

Uniwersytet Warszawski
Wydział Matematyki, Informatyki i Mechaniki

Imię i nazwisko

Nr albumu: nralbumu

Intuicyjny język wyszukiwania TQL (Tablets Query Language)

**Praca magisterska
na kierunku INFORMATYKA**

Praca wykonana pod kierunkiem
dra Roberta Dąbrowskiego
Instytut Informatyki

czerwiec 2010

Oświadczenie kierującego pracą

Potwierdzam, że niniejsza praca została przygotowana pod moim kierunkiem i kwalifikuje się do przedstawienia jej w postępowaniu o nadanie tytułu zawodowego.

Data

Podpis kierującego pracą

Oświadczenie autora (autorów) pracy

Świadom odpowiedzialności prawnej oświadczam, że niniejsza praca dyplomowa została napisana przeze mnie samodzielnie i nie zawiera treści uzyskanych w sposób niezgodny z obowiązującymi przepisami.

Oświadczam również, że przedstawiona praca nie była wcześniej przedmiotem procedur związanych z uzyskaniem tytułu zawodowego w wyższej uczelni.

Oświadczam ponadto, że niniejsza wersja pracy jest identyczna z załączoną wersją elektroniczną.

Data

Podpis autora (autorów) pracy

Streszczenie

Sumerologia jest dziedziną badań nad antycznym językiem Sumerów, w której kluczowym zagadnieniem jest przeszukiwanie dużych zbiorów informacji zapisanych na odnalezionych tabliczkach sumeryjskich.

W pracy przedstawiono definicję przeznaczonego dla sumerologów intuicyjnego języka przeszukiwania zbiorów tabliczek (Tablets Query Language) wraz z jego przykładowymi implementacjami. Jedna z nich jest oparta na relacyjnej bazie danych, a druga na bazie w postaci pliku XML.

Celem tej pracy jest stworzenie języka zapytań intuicyjnego dla sumerologów, stanowiącego znaczące uproszczenie w stosunku do SQL dzięki wprowadzeniu pojęć naturalnych dla rozważanej dziedziny. TQL nadal pozwala na tworzenie skomplikowanych zapytań wyszukiwanych, natomiast nie udostępnia funkcjonalności tworzenia i modyfikowania bazy. Można go rozszerzać i zmieniać tak, by mógł służyć do zastosowań z innych dziedzin.

Słowa kluczowe

tabliczka, sumerologia, kliny, odczyty, język dziedzinowy, DSL, Domain-Specific Languages

Dziedzina pracy (kody wg programu Socrates-Erasmus)

11.3 Informatyka

Klasyfikacja tematyczna

H. INFORMATION SYSTEMS
H.2. DATABASE MANAGEMENT
H.2.3 Languages

Tytuł pracy w języku angielskim

Intuitive query language TQL (Tablets Query Language)

Spis treści

Wprowadzenie	5
I Problem przeszukiwania bazy tabliczek sumeryjskich	9
1. Podstawowe pojęcia	11
1.1. Pojęcia dziedzinowe	11
1.2. Pojęcia informatyczne	11
1.3. Pojęcia paradygmatyczne	12
2. Dziedzina problemu	13
3. Wcześniejsze rozwiązania	15
3.1. The Cuneiform Digital Library Initiative [3]	15
3.2. The Electronic Text Corpus of Sumerian Literature [4]	16
II Definicja języka TQL	17
4. Tablets Query Language	19
5. Gramatyka	21
5.1. Struktura leksykalna	21
5.1.1. Literały	21
5.1.2. Słowa kluczowe	21
5.1.3. Znaki specjalne	21
5.2. Struktura składniowa języka	22
6. Semantyka	23
6.1. Zapytanie proste	23
6.2. Zapytanie złożone	23
6.3. Zapytanie zdefiniowane	24
6.4. Wywołanie zapytania zdefiniowanego	24
6.4.1. Zwyczajne	24
6.4.2. Z dodatkowym warunkiem wyszukiwania	24

III Implementacja	25
7. Moduły podstawowe	29
7.1. Parser	29
7.2. Analizator kontekstowy	29
7.3. Translator	30
7.4. Baza	30
7.5. Pliki pomocnicze	31
8. Moduły wymienne	33
8.1. Baza PostgreSQL	34
8.1.1. Diagram encji	34
8.1.2. Translator_postgres	35
8.1.3. Database_postgres	38
8.2. Baza XML	39
8.2.1. Schemat dokumentu	39
8.2.2. Translator_xml	39
8.2.3. Database_xml	42
Podsumowanie	43
Podsumowanie projektu	43
Możliwości rozwoju	43
Dodatki	47
A. Schemat bazy danych xml	47
Bibliografia	49

Wprowadzenie

Opis problemu

Ok. 3500 lat przed naszą erą starożytni Sumerowie jako prawdopodobnie pierwsi na świecie zaczęli używać pisma. Ze względu na kształt liter odcisniętych za pomocą trzciny w miękkiej glinie zostało ono później nazwane pismem klinowym. Sumerowie używali go głównie w celach administracyjnych. Zapisywali między innymi rachunki za dostawy zwierząt oraz rozliczenia z wykonanej na polu lub w warsztacie pracy i należnej zapłaty. Jednymi z ciekawszych tekstów są tzw. listy królów, na których znajdują się daty panowania kolejnych władców i ich osiągnięcia w poszczególnych latach. Ważniejsze tabliczki wypalano, jednak większość po pewnym czasie niszczone, aby zużytą glinę można było ponownie wykorzystać. Do dzisiaj przetrwały głównie te tabliczki, które zostały wypalone podczas przypadkowego pożaru archiwum. Wiele z nich jest zniszczonych, jednak wciąż stanowią ogromną wartość historyczną.

Na podstawie zachowanych tabliczek można się wiele dowiedzieć o historii Sumeru, a także o życiu zwykłych ludzi w tym okresie, o ich zarobkach, zasiłkach i dniach wolnych. Odczytywaniem tych informacji zajmują się sumerolodzy z całego świata. Aby łatwiej dzielić się zdobytą wiedzą, zaczęli tworzyć cyfrowe bazy tabliczek, dostępne w internecie. Największa z nich to Cuneiform Digital Library Initiative (CDLI), zawierająca prawie 225 tys tekstów.

Jednak wyszukiwanie interesujących tabliczek jest w tym momencie dosyć niewygodne.

Sumerolodzy mogliby skorzystać ze specyficznych dla baz danych języków zapytań (np. SQL), ale po pierwsze większość z nich nie zna tych języków, a po drugie budowanie w ten sposób złożonych zapytań dotyczących danych tekstowych jest czasochłonne. Takie języki były tworzone z myślą o bardziej generycznych zastosowaniach i nie odpowiadają potrzebom sumerologów. Z tych powodów większość istniejących obecnie serwisów internetowych udostępnia formularze ułatwiające wprowadzanie kryteriów wyszukiwania. Jednak mają one ograniczone możliwości i nie pozwalają na konstruowanie bardziej skomplikowanych zapytań. Dlatego istnieje potrzeba stworzenia narzędzia wspomagającego wyszukiwanie, które będzie łączyło w sobie jak największą siłę wyrazu oraz łatwość użycia dla osób znających jedynie dziedzinę problemu.

Przy tego typu problemach z pomocą przychodzi informatyka.

Propozycja rozwiązania

Do tej pory powstało wiele różnych systemów informatycznych, a doświadczenie stopniowo zdobywane przy ich budowie prowadzi do ciągłego rozwijania nauki zwanej inżynierią oprogramowania, zajmującej się praktyczną stroną wytwarzania systemów. Obecnie jest znanych bardzo wiele podejść do tworzenia oprogramowania, a każdemu odpowiada gałąź problemów, w których sprawdza się najlepiej.

Najbardziej popularne jest programowanie zorientowane obiektowo (ang. *Object-Oriented*

Programming, OOP), które polega na modelowaniu świata rzeczywistego w postaci *obiektów*. Obiekty posiadają dane oraz metody, które mogą te dane zmieniać. Program jest zbiorem obiektów komunikujących się między sobą.

Innym podejściem jest architektura zorientowana na usługi (ang. *Service-Oriented Architecture*, SOA), w której definiuje się niezależne od siebie usługi (services) i udostępnia tylko ich interfejsy, ukrywając implementację.

Rozważając problem sumerologów zdecydowaliśmy się na programowanie zorientowane na język (ang. *Language-Oriented Programming*). Polega ono na stworzeniu języka odpowiadającego dziedzinie problemu, czyli tzw. języka dziedzinowego (ang. *Domain-Specific Language*, DSL). Dopiero w tym nowym języku rozwiązuje się konkretny problem (np. znalezienie tekstów sumeryjskich z epoki Ur III dotyczących owiec). Opisywanie problemów dziedziny jest wtedy znacznie prostsze i bardziej naturalne, a dzięki temu wygodne dla ludzi związanych tylko z konkretną dziedziną.

Wyróżnia się dwa rodzaje języków dziedzinowych - “wewnętrzne” (ang. *internal*), które zachowują składnię istniejących już języków i ograniczają tylko ich możliwości, oraz “zewnętrzne” (ang. *external*), które są zupełnie nowymi językami, tłumaczonymi do istniejących wcześniej. W naszym przypadku narzucającym się podejściem było stworzenie nowego, zewnętrznego DSL, zaprojektowanego specjalnie do wyszukiwania tabliczek sumeryjskich, który miałby składnię przejrzystą dla sumerologa.

Do każdej bazy zawierającej tabliczki można skonstruować moduł pozwalający na wyszukiwanie z użyciem takiego DSL. Zatem jest to rozwiązanie, które może zostać wykorzystane przez wiele różnych serwisów, zarówno nowo powstających, jak i tych już istniejących. Dzięki temu wyszukiwarki, internetowe oraz stacjonarne, mogłyby mieć taki sam lub bardzo podobny interfejs, co znacznie ułatwiłoby pracę sumerologom. DSL pozwoliłby usystematyzować pracę nad wszelkimi wyszukiwarkami tabliczek sumeryjskich, a przez to również pracę sumerologów. Jest to niewątpliwa zaleta takiego rozwiązania. Wyraźnie widać jego przewagę nad podejściem polegającym na niezależnym rozwijaniu interfejsów poszczególnych wyszukiwarek, czego efektem są coraz bardziej skomplikowane formularze do wpisywania zapytań.

Tablets Query Language

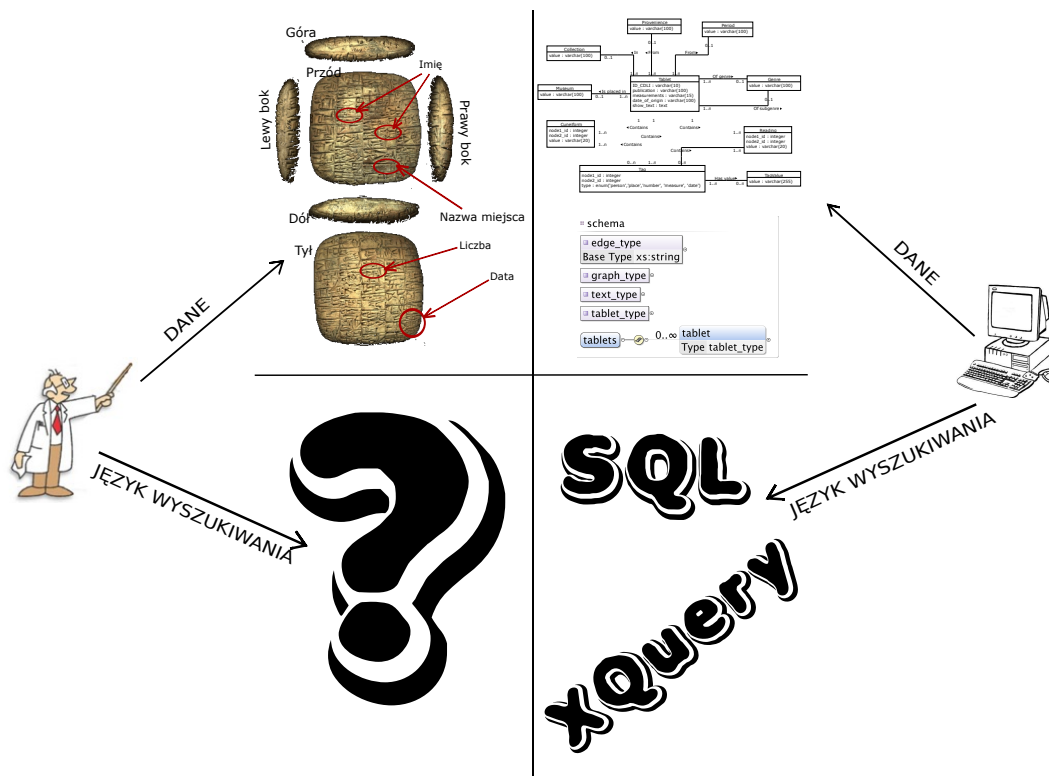
Wychodząc naprzeciw potrzebom sumerologów prezentujemy stworzony przez nas język dziedzinowy do wyszukiwania tabliczek sumeryjskich — *Tablets Query Language* (TQL). Został on zaprojektowany przede wszystkim jako intuicyjny i zrozumiały dla osób z dziedziny problemu. Ma także minimalnie ograniczać siłę wyrażu tzn. pozwalać na konstruowanie jak najbardziej skomplikowanych zapytań. Chcemy również, żeby był niezależny od sposobu reprezentacji tabliczek. W niniejszej pracy udowodnimy, że TQL spełnia wszystkie te wymagania.

Zgodnie z paradygmatem języka dziedzinowego, TQL jest tłumaczony do innych, istniejących wcześniej języków zapytań bardziej ogólnego zastosowania, np. SQL, XQuery, odpowiednio do sposobu reprezentacji danych.

Sposobem reprezentacji danych nazywamy rodzaj bazy danych (np. relacyjna, obiektowa, XML), konkretną implementację serwera (np. PostgreSQL, Zorba) oraz schemat danych. Dla każdego takiego sposobu zapytanie w docelowym języku może się różnić, a zatem różnić się będzie również program tłumaczący TQL na ten język. Oznacza to, że dla każdego sposobu reprezentacji danych należy stworzyć oddzielny program tłumaczący. W ramach niniejszej pracy oprócz języka TQL zaprezentujemy dwie prototypowe implementacje komponentu do wyszukiwarki, zawierającego m.in. program tłumaczący. Ten komponent, nazywany dalej w

translatorem, tłumaczy zapytanie TQL na zapytanie w języku specyficznym dla danej bazy, przekazuje przetłumaczone zapytanie do bazy, odbiera wynik i przedstawia go w odpowiedniej postaci. Pokażemy implementacje translatora dla bazy PostgreSQL (z językiem wyszukiwania SQL) oraz dla bazy XML (z językiem wyszukiwania XQuery).

TQL może być także podstawą do tworzenia podobnych języków wyszukiwań dostosowanych do potrzeb innych grup ludzi, np. językoznawców. Większość programów ułatwiających tworzenie zapytań jest skomplikowana, daje ograniczone możliwości lub jest przystosowana głównie do przetwarzania danych liczbowych. Tablets Query Language rozwiązuje te problemy: jest prosty i intuicyjny, przystosowany głównie do tekstów, minimalnie zmniejsza siłę wyrazu oraz łatwo go rozbudowywać.



Rysunek 1: Przedstawienie problemu

Część I

Problem przeszukiwania bazy tabliczek sumeryjskich

Rozdział 1

Podstawowe pojęcia

1.1. Pojęcia dziedzinowe

Sumerolog - naukowiec zajmujący się odczytywaniem pisma klinowego w języku sumeryjskim. Na potrzeby tej pracy to pojęcie jest rozszerzone do wszystkich ludzi zajmujących się odczytywaniem tabliczek klinowych (także w innych językach) i czerpiących z nich wiedzę historyczną.

Tabliczka (ang. tablet) - w tej pracy tabliczka będzie oznaczała tabliczkę klinową w wersji elektronicznej (chyba, że zostanie zaznaczone inaczej). Dla rozróżnienia, kiedy będziemy mówić o “prawdziwej”, glinianej tabliczce, będziemy używać pojęcia **gliniana tabliczka**

Proweniencja (ang. provenience) - pojęcie używane przez sumerologów, oznacza miejsce pochodzenia/znalezienia glinianej tabliczki

Kliny (ang. cunes) - znaki występujące na glinianych tabliczkach.

Odczyty (ang. readings) - sposób transkrypcji klinów. Tabliczki elektroniczne są zapisane za pomocą odczytów. Na ogół jeden klin odpowiada jednemu odczytowi, ale może odpowiadać wielu różnym sekwencjom odczytów. Dodatkowo jeden odczyt może być zapisany za pomocą sekwencji klinów.

Pieczęć (ang. seal) - część tabliczki zawierająca znak rozpoznawczy autora

1.2. Pojęcia informatyczne

Alfabet (zbiór Σ , zbiór symboli terminalnych) - zbiór symboli (np. słów kluczowych, znaków specjalnych, literalów), z których zbudowane są słowa - konstrukcje języka.

Zbiór symboli nieterminalnych - zbiór symboli pomocniczych, rozłączny z alfabetem.

Słowo nad alfabetem Σ - skończony ciąg symboli należących do zbioru Σ .

Język nad alfabetem Σ - zbiór słów nad alfabetem Σ .

Reguły gramatyki - reguły definiujące sposób tworzenia słów nad danym alfabetem. Każda reguła jest postaci $S1 \rightarrow S2$, gdzie $S1$ i $S2$ to ciągi symboli terminalnych i nieterminalnych, przy czym w ciągu $S1$ musi wystąpić przynajmniej jeden symbol nieterminalny.

Reguły określają możliwe podstawienia symboli w wyprowadzanym słowie - ciąg $S1$ można zastąpić przez $S2$.

Gramatyka - formalny sposób definiowania języka. Składa się z czterech elementów: zbioru symboli terminalnych, zbioru symboli nieterminalnych, symbolu startowego (należącego do zbioru symboli nieterminalnych) oraz zbioru reguł gramatyki. Wyprowadzanie słowa należącego do języka rozpoczynamy od symbolu startowego, przeprowadzamy podstawienia zgodnie z regułami gramatyki i kończymy, gdy wszystkie symbole w słowie należą do zbioru symboli terminalnych. Język określony przez gramatykę jest to zbiór słów, które są możliwe do wyprowadzenia z symbolu startowego za pomocą reguł gramatyki.

Struktura leksykalna języka - definicja symboli terminalnych (alfabetu).

Struktura składniowa języka - opis składni języka, zapisany np. za pomocą reguł gramatyki.

Semantyka - znaczenie i funkcja poszczególnych konstrukcji języka.

1.3. Pojęcia paradygmatyczne

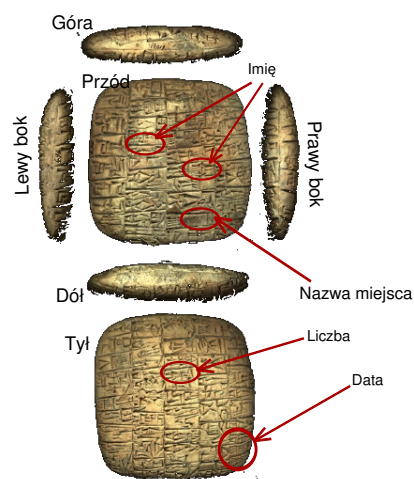
Język dziedzinowy (ang. domain-specific language, DSL) - język programowania szczególnego zastosowania, rozwiązujący specyficzny problem lub zajmujący się wąską dziedziną, stworzony specjalnie na potrzeby danej dziedziny i do niej dostosowany.

Rozdział 2

Dziedzina problemu

Podstawowym elementem rozpatrywanej przez nas dziedziny jest *tabliczka sumeryjska*. Pojęcie to odnosi się zarówno do fizycznej *tabliczki glinianej* jak i do jej cyfrowej reprezentacji, odpowiedniej do przechowywania na komputerze.

Tabliczki gliniane mają różne kształty i rozmiary. Mogą być zapisane z wielu stron (z przodu, z tyłu, od góry, z boku, itp.), a także mogą zawierać pieczęcie, na których również znajduje się tekst. Przez kilka tysięcy lat istnienia ulegały też uszkodzeniom uniemożliwiającym obecnie ich pełne odczytanie. Przykładowa gliniana tabliczka zaprezentowana jest na rysunku 2.1.



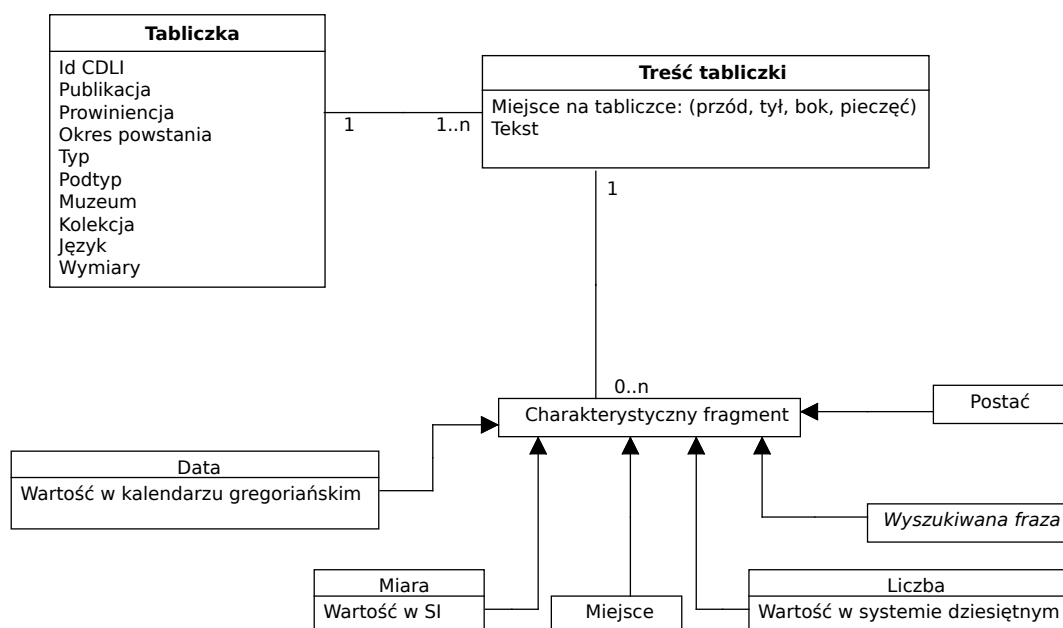
Rysunek 2.1: Gliniana tabliczka - struktura

Tekst na tabliczkach jest zapisany za pomocą *klinów*. Pismo klinowe jest prawdopodobnie najstarszym pismem na świecie (konkurując z nim jedynie hieroglify egipskie). Początkowo składało się z piktogramów z czasem coraz bardziej upraszczanych. Służyło pomocą w rozliczeniach i umowach, zapisywano np. ilość i rodzaj obiecanych ludziom zwierząt.

Treść tabliczek to najczęściej dokumenty urzędowe dotyczące czynności administracyjnych i gospodarczych, produkcji i dystrybucji rozmaitych artykułów, uroczystości religijno-kulturowych, umów cywilno-prawnych itp. [1] Często można spotkać w niej między innymi imiona osób i bóstw, liczby, jednostki (np. mina - jednostka wagi), miejsca, daty. Niektóre z tych elementów można przetłumaczyć na współczesny język, na przykład jednostki można przeliczyć na SI, datę opisową na datę wg obecnego kalendarza gregoriańskiego.

Dla ułatwienia pracy sumerolodzy tłumaczą zapis klinowy na tzw. *odczyty*. Są one sposobem zapisania tekstu sumeryjskiego za pomocą współczesnych znaków - liter alfabetu łacińskiego i cyfr arabskich, na podstawie prawdopodobnego brzmienia słów w języku sumeryjskim. W tłumaczeniu uwzględniony jest także kształt tabliczki i rozmieszczenie poszczególnych fragmentów tekstu (zgodnie z rysunkiem 2.1).

Odczyty zawarte w cyfrowym zapisie tabliczki są tylko jednym z wariantów tłumaczenia z klinów. Przede wszystkim dlatego, że jeden klin może zostać odczytany na wiele różnych sposobów. Ponadto odczytowi może odpowiadać nie tylko jeden klin, ale także ich sekwencja.



Rysunek 2.2: Co powinna zawierać cyfrowa reprezentacja tabliczki

Ponieważ w cyfrowej wersji nie ma klinów, trudno jest zweryfikować ewentualne pomyłki w tłumaczeniach.

W takiej właśnie formie tabliczki elektroniczne są przechowywane w pamięci komputerów. Niestety ten zapis może być mylący ze względu na niejednoznaczności występujące przy tłumaczeniu.

Źródłem dodatkowych problemów są uszkodzenia tabliczek. W wielu przypadkach naukowcy mogą domyślić się jak wyglądał brakujący fragment. Zawsze natomiast informacje o uszkodzeniu oraz domniemaną brakującą treść zapisują w tabliczce cyfrowej. Ponieważ pewne konwencje opisywania tych informacji zostały przyjęte dosyć późno i nie przez wszystkich sumerologów, są one zapisywane na różne sposoby. Również tutaj możliwe są pomyłki i niejednoznaczności.

Tabliczka cyfrowa, oprócz treści, zawiera także informacje, które nie były zawarte na tabliczce glinianej. Są to metadane takie jak miejsce znalezienia tabliczki, okres, w którym powstała, czy nazwa kolekcji, do której obecnie należy. Te informacje są istotne przy przeszukiwaniu, gdyż często pozwalają na określenie o czym jest tabliczka bez dokładnej analizy jej treści. W praktyce, wśród sumerologów, atrybutem, który w znacznym stopniu pomaga zidentyfikować tabliczkę jest informacja o jej publikacji.

Nie ma wątpliwości, że cyfrowa postać tabliczek sumeryjskich wraz z podstawowymi możliwościami wyszukiwania znacząco ułatwia sumerologom ich badanie. Brakuje im jednak możliwości bardziej skomplikowanego i rozbudowanego wyszukiwania. Wychodząc naprzeciw tej potrzebie chcemy stworzyć język, w którym będą oni mogli w łatwy sposób wyrażać, jakich tabliczek potrzebują, i który jednocześnie będzie można wykorzystać do przeszukiwania baz danych. Język ten powinien przede wszystkim umożliwiać wyszukiwanie na podstawie treści tabliczki (odczytów), alternatywnych tłumaczeń (wymaga to przetłumaczenia tabliczki z powrotem na kliny) oraz metadanych. Dodatkową zaletą byłoby wyszukiwanie po specyficznych fragmentach (tagach), takich jak imiona, jednostki, daty. W pierwszej wersji języka implementujemy tylko wyszukiwanie po odczytach i metadanych.

Rozdział 3

Wcześniejsze rozwiązania

Dotychczasowe rozwiązania problemu przeszukiwania bazy tabliczek opierają się o pomysł stworzenia formularza, w którym, wypełniając odpowiednie pola, można określić jakich tabliczek szukamy. Formularze różnią się od siebie przede wszystkim poziomem skomplikowania. Im bardziej rozbudowany, tym więcej można w nim zdefiniować warunków dotyczących parametrów wyszukiwania, ale jednocześnie tym trudniej z niego korzystać. Z punktu widzenia dziedziny naszego problemu podstawową wadą takich rozwiązań jest brak możliwości tworzenia zapytań złożonych oraz brak elastyczności.

Poniżej przedstawimy dwie przykładowe strony internetowe oferujące wyszukiwanie za pomocą formularzy.

3.1. The Cuneiform Digital Library Initiative [3]



The screenshot displays the search interface of the Cuneiform Digital Library Initiative (CDLI). At the top, there are radio buttons for 'Representation of results' (Concise list of results, List of results with Images, Browse) and 'Number of results' (100, 1000, 1, 10, 50 records to be shown). Below this is the 'Transliteration Search' section, which includes a text input field for 'Word/Part of Word' (e.g., 'udu') and a dropdown menu for 'part of word'. A checkbox for 'Advanced search syntax' is also present. The 'Catalogue Search' section follows, featuring several rows of search criteria: 'Primary Publication' (e.g., 'TRU 001'), 'Author(s)' (last name first), 'Date of publication', 'Other Publication(s)' (e.g., 'MVN 18'), 'Collection' (preferred 'contains'), 'Museum Number' (e.g., 'VAT'), and 'Accession Number' (e.g., 'K 00001'). Each criterion has a dropdown menu for search operators (e.g., 'begins with', 'contains', 'does not contain') and a text input field. To the right of each input field is a 'Sort by' radio button.

Rysunek 3.1: Formularz wyszukiwania na stronie CDLI

Największą znaną nam bazą tekstów sumeryjskich jest The Cuneiform Digital Library Initiative (CDLI), która zawiera ok. 225 tys. tabliczek. Posiada stosunkowo rozbudowany formularz, który pozwala na wyszukiwanie po wielu parametrach tabliczki (m. in. po danych dotyczących czasu i miejsca powstania, komentarzach CDLI). Dla każdego parametru jest pole tekstowe z możliwymi opcjami wyszukiwania: "begins with", "contains", "does not contain". Do wyszukiwania na podstawie treści tabliczki jest pole tekstowe z opcjami "word", "part of word". Jest również checkbox "advanced search syntax", jednak brakuje wyjaśnienia

jak go używać. Nie można tworzyć bardziej skomplikowanych warunków dla poszczególnych parametrów (np. nie można wyszukać tabliczek, które powstały w Uruk lub w Garshanie) ani łączyć kilku zapytań w jedno zapytanie złożone.

3.2. The Electronic Text Corpus of Sumerian Literature [4]

The Electronic Text Corpus of Sumerian Literature

Search for as sort by

in category or in comp. no(s).

Encode char. as: Display: ☒ Full text in new window / ☐ Full para. in gloss

Rysunek 3.2: Formularz wyszukiwania na stronie ETCSL

The Electronic Text Corpus of Sumerian Literature (ETCSL) posiada znacznie mniejszą bazę niż CDLI, składającą się z około 400 tekstów, głównie literackich. Ma ograniczone możliwości wyszukiwania po metadanych - tylko po kategorii tekstu, jednak udostępnia tworzenie bardziej skomplikowanych zapytań dotyczących treści tabliczki. Pozwala określić typ wyszukiwanego słowa - do wyboru: form, lemma, label, pos, emesal, sign, a także jego znaczenie w tekście (np. czy jest imieniem bóstwa) oraz część mowy, do której należy. Można również wyszukiwać na podstawie fragmentu słowa.

Część II

Definicja języka TQL

Rozdział 4

Tablets Query Language

Tablets Query Language (TQL) jest naszą propozycją rozwiązania problemu wyszukiwania tabliczek sumeryjskich. Jest to zewnętrzny język dziedziny stworzony od podstaw do tego, aby służyć sumerologom jako język zapytań. Jego składnia jest bardzo prosta i intuicyjna, co będzie widać na przykładach w rozdziale 6. Jednocześnie, dzięki możliwości formułowania skomplikowanych wyrażeń do wyszukiwania na podstawie pojedynczej metadanej lub treści tabliczki, oraz możliwości łączenia wielu zapytań w jedno zapytanie złożone, TQL ma bardzo dużą siłę wyrazu.

Dodatkowo, jednym z głównych założeń przyjętych przy tworzeniu języka TQL jest niezależność od rzeczywistej reprezentacji danych. Przy konstruowaniu go nie brałyśmy pod uwagę sposobu fizycznej reprezentacji tabliczek, czyli rodzaju bazy danych oraz schematu danych. Skupiłyśmy się jedynie na dziedzinie problemu, czyli na tym, co zawiera tabliczka oraz na podstawie jakich informacji chcemy wyszukiwać. Zakładamy, że niezależnie od reprezentacji danych takie wyszukiwanie będzie możliwe, chociaż oczywiście skonstruowanie odpowiedniego zapytania może być skomplikowane. Przetłumaczenie zapytania TQL na zapytanie w języku odpowiednim do reprezentacji danych, np. SQL, XQuery, jest zadaniem programu tłumaczącego. Dzięki temu praca polegająca na przetłumaczeniu zapytania z "języka dziedziny" na "język komputerów" jest wykonana tylko raz dla każdego sposobu reprezentacji danych, i to przez programistów, a nie sumerologów.

Język TQL umożliwia wyszukiwanie na podstawie kryteriów dotyczących następujących danych:

Opis	Nazwa pola w TQL
numer tabliczki w bazie CDLI	cdli_id
miejsce pochodzenia (proweniencja)	provenience
okres powstania	period
typ i podtyp	genre
rok powstania	year
publikacja	publication
treść (odczyty)	text
kolekcja	collection
muzeum	museum

Język można łatwo rozszerzać, aby umożliwić tworzenie kryteriów wyszukiwania w oparciu o inne dane (np. kliny, zawartość pieczęci).

W kolejnych dwóch rozdziałach przedstawimy gramatykę i semantykę zaprojektowanego

przez nas języka TQL.

Rozdział 5

Gramatyka

W tym rozdziale przedstawimy gramatykę zaprojektowanego przez nas języka. W pierwszej części pokażemy strukturę leksykalną TQL, czyli elementy, z których buduje się zapytania. W drugiej części zaprezentujemy reguły tworzenia zapytań, zapisane w formie reguł gramatyki w notacji BNF.

5.1. Struktura leksykalna

5.1.1. Literały

String

Literał **String** jest ciągiem dowolnych znaków w cudzysłowie ("). Nie może zawierać jedynie znaków "\"" niepoprzedzonych "\".

SłowoOdLitery

Literał **SłowoOdLitery** to ciąg liter, cyfr oraz znaków "-", "'", "_", zaczynający się od litery, z wyjątkiem słów kluczowych.

SłowoOdLiczby

Literał **SłowoOdLiczby** to ciąg liter, cyfr oraz znaków "-", "'", "_", zaczynający się od cyfry.

5.1.2. Słowa kluczowe

```
as      define  in
search
```

5.1.3. Znaki specjalne

```
( )      +
/  --    *
:        \n (koniec linii)
```

5.2. Struktura składniowa języka

Poniżej przedstawimy reguły gramatyki TQL w notacji BNF. Nieterminalne są zapisane pomiędzy „ \langle ” a „ \rangle ”. Symbole „ $::=$ ” (produkcja), „ $|$ ” (lub) i „ ϵ ” (pusta reguła) należą do notacji BNF. Wszystkie pozostałe symbole to terminale.

$$\begin{aligned}
 \langle \text{Zapytanie Złożone} \rangle &::= \langle \text{Lista Zapytań} \rangle \\
 \langle \text{Lista Zapytań} \rangle &::= \langle \text{Zapytanie} \rangle \\
 &\quad | \quad \langle \text{Zapytanie} \rangle \langle \text{Lista Zapytań} \rangle \\
 \langle \text{Zapytanie} \rangle &::= \langle \text{Lista Linii Zapytania} \rangle \langle \text{Lista Pustych Linii} \rangle \\
 &\quad | \quad \text{define } \backslash n \langle \text{Zapytanie} \rangle \text{ as } \langle \text{Nazwa} \rangle \langle \text{Lista Pustych Linii} \rangle \\
 &\quad | \quad \text{search } \backslash n \langle \text{Zapytanie} \rangle \text{ in } \langle \text{Nazwa} \rangle \langle \text{Lista Pustych Linii} \rangle \\
 &\quad | \quad \text{search } \langle \text{Nazwa} \rangle \langle \text{Lista Pustych Linii} \rangle \\
 &\quad | \quad \langle \text{Lista Pustych Linii} \rangle \\
 \langle \text{Lista Linii Zapytania} \rangle &::= \langle \text{Linia Zapytania} \rangle \backslash n \\
 &\quad | \quad \langle \text{Linia Zapytania} \rangle \backslash n \langle \text{Lista Linii Zapytania} \rangle \\
 \langle \text{Linia Zapytania} \rangle &::= \langle \text{Nazwa pola} \rangle : \langle \text{Wyrażenie} \rangle \\
 \langle \text{Wyrażenie} \rangle &::= \langle \text{Wyrażenie} \rangle + \langle \text{Wyrażenie1} \rangle \\
 &\quad | \quad \langle \text{Wyrażenie} \rangle / \langle \text{Wyrażenie1} \rangle \\
 &\quad | \quad \langle \text{Wyrażenie1} \rangle \\
 \langle \text{Wyrażenie1} \rangle &::= -- \langle \text{Wyrażenie1} \rangle \\
 &\quad | \quad \langle \text{Wyrażenie2} \rangle \\
 \langle \text{Wyrażenie2} \rangle &::= \langle \text{Wyrażenie3} \rangle * \langle \text{Wyrażenie3} \rangle \\
 &\quad | \quad \langle \text{Tekst} \rangle \\
 &\quad | \quad (\langle \text{Wyrażenie} \rangle) \\
 \langle \text{Wyrażenie3} \rangle &::= \langle \text{Wyrażenie2} \rangle \\
 &\quad | \quad \epsilon \\
 \langle \text{Lista Pustych Linii} \rangle &::= \epsilon \\
 &\quad | \quad \langle \text{Pusta Linia} \rangle \langle \text{Lista Pustych Linii} \rangle \\
 \langle \text{Pusta Linia} \rangle &::= \backslash n \\
 \langle \text{Tekst} \rangle &::= \text{String} \\
 &\quad | \quad \langle \text{Słowo} \rangle \\
 \langle \text{Słowo} \rangle &::= \text{Słowo0dLiter} \\
 &\quad | \quad \text{Słowo0dLiczby} \\
 \langle \text{Nazwa pola} \rangle &::= \text{Słowo0dLiter} \\
 \langle \text{Nazwa} \rangle &::= \text{String}
 \end{aligned}$$

Powyższa gramatyka jest gramatyką bezkontekstową.

Rozdział 6

Semantyka

Poniżej przedstawiamy semantykę wybranych przykładów.

6.1. Zapytanie proste

```
provenience: Gar*  
period: "Ur III"  
genre: Administrative  
text: udu + (masz2/ugula) --szabra
```

Wynikiem zapytania będą wszystkie tabliczki, które:

- pochodzą z miejscowości o nazwie zaczynającej się na “Gar”
- pochodzą z okresu Ur III
- są dokumentami administracyjnymi
- zawierają słowo “udu” oraz conajmniej jedno ze słów “masz2” lub “ugula”
- nie zawierają słowa “szabra”

6.2. Zapytanie złożone

```
provenience: Ur  
period: "Ur III"/"Ur IV"  
text: udu --szabra
```

```
text: masz2/ugula  
publication: *tan  
provenience: Ur
```

Wynikiem zapytania będą wszystkie tabliczki, które:

- pochodzą z miejscowości Ur
- pochodzą z okresu Ur III lub Ur IV
- zawierają słowo “udu”

- nie zawierają słowa “szabra”

oraz wszystkie tabliczki, które:

- zawierają słowo ”masz2“ lub ”ugula“
- zostały opublikowane w pracy, której nazwa kończy się na ”tan“
- pochodzą z miejscowości Ur

6.3. Zapytanie zdefiniowane

```
define
  provenience: Gar*a
  period: Ur III
  text: "udu ban"/mash2
as "zwierzęta w Gar*a"
```

Wynikiem zapytania (po jego wywołaniu) będą wszystkie tabliczki, które:

- pochodzą z miejscowości, których nazwy zaczynają się na ”Gar“ i kończą na ”a”
- pochodzą z okresu Ur III
- zawierają conajmniej jedną z fraz ”udu ban“ lub ”mash2“

6.4. Wywołanie zapytania zdefiniowanego

6.4.1. Zwykle

```
search "zwierzęta w Gar*a"
```

Wynikiem zapytania będą dokładnie te tabliczki, które spełniają wszystkie warunki zapytania ”zwierzęta w Gar*a”.

6.4.2. Z dodatkowym warunkiem wyszukiwania

```
search
  text: adad-tilati
in "zwierzęta w Gar*a"
```

Wynikiem zapytania będą wszystkie tabliczki, które:

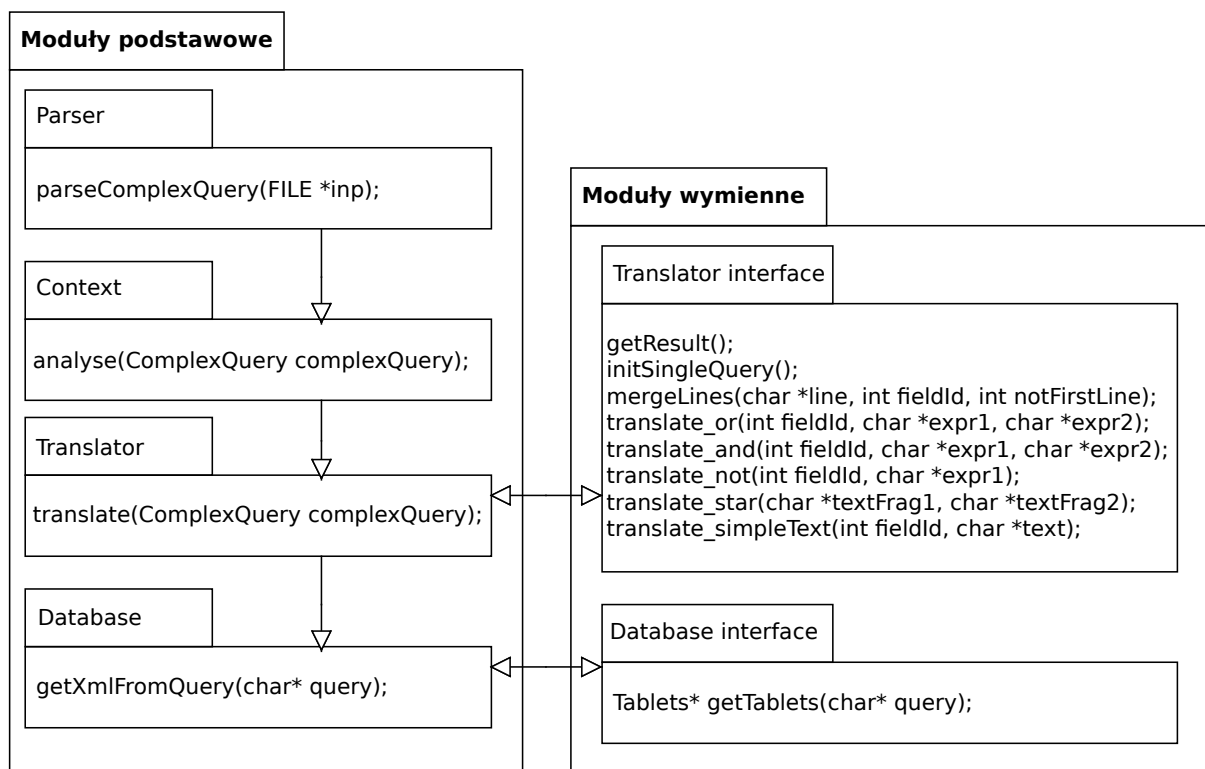
- spełniają wszystkie warunki zapytania ”zwierzęta w Gar*a”
- zawierają słowo ”adad-tilati”

Część III

Implementacja

Jednym z głównych założeń języka TQL, szczególnie ważnym z punktu widzenia implementacji, jest niezależność od struktury danych. W związku z tym istotną cechą programu tłumaczącego zapytania w TQL (translatora) jest możliwość dostosowania go do współpracy z różnymi bazami danych. Wynikiem tego jest podział translatora na 2 rodzaje modułów (rysunek 6.1):

1. **Podstawowe** – niezależne od struktury danych, zajmujące się głównie parsowaniem i analizą składniową zapytania.
2. **Wymienne** – zależne od struktury danych, tłumaczące zapytanie z TQL na język odpowiedni dla używanej bazy danych i wywołujące je.



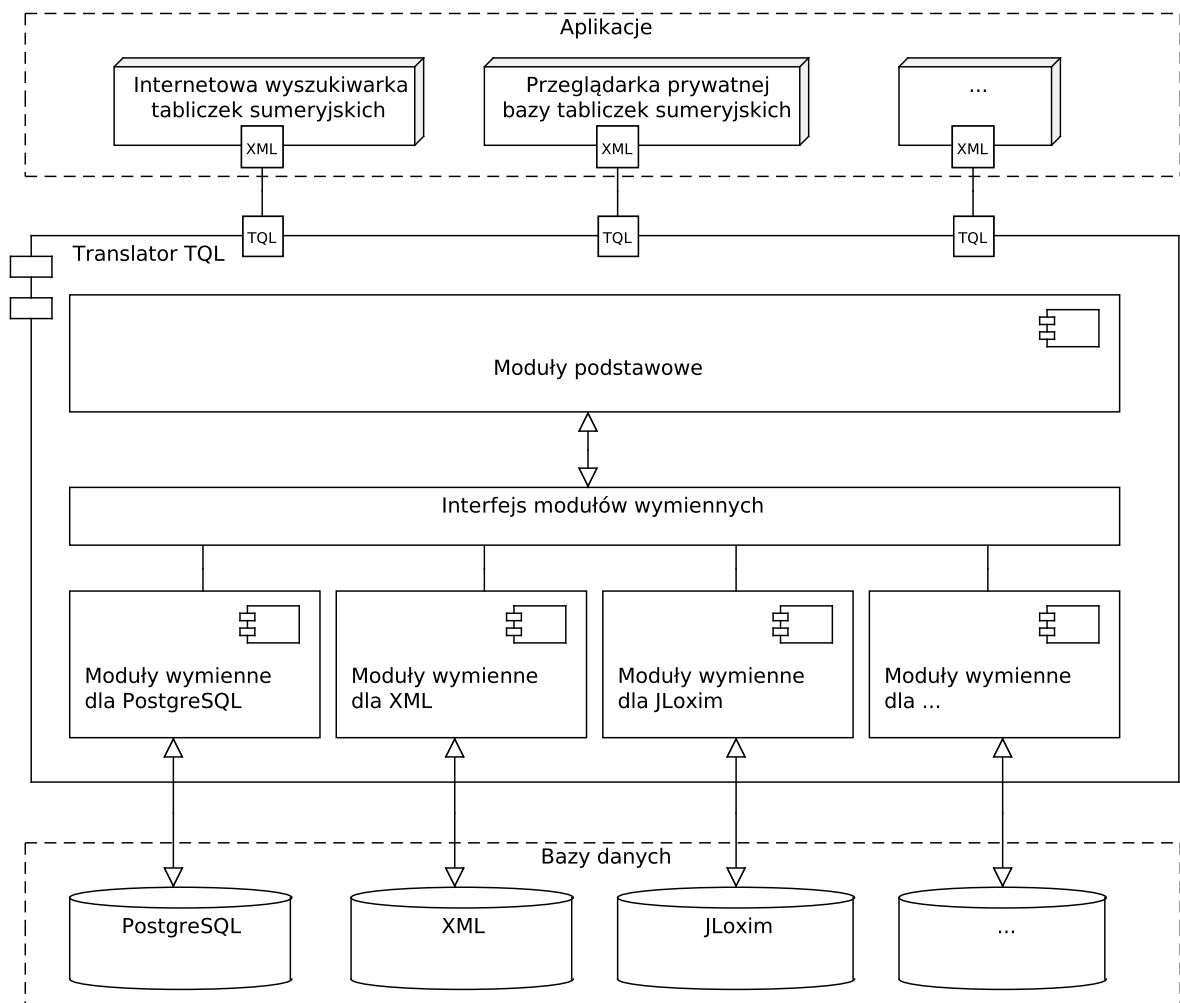
Rysunek 6.1: Podział programu na moduły

Na rysunku 6.2 przedstawiamy ogólnie strukturę systemu, który korzysta z translatora TQL. Taki system, oprócz modułów podstawowych translatora, zawiera bazę tabliczek w dowolnej formie, implemetację modułów wymiennych umożliwiającą korzystanie z tej bazy, oraz interfejs użytkownika - do wprowadzania zapytań i wyświetlania wyników.

W niniejszej pracy prezentujemy dwie prototypowe implementacje translatora, zawierające własne zestawy modułów wymiennych:

- dla bazy PostgreSQL
- dla bazy XML

Wybór modułów wymiennych odbywa się na poziomie kompilacji. Stworzony przez nas Makefile domyślnie buduje obie implementacje translatora TQL.



Rysunek 6.2: Struktura systemu korzystającego z tłumacza

Rozdział 7

Moduły podstawowe

W tym rozdziale przedstawimy dokładniej poszczególne moduły podstawowe tłumacza.

7.1. Parser

Parser został utworzony za pomocą narzędzia BNFC [5]. Następnie zostały w nim wprowadzone modyfikacje:

- poprawienie nazw stałych oznaczających symbole na bardziej intuicyjne,
- dodanie tablicy symboli,
- usunięcie niepotrzebnych funkcji z interfejsu,
- uporządkowanie kodu.

Moduł ten parsuje zapytanie w języku TQL, tworząc drzewo struktury składniowej, które jest zdefiniowane w pliku pomocniczym Absyn.h. Główna funkcja tego modułu to `parseComplexQuery(FILE *inp)`. Argumentem jest wskaźnik do pliku zawierającego zapytanie TQL, a wynikiem odpowiednie drzewo struktury składniowej. Na parser składają się następujące pliki:

- Parser.cpp
- Parser.h
- TQL.y
- TQL.l

7.2. Analizator kontekstowy

Podstawową funkcją w tym module jest `analyse(ComplexQuery complexQuery)`, której argumentem oraz wynikiem jest drzewo struktury składniowej zapytania TQL. Funkcja:

- sprawdza, czy podano prawidłowe nazwy pól, na podstawie których następuje wyszukiwanie,
- upraszcza drzewo - z wywołania zapytania (wywołanie *search in*) tworzy zapytanie proste.

Składa się z następujących plików:

- Context.cpp
- Context.h

7.3. Translator

Zadaniem tłumacza jest przetłumaczenie drzewa składni abstrakcyjnej na zapytanie w docelowym języku. Składa się z następujących plików:

- Translator.cpp
- Translator.h
- Translator.interface.h (interfejs modułu tłumacza zależnego od bazy danych)

Tłumaczenie poszczególnych elementów drzewa zależy od implementacji interfejsu zawartego w pliku Translator.interface.h. Funkcja translate(ComplexQuery complexQuery) przechodzi całą strukturę drzewa, wywołując w razie potrzeby odpowiednie funkcje z Translator.interface. Następnie pobiera przetłumaczone zapytanie za pomocą funkcji getResult() i przekazuje je jako wynik.

7.4. Baza

Moduł bazy jest odpowiedzialny za wywołanie przetłumaczonego zapytania i przekazanie wyniku w określonej formie - jako XML. Składa się z następujących plików:

- Database.cpp
- Database.h
- Database.interface.h (interfejs modułu bazy zależnego od bazy danych)

Główną funkcją w tym module jest getXmlFromQuery(char *query), która wywołuje funkcję getTablets(char *query) z Database.interface.h, jako parametr podając przetłumaczoną treść zapytania. Dostaje w wyniku strukturę danych Tablets, wypełnioną informacjami o wyszukanych tabliczkach. Następnie na podstawie otrzymanej struktury tworzy dokument XML i przekazuje go jako wynik wywołania zapytania.

Struktura Tablets zawiera metadane tabliczki, jej treść oraz informację o frazach, na podstawie których tabliczka została znaleziona. Poniżej przedstawiamy definicję struktury Tablets:

```
typedef struct{
    char* id;
    char* id_cdli;
    char* publication;
    char* measurements;
    char* year;
    char* provenience;
    char* period;
    char* genre;
    char* subgenre;
```



```

    char* collection;
    char* text;
    Tags* tags; // specjalnie oznaczone miejsca w tekście - frazy wyszukiwania
} Tablet;

typedef struct{
    int size;
    Tablet* tabs;
} Tablets;

```

7.5. Pliki pomocnicze

Definicje struktur danych i funkcji służących do budowy drzewa struktury składniowej (wygenerowane za pomocą BNFC [\[5\]](#), następnie uproszczone):

- Absyn.cpp
- Absyn.h

Tablica symboli:

- Symbols.cpp
- Symbols.h

Obsługa błędów:

- Err.cpp
- Err.h

Moduł do dzielenia tekstu względem separatora, pobrany z internetu [\[6\]](#):

- Cexplode.cpp
- Cexplode.h

Rozdział 8

Moduły wymienne

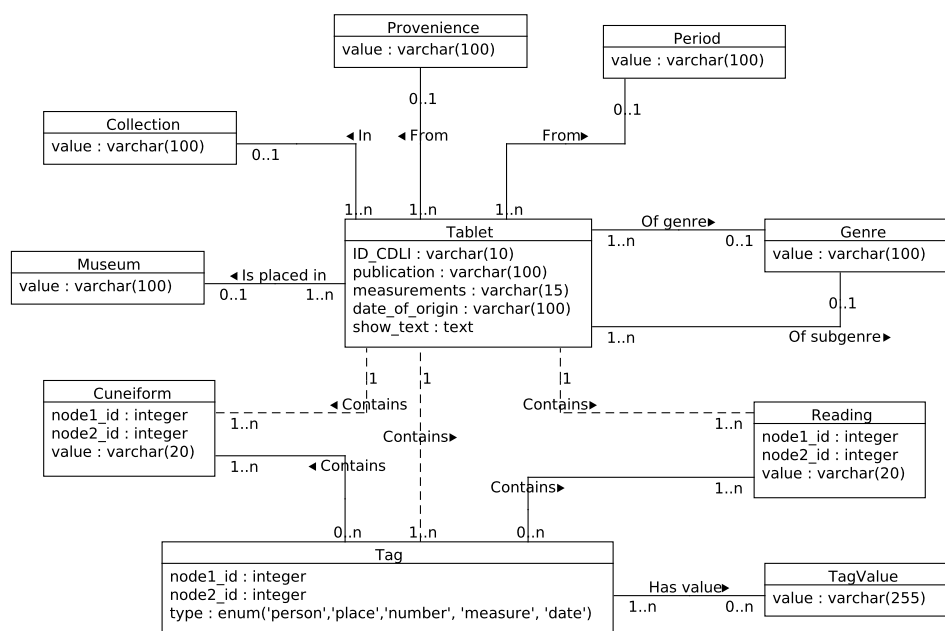
Pliki zależne od wyboru konkretnej bazy danych to:

- Translator_<nazwa>.cpp - dla modułu translatora
- Database_<nazwa>.cpp - dla modułu bazy

Ich interfejsy są wspólne dla wszystkich baz danych.

8.1. Baza PostgreSQL

8.1.1. Diagram encji



Rysunek 8.1: Diagram encji

Jednym z problemów przy projektowaniu bazy danych był wybór takiej reprezentacji treści tabliczki, żeby efektywnie wyszukiwać zarówno treść konkretnej tabliczki jak i tabliczki, których treść spełnia kryteria (wg odczytów lub zapisu klinowego).

W przypadku pierwszego problemu jako rozwiązanie narzuca się przechowywanie treści jako otwarty tekst. Natomiast najlepszym rozwiązaniem drugiego jest reprezentacja treści tabliczki w formie grafu, którego krawędziami są odczyty, kliny i tagi (zgodnie z pomysłem dr Wojciecha Jaworskiego[2, s.13-24]). Zdecydowałyśmy się na połączenie obu sposobów. Odczyty, kliny i tagi przechowujemy w tabelach Reading, Cuneiform i Tag, natomiast otwarty tekst w kolumnie show_text tabeli Tablet. Aby zapewnić możliwość odwzorowania treści tabliczki między reprezentacjami, węzły są liczbami postaci:

$\langle \text{numer wężła w tabliczce} \rangle * 1\,000\,000 + \langle \text{id tabliczki} \rangle$

gdzie numer wężła w tabliczce to numer kolejnego słowa (słowa są oddzielone spacjami i końcem linii) pomnożone przez 10 (żeby umożliwić wstawienie kilku wężłów w jednym słowie np. pozwolić na przetłumaczenie jednego słowa na sekwencję trzech klinów).

Przerywane linie na diagramie encji oznaczają opisany powyżej związek pomiędzy id wężła (node1_id i node2_id) a id tabliczki (Tablet.id).

8.1.2. Translator_postgres

Tłumaczy otrzymane fragmenty drzewa struktury zapytania na język SQL. Przetłumaczone fragmenty zbiera do buforów (*select*, *from*, *where*), które następnie odpowiednio łączy. Każde proste zapytanie TQL jest tłumaczone na pojedyncze zapytanie SQL. Tłumaczenie kilku prostych zapytań łączone jest za pomocą UNION.

Stałe fragmenty zapytania

Tłumaczenie prostego zapytania zaczyna się od inicjalizacji buforów przechowujących poszczególne części wynikowego SQL-a.

select jest inicjowany na

```
SELECT t.id, t.id_cdli, t.publication, t.measurements, t.origin_date,  
       p.value as provenience, pd.value as period,  
       g1.value as genre, g2.value as subgenre,  
       c.value as collection, t.text
```

from jest inicjowany

```
FROM tablet t  
  LEFT JOIN provenience p ON p.id = t.provenience_id  
  LEFT JOIN collection c ON c.id = t.collection_id  
  LEFT JOIN genre g1 ON g1.id = t.genre_id  
  LEFT JOIN genre g2 ON g2.id = t.subgenre_id  
  LEFT JOIN period pd ON pd.id = t.period_id
```

where początkowo zawiera pusty ciąg znaków.

Tłumaczenie zapytań o atrybuty tabliczki

Poniższe tłumaczenia są dodawane do bufora *where* i łączone za pomocą AND.

Konstrukcja	Tłumaczenie na SQL
provenience: wartosc	p.value LIKE 'wartosc'
publication: wartosc	t.publication LIKE 'wartosc'
period: wartosc	pd.value LIKE 'wartosc'
year: wartosc	t.origin_date LIKE 'wartosc'

Konstrukcja	Tłumaczenie na SQL
genre: wartosc	<pre> g1.value LIKE 'wartosc' OR g2.value LIKE 'wartosc' </pre>
cdli_id: wartosc	<pre> t.cdli_id LIKE 'wartosc' </pre>
museum: wartosc	<pre> t.museum LIKE 'wartosc' </pre>
collection: wartosc	<pre> c.value LIKE 'wartosc' </pre>

Tłumaczenie operatorów:

Operator	Tłumaczenie
/	OR
–	NOT
+	AND
*	%

Tłumaczenie zapytań o treść tabliczki

Przy tłumaczeniu zapytań o treść tabliczki korzystamy z przedstawienia treści tabliczki w formie grafu.

Pojawienie się wyszukiwania po treści tabliczki niesie za sobą konieczność dodania do bufora *from*:

```

INNER JOIN (
  <wynikowe zapytanie o treść tabliczki>
) AS sequence ON sequence.id_tab = t.id

```

natomiast do *select* dodajemy:

```
, sequence.nodes as nodes
```

gdzie <wynikowe zapytanie o treść tabliczki> to kombinacja zapytań typu:

```

SELECT
  id_tab,
  CAST(array_accum(nodes) as TEXT) as nodes,
  COUNT(DISTINCT id_seq) AS seq,
  <id_sekw> AS id_seq
FROM (
  SELECT
    t1.node1_id % 1000000 AS id_tab,

```

```

    '{' || t1.node1_id || ',' || t<dl_sekw>.node2_id || '}' AS nodes,
    1 AS id_seq
FROM
    <nazwa_tabeli> t1
    LEFT JOIN <nazwa_tabeli> t2 ON (t2.node1 = t1.node2)
    LEFT JOIN <nazwa_tabeli> t3 ON (t3.node1 = t2.node2)
    ...
    LEFT JOIN <nazwa_tabeli> t<dl_sekw> ON (t<dl_sekw>.node1 = t<dl_sekw-1>.node2)
WHERE
    t1.value LIKE '<sekw[1]>'
    AND
    t2.value LIKE '<sekw[2]>'
    AND
    t3.value LIKE '<sekw[3]>'
    AND
    ...
    AND
    t<dl_sekw>.value LIKE '<sekw[<dl_sekw>]>'
) AS a
GROUP BY id_tab

```

Zmienne użyte w powyższym pseudo-kodzie:

id_sekw - kolejny numer sekwencji (przydatny przy bardziej skomplikowanym zapytaniu - do rozróżniania podzapytań)

dl_sekw - ilość słów składających się na wyszukiwaną sekwencję

sekw - tablica zawierająca słowa składające się na wyszukiwaną sekwencję

nazwa_tabeli - nazwa tabeli, w której wyszukiujemy (Reading lub Cuneiform)

Operator	Tłumaczenie
/	<pre> SELECT id_tab, CAST(array_accum(nodes) as TEXT) as nodes, COUNT(DISTINCT id_seq) as seq, <id_sekw> as id_seq FROM (<zapytanie1> UNION <zapytanie2>) as c GROUP BY id_tab </pre>

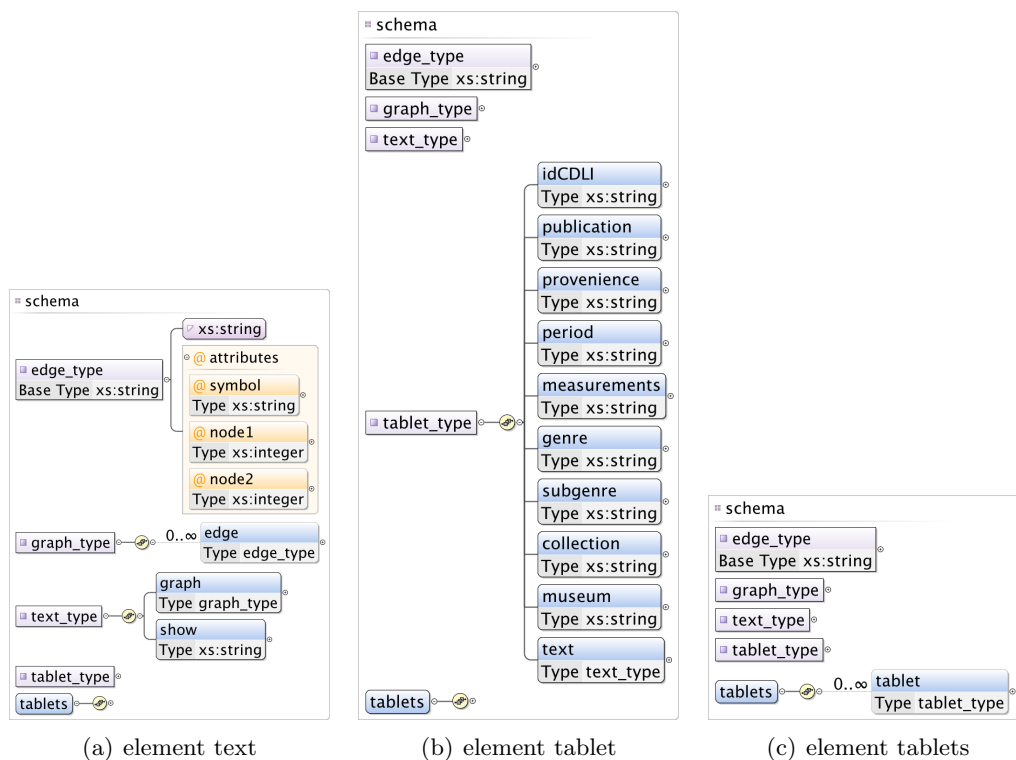
Operator	Tłumaczenie
+	<pre> SELECT * FROM (SELECT id_tab, CAST(array_accum(nodes) as TEXT) as nodes, COUNT(DISTINCT id_seq) as seq, <id_sekw> as id_seq FROM (<zapytanie1> UNION <zapytanie2>)) as c GROUP BY id_tab) as b WHERE b.seq=2 </pre>
—	<pre> SELECT id_tab, '' as nodes, 0 as seq, <id_sekw> as id_seq FROM ((SELECT id as id_tab from tablet) EXCEPT (SELECT id_tab from <zapytanie_negowane> as a)) as b </pre>
*	%

8.1.3. Database_postgres

Odpowiada za wywołanie zapytania w konkretnej bazie i zapisanie wyniku do struktury Tablets. Korzysta z pliku database.conf, który zawiera dane dostępu do bazy (nazwa bazy, host, port, użytkownik, hasło) oraz biblioteki libpq-fe.h do PostgreSQL.

8.2. Baza XML

8.2.1. Schemat dokumentu



Rysunek 8.2: Schematy poszczególnych elementów (kolor pomarańczowy oznacza atrybuty, a niebieski elementy)

Schemat dokumentu zamieszczony powyżej jest graficznym przedstawieniem dokumentu XML Schema (dodatek A) wygenerowanym przez program Oxygen.

Jest on oparty, podobnie jak w bazie postgres, na pomysłe dra Wojciecha Jaworskiego, aby przedstawić treść tabliczki w formie grafu. Każda krawędź tego grafu (odpowiadająca odczytowi) jest oddzielnym elementem, zawierającym atrybuty *node1*, *node2* i *symbol*. Atrybuty *node1* oraz *node2* oznaczają numery węzłów grafu, natomiast atrybut *symbol* to typ symbolu znajdującego się na danej krawędzi. Dodatkowo przechowujemy treść tabliczki w formie napisu (element *show*).

Metadane tabliczek przechowywane są w postaci podelementów elementu *tablet*.

8.2.2. Translator xml

Do przeszukiwania dokumentu XML wykorzystywany jest język XQuery, który jest częścią rekomendacji W3C dotyczącej XML.

Proste zapytanie TQL jest tłumaczone na pojedynczą konstrukcję FLWOR (For Let Where Order by Return).

Stałe fragmenty zapytania

Każde zapytanie w części For zawiera:

```
FOR $tablet IN .//tablet
```

a w części Return:

```
RETURN <tablet>
{$tablet/idCDLI}
{$tablet/publication}
{$tablet/provenience}
{$tablet/period}
{$tablet/measurements}
{$tablet/genre}
{$tablet/subgenre}
{$tablet/collection}
{$tablet/museum}
{$tablet/text/show}
<seq>...</seq>
</tablet>
```

Zawartość elementu seq zależy od ilości sekwencji, po których wyszukujemy.

Tłumaczenie zapytań o atrybuty tabliczki

Konstrukcja	Tłumaczenie na XQuery
provenience: wartosc	<code>fn:matches(\$tablet/provenience,'^wartosc\$')</code>
publication: wartosc	<code>fn:matches(\$tablet/publication,'^wartosc\$')</code>
period: wartosc	<code>fn:matches(\$tablet/period,'^wartosc\$')</code>
genre: wartosc	<code>(fn:matches(\$tablet/genre,'^wartosc\$')</code> <code>or fn:matches(\$tablet/subgenre,'^wartosc\$'))</code>
cdli_id: wartosc	<code>fn:matches(\$tablet/idCDLI,'^wartosc\$')</code>

Tłumaczenie zapytań o treść tabliczki

Każda sekwencja, po której wyszukujemy powoduje dodanie do zapytania następujących konstrukcji:

- do części Let:

```
let $seq <id_sekw> := (  
  for $edge_end in $tablet//edge  
  for $edge_start in $tablet//edge  
  where (  
    fn:matches($edge_start,'^<sekw[0]>$')  
    and (  
      some $edge1 in $tablet//edge[@node1=$edge_start/@node2]  
      satisfies (fn:matches($edge1,'^<sekw[1]>$')  
      and ...  
      and fn:matches($edge_end,'^<sekw[dl_sekw-1]>$')))))  
return <seq<id_sekw>> {$edge_start/@node1} {$edge_end/@node2} </seq<id_sekw>>
```

- do części Where

\$seq<id_sekw>

- do części Return w elemencie seq

\$seq<id_sekw>

Tłumaczenie operatorów

Poniższe tłumaczenia dotyczą zarówno konstrukcji prostych jak i złożonych.

Operator	Tłumaczenie
/	(<zapytanie1> or <zapytanie2>)
—	not (<zapytanie_negowane>)
+	(<zapytanie1> and <zapytanie2>)
*	.*

Zapytania złożone

Zapytanie złożone, składające się z wielu zapytań prostych tłumaczymy na sekwencję zapytań XQuery połączonych znakiem ','.

8.2.3. Database_xml

Odpowiada za wywołanie zapytania i zapisanie wyniku do struktury Tablets. Jako bazę danych wykorzystujemy plik XML, określony w pliku konfiguracyjnym xml.conf. Do wyszukiwania wykorzystujemy procesor XQuery Zorba [8]. Posiada on API m.in. do C++, które pozwala na przekazanie zapytania do bazy oraz przetworzenie wyniku.

Podsumowanie

Podsumowanie projektu

Przedstawiona w niniejszej pracy realizacja języka TQL spełnia najważniejsze postulaty. Po pierwsze jest on intuicyjny i prosty w użyciu dla osób znających jedynie dziedzinę problemu. Po drugie minimalnie ogranicza siłę wyrazu pozwalając na tworzenie skomplikowanych zapytań.

Język TQL ma jednak kilka ograniczeń w stosunku do typowych języków zapytań. Największym z nich jest brak wpływu na postać wyniku, co utrudnia np. zbieranie danych statystycznych (m. in. nie da się spytać o ilość tabliczek spełniających dane kryteria). Można jednak stworzyć narzędzia do samej prezentacji wyników zapytań, które pokonają to ograniczenie.

Dzięki specyficznym cechom języka TQL, można również dostosować go do wykorzystania w innych dziedzinach. Należy tu wspomnieć o braku ograniczenia ilości i nazw pól, po których można wyszukiwać oraz przystosowaniu do wyszukiwania wg kryteriów dotyczących danych tekstowych. Na podstawie TQL można łatwo stworzyć rodzinę języków dla różnych dziedzin zajmujących się przeszukiwaniem tekstów.

W ramach niniejszej pracy zaprezentowane zostały dwie implementacje języka TQL. Po stworzeniu pierwszej nakład czasu poświęconego na drugą był już niewielki. Największym zadaniem było przetłumaczenie poszczególnych konstrukcji TQL na XQuery. Pozwala to sądzić, że każda kolejna implementacja, dla różnych typów baz i różnych schematów danych, nie będzie wymagała dużego wysiłku.

Możliwości rozwoju

W pierwszej kolejności chcemy stworzyć własną stronę internetową z wyszukiwarką tabliczek, która byłaby dostępna dla Wydziału Historii UW i rozwijana z pomocą z naukowców z tej instytucji. Poza tym planujemy nawiązać współpracę z projektem CDLI, (zob. rozdział 3.). Implementacja TQL dla bazy CDLI, wraz z interfejsem www, z pewnością byłaby przydatnym narzędziem dla sumerologów na całym świecie.

W dalszej perspektywie można rozwinąć język TQL między innymi o wyszukiwanie wg klinów i wg tagów. Dodatkowo warto stworzyć narzędzie do analizy uszkodzonych fragmentów i wykrywania błędów w odczytach tekstów. Dzięki niemu można by zmniejszyć ryzyko pominięcia istotnych tabliczek podczas wyszukiwania.

Rozwój projektu będzie zależał również od potrzeb zgłaszanych przez samych sumerologów.

Dodatki

Dodatek A

Schemat bazy danych xml

```

<?xml version="1.0"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">

  <xs:complexType name="edge_type">
    <xs:simpleContent>
      <xs:extension base="xs:string">
        <xs:attribute name="symbol" type="xs:string"/>
        <xs:attribute name="node1" type="xs:integer"/>
        <xs:attribute name="node2" type="xs:integer"/>
      </xs:extension>
    </xs:simpleContent>
  </xs:complexType>

  <xs:complexType name="graph_type">
    <xs:sequence>
      <xs:element name="edge" minOccurs="0" maxOccurs="unbounded" type="edge_type"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="text_type">
    <xs:sequence>
      <xs:element name="graph" type="graph_type"/>
      <xs:element name="show" type="xs:string"/>
    </xs:sequence>
  </xs:complexType>

  <xs:complexType name="tablet_type">
    <xs:sequence>
      <xs:element name="idCDLI" type="xs:string"/>
      <xs:element name="publication" type="xs:string"/>
      <xs:element name="provenience" type="xs:string"/>
      <xs:element name="period" type="xs:string"/>
      <xs:element name="measurements" type="xs:string"/>
      <xs:element name="genre" type="xs:string"/>
      <xs:element name="subgenre" type="xs:string"/>
      <xs:element name="collection" type="xs:string"/>
      <xs:element name="museum" type="xs:string"/>
      <xs:element name="text" type="text_type"/>
    </xs:sequence>
  </xs:complexType>

  <xs:element name="tablets">
    <xs:complexType>
      <xs:sequence>
        <xs:element name="tablet" type="tablet_type" minOccurs="0" maxOccurs="unbounded"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>
</xs:schema>

```

Bibliografia

- [1] Agata Powalka, Marek Stępień, Jerzy Tyszkiewicz, *Automatyczna analiza dokumentów sumeryjskich*, Czas i przestrzeń, <http://www.czasiprzestrzen.wuw.pl/?id=str,sumeryjskich,1,0>
- [2] Wojciech Jaworski, *Modelowanie treści sumeryjskich tekstów gospodarczych z epoki Ur III*, <http://nlp.ipipan.waw.pl/NLP-SEMINAR/071119.pdf>, 19 listopada 2007
- [3] Cuneiform Digital Library Initiative, <http://cdli.ucla.edu>
- [4] The Electronic Text Corpus of Sumerian Literature, <http://etcsl.orinst.ox.ac.uk/>
- [5] BNF Converter, <http://www.ohloh.net/p/12800>
- [6] C-explode, <http://maz-programmersdiary.blogspot.com/2008/09/c-explode.html>
- [7] XQuery, <http://www.w3.org/XML/Query/>
- [8] Zorba - procesor XQuery, <http://www.zorba-xquery.com/>

Spis rysunków

1.	Przedstawienie problemu	7
2.1.	Gliniana tabliczka - struktura	13
2.2.	Co powinna zawierać cyfrowa reprezentacja tabliczki	14
3.1.	Formularz wyszukiwania na stronie CDLI	15
3.2.	Formularz wyszukiwania na stronie ETCSL	16
6.1.	Podział programu na moduły	27
6.2.	Struktura systemu korzystającego z translatora	28
8.1.	Diagram encji	34
8.2.	Schematy poszczególnych elementów (kolor pomarańczowy oznacza atrybuty, a niebieski elementy)	39