozumowania i znalezienie zależności pomiędzy słowami, konieczne było odpowiednie odwzorowanie słownika semantycznego w strukturę grafu skierowanego, ważonego.

Poniżej przedstawiono odwzorowanie elementów pomiędzy modelami:

- o wierzchołki (odwzorowujące pojęcia) – reprezentują byty istniejące w modelu; dla każdego pojęcia istnieje odpowiadające mu słowo lub ich ciąg ze słownika semantycznego; na potrzeby modelu dla pojęć nie rozróżniamy części mowy ani wielkości liter,
- o skierowane krawędzie (odwzorowujące relacje) – reprezentują istniejące zależności pomiędzy pojęciami; pokrywają się z relacjami istniejącymi w słowniku semantycznym; każdej relacji przyporządkowana jest dodatkowo liczbowa wartość określająca wagę relacji.

Dla potrzeb algorytmu wszystkie nazwy przechowywane w modelu występują w pierwszej formie podstawowej, którą zwraca biblioteka CLP. W przypadku wczytania pojęcia z grafu, które nie spełnia tego warunku jest ono sprowadzane do tej postaci.



Rys. 3. Fragment modelu dla słowa wieloznacznego *zamek*



Stworzony model może być przedstawiony w postaci ważonego grafu skierowanego<sup>15</sup>, w którym pojęcia przedstawione są jako wierzchołki, a relacje zachodzące pomiędzy nimi – wyrażane poprzez krawędzie. Relacje są skierowane i łączą ze sobą dwa, niekoniecznie różne (np. relacja zwrotna), pojęcia. Dla dwóch dowolnie wybranych słów ze słownika semantycznego odległością w modelu nazywamy długość ścieżki pomiędzy nimi. Każda relacja posiada określający ją współczynnik ważności, mówiący o tym, jak mocny jest związek pomiędzy dwoma pojęciami. Istnienie tego współczynnika motywowane jest mnogością relacji i potrzebą rozróżnienia wprowadzanych przez nie zależności przy jednoczesnym zachowaniu prostoty samego algorytmu. Jest to istotne, gdyż poszukiwanie w sąsiedztwie słowa wieloznacznego słów współwystępujących z nim w zdaniu testowym odbywać się musi w pewnym promieniu, którego wielkość nie powinna być ustalona arbitralnie, ale uzależniona od rodzaju rozpatrywanych relacji. Im waga relacji jest większa, tym wnosi mniejszy narzut na rozmiar promienia poszukiwań. Najsilniejszą relacją jest relacja synonimii. Przejście pomiędzy dwoma pojęciami z użyciem tej relacji ma zerowy koszt.

Słownik semantyczny z definicji zawiera relacje należące do pierwszej i drugiej strefy konotacyjnej, czyli informacje na temat ogólnie znanych cech obiektów. Nie dostarcza on zatem ani wiedzy, ani mechanizmu wnioskowania o cechach nietypowych, indywidualnych, skojarzeniach i prawdach relatywnych (np. fantazja literacka, operowania pewnymi relacjami tylko w danym zakresie i pod określonymi warunkami). Kolejnym ograniczeniem słownika jest fakt, że o ile informacje na temat instancji poszczególnych obiektów dobrze komponują się w jego strukturę, to nie ma łatwego sposobu na określenie zależności czasowych ani przestrzennych. Przykładem jest zdanie *Średniowieczne zamki były bardzo awaryjne*. – w procesie ujednoznaczniania słowa *zamek* można pominąć *zamek błyskawiczny*, który został wymyślony w XX wieku.

### 3.2. OBRÒBKA ZDANIA

Wejściem algorytmu jest ciąg znaków. W trakcie działania algorytmu wpływ znaków interpunkcyjnych na kontekst w jakim występują słowa jest ignorowany. Wszystkie operacje związane z analizą interpunkcji, podziałem wypowiedzi na akapity i zdania muszą być przeprowadzane przed działaniem algorytmu. Podyktowane jest to chęcią oddzielenia od siebie zakresu działania poszczególnych modułów biorących udział w procesie przetwarzania wypowiedzi. Jedyną operacją podczas której interpunkcja jest uwzględniana, jest podział wejścia na słowa. W pierwszym kroku algorytm dzieli wejście na poszczególne słowa, które w dalszym opisie będziemy oznaczać jako wektor wejściowy. W celu wydajnego wyszukiwania wyrazów wielosegmentowych zostało stworzone dodatkowe mapowanie, które pozwala znaleźć je na podstawie pojedynczego fragmentu. Wyraz jest identyfikowany jako wieloznaczny tylko w przypadku, gdy wszystkie fragmenty wchodzące w jego skład znajdują się we właściwej kolejności w zdaniu.

---

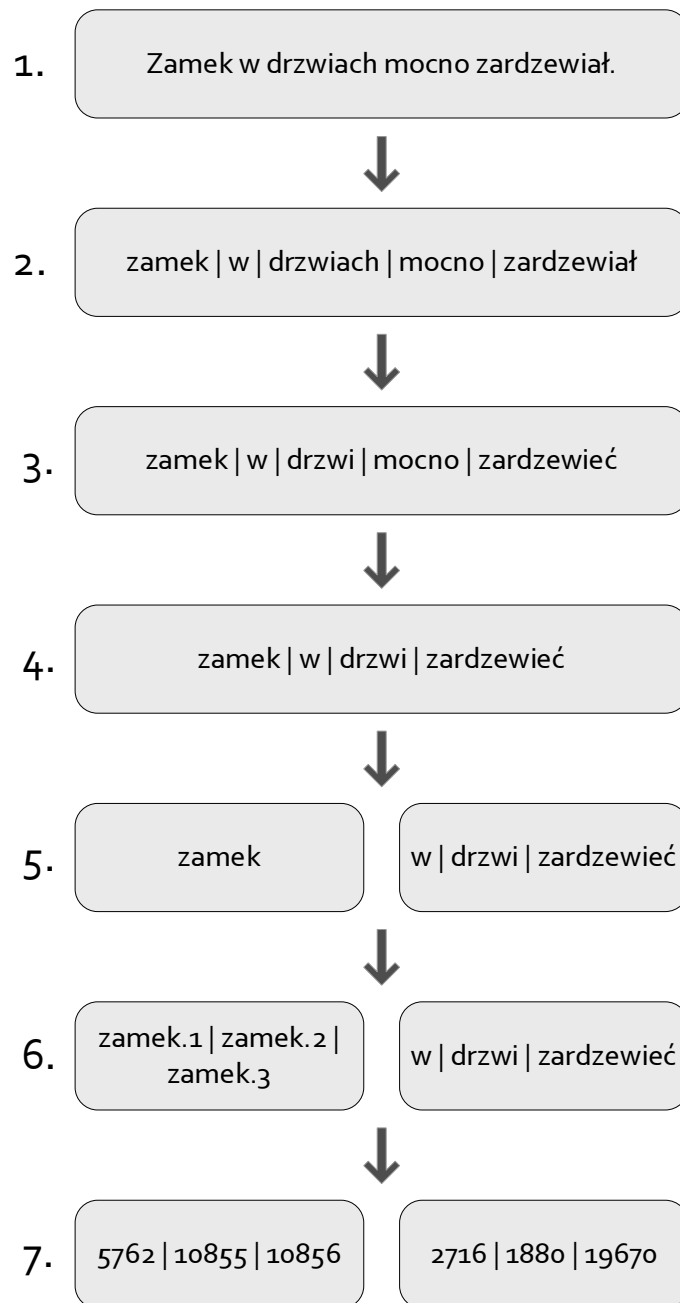
<sup>15</sup> T.Cormen, C.Leiserson, R.Rivest, C.Stein, *Introduction to Algorithms*, MIT Press 2001, s.527-529.

W następnym kroku przeprowadzana jest wstępna identyfikacja wyrazów wieloznacznych. Algorytm pozwala na identyfikację homonimów fleksyjnych (np. *jak* jako zwierzę, partykuła itd.). Pozwala to na ujednoznaczenie pojęcia na podstawie jego formy gramatycznej (np. forma *jakiem* bezpośrednio wskazuje na słowo *jak* w znaczeniu zwierzęcia). Niestety, model słownika semantycznego nie zawiera odwzorowania na identyfikator fleksyjny z biblioteki CLP. Powoduje to brak możliwości wykorzystania tego elementu w procesie działania algorytmu. Poza wyszukiwaniem wyrażań wielosegmentowych, kolejność słów wypowiedzi nie ma znaczenia dla działania algorytmu. Wybór zbioru jako struktury przechowywania słów podyktowany jest nieistnieniem narzędzi umożliwiających uwzględnienie składni w procesie ujednoznaczniania oraz brakiem potrzeby wielokrotnego rozpatrywania tych samych słów (nieumiejętność prawidłowego określenia powiązań między słowami w zdaniu, granic wypowiedzi, itp.).

Przy pomocy biblioteki CLP, słowa z wektora wejściowego sprowadzane są do formy podstawowej. Podyktowane jest to faktem trzymania wszystkich słów modelu w formie podstawowej. Algorytm zakłada poprawność ortograficzną wejściowej wypowiedzi. Na żadnym z etapów algorytmu na zdaniu wejściowym nie jest wykonywana operacja korekty pisowni. Wszelką korektę pisowni należy przeprowadzić przed procesem ujednoznaczniania. Słowa, które nie zostały znalezione przez CLP są traktowane jako błędne i usuwane z dalszej obróbki. Zakłada się, że definicje w bibliotece CLP pokrywają się z definicjami w modelu.

Następnie wyszukiwane w grafie są wszystkie wierzchołki odpowiadające słowom z wektora wejściowego. Jeżeli słowo w nim nie występuje, jest ignorowane (można traktować to jako jego usunięcie z wektora wejściowego). Wprowadzenie takiego ograniczenia powoduje zmniejszenie przestrzeni poszukiwań o słowa, których nie potrafimy interpretować bez wiedzy o składni wypowiedzi. Pozwala to zachować prostotę algorytmu.

Kolejne etapy algorytmu korzystają z dwóch zbiorów wierzchołków. W pierwszym przechowywane są wskazania na te, które interpretowane są jako pojęcia wieloznaczne, w drugim – wskazania na pojęcia jednoznaczne. Dla każdego słowa z wektora wejściowego wyszukiwany jest odpowiadający mu wierzchołek grafu. Jeżeli dla danego słowa w grafie istnieje więcej niż jeden wierzchołek, algorytm automatycznie zakłada, że jest ono wieloznaczne i dodaje je do zbioru do ujednoznacznienia. Jeżeli istnieje tylko jeden wierzchołek, słowo jest dodawane do zbioru jednoznacznych. Algorytm kończy działanie, gdy zbiór pozycji do ujednoznacznienia jest pusty (wszystkie słowa wypowiedzi zostały ujednoznacznione).



**Rys. 4.** Kroki przetwarzania tekstu wejściowego na przykładzie zdania:  
*Zamek w drzwiach mocno zardzewiał.*

### 3.3. ROZUMOWANIE W MODELU

Na tym etapie algorytm ujednoznacza wyrazy homonimiczne i polisemiczne. W procesie ujednoznaczniania, wszystkie słowa ze zdefiniowanych wcześniej zbiorów traktowane są równorzędnie pod względem ważności. W dalszej analizie ignorowany jest wpływ składni języka polskiego na zasięg oddziaływania słów na siebie. Podyktowane jest to trudnościami z prostym opisaniem reguł składni języka polskiego i ich wpływem na proces ujednoznaczniania.

W celu uproszczenia dalszego opisu zostanie wprowadzone pojęcie granicy przeszukiwań przy zadanej wartości. Dla konkretnego wierzchołka grafu, jego granicą przy zadanej wielkości nazywamy wszystkie wierzchołki, które znajdują się w odległości mniejszej bądź równej zadanej wartości, przy czym odległość między wierzchołkami grafu rozumiana jest w sposób standardowy w teorii grafów, czyli jako suma wag najkrótszej skierowanej ścieżki łączącej te dwa wierzchołki.

Parametrami procesu przeszukiwania są:

- o wierzchołek startowy – wierzchołek grafu, od którego liczona jest granica przeszukiwań; wierzchołkami grafu są wierzchołki pojęć wieloznacznych,
- o waga startowa – określająca granicę przeszukiwania grafu od wierzchołka startowego; w momencie uruchomienia algorytmu aktualna wartość granicy jest równa wadze startowej,
- o aktualna waga – określająca aktualną granicę przeszukiwania grafu w otoczeniu wierzchołka startowego,
- o przyrost wagi – określający zmianę granicy przeszukiwania w grafie; aktualna waga jest zwiększana w przypadku braku trafień w danej granicy,
- o maksymalna waga – wyznacza maksymalną granicę, w której wyszukiwane są powiązania pomiędzy pojęciami; jeżeli aktualna waga przekracza tę wartość, aktualny proces przeszukiwania grafu jest przerywany.

W następnym kroku algorytm ustawia wartość dla granicy na wagę startową. Po czym dla wybranego wierzchołka ze zbioru wieloznacznych przegląda graf w przód w poszukiwaniu ścieżek do wierzchołków z obu zbiorów wskazań, pomniejszonego o wierzchołki danego pojęcia. Podyktowane jest to koniecznością uniknięcia sytuacji, w której wierzchołkiem ujednoznaczającym jest inne znaczenie danego pojęcia.

Algorytm nie uwzględnia wpływu słów z wektora wejściowego na zmianę wag relacji między słowami. Spowodowane jest to problemem z określeniem, które relacje z grafu mają wpływ na słowa wektora wejściowego.

Jeżeli w zadanej granicy wierzchołki nie zostaną znalezione, jest ona powiększana o przyrost wagi i proces przeszukiwania jest ponawiany. W przypadku znalezienia wierzchołka lub osiągnięcia granicy maksymalnej, proces jest przerywany. Jeżeli na danej głębokości znaleziono więcej niż jeden wierzchołek, informacja o nich jest zapisywana, a odległości do tych wierzchołków są sumowane. Po zakończeniu przeszukiwania, algorytm jest powtarzany na zbiorze krawędzi wstecz, a wierzchołek

źródłowy staje się wierzchołkiem końcowym ścieżki. Ostatecznie zbiorem wierzchołków wskazujących znaczenie zostają te, które znajdują się najbliżej wierzchołka startowego.

W tym momencie proces przeszukiwania grafu dla danego wierzchołka jest zakończony i ponawiany jest dla kolejnego słowa ze zbioru wieloznacznych. Proces przeszukiwania kończy się w momencie gdy dla wszystkich wierzchołków ze zbioru wieloznacznych przeszukano graf.

### 3.4. ANALIZA WYNIKÓW

Zastosowanie powyższej procedury powoduje znalezienie dla każdego ze znaczeń słowa wieloznacznego, wierzchołków z wypowiedzi, które znajdują się najbliżej niego w grafie i określenie odległości do nich. Algorytm dokonuje podziału otrzymanych wyników na grupy, według słów wieloznacznych i analizuje je osobno. W przypadku gdy tylko dla jednego z tych znaczeń wartość jest niezerowa, algorytm zwraca je w odpowiedzi. Gdy algorytm wskazał więcej niż jedno znaczenie, dla każdego z nich wyliczana jest suma odwrotności odległości poszczególnych wierzchołków. Premiuje to te znaczenia, wokół których w minimalnej możliwej odległości znajduje się najwięcej wierzchołków. Po dodaniu sum algorytm potrafi wskazać procentową część każdego ze znaczeń. Na wyjściu algorytm zwraca wszystkie znalezione znaczenia oraz procentowy udział każdego z nich, przy czym jako najbardziej prawdopodobne wskazuje to, którego suma jest największa. Najgorszym przypadkiem jest sytuacja, w której słowa znajdujące się w wypowiedzi wejściowej znajdują się w granicach wszystkich znaczeń. Próbą poprawienia odpowiedzi dla takiego przypadku jest powiększenie grafu o nowe słowa, lub dodanie dodatkowych słów do wypowiedzi wejściowej (powiększenie kontekstu).

## ROZDZIAŁ IV

### IMPLEMENTACJA

Stworzenie aplikacji składało się z dwóch etapów:

- o wygenerowania pliku modelu na podstawie otrzymanej kopii bazy danych słownika semantycznego,
- o implementacji programu budującego graf sieci semantycznej na podstawie pliku modelu, a następnie ujednoznaczającego wyrazy w zadanej wypowiedzi.

#### 4.1. GENERACJA PLIKU MODELU

Aby umożliwić algorytmowi operowanie na danych słownika semantycznego, konieczne było wygenerowanie pliku zawierającego jego model. W tym celu stworzono zestaw skryptów języka SQL wybierających poprawne dane z udostępnionego nam zrzutu bazy, a następnie formatujących je do właściwej postaci.

Plik modelu zawiera zarówno pojęcia jak i relacje słownika. Każdy element zapisany jest w osobnej linii – najpierw pojęcia, a po nich relacje. Wyrazy jednoznaczne i wieloznaczne rozróżniane są poprzez kolejne numery nadawane im w trakcie formatowania – jako, że w modelu przechowujemy wyłącznie słowo, jest to konieczne dla jednoznacznego ich identyfikowania.

Poniżej wklejono przykładowe linie pliku:

```
akceptowanie
akcesoria do włosów
akcja.1(działanie)
akcja.2(fabuła)
akcja.3(papier wartościowy)
```

Dwie pierwsze pozycje to wyrazy jednoznaczne. Kolejne pojęcie – *akcja* – jest wieloznaczne i słownik zawiera trzy jego definicje. Każde znaczenie pojęcia wieloznacznego określane jest przez

kolejny numer i opis. Algorytm rozpoznaje pojęcie jako wieloznaczne poprzez zakończone sukcesem poszukiwanie numeru porządkowego pojęcia i jego opisu w nawiasach okrągłych.

Po definicjach wszystkich pojęć znajduje się zestaw relacji, czyli wszystkich powiązań między nimi.

```
zamek.1[IS A],automatyczny.1  
zakup.3[SIMILAR TO],towar.1  
zadra[IS A],drzazga  
folwark[IS A KIND OF],wieś
```

Każde powiązanie znajduje się w osobnej linii i składa się z trzech części:

- o wierzchołka początkowego – w przypadku pojęcia jednoznacznego, jego nazwy, a dla wieloznacznego, nazwy wraz z numerem porządkowym, już bez opisu,
- o nazwy relacji w nawiasach prostokątnych,
- o wierzchołka docelowego – także w postaci nazwy pojęcia waz lub bez numeru.

Plik tekstowy w takim formacie jest wczytywany do programu w momencie jego inicjalizacji i przechowywany w pamięci.

## 4.2. IMPLEMENTACJA UJEDNOZNACZNIANIA

Program ujednoznaczniający został zaimplementowany w języku C++. Jego działanie ma charakter sekwencyjny – w danym momencie może być analizowane pojedyncze zdanie testowe, a całość odbywa się w ramach jednego procesu. Kod programu jest przystosowany do pracy zarówno w środowisku Windows (użycie kompilatora firmy Microsoft w ramach Microsoft Visual Studio), jak i w środowisku Unix/Linux (przy użyciu kompilatora gcc).

Program może być używany zarówno jako niezależna aplikacja, jak również moduł biblioteczny. W przypadku uruchomienia programu jako niezależnej aplikacji zostały zaimplementowane dwa sposoby działania:

- o jako program konsolowy,
- o jako webservice.

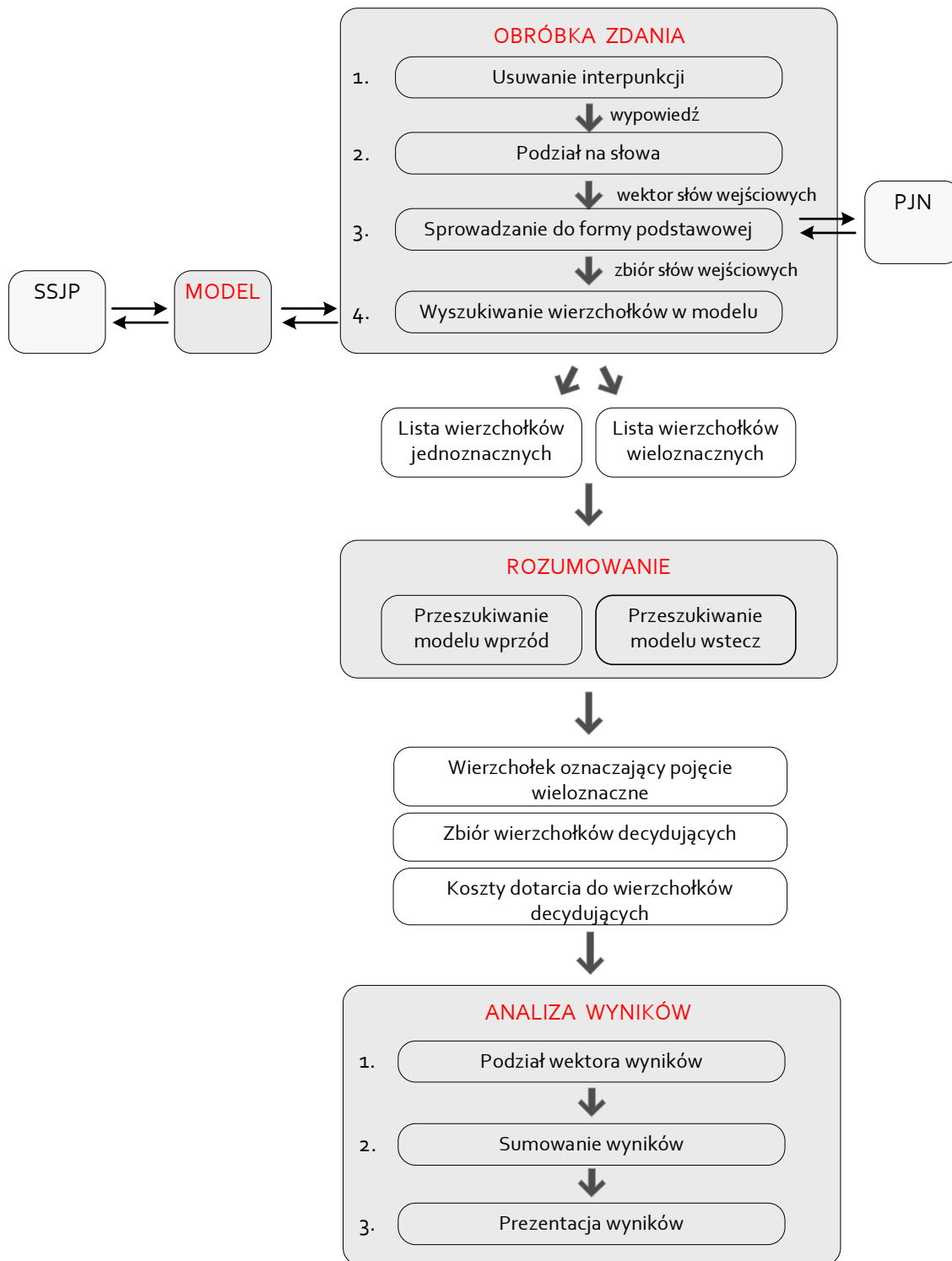
W ramach aplikacji możemy wyróżnić następujące moduły:

- o moduł modelu – odpowiedzialny za wczytywanie, przechowywanie i sprawdzanie spójności struktur potrzebnych do prawidłowego działania algorytmu,
- o moduł odpowiedzialny za obróbkę wejściowego zdania – mapujący słowa wchodzące w skład zdania wejściowego na wierzchołki grafu i identyfikujący w nim słowa wieloznaczne,
- o moduł odpowiedzialny za przeglądanie grafu – wyszukujący kolejno powiązania pomiędzy wierzchołkami wieloznacznymi ze zdania wejściowego a pozostałymi,
- o moduł odpowiedzialny za interpretację wyników – obliczający prawdopodobieństwo jakie algorytm przydziela poszczególnym znaczeniom wieloznacznym słów.

Poza wymienionymi składowymi programu, do wykorzystywane są także dwa zewnętrzne moduły:

- o biblioteka CLP – używana do ujednoznaczniania fleksyjnego słów,
- o Słownik Semantyczny Języka Polskiego – użyty do wygenerowania modelu.





Rys. 5. Moduły aplikacji i ich wzajemne zależności

Moduł odpowiedzialny za przechowywanie danych dla algorytmu korzysta z parametryzowanych struktur biblioteki stl. Podstawową strukturą wykorzystywaną w ramach programu jest klasa string używana zarówno do przechowywania nazw pojęć i relacji w grafie, jak również jako klucze dla map odwzorowujących wyrazy w grafie. W ramach programu sieć semantyczna została zapisana w postaci skierowanego, ważonego grafu. Jego struktura jest przechowywana w postaci list sąsiedztwa.<sup>16</sup> Pojęcia sieci semantycznej zostały odwzorowane na wierzchołki, a relacje pomiędzy nimi – na krawędzie. Informacja o nich jest zapisana przy pomocy dwóch map:

- o `map<SOURCEIDX, EDGES> _relationsForward;`
- o `map<SOURCEIDX, EDGES> _relationsBackward;`

gdzie typy SOURCEIDX, TARGETIDX, RELATIONIDX i EDGES to odpowiednio:

- o `typedef map<RELATIONIDX, set<TARGETIDX>> EDGES;`
- o `typedef unsigned int SOURCEIDX;`
- o `typedef unsigned int TARGETIDX;`
- o `typedef int RELATIONIDX;`

Struktura `_relationsForward` odpowiada za przechowywanie relacji skierowanych od wierzchołka źródłowego. Struktura `_relationsBackward` – za relacje wstecz. Informacje zawarte w nich duplikują się. Taka nadmiarowość pozwala na przyspieszenie procesu przeszukiwania grafu. Poszczególne wierzchołki (odpowiadające pojęciom) grafu są reprezentowane przez kolejne liczby naturalne. Użycie mapy podyktowane jest faktem istnienia wierzchołków, dla których nie istnieją powiązane z nimi relacje wejściowe czy wyjściowe.<sup>17</sup> Z każdym wierzchołkiem powiązana jest struktura EDGES, która zawiera informacje o powiązanych z nim skierowanych krawędziach i wierzchołkach końcowych dla struktury `_relationsForward` i początkowych dla `_relationsBackward`. Ze względu na potrzebę rozróżnienia typów relacji, struktura EDGES została zdefiniowana jako osobna mapa, która separuje od siebie poszczególne relacje i grupy wierzchołków wyjściowych. Wykorzystanie mapy pozwala przyspieszyć sprawdzanie istnienia określonej relacji i, w razie potrzeby, usunięcie jej z procesu przeszukiwania. Obecność zbioru do przechowywania wierzchołków docelowych spowodowana jest potrzebą zmniejszenia narzutu na czas poszukiwania pojedynczego wierzchołka.

---

<sup>16</sup> T.Cormen, C.Leiserson, R.Rivest, C.Stein, *Introduction to Algorithms*, MIT Press 2001, s.529.

<sup>17</sup> U.Breymann, *Designing Components With the C++ STL (2<sup>nd</sup> Edition)*, Addison Wesley 2002, s.78-79.

W celu odwzorowania nazw poszczególnych wierzchołków, a co za tym idzie – również pojęć, zostały zdefiniowane następujące struktury:

```
o vector<string> _vWordIdxToString;  
o map<string, CONCEPTIDX> _mUnambiguousWords;  
o map<string, set<CONCEPTIDX>> _mMultisegmentUniqueConcepts;  
o map<string, vector<DESCRIPTION>> _mAmbiguousWords  
o map<string, set<CONCEPTIDX>> _mMultisegmentAmbigConcepts;  
o vector<string> _vMultisegmentWordsToAnalyze;
```

gdzie typy DESCRIPTION i CONCEPTIDX to:

```
o typedef pair<string, int> Description;  
o typedef int CONCEPTIDX;
```

Struktura `_vWordIdxToString` pozwala na wyszukanie nazwy wierzchołka skojarzonego z indeksem. Natomiast struktury `_mUnambiguousWords` i `_mAmbiguousWords` realizują odwrotną operację, przy czym pozwalają na dodatkowe określenie czy wierzchołek jest jedno- bądź wieloznaczny. Fakt oddzielenia ich od siebie został podyktowany potrzebą zdefiniowania opisów pojęć wieloznacznych i istnieniem więcej niż jednego wierzchołka dla pojęcia wieloznacznego, która to informacja w przypadku wierzchołka jednoznacznego byłaby nadmiarowa. Wartość liczbowa w opisie określa numer porządkowy znaczenia ze słownika semantycznego. Struktura `_mMultisegmentAmbigConcepts` pozwala na podstawie fragmentu słowa określić, czy istnieje wielosegmentowe, wieloznaczne słowo je zawierające i w przypadku odpowiedzi pozytywnej zwraca wektor ich wierzchołków. `_mMultisegmentUniqueConcepts` realizuje podobną funkcjonalność, z tym, że dla słów wielosegmentowych jednoznacznych. Struktura `_vMultisegmentWordsToAnalyze` używana jest na początkowym procesie działania algorytmu. Umieszczane są w niej wszystkie wielosegmentowe słowa, które w dalszym etapie pozwolą na wypełnienie struktur `_mMultisegmentUniqueConcepts` i `_mMultisegmentAmbigConcepts`.

W celu przechowywania informacji o dozwolonych dla algorytmu przeszukiwania relacjach i ich wagach, stworzone zostały następujące struktury:

```
o map<string, RELATIONIDX> _mRelationDescriptionToIdx;  
o map<RELATIONIDX, AllowedRelationProp> _vAllowedRelations;
```

gdzie typ RELATIONIDX to:

```
o typedef int RELATIONIDX;
```

Natomiast AllowedRelationProp to struktura pomocnicza:

```
class AllowedRelationProp {  
    public:  
        int _forwardCost;  
        int _backwardCost;  
};
```

Struktura \_mRelationDescriptionToIdx pozwala powiązać nazwę relacji z odpowiadającym jej indeksem. Natomiast AllowedRelationProp – określić koszt relacji przy przechodzeniu odpowiednio w przód dla \_forwardCost i wstecz dla \_backwardCost.

Definicje dostępnych dla algorytmu relacji znajdują się w osobnym pliku tekstowym. Opis każdej relacji umieszczony jest w oddzielnej linii w nawiasach okrągłych. Na początku linii znajduje się tekstowa nazwa relacji, za nią zdefiniowane są dwie liczby naturalne określające koszt relacji przy przechodzeniu grafu w przód i wstecz. Przykładowa postać pliku wygląda następująco:

```
(action,1,1)  
  
(action negative,1,1)  
  
(action negative rt,1,1)
```

Tworzenie modelu sieci semantycznej polega na wczytaniu zawartości pliku z zapisem grafu połączeń dla odwzorowania słownika semantycznego i pliku z dozwolonymi relacjami. W tym momencie wypełniane są struktury: \_vWordIdxToString, \_mUnambiguousWords, \_mAmbiguousWords, \_vMultisegmentWordsToAnalyze, \_mRelationDescriptionToIdx, \_vAllowedRelations, \_relationsForward, \_relationsBackward. W kolejnym kroku wielosegmentowe słowa zapisywane w strukturze \_vMultisegmentWordsToAnalyze i używane są do wypełnienia struktur \_mMultisegmentAmbigConcepts oraz \_mMultisegmentUniqueConcepts. W trakcie tego procesu, dodatkowo zapisywane są informacje o możliwych błędach wczytywania modelu.

W skład możliwych błędów wchodzi:

- o deklaracja nie jest wieloznaczna; deklaracja została usunięta z listy wieloznacznych – kod błędu 10001,
- o brak opisu wieloznacznego słowa – kod błędu 10002,
- o błędna relacja, relacja pominięta – kod błędu 10003,
- o definicja słowa wieloznacznego w złej kolejności – kod błędu 10004,
- o deklaracja nie jest wieloznaczna, pomimo takiego jej oznaczenia – kod błędu 10005,
- o zła kolejność deklaracji; deklaracja występuje przed deklaracją – kod błędu 10006,
- o brakująca deklaracja została dodana do modelu – kod błędu 10007,
- o w definicji relacji znajduje się niezdefiniowane wcześniej pojęcie – kod błędu 10008,

- o podwójna definicja tej samej krawędzi; drugie wystąpienie zostało zignorowane – kod błędu 10009,
- o wieloznaczna deklaracja dodana drugi raz (przyczyna - konwersja do małych liter) – kod błędu 10010,
- o deklaracja nadpisuje deklarację; dodane dodatkowe mapowanie – kod błędu 10011,
- o jednoznaczne słowo posiada również wersję wieloznaczną; obie dodane do modelu jako wieloznaczne – kod błędu 10012.

Log zawierający informacje o procesie wczytywania modelu znajduje się w katalogu aplikacji w pliku tekstowym *log.txt*. Każda kolejna linia zawiera pojedyncze wystąpienie błędu wraz z określeniem jego typu, numeru linii, w której wystąpił i opisu.

Poniżej znajdują się przykładowe linie z błędami wygenerowanymi w trakcie użytkowania programu:

- o 10002 WARNING(3332): Brak opisu wieloznacznego słowa: "Bartosiak.1()". Dodany opis o treści: "Brak opisu",
- o 10003 WARNING(57459): Błędna relacja: "[0],obiekt fizyczny.1". Relacja pominięta.

Moduł odpowiedzialny za obróbkę zdania wejściowego oparty jest o funkcję biblioteczną `strok`<sup>18</sup> przy użyciu wykluczającego filtra dla znaków `~`!@#%&^&*( )_+{}|: "<>?-=[]\; ',./`. Przy jej udziale dokonywany jest podział wejściowego ciągu znaków na słowa i usunięcie niepotrzebnych lub niedozwolonych znaków. W kolejnym kroku, przy użyciu biblioteki CLP poszukiwana jest forma podstawowa danego słowa. W wersji programu dla systemu Windows, tworzone jest w tym celu sieciowe połączenie socketowe w oparciu o protokół TCP/IP.<sup>19</sup> Dedykowany serwer uruchomiony na maszynie uczelnianej zwraca odpowiednią formę dla zapytania, pośrednicząc w komunikacji z biblioteką CLP. W wersji pod system UNIX/LINUX forma podstawowa jest uzyskiwana bezpośrednio. W przypadku otrzymania kilku form podstawowych, uwzględniane są wszystkie z nich. Następnie, przy użyciu struktur `_mUnambiguousWords`, `_mAmbiguousWords`, `_mMultisegmentAmbigConcepts` oraz `_mMultisegmentUniqueConcepts` określone i zapisane strukturach tymczasowych zostają indeksy wierzchołków słów jedno- i wieloznacznych.

Chodzenie po grafie zostało zrealizowane iteracyjną metodą przeszukiwania wszerz, przy użyciu struktur pomocniczych do przechowywania informacji o aktualnie odwiedzonych wierzchołkach i wartościach odległości do poszczególnych wierzchołków. Ostateczne wyniki zapisywane są w strukturze:

```
class HitCounter {
    public:
        int _travelCost;
        set<PathToVertex> _wordsIndexesWhichCount;
```

<sup>18</sup> Zob. <http://www.cplusplus.com/reference/string/strtok.html>.

<sup>19</sup> Zob. [http://msdn.microsoft.com/en-us/library/ms740673\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms740673(VS.85).aspx).

```
};
```

gdzie PathToVertex to:

```
class PathToVertex {  
    public:  
    vector<CONCEPTIDX> _path;  
    TARGETIDX _targetIdx;  
};
```

Zmienna `_travelCost` określa koszt dotarcia do poszczególnego wierzchołka. Zbiór `_wordsIndexesWhichCount` zawiera informacje o ścieżce (wektor `_path`) do wierzchołka docelowego i jego indeks (zmienna `_targetIdx`).

W fazie analizy następuje podział na wyniki cząstkowe ze względu na różne pojęcia wieloznaczne. Do tego celu używana jest struktura pomocnicza `FormatResultsStruct`.

```
class FormatResultsStruct {  
    public:  
    int _conceptsCounter;  
    int _travelCost;  
    float _wholeSum;  
};
```

W zmiennej `_conceptsCounter` zapisana jest ilość wierzchołków w otoczeniu znaczenia, parametr `_travelCost` określa sumaryczny koszt dotarcia do wszystkich wierzchołków, natomiast `_wholeSum` zawiera zbiorczą sumę odwrotności kosztów dotarcia dla całego pojęcia. W procesie wypisywania wyników wykorzystywane są struktury `_mAmbiguousWords`, `_mUnambiguousWords` i `_vWordIdxToString`, które pozwalają określić nazwy pojęć dla wierzchołków.

W przypadku uruchomienia programu jako niezależnej aplikacji, program wczytuje z dysku model sieci semantycznej i dozwolone relacje, a następnie dla każdej linii z podanego jako parametr pliku ujednoznacznia znalezione w nim pojęcia wieloznaczne. Wynik działania programu zapisywany jest w pliku wyjściowym. W trakcie działania algorytmu nie są modyfikowane pliki tekstowe modelu ani dostępnych relacji.

W przypadku uruchomienia programu jako webservice, zostaje otwarte połączenie sieciowe przy użyciu biblioteki socketów. Za pomocą tego interfejsu następować może komunikacja ze stroną PHP stanowiącą GUI użytkownika, lub innym, dowolnie wybranym modułem.

W przypadku używania programu jako biblioteki, udostępniony został następujący interfejs:

```
typedef int IDX;
typedef float PROBABILITY;
struct Results {
    vector<pair<IDX,PROBABILITY> _results;
};
struct Desambiguator {
    void init();
    Results getResults(const string & input);
    const string & getConceptName(IDX idx);
    const string & getConceptDescription(IDX idx);
    int getConceptRelativeNumber(IDX idx);
};
```

## ROZDZIAŁ V

### ANALIZA ZACHOWANIA ALGORYTMU

W języku naturalnym istnieje wiele rodzajów wypowiedzi. Dla idealnie działającego algorytmu, typ zdania nie powinien wpływać na jakość otrzymanych wyników. W celu przetestowania działania programu, wybranych zostało kilka ich typów, dla których przygotowano zestawy testowe:

- testy autorskie (pkt 5.1.) – mające pokazać specyfikę działania programu (11 zdań testowych),
- testy na zdaniach mowy potocznej (pkt 5.2.) – mające dostarczyć wyników działania dla zdań ubogich w kontekst (6 zdań testowych),
- testy na wyselekcjonowanych fragmentach dzieł literackich (pkt 5.3.) – dzieła literackie stanowią znaczący zasób ludzkiej wiedzy. Zbiór zdań z wybranych utworów literackich prezentuje zachowania dla tej dziedziny (13 wypowiedzi testowych),
- testy na wybranym podzbiórce notatek prasowych PAP zawierających wybrane słowo (pkt 5.4. i 5.5.) – analogicznie jak dla dzieł literackich, przy czym pojedyncza wypowiedź zawiera w sobie terminologię z pojedynczej dziedziny (19 i 181 wypowiedzi testowych).

W celu sprawdzenia działania algorytmu dla szczególnych przypadków przeprowadzono dodatkowe testy:

- testy na zbiorach słów sztucznie wygenerowanych ze słownika semantycznego (pkt 5.6.) – mające sprawdzić wyniki działania w oderwaniu od zupełności modelu (450 testowych ciągów słów),
- powtórzenie wybranych testów dla ograniczonego zbioru relacji (pkt 5.7.) – przeprowadzone w celu sprawdzenia działania algorytmu bez relacji syntagmatycznych (181 zdań testowych).

Pojedynczy test składa się z:

- poddanego analizie zdania,
- zwróconego przez program zbioru znaczeń znalezionych w tekście pojęć wieloznacznych - jeśli znaczenie słowa zostało określone z prawdopodobieństwem 100% (w zadanym



sąsiedztwie nie znaleziono żadnych innych, dopuszczalnych znaczeń) – podawane jest tylko to znaczenie; w przeciwnym wypadku, wskazywane są wszystkie znalezione znaczenia słowa wraz z ich prawdopodobieństwami oraz kluczowymi słowami,

- prawdopodobieństwo każdego z pojęć (jeżeli wskazane zostało więcej niż jedno pojęcie),
- słów wskazujących dla każdego z wskazanych pojęć.

Poniżej przedstawiony został format przykładowego testu:

<b><u>ETYKIETA</u></b>	<b><u>WARTOŚĆ</u></b>	<b><u>ZNACZENIE</u></b>
<i>ZDANIE:</i>	<i>„<u>Zamek</u> w drzwiach zardzewiał.”</i>	<i>Zdanie testowe</i>
<i>ODPOWIEDŹ 1:</i>	<i>zamek (urządzenie do zamykania)</i>	<i>Pierwsza ze zwróconych odpowiedzi</i>
<i>WAGA 1:</i>	<i>0. 666667</i>	<i>Obliczone przez program prawdopodobieństwo pierwszej odpowiedzi</i>
<i>SŁOWA KLUCZOWE 1:</i>	<i>drzwi , zardzewieć</i>	<i>Wskazane przez algorytm słowo wskazujące na pierwszą odpowiedź</i>
<i>ODPOWIEDŹ 2:</i>	<i>zamek (część broni)</i>	<i>Druga ze zwróconych odpowiedzi</i>
<i>WAGA 2:</i>	<i>0. 333333</i>	<i>Obliczone przez program prawdopodobieństwo drugiej odpowiedzi</i>
<i>SŁOWA KLUCZOWE 2:</i>	<i>zardzewieć</i>	<i>Wskazane przez algorytm słowo wskazujące na drugą odpowiedź</i>

**Tabela 3.** Format przykładowego testu

Większość testów wykonana została dla następujących pojęć wieloznacznych:

- zamek w znaczeniu:
  - warownej budowli mieszkalnej
  - części broni
  - urządzenia do zamykania
- ślimak rozumiany jako:
  - motyw dekoracyjny
  - składnik ucha wewnętrznego ssaków
  - element maszynowy w postaci śruby
  - droga w formie linii spiralnej
  - zwierzę z rodziny ślimaków
- kozak jako:
  - żołnierz lekkiej jazdy
  - śmiały, odważny mężczyzna
  - ludowy taniec ukraiński
  - grzyb jadalny
  - ciepłe, damskie lub dziecięce obuwie zimowe o wysokiej, miękkiej cholewie

Słowa te zostały wybrane, ponieważ są one stosunkowo dobrze opisane w słowniku semantycznym (średnia dla modelu to ok. 2 relacje dla pojęcia, a dla każdego ze znaczeń słowa *zamek* występuje odpowiednio 20, 30 i 60 relacji ich dotyczących), a przy tym są popularne więc często występują w zdaniach języka potocznego.

## 5.1. TESTY AUTORSKIE

Test został przeprowadzony na uprzednio przygotowanych zdaniach i miał za zadanie zaprezentować działanie algorytmu. Niektóre spośród przeprowadzonych testów mają charakter przyrostowy, czego celem było pokazanie wpływu zmian w zdaniu wejściowym na otrzymane wyniki.

### 5.1.1. ZBIÓR TESTOWY I OTRZYMANE WYNIKI

Poniżej przedstawiony został zestaw testowy wraz z zwróconymi przez system odpowiedziami i wskazanym prawdopodobieństwem ich występowania, oraz słowa kluczowe, które spowodowały ich wybranie:

[1]

ZDANIE: „Zamek w warszawie był celem wielu ataków.”

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: cel (cel ataku)

[2]

ZDANIE: „Kiedy będę wystarczająco bogaty będę mieszkać w zamku.”

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: mieszkać (zajmować na stałe jakieś pomieszczenie)

[3]

ZDANIE: „Zamek który widzieliśmy był ładny, chociaż zaniedbany.”

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: ładny, zaniedbany

[4]

ZDANIE: „Zamek się zamyka.”

ODPOWIEDŹ: zamek (urządzenie do zamykania)

SŁOWA KLUCZOWE: zamykać

[5]

ZDANIE: „Kupiony przeze mnie pistolet okazał się mieć nie działający zamek.”

ODPOWIEDŹ: zamek (część broni)

SŁOWA KLUCZOWE: pistolet

[6]

ZDANIE: „Źle działający zamek spowodował, że ładunek chybił.”

ODPOWIEDŹ: zamek (część broni)

SŁOWA KLUCZOWE: [działający, działać], ładunek

[7]

ZDANIE: „Będąc na Ukrainie zwiedziłam zamek nazywający się Wysokim Zamkiem.”

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: wysoki zamek

[8]

ZDANIE: „Zamek zarzewiał.”

ODPOWIEDŹ 1: zamek (urządzenie do zamykania)

WAGA 1: 0.500000

SŁOWA KLUCZOWE 1: zarzewieć

ODPOWIEDŹ 2: zamek (część broni)

WAGA 2: 0.500000

SŁOWA KLUCZOWE 2: zarzewieć

[9]

ZDANIE: „Zamek w drzwiach zarzewiał.”

ODPOWIEDŹ 1: zamek (urządzenie do zamykania)

WAGA 1: 0.666667

SŁOWA KLUCZOWE 1: drzwi, zarzewieć

ODPOWIEDŹ 2: zamek (część broni)

WAGA 2: 0.333333

SŁOWA KLUCZOWE 2: zarzewieć

[10]

ZDANIE: „Zamek był metalowy.”

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[11]

ZDANIE: „Zamek był zrobiony z metalu.”

ODPOWIEDŹ 1: zamek (urządzenie do zamykania)

WAGA 1: 0.500000  
SŁOWA KLUCZOWE 1: metal  
ODPOWIEDŹ 2: zamek (część broni)  
WAGA 2: 0.500000  
SŁOWA KLUCZOWE 2: metal

### 5.1.2. ANALIZA

Z rezultatów przeprowadzonych testów można wnioskować, że algorytm dobrze radzi sobie ze zdaniami zawierającymi wyrazy wprost wskazujące na konkretne znaczenie wyrazu. Wyrazy: *cel*, *mieszkać*, *ładny*, *zaniedbany*, występujące w testach [1], [2] i [3], jednoznacznie wskazują na *zamek* w rozumieniu *warownej budowli mieszkalnej* i zostało to poprawnie zinterpretowane przez algorytm. Podobnie jest dla zdań [4], [5] i [6] – słowo *zamykać* wskazuje na *zamek* jako *urządzenie do zamykania*, a *zepsuty*, *działający*, *ładunek* i *pistolet* – na *część broni*. Jak widać na przykładzie zdania numer [7], algorytm równie dobrze radzi sobie z występującym w modelu wyrażeniami wielosegmentowymi (tu np. poprawnie zidentyfikowano znaczenie słowa *zamek* występującego razem z nazwą własną *Wysoki Zamek*).

Dodatkowo, test [6] pokazuje działanie algorytmu, gdy na etapie ujednoznaczniania fleksyjnego zwrócono więcej niż jedną poprawną formę wyrazu (w tym wypadku dla słowa *działający* zwrócono *działać* i *działający*). Algorytm wówczas sprawdza istnienie wszystkich zwróconych form, ale uwzględnia do wyniku tylko jedną – z najniższym kosztem.

Na przykładzie testów [8] i [9] widać, że wraz ze wzrostem ilości informacji w zdaniu, zwiększa się ilość słów stanowiących o kontekście całej wypowiedzi, a co za tym idzie – trafność udzielonej odpowiedzi. W zdaniu [8] zwrócono dwie możliwe odpowiedzi – zarówno *zamek* w rozumieniu *urządzenia do zamykania*, jak i *części broni* może *zardzewieć*. Dostarczenie dodatkowej informacji do zdania (test [9]), czyli jego uszczegółowienie, powoduje zwiększenie prawdopodobieństwa jednego ze znaczeń.

Problemem uwidocznionym w wynikach otrzymanych dla zdań [10] i [11] jest brak relacji pomiędzy różnymi formami pokrewnymi tego samego wyrazu. Pozornie, obejście tego problemu jest stosunkowo łatwe i polega na sprawdzeniu istnienia w modelu każdej możliwej formy rozpatrywanego słowa. Dość prosto można znaleźć kontrprzykłady skuteczności takiego podejścia – po pierwsze, przy braku mechanizmu rozpoznawania wyrażen wielosegmentowych na poziomie sprowadzania słów wypowiedzi do formy podstawowej, łatwo o błędy dla nich. Po drugie, przykładowo, złoty kolor przedmiotu nie jest tożsamy z wykonaniem go ze złota, także zastosowanie tego typu uproszczenia w wielu przypadkach mogłoby odbyć się ze stratą dla kontekstu całości. Umieszczenie w słowniku tylko jednej formy z całej rodziny wyrazów powoduje zatem nadmierne zubożenie słownika. Rozwiązaniem może być uzupełnienie słownika semantycznego o konieczne relacje dla wszystkich poprawnych form wyrazów, lub wskazanie części mowy zachowujących prawdziwość relacji. Wprowadzenie do słownika informacji o więcej niż jednej formie wyrazu, powoduje konieczność wyszukiwania całej rodziny wyrazów wokół danego wierzchołka, co z kolei prowadzi do wzrostu złożoności obliczeniowej procesu przeszukiwania.

## 5.2. TEKSTY JĘZYKA POTOCZNEGO

Kolejny test został przeprowadzony na zdaniach języka potocznego. Problemem w analizie tego typu wypowiedzi jest szczątkowa forma lub wręcz zupełny brak kontekstu.

### 5.2.1. ZBIÓR TESTOWY I WYNIKI

Zestaw zdań wchodzących w skład testu przedstawiał się następująco:

[1]

ZDANIE: „Zepsuł mi się zamek.”  
ODPOWIEDŹ 1: *zamek (urządzenie do zamykania)*  
WAGA 1: 0.500000  
SŁOWA KLUCZOWE 1: *zepsuć się*  
ODPOWIEDŹ 2: *zamek (część broni)*  
WAGA 2: 0.500000  
SŁOWA KLUCZOWE 2: *zepsuć się*

[2]

ZDANIE: „Kupiłem sobie nowy zamek.”  
ODPOWIEDŹ: *(brak wyników)*  
SŁOWA KLUCZOWE: *(brak wyników)*

[3]

ZDANIE: „Zamek który kupiłem okazał się słaby.”  
ODPOWIEDŹ: *(brak wyników)*  
SŁOWA KLUCZOWE: *(brak wyników)*

[4]

ZDANIE: „Zamek mi zarzewiał.”  
ODPOWIEDŹ 1: *zamek (urządzenie do zamykania)*  
WAGA 1: 0.500000  
SŁOWA KLUCZOWE 1: *zarzewieć*  
ODPOWIEDŹ 2: *zamek (część broni)*  
WAGA 2: 0.500000

SŁOWA KLUCZOWE 2:        *zardzewieć*

[5]

ZDANIE:                    „Przeżyliśmy prawdziwą plagę ślimaków.”

ODPOWIEDŹ:            (*brak wyników*)

SŁOWA KLUCZOWE: (*brak wyników*)

[6]

ZDANIE:                    „Znalazłem wczoraj dużo kozaków.”

ODPOWIEDŹ 1:        *kozak (żołnierz lekkiej jazdy)*

WAGA 1:                    *0.500000*

SŁOWA KLUCZOWE 1:    *kozak (grzyb jadalny)*

ODPOWIEDŹ 2:        *kozak (grzyb jadalny)*

WAGA 2:                    *0.500000*

SŁOWA KLUCZOWE 2:    *kozak (żołnierz lekkiej jazdy)*

### 5.2.2. ANALIZA

Zgodnie z przewidywaniami, wypowiedzi języka potocznego są ubogie w słowa, co czyni kontekst trudnym, a czasem niemożliwym do określenia. Należy jednak zauważyć, że bez znajomości kontekstu, często jest to problem niemożliwy do rozwiązania także dla człowieka. Otrzymane wyniki wskazują również na braki pojęć w modelu i relacji między nimi.

## 5.3. TEKSTY DZIEŁ LITERACKICH

Test został przeprowadzony na kilku fragmentach zaczerpniętych z tekstów literackich.

### 5.3.1. ZBIÓR TESTOWY I WYNIKI

Fragmenty książek wchodzące w zbiór testowy:

[1]

**FRAGMENT:** „Ogniem i mieczem” – H. Sienkiewicz  
*Namiestnik jednak nie dostał się tego wieczora do zamku, bo pan Grodzicki zaprowadził taki porządek, że gdy przed zachodem słońca wybito hasło, nie wpuszczano nikogo z zamku i do zamku i gdyby nawet sam król przyjechał, musiałby nocować w Słobódce stojącej pod wałami fortecy.*

**ODPOWIEDŹ:** zamek (warowna budowla mieszkalna)

**SŁOWA KLUCZOWE:** forteca

[2]

**FRAGMENT:** „Ogniem i mieczem” – H. Sienkiewicz  
*Forteca była istotnie nie do zdobycia, bo pròcz armat broniły jej Dnieprowe przepaście i niedostępne skały pionowo zeskakujące w wodę; nie potrzebowała nawet wielkiej załogi. Toteż w zamku nie stało więcej nad sześćset ludzi, ale za to co najprzebrańszego żołnierza, uzbrojonego w muszkiety i samopały. Dniepr, płynąc w tym miejscu ściśniętym korytem, tak był wąski, że rzucona z wałów strzała przelatowała daleko na drugi brzeg. Działa zamkowe panowały nad oboma brzegami i nad całą okolicą. Pròcz tego o pòł mili od zamku stała wysoka wieża, z której osiem mil wokoło widać było, a w niej stu żołnierzy, do których pan Grodzicki każdego dnia zaglądał. Ci, spostrzegłszy w okolicy lud jaki, dawali natychmiast znać do zamku, a wòwczas bito w dzwony i cała załoga wnet stawiała pod bronią.*

**ODPOWIEDŹ 1:** zamek (warowna budowla mieszkalna)

**WAGA 1:** 0.750000

**SŁOWA KLUCZOWE 1:** człowiek, okolica, stać

**ODPOWIEDŹ 2:** zamek (część broni)

**WAGA 2:** 0.250000

**SŁOWA KLUCZOWE 2:** [nie działać, działać, działa]

[3]



FRAGMENT: „Felicitas” – A. Zajdel  
Słyszac trzask zamka w przedpokoju, przymknąłem szybko wieko walizki i odskoczyłem w stronę biurka.

ODPOWIEDŹ: zamek (urządzenie do zamykania)

SŁOWA KLUCZOWE: walizka

[4]

FRAGMENT: „Ferma” – A. Zajdel  
To nowe serce... Okazało się, że ono nie jest w ogóle moje. Ono nie było wszczepione, tylko jakby... zainstalowane. Można je było odłączyć w każdej chwili. Cięcia operacyjne też nie zostały zszyte, tylko jakby... jakby zapięte na zamek błyskawiczny, rozumiesz? Skóra, osierdzie wszystko spojone tak, że w każdej chwili można otworzyć, jak kieszeń!

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[5]

FRAGMENT: „Ferma” – A. Zajdel  
- Nie mogę. Zamykają nas o zmroku, nie mam klucza.  
- Ach, tak... Zaczekaj, spróbuję otworzyć. Sven sięgnął do torby, wyjął kilka drobnych narzędzi i przez chwilę manipulował przy zamku. Bez trudu udało mu się sforsować prosty zamek.

ODPOWIEDŹ: zamek (urządzenie do zamykania)

SŁOWA KLUCZOWE: [zamykać, nie zamykać]

[6]

FRAGMENT: „Gra w zielone” – A. Zajdel  
W pewnej chwili, odebrawszy telefon, profesor przeprosił mnie i wyszedł zapewniając, że niedługo wróci. Rozejrzałem się po gabinecie. W głębi, w ścianie, zauważyłem skrytkę, jakby kasę pancerną. Była nie domknięta, pęk kluczy zwisał z zamka sejfu.

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[7]

FRAGMENT: „Bajka o królu Murdasie” – S. Lem  
Szeregi elektrykarzy i budowniczych jeły wnosić do zamku druty i szpule, a

gdy rozbudowany król wypełnił swoją osobą cały pałac, tak że był jednocześnie z frontu, w piwnicach i oficynie, przyszła kolej na stojące w pobliżu domostwa.

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[8]

FRAGMENT: „Czarna kawa” – A. Christie  
- Doskonale, Tredwell. Możesz teraz zamknąć.  
- Tak, proszę pana - odpowiedział lokaj, wychodząc. Zamknął za sobą drzwi i wszyscy usłyszeli odgłos klucza przekręcanego w zamku.

ODPOWIEDŹ: zamek (urządzenie do zamykania)

SŁOWA KLUCZOWE: drzwi

[9]

FRAGMENT: „Szatański wirus” – A. MacLean  
Odwróciłem się plecami do drzwi, tak długo skrobałem po nich lufą, aż trafiłem na zamek; i pociągnąłem za spust. Po drugim wystrzale drzwi nagle ustąpiły pod moim ciężarem, co może wiele powiedzieć o stanie mojego umysłu, bo nawet nie obejrzałem przedtem zawiasów, by sprawdzić, w którą stronę się otwierają do środka czy na zewnątrz.

ODPOWIEDŹ 1: zamek (urządzenie do zamykania)

WAGA 1: 0.600000

SŁOWA KLUCZOWE 1: drzwi, spust, otwierać

ODPOWIEDŹ 2: zamek (część broni)

WAGA 2: 0.200000

SŁOWA KLUCZOWE 2: lufa

ODPOWIEDŹ 3: zamek (warowna budowla mieszkalna)

WAGA 3: 0.200000

SŁOWA KLUCZOWE 3: stanie

[10]

FRAGMENT: „Szatański wirus” – A. MacLean  
Drzwi na górze również były zamknięte na klucz, ale to nie moje drzwi i miałem jeszcze pięć nabojęw w magazynku. Zamek ustąpił po pierwszym strzale i wytoczyłem się do hallu.

ODPOWIEDŹ 1: zamek (część broni)

WAGA 1: 0.500000

SŁOWA KLUCZOWE 1: magazynek, nabój

ODPOWIEDŹ 2: zamek (urządzenie do zamykania)  
WAGA 2: 0.500000  
SŁOWA KLUCZOWE 2: drzwi

[11]

FRAGMENT: „Elfia krew” – A. Norton  
- Muszę przynieść swoje narzędzia.  
Podniosła się i zniknęła w prywatnych pomieszczeniach, skąd wróciła w oka mgnieniu ze skórzaną saszetką pełną precyzyjnych narzędzi, których kobiety-kowale używają przy tworzeniu biżuterii.  
- Ten zamek jest bardzo stary i delikatny - powiedziała, sadowiąc się obok Shany, po czym otworzyła saszetkę i wyjęła zestaw do próbek. - Sądząc po jego wyglądzie, został wykonany przez kobietę. Nie jest może prosty, lecz sama robiłam bardziej skomplikowane.

ODPOWIEDŹ: (brak wyników)  
SŁOWA KLUCZOWE: (brak wyników)

[12]

FRAGMENT: „Bóg Imperator Diuny” – F. Herbert  
Doszedł do drzwi swojej kwatery, zbliżył dłoń do zamka zaprogramowanego na jego linie papilarne i zawahał się. Czuł się jak zaszczute zwierzę uciekające do swojego legowiska.

ODPOWIEDŹ: zamek (urządzenie do zamykania)  
SŁOWA KLUCZOWE: drzwi

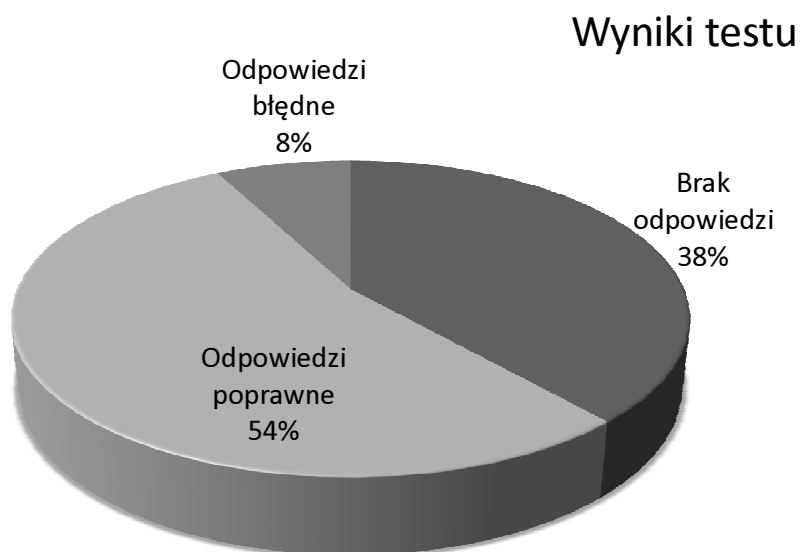
[13]

FRAGMENT: „Koniec wieczności” – I. Asimov  
Błyskawicznie rzucił się do swoich akt. Gdy palcem jednej ręki przyciskał szyfrowy zamek szafki, druga sięgnęła do teczki. Z biurka wysunął się srebrny język folii, jego perforacja była widoczna nawet gołym okiem.

ODPOWIEDŹ: (brak wyników)  
SŁOWA KLUCZOWE: (brak wyników)

### 5.3.2. ANALIZA

Dla zdecydowanej większości fragmentów, dla których program był w stanie udzielić jakakolwiek odpowiedź – była ona poprawna. Tylko w jednym przypadku rozwiązanie udzielone przez program nie było prawidłowe, wskazano bowiem dwie równie prawdopodobne odpowiedzi. Dużym problemem w działaniu programu jest znaczna liczba przypadków, w których program nie jest w stanie wskazać żadnego rozwiązania. Przyczyną takiego stanu jest brak wiedzy w modelu, w słowniku nie istnieją pojęcia lub relacje występujące w wypowiedziach. Jedynym rozwiązaniem jest uzupełnienie słownika o wymagane zależności.



**Rys. 6.** Dzieła literackie – wyniki testów

Świadczy to o brakach w modelu – w słowniku nie istnieją pojęcia lub relacje występujące w wypowiedziach. Jedynym rozwiązaniem jest uzupełnienie słownika o wymagane zależności.

## 5.4. NOTATKI PAP

Test został przeprowadzony na wyselekcjonowanym podzbiorze notatek prasowych PAP.

### 5.4.1. ZBIÓR TESTOWY I WYNIKI

Zbiór testowy stanowiły wybrane notatki prasowe, dotyczące pojęć wieloznacznych:

- *zamek* w znaczeniu:
  - urządzenia do zamykania
  - warownej budowli mieszkalnej
  - części broni
- *ślimak* jako:
  - część ucha wewnętrznego ssaków
  - zwierzę z rodziny ślimaków
  - motyw dekoracyjny
  - element maszynowy w postaci śruby
  - droga w formie linii spiralnej
- *kozak* w rozumieniu:
  - żołnierz lekkiej jazdy
  - grzyb jadalny
  - ukraiński taniec ludowy
  - śmiały, odważny mężczyzna
  - ciepłe, damskie lub dziecięce obuwie zimowe o wysokiej, miękkiej cholewie

Poniżej przedstawione zostały notatki testowe wraz z wynikami zwróconymi przez program:

[1]

NOTATKA #007516:

*W Centrum Sztuki Współczesnej Zamek Ujazdowski w Warszawie otwarto wystawę prac młodego polskiego artysty Dominika Lejmana, zatytułowaną: "Luksus przetrwania".*

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[2]

NOTATKA #016750:

Pozostałości nieznanego dotychczas średniowiecznego zamku rycerskiego odkryli śląscy archeolodzy we wsi Giebło k. Ogrodzieńca. Na fragmenty starych murów natrafiono podczas robót ziemnych na prywatnej posesji. "To prawdziwa rewelacja. Wydawało się, że wszystkie dawne zamki jurajskie są już znane, tymczasem natrafiliśmy na obiekt, o którym nikt nie miał pojęcia" - powiedział PAP Jacek Pierzak, konserwator zabytków archeologicznych w Katowicach. Na razie odstonięto jedynie fragmenty fundamentów budowli. "Na pewno budynek miał charakter rezydencji rycerskiej, być może była to wieża mieszkalno-obronna. Spodziewamy się znaleźć w pobliżu również pozostałości zabudowań gospodarczych oraz obronnych, np. wału czy fosy" - oświadczył Pierzak. Archeolodzy ustalają teraz, z jakiego dokładnie okresu pochodzi budowla. Zgromadzony dotychczas materiał archeologiczny - głównie ceramika - wskazuje na XIV wiek. Niewykluczone jednak, że zamek postawiono już w wieku XIII. O tym, że zamek był posiadłością rycerską, świadczy m.in. fakt, że w pobliżu znajduje się - jeden z trzech w woj. Śląskim - kościołów romańskich." Fundatorami takich niewielkich wiejskich kościołów w tym okresie byli zwykle zamieszkałi w pobliżu rycerze, wewnątrz umieszczano dla nich nawet specjalne łóża" - wytłumaczył Pierzak. W Jurze Krakowsko-Częstochowskiej znajdują się pozostałości ponad 20 średniowiecznych zamków. Najstynniejsze są ruiny w Ogrodzieńcu, oddalonym od Giebła o kilka kilometrów.

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: budowla, obronny, rycerski, średniowieczny, zabytek

[3]

NOTATKA #016166:

Ponad 140 osób z 27 miast Polski przyjedzie do Krakowa na VI Krajowy Zjazd Bractwa Młodych Miłośników Starych Miast. Spotkanie rozpocznie się w piątek Drużyny z 18 miast będą rywalizować w Ogólnopolskim Konkursie Historycznym "Wawel - wszystkich Polaków miejsce najgodniejsze". Jedna trzecia konkursowych pytań ma dotyczyć Biskupstwa Krakowskiego, które w tym roku świętuje jubileusz 1000-lecia, a pozostałe zadania - historii, architektury i zbiorów Zamku Królewskiego na Wawelu. W finale zmierzą się tylko trzy drużyny, które wędrując po Wzgórzu Wawelskim będą szukać odpowiedzi na pytania jurorów. Zwycięzcy konkursu pojadą do Rzymu.

Bractwo Młodych Miłośników Starych Miast powstało w styczniu 1995 r. Jego głównym celem jest upowszechnianie wśród młodzieży wiedzy o zabytkach i historii miast, w których mieszkają. Co roku w innym miejscu Polski organizowane są zjazdy bractwa. Współorganizatorami zjazdu są: Zamek Królewski na Wawelu, Katedra Krakowska, Muzeum Narodowe, Dom Jana Matejki, Muzeum Historyczne Miasta Krakowa, Kopalnia Soli w Wieliczce. Honorowy patronat objął metropolita krakowski ks. kardynał Franciszek Macharski.

**ODPOWIEDŹ:** zamek (warowna budowla mieszkalna)  
**SŁOWA KLUCZOWE:** cel, dom, kròlewski, mieszkać (zajmować na stałe jakieś pomieszczenie), muzeum, zabytek

[4]

NOTATKA #026921:

Gwiazda popu 42-letnia Madonna wyszła w piątek wieczorem za mąż za 32-letniego brytyjskiego reżysera filmowego Guya Ritchie w szkockim zamku Skibo – oświadczyła w sobotę rano miejscowa kobieta-pastor Susan Brown, która udzieliła ślubu. "To się stało" – zapewniła pani Brown dziennikarzy zgromadzonych przed jej domem. Potwierdzenie to ogłoszone zostało po 12 godzinach oczekiwania przedstawicieli mediów, którzy do soboty rano nie byli pewni, czy w piątek wieczorem ślub Madonny z Ritchie rzeczywiście się odbył.

Ceremonia ślubna zorganizowana została tak, aby nie miało do niej dostępu kilkuset dziennikarzy, którzy przybyli do malowniczej miejscowości Dornoch w północnej Szkocji, gdzie znajduje się zbudowany w końcu XIX wieku ekskluzywny zamek Skibo. Ślubu udzieliła ta sama pastor-kobieta, która dzień wcześniej ochrzciła czteromiesięcznego Rocco – syna Madonny i Ritchie.

**ODPOWIEDŹ:** zamek (warowna budowla mieszkalna)

**SŁOWA KLUCZOWE:** dom (bliskie sercu miejsce), stać

[5]

NOTATKA #041103:

Wieczorem, w ramach działań prewencyjnych policjanci wysadzili zamek w drzwiach jednego z samochodów zaparkowanych koło hotelu Sheraton w Warszawie, w którym zatrzymał się Kasjanow. Policyjny pies sygnalizował, że w aucie może być bomba, ale jej nie było. Wg Aleksieja Wolina, członka delegacji rosyjskiej, Kasjanow ze spokojem odniósł się do incydentu. "Obecnie też Michaił Kasjanow nie ma zamiaru wprowadzać żadnych zmian do programu swej wizyty w Polsce" – podkreślił Wolin.

**ODPOWIEDŹ:** zamek (urządzenie do zamykania)

**SŁOWA KLUCZOWE:** drzwi

[6]

NOTATKA #041416:

W Ujeździe koło Sandomierza (woj. świętokrzyskie) odbył się dwudniowy V Turniej Rycerski o szablę Krzysztofa Baldwina Ossolińskiego. Jednym z głównych punktów imprezy była rekonstrukcja bitwy o zamek Krzyżtopór.

**ODPOWIEDŹ:** zamek (warowna budowla mieszkalna)

**SŁOWA KLUCZOWE:** rycerski

[7]

NOTATKA #044864:

*Zamek Kròlewski na Wawelu od wtorku udostępnia zwiedzającym dwie Komnaty na pierwszym pięttrze północnego skrzydła zamku, w których prezentowana jest kolekcja XVIII-wiecznych mebli, portretów, dywanów, porcelany i sreber. W obu salach przez ostatnie 10 lat mieściły się magazyny.*

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: kròlewski, zwiedzający

[8]

NOTATKA #049001:

*Kętrzyn (woj. warmińsko-mazurskie) najechało w czwartek 200 rycerzy i żołnierzy z kraju i zagranicy. Rozpoczął się tam II Mazurski Turniej Rycerski. Imprezy potrwałją do niedzieli, a zakończą je najazd na zamek w Barcianach. Na turniej rycerski do Kętrzyna przyjechali członkowie bractw polskich i Czesi prezentujący wojsko z czasów cesarstwa austro-węgierskiego za panowania cesarza Franciszka Józef.*

ODPOWIEDŹ: zamek (warowna budowla mieszkalna)

SŁOWA KLUCZOWE: rycerski

[9]

NOTATKA #020725:

*Setki tysięcy ślimaków, które wypęłżyły na płytę lotniska w Nicei sprawiły, że przez pięć godzin jeden z dwóch pasów startowych nie przyjmował samolotów. Plagę ślimaków wykryto podczas rutynowej kontroli pasa. Najprawdopodobniej skorupiaki przedostały się w nocy na betonową nawierzchnię, uciekając przed ulewnym deszczem. Ślimaczy śluz, którym pokryty był trzykilometrowy pas, wykluczał starty i lądowania. Zarząd lotniska zdecydował się na rozwiązanie radykalne – nieproszonych gości pozamiatano, a pas umyto. Do najbliższej składnicy odpadów odstawiono około pięciu metrów sześciennych ślimaków.*

ODPOWIEDŹ: ślimak (zwierzę z rodziny ślimaków)

SŁOWA KLUCZOWE: uciekać przed

[10]

NOTATKA #028749:

*Kilkaset eksponatów przyrodniczych i ekologicznych wystąło do Muzeum Narodowego w Szczecinie dwoje naukowców, uczestników wyprawy jachtem "Maria" dookoła świata "Szlakiem Magellana". Wśród eksponatów są m.in. zbiory krabów, ślimaków, małży, koralowców, gąbek, rozgwiazd z wysp na Pacyfiku znalezionych u wybrzeży Chile.*

ODPOWIEDŹ: (brak wyników)



SŁOWA KLUCZOWE: (brak wyników)

[11]

NOTATKA #030828:

*Oddział dziecięcy ośrodka rehabilitacji laryngologicznej Akademii Medycznej otwarto w Poznaniu. Ośrodek opiekuje się m.in. pacjentami z wszczepami implantów ślimakowych.*

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[12]

NOTATKA #036992:

*Powodzeniem zakończyła się operacja wszczepienia implantu pniowego do mózgu 51-letniego poznaniaka. W piątek – 24 godziny po zabiegu – pacjent został wybudzony. Poznaniak utracił słuch na skutek niezłośliwego guza na obu nerwach słuchowych. Z tego powodu nie można było pacjentowi zastosować wszczepienia implantu ślimakowego. Operacja polegała na usunięciu niezłośliwych guzków i wszczepieniu implantu pniowego z 21 elektrodami do zachyłka bocznego komory w IV pniu mózgu.*

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[13]

NOTATKA #037915:

*Na Warmii i Mazurach rozpoczął się zbiór ślimaków. Zbieranie ślimaków w województwie o największej stopie bezrobocia to obok zbierania poroży czy grzybów i jagód jeden ze sposobów zarobkowania, zwłaszcza dla mieszkańców na terenach popegeerowskich.*

ODPOWIEDŹ: (brak wyników)

SŁOWA KLUCZOWE: (brak wyników)

[14]

NOTATKA #038966:

*Rozporządzenie wojewody warmińsko-mazurskiego w sprawie zbioru ślimaków – winniczków obowiązuje w tym roku od 19 kwietnia do 15 maja.*

*Jak powiedziała PAP Maria Mellin - Wojewódzki Konserwator Przyrody, ślimaki można zbierać wszędzie z wyjątkiem parków krajobrazowych, rezerwatów oraz strefy przygranicznej. Nie wolno zbliżyć się do granicy na mniej niż 15 metrów. Ślimaki muszą mieć co najmniej 30 mm długości.*

**ODPOWIEDŹ:** ślimak (zwierzę z rodziny ślimaków)

**SŁOWA KLUCZOWE:** zbierać

[15]

**NOTATKA #011939:**

*Pierwszą w Polsce operację wszczepienia drugiego implantu ślimaka ucha wewnętrznego przeprowadzono w Poznaniu. Operacja wszczepienia powiodła się, implant działa. Za miesiąc pacjent odzyska słuch w obu uszach.*

**ODPOWIEDŹ 1:** ślimak (zwierzę z rodziny ślimaków)

**WAGA 1:** 0.500000

**SŁOWA KLUCZOWE 1:** ucho (narząd słuchu)

**ODPOWIEDŹ 2:** ślimak (składnik ucha wewnętrznego ssaków)

**WAGA 2:** 0.500000

**SŁOWA KLUCZOWE 2:** działać

[16]

**NOTATKA #012846:**

*W lasach Gór Świętokrzyskich rozpoczął się – nie notowany – zazwyczaj w lipcu – obfity wysyp grzybów jadalnych: borowików, czerwonych kozaków, maślaków oraz kani.*

**ODPOWIEDŹ:** kozak (grzyb jadalny)

**SŁOWA KLUCZOWE:** borowik, czerwony, grzyb jadalny, las

[17]

**NOTATKA #019167:**

*Jeśli uwielbiamy maślaka żółtego, szukajmy go wyłącznie pod modrzewiem. Wielbiciele kozaka brunatnego - kierujcie się pod brzozę, a czerwonego - pod osikę. Tajemnicą tej współpracy jest symbioza grzybów i roślin - mikoryza. Mikoryza to symbiotyczna współżycie korzeni roślin ze strzępkami grzybów, podczas której z gleby do rośliny za pośrednictwem grzyba wędrują związki mineralne, a z korzeni do grzybni transportowane są produkty fotosyntezy. Niektóre grzyby wiążą się tylko z jednym drzewem – np. kozak brunatny z brzozą, czerwony – z osiką, a kurka – z sosną. Inne żyją w symbiozie z wieloma gatunkami drzew. I tak olszówkę, czyli krowiaka podwiniętego, możemy znaleźć pod olchą, sosną, świerkiem i brzozą.*

**ODPOWIEDŹ 1:** kozak (grzyb jadalny)

**WAGA 1:** 0.600000

**SŁOWA KLUCZOWE 1:** czerwony, roślina (pojęcie kategoriale), życie

**ODPOWIEDŹ 2:** kozak (żołnierz lekkiej jazdy)

**WAGA 2:** 0.200000

SŁOWA KLUCZOWE 2:       żyć  
ODPOWIEDŹ 3:       kozak (śmiały, odważny mężczyzna)  
WAGA 3:               0.200000  
SŁOWA KLUCZOWE 3:       żyć

[18]

NOTATKA #033110:

*Światowej sławy irlandzki spektakl "Spirit of the Dance" będzie mogła obejrzeć polska publiczność podczas jedyne koncertu w Polsce – 22 marca w Spodku w Katowicach. Eklektyczne widowisko, przygotowane przez Irlandzki Międzynarodowy Zespół Tańca, łączy wiele gatunków tanecznych: irlandzkie stepowanie, flamenco, klasyczny balet, taniec ukraińskich kozaków.*

ODPOWIEDŹ 1:       kozak (ludowy taniec ukraiński)  
WAGA 1:               0.666667  
SŁOWA KLUCZOWE 1:       taniec  
ODPOWIEDŹ 2:       kozak (śmiały, odważny mężczyzna)  
WAGA 2:               0.333333  
SŁOWA KLUCZOWE 2:       zespół

[19]

NOTATKA #047864:

*W piątek w Warszawie złożono akt erekcyjny pod budowę pomnika Tarasa Szewczenki. Stanisław Żurowski, podsekretarz stanu w Ministerstwie Kultury oraz pełnomocnik rządu ds. polskiego dziedzictwa kulturowego za granicą, uznał piątkową uroczystość za wydarzenie "o znaczeniu symbolicznym". Budowę pomnika określił jako "podwaliny pod wzajemne relacje" między Polską a Ukrainą. Pomnik Tarasa Szewczenki, wyrzeźbiony przez artystę z Kijowa, Anatolija Kuszczę, będzie miał 2,8 m wysokości. Stanie na ponad dwumetrowym cokole, zaprojektowanym przez architekta Baltazara Brukalskiego. Na pomniku ma znaleźć się napis: "Podajże rękę Kozakowi i serce swe ku niemu przychyl".*

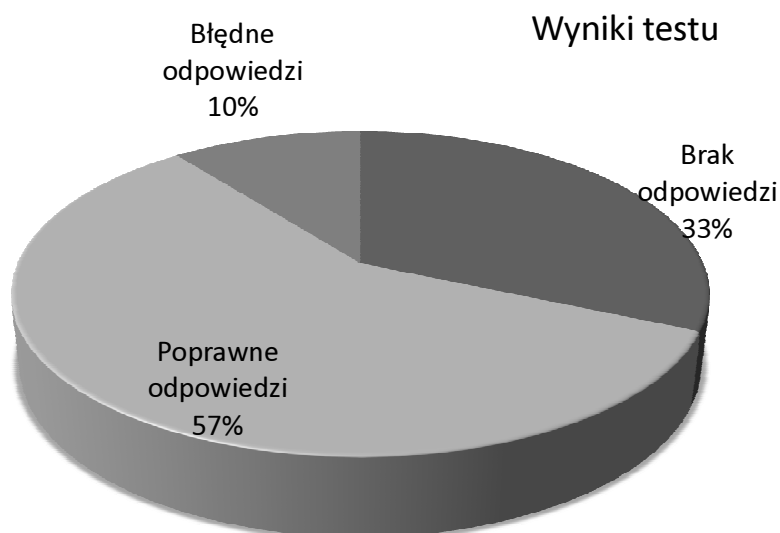
ODPOWIEDŹ 1:       kozak (ludowy taniec ukraiński)  
WAGA 1:               0.400000  
SŁOWA KLUCZOWE 1:       kultura  
ODPOWIEDŹ 2:       kozak (żołnierz lekkiej jazdy)  
WAGA 2:               0.200000  
SŁOWA KLUCZOWE 2:       mieć  
ODPOWIEDŹ 3:       kozak (grzyb jadalny)  
WAGA 3:               0.200000

SŁOWA KLUCZOWE 3:       *mieć*  
ODPOWIEDŹ 4:       *kozak (śmiały, odważny mężczyzna)*  
WAGA 4:               0.200000  
SŁOWA KLUCZOWE 4:       *mieć*

#### 5.4.2. ANALIZA

Podobnie jak dla fragmentów tekstów literackich, jeśli jakkolwiek wynik został zwrócony – w zdecydowanej większości przypadków był on poprawny. Problemem jednak znowu jest duża ilość wypowiedzi, dla których nie został zwrócony żaden wynik, co potwierdza istnienie sporych braków w modelu. Zwrócone dwie błędne wypowiedzi są wynikiem istnienia niepoprawnej relacji w słowniku.

Notatki [11], [12] i [13] potwierdzają istnienie problemu omówionego w punkcie 5.1.2, a mianowicie braku zdefiniowanych w słowniku semantycznym zależności pomiędzy różnymi formami tego samego wyrazu. W tym wypadku, pomiędzy występującym w słowniku słowem *ślimak*, a wyrazami: *ślimakowych* lub *ślimakowego* pochodzących od słowa *ślimakowy* (nieistniejącego w słowniku), oraz *zbierać* a *zbieranie*.



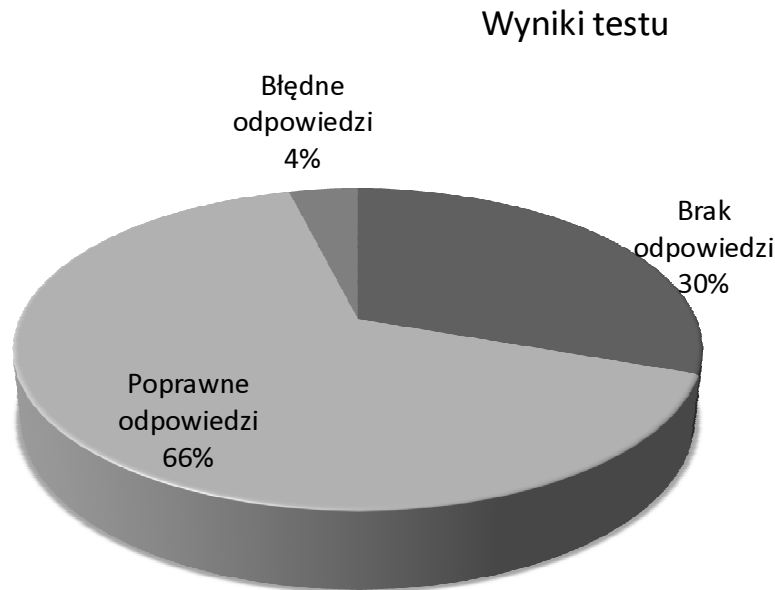
**Rys. 7.** Wybrane notatki PAP – wyniki testów

## 5.5. WSZYSTKIE NOTATKI ZAWIERAJĄCE WYBRANE SŁOWO

Test polegał na sprawdzeniu skuteczności działania algorytmu, dla wszystkich notatek PAP zawierających którąś z form słowa *zamek*. Po przeprowadzeniu testów, porównano najbardziej prawdopodobne według programu wyniki, z poprawnymi odpowiedziami.

### 5.5.1. ZBIÓR TESTOWY I WYNIKI

Zbiór testowy stanowiły wszystkie notatki PAP zawierające słowo *zamek*, czyli 181 notatek.



**Rys. 8.** Wszystkie notatki PAP – wyniki testów

### 5.5.2. ANALIZA

Dla 181 notatek, znaczenia 119 słów zostały przyporządkowane poprawnie. Spośród błędnie sklasyfikowanych notatek, dla 7 wskazano błędne odpowiedzi, a dla pozostałych 55 wypowiedzi nie zaproponowano żadnej odpowiedzi, co świadczy o brakach odpowiednich relacji czy pojęć w słowniku semantycznym, lub zbyt dużej objętości testowanych notatek.

Szczególnie często mylone było słowo *zamek* w znaczeniu *warownej budowli mieszkalnej*, z częścią *broni*. Powodowane jest to zbyt obszernym kontekstem. W celu zapewnienia poprawnego działania algorytmu, konieczne byłoby wcześniejsze wykonanie analizy zdań w celu określenia faktycznych granic wypowiedzi.

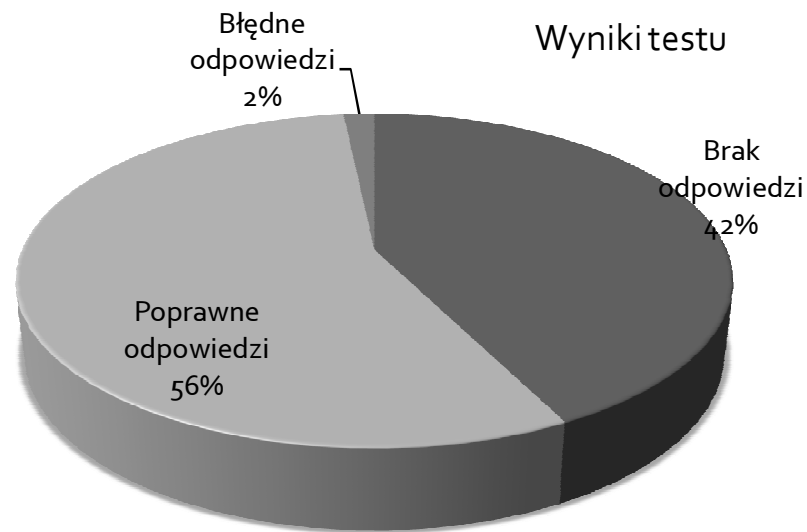
## 5.6. OGRANICZONY ZBIÓR RELACJI

Test polegał na sprawdzeniu działania algorytmu dla zbioru relacji ograniczonego jedynie do relacji paradygmatycznych (zgodnie z tabelą 1), a następnie porównaniu otrzymanych wyników, z uzyskanymi w poprzednich testach (czyli dla pełnego zbioru relacji).

### 5.6.1. ZBIÓR TESTOWY I WYNIKI

Testowi poddano wszystkie notatki PAP zawierające słowo *zamek*, czyli 181 notatek (podobnie jak dla testu nr 5).

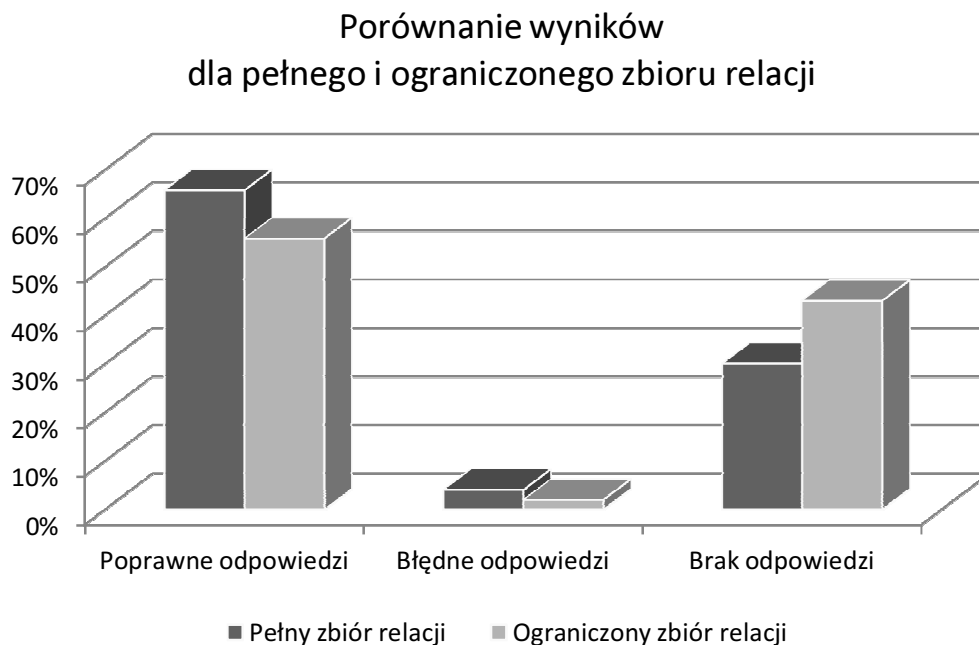
Otrzymane wyniki przedstawiają się następująco:



**Rys. 9.** Ograniczony zbiór relacji – wyniki testów

## 5.6.2. ANALIZA

Okazało się, że dla większości zdań ograniczenie zbioru relacji jedynie do paradygmatycznych skutkowało obniżeniem skuteczności działania algorytmu. Znaczenie zostało określone dla 105 z 181 zdań (poprzednio było to 126 znaczeń), czyli 56%. Dla 42%, czyli 76 notatek nie znaleziono rozwiązania w ogóle.



**Rys. 10.** Ograniczony zbiór relacji – porównanie wyników

Pogorszenie skuteczności działania programu spowodowane jest tym, że zależności paradygmatyczne rzadziej występują w zdaniach języka naturalnego. Prawidłowe odpowiedzi mogły zatem zostać udzielone jedynie dla wypowiedzi typu definicyjnego i takich, w których w jakiś sposób przemycone zostały treści tego typu.

## 5.7. SZTUCZNIE WYGENEROWANE CIĄGI SŁÓW

Test został wykonany na zbiorach słów sztucznie wygenerowanych z bazy słownika semantycznego i miał zweryfikować działanie algorytmu dla istniejących danych, a więc w oderwaniu od problemu kompletności słownika.

### 5.7.1. ZBIÓR TESTOWY I WYNIKI

Jako zbiór testowy wybrano zestawy trzech pojęć połączonych ze sobą relacjami. W każdym z zestawów, przynajmniej jedno słowo było jednym z trzech zawartych w słowniku semantycznym znaczeń słowa *zamek*, czyli:

- urządzenie do zamykania
- część broni
- warowna budowla mieszkalna

Testy przeprowadzono dla 450 spreparowanych zestawów słów, jak na przykład:

[1] *zamek* – IS A PART OF – *drzwi* – IS A KIND OF – *obiekt fizyczny* [*zamek* jako urządzenie do zamykania]

[2] *broń* – IS A – *kusza* – CONSISTS OF – *zamek* [*zamek* jako część broni]

[3] *pistolet* – CONSISTS OF – *zamek* – STATE NEGATIVE – *zepsuty* [*zamek* jako część broni]

[4] *dom* – IS A – *zamek* – CONSISTS OF – *budynki mieszkalne* [*zamek* jako warowna budowla mieszkalna]

[5] *zamek* – SIMILAR TO – *twierdza* – STATE NEGATIVE – *ruiny* [*zamek* jako warowna budowla mieszkalna]

Dla każdego z zestawów zwracano najbardziej prawdopodobną odpowiedź – istnienie dwóch, równie prawdopodobnych odpowiedzi traktowano jako jej brak.



### 5.7.2. ANALIZA

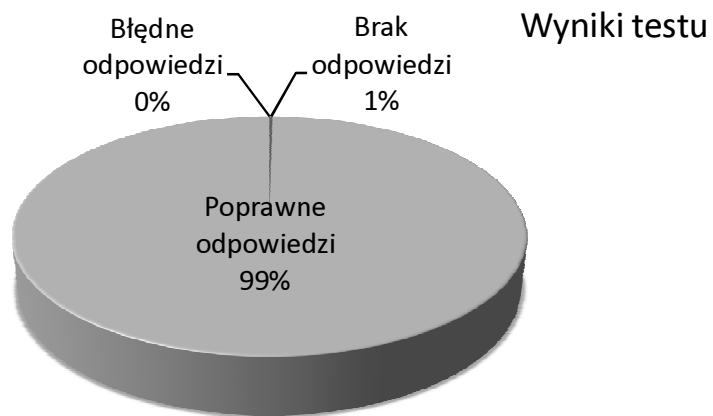
Jak się okazało, jedynie dla 4 z 450 zestawów słów nie zwrócono odpowiedzi. Nastąpiło to dla poniższych zestawów:

[1] Hohenzollern – IS A kind Of – zamek – CATEGORY – obiekt fizyczny [zamek jako warowna budowla mieszkalna]

[2] Hohenzollern – IS A kind Of – zamek – consists of – budynki mieszkalne [zamek jako warowna budowla mieszkalna]

[3] Hohenzollern – IS A kind Of – zamek – consists of – gospodarcze [zamek jako warowna budowla mieszkalna]

[4] Hohenzollern – IS A kind Of – zamek – action negative rt – niechronienie [zamek jako warowna budowla mieszkalna]



**Rys. 11.** Sztucznie wygenerowane ciągi słów – wyniki testów

Dla wszystkich zdań dla których algorytm nie zwrócił odpowiedzi, pierwsze ze słów, czyli *Hohenzollern* nie przeszło etapu ujednoznaczniania fleksyjnego, zatem nie było brane pod uwagę w trakcie procesu ujednoznaczniania semantycznego.

Dla testu [1] pozostałe ze słów, czyli: *obiekt fizyczny*, powiązane było z każdym ze znaczeń słowa *zamek*, dlatego nie można było rozstrzygnąć, o które znaczenie faktycznie chodziło.

W testach [2] i [3] słowa przechowywane w słowniku semantycznym nie są w formie podstawowej, co uniemożliwiło ich zidentyfikowanie.

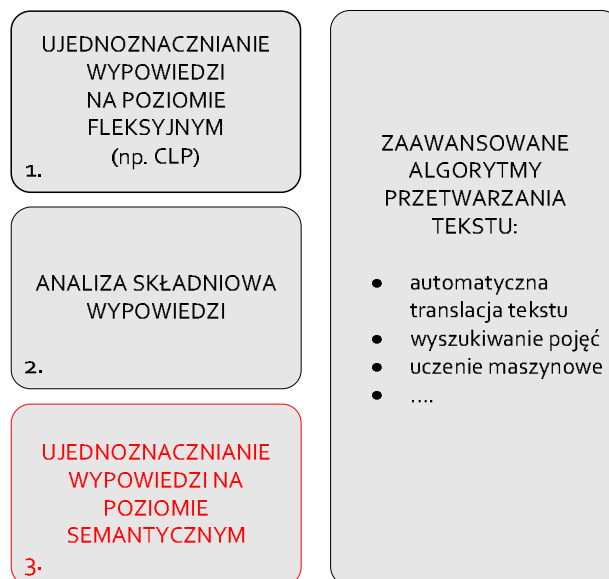
W teście [4] słowo niechronienie nie przeszło etapu ujednoznaczniania fleksyjnego, więc nie mogło pomóc określić właściwego znaczenia wyrazu *zamek*.

Można zatem powiedzieć, że skuteczność algorytmu jest duża i zidentyfikowane błędy wynikają z niekompletności słownika lub pomyłek w nim znalezionych.

## ROZDZIAŁ VI

### PODSUMOWANIE

Celem pracy było stworzenie i przetestowanie na wybranym zbiorze tekstów algorytmu rozstrzygającego wieloznaczność wyrazu za pomocą relacyjnego opisu znaczenia. Postawione zadanie zostało zrealizowane. W ramach pracy sprecyzowany został zakres stawianych programowi wymagań. W następnym kroku powstał teoretyczny opis algorytmu, który następnie posłużył jako punkt wyjściowy do implementacji zaproponowanego rozwiązania. Jako bazę wiedzy dla rozwiązania użyto słownika semantycznego stworzonego przez Katedrę Lingwistyki Komputerowej Uniwersytetu Jagiellońskiego, dla którego wykorzystania stworzony został odpowiedni importer danych. W kolejnym kroku został zaproponowany zestaw testów i procedur ich uruchamiania. Na ich podstawie określono skuteczność algorytmu i jego zachowanie.



**Rys. 12.** Miejsce opracowanego modułu w algorytmach przetwarzania języka

Opracowany i zaimplementowany algorytm stanowi element pośredni etapu analizy tekstu. Jego istotą jest zmierzenie się z integralną cechą języka naturalnego jakim jest wieloznaczność wyrazu w wypowiedzi. Wykorzystane go wraz z modułami analizy fleksyjnej i składniowej w bardziej złożonych algorytmach przetwarzania tekstu, pozwala na odseparowanie w postaci niezależnego narzędzia informacji o wieloznaczności pojęć w wypowiedzi.

Po analizie odpowiedzi generowanych w odpowiedzi na zdania testowe można zauważyć, że skuteczność działania stworzonego algorytmu jest silnie zależna od analizowanej wypowiedzi (kontekstu, jaki dostarcza) oraz ilości i jakości pojęć zgromadzonych w słowniku semantycznym, które dostarczają informacji o kontekście słów w zdaniach. Słownik Semantyczny Katedry Lingwistyki Komputerowej UJ znajduje się w fazie tworzenia i w tym momencie zawiera niecałe 52000 pojęć języka polskiego, z których wiele wymaga poprawy. W porównaniu do ogółu pojęć języka polskiego, jest to stosunkowo niewielka liczba, jednak wystarczająca do przetestowania skuteczności opracowanego algorytmu. Uzyskanie dobrych wyników działania już na tym etapie prac napawa optymizmem, należy jednak podkreślić, że ewentualne zastosowanie rozwiązania na większą skalę wymaga pogłębionego wysiłku i nakładu pracy.

Pomimo, że otrzymane wyniki wskazują na wysoką skuteczność algorytmu, nie jest możliwe wykazanie jego skalowalności, czyli ekstrapolacji jego zachowania na kilkukrotnie powiększony model. Wykorzystywany słownik posiada ok. 52000 pojęć i 90000 relacji pomiędzy nimi, co oznacza, że każde pojęcie stanowi węzeł początkowy dla średnio dwóch relacji. Przy tak małej liczbie powiązań, niewielką liczbę stanowią pojęcia łączące się z więcej niż jednym znaczeniem słowa wieloznacznego. Powoduje to, że w większości przypadków obecność w wypowiedzi wieloznacznego pojęcia i pojęcia, które z nim się łączy determinuje to znaczenie jako wynik. Powoduje to, że algorytm w większości wypadków poszukuje jakiegokolwiek pojęcia mogącego wskazywać znaczenie. Natomiast sytuacja, w której wybór znaczenia jest determinowany wagą jaką wnoszą poszczególne znalezione pojęcia i wagi relacji między nimi jest rzadko spotykana. Analizując przykłady pojęć opisanych w słowniku można przyjąć, że dla każdego z nich jesteśmy w stanie wskazać szacunkowo, nie mniej niż 100 powiązań z innymi pojęciami. Nie jesteśmy także w stanie wskazać jaką ilość słów powinien zawierać słownik, aby prawidłowo odwzorowywać rzeczywistość, gdyż nie umiemy wskazać liczby słów języka polskiego – cytując profesora Jerzego Bralczyka:

*„Nie można też stwierdzić, ile słów jest w całym języku. Mamy wątpliwości, czy przyjąć, że wszystkie słowa zostają w języku na zawsze, czy raczej uważać, że mogą z niego znikać. Nie mamy ustalonych kryteriów, pozwalających uznać, że oto dane słowo już na stałe weszło do języka. Nie wiemy, jakie słowa powinniśmy uważać za polskie, a jakie są tylko chwilowymi, przypadkowymi gośćmi w polszczyźnie. (...) Nie wiemy, ile jest słów, bo granica między słowem i połączeniem słów nie jest wyraźna. Wszystkim wiadomo, że często połączenie dwóch, a czasem i trzech słów daje nową wartość, zasługującą na osobne miejsce w słowniku.”<sup>20</sup>*

---

<sup>20</sup> J.Bralczyk, *Słownik 100 tysięcy potrzebnych słów*, Wydawnictwo Naukowe PWN 2008.

Nawet ograniczone zbliżenie się do takiego stanu modelu doprowadzi do sytuacji, w której akcent działania algorytmu przesunie się na wybór znaczenia na podstawie wielu pojęć. Obecne testy wskazują, że w takiej sytuacji o jakości działania algorytmu będzie decydować ilość relacji pomiędzy pojęciami i ich wagi. Zachowanie algorytmu dla takiego przypadku wymaga dalszych prac.

Można wskazać również inne elementy mogące polepszyć działanie programu. Do najważniejszych można zaliczyć:

- określanie granic wypowiedzi – ograniczenia wymaga rozmiar wypowiedzi, ponieważ w odpowiednio dużym kontekście prawdopodobieństwo pojawienia się dowolnego słowa wzrasta,
- uwzględnienie w ramach słownika i algorytmu informacji o instancjach poszczególnych pojęć i powiązanie ich z nimi; pozwoli to przykładowo na przypisanie artykułu o *Zamku Królewskim w Warszawie* do pojęcia *zamek* w znaczeniu *budowla* – problemem dla takiego zagadnienia jest ogromna ilość informacji, które wymagają obróbki dla otrzymania zadowalającej bazy wiedzy,
- powiększenie słownika semantycznego doprowadzi do sytuacji, w której pojawi się konieczność zmian w metryce relacji, sposobu w jaki relacje oddziałują na odległość pojęć między sobą,
- uwzględnienie reguł składni zdania – pozwoli to na etapie analizy na określenie granic kontekstu niezależnych fragmentów wypowiedzi,
- na etapie pobierania informacji o istniejących pojęciach w wypowiedzi uwzględnienie często występujących w języku zwrotów, kalek językowych,
- ujednoznacznianie na poziomie leksykalnym poprzez powiązanie pojęć słownika semantycznego z odpowiednimi pojęciami słownika fleksyjnego – wymaga synchronizacji pracy przy tworzeniu słownika fleksyjnego i semantycznego.

## BIBLIOGRAFIA

- Anusiewicz Janusz, *Lingwistyka kulturowa. Zarys problematyki*, Warszawa 1995.
- Bartmiński Jerzy, *Definicja kognitywna jako narzędzie opisu konotacji słowa*, Lublin 1988.
- Beckwith Richard, Miller George A., Tengri Randee, *Design and Implementation of the WordNet Lexical Database and Searching Software*, The MIT Press 1998.
- Bralczyk Jerzy, *Słownik 100 tysięcy potrzebnych słów*, Wydawnictwo Naukowe PWN 2008.
- Breyman Ulrich, *Designing Components With the C++ STL (2nd Edition)*, Addison Wesley 2002.
- Brzozowska Małgorzata, *Etymologia a konotacja wybranych nazw kamieni w: Etnolingwistyka. Problemy języka i kultury nr 12.*
- Chomsky Noam, *Syntactic Structures*, Mouton, The Hague 1957.
- Cormen Thomas H., Leiserson Charles E., Rivest Ronald L., Stein Clifford *Introduction to Algorithms*, MIT Press 2001.
- Jurafsky Daniel, Martin James H., *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*, Prentice Hall 2000.
- Mykowiecka Agnieszka, *Inżynieria lingwistyczna, Komputerowe przetwarzanie tekstów w języku naturalnym*, Warszawa 2007.
- Ogden Charles K., *The Meaning of Meaning: A Study of the Influence of Language Upon Thought and of the Science of Symbolism*, London 1946.
- Olinkiewicz Elżbieta, Radzyńska Katarzyna, Styś Halina, *Słownik Encyklopedyczny – Język polski*, Wydawnictwo Europa 1999.
- Quine Willard Van Orman, *Word and Object*, The MIT Press 1960.
- Ravin Yael, Leacock Claudia, *Polysemy: Theoretical and Computational Approaches*, Oxford University Press 2000.
- Sapir Edward, *Kultura, język, osobowość*, Warszawa 1978.

Wiesław Lubaszewski, Henryk Wróbel, Marek Gajęcki, Barbara Moskal, Alicja Orzechowska, Paweł Pietras, Piotr Pisarek, Teresa Rokicka, *Słownik Fleksyjny Języka Polskiego*, LexisNexis 2001

WordNet a lexical database for the English language: <http://wordnet.princeton.edu/>

Princeton University Cognitive Science Laboratory: <http://cogsci.princeton.edu/>

Introduction to Information Extraction Technology – Douglas E. Appelt, David Israel:  
<http://www.ai.sri.com/~appelt/ie-tutorial/>

Gramatyka języka polskiego – Grzegorz Jagodziński: <http://grzegorz.w.interia.pl/gram/gram00.html>

Wstęp do kognitywistyki – Włodzisław Duch:  
[http://www.fizyka.umk.pl/~duch/Wyklady/Cog\\_plan.html](http://www.fizyka.umk.pl/~duch/Wyklady/Cog_plan.html)

Word Sense Disambiguation Improves Statistical Machine Translation – Yee Seng Chang, Hwee Tou Ng, David Chiang: [http://www.isi.edu/~chiang/papers/chan\\_wsd\\_in\\_mt.pdf](http://www.isi.edu/~chiang/papers/chan_wsd_in_mt.pdf)

Frequency Asked Questions for Google Suggest: <http://labs.google.com/suggestfaq.html>

Pozycjonowanie w AdWords:  
<http://adwords.google.pl/support/bin/answer.py?answer=65133&topic=9352>

Strona Słownika Fleksyjnego, KI AGH: [http://winnie.ics.agh.edu.pl/proj\\_uk/fleksbaz/](http://winnie.ics.agh.edu.pl/proj_uk/fleksbaz/)

Strona Słownika Semantycznego Języka Polskiego, KI AGH:  
[http://winnie.ics.agh.edu.pl/proj\\_re/slse/index.html](http://winnie.ics.agh.edu.pl/proj_re/slse/index.html)

WordNet: What is to be Done? Hans Patrick,  
<http://www.fi.muni.cz/gwc2004/pres/panel/Hanks/hanks-panel.pdf>.

Sockety dla systemu operacyjnego Windows:  
[http://msdn.microsoft.com/en-us/library/ms740673\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms740673(VS.85).aspx)

Opis funkcji bibliotecznej strtok: <http://www.cplusplus.com/reference/cstring/strtok.html>