

Języki formalne i automaty

JFA

Marcin Kubica

e-mail: kubica@mimuw.edu.pl

2006

Niniejsze materiały powinny być pierwszym źródłem informacji dotyczących przedmiotu *Języki formalne i automaty* (JFA). Czytelnikom, którzy oprócz lektury tych notatek, chcieliby sięgnąć do podręcznika, polecam w pierwszej kolejności książkę D. Kozena, *Automata and Computability*. Jest ona napisana przystępnie i zawiera mnóstwo zadań, choć nie jest (jak dotąd) dostępna w języku polskim. W drugiej kolejności polecam książkę J. Hopcrofta i J. Ullmana, *Wprowadzenie do teorii automatów, języków i obliczeń*. Oprócz tematów wchodzących w zakres tego kursu można w niej znaleźć też wiele informacji na zaawansowane pokrewne tematy. Ponadto została ona przetłumaczona na język polski. Obie te książki nie zawierają informacji na temat generatorów analizatorów leksykalnych i składniowych. W pierwszej kolejności polecam tutaj dokumentację Flex'a i Bison'a. Dodatkowo, można przeczytać odpowiednie rozdziały z książki A. Aho, R. Sethi i J. D. Ullmana, *Kompilatory*.

Do każdego wykładu dołączono szereg zadań. W przypadku studiów internetowych, rozwiązania zadań z p. *Praca domowa* należy wykonywać w określonych terminach i zgłaszać poprzez serwis [edu](http://edu.mimuw.edu.pl). Pozostałe zadania mają charakter uzupełniający i są przeznaczone do samodzielnego rozwiązywania.

Materiały te są również dostępne w formacie PS i PDF.

Uwagi dotyczące niniejszych materiałów proszę przysyłać na adres: kubica@mimuw.edu.pl.

Literatura

- [HU] John E. Hopcroft, Rajeev Motwani, Jeffrey D. Ullman *Wprowadzenie do teorii automatów, języków i obliczeń*, PWN, 2005r.
- [K] Dexter C. Kozen, *Automata and Computability*,
- [ASU] A. V. Aho, R. Sethi, J. D. Ullman, *Kompilatory*, WNT, 2002,

Podziękowania

Niniejsze materiały powstały na podstawie notatek do prowadzonych przeze mnie, na przestrzeni kilku lat, wykładów z teorii języków i automatów oraz budowy kompilatorów. Chciałbym gorąco podziękować moim kolegom, którzy w tym czasie prowadzili ćwiczenia z tych przedmiotów, a w szczególności: Łukaszowi Krzeszczakowskiemu, Łukaszowi Maśko, Marcinowi Sawickiemu i Tomkowi Waleniowi. Ich uwagi miały wpływ na postać prowadzonych zajęć, a więc i również na te materiały. W szczególności część zamieszczonych tu zadań pochodzi od nich.

Szczególnie gorąco chciałbym podziękować Tomkowi Waleniowi za pomoc przy konwertowaniu niniejszych materiałów na format html i przygotowanie umieszczonych tu elementów multimedialnych.

Materiały te zostały opracowane w PJWSTK w projekcie współfinansowanym ze środków EFS.



1 Sylabus

W ramach kursu przedstawiane są klasy języków tworzące hierarchię Chomsky'ego: języki regularne, bezkontekstowe, kontekstowe i częściowo obliczalne (rekurencyjnie przeliczalne). Dodatkowo przedstawiona jest klasa języków obliczalnych.

Dla każdej z tych klas przedstawione są formalizmy służące do opisu języków: automaty skończone deterministyczne, niedeterministyczne i z ε -przejściami, wyrażenia regularne, automaty stosowe, gramatyki bezkontekstowe, gramatyki kontekstowe, maszyny Turinga oraz gramatyki ogólne (typu 0). Przedstawione są też fakty dotyczące przynależności języków do poszczególnych klas, w tym: lematy o pompowaniu dla języków regularnych i bezkontekstowych.

Program wykładów

1. Wprowadzenie, podstawowe pojęcia.
2. Wyrażenia regularne i wzorce.
3. Analiza leksykalna i Flex/Lex.
4. Deterministyczne automaty skończone.
5. Niedeterministyczne automaty skończone.
6. Równoważność automatów skończonych i wyrażeń regularnych.

7. Lemat o pompowaniu dla języków regularnych.
8. Minimalizacja deterministycznych automatów skończonych.
9. Języki i gramatyki bezkontekstowe.
10. Postać normalna Chomsky'ego i algorytm Cocke-Younger'a-Kasami.
11. Lemat o pompowaniu dla języków bezkontekstowych.
12. Automaty stosowe.
13. Analiza składniowa i Bison/Yacc.
14. Maszyny Turinga i obliczalność.
15. Języki obliczalne, częściowo obliczalne i nieobliczalne.

Program ćwiczeń

1–2 Flex/Lex - ćwiczenia.

3–4 Bison/Yacc - ćwiczenia.

Kryteria zaliczeń (dla studiów internetowych)

- prace domowe (przez Internet),
- egzamin (w uczelni).

Podstawą do wystawienia oceny będą punkty zdobyte za prace domowe i punkty zdobyte w egzaminie. Udział w egzaminie jest obowiązkowy, a robienie prac domowych gorąco zalecane.

Wymagane oprogramowanie

- Flex/Lex,
- Bison/Yacc,
- kompilator C/C++.

W przypadku Linuksa (zalecany) są to: `flex`, `bison`, `gcc`, `g++`. W przypadku Windows są to: DJGPP (z `gcc`), pakiety `flex` i `bison` z GnuWin32.

Powiązania merytoryczne

Nazwa przedmiotu poprzedzającego	Wymagany materiał
Matematyka dyskretna	Podstawy teorii mnogości, ciągi, relacje i funkcje.
Programowanie I i II	Podstawy programowania, podstawowe typy danych, programowanie w C lub C++.
Algorytmy i struktury danych	Złożoność asymptotyczna, podstawowe struktury danych (stosy), programowanie dynamiczne.

Wykład 2. Wprowadzenie

Celem niniejszego wykładu jest:

- poznanie podstawowych środków opisu składni (wyrażenia regularne, gramatyki bezkontekstowe),
- poznanie narzędzi wspomagających przetwarzanie danych o określonej składni (Lex, Yacc).
- zrozumienie podstawowych pojęć dotyczących obliczalności i odpowiedzenie sobie na kilka podstawowych pytań jej dotyczących:
 - Co to znaczy, że dany problem może zostać rozwiązany za pomocą programu komputerowego?
 - Czy istnieją problemy obliczeniowe, których nie można rozwiązać za pomocą programu komputerowego?
 - W jakim stopniu to co możemy obliczyć zależy od modelu obliczeń czy dostępnych konstrukcji programistycznych?

Nie są to proste pytania.

W trakcie kolejnych wykładów będziemy poznawać kolejne klasy języków tworzące tzw. hierarchię Chomsky'ego. Hierarchia ta jest owocem prac Noama Chomsky'ego, lingwisty, który próbował sformalizować pojęcia *gramatyki* i *języka*. Hierarchia ta składa się z czterech rodzajów gramatyk:

0. gramatyki ogólne,
1. gramatyki kontekstowe,
2. gramatyki bezkontekstowe,

3. gramatyki liniowe.

W trakcie badań nad pojęciem obliczalności pojawiło się wiele modeli, o różnej sile wyrazu:

1. automaty skończone, wyrażenia regularne (stała skończona pamięć),
2. automaty stosowe (stała skończona pamięć + nieograniczony stos),
3. bez ograniczeń:
 - maszyny Turinga,
 - systemy Posta,
 - rachunek λ , i inne.

Niektóre z nich powstały na długo przed tym zanim powstały komputery i informatyka. Okazało się, że modele te w zadziwiający sposób odpowiadają kolejnym klasom w hierarchii Chomsky'ego:

- gramatyki liniowe = automaty skończone = wyrażenia regularne,
- gramatyki bezkontekstowe = automaty stosowe,
- gramatyki kontekstowe = liniowo ograniczone maszyny Turinga,
- gramatyki ogólne = maszyny Turinga.

Nie są to czysto teoretyczne formalizmy. Niektóre z nich mają wiele praktycznych zastosowań. Szczególny nacisk położymy na dwa najważniejsze z tych formalizmów: wyrażenia regularne i gramatyki bezkontekstowe. W szczególności służą do specyfikowania analizatorów leksykalnych i składniowych — typowych modułów pojawiających się wszędzie tam, gdzie program wczytuje dane o określonej składni. Poznamy też narzędzia wspomagające tworzenie takich analizatorów — generatory, które same tworzą kod źródłowy analizatorów na podstawie ich specyfikacji.

2.1 Języki jako problemy decyzyjne

We wprowadzeniu wspomnieliśmy o *obliczeniach* i *językach*. Gdy mówimy o obliczaniu, to zwykle mamy na myśli obliczenie wartości takiej czy innej funkcji. Gdy zaś mówimy o języku, to mamy na myśli (potencjalnie nieskończony) zbiór napisów złożonych ze znaków ustalonego (skończonego) alfabetu. (Napisy takie będziemy nazywać *słowami*.) Jak połączyć te dwa pojęcia?

Def. 1. *Problem decyzyjny* to funkcja, która możliwym danym wejściowym przyporządkowuje wartości logiczne tak/nie. □

Inaczej mówiąc problem decyzyjny, to taki problem, w którym szukany wynik to wartość logiczna. Na problem decyzyjny możemy też patrzeć jak na zbiór danych — zbiór tych danych, dla których odpowiedź brzmi „tak”. I odwrotnie, na dowolny zbiór danych możemy patrzeć jak na problem decyzyjny — czy dane należą do określonego zbioru?

Jak to jednak ma się do języków? Jeśli interesuje nas tylko składnia języka (a nie jego semantyka, czyli znaczenie), to każdy język możemy przedstawić sobie jako zbiór (być może nieskończony) słów. Tak więc każdy język nie jest niczym więcej niż problemem decyzyjnym określonym dla słów. Powyższe intuicje są przedstawione bardziej formalnie poniżej.

W trakcie tego kursu będziemy zajmować się m.inn. tym, dla jakich języków (gdy spojrzymy na nie jak na problemy decyzyjne) można skonstruować automatyczne mechanizmy odpowiadające na pytanie, czy dane słowo należy do języka.

2.2 Podstawowe pojęcia dotyczące słów

Def. 2. *Alfabet* to dowolny niepusty, skończony zbiór. Elementy alfabetu nazywamy *znakami*. □

Alfabetem może to być np. zbiór cyfr dziesiętnych, zbiór znaków ASCII, czy zbiór bitów $\{0, 1\}$. Alfabet będziemy zwykle oznaczać przez Σ , a elementy alfabetu przez a, b, c, \dots

Def. 3. *Słowo* (lub *napis*) nad alfabetem Σ , to dowolny skończony ciąg znaków z Σ . W szczególności, pusty ciąg jest też słowem (i to nad dowolnym alfabetem). Oznaczamy go przez ε . □

Słowa będziemy zwykle reprezentować przez zmienne x, y, z, \dots

W przypadku języków programowania mówimy zwykle o napisach. Natomiast w teorii języków używa się raczej terminu słowo. Dodatkowo, w odniesieniu do języków naturalnych bardziej właściwe niż „słowo” wydaje się „zdanie”. Tak naprawdę jednak chodzi o to samo pojęcie.

Def. 4. *Długość* słowa oznaczamy przez $|\cdot|$. □

Na przykład $|ala| = 3$, $|\varepsilon| = 0$.

Def. 5. *Sklejenie* (konkatenacja) słów to słowo powstałe z połączenia słów, tak jakby zostały one zapisane kolejno po sobie. Sklejanie zapisujemy pisząc po sobie sklejaną słowa. □

Na przykład, jeżeli $x = ala$, $y = ma$, $z = kota$, to $xyz = alamakota$. Sklejanie słów jest łączne, tzn. $x(yz) = (xy)z$, a ε jest elementem neutralnym sklejaną, tzn. $x\varepsilon = \varepsilon x = x$.

Sposób zapisu sklejaną słów przypomina zwyczajowy zapis iloczynów (bez znaku mnożenia). Przez analogię, możemy wprowadzić operacje „potęgowania” słów:

Def. 6. Przez a^n oznaczamy n -krotne powtórzenie znaku a . Podobnie, przez x^n oznaczamy n -krotne powtórzenie słowa x . □

Na przykład, $a^5 = aaaaa$, $a^0 = \varepsilon$, $a^{n+1} = a^n a$, oraz $x^0 = \varepsilon$, $x^{n+1} = x^n x$. Jeżeli $x = ala$, to $x^3 = alaalaala$.

Def. 7. Przez $\#_a(x)$ oznaczamy liczbę wystąpień znaku a w słowie x . □

Na przykład, $\#_a(ala) = 2$.

Def. 8. Prefiksem słowa nazywamy dowolny jego początkowy fragment, tzn. x jest prefiksem y wtw., gdy istnieje takie z , że $xz = y$. Jeżeli ponadto $x \neq \varepsilon$ i $x \neq y$, to mówimy, że x jest właściwym prefiksem y .

Analogicznie sufiksem słowa nazywamy dowolny jego końcowy fragment, tzn. x jest sufiksem y wtw., gdy istnieje takie z , że $zx = y$. Jeżeli ponadto $x \neq \varepsilon$ i $x \neq y$, to mówimy, że x jest właściwym sufiksem y .

Podstawem danego słowa nazywamy dowolny jego spójny fragment, tzn. powiemy, że x jest podstawem y wtw., gdy istnieją takie v i w , że $vxw = y$. □

Na przykład, *kaj* i *kaja* są prefiksami (właściwymi) słowa *kajak*, natomiast *jak* jest jego sufiksem (właściwym).

Def. 9. Przez $\text{rev}(x)$ oznaczamy słowo x czytane wspak. Na przykład $\text{rev}abbab = babba$.

Jeżeli słowo x czytane wprost i wspak jest takie samo (tzn. $x = \text{rev}(x)$), to mówimy, że x jest *palindromem*. □

2.3 Podstawowe pojęcia dotyczące języków

Przejdziemy teraz o jeden poziom wyżej i zajmijmy się zbiorami słów czyli językami.

Def. 10. *Język* (nad alfabetem Σ), to dowolny zbiór słów (nad alfabetem Σ). □

Jeżeli alfabet Σ jest znany z kontekstu, to będziemy go czasami pomijać. Języki będziemy zwykle oznaczać przez A, B, C, \dots . Język złożony ze wszystkich możliwych słów nad alfabetem Σ będziemy oznaczać przez Σ^* .

Język, to podzbiór zbioru Σ^* . Będziemy więc (w kolejnych wykładach) utożsamiać język z problemem decyzyjnym, dla zbioru danych Σ^* .

Skoro języki to zbiory słów (napisów), a więc określone są na nich wszystkie podstawowe operacje na zbiorach:

- \emptyset to pusty język, który nie zawiera żadnego słowa,
- $A \cup B$ to suma języków A i B ,
- $A \cap B$ to część wspólna (przecięcie) języków A i B ,
- $A \setminus B$ to różnica języków A i B .

Na przykład, $(\{ala, ola, ula\} \cap \{abba, ala, ula\}) \setminus \{ula, bula\} = \{ala\}$.

Przyda nam się też kilka operacji charakterystycznych dla języków:

Def. 11. \bar{A} to dopełnienie języka A , czyli $\bar{A} = \Sigma^* \setminus A$. □

Na przykład $\overline{\Sigma^*} = \emptyset$.

Def. 12. AB oznacza sklejanie (konkatenację) języków A i B , czyli język zawierający wszystkie możliwe sklejania słów z A ze słowami z B , $AB = \{xy : x \in A, y \in B\}$. □

Na przykład $\{a, ab\}\{b, bb\} = \{ab, abb, abbb\}$.

Zauważmy, że sklejanie języków jeszcze bardziej przypomina mnożenie, niż to miało miejsce w przypadku sklejania słów. Elementem neutralnym (czyli „jedyką”) sklejania języków jest język zawierający tylko słowo puste $\{\varepsilon\}$, $\{\varepsilon\}A = A\{\varepsilon\} = A$. Natomiast „zerem” sklejania języków jest język pusty, $\emptyset A = A\emptyset = \emptyset$. W odróżnieniu od operacji mnożenia, sklejanie języków nie jest przemienne.

Sklejanie języków jest rozdzielne względem sumowania języków:

$$A(B \cup C) = AB \cup AC$$

Jest tak dlatego, że każde sklejanie słowa należącego do A ze słowem należącym do $(B \cup C)$ jest sklejaniem słowa należącego do A ze słowem należącym do B , lub sklejaniem słowa należącego do A ze słowem należącym do C , i vice versa. Analogicznie:

$$(A \cup B)C = AC \cup BC$$

Mając zdefiniowane sklejanie języków, możemy analogicznie do potęgowania słów zdefiniować potęgowanie języków:

Def. 13. Język A^n definiujemy rekurencyjnie:

- $A^0 = \{\varepsilon\}$,
- $A^{n+1} = A^n A$.

Czyli $A^n = \underbrace{A \dots A}_{n \text{ razy}}$. □

Inaczej mówiąc, A^n to język zawierający wszystkie możliwe sklejania n (niekoniecznie różnych) słów wziętych z A .

Języki, jako zbiory, mogą być nieskończone. Możemy więc rozważać sklejanie dowolnej liczby słów pochodzących z A . Prowadzi to do tzw. domknięcia Kleene’ego:

Def. 14. Domknięcie Kleene’ego języka A , oznaczane jako A^* , to $A^* = \bigcup_{n \geq 0} A^n$. Inaczej mówiąc, A^* to język zawierający wszystkie możliwe sklejania dowolnej liczby (łącznie z 0) słów należących do A . □

Dla każdego języka A mamy $\varepsilon \in A^*$, gdyż $\{\varepsilon\} = A^0 \subseteq A^*$. Jeżeli tylko A zawiera jakieś niepuste słowo, to A^* jest zbiorem nieskończonym. Na przykład, $\{ab, aa\}^* = \{\varepsilon, ab, aa, abab, abaa, aaab, aaaa, \dots\}$.

Jeżeli spojrzymy na Σ nie jak na alfabet, ale jak na zbiór słów długości 1, to wyjaśni się dlaczego Σ^* jest zbiorem wszystkich słów nad alfabetem Σ — po prostu każde słowo jest sklejaniem pewnej liczby znaków.

Czasami będzie nas interesowało sklejanie dowolnej, ale dodatniej, liczby słów z A :

Def. 15. Przez A^+ oznaczamy język zawierający wszystkie możliwe sklejania dowolnej dodatniej liczby słów należących do A , $A^+ = AA^*$. \square

2.4 Powtórzenie pojęć dotyczących relacji

Przypomnijmy sobie kilka pojęć dotyczących relacji.

Def. 16. Niech X będzie dowolnym zbiorem, a $\rho \subseteq X \times X$ relacją (binarną) określoną na tym zbiorze. Powiemy, że relacja ρ jest:

- zwrotna, jeżeli dla każdego $x \in X$ zachodzi $x\rho x$,
- symetryczna, gdy dla dowolnych $x, y \in X$ jeżeli mamy $x\rho y$, to mamy również $y\rho x$,
- przechodnia, gdy dla dowolnych $x, y, z \in X$ jeżeli mamy $x\rho y$ i $y\rho z$, to mamy również $x\rho z$,
- antysymetryczna, gdy dla dowolnych $x, y \in X$ jeśli $x\rho y$ i $y\rho x$, to $x = y$,
- relacją równoważności, jeśli jest zwrotna, symetryczna i przechodnia,
- częściowym porządkiem, jeśli jest zwrotna, przechodnia i antysymetryczna.

\square

Relacja równoważności dodatkowo dzieli zbiór X na rozłączne klasy abstrakcji.

Def. 17. Jeśli ρ jest relacją równoważności, $x \in X$, to przez $[x]_\rho$ oznaczamy klasę abstrakcji x :

$$[x]_\rho = \{y \in X : x\rho y\}$$

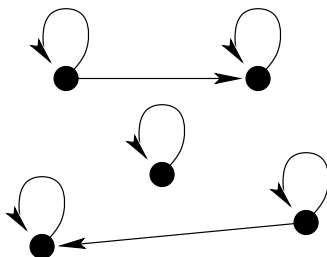
Przez X/ρ oznaczamy zbiór klas abstrakcji elementów zbioru X :

$$X/\rho = \{[x]_\rho : x \in X\}$$

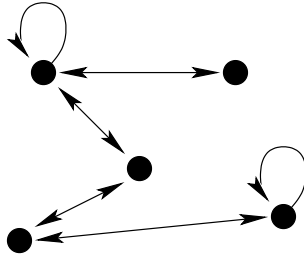
\square

Pojęcia te są bardziej intuicyjne, gdy przedstawimy sobie relację jako graf skierowany (z pętelkami). Wyobraźmy sobie, że elementy zbioru X , to wierzchołki grafu. Z wierzchołka x prowadzi krawędź do wierzchołka y wtedy i tylko wtedy, gdy $x\rho y$. (W szczególności, gdy $x\rho x$, to mamy „pętelkę” prowadzącą z x do x .)

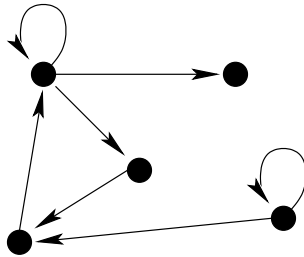
Relacja jest zwrotna, gdy w każdym wierzchołku jest „pętelka”.



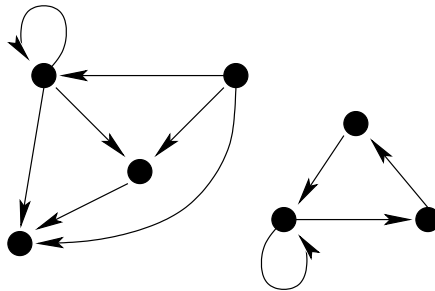
Relacja jest symetryczna, gdy krawędzie między (różnymi) wierzchołkami są dwukierunkowe (tzn. jeżeli jest krawędź w jednym kierunku, to jest i w drugim).



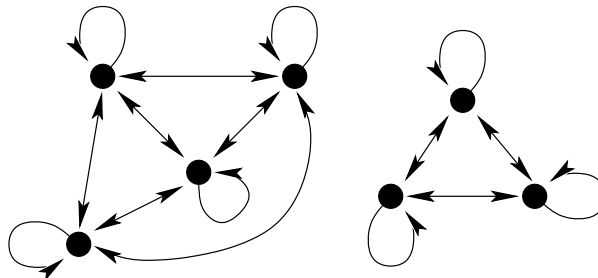
Relacja jest antysymetryczna, gdy krawędzie między (różnymi) wierzchołkami mogą być tylko jednokierunkowe.



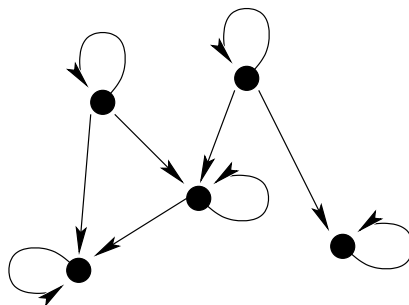
Relacja jest przechodnia, jeżeli dla każdej ścieżki w grafie (długości przynajmniej 1) istnieje w nim krawędź łącząca początek ścieżki z jej końcem. Inaczej mówiąc, graf musi zawierać wszystkie możliwe krawędzie idące „na skróty”.



Relacja jest relacją równoważności, jeżeli graf jest podzielony na ileś spójnych składowych, każda składowa to klika (tzn. między każdymi dwoma wierzchołkami tej samej składowej mamy krawędź), natomiast między wierzchołkami należącymi do różnych składowych nie mamy krawędzi. Spójne składowe takiego grafu są nazywane *klasami abstrakcji*.



Relacja jest częściowym porządkiem, jeżeli nie zawiera cykli (z wyjątkiem pęteli, które są we wszystkich wierzchołkach).



Pojęcia zwrotności, przechodniości i symetryczności wszystkie polegają na tym, że pewne pary/krawędzie muszą być obecne. Będziemy mówić o odpowiednim *domknięciu* relacji, jako o relacji powstałej przez dodanie odpowiednich par/krawędzi, niezbędnych do spełnienia określonej własności. Szczególnie przydatne będzie nam domknięcie zwrotno-przechodnie.

Def. 18. Niech $\rho \subseteq X \times X$ będzie dowolną relacją binarną. *Domknięciem zwrotno-przechodnim* relacji ρ nazywamy najmniejszą taką relację $\rho' \subseteq X \times X$, która:

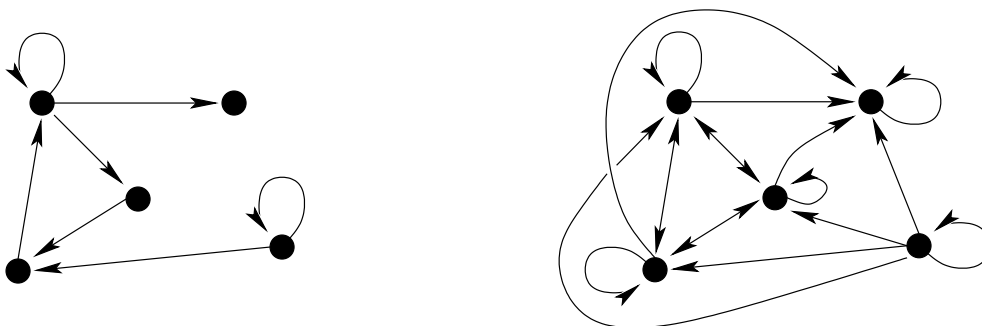
- zawiera relację ρ , $\rho \subseteq \rho'$,
- jest zwrotna i przechodnia.

□

Podobnie definiujemy domknięcie przechodnie, czy symetryczne.

Intuicyjnie, domknięcie zwrotno-przechodnie polega na dodaniu wszystkich takich krawędzi, które w oryginalnym grafie mogły być realizowane przez ścieżki — łącznie z pętelkami, które odpowiadają ścieżkom długości 0.

Przykład: Relacja i jej domknięcie zwrotno-przechodnie:



2.5 Podsumowanie

W tym wykładzie poznaliśmy podstawowe pojęcia dotyczące słów i języków. Przypomnieliśmy też podstawowe informacje na temat relacji, które będą nam potrzebne w dalszych wykładach.

2.6 Skorowidz

- **Alfabet** to dowolny niepusty, skończony zbiór. Elementy alfabetu nazywamy znakami.
- **Domknięcie Kleene’ego** języka A , oznaczane jako A^* , to $A^* = \bigcup_{n \geq 0} A^n$. Inaczej mówiąc, A^* to język zawierający wszystkie możliwe sklejania dowolnej liczby (łącznie z 0) słów należących do A .
- **Domknięcie zwrotno-przechodnie** relacji ρ to najmniejsza taka relacja $\rho' \subseteq X \times X$, która: zawiera relację ρ , $\rho \subseteq \rho'$, oraz jest zwrotna i przechodnia.
- **Dopełnienie języka** A , oznaczane jako \bar{A} , to $\bar{A} = \Sigma^* \setminus A$.
- **Język** (nad alfabetem Σ), to dowolny zbiór słów (nad alfabetem Σ).
- **Palindrom** to takie słowo x , które czytane wprost i wspak jest takie samo, $x = \text{rev}(x)$.
- **Podsłowo** danego słowa to dowolny jego spójny fragment, tzn. powiemy, że x jest pod słowem y wtw., gdy istnieją takie v i w , że $v x w = y$.
- **Prefiks** słowa to dowolny jego początkowy fragment, tzn. x jest prefiksem y wtw., gdy istnieje takie z , że $x z = y$. Jeżeli ponadto $x \neq \varepsilon$ i $x \neq y$, to mówimy, że x jest *właściwym* prefiksem y .
- **Problem decyzyjny** to funkcja, która możliwym danym wejściowym przyporządkowuje wartości logiczne tak/nie.
- **Sklejenie (konkatenację) języków** to język zawierający wszystkie możliwe sklejania słów z A ze słowami z B , $AB = \{xy : x \in A, y \in B\}$.
- **Sklejenie (konkatenacja) słów** to słowo powstałe z połączenia słów, tak jakby zostały one zapisane kolejno po sobie.
- **Słowo** (lub napis) nad alfabetem Σ , to dowolny skończony ciąg znaków z Σ .
- **Sufiks** słowa to dowolny jego końcowy fragment, tzn. x jest sufiksem y wtw., gdy istnieje takie z , że $z x = y$. Jeżeli ponadto $x \neq \varepsilon$ i $x \neq y$, to mówimy, że x jest *właściwym* sufiksem y .

2.7 Praca domowa

1. Jaka jest różnica między \emptyset i ε ? Ile słów zawierają te języki?
2. Podaj które prefiksy słowa $ababbbabab$ są równocześnie jego sufiksami.
3. Podaj wszystkie słowa należące do języka $\{aba, ba, a\}\{ba, baba\}$.
4. Podaj wszystkie słowa należące do języka $\{a, ab, ba, baba\} \cap \{ab, ba\}^2$.
5. Podaj wszystkie słowa należące do języka $\{a, ab, ba\} \setminus \{a\}^*\{bb\}^*\{a\}^*$.

2.8 Zadania

1. Operacje na zbiorach. Niech $A, B \subset X$. Narysować diagram przedstawiający te zbiory. Które są sobie równe? Które z nich można uszeregować zgodnie z zawieraniem.

- $A \cup B$,
- $A \cap B$,
- X ,
- A ,
- B ,
- \emptyset ,
- $X \setminus A$,
- $A \cap B \cap (X \setminus A)$,
- $(A \cup B) \setminus (A \cap B)$,
- $X \setminus (A \cup B)$,
- $X \setminus (A \cap B)$,
- $(A \setminus B) \cup (B \setminus A)$,
- $(X \setminus A) \cup (X \setminus B)$,
- $(X \setminus A) \cap (X \setminus B)$,
- $X \setminus ((X \setminus A) \cup (X \setminus B))$,
- $X \setminus ((X \setminus A) \cap (X \setminus B))$.

2. Wybieramy dowolny numer telefonu i traktujemy go jak słowo złożone z cyfr. Podaj:

- wszystkie prefiksy tego ciągu,
- wszystkie sufiksy tego ciągu,
- kilka podciągów (niekoniecznie spójnych); ile jest wszystkich?

- kilka spójnych podciągów (podslów); ile jest wszystkich?
3. Wybieramy zbiór 2–3 elementowy. Jakie są jego podzbiory? Ile podzbiorów ma zbiór n -elementowy?
4. Które z poniższych tożsamości są prawdziwe? W przypadku prawdziwych uzasadnij, a w przypadku fałszywych podaj przykład słowa, które przeczy tożsamości.
- $(A \cup B) \cup C = A \cup (B \cup C)$,
 - $(A \cup B) \cap C = A \cup (B \cap C)$,
 - $(A \cup B) \cap C \subseteq A \cup (B \cap C)$,
 - $A(B \cup C) = AB \cup AC$,
 - $A(B \cap C) = AB \cap AC$,
 - $A(B \cap C) \subseteq AB \cap AC$,
 - $(AB)^* = A^*B^*$,
 - $(A \cup B)^* = (A^*B^*)^*$,
 - $A^{*+} = A^{+*}$,
 - $\overline{(A \cup B)} = \overline{A} \cap \overline{B}$.
5. Co to za języki:

- $(A \cap A^*)^*$, Rozwiązanie¹
- $A^* \setminus A^+$, Rozwiązanie²
- $(A \cap B)^* \setminus (A^* \cup B^*)$, Rozwiązanie³
- $\{aa, ba, bb, ab\}^*$, Rozwiązanie⁴
- $(\Sigma\Sigma)^*$, Rozwiązanie⁵
- $\overline{(\Sigma\Sigma)^*}$, Rozwiązanie⁶
- $\Sigma(\Sigma\Sigma)^*$ Rozwiązanie⁷
- $\text{III}\Sigma\Sigma\text{A}$ Rozwiązanie⁸

¹ A^*

² Jeśli $\varepsilon \in A$, to \emptyset . W przeciwnym przypadku jest to $\{\varepsilon\}$.

³ Słowa powstałe ze sklejania słów z A i B , przy czym konieczne jest wzięcie słów z obu tych języków.

⁴ Słowa parzystej długości

⁵ Słowa parzystej długości

⁶ Słowa nieparzystej długości.

⁷ Również słowa nieparzystej długości.

⁸ A to jest po grecku pizza. :-)

6. Używając symboli i operacji wprowadzonych na pierwszym wykładzie (czyli operacji na słowach, zbiorach, językach i definicji zbiorów) zdefiniować języki:
- słowa zawierające tyle samo liter a i b ,
 - słowa nad alfabetem $\{a, b\}$ (nie)parzystej długości,
 - inne ...
7. Podać domknięcie zwrotno-przechodnie i domknięcie symetryczne relacji $\{(1, 2), (2, 3), (3, 1), (2, 4), (5, 6)\}$ określonej na zbiorze $\{1, 2, \dots, 6\}$.
8. Pokazać, że następujące relacje są relacjami równoważności:
- słowa x i y są w relacji wtw., gdy $|x| = |y|$,
 - słowa x i y są anagramami,
 - języki A i B są w relacji wtw., gdy $A^* = B^*$.
9. Pokazać, że następujące relacje są częściowymi porządkami:
- relacja zawierania \subseteq na językach nad ustalonym alfabetem,
 - relacja porządku leksykograficznego (alfabetycznego) na słowach,
 - relacja bycia prefiksem.
10. Relacja \prec jest określona na językach nad ustalonym alfabetem w następujący sposób: $A \prec B$ wtw., gdy istnieje taki język C , że $AC = B$. Udowodnij, że relacja \prec jest częściowym porządkiem.

Wykład 3. Wzorce i wyrażenia regularne

3.1 Wstęp

Wyrażenia regularne to rodzaj wzorców, do których dopasowuje się fragmenty różnych tekstów. Nazwy tej używa się w dwóch kontekstach. Pierwszy z nich dotyczy wyrażeń regularnych, jakie można znaleźć np. przy wyszukiwaniu tekstu, w bardziej rozbudowanych edytorach tekstów. Są one również szeroko stosowane w Uniksie, w różnego rodzaju programach do wyszukiwania i przetwarzania informacji tekstowych, np. `grep`, `awk` i `sed`, a także w generatorze analizatorów leksykalnych `Lex`.

Terminu tego używa się również w teorii języków formalnych w podobnym znaczeniu. Jednak w tym przypadku liczba konstrukcji jakie mogą być używane do budowy wyrażeń regularnych jest ograniczona do absolutnego minimum.

Jest to poniekąd naturalne, gdyż oba rodzaje wyrażeń regularnych służą innym celom. W pierwszym przypadku ma to być wygodny w użyciu język, pozwalający na zwięzłe zapisywanie wzorców. Stąd jest od dosyć rozbudowany i zawiera wiele konstrukcji. W drugim przypadku bardziej nas interesuje badanie samych wyrażeń regularnych i ich siły wyrazu. Stąd język tych wyrażeń regularnych jest maksymalnie uproszczony, a wszystkie konstrukcje, które można wyrazić w inny sposób zostały z niego usunięte.

W trakcie tego kursu będziemy używali obu rodzajów „wyrażeń regularnych”. Musimy więc nadać im różne nazwy. Wyrażenia regularne takie, jak są używane w `Lex`'ie, będziemy nazywać *wzorcami*. Natomiast sam termin „wyrażenia regularne” zarezerwujemy dla takiego jego znaczenia, jakie ma w teorii języków formalnych.

Na potrzeby tego rozdziału ustalamy jako domyślny alfabet Σ zbiór znaków ASCII.

3.2 Wzorce

Poniżej podajemy definicję wzorców. Z jednej strony jest ona wzorowana na wyrażeniach regularnych w `Lex`'ie, ale nie wymieniliśmy tutaj wszystkich dostępnych tam konstrukcji. Z drugiej strony dodaliśmy dwie konstrukcje, które nie są dostępne w `Lex`'ie. Definicja ta, oprócz składni wzorców, opisuje też nieformalnie ich semantykę. Semantyka ta jest bardziej formalnie zdefiniowana dalej.

Def. 19. *Wzorce*, to wyrażenia, które możemy budować w podany poniżej sposób.

- dla $a \in \Sigma$, a jest wzorcem, do którego pasuje tylko słowo a (chyba, że a ma jakieś specjalne znaczenie opisane poniżej),
- ε jest wzorcem, do którego pasuje tylko słowo puste (ta konstrukcja nie jest dostępna w `Lex`'ie),
- \emptyset jest wzorcem, do którego nie pasuje żadne słowo (ta konstrukcja nie jest dostępna w `Lex`'ie),
- $.$ to wzorzec, do którego pasuje dowolny znak (oprócz końca wiersza),

- “ x ” to wzorec, do którego pasuje tylko słowo x (nawet jak zawiera jakieś znaki o specjalnym znaczeniu),
- $\backslash a$ — ma analogiczne znaczenie jak w C, np. $\backslash n$ oznacza znak nowej linii,
- $[\dots]$ to wzorec, do którego pasuje dowolny ze znaków wymienionych w kwadratowych nawiasach, np. do $[abc]$ pasuje dowolny ze znaków a , b i c ,
- $[a - b]$ to wzorec, do którego pasuje dowolny znak od a do b , zakresy takie można łączyć, np. $[a - zA - Z]$ to wzorec, do którego pasuje dowolna litera, mała lub wielka,
- $[\hat{ \dots }]$ to wzorec, do którego pasuje dowolny znak oprócz znaków wymienionych wewnątrz kwadratowych nawiasów, np. do $[\hat{xyz}]$ pasuje dowolny znak oprócz x , y i z ,
- jeśli α jest wzorcem, to $\alpha?$ jest wzorcem, do którego pasuje ε oraz wszystkie te słowa, które pasują do α ,
- jeśli α jest wzorcem, to α^* jest wzorcem, do którego pasuje sklejenie zera lub więcej słów pasujących do α ,
- jeśli α jest wzorcem, to α^+ jest wzorcem, do którego pasuje sklejenie jednego lub więcej słów pasujących do α ,
- jeśli α i β są wzorcami, to $\alpha\beta$ jest wzorcem, do którego pasują sklejenia słów pasujących do α i słów pasujących do β ,
- jeśli α i β są wzorcami, to $\alpha | \beta$ jest wzorcem, do którego pasują te słowa, które pasują do α lub do β ,
- jeśli α i β są wzorcami, to $\alpha \cap \beta$ jest wzorcem, do którego pasują te słowa, które pasują równocześnie do α i β (ta konstrukcja nie jest dostępna w Lex’ie),
- jeśli α jest wzorcem, to $\bar{\alpha}$ jest wzorcem, do którego pasują wszystkie te słowa, które nie pasują do α (ta konstrukcja nie jest dostępna w Lex’ie),
- jeśli α jest wzorcem, to jest nim również (α) i pasują do niego te same słowa, co do α — inaczej mówiąc, do budowy wzorców możemy używać nawiasów,

□

Przykład: Oto garść przykładowych wzorców:

- $[0 - 9]^+$ opisuje (niepuste) ciągi cyfr, czyli zapisy dziesiętne liczb naturalnych.
- Do wzorca $(a^*bba^*) \cap (ab | ba)^*$ pasuje tylko słowo $abba$. Do a^*bba^* pasują słowa zawierające dokładnie dwie litery b i to położone obok siebie, a do $(ab | ba)^*$ pasują słowa zbudowane z „cegiełek” ab i ba . Jedynym słowem, które pasuje do obydwu tych wzorców jest właśnie $abba$.

- Wzorzec $[A - Z][A - Z][0 - 9][0 - 9][0 - 9][0 - 9][0 - 9]$ opisuje jedną z form numerów rejestracyjnych, złożonych z dwóch liter i pięciu cyfr.

Zdefiniujemy teraz bardziej formalnie semantykę wzorców.

Def. 20. Niech α i β będą wzorcami. Przez $L(\alpha)$ oznaczamy język opisywany przez wzorzec α , czyli język złożony z tych słów, które pasują do wzorca α . $L(\alpha)$ definiujemy indukcyjnie ze względu na budowę wzorca α :

- $L(a) = \{a\}$, dla $a \in \Sigma$, o ile a nie ma specjalnego znaczenia opisanego niżej,
- $L(\varepsilon) = \{\varepsilon\}$,
- $L(\emptyset) = \emptyset$,
- $L(\cdot) = \Sigma \setminus \{\backslash n\}$,
- $L("x") = \{x\}$,
- $L([a_1 a_2 \dots a_k]) = \{a_1, a_2, \dots, a_k\}$,
- $L([a - b]) = \{a, \dots, b\}$ — zależy od kolejności znaków w kodowaniu ASCII,
- $L([\hat{a} - b]) = \Sigma \setminus L([a - b])$
- $L(\alpha?) = L(\alpha) \cup \{\varepsilon\}$,
- $L(\alpha^*) = L(\alpha)^*$,
- $L(\alpha^+) = L(\alpha)^+ = L(\alpha)L(\alpha)^*$,
- $L(\alpha\beta) = L(\alpha)L(\beta)$,
- $L(\alpha \mid \beta) = L(\alpha) \cup L(\beta)$,
- $L(\alpha \cap \beta) = L(\alpha) \cap L(\beta)$,
- $L(\bar{\alpha}) = \overline{L(\alpha)}$,
- $L((\alpha)) = L(\alpha)$.

□

Badając wzorce często będziemy chcieli porównywać nie tyle je same, ale ich znaczenia. Przyda nam się do tego pojęcie równoważności wzorców.

Def. 21. Powiemy, że dwa wzorce α i β są równoważne, $\alpha \equiv \beta$, wtw., gdy $L(\alpha) = L(\beta)$. □

Oto garść tożsamości dotyczących wzorców. Niektóre z nich, tak naprawdę, poznaliśmy już w poprzednim wykładzie jako tożsamości dotyczące języków.

- $\alpha(\beta \mid \gamma) \equiv \alpha\beta \mid \alpha\gamma$ — to nic innego jak rozdzielność sklejanania języków względem ich sumowania, zastosowana do wzorców,
- $(\beta \mid \gamma)\alpha \equiv \beta\alpha \mid \gamma\alpha$ — j.w.,
- $\emptyset\alpha \equiv \alpha\emptyset \equiv \emptyset$ — \emptyset jest „zerem” sklejanania, również dla wzorców,
- $\alpha^+ \equiv \alpha \mid \alpha\alpha^+ \equiv \alpha \mid \alpha^+\alpha$ — oto dwie rekurencyjne definicje operacji $^+$,
- $(\alpha\beta)^*\alpha \equiv \alpha(\beta\alpha)^*$, — oba wzorce opisują naprzemienne ciągi α i β zaczynające i kończące się α -ą,
- $(\alpha^*)^* \equiv \alpha^*$ — domykanie domknięcia Kleene’ego nie wnosi nic dodatkowego,
- $(\alpha \mid \beta)^* \equiv (\alpha^*\beta^*)^*$ — ciąg α i β można zawsze przedstawić jako poprzepłatane ze sobą ciągi α i ciągi β (wliczając w to ciągi długości zero), i vice versa.

3.3 Wyrażenia regularne

Okazuje się, że wiele z konstrukcji występujących we wzorcach może być zastąpionych prostszymi konstrukcjami. Na przykład zamiast $[a - c]$ można napisać $a \mid b \mid c$. Postępując w ten sposób, każdy wzorec można przerobić na taki, który został zapisany przy użyciu tylko pewnego minimalnego zestawu konstrukcji. Wzorce, które możemy zapisać przy użyciu owego minimalnego zestawu operacji, nazywamy *wyrażeniami regularnymi*.

Def. 22. Wyrażenia regularne, to takie wzorce, które są zbudowane tylko przy użyciu:

- ε ,
- \emptyset ,
- a , dla $a \in \Sigma$,
- $\alpha\beta$, gdzie α i β to wyrażenia regularne,
- $\alpha \mid \beta$, gdzie α i β to wyrażenia regularne,
- α^* , gdzie α to wyrażenie regularne,
- (α) , gdzie α to wyrażenie regularne.

□

Jak zobaczymy dalej, wszystkie pozostałe konstrukcje są redundantne, choć w praktyce są przydatne. Wyrażenia regularne będą nam potrzebne wówczas, gdy będziemy się zajmowali tym, co można wyrazić za pomocą wzorców i ogólnymi właściwościami języków, które można opisywać za pomocą wzorców. Dzięki ograniczeniu zestawu możliwych konstrukcji do niezbędnego minimum, nasze rozważania będą prostsze.

3.4 Języki regularne

Klasę języków, dla których istnieją opisujące je wzorce nazywamy *językami regularnymi*.

Def. 23. Powiemy, że język A jest *regularny*, wtw., gdy istnieje wzorec α opisujący A , czyli $A = L(\alpha)$. \square

Jak później pokażemy, wzorce mają taką samą siłę wyrazu co wyrażenia regularne, tzn. jeśli dla danego języka istnieje opisujący go wzorec, to istnieje również opisujące go wyrażenie regularne.

Jeden z problemów, jakim będziemy się zajmować w trakcie tego kursu, to: jakie języki są regularne? Okazuje się, że nie wszystkie. Poniżej podajemy kilka prostych własności klasy języków regularnych.

Fakt 1. *Każdy język skończony jest regularny.*

Dowód: Niech $\{x_1, x_2, \dots, x_n\}$ będzie niepustym skończonym językiem. Język ten można opisać wzorcem (a dokładniej wyrażeniem regularnym) postaci $x_1 | x_2 | \dots | x_n$. Język pusty można opisać wzorcem \emptyset . \square

Fakt 2. *Niech A i B będą językami regularnymi. Języki $A \cup B$, $A \cap B$, AB , \bar{A} i A^* są też regularne.*

Dowód: Skoro A i B są regularne, to istnieją opisujące je wzorce α i β , $L(\alpha) = A$, $L(\beta) = B$. Interesujące nas języki są regularne, bo opisują je wzorce: $\alpha | \beta$, $\alpha \cap \beta$, $\alpha\beta$, $\bar{\alpha}$ i α^* . \square

3.5 Podsumowanie

W tym wykładzie poznaliśmy wzorce oraz wyrażenia regularne. Języki jakie można za ich pomocą opisać tworzą klasę języków regularnych.

3.6 Skorowidz

- **Język opisywany** przez wzorec to język złożony z tych słów, które pasują do wzorca α .
- **Język regularny**, to taki język, dla którego istnieje opisujący go wzorec.
- **Wyrażenia regularne** to szczególny przypadek wzorców, do których zapisu użyto wyłącznie: symboli alfabetu, ε , \emptyset , sklejania, $|$, $*$ i nawiasów.
- **Wzorce** to rodzaj wyrażeń, do których dopasowujemy słowa nad ustalonym alfabetem. Wzorec opisuje język złożony ze słów, które do niego pasują.

3.7 Praca domowa

1. Podaj wyrażenie regularne równoważne wzorcowi: $([a - c] | a(b?))^+$.
2. Podaj wzorzec opisujący poprawne numery indeksów studentów, np. postaci s0124.
3. Uprość następujący wzorzec $(a | b)(aa | ab | bb | ba)^*(a | b)$.

3.8 Ćwiczenia

1. Podaj wzorce/wyrażenia regularne opisujące język złożony ze słów:
 - (a) nad alfabetem $\{a, b\}$, które zawierają podsłowo $bbab$,
 - (b) nad alfabetem $\{a, b, c\}$, które zaczynają się i kończą tym samym znakiem,
 - (c) nad alfabetem $\{a, b\}$, które nie zawierają podsłowa aa ,
 - (d) nad alfabetem $\{a, b\}$, które nie zawierają podsłowa ab ,
 - (e) nad alfabetem $\{a, b, c\}$, które nie zawierają podsłowa aa ,
 - (f) nad alfabetem $\{a, b, c\}$, które nie zawierają podsłowa ab ,
 - (g) nad alfabetem $\{a, b, c\}$ złożonych tylko z jednego rodzaju symboli,
 - (h) nad alfabetem $\{a, b, c\}$ złożonych co najwyżej z dwóch rodzajai symboli,
 - (i) nad alfabetem $\{a, b\}$, które zawierają parzystą liczbę liter a .
2. Porównaj podane wzorce pod kątem równoważności / zawierania się odpowiadających im języków. Jeśli są równoważne, to uzasadnij to. Jeżeli nie, to podaj przykład słowa, które pasuje do jednego wzorca, ale nie do drugiego.
 - (a) $(a?b?)^*$ i $(a | b)^*$,
 - (b) $(a | b)^*$ i $a^* | b^*$,
 - (c) $(a^*b)^*$ i $(b^*a)^*$,
 - (d) $(a | b)^*$ i $(aa | ab | ba | bb)$,
 - (e) $aa(bbaa | baa)^*$ i $(aabb | aab)^*a$.

Wykład 4. Analiza leksykalna i generatory analizatorów leksykalnych

Wzorce możemy spotkać w wielu narzędziach — wszędzie tam, gdzie chcemy wyszukiwać pewne fragmenty tekstu, wyszukiwany fragment możemy opisać właśnie za pomocą wzorca:

- w edytorach tekstów często możemy podać nie tylko wyszukiwany tekst, ale właśnie wzorzec,
- programy systemowe służące do wyszukiwania informacji, np. `grep`, `sed`, czy `awk`,
- biblioteki programistyczne zawierające procedury rozpoznające wystąpienia wzorców w tekście,
- generatory analizatorów leksykalnych (skanerów).

Tym ostatnim zastosowaniem zajmiemy się bliżej. Najpierw jednak musimy przybliżyć pojęcie analizy leksykalnej.

4.1 Analiza leksykalna

Miejsce na analizę leksykalną jest wszędzie tam, gdzie wczytujemy dane o określonej składni. Wczytując takie dane, zanim będziemy mogli je przetwarzać, musimy rozpoznać ich składnię. Pomysł polega na tym, aby najpierw wczytywany ciąg znaków podzielić na elementarne cegiełki składniowe — nazywane *leksemami* — a dopiero dalej analizować ciąg leksemów. Analizator leksykalny (nazywany też skanerem) to wyodrębniony moduł zajmujący się tym zadaniem.

Analiza leksykalna powstała i rozwinęła się w ramach prac nad budową kompilatorów. Wszak kompilator musi najpierw wczytać i zanalizować składnię wczytywanego programu. Jej zastosowania są jednak dużo szersze. Praktycznie analizę leksykalną można zastosować w każdym programie, który wczytuje dane posiadające jakąś składnię, na przykład w: przeglądarkach internetowych, edytorach, systemach składu tekstu, programach konwertujących, czy jakichkolwiek aplikacjach posiadających pliki konfiguracyjne o określonej składni.

Czemu wyodrębnić analizę leksykalną jako osobny moduł? Jest kilka powodów:

- uproszczenie konstrukcji programu — łatwiej opracować osobno analizę leksykalną i resztę analizy składniowej,
- zwiększenie przenośności — w module analizie leksykalnej można ukryć wszystkie szczegóły związane ze sposobem wczytywania plików, reprezentacją znaków, itp.,
- zwiększenie efektywności — analiza leksykalna, choć koncepcyjnie prosta, może zajmować dużą część czasu pracy programu; właściwe jej zaimplementowanie może poprawić efektywność programów,

Napisanie efektywnego skanera nie jest proste. Nie musimy jednak tego robić ręcznie. Praktycznie dla każdego języka programowania istnieją narzędzia, tzw. generatory analizatorów leksykalnych, które zrobią to za nas. Musimy im jedynie dostarczyć specyfikacji opisującej jak wyglądają leksemy, owe elementarne cegiełki, na które chcemy podzielić wczytywany ciąg znaków i jakie jest ich znaczenie.

4.2 Leksemy, żetony i atrybuty

W trakcie analizy leksykalnej wczytywany ciąg znaków jest dzielony na leksemy. Jednak to co jest przekazywane dalej, to nie są dokładnie leksemy. Formalnie, leksem to ciąg znaków. To co jest przekazywane, to informacja reprezentująca znaczenie leksemu. Informacja ta jest reprezentowana za pomocą tzw. *żetonu* i opcjonalnego *atrybutu*. Żeton niesie informację o rodzaju leksemu. Jeżeli leksemy danego rodzaju niosą ze sobą pewną „wartość”, to żetonowi towarzyszy atrybut i jest on równy tej wartości. Podsumujmy więc znaczenie tych trzech terminów:

leksem to ciąg kolejnych znaków stanowiących semantycznie niepodzielną całość,

żeton (ang. *token*) to stała (całkowita) reprezentująca rodzaj wczytanego leksemu,

atrybut to opcjonalna wartość reprezentująca znaczenie leksemu.

Typowe leksemy, to:

- identyfikatory,
- słowa kluczowe,
- napisy,
- liczby całkowite i rzeczywiste,
- operacje arytmetyczne i relacje,
- nawiasy i innego rodzaju znaki „interpunkcyjne”.

Najlepiej powyższe pojęcia zilustrować na przykładach.

Przykład: Rozważmy instrukcję (z kompilowanego programu) postaci:

$$E := m * c^2;$$

Możemy ją rozbić na następujący ciąg leksemów:

$$E, :=, m, *, c, ^, 2, ;$$

Natomiast ich reprezentacja za pomocą żetonów i atrybutów będzie następująca:

Leksem	Zeton	Atrybut
E	identyfikator	„E”
:=	przypisanie	
m	identyfikator	„m”
*	mnożenie	
c	identyfikator	„c”
^	potęgowanie	
2	liczba całkowita	2
;	średnik	

Przykład: Rozważmy fragment źródeł HTML:

```
<head>
  <meta http-equiv="Content-Type" content="text/html; charset=iso-8859-2">
  <title>PJWSTK - JFA</title>
```

Jak ten fragment należy podzielić na leksemy? Pamiętajmy, że leksemy muszą być niepodzielne semantycznie. Oto jeden z możliwych podziałów:

Leksem	Zeton	Atrybut
<head	open-start-tag	„head”
>	close-tag	
<meta	open-start-tag	„meta”
http-equiv	identyfikator	„http-equiv”
=	równość	
"Content-Type"	napis	„Content-Type
content	identyfikator	„content”
=	równość	
"text/html; charset=iso-8859-2"	napis	„text/html; charset=iso-8859-2”
>	close-tag	
<title	open-start-tag	„title”
>	close-tag	
PJWSTK - JFA	text	„PJWSTK - JFA”
</title	open-end-tag	„title”
>	close-tag	

Jak widać z powyższych przykładów, typowe atrybuty leksemów to

- identyfikator — nazwa identyfikatora (jeśli nie odróżnia się małych i wielkich liter, to w znormalizowanej postaci, np. pisana samymi wielkimi literami),
- napis — treść napisu,

- liczba — jej wartość,

Pisząc specyfikację dla Lex'a musimy dla każdego żetonu opisać jaką postać mogą mieć leksemy odpowiadające temu żetonowi. Robimy to podając dla każdego żetonu wzorzec. Następnie dla każdego żetonu podajemy fragment kodu, który na podstawie leksemu oblicza wartość atrybutu.

Przykład: Oto kilka przykładów żetonów i wzorców opisujących odpowiadające im leksemy:

Żeton	Przykład leksemu	Wzorzec
sł_klucz_if	if	if
identyfikator (złożony tylko z liter)	Pom	[a - zA - Z] ⁺
liczba_całkowita	-42	-?[0 - 9] ⁺

Projektując podział wejścia na leksemy i dobierając żetony powinniśmy kierować się następującymi zasadami:

- Leksemy powinny być niepodzielne semantycznie.
- Leksemom, które mają tego samego rodzaju semantykę powinien odpowiadać ten sam żeton.
- Leksemy odpowiadające jednemu żetonowi powinno się dać opisać wzorcem.

Skaner jest zwykle zrealizowany w postaci modułu, który udostępnia procedurę „daj kolejny leksem”. Procedura ta rozpoznaje jeden leksem i zwraca odpowiadający mu żeton i atrybut. Można więc powiedzieć, że skaner przetwarza strumień znaków w strumień par: żeton, atrybut.

Proste atrybuty, takie jak liczby, są przekazywane wprost. Bardziej złożone atrybuty, takie jak nazwy identyfikatorów, mogą być przez skaner umieszczane w specjalnym słowniku, nazywanym tablicą symboli. Wówczas jako atrybut przekazywany jest wskaźnik lub pozycja w tablicy symboli. Z jednej strony, dzięki temu oszczędzamy na pamięci. Z drugiej strony, sama pozycja w tablicy symboli jest wystarczająca, gdyż zwykle nie jest istotne, jak się dany identyfikator nazywa, tylko który z identyfikatorów to jest.

4.3 Generator analizatorów leksykalnych Lex

Generator analizatorów leksykalnych na podstawie specyfikacji skanera sam generuje jego kod. Jest wiele różnych generatorów analizatorów leksykalnych, różnego pochodzenia i przeznaczonych dla różnych języków programowania. Skupiamy się na generatorze Flex/Lex przeznaczonym dla C/C++. Nazwa [F]lex to skrót od *[Fast] LEXical analyzer generator*.

Specyfikacja taka zawiera z jednej strony opis składni leksemów różnych rodzajów, a z drugiej zawiera fragmenty kodu, które mają być wykonane w przypadku napotkania leksemu określonego rodzaju.

Dla każdego rodzaju leksemów (tj. dla każdego żetonu) specyfikacja skanera zawiera:

- wzorzec opisujący postać leksemów (danego rodzaju),
- fragment kodu wykonywany w momencie napotkania leksemu danego rodzaju.

Fragmenty kodu wykonywane w momencie rozpoznania leksemu mają oczywiście dostęp do słowa (napisu) stanowiącego rozpoznany leksem. Ich głównym zadaniem jest przekazanie żetonu i obliczenie ew. atrybutu. Jednak z punktu widzenia Lex'a mogą one robić cokolwiek. Jeżeli np. naszym zadaniem jest przetłumaczenie jednej postaci danych na inną, takie fragmenty mogą wypisywać przetłumaczone odpowiedniki wczytanych leksemów.

4.3.1 Definicje nazwanych wzorców

W specyfikacji dla Lex'a zwykle pewne wzorce pojawiają się wielokrotnie. Żeby nie powtarzać ich, wprowadzono możliwość definiowania pomocniczych nazwanych wzorców. Definicje takie mają postać:

```
nazwa      wzorzec
```

Od momentu zdefiniowania nazwanego wzorca, można się do niego odwoływać w innych wzorcach poprzez: `{nazwa}`. W definicjach nazwanych wzorców można się odwoływać do nazwanych wzorców zdefiniowanych wcześniej, ale tylko wcześniej — tym samym definicje wzorców nie mogą być rekurencyjne.

Przykład:

```
cyfra      [0-9]
cyfry      {cyfra}+
cz_ułamkowa  "."{cyfry}
znak_opc   ("+"|-")?
wykładnik_opc (E{znak_opc}{cyfry})?
liczba     {znak_opc}({cyfry}{cz_ułamkowa}|{cyfry}|{cz_ułamkowa}){wykładnik_opc}
```

Definicje nazwanych wzorców pozwalają w zwięzły sposób powiązać nazwy z określonymi wzorcami. Nie wnoszą one jednak nic do siły wyrazu wzorców. Definicji nazwanych wzorców można się zawsze pozbyć rozwijając wszystkie definicje.

4.3.2 Specyfikacje skanerów dla Lex'a

Język specyfikacji dla Lex'a zawiera mnóstwo szczegółów, których nie będziemy tutaj omawiać. W tej chwili naszym zadaniem jest wprowadzenie Czytelnika w typowy sposób wykorzystania Lex'a. Gdy zawarte tu informacje okażą się niewystarczające, odsyłamy Czytelnika do dokumentacji generatora skanerów (`man lex`).

Specyfikacja dla Lex'a składa się z trzech części:

- deklaracji programistycznych (stałych i zmiennych), oraz definicji regularnych,
- reguł rozpoznawania leksemów,
- dodatkowych procedur zdefiniowanych przez użytkownika.

Części te są oddzielone od siebie wierszami postaci: „`%%`”.

Deklaracje programistyczne Dowolne globalne **deklaracje** programistyczne można umieścić ujmując je w nawiasy `{ ... }`. Dotyczy to również dyrektyw `#include`.

Definicje nazwanych wzorców Zapisane zgodnie z opisaną powyżej składnią.

Reguły rozpoznawania leksemów Reguły rozpoznawania leksemów mają postać:

wzorzec fragment kodu

Jeżeli fragment kodu jest dłuższy niż jedna instrukcja, to musi być ujęty w nawiasy `{ ... }`. Wzorzec musi zaczynać się na początku wiersza i nie może zawierać żadnych (nie ujętych w cudzysłowy) odstępów.

Dodatkowe procedury W trzeciej części można umieścić deklaracje własnych procedur. W szczególności można tam zdefiniować procedurę `main`.

4.3.3 Rozpoznawanie leksemów

Zanim zobaczymy przykłady specyfikacji, musimy zrozumieć, jak działają wygenerowane przez Lex'a skanery. Skaner udostępnia funkcję (`int yylex()`) służącą do wczytania jednego leksemu. Jeżeli plik został wczytany, to `yylex` zwraca 0. Wpp. stara się dopasować pewien fragment początkowy wejścia do jednego z podanych wzorców, wg. następujących zasad:

- Jeżeli istnieje kilka możliwych dopasowań leksemu, przede wszystkim wybierany jest jak najdłuższy leksem.

- Jeśli jest kilka wzorców pasujących do najdłuższego możliwego leksemu, to wybierany jest pierwszy z nich (zgodnie z kolejnością w specyfikacji). Następnie wykonywany jest fragment kodu odpowiadający dopasowanemu wzorcowi. Dopasowany leksem jest dostępny poprzez zmienną `ytext` i ma długość `ylen`. (Zmienna `ytext` może być typu `char ytext[]` lub `char *ytext`.)
- Jeśli wykonywany fragment kodu zakończy się poprzez `return`, to jest to powrót z wywołania procedury `yylex`. Wpp. rozpoczyna się dopasowywanie kolejnego leksemu itd.
Żeton powinien być zwracany właśnie jako wynik procedury `yylex`. Jeśli chcemy przekazywać atrybuty, to powinniśmy je przypisać (zadeklarowanym w pierwszej części specyfikacji) zmiennym.
- Jeżeli nic się nie udało dopasować, to jeden znak z wejścia jest przepisywany na wyjście i rozpoznawanie leksemów zaczyna się od początku.

Przykład: Oto specyfikacja skanera, który każdy napis „FOO” zamienia na „Foo”, „BAR” na „Bar” i „FOOBAR” na „Foobar”. Wszystko inne jest przepisywane z wejścia na wyjście bez zmian.

```
%{
#include <stdio.h>
%}

%%

FOO    printf("Foo");
BAR    printf("Bar");
FOOBAR printf("Foobar");

%%

int main() {
    yylex();
    return 0;
}
```

Przykład ten znajduje się w pliku `foobar.l`.

Przykład: Oto specyfikacja skanera, który wyszukuje w wejściu liczby zmiennopozycyjne i wypisuje je. Wszystko inne jest ignorowane.

```

%{
    #include <stdio.h>
    #include <math.h>
}%

cyfra      [0-9]
cyfry      {cyfra}+
cz_ułamkowa  "."{cyfry}
znak_opc   ("+-")?
wykładnik_opc (E{znak_opc}{cyfry})?
liczba     {znak_opc}({cyfry}{cz_ułamkowa}|{cyfry}|{cz_ułamkowa}){wykładnik_opc}

%%

{liczba}   { printf("%g", atof(yytext)); }
.          { /* nic */ }

%%

int main() {
    yylex();
    return 0;
}

```

Przykład ten znajduje się w pliku `liczby.l`.

Przykład: Następująca specyfikacja opisuje skaner, który usuwa komentarze postaci `/* ...*/`. Uwaga: nie uwzględnia ona innych rodzajów komentarzy, ani napisów.

```

%%

"/*"([^\*]|\n|("/*"[/]))*"*/"    { /* Nic */ }

"/*"                                { printf("Nieprawidłowy komentarz."); }

%%

int main() {
    yylex();
    return 0;
}

```

Przykład ten znajduje się w pliku `uncomment.l`.

W powyższych przykładach nie używaliśmy żetonów ani atrybutów. Jedno wywołanie funkcji `yylex` przetwarzało całe wejście. Poniższy przykład ilustruje wykorzystanie żetonów i atrybutów.

Przykład: Oto specyfikacja prostego skanera, który sumuje wszystkie liczby całkowite znajdujące się na wejściu, pomijając wszystko inne.

```
%{
    #include <stdio.h>
    #define LICZBA 1
    int attr;
}%

INTEGER -?[0-9]+

%%

{INTEGER} { attr = atoi(yytext);
           return LICZBA; }

.|\\n      { /* Nic */ }

%%

int main() {
    int sum=0;

    while (yylex() == LICZBA) sum += attr;
    printf("suma = %d\\n", sum);
    return 0;
}
```

Przykład ten znajduje się w pliku `sum.1`.

4.4 Podsumowanie

Wykład ten wyjaśnia czym jest analiza leksykalna i generatory analizatorów leksykalnych, takie jak `Lex`. Pokazuje też jak pisać specyfikacje analizatorów leksykalnych dla `Lex`'a. Umiejętność korzystania z `Lex`'a będzie jeszcze przedmiotem ćwiczeń laboratoryjnych.

4.5 Skorowidz

- **Analiza leksykalna** polega na wczytaniu strumienia znaków, podzieleniu go na leksemy i przekształceniu na żetony i ew. atrybuty.
- **Atrybut** (leksemu) to opcjonalna wartość reprezentująca znaczenie/wartość leksemu.
- **Definicje nazwanych wzorców** to część specyfikacji dla Lex'a wiążąca z określonymi nazwami określone wzorce.
- **Generator analizatorów leksykalnych** to program generujący, na podstawie specyfikacji, kod analizatora leksykalnego (skanera).
- **Leksem** to ciąg kolejnych znaków stanowiących semantycznie niepodzielną całość.
- **Lex** to generator analizatorów leksykalnych.
- **Skaner** (analizator leksykalny) to moduł zajmujący się analizą leksykalną.
- **Żeton** to stała reprezentująca rodzaj leksemu.

4.6 Praca domowa

1. Dla następującego fragmentu pliku ini:

```
[info]
drivername=ContentIndex
symbolfile=perfci.h
[languages]
009=English
015=Polish
[objects]
CIOBJECT_009_NAME=Indexing Service
CIOBJECT_015_NAME=Usługa indeksowania
```

podaj:

- jak byś go podzielił na leksemy,
 - jakie żetony byś wyodrębnił.
2. Napisz specyfikację skanera, który zamienia wszystkie słowa tak, że pierwsza litera jest wielka, a kolejne są małe. Wszystkie inne znaki, które nie tworzą słów, powinny być wypisywane bez zmian.

4.7 Ćwiczenia

1. Podaj leksemy pojawiające się w poniższych programach.

- ```
function abs (i : integer): integer;
{ Oblicza wartość bezwzględną i }
begin
 if i > 0 then abs := i else abs := -i
end;
```
- ```
int abs (i)
int i;
/* Oblicza wartość bezwzględną i */
{
  return i>0?i:-i;
}
```

2. Podaj wyrażenia wzorce opisujące następujące leksemy:

- numery telefoniczne: alarmowe i skrócone (trzy i czterocyfrowe), miejscowe, międzymiastowe, na komórki,
- identyfikatory (w Twoim ulubionym języku programowania)
- napisy (w Twoim ulubionym języku programowania) wraz z możliwością umieszczenia ogranicznika napisu w napisie,
- komentarze (w Twoim ulubionym języku programowania).

3. Napisz specyfikację skanera, który:

- Usuwa białe znaki z początku/końca każdej linii. (Przeczytaj w dokumentacji Lex'a o znaczeniu \wedge i $\$$ we wzorcach.)
- Zamienia ciąg białych znaków na jeden odstęp.
- Wstawia odstęp po każdej kropce i przecinku, chyba że taki odstęp już tam jest.
- Usuwa komentarze typu `//` z programu w C/C++. (Ignorujemy problemy z komentarzami postaci `/*...//...*/`):
 - nie uwzględniając stałych napisowych,
 - uwzględniając je.

Wykład 5. Deterministyczne automaty skończone

5.1 Wstęp

Przypomnijmy, że na język można patrzeć jak na problem decyzyjny. W ramach tego wykładu zajmiemy się modelowaniem takich mechanizmów liczących, które są przystosowane do rozstrzygania, czy dane słowo należy do języka i są wyposażone w pomocniczą pamięć stałej wielkości. Takie modele to właśnie automaty skończone. Języki, dla których takie mechanizmy istnieją, jak się później okaże, to języki regularne. Zaczniemy jednak od ustalenia pewnych, być może oczywistych i intuicyjnych pojęć.

5.2 Intuicje

Ustalmy najpierw pewne intuicyjne pojęcia: *stanu* i *przejścia*. *Stan* jakiegokolwiek systemu (czy mechanizmu), to pełna informacja na jego temat w danym momencie czasu. W szczególności stan zawiera wszelką „pamięć” w jaką wyposażony jest dany system. Również wszystko to co zdarzyło się do danego momentu ma wpływ na przyszłość tylko o tyle, o ile wpłynęło na stan systemu w danej chwili.

Będziemy rozpatrywać takie systemy i mechanizmy, których stan zmienia się w sposób dyskretny (tj. skokowy), a nie ciągły. Podejście takie ma zastosowanie np. do komputerów i urządzeń cyfrowych. Praca podzespołów komputerowych jest taktowana odpowiednimi zegarami. Ich stan może się zmieniać tylko wraz z kolejnymi taktami odpowiednich zegarów. Na cały komputer możemy spojrzeć jak na system, w którym czas biegnie nie w sposób ciągły, ale drobnymi skokami. W przypadku wielu innych systemów takie podejście też może być poprawne, jeżeli ich stan, który będzie nas interesował, będzie miał charakter dyskretny. Wówczas jego zmiany będą musiały mieć charakter skokowy. Dlatego też w naszych modelach czas będzie *dyskretny*, tzn. będzie on płynął drobnymi kroczkami.

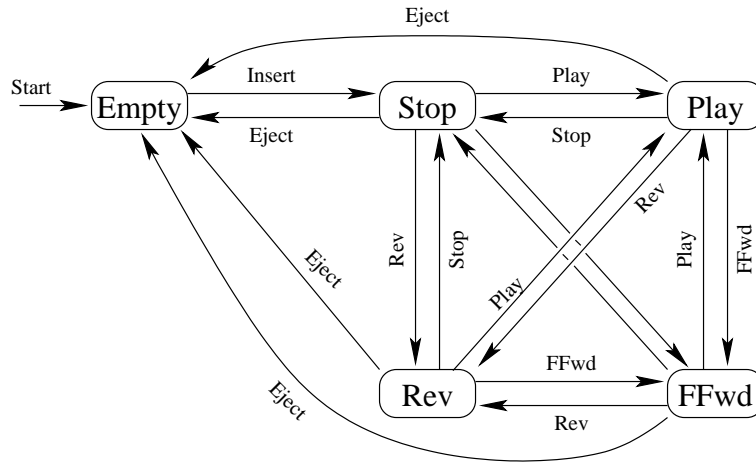
Zmianę stanu będziemy nazywać *przejściem* — od stanu w chwili poprzedniej do stanu w chwili następnej. Zwykle przejście ma miejsce jako reakcja na pewne zdarzenie. Założymy przy tym, że przejścia są natychmiastowe — jest to oczywiście uproszczenie, tak samo jak założenie, że czas jest dyskretny.

Najlepiej powyższe intuicje zilustrować na przykładach:

- Kostka Rubik’a — jej stany to możliwe ułożenia, a przejścia zachodzą na skutek wykonywania pojedynczych ruchów.
- Zegarek cyfrowy — Stan zegarka to aktualny czas (z dokładnością do jego pomiaru przez zegarek) oraz cała pamięć zegarka (np. godzina, na którą ustawiony jest budzik). Przejścia zachodzą na skutek upływu czasu lub naciskania przycisków zegarka.
- Sterownik windy — stany tego sterownika zawierają informację m.inn. nt. położenia windy (z pewną dokładnością), kierunku jej jazdy, otwarcia lub zamknięcia drzwi, oraz tego, jakie guziki (w windzie i na zewnątrz) były naciśnięte.

- Wszechświat(?) — pytanie czy wszechświat jest systemem dyskretnym pozostawimy filozofom i fizykom.

Przykład: Na poniższym diagramie przedstawiono uproszczony schemat działania bardzo prostego odtwarzacza video. Stany zaznaczono owalami. W momencie włączenia system znajduje się w stanie „Empty”. Przejścia przedstawiono jako strzałki i umieszczono przy nich etykiety reprezentujące zdarzenia powodujące określone przejścia.



5.3 Determisticzne automaty skończone

Powróćmy do języków. Interesuje nas modelowanie mechanizmów obliczeniowych odpowiadających na pytanie, czy dane słowo należy do ustalonego języka. Automat skończony wczytuje dane słowo znak po znaku i po wczytaniu całego słowa udziela odpowiedzi, czy słowo to należy do ustalonego języka. Przejścia między stanami zachodzą w automacie skończonym na skutek wczytywania kolejnych znaków analizowanego słowa. W automacie deterministycznym, dla każdego znaku jaki może się pojawić na wejściu i dla każdego stanu automatu, ze stanu tego wychodzi dokładnie jedno przejście odpowiadające danemu znakowi — w każdej sytuacji działanie automatu musi być określone jednoznacznie. (W tym wykładzie będziemy się zajmować wyłącznie automatami deterministycznymi. Automaty niedeterministyczne pojawią się w kolejnym wykładzie i dopiero wówczas będziemy je odróżniać.)

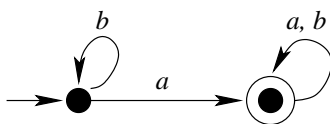
To, jaka jest odpowiedź udzielana przez automat zależy od stanu osiąganego po wczytaniu całego słowa. Jeżeli odpowiedź jest pozytywna, to mówimy, że automat *akceptuje* słowo, a w przeciwnym przypadku mówimy, że je *odrzuca*. Stany, w których automat akceptuje nazywamy *stanami akceptującymi*. Pamiętajmy jednak, że o tym, czy automat akceptuje słowo decyduje wyłącznie ostatni stan — ten, w którym automat jest po wczytaniu całego słowa.

Dodatkowo, automat skończony, jak sama nazwa wskazuje, ma skończoną liczbę stanów. Oznacza to, że dysponuje on stałą pamięcią dodatkową, tzn. wielkość pamięci dodatkowej

nie zależy od długości wczytywanego słowa i jest ograniczona przez stałą. Jeżeli do stwierdzenia, czy słowo należy do ustalonego języka potrzebna jest pamięć pomocnicza, której nie da się ograniczyć przez stałą, to taki język nie jest akceptowany przez żaden automat skończony. Przykłady takich języków, oraz techniki dowodzenia tego poznamy w dalszych wykładach.

Przykład: Zanim podamy formalną definicję automatu skończonego zobaczymy prosty przykład takiego automatu.

Do prezentowania automatów skończonych będziemy używać diagramów podobnych do diagramu odtwarzacza wideo z poprzedniego punktu. Stany automatu oznaczamy na diagramie punktami, a przejścia strzałkami. Przejścia etykietujemy znakami, których wczytanie powoduje dane przejście. Początkowy stan automatu oznaczamy strzałką „znikąd”. Stany akceptujące zaznaczamy otaczając je dodatkowym kółkiem.



Powyższy automat wczytuje słowa nad alfabetem $\Sigma = \{a, b\}$. Akceptuje on wszystkie słowa zawierające choć jedną literę a — tylko takie słowa mogą spowodować przejście do stanu po prawej stronie.

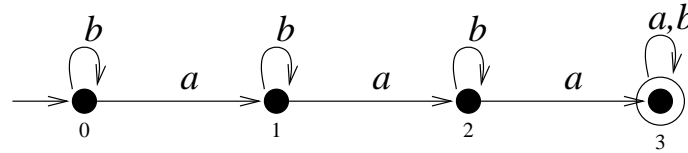
Def. 24. Automat skończony (deterministyczny), to dowolna taka piątka $M = \langle Q, \Sigma, \delta, s, F \rangle$, w której:

- Q to skończony zbiór *stanów*,
- Σ to (skończony) *alfabet wejściowy*,
- $\delta : Q \times \Sigma \rightarrow Q$ to *funkcja przejścia* — $\delta(q, a)$ to stan do którego przechodzimy ze stanu q na skutek wczytania znaku a ,
- s to stan początkowy,
- $F \subseteq Q$ to zbiór stanów akceptujących.

□

Przedstawiając automat musimy podać informacje o wszystkich pięciu elementach. Formalne podawanie piątki, takiej jak w definicji, jest mało czytelne. Dlatego też zwykle przedstawia się automat w postaci tabelki lub diagramu.

Przykład: Następująca diagram i tabelka przedstawiają ten sam automat.



	a	b
$\rightarrow 0$	1	0
1	2	1
2	3	2
F 3	3	3

Tabelka przedstawia przede wszystkim funkcję przejścia δ . Wiersze są indeksowane stanami. Kolumny są indeksowane znakami z alfabetu wejściowego. Komórki tabeli zawierają stany, do których przechodzimy na skutek wczytania pojedynczych znaków. Dodatkowo, stan początkowy jest zaznaczony za pomocą strzałki \rightarrow , a stany akceptujące są zaznaczone literami F .

Zaletą diagramów jest ich przejrzystość oraz to, że nie ma konieczności nazywania stanów. Natomiast zaletą tabel jest ich zwarta postać.

Dane dla automatu to dowolne słowo $x \in \Sigma^*$. Działanie automatu możemy zasymulować w następujący sposób:

1. Na początku automat jest w stanie początkowym.
2. Wczytujemy kolejne znaki słowa na wejściu i zmieniamy stan zgodnie z przejściami odpowiadającymi tym znakom.
3. Po wczytaniu całego słowa i wykonaniu wszystkich przejść sprawdzamy, czy automat jest w stanie akceptującym. Jeżeli tak, to słowo jest akceptowane przez automat, wpp. nie.

Mając automat przedstawiony w formie diagramu możemy potraktować ten diagram jak planszę. Na początku stawiamy pionek na polu / w stanie początkowym. Następnie przesuujemy pionek zgodnie z wczytywanymi znakami. Jeśli na koniec jesteśmy na polu / w stanie akceptującym, to słowo jest akceptowane.

Przykład:

Czas formalnie zdefiniować działanie deterministycznych automatów skończonych. Niech $M = \langle Q, \Sigma, \delta, s, F \rangle$ będzie ustalonym deterministycznym automatem skończonym.

Def. 25. Funkcję przejścia $\delta : Q \times \Sigma \rightarrow Q$ rozszerzamy do funkcji $\delta^* : Q \times \Sigma^* \rightarrow Q$ w następujący sposób:

$$\delta^*(q, \varepsilon) = q$$

$$\delta^*(q, xa) = \delta(\delta^*(q, x), a)$$

□

Podczas gdy funkcja δ określa działanie automatu dla pojedynczego znaku, funkcja δ^* określa jego działanie dla całych słów.

Def. 26. Automat akceptuje słowo x wtw., gdy $\delta^*(s, x) \in F$.

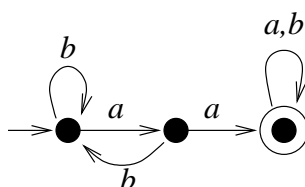
Język akceptowany przez automat M , to:

$$L(M) \equiv \{x \in \Sigma^* : \delta^*(s, x) \in F\}$$

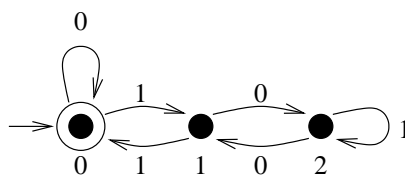
□

Inaczej mówiąc, język akceptowany przez automat składa się dokładnie z tych słów, które ten automat akceptuje.

Przykład: Automat akceptujący słowa nad alfabetem $\{a, b\}$ zawierające pod słowo aa .



Przykład: Automat akceptujący $\{x \in \{0,1\}^* : x \text{ w układzie binarnym jest podzielne przez } 3\}$.



Stany reprezentują reszty modulo 3. Po wczytaniu dowolnego ciągu cyfr automat jest w stanie odpowiadającym reszcie z dzielenia wczytanej liczby przez 3. Stan początkowy (odpowiadający wczytaniu ε) to 0. Budując przejścia automatu musimy się zastanowić, jak zmienia się reszta z dzielenia danej liczby przez 3, gdy do jej zapisu binarnego dopiszemy (z prawej strony) 0 lub 1. Korzystamy tu z następujących tożsamości:

$$x0 = 2 \cdot x \equiv (2(x \bmod 3)) \bmod 3$$

$$x1 = 2 \cdot x + 1 \equiv (2(x \bmod 3) + 1) \bmod 3$$

Czyli $\delta(q, x) = (2 \cdot q + x) \bmod 3$.

5.4 Własności klasy języków akceptowanych przez deterministyczne automaty skończone

Zajmiemy się teraz własnościami klasy języków akceptowanych przez deterministyczne automaty skończone. Jak się później okaże, jest to dokładnie klasa języków regularnych. Zanim się to jednak stanie, pokażemy, że automaty skończone mają niektóre z cech wzorców, które pokazaliśmy w poprzednim wykładzie.

Fakt 3. *Każdy język skończony jest akceptowany przez pewien deterministyczny automat skończony.*

Dowód: Niech A będzie skończonym językiem. Wszystkie słowa z A razem wzięte mają skończoną liczbę prefiksów. Stanami automatu akceptującego A będą właśnie te prefiksy, plus jeden dodatkowy stan, który będziemy nazywać „śmietnikiem”. Po wczytaniu dowolnego słowa, jeżeli jest ono prefiksem pewnego słowa z A , nasz automat będzie dokładnie w takim stanie jak to słowo. Natomiast po wczytaniu słowa, które nie jest prefiksem żadnego słowa z A automat przejdzie do śmietnika.

Stan początkowy to ε . Jeśli q i qa są prefiksami pewnego słowa z A , to $\delta(q, a) = qa$, wpp. $\delta(q, a) = \text{śmietnik}$. Stany akceptujące to te prefiksy, które należą do A .

W oczywisty sposób, język akceptowany przez automat to A . Natomiast dzięki temu, że słowa z A mają skończoną liczbę prefiksów, jest to poprawny automat skończony. \square

Twierdzenie: Niech A i B będą językami (nad tym samym alfabetem Σ) akceptowanymi odpowiednio przez automaty deterministyczne $M_1 = \langle Q_1, \Sigma, \delta_1, s_1, F_1 \rangle$ i $M_2 = \langle Q_2, \Sigma, \delta_2, s_2, F_2 \rangle$, $L(M_1) = A$, $L(M_2) = B$. Wówczas następujące języki są również akceptowane przez pewne deterministyczne automaty skończone:

- \bar{A} ,
- $A \cap B$,
- $A \cup B$.

Dowód:

- Dopełnienie: Weźmy $M'_1 = \langle Q_1, \Sigma, \delta_1, s_1, Q \setminus F_1 \rangle$. Mamy:

$$\begin{aligned} L(M'_1) &= \{x \in \Sigma^* : \delta_1^*(s_1, x) \in Q \setminus F_1\} = \\ &= \{x \in \Sigma^* : \delta_1^*(s_1, x) \notin F_1\} = \\ &= \Sigma^* \setminus \{x \in \Sigma^* : \delta_1^*(s_1, x) \in F_1\} = \\ &= \Sigma^* \setminus L(M_1) = \Sigma^* \setminus A = \bar{A} \end{aligned}$$

Tak więc \bar{A} jest akceptowany przez M_1 .

- Weźmy $M_3 = \langle Q_1 \times Q_2, \Sigma, \delta_3, (s_1, s_2), F_1 \times F_2 \rangle$, przy czym $\delta_3((q_1, q_2), a) = (\delta_1(q_1, a), \delta_2(q_2, a))$.

Automat taki nazywamy *automatem produktowym*. Jego stany to stany automatów M_1 i M_2 . Symuluje on równocześnie działanie obydwóch tych automatów. Jednocześnie akceptuje tylko takie słowa, które byłyby zaakceptowane równocześnie przez M_1 i M_2 .

Można pokazać (przez indukcję po długości słowa x), że $\delta_3^*((p, q), x) = (\delta_1^*(p, x), \delta_2^*(q, x))$. Stąd uzyskujemy:

$$\begin{aligned} L(M_3) &= \{x \in \Sigma^* : \delta_3^*((s_1, s_2), x) \in F_1 \times F_2\} = \\ &= \{x \in \Sigma^* : (\delta_1^*(s_1, x), \delta_2^*(s_2, x)) \in F_1 \times F_2\} = \\ &= \{x \in \Sigma^* : \delta_1^*(s_1, x) \in F_1 \wedge \delta_2^*(s_2, x) \in F_2\} = \\ &= \{x \in \Sigma^* : \delta_1^*(s_1, x) \in F_1\} \cap \{x \in \Sigma^* : \delta_2^*(s_2, x) \in F_2\} = \\ &= L(M_1) \cap L(M_2) = A \cap B \end{aligned}$$

Tak więc $A \cap B$ jest akceptowany przez M_3 .

- Z praw De Morgana mamy $A \cup B = \overline{\overline{A \cap B}}$. Stąd $A \cup B$ jest akceptowany przez pewien deterministyczny automat skończony. \square

5.5 Podsumowanie

Ten wykład był poświęcony deterministycznym automatom skończonym. Poznaliśmy ich budowę i sposób działania. Poznaliśmy też podstawowe operacje, na które domknięta jest klasa języków akceptowanych przez automaty deterministyczne. Jak się okaże w dalszej części kursu, jest to dokładnie klasa języków regularnych.

5.6 Skorowidz

- **Automat skończony, deterministyczny** to dowolna taka piątka $M = \langle Q, \Sigma, \delta, s, F \rangle$, w której: Q to skończony zbiór *stanów*, Σ to (skończony) *alfabet wejściowy*, $\delta : Q \times \Sigma \rightarrow Q$ to *funkcja przejścia*, s to stan początkowy, $F \subseteq Q$ to zbiór stanów akceptujących.
- **Automat produktowy** to automat symulujący równocześnie działanie dwóch danych automatów deterministycznych — jego zbiór stanów to produkt kartezjański zbiorów stanów danych automatów.
- **Język akceptowany** przez deterministyczny automat skończony $M = \langle Q, \Sigma, \delta, s, F \rangle$ to język postaci:

$$L(M) \equiv \{x \in \Sigma^* : \delta^*(s, x) \in F\}$$

- **Przejście** to zmiana stanu — od stanu w chwili poprzedniej do stanu w chwili następczej. Zwykle przejście ma miejsce jako reakcja na pewne zdarzenie.

- **Stan** jakiegokolwiek systemu (czy mechanizmu), to pełna informacja na jego temat w danym momencie czasu. W szczególności stan zawiera wszelką „pamięć” w jaką wyposażony jest dany system.

5.7 Praca domowa

1. Podaj deterministyczny automat skończony akceptujący te słowa nad alfabetem $\{0, 1\}$, w których na parzystych pozycjach występują 1-ki.
2. Podaj deterministyczny automat skończony akceptujący te słowa nad alfabetem $\{a, b\}$, które zawierają pod słowo $abba$.
3. Podaj deterministyczny automat skończony akceptujący przecięcie języków akceptowanych przez automaty:

	a	b
\rightarrow 1	2	2
F 2	1	1

	a	b
\rightarrow 1	1	2
F 2	1	2

5.8 Ćwiczenia

Podaj deterministyczne automaty skończone akceptujące następujące języki:

1. słowa nad alfabetem $\{0, 1\}$, które zawierają parzystą liczbę zer,
2. słowa nad alfabetem $\{a, b\}$ zakończone na aa ,
3. słowa nad alfabetem $\{1, 2, 3, 4\}$ zawierające pod słowo 42,
4. słowa nad alfabetem $\{a\}$ o długości podzielnej przez 2 lub 3,
5. słowa nad alfabetem $\{a, b\}$ zawierające pod słowo abb ,
6. słowa nad alfabetem $\{a, b\}$ nie zawierające pod słowa aba ,
7. słowa nad alfabetem $\{a, b\}$ zawierające przynajmniej 2 wystąpienia słowa aba (być może zachodzące na siebie),
8. te słowa nad alfabetem $\{a, b\}$, w których pierwszy i ostatni znak są takie same (uwaga, słowa a i b powinny zostać zaakceptowane),
9. te słowa nad alfabetem $\{0, 1, \dots, 9\}$, które zawierają jako pod słowo wybrany przez Ciebie numer telefonu,
10. zapisy 10-tne liczb podzielnych przez 3,
11. te słowa nad alfabetem $\{0, 1\}$, w których każde dwa kolejne znaki, licząc od początku słowa, zawierają choć jedną 1-kę,

12. przecięcie/sumę języków akceptowanych przez automaty:

			<i>a</i>	<i>b</i>
→ <i>F</i>	1		2	2
	2		1	1

			<i>a</i>	<i>b</i>
→	1		1	2
<i>F</i>	2		2	1

13. przecięcie/sumę języków akceptowanych przez automaty:

			<i>a</i>	<i>b</i>
→	1		1	2
<i>F</i>	2		2	1

			<i>a</i>	<i>b</i>
→ <i>F</i>	1		2	1
	2		1	2

14. przecięcie/sumę języków akceptowanych przez automaty:

			<i>a</i>	<i>b</i>
→	1		2	2
<i>F</i>	2		1	1

			<i>a</i>	<i>b</i>
→ <i>F</i>	1		2	3
	2		3	1
	3		1	2

Wykład 6. Niedeterministyczne automaty skończone

6.1 Wstęp

W poprzednim wykładzie poznaliśmy deterministyczne automaty skończone. Czym jednak jest ów *determinizm*? Determinizm to koncepcja filozoficzna, według której w Świecie nie ma miejsca na przypadek — gdyby mieć pełną informację o stanie świata, można by przewidzieć jego przyszłość. W determinizmie brak też miejsca na wolną wolę — nasze decyzje są określone przez stan Świata, którego jesteśmy częścią.

Niedeterminizm stanowi przeciwną koncepcję. Nawet gdybyśmy mieli pełną wiedzę na temat stanu Świata w danej chwili, to i tak nie bylibyśmy w stanie przewidzieć przyszłości. Co najwyżej moglibyśmy określić jakie są możliwe przyszłe dzieje.

W informatyce determinizm oznacza, że jeżeli będziemy wielokrotnie uruchamiać ten sam program czy mechanizm obliczeniowy, podając te same dane, to zawsze uzyskamy te same wyniki. Niedeterminizm zaś przeciwnie — dla tych samych danych możemy uzyskać różne wyniki. Prosty programy komputerowe (bez randomizacji) zwykle działają w sposób deterministyczny. Niedeterminizm pojawia się np. w algorytmach randomizacyjnych czy programowaniu współbieżnym, czyli tam, gdzie jest jakieś źródło przypadkowości. Wówczas zwykle wymagamy, aby bez względu na przebieg obliczeń była zachowana jakaś wymagana przez nas własność. Niedeterminizm jest też używany przy specyfikowaniu systemów — specyfikacja nie określa jednoznacznie jak powinien zachowywać się system, każde zachowanie zgodne ze specyfikacją jest dopuszczalne.

Automat niedeterministyczny jest więc automatem, którego działania nie da się w pełni przewidzieć. Będąc w tym samym stanie i wczytując ten sam znak może wykonać różne przejścia. Przestaje więc obowiązywać wymóg, że dla danego znaku, z każdego stanu wychodzi dokładnie jedno przejście odpowiadające wczytaniu tego znaku. Takich przejść może być dowolnie wiele, a nawet może ich nie być wcale (co to znaczy wyjaśni się za chwilę). Automat taki może również mieć wiele stanów początkowych — każde uruchomienie automatu może rozpocząć się w dowolnym z tych stanów.

Może więc się zdarzyć, że ten sam automat raz zaakceptuje dane słowo, a raz nie. Czy więc takie słowo należy do języka akceptowanego przez automat, czy nie? Otóż, jeśli tylko automat **może** zaakceptować słowo, to powiemy, że należy ono do języka akceptowanego przez automat. Jest to dosyć nietypowa interpretacja niedeterminizmu. Nie odpowiada ona intuicyjnemu rozumieniu przypadku.

Na taki niedeterministyczny automat możemy spojrzeć jak na mechanizm o dwoistej budowie: jedna jego część to automat skończony określający jakie przejścia są możliwe, druga część, to tzw. *wyrocznia*, która w przypadku wielu możliwych do wykonania przejść wybiera przejście, które prowadzi do zaakceptowania słowa. (Zapachniało metafizyką? Chwila cierpliwości, zaraz wszystko będzie ściśle zdefiniowane.) Jedno z zagadnień, którymi się zajmiemy, to czy taka wyrocznia zmienia się obliczeniową automatów, czyli czy dla tych samych języków potrafimy zbudować automaty deterministyczne i nie, czy też są jakieś różnice? Jak się okaże, w przypadku automatów skończonych nie ma żadnej różnicy, a wyrocznia i niedeterminizm można zastąpić większą liczbą stanów.

6.2 Działanie niedeterministycznych automatów skończonych

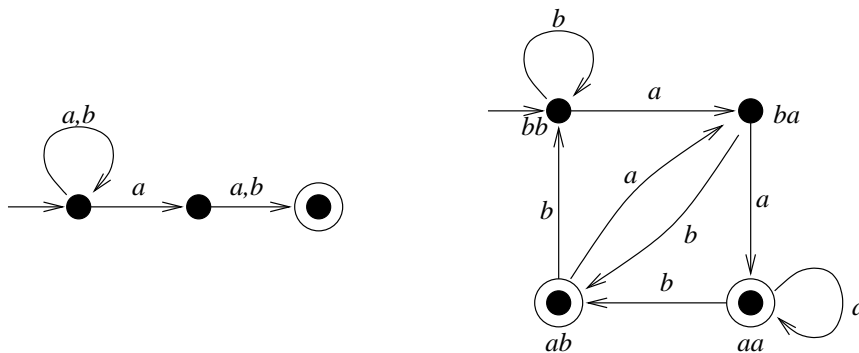
Jak już powiedzieliśmy, w automacie niedeterministycznym, dla danego znaku, nie musi z każdego stanu wychodzi dokładnie jedno przejście odpowiadające wczytaniu tego znaku. Takich przejść może nie być, może być jedno, lub może być wiele. Podobnie, automat niedeterministyczny może mieć wiele stanów początkowych.

Działanie automatu niedeterministycznego możemy symulować w następujący sposób. Spójrzmy na diagram automatu jak na planszę, po której będziemy przesuwać pionki — stany to pola, na których będziemy stawiać pionki, a przesuwać będziemy je zgodnie z przejściami.

1. Na początku symulacji ustawiamy pionki we wszystkich stanach początkowych.
2. Wczytując pojedynczy znak przesuujemy zgodnie z przejściami wszystkie pionki na planszy. Przy czym, jeżeli z pola, na którym stoi pionek wychodzi wiele przejść odpowiadających wczytanemu znakowi, to stawiamy pionki na wszystkich polach docelowych. W szczególności, jeżeli z pola, na którym stoi pionek nie wychodzi żadne takie przejście, to taki pionek zdejmujemy z planszy. Jeżeli na jednym polu znajdzie się kilka pionków, to zostawiamy jeden, a resztę zdejmujemy.
3. Automat akceptuje słowo jeśli po wczytaniu całego słowa w którymkolwiek z pól/stanów akceptujących znajduje się pionek.

Pionki reprezentują wszystkie te stany, w których automat może się znaleźć po wczytaniu odpowiedniego (prefiksu) słowa. Automat akceptuje słowo jeśli choć jeden pionek dotrze do stanu akceptującego, czyli może się zdarzyć takie obliczenie, które prowadzi do jego zaakceptowania.

Przykład: Oto automat niedeterministyczny, który akceptuje te słowa nad alfabetem $\{a, b\}$, w których przedostatni znak to a . Obok przedstawiono dla porównania deterministyczny automat skończony akceptujący taki sam język. Automat deterministyczny musi pamiętać ostatnie dwa znaki wczytywanego słowa — dlatego też ma cztery stany.



Def. 27. Automat niedeterministyczny, to dowolna taka piątka $M = \langle Q, \Sigma, \delta, S, F \rangle$, gdzie:

- Q to skończony zbiór stanów,
- Σ to (skończony) alfabet wejściowy,
- δ to funkcja⁹ przejścia $\delta : Q \times \Sigma \rightarrow 2^Q$,
- $S \subseteq Q$ to zbiór stanów początkowych,
- $F \subseteq Q$ to zbiór stanów akceptujących.

□

Def. 28. Rozszerzenie funkcji przejścia $\delta^* : 2^Q \times \Sigma^* \rightarrow 2^Q$ definiujemy następująco:

$$\delta^*(A, \varepsilon) = A$$

$$\delta^*(A, ax) = \delta^*\left(\bigcup_{q \in A} \delta(q, a), x\right)$$

□

Funkcja δ^* zamiast na pojedynczych stanach operuje na zbiorach stanów. Dokładnie, opisuje ona jak zmienia się ustawienie pionków pod wpływem wczytywanych słów.

Def. 29. Niedeterministyczny automat skończony $M = \langle Q, \Sigma, \delta, S, F \rangle$ akceptuje słowo x wtw., gdy $\delta^*(S, x) \cap F \neq \emptyset$. Język akceptowany przez automat M , to:

$$L(M) = \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\}$$

□

Formalnie automaty niedeterministyczne są czymś innym niż deterministyczne. Jednak my będziemy traktować automaty deterministyczne jak szczególny przypadek automatów niedeterministycznych. Patrząc na diagramy reprezentujące automaty, automaty deterministyczne muszą spełnić więcej wymagań niż automaty niedeterministyczne. Następujący fakt utwierdza nas w takim spojrzeniu:

Fakt 4. Niech $M = \langle Q, \Sigma, \delta, s, F \rangle$ będzie automatem deterministycznym. Wówczas automat niedeterministyczny $M' = \langle Q, \Sigma, \delta', \{s\}, F \rangle$, gdzie $\delta'(q, a) = \{\delta(q, a)\}$ akceptuje dokładnie ten sam język, $L(M) = L(M')$.

Kilka własności funkcji δ^* , które przydadzą nam się w dalszej części wykładu:

⁹W literaturze można spotkać definicje, w których niedeterministyczny automat skończony ma *relacje* przejścia postaci $\delta \subseteq Q \times \Sigma \times Q$. Relacje takie są jednak izomorficzne z funkcjami podanej przez nas postaci.

Fakt 5. Dla dowolnych $x, y \in \Sigma^*$ mamy $\delta^*(A, xy) = \delta^*(\delta^*(A, x), y)$.

Dowód przebiega przez indukcję ze względu na $|y|$.

Fakt 6. Funkcja δ^* jest rozłączna ze względu na sumowanie zbiorów, tzn.:

$$\delta^*\left(\bigcup_i A_i, x\right) = \bigcup_i \delta^*(A_i, x)$$

Dowód przebiega przez indukcję ze względu na $|x|$.

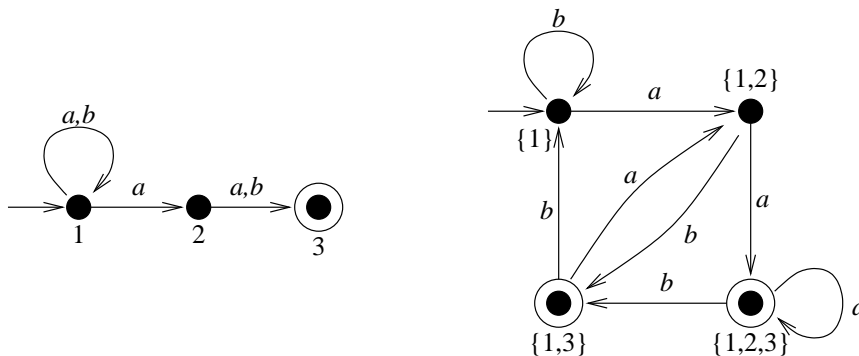
6.3 Determinizacja automatów skończonych

Pokażemy teraz, że siła wyrazu automatów niedeterministycznych jest dokładnie taka sama, jak deterministycznych. W tym celu musimy pokazać, że dla każdego niedeterministycznego automatu skończonego istnieje równoważny mu automat deterministyczny. To, że automaty deterministyczne możemy traktować jak szczególny przypadek automatów niedeterministycznych powiedzieliśmy już wcześniej.

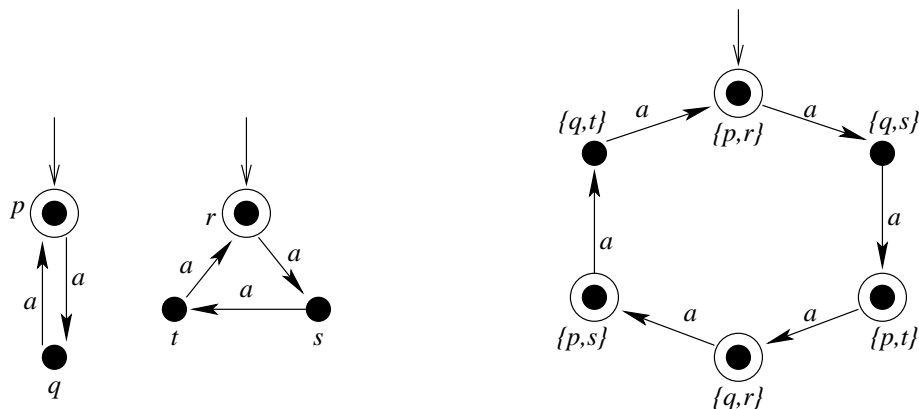
W trakcie symulacji działania automatu niedeterministycznego istotny jest zbiór pól, na których w danym momencie znajdują się pionki. Nasza konstrukcja będzie polegała na zbudowaniu automatu deterministycznego, którego stanami będą *zbiory* stanów automatu niedeterministycznego (mogące się pojawić w czasie symulacji), jego stanem początkowym będzie zbiór stanów początkowych automatu niedeterministycznego, a akceptujące będą te stany, które zawierają choć jeden stan akceptujący automatu niedeterministycznego.

Konstruując automat deterministyczny liczba stanów może wzrosnąć nawet wykładniczo. Jednak nie zawsze tak musi być, gdyż nie zawsze wszystkie możliwe zbiory stanów automatu niedeterministycznego mogą się pojawić w czasie symulacji. Zobaczmy kilka przykładów.

Przykład: Przypomnijmy sobie automat niedeterministyczny akceptujący słowa nad alfabetem $\{a, b\}$, w których przedostatni znak to a . W trakcie dowolnego obliczenia, na polu początkowym nr 1 będzie cały czas stał pionek. Tak więc mogą się pojawić cztery różne konfiguracje pionków (czy też cztery różne zbiory pól zajmowanych przez pionki). Ponieważ konfiguracje nieosiągalne nie mają znaczenia dla działania automatu, możemy się ograniczyć się do tych czterech konfiguracji. Odpowiada to automатовi zapamiętującemu dwa ostatnie znaki w słowie.



Przykład: Automat przedstawiony na poniższym rysunku akceptuje słowa (nad alfabetem $\{a\}$) o długości podzielnej przez 2 lub 3. Jedyny niedeterminizm polega tu na wybraniu stanu początkowego. Automat deterministyczny ma 6 stanów.



Konstruowany automat deterministyczny, równoważny automатовi niedeterministycznemu, jest nazywany *automatem potęgowym* — gdyż jego zbiór stanów to zbiór potęgowy zbioru stanów automatu niedeterministycznego. Formalnie, jego konstrukcja przebiega następująco.

Konstrukcja automatu potęgowego Niech $M = \langle Q, \Sigma, \delta, S, F \rangle$ będzie niedeterministycznym automatem skończonym. Jego determinizacją będziemy nazywać automat deterministyczny $M' = \langle 2^Q, \Sigma, \delta', S, F' \rangle$, gdzie:

$$\delta'(A, a) = \delta^*(A, a)$$

$$F' = \{A \subseteq Q : A \cap F \neq \emptyset\}$$

Pokażemy, że oba automaty akceptują te same języki.

Lemat 1. Dla dowolnych $A \subseteq Q$ i $x \in \Sigma^*$ zachodzi

$$\delta^*(A, x) = \delta'^*(A, x)$$

Dowód: Dowód przebiega przez indukcja ze względu na $|x|$.

1. $\delta^*(A, \varepsilon) = A = \delta'^*(A, \varepsilon)$,
2. Załóżmy, że lemat zachodzi dla słowa x . Pokażemy, że zachodzi również dla słowa postaci ax :

$$\delta^*(A, ax) = \delta^*(\delta^*(A, a), x) = \delta^*(\delta'(A, a), x) = \delta^*(\delta'^*(A, a), x)$$

Z założenia indukcyjnego:

$$\delta^*(\delta'^*(A, a), x) = \delta'^*(\delta'^*(A, a), x) = \delta'^*(A, ax)$$

□

Korzystając z lematu, dowód, że jest natychmiastowy:

$$\begin{aligned} L(M) &= \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\} = \\ &= \{x \in \Sigma^* : \delta'^*(S, x) \cap F \neq \emptyset\} = \\ &= \{x \in \Sigma^* : \delta'^*(S, x) \in F'\} = L(M') \end{aligned}$$

W tak skonstruowanym automacie potęgowym możemy, bez zmiany akceptowanego przez niego języka, usunąć wszystkie stany nieosiągalne.

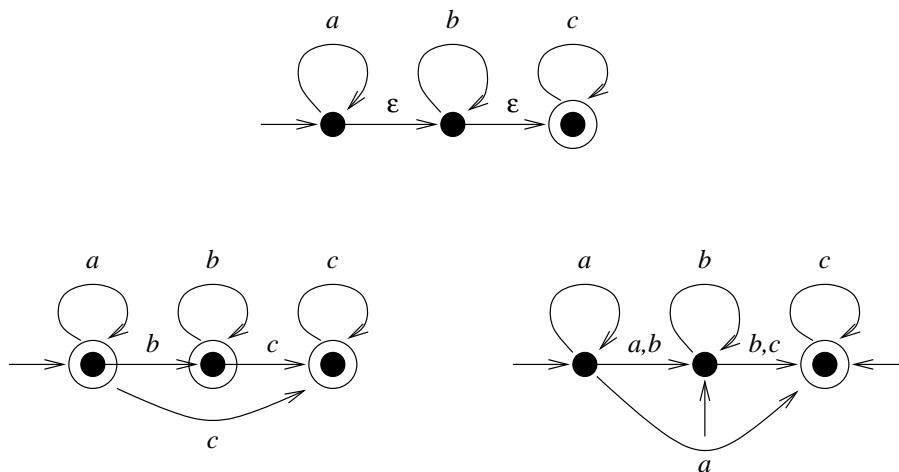
6.4 Automaty z ε -przejdściami

Automaty z ε -przejdściami stanowią rozszerzenie automatów niedeterministycznych. Mogą one dodatkowo zawierać tzw. ε -przejdścia. Są to przejdścia, które automat może w każdej chwili wykonywać, nie wczytując nic z wejścia. Tak więc automat z ε -przejdściami nie musi w każdym kroku wczytywać jednego znaku.

Podobnie, jak automaty deterministyczne stanowią szczególny przypadek automatów niedeterministycznych, tak automaty niedeterministyczne stanowią szczególny przypadek automatów z ε -przejdściami — są to automaty z ε -przejdściami, w których nie ma ε -przejdść.

Okazuje się, że automaty z ε -przejdściami mają taką samą siłę wyrazu co automaty niedeterministyczne. Czasami jednak ε -przejdścia pozwalają uprościć budowę automatu.

Przykład: Poniższy rysunek przedstawia automat z ε -przejdściami akceptujący język $a^*b^*c^*$. Poniżej przedstawiono dwa równoważne mu automaty niedeterministyczne. Jak widać, ε -przejdścia można usunąć zachowując ten sam zbiór stanów, ale kosztem zwiększenia liczby przejdść oraz większej liczby stanów początkowych lub akceptujących.



Def. 30. Automat z ε -przejściami, to dowolna taka piątka $M = \langle Q, \Sigma, \delta, S, F \rangle$, w której:

- Q jest skończonym zbiorem stanów,
- Σ to (skończony) alfabet,
- $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ to funkcja przejścia,
- $S \subseteq Q$ to zbiór stanów początkowych, oraz
- $F \subseteq Q$ to zbiór stanów akceptujących.

□

Jak widać, jedyna różnica między definicją automatu niedeterministycznego i ε -przejściami polega na dodatkowej etykiecie jaka może towarzyszyć przejściom: ε . Zobaczmy jak definiujemy wpływ ε -przejść na działanie automatów:

Def. 31. Niech A będzie zbiorem stanów, $A \subseteq Q$. Przez $D(A)$ oznaczamy zbiór stanów, które można osiągnąć z A poprzez ε -przejścia.

Rozszerzenie funkcji przejścia $\delta^* : 2^Q \times \Sigma^* \rightarrow 2^Q$ definiujemy następująco:

$$\delta^*(A, \varepsilon) = D(A)$$

$$\delta^*(A, ax) = \delta^*\left(\bigcup_{q \in D(A)} \delta(q, a), x\right)$$

□

Każde obliczenie automatu z ε -przejściami to przeplot ε -przejść i zwykłych przejść. Odpowiada temu w definicji δ^* przeplatanie funkcji D i δ .

Def. 32. Język akceptowany przez automat z ε -przejściami to:

$$L(M) = \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\}$$

□

Powyższa definicja jest identyczna jak w przypadku automatów niedeterministycznych. Pamiętajmy jednak, że funkcja δ^* , opisująca działanie automatu, jest zdefiniowana inaczej.

Pokażemy teraz, że automaty z ε -przejściami i automaty niedeterministyczne mają taką samą siłę wyrazu. (Tym samym, automaty z ε -przejściami mają taką samą siłę wyrazu co automaty deterministyczne.) Powiedzieliśmy już wcześniej, że automaty niedeterministyczne są szczególnym przypadkiem automatów z ε -przejściami. Musimy pokazać, że ε -przejścia nie zwiększają siły wyrazu i że zawsze można je wyeliminować, zamieniając automat z ε -przejściami w automat niedeterministyczny.

Eliminacja ε -przejęć Eliminacja ε -przejęć opiera się na dodaniu dodatkowych przejęć i stanów początkowych, tak aby definicje funkcji δ^* dla automatów niedeterministycznych i z ε -przejęciami były sobie równoważne. W tym celu dodamy dodatkowe stany początkowe oraz dodatkowe przejęcia.

Niech $M = \langle Q, \Sigma, \delta, S, F \rangle$ będzie automatem skończonym z ε -przejęciami. Równoważny mu niedeterministyczny automat skończony M' ma postać $M' = \langle Q, \Sigma, \delta', S', F \rangle$, gdzie $S' = D(S)$, $\delta'(q, a) = D(\delta(q, a))$ (dla $a \neq \varepsilon$).

Dowód poprawności Zbiory stanów i stanów akceptujących są w obu automatach takie same. Pokażemy przez indukcję (ze względu na $|x|$), że dla dowolnego $A \subseteq Q$ zachodzi $\delta^*(A, x) = \delta'^*(D(A), x)$.

1. Dla $x = \varepsilon$ mamy:

$$\delta^*(A, \varepsilon) = D(A) = \delta'^*(D(A), \varepsilon)$$

2. Załóżmy, że teza jest prawdziwa dla x . Pokażemy jej prawdziwość dla ax (dla dowolnego $a \in \Sigma$).

$$\begin{aligned} \delta^*(A, ax) &= \delta^*\left(\bigcup_{q \in D(A)} \delta(q, a), x\right) = \\ &= \delta'^*(D\left(\bigcup_{q \in D(A)} \delta(q, a)\right), x) = \\ &= \delta'^*\left(\bigcup_{q \in D(A)} D(\delta(q, a)), x\right) = \\ &= \delta'^*\left(\bigcup_{q \in D(A)} \delta'(q, a), x\right) = \delta'^*(D(A), ax) \end{aligned}$$

Stąd mamy:

$$L(M') = \{x \in \Sigma^* : \delta'^*(S', x) \cap F \neq \emptyset\} = \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\} = L(M)$$

Czyli oba automaty, M i M' , akceptują ten sam język. □

6.5 Podsumowanie

W tym wykładzie przedstawiliśmy automaty niedeterministyczne i automaty z ε -przejęciami. Pokazaliśmy, że mają one taką samą siłę wyrazu. Siła wyrazu automatów deterministycznych, niedeterministycznych i z ε -przejęciami jest dokładnie taka sama.

6.6 Skorowidz

- **Automat potęgowy** to deterministyczny automat skończony równoważny danemu niedeterministycznemu automatowi skończonemu $M = \langle Q, \Sigma, \delta, S, F \rangle$ następującej postaci: $M' = \langle 2^Q, \Sigma, \delta', S, F' \rangle$, gdzie $\delta'(A, a) = \delta^*(A, a)$, $F' = \{A \subseteq Q : A \cap F \neq \emptyset\}$.

- **Automat skończony, niedeterministyczny** to dowolna taka piątka $M = \langle Q, \Sigma, \delta, S, F \rangle$, w której: Q to skończony zbiór stanów, Σ to (skończony) alfabet wejściowy, δ to funkcja przejścia $\delta : Q \times \Sigma \rightarrow 2^Q$, $S \subseteq Q$ to zbiór stanów początkowych, $F \subseteq Q$ to zbiór stanów akceptujących.
- **Automat skończony z ε -przejściami** to dowolna taka piątka $M = \langle Q, \Sigma, \delta, S, F \rangle$, w której: Q jest skończonym zbiorem stanów, Σ to (skończony) alfabet, $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ to funkcja przejścia, $S \subseteq Q$ to zbiór stanów początkowych, oraz $F \subseteq Q$ to zbiór stanów akceptujących.
- **Determinizacja** niedeterministycznego automatu skończonego polega na skonstruowaniu równoważnego mu automatu potęgowego.
- **Eliminacja ε -przejęć** to konstrukcja niedeterministycznego automatu skończonego równoważnego danemu automatowi skończonemu z ε -przejściami.
- **Język akceptowany** przez niedeterministyczny automat skończony $M = \langle Q, \Sigma, \delta, S, F \rangle$ to:

$$L(M) = \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\}$$

- **Język akceptowany** przez automat skończony z ε -przejściami to:

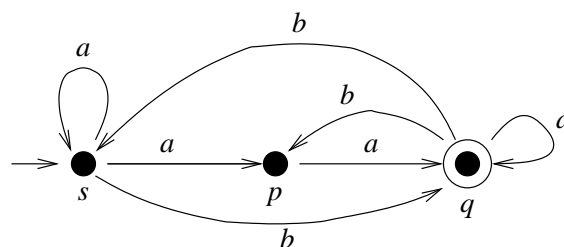
$$L(M) = \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\}$$

- **Niedeterminizm** stanowi koncepcję przeciwną do determinizmu: Nawet gdybyśmy mieli pełną wiedzę na temat stanu Świata w danej chwili, to i tak nie byłibyśmy w stanie przewidzieć przyszłości. Co najwyżej moglibyśmy określić jakie są możliwe przyszłe dzieje.

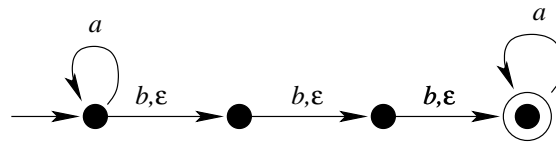
W informatyce niedeterminizm oznacza, że uruchamiając kilka razy ten sam program czy inny mechanizm obliczeniowy, dla tych samych danych możemy uzyskać różne wyniki.

6.7 Praca domowa

1. Podaj automat niedeterministyczny akceptujący słowa nad alfabetem $\{a, b\}$ zawierające podslowo $ababb$.
2. Zdeterminizuj następujący automat niedeterministyczny:

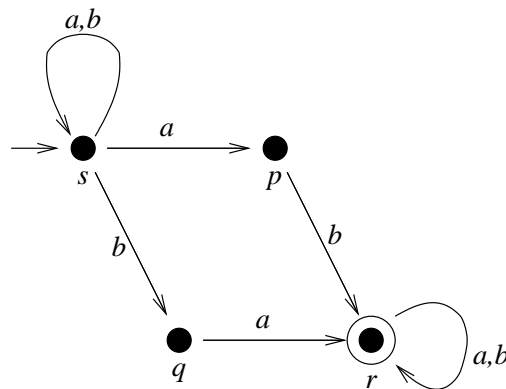
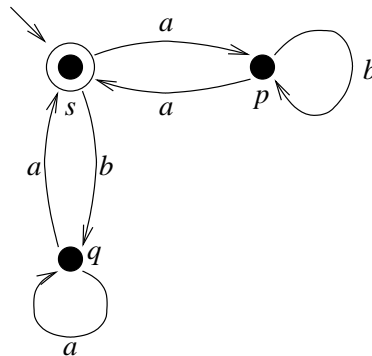
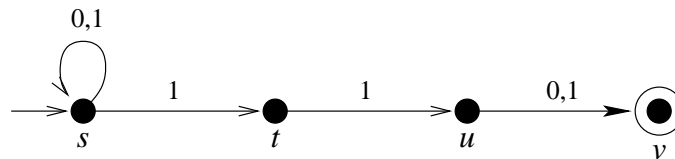


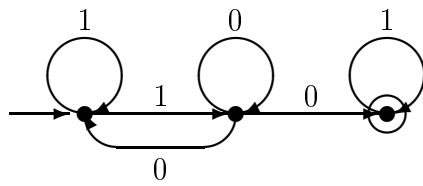
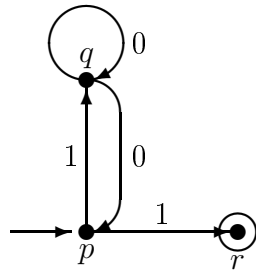
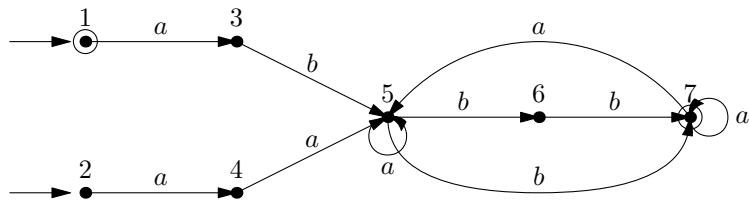
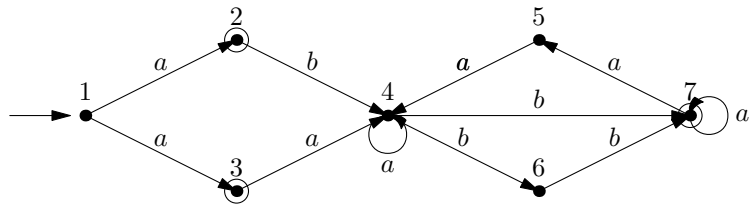
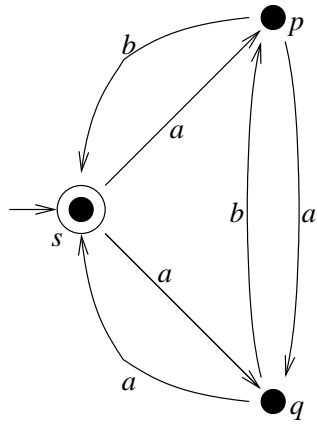
3. Wyeliminuj ε -przejścia w poniższym automacie. Jaki język on akceptuje?



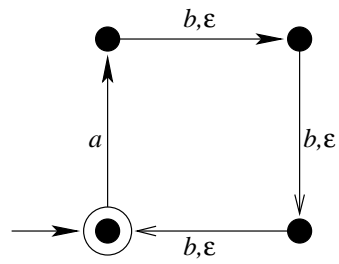
6.8 Ćwiczenia

1. Podaj automat niedeterministyczny akceptujący słowa nad alfabetem $\{a, b\}$ zawierające pod słowo *bababba*. Porównaj go z automatem deterministycznym stworzonym ręcznie i w wyniku determinizacji automatu niedeterministycznego.
2. Zdeterminizuj następujące automaty niedeterministyczne. Pomiń stany nieosiągalne.





3. Wyliminuj ε -przejścia, a następnie zdeterminizuj poniższy automat. Możesz pominąć stany nieosiągalne.



Wykład 7. Równoważność wzorców, wyrażeń regularnych i automatów skończonych

Tematem tego wykładu jest pokazanie, że wzorce, wyrażenia regularne i automaty skończone mają taką samą siłę wyrazu. Wszystkie one opisują dokładnie języki regularne.

Tw. 1. Niech Σ będzie ustalonym alfabetem, a A niech będzie językiem nad tym alfabetem, $A \subseteq \Sigma^*$. Następujące stwierdzenia są równoważne:

- a) A istnieje taki automat skończony M , że $A = L(M)$,
- b) A jest regularny, tzn. $A = L(\alpha)$ dla pewnego wzorca α ,
- c) $A = L(\alpha)$ dla pewnego wyrażenia regularnego α .

Dowód: Pokażemy, że $c) \Rightarrow b)$, $b) \Rightarrow a)$ i $a) \Rightarrow c)$.

$c) \Rightarrow b)$ Oczywiście, bo każde wyrażenie regularne jest wzorcem.

$b) \Rightarrow a)$ Indukcyjnie ze względu na strukturę wzorca budujemy automat z ε -przejściami, o jednym stanie początkowym i jednym stanie akceptującym, który akceptuje dokładnie słowa pasujące do wzorca.

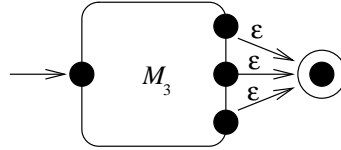
- Dla wzorców postaci ε , \emptyset i a (dla $a \in \Sigma$) odpowiednie automaty przedstawiono na poniższym rysunku.



- Dla wzorca postaci $\bar{\alpha}$ konstrukcja przebiega w kilku krokach:
 1. Niech M_1 będzie automatem skonstruowanym dla wzorca α , $l(M_1) = L(\alpha)$.
 2. Niech $M_2 = \langle Q, \Sigma, \delta, s, F \rangle$ będzie automatem powstałym w wyniku eliminacji ε -przejęć i determinizacji automatu M_1 .
 3. Niech $M_3 = \langle Q, \Sigma, \delta, s, Q \setminus F \rangle$. Analogicznie jak w dowodzie twierdzenia 5.4, M_3 akceptuje dopełnienie języka akceptowanego przez M_2 :

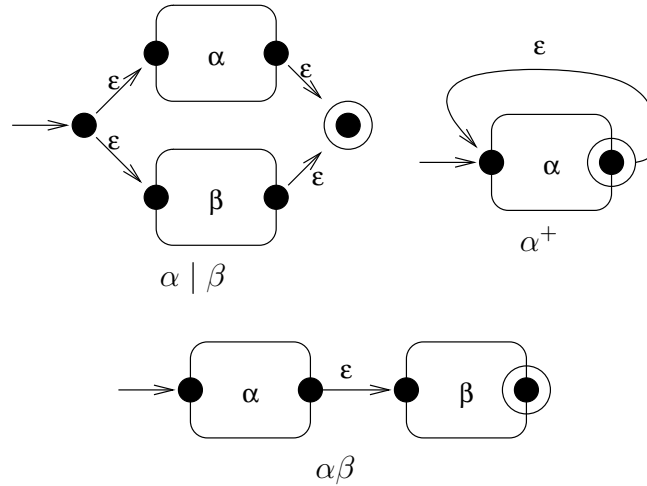
$$L(M_3) = \overline{L(M_2)} = \overline{L(\alpha)} = L(\bar{\alpha})$$

4. Niech M_4 będzie automatem powstałym z M_3 poprzez dodanie: nowego stanu akceptującego, ε -przejęć z dotychczasowych stanów akceptujących do nowego stanu akceptującego i zmianą dotychczasowych stanów akceptujących na nieakceptujące. Konstrukcję tę przedstawia poniższy rysunek:



Konstrukcja ta nie zmienia języka akceptowanego przez automat, $L(M_4) = L(M_3) = L(\bar{\alpha})$, natomiast automat ten spełnia przyjęte założenia o jednym stanie początkowym i jednym stanie akceptującym.

- Jeśli wzorec jest postaci: $\alpha \mid \beta$, α^+ lub $\alpha\beta$, to konstruujemy automat z automatów akceptujących $L(\alpha)$ i $L(\beta)$, w sposób przedstawiony na rysunku. Stany akceptujące zmieniamy tak, że tylko stan zaznaczony na rysunku pozostaje akceptujący:



- Dla wzorców postaci $\alpha = \beta \cap \gamma$ i $\alpha = \beta^*$ wykorzystujemy tożsamości:

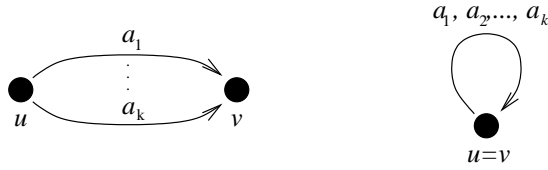
$$\beta \cap \gamma \equiv \overline{\overline{\beta} \mid \overline{\gamma}}$$

$$\beta^* \equiv \varepsilon \mid \beta^+$$

i wcześniej opisane konstrukcje.

a) \Rightarrow c) Niech $M = \langle Q, \Sigma, \delta, S, F \rangle$ będzie danym automatem skończonym (deterministycznym lub nondeterministycznym). Konstrukcja wyrażenia regularnego opisującego język akceptowany przez M przebiega rekurencyjnie. Budujemy wyrażenia regularne opisujące coraz większe fragmenty automatu. Niech u i v będą stanami, a X będzie zbiorem stanów, $u, v \in Q$, $X \subseteq Q$. Wyrażenie regularne $\alpha_{u,v}^X$ opisuje te słowa, za pomocą których można przejść z u do v zatrzymując się po drodze jedynie w stanach z X . Wyrażenie $\alpha_{u,v}^X$ definiujemy indukcyjnie ze względu na $|X|$.

1. Załóżmy, że $X = \emptyset$. Oznaczmy przez a_1, \dots, a_k wszystkie znaki powodujące przejście z u do v . Możliwe są dwa przypadki: $u = v$ lub $u \neq v$:



Jeśli $u \neq v$ i nie ma żadnych przejść z u do v , to wczytanie żadnego słowa nie spowoduje przejścia z u do v . Tak więc:

$$\alpha_{u,v}^{\emptyset} = \begin{cases} a_1 \mid \dots \mid a_k \mid \varepsilon & ; u = v \\ a_1 \mid \dots \mid a_k & ; u \neq v \wedge k > 0 \\ \emptyset & ; u \neq v \wedge k = 0 \end{cases}$$

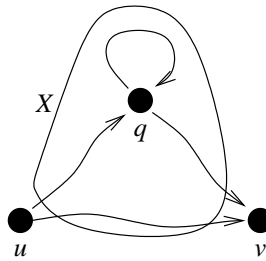
2. Załóżmy, że $X \neq \emptyset$. Niech q będzie dowolnym stanem należącym do X , $q \in X$. Wówczas możliwe są dwa rodzaje przejścia z u do v :

- Przechodzimy z u do v zatrzymując się po drodze w stanach należących do X , ale nie w q . Czyli po drodze zatrzymujemy się w stanach należących do $X \setminus \{q\}$. Słowa, które możemy wczytać w ten sposób są opisane przez $\alpha_{u,v}^{X \setminus \{q\}}$.
- Przechodząc z u do v przynajmniej raz, a być może więcej razy, przechodzimy przez q . Całą trasę możemy podzielić na odcinki, na każdym przejściu przez stan q .

Pierwszy odcinek to przejście z u do q (pierwsza wizyta w q), przechodząc po drodze przez stany z $X \setminus \{q\}$. To co możemy wczytać przechodząc przez ten pierwszy odcinek jest opisane przez $\alpha_{u,q}^{X \setminus \{q\}}$.

Następnie pewną liczbę razy (być może zero) wychodzimy z q i wracamy do q . To co możemy wczytać w trakcie pojedynczego przejścia z q do q jest opisane przez $\alpha_{q,q}^{X \setminus \{q\}}$.

Na koniec, wychodzimy z q (ostatnia wizyta w q) i przechodzimy do v . To co możemy w trakcie tego przejścia wczytać opisuje $\alpha_{q,v}^{X \setminus \{q\}}$.



Łącząc obie powyższe możliwości, wszystko co może zostać wczytane w trakcie przejścia z u do v , po drodze zatrzymując się w stanach należących do X , jest opisane przez:

$$\alpha_{u,v}^X = \alpha_{u,v}^{X \setminus \{q\}} \mid \alpha_{u,q}^{X \setminus \{q\}} (\alpha_{q,q}^{X \setminus \{q\}})^* \alpha_{q,v}^{X \setminus \{q\}}$$

Słowa, które są akceptowane przez M to te, które mogą zostać wczytane przy przejściu z jednego ze stanów początkowych do jednego ze stanów akceptujących. Niech $S = \{s_1, \dots, s_k\}$, $F = \{f_1, \dots, f_l\}$. Wówczas:

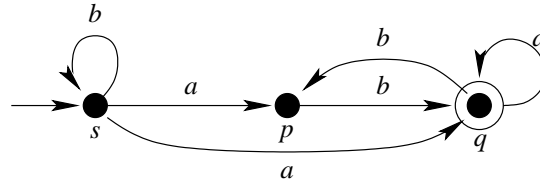
$$L(M) = L(\alpha_{s_1, f_1}^Q \mid \dots \mid \alpha_{s_1, f_l}^Q \mid \dots \mid \alpha_{s_k, f_1}^Q \mid \dots \mid \alpha_{s_k, f_l}^Q)$$

□

Tym samym pokazaliśmy, że: wzorce, wyrażenia regularne, oraz automaty skończone (wszystkie trzy poznane rodzaje), wszystkie opisują dokładnie klasę języków regularnych.

Konstrukcja wyrażenia regularnego odpowiadającego automatowi skończonemu, podana w powyższym dowodzie daje dosyć skomplikowane wyrażenie regularne. Zauważmy, że rozmiar wynikowego wyrażenia regularnego rośnie czterokrotnie z każdym kolejnym stanem automatu! Jednocześnie wiele powstających podwyrażeń można pominąć, gdyż nie wnoszą nic do wyniku. Jeśli więc będziemy chcieli przekształcić jakiś automat na wyrażenie regularne, to nie musimy sztywno trzymać się pokazanej powyżej konstrukcji. Gdy jednak zdrowy rozsądek nie podpowiada nam jak sobie poradzić, wówczas powinniśmy zastosować konstrukcję z dowodu, zmniejszając rozmiar rozważanego automatu. Ilustruje to poniższy przykład:

Przykład: Przekształćmy następujący automat na wyrażenie regularne:



Język akceptowany przez automat jest opisany wyrażeniem.

$$\alpha_{s,q}^{\{s,p,q\}} = \alpha_{s,q}^{\{s,q\}} \mid \alpha_{s,p}^{\{s,q\}} (\alpha_{p,p}^{\{s,q\}})^* \alpha_{p,q}^{\{s,q\}}$$

Patrząc na automat możemy wywnioskować:

$$\alpha_{s,q}^{\{s,q\}} \equiv b^* a a^*$$

$$\alpha_{p,p}^{\{s,q\}} \equiv \varepsilon \mid b a^* b$$

$$\alpha_{p,q}^{\{s,q\}} \equiv b a^*$$

Stąd:

$$\alpha_{s,q}^{\{s,p,q\}} \equiv b^* a a^* \mid \alpha_{s,p}^{\{s,q\}} (\varepsilon \mid b a^* b)^* b a^*$$

Dalej:

$$\alpha_{s,p}^{\{s,q\}} = \alpha_{s,p}^{\{s\}} \mid \alpha_{s,q}^{\{s\}} (\alpha_{q,q}^{\{s\}})^* \alpha_{q,p}^{\{s\}}$$

Zauważmy, że:

$$\alpha_{s,p}^{\{s\}} \equiv \alpha_{s,q}^{\{s\}} \equiv b^* a$$

$$\alpha_{q,q}^{\{s\}} \equiv \varepsilon \mid a$$

$$\alpha_{q,p}^{\{s\}} \equiv b$$

Czyli:

$$\alpha_{s,p}^{\{s,q\}} \equiv b^*a \mid b^*a(\varepsilon \mid a)^*b$$

Stąd:

$$\begin{aligned} \alpha_{s,q}^{\{s,p,q\}} &\equiv b^*aa^* \mid (b^*a \mid b^*a(\varepsilon \mid a)^*b)(\varepsilon \mid ba^*b)^*ba^* \\ &\equiv b^*aa^* \mid (b^*a \mid b^*aa^*b)(ba^*b)^*ba^* \end{aligned}$$

7.1 Zastosowania automatów skończonych

Automaty skończone nie są tylko teoretycznym narzędziem do badania języków regularnych. Mają one również wiele zastosowań praktycznych. Na przykład, stosuje się je do wyszukania wzorca lub pod słowa w tekście — dla danego wzorca czy pod słowa należy skonstruować odpowiedni automat (deterministyczny), a następnie wczytać analizowany tekst; każde przejście takiego automatu przez stan akceptujący reprezentuje jedno wystąpienie wzorca czy pod słowa w tekście.

Dalej, analizatory leksykalne są realizowane jako automaty skończone. Lex na podstawie specyfikacji tworzy odpowiedni automat skończony, który wyszukuje wystąpienia leksemów i równocześnie określa które fragmenty kodu ze specyfikacji należy zastosować.

Automaty skończone mają też zastosowanie w bardzo młodej dziedzinie informatyki, weryfikacji modelowej (ang. *model-checking*). Weryfikacja modelowa jest stosowana do badania, czy określony mechanizm (np. układ elektroniczny, protokół komunikacyjny, prosty program o ograniczonej liczbie stanów) posiada wymagane własności. W tym celu tworzy się automat skończony będący modelem systemu. Język akceptowany przez ten automat reprezentuje możliwe przebiegi zdarzeń. Następnie bada się, czy język akceptowany przez automat zawiera się w języku złożonym z dopuszczalnych przez specyfikację przebiegów zdarzeń. Zaletą weryfikacji modelowej w porównaniu z testowaniem jest to, że weryfikacja modelowa obejmuje **wszystkie** możliwe przebiegi zdarzeń, a nie tylko niektóre.

7.2 Podsumowanie

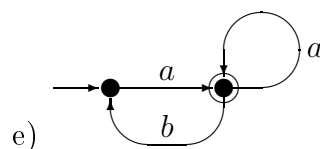
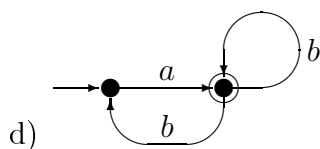
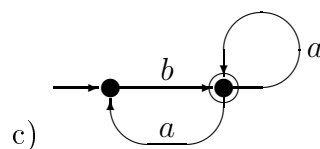
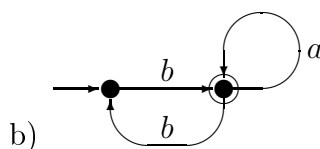
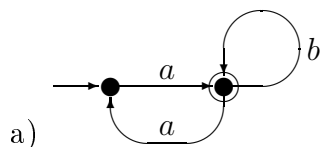
W niniejszym wykładzie pokazaliśmy, że: wzorce, wyrażenia regularne, oraz automaty skończone (wszystkie trzy poznane rodzaje), wszystkie opisują dokładnie klasę języków regularnych. Ponieważ dowód był konstruktywny, widać z niego jak można tłumaczyć wzorce na automaty i odwrotnie.

7.3 Skorowidz

- **Weryfikacja modelowa** (ang. *model-checking*) to metoda weryfikacji polegająca na stworzeniu modelu weryfikowanego systemu w postaci np. automatu skończonego i wykazaniu jego poprawności poprzez badanie stworzonego modelu.

7.4 Praca domowa

1. Dopasuj do siebie automaty i wyrażenia regularne:



- 1) $a(ba \mid b)^*$,
- 2) $b(bb \mid a)^*$,
- 3) $a(aa \mid b)^*$,
- 4) $a(ba \mid a)^*$,
- 5) $b(ab \mid a)^*$.

2. Podaj automat skończony akceptujący język opisany wzorcem: $b(aba \mid aa)^+a$.

3. Podaj wyrażenie regularne opisujące język akceptowany przez następujący automat niedeterministyczny:

	a	b
$\rightarrow F1$	2, 3	4
2		1
3	1	
4	1	1

7.5 Ćwiczenia

1. Dopasuj do siebie automaty i wyrażenia regularne:

a)

	a	b
$\rightarrow 1$	2	2
$F2$		2

b)

	a	b
$\rightarrow F1$	2	2
2	1	1

c)

	a	b
$\rightarrow F1$		2
$F2$	1	2

d)

	a	b
$\rightarrow 1$	2	2
$F2$		1

e)

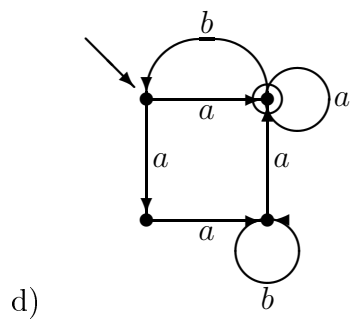
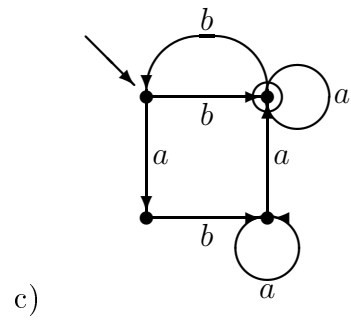
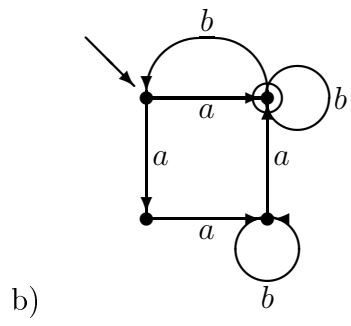
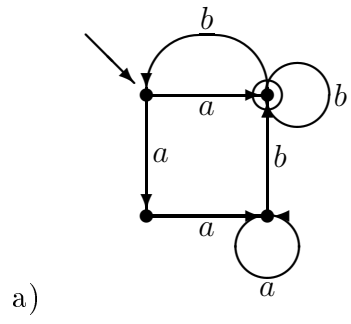
	a	b
$\rightarrow F1$	1	2
$F2$	1	2

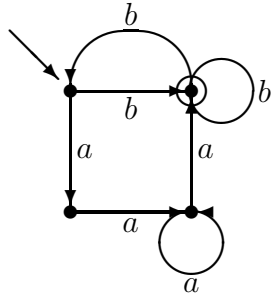
- 1) $(b \mid ba)^*$,
- 2) $(a \mid b)(ba \mid bb)^*$,
- 3) $(a \mid b)^*$,
- 4) $(a \mid b)b^*$,
- 5) $((a \mid b)(a \mid b))^*$.

Rozwiązanie¹⁰

2. Dopasuj automaty do podanych wyrażeń regularnych:

¹⁰ 1-c, 2-d, 3-e, 4-a, 5-b





e)

- 1) $(aaa^*a \cup b)(b^* \cup bb \cup baaa^*a)^*$,
- 2) $(b \cup aba^*a)(a^* \cup bb \cup baba^*a)^*$,
- 3) $(aab^*a \cup a)(b^* \cup ba \cup baab^*a)^*$,
- 4) $(a \cup aaa^*b)(b^* \cup ba \cup baaa^*b)^*$,
- 5) $(a \cup aab^*a)(a^* \cup ba \cup baab^*a)^*$.

Rozwiązanie¹¹

3. Podaj automaty skończone (deterministyczne lub nie) odpowiadające następującym wzorcom:

- $(a(ab)^* \mid a^*bb)^*$,
- $(ab^*a)^* \mid (ba^*b)^*$,
- $(ab \mid abb \mid aabbb)^+$,
- $aa \mid (a \mid bab)b^*$,
- $a(bab)^* \mid a((aa)^* \mid (bb)^*)$,
- $a * [((ba)^* \mid bb)^* \mid aa]^*ab$.

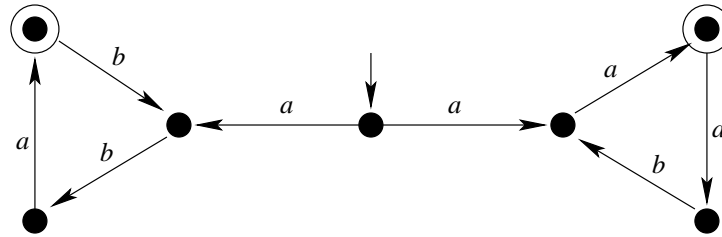
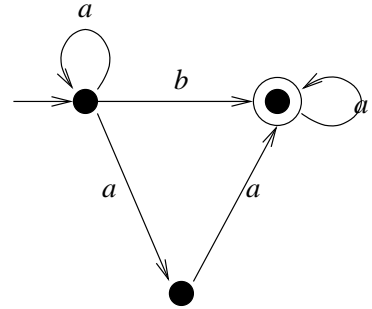
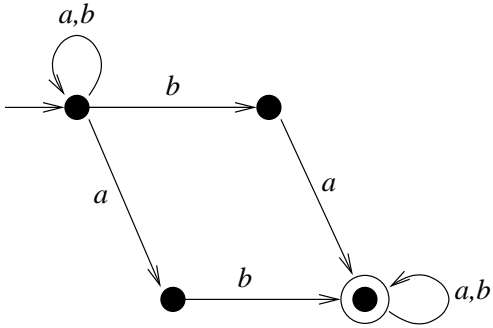
4. Podaj deterministyczne automaty skończone akceptujące języki opisane wyrażeniami regularnymi:

- $a^*b^*a^*$,
- $a^*b^* \mid b^*a^*$,
- $(ab)^* \mid (ba)^*$,
- $a(baba)^+a$,
- $(bab \mid abba)^*aba$,
- $(aba)^* \mid (bab)^*$,

¹¹ 1-e, 2-c, 3-b, 4-a, 5-d.

- $(aba \mid abba)^*bab$,
- $((bab \mid bba)(bab \mid baba))^+$.

5. Podaj wyrażenia regularne odpowiadające automatom:



Wykład 8. Lemat o pompowaniu dla języków regularnych

8.1 Wstęp

Okazuje się, że nie każdy język jest regularny. Standardowy przykład języka, który nie jest regularny, to $A = \{a^n b^n : n \geq 0\}$. Intuicyjnie jest to dosyć oczywiste. Żeby język był regularny, musi być możliwe stwierdzenie czy słowo należy do języka przy użyciu stałej pamięci. Jednak żeby sprawdzić, czy słowo jest postaci $a^n b^n$, musimy najpierw wczytać litery a i zapamiętać ich dokładną liczbę. Zapamiętanie liczby n wymaga przynajmniej $\log_2 n$ bitów pamięci, a to nie jest stała. Tak więc, jaką byśmy pamięcią nie dysponowali, zawsze znajdzie się dostatecznie duże n , takie że po wczytaniu a^n nie będziemy dokładnie wiedzieli ile tych liter a było. Tym samym nie będziemy w stanie sprawdzić, czy liter b jest dokładnie tyle samo.

Spróbujmy sformalizować powyższe rozumowanie. Dowodzimy nie wprost, że A nie jest regularne. Założmy przeciwnie, że A jest regularne. Niech $M = (Q, \Sigma, \delta, s, F)$ będzie automatem akceptującym A i niech $n = |Q|$. Wówczas z zasady szufladkowej Dirichleta istnieje taki stan q , $u \geq 0$ oraz $v > 0$, że $n = u + v$, $\delta^*(s, a^u) = \delta^*(q, a^v) = q$. Wówczas tak samo jest akceptowane, lub nie, słowo $a^u b^n$ i $a^n b^n$. Sprzeczność. A więc założenie, że A jest regularny było fałszywe.

8.2 Lemat o pompowaniu

Powyższe rozumowanie można uogólnić, uzyskując bardziej ogólne narzędzie do wykazywania, że pewne języki nie są regularne.

Tw. 2 (Lemat o pompowaniu dla języków regularnych). *Niech A będzie językiem regularnym. Wówczas istnieje takie k , że dla dowolnych słów x, y i z takich, że $xyz \in A$ oraz $|y| \geq k$, można y podzielić na słowa u, v i w , $y = uvw$ w taki sposób, że $v \neq \varepsilon$ i dla wszystkich $i \geq 0$ $xuv^i w \in A$.*

Dowód: Skoro A jest językiem regularnym, to istnieje deterministyczny automat skończony rozpoznający A — oznaczmy go przez $M = \langle Q, \Sigma, \delta, s, F \rangle$. Wybieramy $k = |Q|$.

Niech x, y i z będą takimi słowami, że $xyz \in A$ oraz $|y| \geq k$. Zauważmy, że w obliczeniu automatu M akceptującego słowo xyz , w trakcie wczytywania słowa y , pojawia się przynajmniej $k + 1$ stanów.

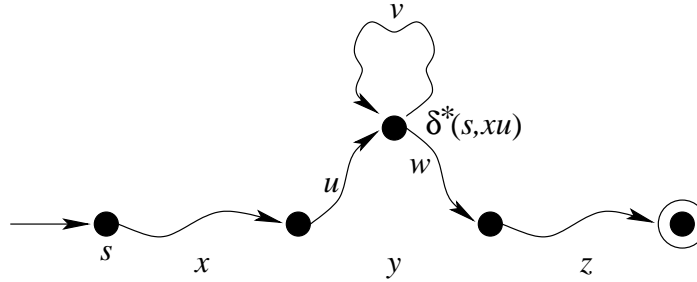
Z zasady szufladkowej Dirichleta wynika, że przynajmniej jeden z tych stanów musi się powtarzać. Weźmy pierwsze powtórzenie się stanu w trakcie wczytywania słowa y .

Możemy więc podzielić słowo y na trzy części $y = uvw$:

- u — od początku słowa y do pierwszego wystąpienia powtarzającego się stanu,
- v — od pierwszego do drugiego wystąpienia powtarzającego się stanu, oraz
- w — od drugiego wystąpienia powtarzającego się stanu do końca słowa y .

Zauważmy, że słowu v odpowiada „pętelka”, tzn. $v \neq \varepsilon$, oraz:

$$\delta^*(s, xu) = \delta^*(s, xuv)$$



Stąd, dla dowolnego $i \geq 0$ mamy:

$$\begin{aligned} \delta^*(s, xuv^i) &= \delta^*(\delta^*(s, xuv), v^{i-1}) = \\ &= \delta^*(\delta^*(s, xu), v^{i-1}) = \\ &= \delta^*(s, xuv^{i-1}) = \\ &\vdots \\ &= \delta^*(s, xu) \end{aligned}$$

Stąd zaś otrzymujemy:

$$\delta^*(s, xuv^i wz) = \delta^*(s, xuvwz) \in F$$

Czyli $xuv^i wz \in A$. □

8.3 Zastosowanie lematu o pompowaniu

Lemat o pompowaniu służy pokazywaniu, że określone języki **nie są** regularne. Schemat takiego dowodu przebiega nie wprost. Lemat o pompowaniu ma postać implikacji: jeśli język jest regularny, to ma pewne własności. Dowód, że język nie jest regularny przebiega według schematu: skoro dany język nie posiada tych własności, to nie może być regularny.

Specjalnie z myślą o takim zastosowaniu powstało alternatywne sformułowanie lematu o pompowaniu:

Lemat o pompowaniu dla języków regularnych — sformułowanie alternatywne:

Niech A będzie językiem. Jeżeli dla każdego $k \geq 0$ istnieją takie słowa x, y, z , że $xyz \in A$, $|y| \geq k$ oraz dla dowolnego podziału słowa y na słowa u, v, w , $y = uvw$ takiego, że $v \neq \varepsilon$, istnieje takie $i \geq 0$, że $xuv^i wz \notin A$, wówczas A nie jest regularny.

Nawet w postaci odwróconej lemat o pompowaniu jest trudny do zrozumienia, gdyż mamy tu do czynienia z czterema poziomami zagnieżdżonych naprzemiennych kwantyfikatorów uniwersalnych i egzystencjalnych. Takie zagnieżdżenie naprzemiennych kwantyfikatorów możemy sobie wyobrazić jako ruchy dwóch graczy — gracze na przemian wskazują

jakie wartości powinny przyjąć zmienne spod kwantyfikatorów, jeden gracz spod uniwersalnych, a drugi spod egzystencjalnych. Pierwszy gracz stara się pokazać, że dany język jest regularny, a drugi, że nie. O wygranej decyduje, czy dla wybranych wartości zmiennych zachodzi własność objęta kwantyfikatorami. Gra jest skończona, więc jeden z graczy ma strategię wygrywającą. W zależności od tego, który z graczy to jest, język może być regularny lub nie.

Lemat o pompowaniu dla języków regularnych — gra z Demonem: Prowadzimy z Demonem grę. Staramy się pokazać, że język nie jest regularny, a Demon stara się pokazać, że język mimo wszystko jest regularny. Zasady gry są następujące:

1. Demon wybiera $k \geq 0$. (Jeśli A jest faktycznie regularny, to Demon może wybrać za k liczbę stanów automatu akceptującego A).
2. Wybieramy takie słowa x, y, z , że $xyz \in A$ i $|y| \geq k$.
3. Demon dzieli y na takie słowa u, v, w , $y = uvw$, że $v \neq \varepsilon$. (Jeśli A jest faktycznie regularny, to Demon może podzielić y na pierwszych dwóch wystąpieniach stanu, który jako pierwszy powtarza się w trakcie wczytywania słowa y .)
4. Wybieramy $i \geq 0$.

Jeśli $xuv^i w z \notin A$, to wygraliśmy, wpp. wygrał Demon.

Oryginalne sformułowanie lematu o pompowaniu mówi, że jeżeli język jest regularny, to Demon ma strategię wygrywającą. Natomiast odwrotne sformułowanie mówi, że jeżeli my mamy strategię wygrywającą, to język nie jest regularny. Pamiętajmy, że lemat o pompowaniu nie służy do pokazywania, że język jest regularny. Nie zawsze pozwala on na pokazanie, że język, który nie jest regularny faktycznie taki nie jest — istnieją języki, które nie są regularne i spełniają lemat o pompowaniu.

Zobaczmy na przykładach jak można zastosować sformułowanie lematu o pompowaniu jak gry z Demonem:

Przykład: Pokażemy, że język $A = \{a^n b^n : n \geq 0\}$ nie jest regularny.

- Demon wybiera k .
- Wybieramy $x = \varepsilon$, $y = a^k$ i $z = b^k$.
- Jakkolwiek Demon by nie podzielił słowa y , to słowo v składa się wyłącznie z literek a i jest niepuste. Możemy wybrać $i = 2$.

Słowo $xuv^2 w z \notin A$, bo ma więcej literek a niż b , czyli wygraliśmy.

Powyższa strategia jest dla nas strategią wygrywającą, a zatem język A nie jest regularny.

Przykład: Pokażemy, że język $B = \{a^{2^n} : n \leq 0\}$ nie jest regularny.

- Demon wybiera k .
- Wybieramy $x = \varepsilon$, $y = a^{2^k-1}$, $z = a$.
- Jakkolwiek Demon by nie podzielił słowa y , to słowo v składa się wyłącznie z literek a oraz $0 < |v| < 2^k$. Możemy wybrać $i = 2$.

Słowo $xuv^2wz \notin B$, bo $2^k < |xuv^2wz| < 2 \cdot 2^k = 2^{k+1}$ a zatem wygraliśmy.

Powyższa strategia jest dla nas strategią wygrywającą, a zatem język B nie jest regularny.

Przykład: Pokażemy, że język $C = \{a^{n!} : n \geq 0\}$ nie jest regularny.

- Demon wybiera k .
- Wybieramy $x = \varepsilon$, $y = a^{(k+1)!-1}$, $z = a$.
- Demon wybiera słowa u, v, w . Mamy $0 < |v| < (k+1)!$.
- Wybieramy $i = k + 2$.

Mamy $|xuv^i w z| = (k+1)! + (i-1) \cdot |v|$. Tak więc, z jednej strony mamy $|xuv^i w z| > (k+1)!$, ale z drugiej:

$$|xuv^i w z| = (k+1)! + (k+1) \cdot |v| < (k+1)! + (k+1) \cdot (k+1)! = (k+2)(k+1)! = (k+2)!$$

Stąd, $|xuv^i w z| \notin C$, czyli wygraliśmy.

Powyższa strategia jest dla nas strategią wygrywającą, a zatem język C nie jest regularny.

8.4 Inne metody dowodzenia nieregularności języków

Wykorzystanie lematu o pompowaniu nie jest jedynym sposobem pokazywania, że jakiś język nie jest regularny. Jeśli już wiemy o jakimś języku (B), że nie jest regularny, to możemy pokazać, że jakiś inny język (A) nie jest regularny korzystając z własności domknięcia klasy języków regularnych. Rozumujemy przy tym nie wprost — gdyby język A był regularny, to B też, a że B nie jest, to i A nie może być.

Przykład: Pokażemy, że język $D = \{a^n b^m : n \neq m\}$ nie jest regularny. Gdyby był on regularny, to (ze względu na domknięcie klasy języków regularnych na dopełnienie) język $A = \{a^n b^n : n \geq 0\}$ byłby też regularny, a nie jest. Tak więc język D nie może być regularny.

Przykład: Rozważmy język $E \subseteq \{[,]\}$ złożonych z poprawnych wyrażeń nawiasowych. Pokażemy, że nie jest on regularny. Gdyby był on regularny, to $E \cap L([*]^*) = \{[n]^n : n \geq 0\}$ też byłby regularny. Natomiast język $\{[n]^n : n \geq 0\}$, tak jak $A = \{a^n b^n : n \geq 0\}$, nie jest regularny. Tak więc język E nie może być regularny.

8.5 Podsumowanie

Wykład ten był poświęcony przede wszystkim lematowi o pompowaniu dla języków regularnych. Poznaliśmy trzy równoważne sformułowania tego lematu, oraz przykłady jego zastosowania do wykazywania, że określone języki nie są regularne. Poznaliśmy też metodę wykazywania, że dane języki nie są regularne, w oparciu o własności domknięcia klasy języków regularnych.

8.6 Skorowidz

- **Lemat o pompowaniu dla języków regularnych** mówi, że dla każdego języka regularnego, każde dostatecznie długie słowo w tym języku można „pompować” dowolnie je wydłużając. Wykazywanie braku tej własności jest metodą pokazywania, że dany język nie jest regularny.
- **Gra z Demonym** — alternatywne sformułowanie lematu o pompowaniu dla języków regularnych jako gry.

8.7 Praca domowa

1. Które z poniższych języków są regularne, a które nie? (Odpowiedzi nie musisz formalnie uzasadniać.)
 - a) $\{a^{n^3} : n \geq 0\}$,
 - b) $\{a^i b^j : i \cdot j \leq 42\}$,
 - c) $\{a^i b^j : i \equiv j \pmod{42}\}$,
 - d) $\{ww : w \in \{a, b\}^*\}$,
 - e) $\{a^n b^m : |n - m| \leq 42\}$
2. Pokaż, że język $\{a^m b^n a^{m \cdot n} : m, n \geq 0\}$ nie jest regularny.
3. Pokaż, że język zawierający te słowa nad alfabetem $\{a, b\}$, które zawierają tyle samo liter a i b , $\{x \in \{a, b\}^* : \#_a(x) = \#_b(x)\}$, nie jest regularny.

8.8 Ćwiczenia

1. Pokaż, że następujące języki nie są regularne:

- $\{a^n b^{2^n} : n \geq 0\}$,
- język złożony z palindromów nad alfabetem $\{a, b\}$, $\{x \in \{a, b\}^* : x = \text{rev}(x)\}$,
- $\{a^p : p \text{ jest liczbą pierwszą}\}$,
- $\{a^{n^2} : n \geq 0\}$,
- $\{a^i b^j c^k : i = j \vee j = k\}$

2. Które z podanych języków są regularne, a które nie? Odpowiedzi uzasadnij.

- $\{a^i b^j : 0 \leq i \leq j \leq 42\}$,
- komentarze w Pascalu,
- $\{a^i b^j c^k : i \neq j \vee j \neq k\}$,
- $\{a^i b^j a^{2(i+j)} : i, j \geq 0\}$,
- $\{x \in \{a, b\}^* : \#_a(x) \neq \#_b(x)\}$,
- $\{a^{2^n} : n \geq 0\}$,
- $\{a^i b^j c^k : i + j + k \leq 42\}$,
- poprawne programy w Pascalu,
- $\{a^{42^n} : n \geq 0\}$,
- $\{x \in \{a, b\}^* : \#_a(x) \text{ jest podzielne przez } 42\}$.

Rozwiązanie¹²

¹² Regularne są: a), b), f), g), j).

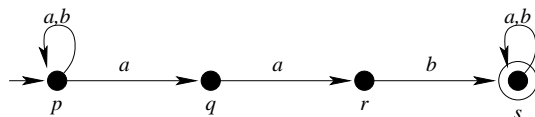
Wykład 9. Minimalizacja deterministycznych automatów skończonych

9.1 Wstęp

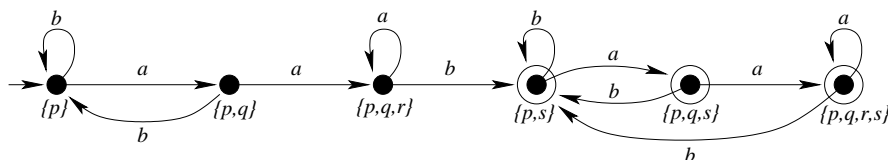
W tym wykładzie zajmiemy się problemem minimalizacji liczby stanów w deterministycznych automatach skończonych. Konstruując automaty skończone często możemy dość do kilku rozwiązań o różnej liczbie stanów. Jak znaleźć automat o najmniejszej możliwej liczbie stanów? W przypadku deterministycznych automatów skończonych znany jest efektywny algorytm przekształcający dowolny automat akceptujący dany język w automat o minimalnej liczbie stanów.

Zanim poznamy ten algorytm i udowodnimy jego poprawność, zobaczymy na przykładzie jak można zminimalizować liczbę stanów.

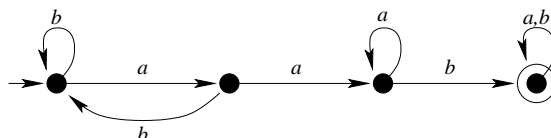
Przykład: Weźmy na początek niedeterministyczny automat skończony akceptujący te słowa (nad alfabetem $\{a, b\}$), które zawierają pod słowo aab .



Automat ten ma 4 stany, a więc odpowiadający mu automat potęgowy ma 16 stanów. Oczywiście możemy się ograniczyć do stanów osiągalnych, a tych jest 6.



Zauważmy, że możemy skleić ze sobą wszystkie stany akceptujące — po ich osiągnięciu nie można już ich opuścić, a więc słowo musi zostać zaakceptowane.



Uzyskaliśmy automat deterministyczny o 4 stanach. Czy to jest minimum? Zastanówmy się ile stanów jest niezbędnych. Potrzebny jest jeden stan akceptujący, do którego wchodzimy, gdy wczytamy pod słowo aab . Dodatkowo potrzebujemy przynajmniej 3 inne stany. Możemy je rozróżnić m.inn. po najkrótszych słowach koniecznych do osiągnięcia stanu akceptującego. Dla stanu początkowego jest to oczywiście słowo aab . Dla stanu w jakim jesteśmy po wczytaniu a jest to ab , a dla stanu, w którym jesteśmy po wczytaniu aa jest to b . Razem uzyskujemy 4 stany, a więc nasz automat jest minimalny.

9.2 Dowód istnienia automatu minimalnego

Jak widać z powyższego przykładu niektóre stany w automacie deterministycznym mogą być sobie „równoważne” i możemy je skleić. Natomiast badając, czy stany są równoważne powinniśmy zwrócić uwagę na to jakie słowa z danych stanów prowadzą do stanów akceptujących. Jeśli są tu jakieś różnice, to stanów nie możemy sklejać. Narazie są to mało precyzyjne intuicje, ale za chwilę sformalizujemy je.

Nadal jednak pozostaje otwarte pytanie, czy usunięcie stanów nieosiągalnych i sklejenie stanów „równoważnych” prowadzi do jednego automatu minimalnego? Czy też zaczynając od różnych automatów możemy uzyskać różne automaty minimalne? Pokażemy, że dla każdego języka regularnego istnieje jeden (z dokładnością do izomorfizmu) minimalny deterministyczny automat skończony.

Def. 33. Niech A będzie ustalonym językiem. Przez A/x oznaczamy język postaci:

$$A/x = \{y \in \Sigma^* : xy \in A\}$$

Języki postaci A/x dla $x \in \Sigma^*$ nazywamy *językami ilorazowymi*. □

Inaczej mówiąc, A/x to język złożony ze wszystkich tych ciągów dalszych, które uzupełniają wystąpienia prefiksów x w słowach z A . W szczególności $A/\varepsilon = A$. Wprost z definicji języka A/x wynika następujący fakt:

Fakt 7. Dla dowolnych $x, y \in \Sigma^*$ mamy $A/xy = (A/x)/y$.

Stąd zaś wynika kolejny fakt:

Fakt 8. Dla dowolnych $x, y, z \in \Sigma^*$, jeśli $A/x = A/y$, to $A/xz = A/yz$.

Jest tak gdyż:

$$A/xz = (A/x)/z = (A/y)/z = A/yz$$

Jeżeli dodatkowo A jest regularny to pojawiają się dodatkowe intuicje. Niech $M = \langle Q, \Sigma, \delta, s, F \rangle$ będzie deterministycznym automatem skończonym akceptującym A , $L(M) = A$, w którym wszystkie stany są osiągalne. A/x to język złożony z tych słów, które prowadzą ze stanu $\delta^*(s, x)$ do stanów akceptujących, (czyli jest to język akceptowany przez automat postaci $\langle Q, \Sigma, \delta, \delta^*(s, x), F \rangle$). Ponieważ to, jakie słowa prowadzą z danego stanu do stanów akceptujących zależy od tego stanu, ale nie od tego jak dostaliśmy się do tego stanu, więc zachodzi następujący fakt:

Fakt 9. Niech $x, y \in \Sigma^*$. Jeśli $\delta^*(s, x) = \delta^*(s, y)$, to $A/x = A/y$.

Możemy więc każdemu stanowi q automatu M przyporządkować jeden z języków ilorazowych A/x języka A — dokładnie taki, dla którego $\delta^*(s, x) = q$. Wynika stąd kolejny fakt:

Fakt 10. Jeśli A jest regularny, to istnieje skończenie wiele języków ilorazowych języka A . Dokładniej, jest ich co najwyżej tyle, co stanów osiągalnych w deterministycznym automacie skończonym akceptującym A .

Okazuje się, że odwrotna implikacja jest również prawdziwa.

Tw. 3 (uproszczone twierdzenie Myhill-Nerode'a). *Niech A będzie ustalonym językiem. Następujące stwierdzenia są równoważne:*

- A jest regularny,
- istnieje skończenie wiele języków ilorazowych języka A .

Dowód: Implikacja w jedną stronę wynika z poprzedniego faktu. Musimy pokazać, że jeżeli istnieje skończenie wiele języków ilorazowych języka A , to A jest regularny. Skonstruujemy deterministyczny automat skończony akceptujący A . Stanami tego automatu będą języki ilorazowe języka A , $Q = \{A/x : x \in \Sigma^*\}$. Nasz automat skonstruujemy tak, że każdy stan jako język będzie zawierał dokładnie te słowa, które z tego stanu prowadzą do stanów akceptujących:

$$A/x = \{y \in \Sigma^* : \delta^*(A/x, y) \in F\}$$

Stanem początkowym jest oczywiście $A = A/\varepsilon$, gdyż nasz automat ma akceptować właśnie język A . Stanami akceptującymi są te stany, które jako języki zawierają słowo puste,

$$F = \{A/x : \varepsilon \in A/x\}$$

Czyli:

$$F = \{A/x : x \in A\}$$

Natomiast funkcja przejścia jest postaci:

$$\delta(A/x, a) = A/xa$$

Powyższa definicja jest poprawna, gdyż jeśli $A/x = A/y$, to $A/xa = A/ya$.

Trzeba jeszcze pokazać, że automat $M = \langle Q, \Sigma, \delta, A, F \rangle$ akceptuje język A . Niech $y \in \Sigma^*$, $y = a_1 a_2 \dots a_k$. Zauważmy, że:

$$\begin{aligned} \delta^*(A/x, y) &= \delta^*(A/x, a_1 a_2 \dots a_k) = \\ &= \delta^*(\delta(A/x, a_1), a_2 \dots a_k) = \\ &= \delta^*(A/x a_1, a_2 \dots a_k) = \\ &\quad \vdots \\ &= \delta^*(A/x a_1, a_2 \dots a_k, \varepsilon) = \\ &= A/x a_1, a_2 \dots a_k = A/xy \end{aligned}$$

Stąd, język akceptowany przez M to:

$$\begin{aligned} L(M) &= \{x \in \Sigma^* : \delta^*(A, x) \in F\} = \\ &= \{x \in \Sigma^* : \delta^*(A/\varepsilon, x) \in F\} = \\ &= \{x \in \Sigma^* : A/x \in F\} = \\ &= \{x \in \Sigma^* : \varepsilon \in A/x\} = \\ &= \{x \in \Sigma^* : x \in A\} = A \end{aligned}$$

□

9.3 Konstrukcja automatu minimalnego

Zauważmy, że automat skonstruowany w powyższym dowodzie ma minimalną liczbę stanów — gdyż ma ich dokładnie tyle ile jest języków ilorazowych języka A . Pokażemy jak z dowolnego deterministycznego automatu skończonego akceptującego język A zrobić automat izomorficzny ze skonstruowanym powyżej.

Niech $M = \langle Q, \Sigma, \delta, s, F \rangle$ będzie dowolnym ustalonym deterministycznym automatem skończonym akceptującym język A , $L(M) = A$, w którym wszystkie stany są osiągalne. (Jeśli M zawiera stany nieosiągalne, to najpierw je usuwamy.) Każdemu stanowi $q \in Q$ możemy przypisać język złożony z tych słów, które prowadzą z q do stanów akceptujących. Język taki to A/x , gdzie x prowadzi z s do q , $\delta^*(s, x) = q$.

Niech $M_{/\approx}$ będzie automatem powstałym z M przez sklejenie ze sobą stanów, którym przypisaliśmy takie same języki ilorazowe. Ścisłe rzecz biorąc, na stanach automatu M definiujemy relację równoważności $\approx \subseteq Q \times Q$ określoną w następujący sposób:

$$p \approx q \text{ w.t.w. } \forall x \in \Sigma^* \delta^*(p, x) \in F \Leftrightarrow \delta^*(q, x) \in F$$

Inaczej mówiąc, niech A/x i A/y będą językami ilorazowymi przypisanymi stanom p i q , wówczas $p \approx q$ wtw., gdy $A/x = A/y$.

Sklejamy te stany automatu M , które są sobie równoważne:

$$M_{/\approx} = \langle Q_{/\approx}, \Sigma, \delta_{/\approx}, [s]_{/\approx}, F_{/\approx} \rangle$$

gdzie:

$$\delta_{/\approx}([q]_{/\approx}, a) = [\delta(q, a)]_{/\approx}$$

Automat $M_{/\approx}$ jest nazywany automatem *ilorazowym*.

Po pierwsze musimy się upewnić, czy $\delta_{/\approx}$ jest poprawnie zdefiniowana. Musimy w tym celu pokazać, że jeżeli sklejamy ze sobą dwa stany $p, q \in Q$, $p \approx q$, to dla dowolnego $a \in \Sigma$ mamy również $\delta(p, a) \approx \delta(q, a)$. Skoro $p \approx q$, to $A/x = A/y$, a zatem $A/xa = A/ya$, czyli $\delta(p, a) \approx \delta(q, a)$.

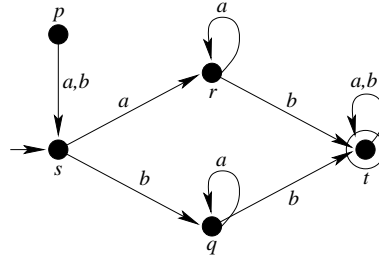
Zauważmy, że nigdy nie sklejimy stanu akceptującego z nieakceptującym, gdyż stany akceptujące charakteryzują się tym, że przypisane im języki zawierają słowa puste.

Czy automat $M_{/\approx}$ akceptuje ten sam język co M , $L(M_{/\approx}) = L(M)$? Można pokazać (przez indukcję), że dla dowolnego stanu $q \in Q$ i słowa $x \in \Sigma^*$ mamy $\delta_{/\approx}^*([q]_{/\approx}, x) = [\delta^*(q, x)]_{/\approx}$. Ponieważ jednak nie skleiliśmy żadnego stanu akceptującego z żadnym ze stanów nieakceptujących, mamy:

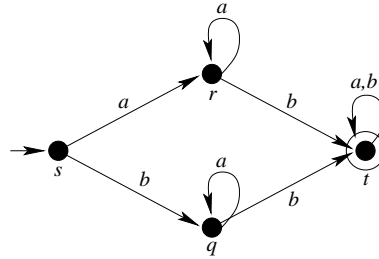
$$\begin{aligned} L(M_{/\approx}) &= \{x \in \Sigma^* : \delta_{/\approx}^*([s]_{/\approx}, x) \in F_{/\approx}\} = \\ &= \{x \in \Sigma^* : [\delta^*(s, x)]_{/\approx} \in F_{/\approx}\} = \\ &= \{x \in \Sigma^* : \delta^*(s, x) \in F\} = L(M) = A \end{aligned}$$

Tak więc usuwając stany nieosiągalne i sklejając ze sobą stany, którym odpowiadają te same języki ilorazowe otrzymujemy ten sam minimalny deterministyczny automat skończony akceptujący dany język regularny.

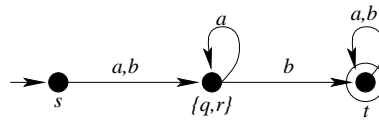
Przykład: Weźmy następujący automat deterministyczny:



Stan p jest nieosiągalny. Po jego usunięciu mamy:



Tylko dwa stany możemy skleić ze sobą: p i q , $p \approx q$. Po ich sklejeniu uzyskujemy;



9.4 Algorytm minimalizacji

Pokażemy jak można zaimplementować opisaną w poprzednim punkcie konstrukcję. Algorytm minimalizacji składa się z dwóch faz:

1. usunięcia stanów nieosiągalnych,
2. wyznaczenia relacji równoważności \approx i sklejenia ze sobą stanów równoważnych.

Znalezienie stanów nieosiągalnych jest prostsze. Trudniej natomiast znaleźć stany równoważne. Przypomnijmy, że dwa stany p i q są sobie równoważne, gdy:

$$\forall x \in \Sigma^* \delta^*(p, x) \in F \Leftrightarrow \delta^*(q, x) \in F$$

Pokażemy algorytm, który znajduje wszystkie pary nierównoważnych sobie stanów. Siłą rzeczy, pozostałe pary stanów są sobie równoważne i można je skleić ze sobą.

Jeśli stany p i q nie są sobie równoważne, to istnieje rozróżniające je słowo x , takie że:

$$\delta^*(p, x) \in F \not\approx \delta^*(q, x) \in F$$

Początkowo zakładamy, że wszystkie stany można skleić ze sobą. Następnie sukcesywnie wyznaczamy pary stanów które nie są sobie równoważne — w kolejności wg. rosnącej długości najkrótszych rozróżniających je słów.

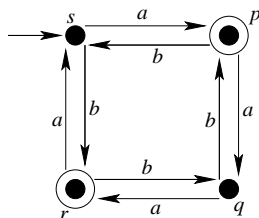
Algorytm

1. Tworzymy tablicę wartości logicznych $T_{\{p,q\}}$ indeksowaną nieuporządkowanymi parami $\{p, q\}$ stanów $p, q \in Q$. Początkowo $T_{\{p,q\}} = \text{true}$ dla wszystkich $p, q \in Q$.
2. Dla wszystkich takich par stanów $\{p, q\}$, że $p \in F$ i $q \notin F$ zaznaczamy $T_{\{p,q\}} = \text{false}$. Jeśli p jest akceptujący, a q nie, to stany te rozróżnia słowo puste ε .
3. Dla wszystkich par stanów $\{p, q\}$ i znaków $a \in \Sigma$ takich, że $T_{\{\delta(p,a), \delta(q,a)\}} = \text{false}$, zaznaczamy również $T_{\{p,q\}} = \text{false}$.
4. Jeżeli w kroku 3 zmieniliśmy choć jedną komórkę tablicy T z **true** na **false**, to powtarzamy krok 3 — tak długo, aż nie będzie on powodował żadnych zmian w tablicy T .
5. Na koniec mamy: $p \approx q \Leftrightarrow T_{\{p,q\}}$.

Dowód poprawności algorytmu: Poprawność algorytmu wynika stąd, że dla każdej pary stanów $\{p, q\}$, które nie są sobie równoważne istnieje pewne słowo, które je rozróżnia. Weźmy najkrótsze takie słowo. W kroku 2 wyznaczamy wszystkie takie pary, które są rozróżniane przez ε . Z każdym powtórzeniem kroku 3 wyznaczamy wszystkie takie pary, które są rozróżniane przez słowa o jeden dłuższe.

Ponieważ tablica T ma skończoną liczbę komórek, więc krok 3 jest powtarzany ograniczoną liczbę razy, a powyższy algorytm zawsze się zatrzyma.

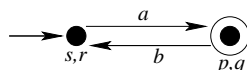
Przykład: Weźmy następujący automat deterministyczny:



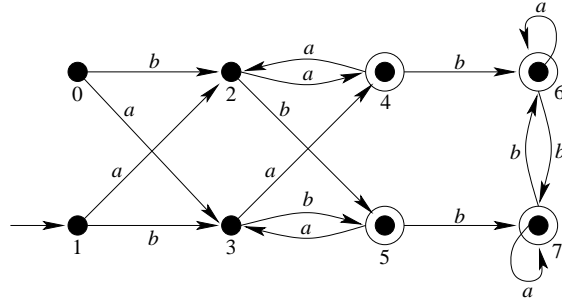
Wszystkie stany są osiągalne. Tablica T obliczona dla tego automatu uzyskuje swą postać już w kroku 2:

p			
F	q		
T	F	r	
F	T	F	s

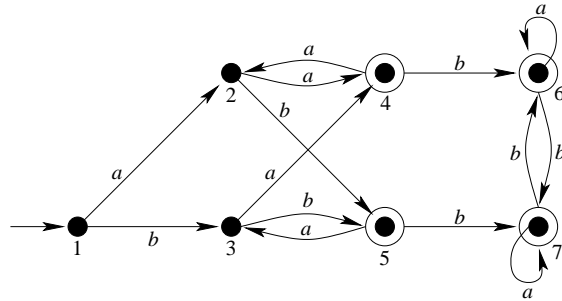
Co oznacza, że możemy skleić ze sobą stany s i q oraz p i r :



Przykład: Weźmy następujący automat deterministyczny:



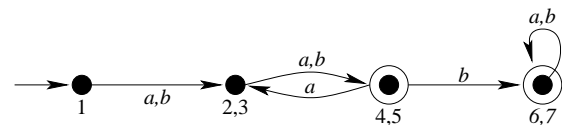
Stan 0 jest nieosiągalny. Po jego usunięciu mamy:



Tablica T obliczona dla tego automatu uzyskuje swą postać już w kroku 2:

1							
F	2						
F	T	3					
F	F	F	4				
F	F	F	T	5			
F	F	F	F	F	6		
F	F	F	F	F	T	7	

Co oznacza, że możemy skleić ze sobą stany: 2 z 3, 4 z 5, oraz 6 z 7.



9.5 Podsumowanie

W tym wykładzie poznaliśmy (uproszczone) twierdzenie Myhilla-Nerode'a. Twierdzenie to dostarcza nam jeszcze jednego kryterium do badania czy języki są regularne — każdy język regularny ma skończenie wiele języków ilorazowych. Jako wniosek z dowodu tego twierdzenia uzyskaliśmy, że istnieje jeden (z dokładnością do izomorfizmu) minimalny automat deterministyczny akceptujący dany język. Poznaliśmy też algorytm wyznaczania tego minimalnego automatu.

9.6 Skorowidz

- **Automat ilorazowy** to minimalny deterministyczny automat skończony powstały przez sklejenie stanów równoważnych (zgodnie z daną relacją równoważności na stanach, spełniającą pewne dodatkowe warunki).
- **Algorytm minimalizacji** deterministycznych automatów skończonych — algorytm wyznaczania minimalnego deterministycznego automatu skończonego na podstawie dowolnego deterministycznego automatu skończonego akceptującego dany język.
- **Język ilorazowy** języka A , to każdy z języków postaci $A/x = \{y \in \Sigma^* : xy \in A\}$ (dla $x \in \Sigma^*$).

9.7 Praca domowa

Zminimalizuj następujące deterministyczne automaty skończone:

1.

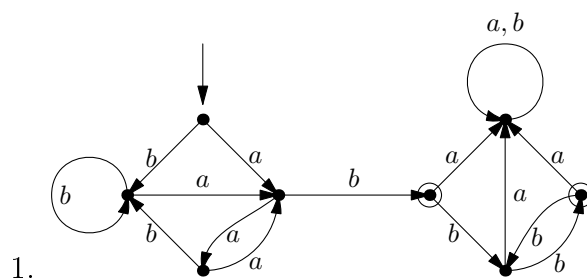
	a	b
$\rightarrow 1$	3	4
2	4	3
$F3$	2	1
$F4$	1	2

2.

	a	b
$\rightarrow 1$	3	4
2	5	6
3	2	3
4	3	1
$F5$	1	5
6	5	2

9.8 Ćwiczenia

Zminimalizuj następujące deterministyczne automaty skończone:



2.

	a	b
$\rightarrow 1$	4	2
$F2$	3	2
3	4	2
$F4$	4	3

3.

	a	b
$\rightarrow F1$	2	1
2	3	2
$F3$	4	3
4	1	4

4.

	a	b
$\rightarrow 1$	3	4
2	4	3
3	4	1
$F4$	1	3

5.

	a	b
$\rightarrow 1$	6	3
2	2	6
$F3$	5	4
4	5	1
5	1	4
6	5	1

6.

	a	b
1	2	4
$F2$	5	1
$\rightarrow 3$	2	3
4	5	1
$F5$	2	4
6	5	6

7.

	a	b
$\rightarrow 1$	6	6
$F2$	3	2
3	6	1
4	1	6
5	6	1
$F6$	5	3

8.

	<i>a</i>	<i>b</i>
0	4	4
1	0	6
2	3	6
3	0	6
4	0	0
→ 5	3	1
<i>F6</i>	5	4

9.

	<i>a</i>	<i>b</i>
0	2	4
1	4	5
2	4	3
→ 3	1	5
4	6	3
<i>F5</i>	2	5
6	6	3

10.

	<i>a</i>	<i>b</i>
→ 0	2	5
1	2	4
<i>F2</i>	5	1
3	2	5
4	5	0
5	1	1
6	5	2
7	5	3

11.

	<i>a</i>	<i>b</i>
→ 0	2	4
<i>F1</i>	2	3
2	3	0
3	0	0
4	3	1
<i>F5</i>	2	3
6	3	2
7	3	5

12.

	a	b
$\rightarrow F1$	2	5
$2F$	1	4
3	7	2
4	5	7
5	4	3
6	3	6
7	3	1

Wykład 10. Języki bezkontekstowe

10.1 Wstęp

W tym wykładzie poznamy *gramatyki bezkontekstowe* i opisywaną przez nie klasę *języków bezkontekstowych*. Gramatyki bezkontekstowe są formalizmem służącym do opisu składni — podobnie jak wyrażenia regularne. Podobnie również jak wyrażenia regularne są maksymalnie uproszczonym formalizmem, a wzorce są ich rozbudowaną wersją przeznaczoną do praktycznych zastosowań, tak i gramatyki bezkontekstowe mają swoją rozbudowaną wersję przeznaczoną do praktycznych zastosowań. Są to: notacja Backusa-Naura (BNF, ang. *Backus-Naur form*) oraz rozszerzona notacja Backusa-Naura (EBNF, ang. *extended Backus-Naur form*). Jest bardzo prawdopodobne, że Czytelnik zetknął się z nimi. Notacje te są zwykle używane w podręcznikach do ścisłego opisu składni języków programowania.

Będziemy również badać właściwości klasy języków, które można opisać gramatykami bezkontekstowymi — tzw. klasy języków bezkontekstowych. Jak zobaczymy, siła wyrazu gramatyk bezkontekstowych będzie większa niż wyrażen regularnych i automatów skończonych. Inaczej mówiąc, klasa języków bezkontekstowych zawiera klasę języków regularnych.

10.2 BNF

W pierwszym przybliżeniu BNF jest formalizmem przypominającym definicje nazwanych wzorców ze specyfikacji dla Lex'a. Główna różnica polega na tym, że zależności między poszczególnymi nazwami mogą być rekurencyjne.

Specyfikacja w BNF-ie określa składnię szeregu *konstrukcji składniowych*, nazywanych też *nieterminalami*. Każda taka konstrukcja składniowa ma swoją nazwę. Specyfikacja składa się z szeregu reguł (nazywanych *produkcjami*), które określają jaką postać mogą przybierać konstrukcje składniowe. Nazwy konstrukcji składniowych ujmujemy w trójkątne nawiasy $\langle \dots \rangle$. Produkcje mają postać:

$$\langle \text{definiowana konstrukcja} \rangle ::= \text{wyrażenie opisujące definiowaną konstrukcję}$$

Wyrażenie opisujące definiowaną konstrukcję może zawierać:

- $\langle \text{konstrukcja} \rangle$ — nazwy konstrukcji, w tym również definiowanej, rekurencyjne zależności są dozwolone,
- tekst, ujęty w cudzysłów, który pojawia się dosłownie w danej konstrukcji,
- $\dots | \dots$ — oddziela różne alternatywne postaci, jakie może mieć dana konstrukcja,
- $[\dots]$ — fragment ujęty w kwadratowe nawiasy jest opcjonalny.

Produkcje traktujemy jak możliwe reguły przepisywania — nazwę konstrukcji stojącej po lewej stronie produkcji możemy zastąpić dowolnym napisem pasującym do wyrażenia podanego po prawej stronie. Każdej konstrukcji składniowej odpowiada pewien język. Składa

sie on z wszystkich tych napisów, które można uzyskać zaczynając od nazwy konstrukcji i stosując produkcje, aż do momentu uzyskania napisu, który nie zawiera już żadnych nazw konstrukcji.

Przykład: BNF może służyć do opisywania składni najrozmaitszych rzeczy. Oto opis składni adresów pocztowych:

```

⟨adres⟩           ::= ⟨adresat⟩ ⟨adres lokalny⟩ ⟨adres miasta⟩ ⟨adres kraju⟩
⟨adresat⟩         ::= ["W.P." | "Sz.Pan."] ⟨napis⟩ ", "
⟨adres lokalny⟩  ::= ⟨ulica⟩ ⟨numer⟩ ["/" ⟨numer⟩ ] ["m" ⟨numer⟩ | "/" ⟨numer⟩ ], "
⟨adres miasta⟩   ::= [ ⟨kod⟩ ] ⟨napis⟩
⟨adres kraju⟩    ::= [ ", " ⟨napis⟩ ]
⟨kod⟩            ::= ⟨cyfra⟩ ⟨cyfra⟩ "-" ⟨cyfra⟩ ⟨cyfra⟩ ⟨cyfra⟩
⟨cyfra⟩          ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9"
⟨numer⟩          ::= ⟨cyfra⟩ [ ⟨numer⟩ ]

```

Specyfikacja ta nie jest oczywiście kompletna, gdyż nie definiuje czym jest napis. Zwróćmy uwagę na to, że definicja numer jest rekurencyjna.

10.3 Gramatyki bezkontekstowe

Gramatyki bezkontekstowe są uproszczoną wersją notacji BNF. Konstrukcje składniowe nazywamy tu *nieterminalami* lub *symbolami nieterminalnymi* i będziemy je oznaczać wielkimi literami alfabetu łacińskiego: A, B, \dots, X, Y, Z . Wśród nieterminali jest jeden wyróżniony symbol, nazywany *aksjomatem*. To właśnie język związany z tym nieterminalem opisuje gramatyka.

Oprócz nieterminali używamy również *terminali* (*symboli terminalnych*) — są to znaki stanowiące alfabet opisywanego przez gramatykę języka. Będziemy je oznaczać małymi literami alfabetu łacińskiego: a, b, \dots .

Produkcje w gramatykach bezkontekstowych mają bardzo prostą postać: $A \rightarrow \alpha$, gdzie A to nieterminal, a α to słowo, które może zawierać terminale i nieterminale. Produkcja taka oznacza, że nieterminal A **możemy** zastąpić napisem α . Dla jednego nieterminala możemy mieć wiele produkcji. Oznacza to, że możemy je dowolnie stosować. Zwykle zamiast pisać:

$$\begin{aligned}
 A &\rightarrow \alpha \\
 A &\rightarrow \beta \\
 &\vdots \\
 A &\rightarrow \gamma
 \end{aligned}$$

Będziemy pisać krócej:

$$A \rightarrow \alpha \mid \beta \mid \dots \mid \gamma$$

Przykład: Zanim formalnie zdefiniujemy gramatyki bezkontekstowe i opisywane przez nie języki, spójrzmy na gramatykę opisującą język: $\{a^n b^n : n \geq 0\}$.

$$A \rightarrow aAb \mid \varepsilon$$

Oto sekwencja zastosować produkcji, która prowadzi do uzyskania słowa $aaabbb$:

$$A \rightarrow aAb \rightarrow aaAbb \rightarrow aaaAbbb \rightarrow aaabbb$$

Trzykrotnie stosowaliśmy tutaj produkcję $A \rightarrow aAb$, aby na koniec zastosować $A \rightarrow \varepsilon$.

Przykład ten pokazuje, że gramatyki bezkontekstowe mogą opisywać języki, które nie są regularne.

Def. 34. *Gramatyka bezkontekstowa*, to dowolna taka czwórka $G = \langle N, \Sigma, P, S \rangle$, gdzie:

- N to (skończony) alfabet symboli *nieterminalnych*,
- Σ to (skończony) alfabet symboli *terminalnych*, $\Sigma \cap N = \emptyset$,
- P to skończony zbiór *produkcji*, reguł postaci $A \rightarrow \alpha$ dla $A \in N$, $\alpha \in (N \cup \Sigma)^*$,
- $S \in N$ to *aksjomat* wyróżniony symbol nieterminalny.

□

Zwykle nie będziemy podawać gramatyki jako formalnej czwórki. Zamiast tego będziemy podawać zbiór produkcji. Zwykle, w sposób niejawni z produkcji wynika zbiór nieterminali i terminali. Jeżeli nie będzie oczywiste, który nieterminal jest aksjomatem, będziemy to dodatkowo zaznaczać.

Def. 35. Niech $G = \langle N, \Sigma, P, S \rangle$ będzie ustaloną gramatyką bezkontekstową. Przez \rightarrow będziemy oznaczać zbiór produkcji P traktowany jako relacja, $\rightarrow \subseteq N \times (N \cup \Sigma)^*$. Przez \rightarrow^* będziemy oznaczać zwrotno-przechodnie domknięcie \rightarrow , $\rightarrow^* \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$. □

Relacja \rightarrow opisuje pojedyncze zastosowanie produkcji, jako relacja między nieterminalem, a zastępującym go słowem. Relacje \rightarrow^* opisuje co można zrobić stosując dowolną liczbę (łącznie z zero) dowolnych produkcji. Możemy też przedstawić sobie relację \rightarrow^* jako skrót do wielokrotnego iterowania relacji \rightarrow :

Fakt 11. *Jeżeli $x \rightarrow^* y$, to istnieje ciąg słów (z_i) taki, że:*

$$x = z_0 \rightarrow z_1 \rightarrow \cdots \rightarrow z_k = y$$

Gramatyka opisuje język złożony z tych wszystkich słów (nad alfabetem terminalnym), które możemy uzyskać z aksjomatu stosując produkcje.

Def. 36. Niech $G = \langle N, \Sigma, P, S \rangle$ będzie ustaloną gramatyką bezkontekstową. Język generowany (lub opisywany) przez gramatykę G , to:

$$L(G) = \{x \in \Sigma^* : S \rightarrow^* x\}$$

Alternatywnie, $L(G)$ to zbiór takich słów $x \in \Sigma^*$, dla których istnieją ciągi słów (z_i) , $z_i \in (N \cup \Sigma)^*$ takie, że:

$$S = z_0 \rightarrow z_1 \rightarrow \dots \rightarrow z_k = x$$

Ciąg (z_i) nazywamy *wyprowadzeniem* słowa x (w gramatyce G). □

Def. 37. Powiemy, że język jest bezkontekstowy, jeżeli istnieje generująca go gramatyka. □

Wyprowadzenie stanowi coś w rodzaju dowodu, że dane słowo faktycznie należy do języka generowanego przez gramatykę. Istnieje jeszcze inny, bardziej strukturalny, sposób ilustrowania jak dane słowo można uzyskać w danej gramatyce. Jest to *drzewo wyprowadzenia*.

Def. 38. *Drzewo wyprowadzenia*, to sposób ilustrowania jak dane słowo może być wyprowadzone w danej gramatyce, w postaci drzewa. Drzewo wyprowadzenia musi spełniać następujące warunki:

- W korzeniu drzewa znajduje się aksjomat.
- W węzłach wewnętrznych drzewa znajdują się nieterminale, a w liściach terminale lub słowa puste ε .
- Jeśli w danym węźle mamy nieterminał X , a w jego kolejnych synach mamy x_1, x_2, \dots, x_k , to musi zachodzić $X \rightarrow x_1x_2 \dots x_k$.
- Terminale umieszczone w liściach, czytane od lewej do prawej, tworzą wyprowadzone słowo.

□

Przyjrzyjmy się zdefiniowanym powyżej pojęciom na przykładach:

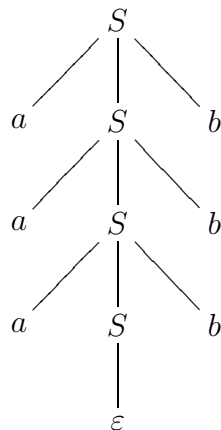
Przykład: Weźmy gramatykę generującą język $\{a^n b^n : n \geq 0\}$:

$$A \rightarrow aAb \mid \varepsilon$$

Wyprowadzenie słów $aaabbb$ ma postać:

$$A \rightarrow aAb \rightarrow aaAbb \rightarrow aaaAbbb \rightarrow aaabbb$$

a jego drzewo wyprowadzenia wygląda następująco:

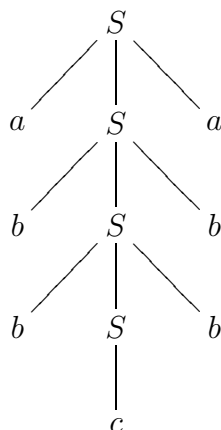


Przykład: Oto gramatyka generująca język palindromów nad alfabetem $\{a, b, c\}$:

$$S \rightarrow aSa \mid bSb \mid cSc \mid a \mid b \mid c \mid \varepsilon$$

Weźmy słowo *abcbba*. Oto jego wyprowadzenie i drzewo wyprowadzenia:

$$S \rightarrow aSa \rightarrow abSba \rightarrow abbSbba \rightarrow abcbba$$



Przykład: Przyjrzyjmy się językowi złożonemu z poprawnych wyrażeń nawiasowych (dla czytelności zbudowanych z nawiasów kwadratowych). Język ten możemy zdefiniować na dwa sposoby. Po pierwsze, wyrażenia nawiasowe, to takie ciągi nawiasów, w których:

- łączna liczba nawiasów otwierających i zamykających jest taka sama,
- każdy prefiks wyrażenia nawiasowego zawiera przynajmniej tyle nawiasów otwierających, co zamykających.

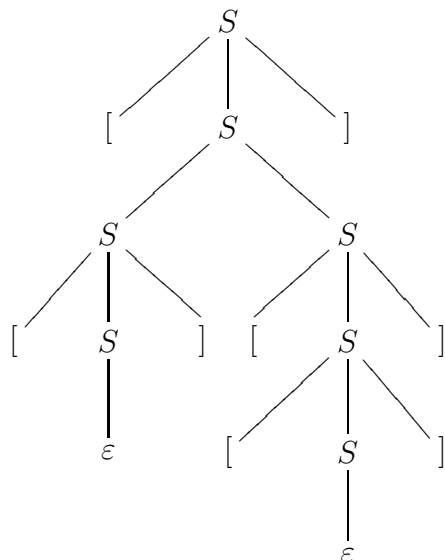
Jeżeli zdefiniujemy funkcje $L(x) = \#_{[(x)}$ i $R(x) = \#_{]}(x)$, to język wyrażeń nawiasowych możemy zdefiniować następująco:

$$W = \{x \in \{[,]\}^* : L(x) = R(x) \text{ i dla każdego prefiksu } y \text{ słowa } x \text{ mamy } L(y) \geq R(y)\}$$

Drugi sposób zdefiniowania języka wyrażeń nawiasowych, to podanie gramatyki bezkontekstowej.

$$S \rightarrow [S] \mid SS \mid \varepsilon$$

Oznaczmy tę gramatykę przez G . Oto przykładowe drzewo wyprowadzenia dla wyrażenia nawiasowego $[[[]]]$:



Pokażemy, że obie definicje są sobie równoważne. W tym celu pokażemy najpierw, że każde słowo, które można wyprowadzić w G należy do W , a następnie, że każde słowo z W można wyprowadzić w G .

To, że $L(G) \subseteq W$ będziemy dowodzić przez indukcję ze względu na długość wyprowadzenia.

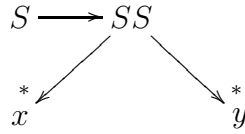
1. Jedynym słowem jakie można wyprowadzić w jednym kroku jest ε . Oczywiście $\varepsilon \in W$.
2. Załóżmy, że wyprowadzenie ma postać $S \rightarrow [S] \rightarrow^* [x]$. Z założenia indukcyjnego x jest wyrażeniem nawiasowym. W takim razie $[x]$ też, gdyż:

$$L([x]) = L(x) + 1 = R(x) + 1 = R([x])$$

oraz dla każdego prefiksu y słowa x mamy:

$$L([y]) = L(y) + 1 > L(y) \geq R(y) = R([y])$$

3. Załóżmy, że wyprowadzenie ma postać:



Z założenia indukcyjnego x i y są wyrażeniami nawiasowymi. W takim razie xy też, gdyż:

$$L(xy) = L(x) + L(y) = R(x) + R(y) = R(xy)$$

oraz dla każdego prefiksu z słowa y mamy:

$$L(xz) = L(x) + L(z) \geq L(x) + R(z) = R(x) + R(z) = R(xz)$$

Z zasady indukcji otrzymujemy $L(G) \subseteq W$.

Pokażemy teraz, że $W \subseteq L(G)$. Dowód będzie przebiegał przez indukcję po długości słowa.

1. Mamy $\varepsilon \in W$ oraz $S \rightarrow \varepsilon$.
2. Niech $x \in W$, $|x| > 0$, $x = a_1 a_2 \dots a_k$. Określmy funkcję $f : \{0, 1, \dots, k\} \rightarrow \mathcal{N}$:

$$f(i) = L(a_1 \dots a_i) - R(a_1 \dots a_i)$$

Zauważmy, że z tego, że x jest wyrażeniem nawiasowym wynika, że:

$$f(0) = f(k) = 0$$

$$f(i) \geq 0$$

Mamy dwa możliwe przypadki:

- Dla pewnego $0 < i < k$ mamy $f(i) = 0$. Wówczas $a_1 \dots a_i$ oraz $a_{i+1} \dots a_k$ są wyrażeniami nawiasowymi. Stąd:

$$S \rightarrow SS \xrightarrow{*} a_1 \dots a_i a_{i+1} \dots a_k = x$$

- Dla wszystkich $0 < i < k$ mamy $f(i) > 0$. Wówczas $a_2 \dots a_{k-1}$ jest wyrażeniem nawiasowym. Stąd:

$$S \rightarrow [S] \xrightarrow{*} [a_2 \dots a_{k-1}] = x$$

Tak więc $x \in L(G)$.

Na mocy zasady indukcji uzyskujemy $W \subseteq L(G)$. Tak więc $W = L(G)$.

Często chcemy, aby gramatyka nie tylko opisywała język, ale również jego strukturę. Wymagamy wówczas, aby gramatyka była dodatkowo *jednoznaczna*.

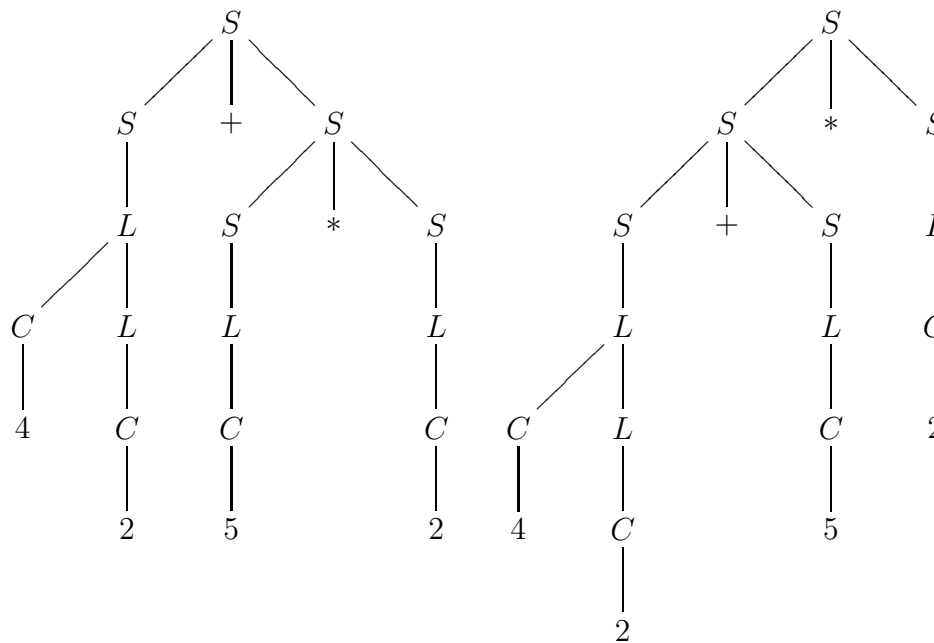
Def. 39. Powiemy, że gramatyka G jest *jednoznaczna*, gdy dla każdego słowa $x \in L(G)$ istnieje tylko jedno drzewo wyprowadzenia słowa x w gramatyce G . \square

Jeśli gramatyka jest jednoznaczna, to drzewa wyprowadzeń jednoznacznie określają strukturę (składniową) słów z języka. Nadal oczywiście słowa mogą mieć wiele wyprowadzeń, ale różnią się one tylko kolejnością stosowanych produkcji, a nie strukturą wyprowadzania słowa.

Przykład: Oto gramatyka opisująca proste wyrażenia arytmetyczne:

$$\begin{aligned} S &\rightarrow S + S \mid S - S \mid S * S \mid S / S \mid L \mid (S) \\ L &\rightarrow C \mid CL \\ C &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

Gramatyka ta jest niestety niejednoznaczna. Na przykład, wyrażenie $42 + 5 * 2$ ma dwa różne drzewa wyprowadzenia:



Struktura wyrażenia ma wpływ na sposób jego interpretacji. Drzewo wyprowadzenia po lewej prowadzi do wartości wyrażenia 52, a to po prawej do 94. Wszędzie tam, gdzie gramatyka opisuje składnię języka, któremu towarzyszy ściśle zdefiniowana semantyka, będziemy chcieli, aby była to jednoznaczna gramatyka.

W przypadku wyrażen arytmetycznych musimy ustalić kolejność wiązania operatorów, oraz wymusić, aby operatory były albo lewostronnie, albo prawostronnie łączne. Oto

jednoznaczna gramatyka wyrażeń arytmetycznych:

$$E \rightarrow E + S \mid E - S \mid S$$

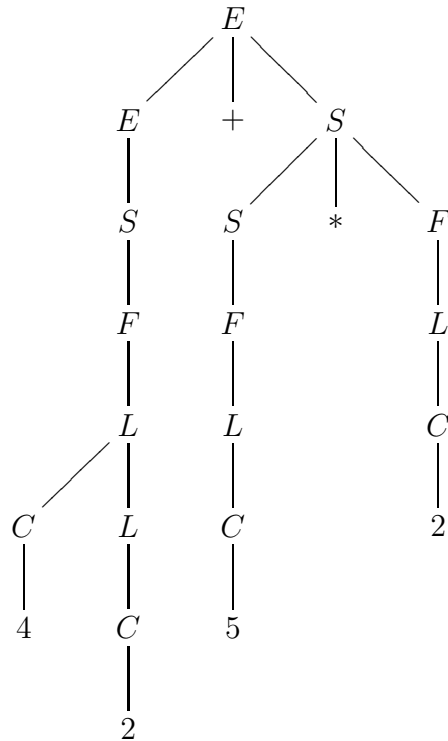
$$S \rightarrow S * F \mid S / F \mid F$$

$$F \rightarrow L \mid (E)$$

$$L \rightarrow C \mid CL$$

$$C \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

W tej gramatyce analizowane wyrażenie ma tylko jedno drzewo wyprowadzenia:



10.4 Języki regularne a bezkontekstowe

Pokażemy teraz, że wszystkie regularne są bezkontekstowe. Jak już wcześniej widzieliśmy, odwrotne stwierdzenie nie jest prawdziwe, gdyż $a^n b^n$ jest językiem bezkontekstowym, ale nie regularny.

Def. 40. Gramatyka (prawostronnie) liniowa, to taka, w której wszystkie produkcje mają postać: $A \rightarrow \alpha B$ lub $A \rightarrow \alpha$, dla $A, B \in N$ i $\alpha \in \Sigma^*$.

Gramatyka silnie (prawostronnie) liniowa, to taka, w której wszystkie produkcje mają postać: $A \rightarrow aB$ lub $A \rightarrow \varepsilon$, dla $A, B \in N$ i $a \in \Sigma$. \square

Okazuje się, że gramatyki liniowe i silnie liniowe opisują dokładnie klasę języków regularnych.

Fakt 12. Gramatyki (silnie) liniowe opisują dokładnie klasę języków regularnych.

Dowód: Pokażemy, jak gramatykę liniową przerobić na równoważną jej gramatykę silnie liniową. Następnie pokażemy jak można gramatykę silnie liniową przekształcić na równoważny jej niedeterministyczny automat skończony i vice versa.

Żeby gramatykę liniową przerobić na liniową, dla każdej produkcji postaci (dla $k > 1$):

$$A \rightarrow a_1 \dots a_k B$$

dodajemy do gramatyki nowe terminale A_1, \dots, A_{k-1} , a samą produkcję zastępujemy produkcjami:

$$\begin{aligned} A &\rightarrow a_1 A_1 \\ A_1 &\rightarrow a_2 A_2 \\ &\vdots \\ A_{k-1} &\rightarrow a_k B \end{aligned}$$

Podobnie, dla każdej produkcji postaci (dla $k > 1$):

$$A \rightarrow a_1 \dots a_k$$

dodajemy do gramatyki nowe terminale A_1, \dots, A_k , a samą produkcję zastępujemy produkcjami:

$$\begin{aligned} A &\rightarrow a_1 A_1 \\ A_1 &\rightarrow a_2 A_2 \\ &\vdots \\ A_{k-1} &\rightarrow a_k A_k \\ A_k &\rightarrow \varepsilon \end{aligned}$$

1. Gramatykę liniową można sprowadzić do silnie liniowej.
2. Gramatykę silnie liniową można przerobić na automat niedeterministyczny.
3. Automat niedeterministyczny można przerobić na gramatykę silnie liniową.

W rezultacie uzyskujemy gramatykę akceptującą ten sam język, ale silnie liniową.

Zauważmy, że jeśli gramatyka jest silnie liniowa (lub liniowa), to w wyprowadzeniu pojawia się na raz tylko jeden nieterminal, i to zawsze na samym końcu słowa. Konstruując niedeterministyczny automat skończony równoważny danej gramatyce silnie liniowej opieramy się na następującej analogii: Miec $G = \langle N, \Sigma, P, S \rangle$ będzie daną gramatyką silnie liniową. Stanami automatu będą nieterminale gramatyki N . Dla każdego słowa $x \in \Sigma^*$ i $A \in N$, takich, że $S \xrightarrow{*} xA$ nasz automat po wczytaniu słowa x będzie mógł przejść do stanu A . Formalnie konstrukcja naszego automatu M wygląda następująco $M = (N, \Sigma, \delta, S, F)$, gdzie:

$$F = \{A : A \rightarrow \varepsilon \in P\}$$

oraz:

$$\delta(A, a) = \{B \in N : A \rightarrow aB \in P\}$$

Można pokazać (przez indukcję po długości słowa), że:

$$\delta^*(A, x) = \{B \in N : A \rightarrow^* xB\}$$

Stąd:

$$\begin{aligned} L(M) &= \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\} = \\ &= \{x \in \Sigma^* : \exists A \in N A \in \delta^*(S, x) \wedge A \in F\} = \\ &= \{x \in \Sigma^* : \exists A \in N S \rightarrow^* xA \wedge A \rightarrow \varepsilon \in P\} = \\ &= \{x \in \Sigma^* : S \rightarrow^* x\} = L(G) \end{aligned}$$

Tak więc faktycznie języki generowane przez gramatyki silnie liniowe są regularne.

Co więcej, powyższą konstrukcją można odwrócić. Jeśli $M = (Q, \Sigma, \delta, s, F)$ jest danym niedeterministycznym automatem liniowym, to równoważna mu gramatyka silnie liniowa ma postać $G = \langle Q, \Sigma, P, s \rangle$, gdzie P zawiera produkcje postaci:

$$\begin{aligned} q &\rightarrow \varepsilon && \text{dla } q \in F \\ q &\rightarrow ap && \text{dla } p, q \in Q, a \in \Sigma, p \in \delta(q, a) \end{aligned}$$

Ponownie można pokazać (przez indukcję po długości słowa), że:

$$\delta^*(q, x) = \{p \in Q : q \rightarrow^* xp\}$$

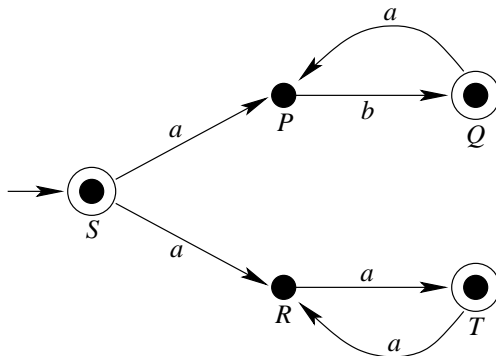
Stąd:

$$\begin{aligned} L(G) &= \{x \in \Sigma^* : s \rightarrow^* x\} = \\ &= \{x \in \Sigma^* : \exists q \in Q s \rightarrow^* xq \wedge q \rightarrow \varepsilon \in P\} = \\ &= \{x \in \Sigma^* : \exists q \in Q q \in \delta^*(s, x) \wedge q \in F\} = \\ &= \{x \in \Sigma^* : \delta^*(S, x) \cap F \neq \emptyset\} = L(M) \end{aligned}$$

□

Tak więc dla wszystkich języków regularnych istnieją generujące je gramatyki silnie liniowe. Tym samym pokazaliśmy, że klasa języków regularnych zawiera się ściśle w klasie języków bezkontekstowych.

Przykład: Rozważmy automat skończony akceptujący język $(ab)^* \mid (aa)^*$:



Jest on równoważny gramatyce silnie liniowej postaci:

$$\begin{aligned} S &\rightarrow aP \mid aR \mid \varepsilon \\ P &\rightarrow bQ \\ Q &\rightarrow aP \mid \varepsilon \\ R &\rightarrow aT \\ T &\rightarrow aR \mid \varepsilon \end{aligned}$$

Oto wyprowadzenie przykładowego słowa $abab$:

$$S \rightarrow aP \rightarrow abQ \rightarrow abaP \rightarrow ababQ \rightarrow abab$$

Jak widać, odpowiada ono dokładnie obliczeniu automatu akceptującemu słowo $abab$.

Jeśli wystarczy nam, aby gramatyka była tylko liniowa, to wystarczy mniejsza liczba nieterminali:

$$\begin{aligned} S &\rightarrow abQ \mid aaT \mid \varepsilon \\ Q &\rightarrow abQ \mid \varepsilon \\ T &\rightarrow aaT \mid \varepsilon \end{aligned}$$

A oto wyprowadzenie przykładowego słowa $abab$:

$$S \rightarrow abQ \rightarrow ababQ \rightarrow abab$$

10.5 Podsumowanie

W tym wykładzie poznaliśmy języki i gramatyki bezkontekstowe wraz z podstawowymi związanymi z nimi pojęciami. Pokazaliśmy też, że klasa języków bezkontekstowych zawiera ściśle w sobie klasę języków regularnych.

10.6 Skorowidz

- **BNF** — notacja Backusa-Naura, formalizm służący do opisu składni, oparty na gramatykach bezkontekstowych. Zobacz też artykuł w Wikipedii.
- **Drzewo wyprowadzenia** — drzewo potwierdzające, że dane słowo należy do języka generowanego przez daną gramatykę; odzwierciedla również strukturę składniową słowa.
- **Gramatyka bezkontekstowa** — formalizm do opisu języków.
- **Gramatyka liniowa** — gramatyka bezkontekstowa, w której występują produkcje jedynie postaci $A \rightarrow \alpha B$ lub $A \rightarrow \alpha$ (dla $A, B \in N$ i $\alpha \in \Sigma^*$). Gramatyki liniowe opisują języki regularne.

- **Gramatyka silnie liniowa** — gramatyka bezkontekstowa, w której występują produkcje jedynie postaci $A \rightarrow aB$ lub $A \rightarrow \varepsilon$ (dla $A, B \in N$ i $a \in \Sigma$). Gramatyki silnie liniowe opisują języki regularne.
- **Język generowany przez gramatykę bezkontekstową** — to język złożony z tych słów $x \in \Sigma^*$, które można wyprowadzić z aksjomatu $S, S \rightarrow^* x$.
- **Wyprowadzenie** — ciąg słów potwierdzający, że dane słowo należy do języka generowanego przez daną gramatykę.

10.7 Praca domowa

1. Podaj gramatykę bezkontekstową generującą język $\{a^n b^{n+k} c^k : n, k \geq 1\}$. Podaj drzewo wyprowadzenia dla słowa $aabbbbcccc$.
2. Podaj liniową gramatykę generującą język $ab(b^* | a^*)$.

10.8 Ćwiczenia

1. Dla gramatyki:

$$\begin{aligned} S &\rightarrow aB|Ab|SS \\ A &\rightarrow a|aS \\ B &\rightarrow b|Sb \end{aligned}$$

i słowa $aabbab$:

- podaj jego wyprowadzenie,
 - podaj jego drzewo wyprowadzenia,
 - czy jest to gramatyka jednoznaczna (podaj kontrprzykład, lub krótko uzasadnij)?
2. Podaj gramatykę bezkontekstową generującą język: $\{a^i b^{i+2 \cdot j} c^j\}$. Narysuj drzewo wyprowadzenia dla słowa $aabbbbbbbcccc$.
 3. Dla podanej poniżej gramatyki, podaj wyprowadzenie i drzewo wyprowadzenia słowa $baabab$.

$$\begin{aligned} S &\rightarrow AB | BA \\ A &\rightarrow a | BBA \\ B &\rightarrow b | AAB \end{aligned}$$

Czy podana gramatyka jest jednoznaczna (odpowiedź uzasadnij lub podaj kontrprzykład)?

4. Dla podanej poniżej gramatyki, podaj wyprowadzenie i drzewo wyprowadzenia słowa *abbabbaa*.

$$\begin{aligned} S &\rightarrow A \mid B \mid AB \\ A &\rightarrow BA \mid aaA \mid abaA \mid aa \\ B &\rightarrow abbB \mid \varepsilon \end{aligned}$$

Czy podana gramatyka jest jednoznaczna (odpowiedź uzasadnij lub podaj kontrprzykład)?

5. Dla podanej poniżej gramatyki, podaj wyprowadzenie i drzewo wyprowadzenia słowa *abaabb*.

$$\begin{aligned} S &\rightarrow aSB \mid ASb \mid ab \\ A &\rightarrow Sa \mid bAA \mid a \\ B &\rightarrow bS \mid BBa \mid b \end{aligned}$$

Czy podana gramatyka jest jednoznaczna (odpowiedź uzasadnij lub podaj kontrprzykład)?

6. Podaj gramatykę bezkontekstową generującą język:

- $\{a^n b^{2n} : n \geq 1\}$,
- $\{a^k b^n : k \geq 2n\}$,
- $\{a^n b^k : 2n \geq 3k\}$.
- $\{a^i b^j c^i\}$,
- $\{a^i b^j a^j b^i\}$,
- $\{a^i b^i a^j b^j\}$,
- $\{a^i b^j c^k : i + 2 \cdot j = k\}$,
- $\{a^i b^j c^k : 3i = 2j + k\}$,
- $\{a^n b^m : n \geq 2m \geq 0\}$,
- $\{a^n b^m c^k d^l : n + m = k + l\}$,
- $\{a^n b^m c^k d^l : n + l = k + m\}$,
- $\{a^n b^m c^k d^l : n + k = m + l\}$,
- $\{xc \operatorname{rev}(x) : x \in \{a, b\}^*\}$,
- $\{xycy : x, y \in \{a, b\}^* \wedge \#_a(x) = \#_b(y)\}$,
- $\{a^n b^m c^k : n + k = 2 * m\}$,

7. Rozszerz podaną w treści wykładu gramatykę wyrażeń arytmetycznych o unarny minus. Zapewnij aby była jednoznaczna. Jak silnie powinien wiązać unarny minus?

8. Podać gramatykę generującą język złożony z takich wyrażeń arytmetycznych składających się z symboli $+$, $*$, 0 , 1 , $($, $)$, których wartość jest większa niż 2.
9. Skonstruuj gramatykę bezkontekstową generującą język złożony ze wszystkich słów nad alfabetem $\{a, b\}^*$, które nie są palindromami.
10. Dla zadanego automatu/wyrażenia regularnego podaj gramatyki (prawostronnie) liniowe / silnie liniowe opisujące odpowiedni język:

- $a(a | b)^*aa$,
- $(a(ab | ba)^* | (b(bb | aa)^*))$,

•

	0	1
$\rightarrow 1$	1	2
$F2$	1	3
$F3$	3	1

•

	a	b
$\rightarrow 1$	2	-
$F2$	2	1

11. Podaj gramatykę wyrażeń regularnych (nad alfabetem $\{a, b\}$).

Wykład 11. Postać normalna Chomsky'ego, algorytm Cocke-Younger'a-Kasami

11.1 Wstęp

W tym wykładzie zajmiemy się problemem jak efektywnie rozpoznawać, czy dane słowo należy do języka generowanego przez daną gramatykę. W tym celu poznamy najpierw szczególną postać gramatyk bezkontekstowych, tzw. postać normalną Chomsky'ego i dowiemy się jak przekształcać gramatyki do tej postaci. Następnie poznamy dynamiczny algorytm rozpoznający czy dane słowo można wyprowadzić w danej gramatyce w postaci normalnej Chomsky'ego.

11.2 Postać normalna Chomsky'ego

Def. 41. Powiemy, że gramatyka jest w *postaci normalnej Chomsky'ego*, jeśli wszystkie produkcje w gramatyce są postaci: $A \rightarrow BC$ lub $A \rightarrow a$ (dla pewnych nieterminali A, B i C oraz terminala a). \square

Zauważmy, że jeśli gramatyka jest w postaci normalnej Chomsky'ego, to każda produkcja przekształca jeden nieterminal na dwa nieterminale, lub jeden nieterminal na jeden nieterminal. Ułatwia to rozstrzygnięcie, czy dane słowo x da się wyprowadzić w danej gramatyce — wiadomo, że w takim wyprowadzeniu $|x| - 1$ razy należy zastosować produkcje przekształcające nieterminal na dwa nieterminale i $|x|$ razy zastosować produkcje przekształcające nieterminal na terminal.

Podstawowe pytanie, jakie można sobie zadać brzmi: czy dla każdego języka bezkontekstowego istnieje generująca go gramatyka w postaci normalnej Chomsky'ego? Odpowiedź brzmi: nie. Zauważmy, że zastosowanie dowolnej produkcji nie zmniejsza liczby symboli. Ponieważ zaczynamy od słowa złożonego z jednego symbolu (aksjomatu), więc nigdy nie uzyskamy słowa pustego. Jeśli więc język A jest bezkontekstowy i $\varepsilon \in A$, to A nie jest generowany przez żadną gramatykę w postaci normalnej Chomsky'ego.

Na szczęście słowo puste jest jedynym elementem, którego nie możemy uzyskać rozważając gramatyki w postaci normalnej Chomsky'ego. Jeżeli dopuścimy, aby słowo puste zniknęło z języka generowanego przez gramatykę, to każdą gramatykę można przekształcić do postaci normalnej Chomsky'ego.

Fakt 13. Dla dowolnej gramatyki bezkontekstowej G istnieje taka gramatyka bezkontekstowa w postaci normalnej Chomsky'ego G' , że $L(G') = L(G) \setminus \{\varepsilon\}$.

Dowód: Niech $G = \langle N, \Sigma, P, S \rangle$. Konstrukcja gramatyki G' przebiega w kilku krokach:

1. Niech gramatyka $G_1 = \langle N, \Sigma, P_1, S \rangle$ będzie gramatyką powstałą z G przez domknięcie zbioru produkcji zgodnie z następującymi dwiema regułami:

- Jeśli mamy produkcje postaci $X \rightarrow \alpha Y \beta$ i $Y \rightarrow \varepsilon$ (dla $X, Y \in N$ i $\alpha, \beta \in (N \cup \Sigma)^*$), to dodajemy produkcję $X \rightarrow \alpha \beta$. Tak dodana produkcja realizuje w jednym kroku to, co można było zrobić w dwóch krokach: $X \rightarrow \alpha Y \beta \rightarrow \alpha \beta$.
- Jeśli mamy produkcje postaci $X \rightarrow Y$ i $Y \rightarrow \gamma$ (dla $X, Y \in N$ i $\gamma \in (N \cup \Sigma)^*$), to dodajemy produkcję $X \rightarrow \gamma$. Tak dodana produkcja realizuje w jednym kroku to, co można było zrobić w dwóch krokach: $X \rightarrow Y \rightarrow \gamma$.

Stosowanie powyższych reguł powtarzamy tak długo, jak długo wprowadzają one jakiegokolwiek nowe produkcje. Ponieważ dodawane produkcje nie są dłuższe od już istniejących, więc cały proces doda tylko ograniczoną liczbę produkcji i w pewnym momencie się zakończy.

Zauważmy, że $L(G_1) = L(G)$, gdyż dodawane produkcje nie pozwalają wyprowadzić niczego, czego i tak nie dałoby się wcześniej wyprowadzić.

2. Dla dowolnego słowa $x \in L(G_1)$, $x \neq \varepsilon$ rozważymy jego najkrótsze wyprowadzenie. W takim wyprowadzeniu nie będą się pojawiały produkcje postaci $X \rightarrow \varepsilon$ ani $X \rightarrow Y$. Niech G_2 będzie gramatyką powstałą z G_1 przez usunięcie tych produkcji, $G_2 = \langle N, \Sigma, P_2, S \rangle$, gdzie

$$P_2 = P_1 \setminus \{X \rightarrow \varepsilon, X \rightarrow Y : X, Y \in N\}$$

Wówczas:

$$L(G_2) = L(G_1) \setminus \{\varepsilon\} = L(G) \setminus \{\varepsilon\}$$

3. Gramatyka G_3 powstaje z G_2 przez dodanie do niej, dla każdego terminala $a, b, c, \dots \in N$, odpowiadającego mu nowego nieterminala (A, B, C, \dots) , oraz produkcji: $A \rightarrow a$, $B \rightarrow b$, $C \rightarrow c, \dots$. We wszystkich pozostałych produkcjach $X \rightarrow \alpha$, gdzie $|\alpha| > 1$, zamieniamy wszystkie terminale na odpowiadające im nieterminale. Nie zmienia to języka generowanego przez gramatykę:

$$L(G_3) = L(G_2) = L(G) \setminus \{\varepsilon\}$$

Natomiast w ten sposób w G_3 mamy tylko dwa rodzaje produkcji: $X \rightarrow a$ gdzie a jest terminalem, oraz $X \rightarrow Y_1 Y_2 \dots Y_k$ gdzie Y_i to nieterminale i $k > 1$.

4. Gramatyka G_4 powstaje z G_3 przez zastąpienie każdej produkcji postaci $X \rightarrow Y_1 Y_2 \dots Y_k$ (dla $k > 2$) przez zestaw produkcji:

$$\begin{aligned} X &\rightarrow Y_1 X_1 \\ X_1 &\rightarrow Y_2 X_2 \\ &\vdots \\ X_{k-2} &\rightarrow Y_{k-1} Y_k \end{aligned}$$

gdzie X_1, X_2, \dots, X_{k-2} to nowe nieterminale dodane do G_4 . Tak uzyskana gramatyka generuje dokładnie ten sam język — jedyne co się zmieniło to to, że zamiast zastosowania jednej produkcji, trzeba ich zastosować kilka,

$$L(G_4) = L(G_3) = L(G) \setminus \{\varepsilon\}$$

Natomiast G_4 jest już w postaci normalnej Chomsky'ego. □

Przykład: Zobaczmy, jak powyższy algorytm działa na przykładzie gramatyki generującej język $\{a^n b^n : n \geq 0\}$.

$$S \rightarrow aSb \mid \varepsilon$$

Pierwszy krok powoduje dodanie produkcji $S \rightarrow ab$:

$$S \rightarrow aSb \mid ab \mid \varepsilon$$

W drugim kroku usuwamy produkcję $S \rightarrow \varepsilon$.

$$S \rightarrow aSb \mid ab$$

W trzecim kroku dodajemy nieterminale A i B , zastępujemy nimi w istniejących produkcjach terminale a i b , oraz dodajemy produkcje $A \rightarrow a$ i $B \rightarrow b$.

$$S \rightarrow ASB \mid AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

W ostatnim kroku algorytmu rozbijamy produkcję $S \rightarrow ASB$ na dwie produkcje.

$$S \rightarrow AX \mid AB$$

$$X \rightarrow SB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Tak uzyskana gramatyka jest w postaci normalnej Chomsky'ego. Na przykład, wyprowadzenie słowa $aaabbb$ ma postać:

$$\begin{aligned} S &\rightarrow AX \rightarrow ASB \rightarrow aSB \rightarrow aSb \rightarrow \\ &\rightarrow aAXb \rightarrow aASBb \rightarrow aaSBb \rightarrow aaSbb \rightarrow \\ &\rightarrow baABbb \rightarrow aaaBbb \rightarrow aaabbb \end{aligned}$$

Przykład: Zobaczmy, jak powyższy algorytm działa na przykładzie gramatyki generującej wyrażenia nawiasowe.

$$S \rightarrow [S] \mid SS \mid \varepsilon$$

Pierwszy krok powoduje dodanie produkcji: $S \rightarrow []$ i $S \rightarrow S$:

$$S \rightarrow [S] \mid [] \mid SS \mid S \mid \varepsilon$$

W drugim kroku usuwamy produkcje $S \rightarrow \varepsilon$ i $S \rightarrow S$.

$$S \rightarrow [S] \mid [] \mid SS$$

W trzecim kroku dodajemy nieterminale L i P , zastępujemy nimi w istniejących produkcjach odpowiednio terminale $[$ i $]$, oraz dodajemy produkcje $L \rightarrow [$ i $P \rightarrow]$.

$$\begin{aligned} S &\rightarrow LSP \mid LP \mid SS \\ L &\rightarrow [\\ P &\rightarrow] \end{aligned}$$

W ostatnim kroku algorytmu rozbijamy produkcję $S \rightarrow LSP$ na dwie produkcje.

$$\begin{aligned} S &\rightarrow LX \mid LP \mid SS \\ X &\rightarrow SP \\ L &\rightarrow [\\ P &\rightarrow] \end{aligned}$$

Tak uzyskana gramatyka jest w postaci normalnej Chomsky'ego. Na przykład, wyprowadzenie słowa $[[[]]]$ w tej gramatyce ma postać:

$$\begin{aligned} S &\rightarrow LX \rightarrow LSP \rightarrow [SP \rightarrow [S] \rightarrow \\ &\rightarrow [SS] \rightarrow [LPS] \rightarrow [[PS] \rightarrow [[]S] \rightarrow \\ &\rightarrow [[]LX] \rightarrow [[]LSP] \rightarrow [[][SP] \rightarrow [[][S]] \rightarrow \\ &\rightarrow [[][LP]] \rightarrow [[][[P]] \rightarrow [[][[[]]] \end{aligned}$$

11.3 Algorytm rozpoznawania Cocke-Younger'a-Kasami (CYK)

Algorytm CYK pozwala stwierdzić, czy dane słowo jest wyprowadzalne w danej gramatyce (w postaci normalnej Chomsky'ego). Jeśli gramatyka nie jest w postaci normalnej Chomsky'ego, to można ją sprowadzić do tej postaci zgodnie z opisaną wyżej konstrukcją.

Tak jak to zwykle bywa w programowaniu dynamicznym, żeby rozwiązać dany problem musimy go najpierw uogólnić i zamiast szukać odpowiedzi na jedno pytanie, znajdziemy odpowiedzi na całe mnóstwo powiązanych ze sobą pytań.

Niech $G = \langle N, \Sigma, P, S \rangle$ będzie daną gramatyką w postaci normalnej Chomsky'ego, a dane słowo $x \in \Sigma^*$ będzie postaci $x = a_1 a_2 \dots a_n$. Dla każdej pary indeksów $1 \leq i \leq j \leq n$ wyznaczmy wszystkie takie nieterminale X , z których można wyprowadzić pod słowo $a_i \dots a_j$. Dokładniej, wypełnimy tablicę:

$$T[i, j] = \{X \in N : X \rightarrow^* a_i \dots a_j\}$$

Gdy już będziemy mieli obliczoną tablicę T , to słowo x można wyprowadzić w gramatyce G wtw., gdy $S \in T[1, n]$.

W jaki sposób możemy wypełniać tę tablicę?

- Zaczynamy od przekątnej. $T[i, i]$ to zbiór nieterminali, z których można wyprowadzić a_i . Jednak takie wyprowadzenie może składać się tylko z zastosowania jednej produkcji zastępującej nieterminal terminalem:

$$T[i, i] = \{X : X \rightarrow a_i \in P\}$$

- Następnie wypełniamy kolejne przy-przekątne. Zauważmy, że jeżeli $X \rightarrow^* a_i \dots a_j$, $i < j$, to wyprowadzenie to musi mieć postać: $X \rightarrow YZ$, $Y \rightarrow^* a_i \dots a_k$, $Z \rightarrow^* a_{k+1} \dots a_j$, dla pewnego k , $i \leq k < j$.

Pomocna nam będzie następująca operacja na zbiorach nieterminali, która na podstawie zbioru wszystkich możliwych Y -ów i Z -ów wyznacza zbiór wszystkich możliwych X -ów:

$$\mathcal{A} \otimes \mathcal{B} = \{X : X \rightarrow YZ \in P, Y \in \mathcal{A}, Z \in \mathcal{B}\}$$

Operacja \otimes może być skomplikowana, ale zauważmy, że jej koszt jest stały (tzn. zależy od gramatyki, ale nie od długości słowa x).

Obliczając $T[i, j]$ uwzględniamy wszystkie możliwe punkty podziału k pod słowa $a_i \dots a_j$:

$$T[i, j] = \bigcup_{i \leq k < j} T[i, k] \otimes T[k + 1, j]$$

Tablica ma rozmiar $\Theta(n^2)$ przy czym aby wypełnić jeden element tablicy trzeba wykonać liniową liczbę operacji \otimes , tak więc złożoność czasowa tego algorytmu to $\Theta(n^3)$.

Przykład: Rozważmy gramatykę w postaci normalnej Chomsky'ego:

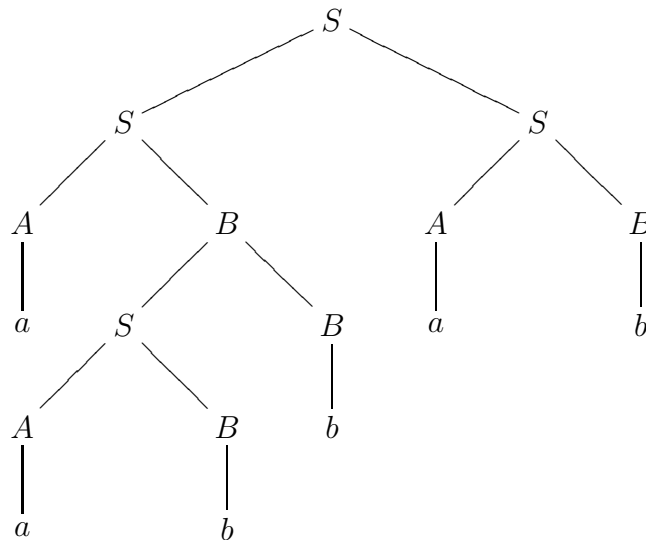
$$\begin{aligned} S &\rightarrow SS \mid AB \\ A &\rightarrow AS \mid AA \mid a \\ B &\rightarrow SB \mid BB \mid b \end{aligned}$$

oraz słowo $aabbab$. Algorytm CYK da nam w wyniku następującą tablicę T :

	1	2	3	4	5	6
1	A	A	A, S	A, B, S	A	A, S
2		A	S	B, S	–	S
3			B	B	–	–
4				B	–	–
5					A	S
6						B

Jeżeli spróbujemy odtworzyć skąd w $T[1, 6]$ wzięło się S , to uzyskamy informacje o drzewach wyprowadzenia badanego słowa. Poniżej podano przykładowe drzewo wyprowadzenia i zaznaczono odpowiadające mu nieterminale w tablicy T .

	1	2	3	4	5	6
1	A	A	A, S	A, B, S	A	A, S
2		A	S	B, S	–	S
3			B	B	–	–
4				B	–	–
5					A	S
6						B



11.4 Podsumowanie

W tym wykładzie poznaliśmy postać normalną Chomsky'ego gramatyk bezkontekstowych i pokazaliśmy, że każdą gramatykę bezkontekstową można sprowadzić do tej postaci (z dokładnością do akceptowania słowa pustego ε). Poznaliśmy też algorytm CKY rozpoznający, czy dane słowo należy do języka generowanego przez daną gramatykę bezkontekstową w postaci normalnej Chomsky'ego.

11.5 Skorowidz

- **Postać normalna Chomsky'ego** — taka postać gramatyk bezkontekstowych, w której gramatyka może zawierać tylko produkcje postaci: $A \rightarrow BC$ lub $A \rightarrow a$ (dla pewnych nieterminali A, B i C oraz terminala a). Zobacz też artykuł w Wikipedii.
- **Algorytm CYK** — algorytm rozpoznawania czy dane słowo jest wyprowadzalne w danej gramatyce w postaci normalnej Chomsky'ego. Zobacz też artykuł w Wikipedii.

11.6 Praca domowa

1. Sprowadź do postaci normalnej Chomsky'ego następującą gramatykę bezkontekstową:

$$\begin{aligned} S &\rightarrow ab \mid ba \mid T \\ T &\rightarrow SS \mid aSb \mid bSa \end{aligned}$$

2. Dla podanej gramatyki bezkontekstowych w postaci normalnej Chomsky'ego i słowa $abaabb$ zasymuluj działanie algorytmu CYK. Podaj zawartość tablicy T oraz przykładowe drzewo wyprowadzenia.

$$\begin{aligned} S &\rightarrow SS \mid AB \\ A &\rightarrow AA \mid a \\ B &\rightarrow BB \mid b \end{aligned}$$

11.7 Ćwiczenia

1. Sprowadź do postaci normalnej Chomsky'ego następujące gramatyki bezkontekstowe:

•

$$\begin{aligned} S &\rightarrow aSb \mid cS \mid Sc \mid T \\ T &\rightarrow ab \mid c \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow aSc \mid B \\ B &\rightarrow bBc \mid \varepsilon \end{aligned}$$

•

$$\begin{aligned} X &\rightarrow aXa \mid Y \\ Y &\rightarrow Yc \mid bX \mid a \end{aligned}$$

•

$$\begin{aligned} X &\rightarrow aXYb \mid Y \\ Y &\rightarrow \varepsilon \mid cY \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow Y \mid aY \mid bbY \\ Y &\rightarrow \varepsilon \mid cS \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow yY \mid xxSy \mid YSY \mid xy \\ X &\rightarrow Y \mid xYX \mid \varepsilon \\ Y &\rightarrow X \mid yX \mid \varepsilon \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow aSbb \mid ASB \mid B \mid \varepsilon \\ A &\rightarrow ABA \mid a \mid \varepsilon \\ B &\rightarrow aB \mid bA \mid S \mid a \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow PSP \mid qqSq \mid Q \mid qq \mid \varepsilon \\ P &\rightarrow pQp \mid xY \mid p \mid q \\ Q &\rightarrow QP \mid PQ \mid S \mid \varepsilon \end{aligned}$$

•

$$\begin{aligned} S &\rightarrow mSnn \mid mNm \mid m \mid \varepsilon \\ M &\rightarrow MN \mid NM \mid N \\ N &\rightarrow mM \mid nM \mid M \mid \varepsilon \end{aligned}$$

2. Dla każdej z poniższych gramatyk bezkontekstowych w postaci normalnej Chomsky'ego zasymuluj działanie algorytmu CYK dla podanego słowa. Czy podane słowo faktycznie można wyprowadzić w podanej gramatyce? Jeśli tak, to podaj przykładowe drzewo wyprowadzenia.

• *aabbab*

$$\begin{aligned} S &\rightarrow AB \mid AA \\ A &\rightarrow SS \mid BB \mid a \\ B &\rightarrow SB \mid AS \mid b \end{aligned}$$

- *aabaaab*

$$\begin{aligned} S &\rightarrow BB \mid AA \\ A &\rightarrow AS \mid BS \mid a \\ B &\rightarrow BA \mid SB \mid b \end{aligned}$$

- *aaaaaba.*

$$\begin{aligned} S &\rightarrow SS \mid AA \\ A &\rightarrow AS \mid SB \mid a \\ B &\rightarrow AB \mid BB \mid b \end{aligned}$$

- *aabbaba*

$$\begin{aligned} S &\rightarrow SS \mid BB \\ A &\rightarrow SA \mid BS \mid a \\ B &\rightarrow BA \mid AA \mid b \end{aligned}$$

- *aacacab*

$$\begin{aligned} S &\rightarrow AC \mid CB \mid AB \\ A &\rightarrow a \mid AA \mid AC \\ B &\rightarrow b \mid BB \mid CB \\ C &\rightarrow c \mid CB \mid CC \end{aligned}$$

Wykład 12. Lemat o pompowaniu dla języków bezkontekstowych

12.1 Wstęp

Okazuje się, że nie każdy język jest bezkontekstowy. Poznamy w tym wykładzie lemat o pompowaniu dla języków bezkontekstowych. Podobnie jak w przypadku lematu o pompowaniu dla języków regularnych, jest to narzędzie umożliwiające wykazanie, że dany język **nie jest** bezkontekstowy.

12.2 Lemat o pompowaniu dla języków bezkontekstowych

Niech A będzie językiem bezkontekstowym. Jeżeli A zawiera dowolnie długie słowa, to okazuje się, że drzewach wyprowadzeń tych słów muszą istnieć fragmenty, które można powielać. Ich powielanie „pompuje” słowa należące do języka.

Tw. 4 (lemat o pompowaniu dla języków bezkontekstowych). *Niech A będzie językiem bezkontekstowym. Istnieje wówczas taka stała $k > 0$, że dla każdego słowa $z \in A$ o długości $|z| \geq k$ można to słowo podzielić na pięć części: $z = uvwxy$ w taki sposób, że: $vx \neq \varepsilon$, $|vwx| \leq k$ oraz dla każdego $i \geq 0$ mamy $uv^iwx^iy \in A$.*

Dowód: Jeśli język A jest bezkontekstowy, to istnieje opisująca go gramatyka bezkontekstowa. Zauważmy, że nie ma znaczenia, czy $\varepsilon \in A$ (gdyż $k > 0$). Załóżmy więc, że $\varepsilon \notin A$. Dzięki temu możemy założyć, że istnieje gramatyka bezkontekstowa w postaci normalnej Chomsky’ego $G = \langle N, \Sigma, P, S \rangle$ generująca język A , $L(G) = A$.

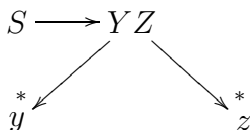
Zauważmy, że drzewa wyprowadzeń dla gramatyk w postaci normalnej Chomsky’ego mają postać regularnych drzew binarnych z „frędzelkami” długości 1. Produkcje przekształcające nieterminal na dwa nieterminale dają węzły wewnętrzne o dwóch synach. Natomiast produkcje przekształcające nieterminal na terminal odpowiadają ojcom liści — każdy ojciec liścia ma tylko jednego syna.

Zachodzi następujący fakt:

Fakt 14. *Jeśli drzewo wyprowadzenia ma wysokość h , to długość wyprowadzanego słowa nie przekracza 2^{h-1} .*

Dowód faktu: Niech x będzie wyprowadzonym słowem. Dowód przebiega przez indukcję po wysokości drzewa wyprowadzenia h :

1. Jeśli $h = 1$, to $|x| = 1 = 2^0 = 2^{h-1}$.
2. Jeśli $h > 1$, to można podzielić x na dwa takie słowa y i z , $x = yz$, że wyprowadzenie x wygląda następująco:

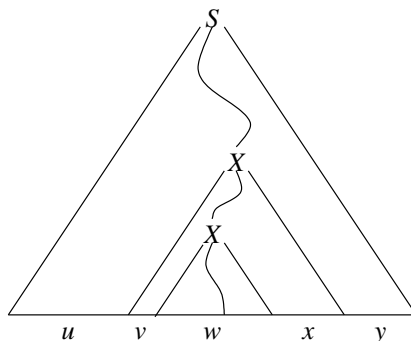


dla pewnych $Y, Z \in N$. Poddrzewa wyprowadzenia y z Y i z z Z mają wysokość co najwyżej $h - 1$. Tak więc, z założenia indukcyjnego mamy:

$$|x| = |xy| = |x| + |y| \leq 2^{h-2} + 2^{h-2} = 2^{h-1}$$

Na mocy zasady indukcji uzyskujemy lemat. □

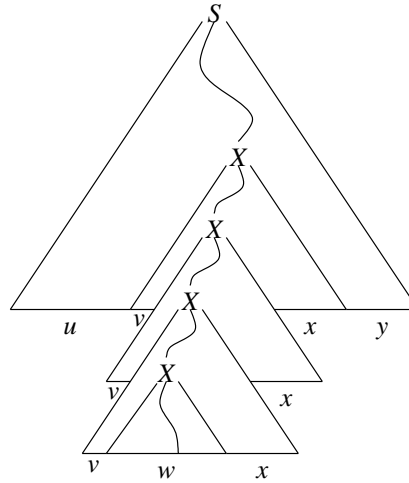
Niech $|N| = n$. Pokażemy, że można wziąć $k = 2^n$. Z powyższego faktu wynika, że jeżeli pewne słowo $z \in A$ ma długość przynajmniej $k = 2^n$, to jego drzewo wyprowadzenia ma wysokość przynajmniej $n + 1$. W takim drzewie wyprowadzenia istnieje ścieżka, na której jest przynajmniej $n + 1$ nieterminali. Wybierzmy i ustalmy jedną taką ścieżkę. Na takiej ścieżce pewien nieterminal musi się powtarzać. Niech X będzie pierwszym nieterminalem, który powtarza się na tej ścieżce idąc do góry (od liścia do korzenia). Jego powtórzenia wyznaczają podział słowa z na słowa u, v, w, x i y , $uvwxy = z$, w sposób przedstawiony na rysunku.



Drzewo wyprowadzenia jest regularne, a oba zaznaczone na rysunku wystąpienia nieterminala X to różne węzły drzewa. Wynika stąd, że v i x nie mogą być oba pustymi słowami. Czyli $vx \neq \varepsilon$.

Ponieważ X jest pierwszym terminalem, który powtarza się na wybranej ścieżce idąc od liścia do korzenia, więc poddrzewo wyprowadzenia vw z X jest wysokości co najwyżej $n + 1$. Stąd $|vw| \leq 2^n = k$.

Nieterminal X występuje (przynajmniej) w dwóch miejscach na ścieżce. Zauważmy, że możemy powielić fragment drzewa wyprowadzenia między pierwszym i drugim wystąpieniem nieterminala X na wybranej ścieżce, z którego pochodzą słowa v i x .



Jak widać, powielając ów fragment drzewa i razy uzyskujemy drzewo wyprowadzenia słowa uv^iwx^iy . Czyli $uv^iwx^iy \in A$. \square

12.3 Zastosowanie lematu o pompowaniu

Lemat o pompowaniu służy pokazywaniu, że określone języki **nie są** bezkontekstowe. Schemat takiego dowodu przebiega nie wprost. Lemat o pompowaniu ma postać implikacji: jeśli język jest bezkontekstowy, to ma pewne własności. Dowód, że język nie jest bezkontekstowy przebiega według schematu: skoro dany język nie posiada tych własności, to nie może być bezkontekstowy. Myśląc o takim zastosowaniu, możemy odwrócić lemat o pompowaniu:

Lemat o pompowaniu dla języków bezkontekstowych — sformułowanie alternatywne: Niech A będzie językiem. Jeśli dla każdego $k > 0$ istnieje takie słowo z , że $|z| \geq k$ oraz dla dowolnego podziału słowa z na słowa u, v, w, x i y , $uvwxy = z$ takiego, że $vx \neq \varepsilon$ i $|vwx| \leq k$ istnieje takie $i \geq 0$, że $uv^iwx^iy \notin A$, wówczas A nie jest bezkontekstowy.

Nawet w postaci odwróconej lemat o pompowaniu jest trudny do zrozumienia, gdyż mamy tu do czynienia z czterema poziomami zagnieżdżonych naprzemiennych kwantyfikatorów uniwersalnych i egzystencjalnych. Takie zagnieżdżenie naprzemiennych kwantyfikatorów możemy sobie wyobrazić jako ruchy dwóch graczy — gracze na przemian wskazują jakie wartości powinny przyjąć zmienne spod kwantyfikatorów, jeden gracz spod uniwersalnych, a drugi spod egzystencjalnych. Pierwszy gracz stara się pokazać, że dany język jest bezkontekstowy, a drugi, że nie. O wygranej decyduje, czy dla wybranych wartości zmiennych zachodzi własność objęta kwantyfikatorami. Gra jest skończona, więc jeden z graczy ma strategię wygrywającą. W zależności od tego, który z graczy to jest, język może być bezkontekstowy lub nie.

Lemat o pompowaniu dla języków bezkontekstowych — gra z Demonem: Prowadzimy z Demonem grę. Staramy się pokazać, że język nie jest bezkontekstowy, a Demon stara się pokazać, że język mimo wszystko jest bezkontekstowy. Zasady gry są następujące:

1. Demon wybiera $k > 0$. (Jeśli A jest faktycznie bezkontekstowy, to Demon może wybrać za $k = 2^{|N|}$.)
2. Wybieramy takie słowo $z \in A$, że $|z| \geq k$.
3. Demon dzieli z na takie słowa u, v, w, x, y , $uvwxy = z$, że $vx \neq \varepsilon$ i $|vwx| \leq k$. (Jeśli A jest faktycznie bezkontekstowy, to może to być podział wyznaczony przez powtarzający się nieterminal na pewnej ścieżce od liścia do korzenia.)
4. Wybieramy $i \geq 0$.

Jeśli $uv^iwx^iy \notin A$, to wygraliśmy, wpp. wygrał Demon.

Oryginalne sformułowanie lematu o pompowaniu mówi, że jeżeli język jest bezkontekstowy, to Demon ma strategię wygrywającą. Natomiast odwrotne sformułowanie mówi, że jeżeli my mamy strategię wygrywającą, to język nie jest bezkontekstowy. Pamiętajmy, że lemat o pompowaniu nie służy do pokazywania, że język jest bezkontekstowy. Nie zawsze pozwala on na pokazanie, że język, który nie jest bezkontekstowy faktycznie taki nie jest — istnieją języki, które nie są bezkontekstowe i spełniają lemat o pompowaniu.

Przykład: Standardowym przykładem języka, który nie jest bezkontekstowy to $A = \{a^n b^n c^n : n \geq 0\}$. Udowodnimy to korzystając z lematu o pompowaniu sformułowanego jako gra z Demonem. Nasza strategia wygląda następująco:

- Demon wybiera k .
- Wybieramy słowo $z = a^k b^k c^k$.
- Demon dzieli słowo $z = uvwxy$.
- Wybieramy $i = 2$.

Wiemy, że $|vwx| \leq k$ i $vx \neq \varepsilon$. Tak więc słowa v i x mogą zawierać jeden rodzaj znaków lub dwa, ale nie wszystkie trzy. Czyli nie może zachodzić $\#_a(uv^iwx^iy) = \#_b(uv^iwx^iy) = \#_c(uv^iwx^iy)$. Stąd wynika, że $uv^iwx^iy \notin A$.

Tak więc powyższa strategia jest wygrywająca, czyli A nie jest bezkontekstowy. \square

Przykład: Pokażemy, że język $B = \{a^i b^j a^i b^j : i, j \geq 0\}$ nie jest bezkontekstowy. Nasza strategia wygląda następująco:

- Demon wybiera k .
- Wybieramy słowo $z = a^k b^k a^k b^k$.
- Demon dzieli słowo $z = uvwxy$.
- Wybieramy $i = 2$.

Wiemy, że $|vwx| \leq k$ i $vx \neq \varepsilon$. Możliwe są trzy przypadki:

Jeżeli słowo v lub x nie jest w całości zbudowane z jednego rodzaju znaków, to $uv^iwx^iy \notin L(a^*b^*a^*b^*)$ i tym samym $uv^iwx^iy \notin B$.

Jeżeli słowa v i x są zbudowane w całości z tych samych znaków, to znaczy że oba pochodzą z tego samego pod słowa postaci a^k lub b^k . Ale wówczas $uv^iwx^iy \notin B$.

Pozostaje jeszcze przypadek, gdy słowo v jest w całości złożone z liter a , a x z liter b (lub odwrotnie). Wówczas jednak również $uv^iwx^iy \notin B$.

Tak więc powyższa strategia jest wygrywająca, czyli B nie jest bezkontekstowy. \square

12.4 Inne metody dowodzenia niebezkontekstowości języków

Zachodzi następujący fakt, którego nie będziemy teraz dowodzić:

Fakt 15. *Przecięcie języka bezkontekstowego i regularnego jest bezkontekstowe.*

Zachodzi też inny fakt:

Fakt 16. *Suma języków bezkontekstowych jest bezkontekstowa.*

Dowód: Jeżeli A i B są językami bezkontekstowymi, to istnieją generujące je gramatyki $G_A = \langle N_A, \Sigma, P_A, S_A \rangle$ i $G_B = \langle N_B, \Sigma, P_B, S_B \rangle$. Co więcej, możemy założyć, że gramatyki te mają rozłączne zbiory nieterminali $N_A \cap N_B = \emptyset$. Niech S będzie nowym symbolem nieterminalnym nie występującym ani w G_A , ani w G_B . Wówczas gramatyka:

$$G = \langle N_A \cup N_B \cup \{S\}, \Sigma, P_A \cup P_B, S \rangle$$

generuje język $A \cup B$. \square

Korzystając z tych faktów możemy pokazywać nie wprost, że języki nie są bezkontekstowe. Schemat rozumowania jest następujący: Zakładamy, że dany język jest bezkontekstowy. Z powyższych faktów wynika, że jakiś inny język, który nie jest bezkontekstowy, musiałby taki być. A więc uprzednie założenie musiało być fałszywe. Stąd dany język nie jest bezkontekstowy.

Przykład: Pokażemy, że język $C = \{ww : w \in \{a, b\}^*\}$ nie jest bezkontekstowy. Założmy przeciwnie, że jest. Wówczas:

$$C \cap L(a^*b^*a^*b^*) = \{a^ib^ja^ib^j : i, j \geq 0\} = B$$

Wiemy wszak, że B nie jest bezkontekstowy, a więc C też nie może być. \square

12.5 Podsumowanie

Wykład ten był poświęcony przede wszystkim lematowi o pompowaniu dla języków bezkontekstowych. Poznaliśmy trzy równoważne sformułowania tego lematu, oraz przykłady jego zastosowania do wykazywania, że określone języki nie są bezkontekstowe. Poznaliśmy też metodę wykazywania, że dane języki nie są bezkontekstowe, w oparciu o własności klasy języków bezkontekstowych.

12.6 Skorowidz

- **Lemat o pompowaniu dla języków bezkontekstowych** mówi, że dla każdego języka bezkontekstowego, każde dostatecznie długie słowo w tym języku można „pompować” (w dwóch miejscach równocześnie) dowolnie je wydłużając. Wykazywanie braku tej własności jest metodą pokazywania, że dany język nie jest bezkontekstowy.
- **Gra z Demonym** — alternatywne sformułowanie lematu o pompowaniu dla języków bezkontekstowych jako gry.

12.7 Praca domowa

1. Udowodnij, że język $D = \{a^{2^i} : i \geq 0\}$ nie jest bezkontekstowy.
2. Dla każdego z poniższych języków określ, czy jest on: (i) regularny, (ii) bezkontekstowy, ale nie regularny, (iii) nie jest bezkontekstowy:
 - $\{a^i b^j c^k : 2i + j = k\}$,
 - $\{a^i b^j c^k : i + j + k = 42\}$,
 - $\{w \in \{a, b, c\}^* : \#_a(w) = \#_b(w) = \#_c(w)\}$,
 - $\{a^i b^j c^k : i + 2j = k \vee 3i + j = 2k\}$,
 - $\{a^{n^2} : n \geq 0\}$.

12.8 Ćwiczenia

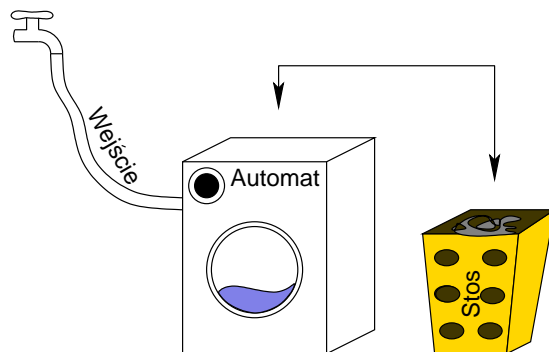
- Udowodnij, że następujące języki nie są bezkontekstowe:
 - $\{a^{n^5} : n \geq 0\}$,
 - $\{a^i b^j c^k : k = i \cdot j\}$,
 - $\{a^p : p \text{ jest liczbą Fibonacciego}\}$,
 - $\{w \in \{a, b, c, d\}^* : \#_a(w) = \#_b(w) \wedge \#_c(w) = \#_d(w)\}$,
 - $\{a^n b^{n!} : n \geq 0\}$.
- Dla każdego z poniższych języków określ, czy jest on: (i) regularny, (ii) bezkontekstowy, ale nie regularny, (iii) nie jest bezkontekstowy:

- $\{a^i b^j c^k : 3i = 2j + k\}$,
- $\{a^i b^j : i + j = k^3 \text{ dla pewnego } k \geq 0\}$,
- $\{a^i b^i c^i : i \leq 42\}$,
- $\{aw_1 b w_2 c : w_1, w_2 \in \{a, b, c\}^*\}$,
- $\{a^i b^j : i > j\}$,
- $\{a^i b^j : i \neq j\}$,
- $\{awbwc : w \in \{a, b, c\}^*\}$,
- $\{a^i b^j : i + 2j = 42\}$,
- $\{a^i b^j c^k : i \cdot j = 2^k\}$,
- $\{a^i b^j c^k : i \cdot j \cdot k = 6\}$.

Wykład 13. Automaty stosowe

13.1 Wstęp

W tym wykładzie poznamy model obliczeń odpowiadający gramatykom bezkontekstowym (tak jak automaty skończone odpowiadały wzorcom). Ten model, to automaty stosowe. Automaty stosowe możemy sobie wyobrazić jako automaty skończone wyposażone w dodatkową nieograniczoną pamięć, ale działającą na zasadzie stosu.



To znaczy, na wierzch stosu można zawsze włożyć nowy element. Jeśli stos nie jest pusty, to można zawsze zobaczyć jaki element jest na wierzchu i ew. zdjąć ten element. Nie można natomiast oglądać ani wyjmować elementów z wnętrza stosu.

Zobaczmy na przykładach jak można konstruować automaty stosowe i pokażemy, że faktycznie ich siła wyrazu jest taka sama jak gramatyk bezkontekstowych.

Pod względem siły wyrazu automaty stosowe są też równoważne imperatywnym językom programowania bez plików (pomocniczych) dynamicznych struktur danych, wskaźników i parametrów proceduralnych. Przy tych ograniczeniach dane globalne w programie są stałego rozmiaru, a danych dynamicznych nie ma. Można korzystać z rekurencji i nie ma ograniczenia na rozmiar stosu, jednak wielkość danych związanych z każdym poziomem rekurencji jest ograniczona przez stałą i nie ma możliwości wyciągania danych z głębi stosu. Jak zobaczymy dokładnie odpowiada to automatom stosowym.

13.2 Automaty stosowe

Automat stosowy, podobnie jak automat skończony, ma skończony zbiór stanów i stan początkowy. Dodatkowo jest on wyposażony w stos, na który może wkładać elementy, podglądać element znajdujący się na wierzchu i zdejmować elementy. Przejścia automatu stosowego są trochę bardziej skomplikowane niż przejścia w automacie skończonym, gdyż oprócz wczytywania znaków z wejścia i zmiany stanów obejmują również operacje na stosie. Automat taki może być niedeterministyczny i może zawierać ε -przejścia, tzn. wykonaniu przejścia nie musi towarzyszyć wczytanie znaku z wejścia. Natomiast przedstawione tutaj automaty stosowe będą pozbawione stanów akceptujących — automat będzie sygnalizował

akceptację wejścia opróżniając stos¹³.

Chcąc opisać działanie automatów stosowych musimy wprowadzić pojęcie *konfiguracji*. Konfiguracja będzie obejmować stan automatu (ze skończonego zbioru stanów), zawartość stosu i pozostałe do wczytania wejście. Automat stosowy w jednym kroku obliczeń wykonuje następujące czynności:

1. Podgląda znak czekający na wczytanie na wejściu.
2. Podgląda element na wierzchołku stosu.
3. Na podstawie tych informacji wybiera jedno z przejść do wykonania.
4. Jeśli to nie jest ε -przejście, to wczytywany jest jeden znak z wejścia.
5. Zdejmowany jest element z wierzchołku stosu.
6. Pewna liczba określonych elementów może być włożona na stos.
7. Zgodnie z przejściem zmieniany jest stan.

Może wydawać się nienaturalne, że w każdym kroku zdejmujemy element z wierzchołka stosu. Jednak dzięki temu, opisany powyżej krok automatu jest na tyle elastyczny, że może zdejmować elementy ze stosu, wkładać nowe, podmieniać element na wierzchołku lub nie zmieniać stosu. Jeżeli chcemy zdjąć element z wierzchołku stosu, to wystarczy, że nie będziemy nic wkładać. Jeżeli chcemy podmienić element na wierzchołku, to wystarczy, że włożymy jeden element, który zastąpi element zdjęty z wierzchołka. Jeżeli chcemy włożyć elementy na stos, to wkładamy na stos właśnie zdjęty element, wraz z elementami, które mają być umieszczone na stosie. Jeśli nie chcemy zmieniać zawartości stosu, to po prostu wkładamy dokładnie ten sam element, który właśnie został zdjęty. Możemy włożyć na stos kilka elementów w jednym kroku, ale niestety nie możemy w jednym kroku zdjąć więcej niż jeden element.

Zauważmy, że powyższy schemat postępowania wymaga, żeby stos nie był pusty. Inaczej nie ma czego podglądać ani zdejmować ze stosu. Dlatego też, w momencie uruchomienia automatu stos nie jest pusty. Zawiera jeden wyróżniony element, który będziemy nazywać „pinezką” i oznaczać przez \perp . Jak tylko stos zostaje opróżniony, automat zatrzymuje się i dalsze jego działanie nie jest możliwe. Jeżeli automat zatrzymuje się z pustym stosem po wczytaniu **całego** słowa z wejścia, to przyjmujemy, że słowo to zostało zaakceptowane.

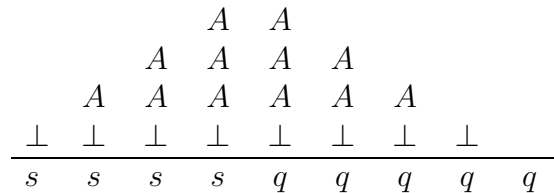
Zanim formalnie zdefiniujemy automaty stosowe i akceptowane przez nie języki, zobaczmy prosty nieformalny przykład:

¹³W literaturze można spotkać trzy różne definicje, wzajemnie równoważnych automatów stosowych: akceptujących pustym stosem, akceptujących stanem i akceptujących pustym stosem i stanem równocześnie. Tutaj przedstawiamy najprostszy z tych rodzajów automatów stosowych — akceptujące pustym stosem.

Przykład: Zbudujemy automat akceptujący słowa postaci $a^n b^n$, dla $n \geq 0$. Automat taki będzie miał dwa stany, s i q , przy czym s będzie stanem początkowym. Na stosie będziemy trzymać dwa rodzaje elementów: \perp i A . Będąc w stanie s automat ma dwa przejścia do wyboru: może wczytać z wejścia znak a i włożyć na stos A , lub nic nie wczytywać, ani nie zmieniać zawartości stosu i przejść do stanu q . W stanie q automat ma również dwa przejścia do wyboru: jeżeli na wejściu czeka na wczytanie znak b i na wierzchołku stosu jest A , to może wczytać z wejścia b i zdjąć ze stosu A . Jeżeli natomiast na wierzchołku stosu jest \perp , to automat może zdjąć ze stosu \perp i tym samym zakończyć działanie.

Jedyna możliwość, żeby zaakceptować słowo wygląda tak: Automat wczytuje sekwencję a^n i umieszcza na stosie A^n . Następnie przechodzi do stanu q i wczytuje znaki b zdejmując ze stosu A . Jednych i drugich musi być tyle samo, czyli wczytana zostaje sekwencja b^n . Po wczytaniu tej sekwencji odsłania się \perp i jest zdejmowana ze stosu. Jeżeli automat ma zaakceptować słowo, to musi ono być w tym momencie całe wczytane, czyli było to słowo $a^n b^n$.

Poniższy diagram przedstawia obliczenie akceptujące słowo $aaabbb$. Widać na nim kolejne zawartości stosu i stany, w których jest automat.



Def. 42. Automat stosowy to dowolna taka piątka $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$, w której:

- Q to skończony zbiór stanów,
- Σ to (skończony) alfabet wejściowy,
- Γ to (skończony) alfabet stosowy,
- δ to relacja przejścia, $\delta \subseteq (Q \times \Gamma \times (\Sigma \cup \{\varepsilon\})) \times (Q \times \Gamma^*)$; określa ona dla danego aktualnego stanu, znaku na wierzchołku stosu i znaku czekającego na wczytanie (lub ε jeżeli nic nie jest wczytywane), do jakiego stanu należy przejść i co należy włożyć na stos (uprzednio zdjąwszy element z wierzchołka stosu),
- $s \in Q$ to stan początkowy,
- $\perp \in \Gamma$ to symbol początkowy na stosie, tzw. „pinezka”.

□

Notacja: W przypadku automatów stosowych, przedstawienie przejść w postaci diagramu nie jest czytelne. Nie będziemy jednak pisać:

$$\delta = \{((s, a, A), (s, AA)), ((s, a, \perp), (s, A \perp)), \dots\}$$

jakby wynikało z definicji, gdyż byłoby to zupełnie nieczytelne. Zamiast tego będziemy przedstawiać przejścia jako reguły postaci:

$$\begin{array}{ccc} (s, A) & \xrightarrow{a} & (s, AA) \\ (\perp) & \xrightarrow{a} & (s, A \perp) \end{array}$$

lub w postaci skrótów:

$$(s, x) \xrightarrow{a} (s, Ax) \text{ dla } x = A, \perp$$

Konfiguracje automatów stosowych definiujemy następująco:

Def. 43. Niech $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$ będzie ustalonym automatem stosowym. Konfiguracja takiego automatu stosowego, to dowolna trójka:

$$\langle q, \alpha, \beta \rangle \in Q \times \Sigma^* \times \Gamma^*$$

gdzie q reprezentuje aktualny stan, α pozostałe do wczytania wejście, a β zawartość stosu (czytaną od wierzchołka w głąb).

Konfiguracja początkowa dla słowa x , to $\langle s, x, \perp \rangle$. □

Jeden krok obliczeń automatu stosowego opisujemy relacją \rightarrow określoną na konfiguracjach.

Def. 44. Niech $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$ będzie ustalonym automatem stosowym. Relacja przejścia między konfiguracjami tego automatu \rightarrow to najmniejsza relacja spełniająca następujące warunki:

- jeśli δ zawiera przejście postaci $(p, A) \xrightarrow{a} (q, \gamma)$ (dla $p, q \in Q$, $a \in \Sigma$, $A \in \Gamma$ i $\gamma \in \Gamma^*$), to dla dowolnych $y \in \Sigma^*$ i $\beta \in \Gamma^*$ mamy: $\langle p, ay, A\beta \rangle \rightarrow \langle q, y, \gamma\beta \rangle$,
- jeśli δ zawiera przejście postaci $(p, A) \xrightarrow{\varepsilon} (q, \gamma)$, (dla $p, q \in Q$, $A \in \Gamma$ i $\gamma \in \Gamma^*$) to dla dowolnych $y \in \Sigma^*$, $\beta \in \Gamma^*$ mamy: $\langle p, y, A\beta \rangle \rightarrow \langle q, y, \gamma\beta \rangle$.

Przez \rightarrow^* będziemy oznaczać zwrotnio-przechodnie domknięcie relacji \rightarrow . □

Inaczej mówiąc, $(p, A) \xrightarrow{a} (q, \gamma)$ oznacza, że możemy będąc w stanie p i mając na wierzchołku stosu A przejść do stanu q wczytując z wejścia a i zastępując na stosie A przez γ . Podobnie $(p, A) \xrightarrow{\varepsilon} (q, \gamma)$ oznacza, że możemy będąc w stanie p i mając na wierzchołku stosu A przejść do stanu q nie wczytując nic z wejścia i zastępując na stosie A przez γ .

Relacja \rightarrow opisuje pojedyncze kroki obliczenia, a relacja \rightarrow^* dowolnie długie obliczenia.

Def. 45. Niech $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$ będzie ustalonym automatem stosowym. Język akceptowany przez automat stosowy M jest określony następująco:

$$L(M) = \{x \in \Sigma^* : \exists_{q \in Q} \langle s, x, \perp \rangle \rightarrow^* \langle q, \varepsilon, \varepsilon \rangle\}$$

□

Inaczej mówiąc, automat stosowy akceptuje takie słowa x , dla których istnieją obliczenia prowadzące od konfiguracji początkowej (dla słowa x) do konfiguracji, w której całe słowo zostało wczytane, a stos opróżniony. Pamiętajmy, że automat stosowy może być niedeterministyczny, a więc wystarczy żeby istniało jedno obliczenie akceptujące, żeby słowo było akceptowane przez automat.

Przykład: Skonstruujemy automat akceptujący słowa zawierające tyle samo liter a , co b . Automat ten będzie używać stosu jako licznika. Z każdym wczytaniem a licznik zwiększa się o 1, a z każdym wczytaniem b zmniejsza się o 1. Dodatkowo wartości licznika reprezentujemy jako \perp i stos A , a ujemne, jako \perp i stos B . Dodatkowo, w momentach gdy odsłonięta jest pinezka, automat może ją zdjąć ze stosu. Automat akceptuje tylko takie słowa x , po których wczytaniu licznik jest wyzerowany, czyli $\#_a(x) = \#_b(x)$.

Automat ten ma jeden stan s i działa w następujący sposób:

$$\begin{array}{ll} (s, \perp) \xrightarrow{a} (s, A \perp) & (s, \perp) \xrightarrow{b} (s, B \perp) \\ (s, A) \xrightarrow{a} (s, AA) & (s, B) \xrightarrow{b} (s, BB) \\ (s, B) \xrightarrow{a} (s, \varepsilon) & (s, A) \xrightarrow{b} (s, \varepsilon) \\ (s, \perp) \xrightarrow{\varepsilon} (s, \varepsilon) & \end{array}$$

Przykładowe obliczenie dla $aabbba$:

$$\begin{array}{cccccccc} & & & & A & & & \\ & & & & A & A & A & B \\ & & & & \perp & \perp & \perp & \perp \\ \hline s & s & s & s & s & s & s & s \end{array}$$

Oczywiście taki automat może usunąć pinezkę gdy jeszcze nie całe słowo jest wczytane. Takie obliczenie nie doprowadzi do zaakceptowania słowa. Pamiętajmy jednak, że automaty stosowe są niedeterministyczne i wystarczy, aby istniało choć jedno obliczenie akceptujące dane słowo, aby należało ono do języka akceptowanego przez automat stosowy.

Przykład: Skonstruujemy automat akceptujący wyrażenia nawiasowe zbudowane z kwadratowych nawiasów [i]. Automat ten działa podobnie jak poprzedni, ale licznik może przyjmować tylko dodatnie wartości. Na stosie trzymamy pinezkę i tyle nawiasów zamykających ile wynosi wartość licznika. Z wczytaniem każdego nawiasu otwierającego licznik

jest zwiększany o 1, a z wczytaniem każdego nawiasu zamykającego zmniejszany o 1.

$$\begin{aligned} (s, x) &\xrightarrow{\downarrow} (s,]x) \text{ dla } x =], \perp \\ (s,]) &\xrightarrow{\downarrow} (s, \varepsilon) \\ (s, \perp) &\xrightarrow{\varepsilon} (s, \varepsilon) \end{aligned}$$

Przykładowe obliczenie dla $[[[]]]$:

$$\begin{array}{cccccccccc} & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp \\ \hline s & s & s & s & s & s & s & s & s & s & s \end{array}$$

Przykład: Automat akceptujący palindromy nad alfabetem $\{a, b\}$. Automat ten musi niedeterministycznie zgadnąć gdzie jest środek palindromu i czy jest on parzystej, czy nieparzystej długości.

Nasz automat będzie miał dwa stany: s i p . Na stosie będziemy przechowywać pinezkę i znaki alfabetu wejściowego. W pierwszej fazie (stan s) automat zapamiętuje pierwszą połowę słowa na stosie. W drugiej fazie (stan p) zdejmuję znaki ze stosu, sprawdzając, czy są one takie same jak wczytywane znaki. Jeżeli słowo jest nieparzystej długości, to przechodząc z pierwszej fazy do drugiej wczytujemy środkowy znak nie zmieniając zawartości stosu. Jeżeli zaś jest parzystej długości, to przechodząc do drugiej fazy nic nie wczytujemy.

$$\begin{aligned} (s, x) &\xrightarrow{y} (s, yx) \quad x = \perp, a, b \quad y = a, b \\ (s, x) &\xrightarrow{y} (p, x) \quad x = \perp, a, b \quad y = a, b, \varepsilon \\ (p, x) &\xrightarrow{x} (p, \varepsilon) \quad x = a, b \\ (p, \perp) &\xrightarrow{\varepsilon} (s, \varepsilon) \end{aligned}$$

Przykładowe obliczenia akceptujące słowa $abba$ i $baabaab$:

$$\begin{array}{cccccccc} & & & & b & b & & & \\ & & & & a & a & a & a & \\ & & & & \perp & \perp & \perp & \perp & \perp & \perp \\ \hline s & s & s & p & p & p & p & & \\ & & & & a & a & & & \\ & & & & a & a & a & a & \\ & & & & b & b & b & b & b & \\ \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp \\ \hline s & s & s & s & p & p & p & p & p & p \end{array}$$

Przykład: W poprzednim wykładzie pokazaliśmy, że język $\{ww : w \in \{a, b\}^*\}$ nie jest bezkontekstowy. Za chwilę pokażemy, że automaty stosowe akceptują dokładnie języki bezkontekstowe. Teraz jednak skonstruujemy automat akceptujący język $\overline{\{ww : w \in \{a, b\}^*\}}$, pokazując tym samym, że jest on bezkontekstowy.

Słowo $x = a_1a_2 \dots a_k$ należy do języka $\overline{\{ww : w \in \{a, b\}^*\}}$ w dwóch przypadkach. Albo jest nieparzystej długości. Albo jest parzystej długości i istnieje taka pozycja i ($1 \leq i \leq \frac{k}{2}$) w pierwszej połowie tego słowa, że na i -tej pozycji w słowie x jest inny znak niż na odpowiadającej jej pozycji w drugiej połowie słowa, $a_i \neq a_{i+\frac{k}{2}}$.

Nasz automat będzie miał szereg stanów, przy czym stanem początkowym będzie s . Na stosie będziemy trzymać dwa rodzaje symboli: \perp i \circ . Automat zgaduje najpierw, czy słowo jest parzystej (stan p), czy nieparzystej (stan n) długości. Jeśli nieparzystej, to tylko to sprawdza, że faktycznie tak jest.

$$\begin{aligned} (s, \perp) &\xrightarrow{\varepsilon} (n, \perp) \\ (s, \perp) &\xrightarrow{\varepsilon} (p, \perp) \\ (n, \perp) &\xrightarrow{x} (m, \perp) \text{ dla } x = a, b \\ (m, \perp) &\xrightarrow{x} (n, \perp) \text{ dla } x = a, b \\ (m, \perp) &\xrightarrow{\varepsilon} (m, \varepsilon) \end{aligned}$$

Jeśli słowo jest parzystej, to automat stosowy musi zgadnąć niedeterministycznie na jakiej pozycji i istnieje niezgodność między pierwszą i drugą połową słowa, a następnie sprawdzić, że faktycznie tak jest. Inaczej mówiąc, słowo x musi być postaci $uavbw$ lub $ubvaw$, przy czym $|v| = |u| + |w|$. To czy w pierwszej połowie na i -tej pozycji jest a czy b można zapamiętać w stanie automatu. Problemem jest natomiast odliczenie odpowiadających sobie pozycji w obu połówkach słowa. Zauważmy jednak, że automat nie musi wiedzieć gdzie jest połowa długości słowa. Wystarczy, że będzie zachodzić $|v| = |u| + |w|$.

Nasz automat najpierw wczytuje słowo u , wkładając na stos $|u|$ znaków \circ .

$$(p, i) \xrightarrow{j} (p, \circ i) \text{ dla } i = \perp, \circ, j = a, b$$

Następnie niedeterministycznie zgaduje kiedy jest koniec u i wczytuje jeden znak, zapamiętując go na stanie.

$$(p, i) \xrightarrow{j} (q_j, i) \text{ dla } i = \perp, \circ, j = a, b$$

Następnie nasz automat wczytuje tyle znaków z wejścia, ile \circ jest na stosie. W ten sposób wczytujemy prefiks v długości $|u|$. Po zdjęciu ze stosu wszystkich \circ przechodzimy do kolejnej fazy.

$$\begin{aligned} (q_i, \circ) &\xrightarrow{j} (q_i, \varepsilon) \text{ dla } i = a, b, j = a, b \\ (q_i, \perp) &\xrightarrow{\varepsilon} (r_i, \perp) \text{ dla } i = a, b \end{aligned}$$

Następnie nasz automat powinien wczytać z wejścia $|w| = |v| - |u|$ znaków słowa v . Liczbę tę zgaduje niedeterministycznie, wkładając równocześnie $|w|$ znaków \circ na stos. W momencie gdy automat zgaduje niedeterministycznie, że całe słowo v zostało wczytane, wczytuje

jeden znak z wejścia sprawdzając równocześnie, że jest on różny od znaku zapamiętanego na stanie.

$$\begin{aligned} (r_i, j) &\xrightarrow{l} (r_i, \circ j) \text{ dla } i = a, b, j = \perp, \circ, l = a, b \\ (r_i, j) &\xrightarrow{l} (t, j) \text{ dla } i = a, b, j = \perp, \circ, l = a, b, l \neq i \end{aligned}$$

Pozostało jeszcze do wczytania w , przy czym na stosie powinno być $|w|$ znaków \circ . Po wczytaniu słowa w powinna odsłonić się pinezka, którą zdejmujemy.

$$\begin{aligned} (t, \circ) &\xrightarrow{i} (t, \varepsilon) \text{ dla } i = a, b \\ (t, \perp) &\xrightarrow{\varepsilon} (t, \varepsilon) \end{aligned}$$

Oto przykładowe obliczenie akceptujące słowo $abbaabaa$.

$$\begin{array}{cccccccccccc} & & & & \circ & \circ & & & & \circ & \circ & & \\ & & & \circ & \circ & \circ & \circ & & & & & & \\ \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \perp & \\ \hline s & p & p & p & q_b & q_b & q_b & r_b & r_b & t & t & t & \end{array}$$

Skonstruowany automat jest wysoce niedeterministyczny. Jeśli gdziekolwiek cokolwiek zostanie źle zgadnięte, to obliczenie nie doprowadzi do zaakceptowania słowa. Pamiętajmy jednak, iż wystarczy, że *można* tak zgadnąć, aby obliczenie było akceptujące, aby słowo było akceptowane przez automat.

13.3 Równoważność gramatyk bezkontekstowych i automatów stosowych

Pokażemy teraz, że siła wyrazu automatów stosowych jest dokładnie taka sama, jak gramatyk bezkontekstowych — tzn. akceptują one dokładnie języki bezkontekstowe.

Tw. 5. Niech A będzie językiem nad alfabetem Σ . Następujące stwierdzenia są równoważne:

- A jest bezkontekstowy, tzn. A jest generowany przez pewną gramatykę bezkontekstową G , $A = L(G)$,
- istnieje automat stosowy M akceptujący A , $A = L(M)$.

Przedstawimy tu jedynie szkic dowodu — oprzemy się na trzech faktach. Pokażemy jak dla danej gramatyki bezkontekstowej skonstruować równoważny jej automat stosowy (z jednym stanem). Następnie pokażemy, że konstrukcję tę można odwrócić, tzn. dla danego automatu stosowego z jednym stanem można skonstruować równoważną mu gramatykę bezkontekstową. Na koniec pokażemy, jak zredukować liczbę stanów w automacie stosowym do jednego.

13.3.2 Gramatyka bezkontekstowa równoważna danemu automатовi stosowemu z jednym stanem

Przedstawioną w poprzednim punkcie konstrukcję można odwrócić. Załóżmy, że mamy dany automat stosowy z jednym stanem $M = \langle \{s\}, \Sigma, \Gamma, \delta, s, \perp \rangle$. Bez zmniejszenia ogólności możemy założyć, że $\Gamma \cap \Sigma = \emptyset$. Nasza gramatyka ma postać $G = \langle \Gamma, \Sigma, P, \perp \rangle$, przy czym jeśli δ zawiera przejście $(s, A) \xrightarrow{a} (s, \gamma)$ (dla $a \in \Sigma \cup \{\varepsilon\}$), to w P mamy produkcję $A \rightarrow a\gamma$.

Odpowiedniość jest taka jak poprzednio. Jeśli wczytywane słowo to x , $x = yz$, to konfiguracji $\langle s, z, \alpha \rangle$ odpowiada w wyprowadzeniu lewostronne $\perp \rightarrow^* y\alpha$. Konfiguracja początkowa $\langle s, x, \perp \rangle$ odpowiada aksjomatowi \perp , a konfiguracja akceptująca $(s, \varepsilon, \varepsilon)$ odpowiada w wyprowadzeniu $\perp \rightarrow^* x$. Stąd $L(M) = L(G)$.

Przykład: Weźmy, przedstawiony powyżej, automat stosowy akceptujący słowa zawierające tyle samo liter a , co b .

$$\begin{array}{ll} (s, \perp) \xrightarrow{a} (s, A \perp) & (s, \perp) \xrightarrow{b} (s, B \perp) \\ (s, A) \xrightarrow{a} (s, AA) & (s, B) \xrightarrow{b} (s, BB) \\ (s, B) \xrightarrow{a} (s, \varepsilon) & (s, A) \xrightarrow{b} (s, \varepsilon) \\ (s, \perp) \xrightarrow{\varepsilon} (s, \varepsilon) & \end{array}$$

Stosując przedstawioną powyżej metodę, uzyskujemy następującą równoważną mu gramatykę (dla czytelności \perp zastąpiono przez S):

$$\begin{array}{l} S \rightarrow aAS \mid bBS \mid \varepsilon \\ A \rightarrow b \mid aAA \\ B \rightarrow a \mid bBB \end{array}$$

Oto wyprowadzenie w tej gramatyce słowa $aabbba$:

$$S \rightarrow aAS \rightarrow aaAAS \rightarrow aabAS \rightarrow aabbS \rightarrow aabbbBS \rightarrow aabbbaS \rightarrow aabbba$$

13.3.3 Redukcja zbioru stanów do jednego

Pozostało nam pokazać, że dla każdego automatu stosowego $M = \langle Q, \Sigma, \Gamma, \delta, s, \perp \rangle$ istnieje równoważny mu automat z jednym stanem $M' = \langle \{s\}, \Sigma, \Gamma', \delta', s \rangle$. Bez zmniejszenia ogólności możemy założyć, że automat M w momencie opróżnienia stosu będzie zawsze w jednym stanie q .

Na stosie automaty M' będziemy trzymać trójki postaci $\langle p, a, r \rangle \in \Gamma' = Q \times \Gamma \times Q$. Jeżeli na wierzchołku stosu M' jest trójka $\langle p, A, r \rangle$, to odpowiada to sytuacji, gdy na wierzchołku stosu w M jest A , M jest w stanie p i po zdjęciu symbolu A ze stosu przejdzie do stanu r . Stany pamiętane na stosie M' muszą się „zazębiać”, tzn. pod trójką $\langle p, A, r \rangle$ musi być trójka postaci $\langle r, B, t \rangle$ (dla pewnych $B \in \Gamma$ i $t \in Q$) itd.

Automat M' na początku ma na stosie trójkę $\langle s, \perp, q \rangle$ i akceptuje słowo x wtw., gdy M zaczynając w stanie s i z symbolem \perp na stosie, wczytuje x i kończy w stanie q z pustym stosiem.

Relacja δ' jest określona następująco:

- jeśli δ zawiera przejście:

$$(p, A) \xrightarrow{a} (r, B_1 B_2 \dots B_k)$$

dla pewnych $a \in \Sigma \cup \{\varepsilon\}$, $p, r \in Q$, $k > 0$, $A, B_1, B_2, \dots, B_k \in \Gamma$, to δ' (dla dowolnych $q_1, \dots, q_k \in Q$) zawiera przejścia:

$$(s, \langle p, A, q_k \rangle) \xrightarrow{a} (s, \langle r, B_1, q_1 \rangle \langle q_1, B_2, q_2 \rangle \dots \langle q_{k-1}, B_k, q_k \rangle)$$

- jeśli δ zawiera przejście

$$(p, A) \xrightarrow{a} (r, \varepsilon)$$

dla pewnych $a \in \Sigma \cup \{\varepsilon\}$, $p, r \in Q$, $A \in \Gamma$, to δ' zawiera przejście:

$$(s, \langle p, A, r \rangle) \xrightarrow{a} (s, \varepsilon)$$

Automat ten jest wysoce niedeterministyczny. Wkładając trójki na stos musimy z góry przewidzieć w jakim stanie będzie automat M w momencie zdjęcia odpowiedniej trójki w M' . Oczywiście zgadnięcie odpowiednich stanów jest możliwe, a więc M' będzie mógł zasymulować każde obliczenie M (i tylko obliczenia M).

13.4 Podsumowanie

W tym wykładzie poznaliśmy automaty stosowe. Zobaczyliśmy na przykładach jak je można konstruować. Dowiedzieliśmy się też, że są one równoważne gramatykom bezkontekstowym.

13.5 Skorowidz

- **Automat stosowy** to dowolna taka piątka $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$, w której: Q to skończony zbiór stanów, Σ to alfabet wejściowy, Γ to alfabet stosowy, $\delta \subseteq (Q \times \Gamma \times (\Sigma \cup \{\varepsilon\})) \times (Q \times \Gamma^*)$ to relacja przejścia, $s \in Q$ to stan początkowy, $\perp \in \Gamma$ to symbol początkowy na stosie.
- **Konfiguracja automatu stosowego** to dowolna trójka $\langle q, \alpha, \beta \rangle \in Q \times \Sigma^* \times \Gamma^*$ gdzie $M = \langle Q, \Sigma, \Gamma, \delta, s \rangle$ jest danym automatem stosowym.

- Język akceptowany przez automat stosowy to

$$L(M) = \{x \in \Sigma^* : \exists_{q \in Q} \langle s, x, \perp \rangle \xrightarrow{*} \langle q, \varepsilon, \varepsilon \rangle\}$$

13.6 Praca domowa

1. Skonstruuj automat stosowy akceptujący wyrażenia nawiasowe zawierające trzy rodzaje nawiasów: $()\{\}\{\}$,
2. Przekształć gramatykę generującą palindromy (nad alfabetem $\{a, b\}$), na automat stosowy (stosując konstrukcję podaną w tym wykładzie).
3. Przekształć automat stosowy akceptujący wyrażenia nawiasowe:

$$\begin{aligned} (s, x) &\xrightarrow{\downarrow} (s,]x) \text{ dla } x =], \perp \\ (s,]) &\xrightarrow{\downarrow} (s, \varepsilon) \\ (s, \perp) &\xrightarrow{\varepsilon} (s, \varepsilon) \end{aligned}$$

na równoważną mu gramatykę bezkontekstową (stosując konstrukcję podaną w tym wykładzie).

13.7 Ćwiczenia

Skonstruuj automaty stosowe akceptujące poniższe języki. Czy potrafisz skonstruować takie automaty, które mają tylko jeden stan? Jeśli tak, to przekształć je na równoważne im gramatyki bezkontekstowe.

1. wyrażenia nawiasowe zawierające trzy rodzaje nawiasów: $()\{\}\{\}$; przy tym: nawiasy okrągłe mogą otaczać tylko nawiasy okrągłe, nawiasy kwadratowe mogą otaczać nawiasy kwadratowe lub okrągłe, a nawiasy wąsate mogą otaczać nawiasy wąsate lub kwadratowe,
2. palindromy nad alfabetem $\{a, b, c\}$,
3. $\{a^i b^{2i}\}$,
4. $\{a^i b^j a^j b^i\}$,
5. $\{w : \#_a(w) \geq 2\#_b(w)\}$,
6. $\{a^i b^j c^k : i \neq j \vee j \neq k\}$,
7. wyrażenia arytmetyczne w notacji polskiej (dla uproszczenia liczby mogą być tylko jednocyfrowe),
8. wyrażenia arytmetyczne w odwrotnej notacji polskiej (dla uproszczenia liczby mogą być tylko jednocyfrowe).

Wykład 14. Analiza składniowa

14.1 Wstęp

Analiza składniowa ma miejsce wówczas, gdy wczytujemy dane o pewnej określonej składni. Zadaniem analizatora składniowego (tzw. *parsera*) jest sprawdzenie, czy dane są poprawne składniowo i rozpoznanie struktury składniowej wczytywanych danych.

Analiza składniowa powstała i rozwinęła się w ramach prac nad budową kompilatorów. Wszak kompilator musi najpierw wczytać i zanalizować składnię wczytywanego programu. Jej zastosowania są jednak dużo szersze. Praktycznie analizę składniową można zastosować w każdym programie, który wczytuje dane posiadające bardziej złożoną składnię, na przykład w: przeglądarkach internetowych, edytorach, systemach składu tekstu, programach konwertujących, czy jakichkolwiek aplikacjach posiadających pliki konfiguracyjne o określonej składni.

Konstrukcja parsera, zwłaszcza w przypadku kompilatorów lub bardziej skomplikowanych języków, to skomplikowane zadanie. Dlatego też powstały narzędzia wspomagające tworzenie parserów. Są to tzw. generatory analizatorów składniowych. Przykładem takiego generatora jest *Bison*, którego poznamy bliżej w tym wykładzie. *Bison* jest następcą generatora *Yacc* (ang. *yet another compiler compiler* — jeszcze jeden kompilator kompilatorów). *Bison* jest przeznaczony do generowania parserów w C/C++. Jednak istnieje wiele generatorów parserów i w zasadzie dla każdego języka programowania można znaleźć generator przeznaczony dla danego języka.

Generator parserów na podstawie odpowiedniej specyfikacji generuje kod źródłowy analizatora składniowego. Specyfikacja taka zawiera opis składni wczytywanego języka oraz fragmenty kodu, które są wykonywane dla poszczególnych konstrukcji składniowych. Do opisu składni używa się jednoznacznych gramatyk bezkontekstowych.

Analiza składniowa stanowi kolejną fazę przetwarzania wczytywanych danych po analizie leksykalnej. Tak więc wejściem dla analizy składniowej jest strumień żetonów i atrybutów. Żetony stanowią symbole terminalne w gramatyce bezkontekstowej opisującej składnię. Efektem pracy parsera jest odtworzenie drzewa wyprowadzenia (a ściślej mówiąc jego obejścia) dla wczytywanego tekstu.

Generowane automatycznie parsery mają postać odpowiednich, deterministycznych automatów stosowych. Sposób generowania parserów wykracza poza ramy tego kursu.

14.2 Gramatyki S-atrybutywne

Zanim poznamy język specyfikacji dla *Bison*'a i *Yacc*'a, musimy zapoznać się z pojęciem gramatyk S-atrybutywnych. (Nie jest to nic skomplikowanego, a ta odstraszaająca nazwa pochodzi stąd, że jest wiele różnych gramatyk atrybutywnych. My jednak nie będziemy ich tutaj omawiać.)

Gramatyki S-atrybutywne to rozszerzenie gramatyk bezkontekstowych — w takim sensie, że podstawowym składnikiem gramatyki S-atrybutywnej jest jednoznaczna gramatyka bezkontekstowa. (Gramatyka musi być jednoznaczna, żeby dla każdego słowa z języka

istniało tylko jedno możliwe drzewo wyprowadzenia.)

Natomiast dodatkowo w gramatyce S-atrybutywnej z każdym węzłem w drzewie wyprowadzenia wiążemy pewną obliczaną wartość nazywaną *atrybutem*. (Jeśli potrzebujemy wielu wartości, to możemy pomyśleć o atrybucie jak o n -ce, rekordzie czy strukturze złożonej z tych wartości.) Liśćmi drzewa wyprowadzenia są symbole terminalne, czyli żetony. Żetonom mogą towarzyszyć atrybuty, jeśli zostały wyznaczone przez skaner. Natomiast atrybuty towarzyszące nieterminalom, czyli węzłom wewnętrznym, są obliczane w trakcie analizy leksykalnej.

Atrybuty nieterminali obliczamy na podstawie atrybutów ich synów w drzewie wyprowadzenia. (Atrybuty takie są nazywane *syntezowanymi* — stąd „S” w nazwie tego rodzaju gramatyk.) Dla każdej produkcji podajemy w specyfikacji fragment kodu, który oblicza atrybut ojca na podstawie atrybutów synów. W kontekście produkcji postaci $X \rightarrow B_1 B_2 \dots B_k$, atrybut X -a oznaczamy przez \$\$, a atrybut B_i przez \$ i .

Parser wywołuje te fragmenty kodu obliczając atrybuty w odpowiedniej kolejności, od liści do korzenia (ang. *bottom-up*). Wynikiem całości obliczeń jest atrybut korzenia drzewa wyprowadzenia.

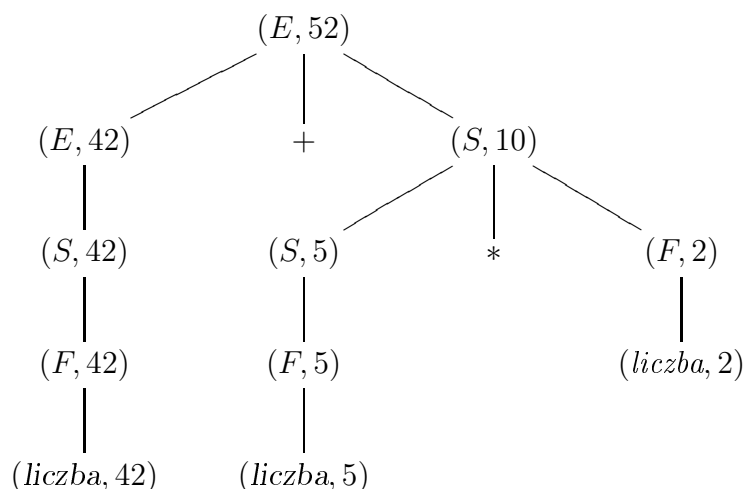
Przykład: Przypomnijmy sobie jednoznaczłą gramatykę wyrażeń arytmetycznych:

$$\begin{aligned} E &\rightarrow E + S \mid E - S \mid S \\ S &\rightarrow S * F \mid S / F \mid F \\ F &\rightarrow \text{liczba} \mid (E) \end{aligned}$$

W tym przypadku liczba jest dla nas terminalem, gdyż jest to leksem rozpoznawany przez skaner. Oto gramatyka S-atrybutywna obliczająca wartość wyrażenia:

$$\begin{aligned} E &\rightarrow E + S & \$\$ = \$1 + \$3 \\ E &\rightarrow E - S & \$\$ = \$1 - \$3 \\ E &\rightarrow S & \$\$ = \$1 \\ S &\rightarrow S * F & \$\$ = \$1 * \$3 \\ S &\rightarrow S / F & \$\$ = \$1 / \$3 \\ S &\rightarrow F & \$\$ = \$1 \\ F &\rightarrow \text{liczba} & \$\$ = \$1 \\ F &\rightarrow (E) & \$\$ = \$2 \end{aligned}$$

A oto drzewo wyprowadzenia dla wyrażenia $42 + 5 * 2$ wzbogacone o obliczone atrybuty



Tak więc obliczona wartość całego wyrażenia to 52.

14.3 Sposób działania parsera

Bison został zaprojektowany tak, aby współpracował razem z Flex'em. Kod wygenerowany przez Bison'a na podstawie specyfikacji parsera może być kompilowany razem z kodem wygenerowanym przez Flex'a na podstawie specyfikacji skanera. W szczególności żetony są definiowane w specyfikacji parsera, a ich definicje w języku programowania są umieszczane w kodzie parsera. Natomiast w kodzie skanera można z nich korzystać.

Parser wywołuje procedurę `yylex` żeby wczytywać kolejne leksemy. Wynikiem tej procedury powinien być żeton odpowiadający wczytanemu leksemowi. Zakłada się, że atrybut żetonu, jeśli jest obecny, to jest przypisany zmiennej `yylval`. Procedura `yylex` może również zwrócić znak ASCII — wówczas leksemem jest pojedynczy wczytany znak.

Parser udostępnia funkcję `int yyparse()`. Wywołanie tej funkcji powoduje próbę wczytania słowa należącego do języka opisywanego przez gramatykę podaną w specyfikacji. Parser nie konstruuje drzewa wyprowadzenia jako struktury danych, ale wykonuje takie czynności, jakby obchodził je w porządku postfiksowym. Dla każdego odwiedzanego nieterminala w drzewie wyprowadzenia wykonywany jest fragment kodu odpowiadający zastosowanej dla niego produkcji. Jako ostatni wykonywany jest fragment kodu dla korzenia. Jeżeli w tym fragmencie kodu nie nastąpi powrót z procedury (`return`), to wejście ponownie jest parsowane. W przypadku wystąpienia błędu składniowego wywoływana jest procedura `void yyerror (char const *s)`. Procedura ta musi być zdefiniowana.

Fragmenty kodu umieszczone w specyfikacji nie muszą tylko obliczać atrybuty. Mogą wykonywać dowolne czynności, np. wypisywać informacje na wyjście czy modyfikować globalne struktury danych. Trzeba jednak zdawać sobie sprawę, w jakiej kolejności różne fragmenty kodu są wykonywane.

14.4 Specyfikacje parserów dla Bison'a

Specyfikacja dla Bison'a składa się z trzech części, oddzielonych od siebie wierszami postaci „%%”:

- Pierwsza część zawiera deklaracje programistyczne (stałych i zmiennych) oraz deklaracje żetonów (postaci %token ŻETON).
- Druga część zawiera gramatykę S-atrybutywną. Gramatyka ta jest zapisana w następujący sposób:

```
nieterminal : prawa strona produkcji {fragment kodu }  
              | prawa strona produkcji {fragment kodu }  
              :  
              | prawa strona produkcji {fragment kodu }  
              ;
```

Czyli produkcje dla jednego nieterminala są zgrupowane razem, a dla każdej z tych produkcji mamy oddzielny fragment kodu ujęty w wąsate nawiasy. Fragmenty kodu są opcjonalne. Jeżeli dany fragment kodu jest pominięty, to tak jakby zostało podane: { \$\$ = \$1 }. Aksjomatem jest ten nieterminal, dla którego produkcje są podane jako pierwsze.

- Trzecia część to definicje procedur pomocniczych. W szczególności można tu umieścić:
 - definicję procedury `yylex`, jeżeli nie korzystamy ze skanera wygenerowanego przez Flex'a,
 - definicję procedury `main`,
 - powinna tu zostać zdefiniowana procedura `yerror`, która jest wywoływana przez parser w przypadku wystąpienia błędów składniowych.

Przykład: Zbudujemy prosty kalkulator wzbogacony o zmienne. Będzie on przyjmował dwa rodzaje poleceń:

```
SET identyfikator = wyrażenie  
PRINT wyrażenie
```

Polecenie SET powoduje przypisanie wartości wyrażenia na zmienną. Polecenie PRINT powoduje wypisanie wartości wyrażenia. Każde polecenie musi znajdować się w osobnym wierszu. Wyrażenia mogą zawierać nazwy zmiennych. Nazwy zmiennych mają postać $[A - Za - z][A - Za - z0 - 9]^*$. Poza tym, składnia wyrażeń jest taka jak w poprzednich przykładach. Początkowo wszystkie zmienne mają wartość 0. W danych może pojawić się co najwyżej 1000 różnych nazw zmiennych.

Skaner odróżnia słowa kluczowe SET i PRINT, liczby i identyfikatory. Wszystkie inne znaki są przekazywane dosłownie. Kod źródłowy specyfikacji skanera znajduje się w pliku calc.l.

```
%{
// Definicje żetonów pochodzące z Bison'a
#include "skalk.tab.h"

// Tablica zmiennych
char *dict[1000];
int dict_size=0;

// Odnajduje nazwę zmiennej w tablicy, lub ją dodaje
int dict_find(const char *key) {
    int i;
    for(i=0; i<dict_size; i++) if (strcmp(key,dict[i])==0) return i;

    i=dict_size; dict_size++;
    dict[i]=(char *)malloc(strlen(key)+1);
    strncpy(dict[i],key,strlen(key)+1);
    return i;
}
%}

INT      [0-9]+
ID       [A-Za-z][A-Za-z0-9]*

%%

PRINT { return CMD_PRINT; }
SET    { return CMD_SET;   }
{ID}   {
        yynval=dict_find(ytext);
        return ID;
      }
{INT}  {
        yynval=atoi(ytext);
        return NUM;
      }
[ \t] ;
[\n] |. { return ytext[0]; }

%%
```


Parser pamięta w tablicy wartości zmiennych, oblicza wartości wyrażeń i przypisuje zmiennym wartości. Kod źródłowy specyfikacji parsera znajduje się w pliku `calc.y`.

```
%{
#include <stdio.h>

// Tablica wartości zmiennych
int d_value[1000];

// Odszukanie wartości zmiennej
int dict_value(int key) {
    return d_value[key];
}

// Przypisanie wartości zmiennej
int dict_set(int key,int value) {
    d_value[key]=value;
}

%}

%token NUM CMD_PRINT CMD_SET ID

%%

input: /* nic */
      | input line
      ;

line: '\n'
     | CMD_SET ID '=' exp '\n' { dict_set($2,$4); }
     | CMD_PRINT exp '\n' { printf("%d\n", $2); }
     ;

exp : exp '+' sk1 { $$ = $1 + $3 }
    | exp '-' sk1 { $$ = $1 - $3 }
    | sk1
    ;

sk1 : sk1 '*' czy { $$ = $1 * $3 }
    | sk1 '/' czy { $$ = $1 / $3 }
    | czy
    ;
```

```

czy : NUM
    | ID          { $$ = dict_value($1); }
    | '(' exp ')' { $$ = $2 }
    ;

%%

main()
{
    // Wyzerowanie wartości zmiennych
    bzero(&d_value, sizeof(d_value));

    yyparse();
}

yyerror(s)
char *s;
{
    fprintf(stderr, "%s\n", s);
}

```

14.5 Podsumowanie

W tym wykładzie poznaliśmy generator analizatorów składniowych Bison oraz związane z nim pojęcie gramatyk S-atrybutywnych. Umiejętność korzystania z Bison'a będzie jeszcze przedmiotem ćwiczeń laboratoryjnych.

14.6 Skorowidz

- **Analiza składniowa** polega na sprawdzeniu poprawności składniowej i odtworzeniu (obejścia) drzewa wyprowadzenia dla wczytywanych danych.
- **Bison** to jeden z generatorów analizatorów składniowych.
- **Generator analizatorów składniowych** to program generujący, na podstawie specyfikacji, kod analizatora składniowego (parsera).
- **Gramatyki S-atrybutywne** to gramatyki atrybutywne, w których wszystkie atrybuty są syntezowane, tzn. węzła w drzewie wyprowadzenia jest obliczany na podstawie atrybutów jego synów.
- **Parser** to inaczej analizator składniowy, moduł zajmujący się analizą składniową.

14.7 Praca domowa

Odwrotna notacja polska (ONP) to sposób zapisu wyrażeń, w którym najpierw podajemy argumenty, a potem operację. Jeżeli wiadomo ile argumentów mają operacje (a w przypadku operacji arytmetycznych tak jest — mają one po dwa argumenty), to w ONP nie są potrzebne nawiasy. Na przykład, wyrażenie $2*(3+5)$ zapisane w ONP ma postać $2\ 3\ 5\ +\ *$. Nota bene, wyrażenie zapisane w ONP to gotowy program dla maszyny stosowej, patrz Ćwiczenia.

Napisz specyfikację (dla Lex'a i Bison'a) analizatora, który wczyta wyrażenie arytmetyczne i wypisze je w ONP.

14.8 Ćwiczenia

1. Napisz specyfikację (dla Lex'a i Bison'a) analizatora, który wczyta wyrażenie arytmetyczne i wyznaczy wielkość stosu potrzebnego do obliczenia wyrażenia przez maszynę stosową.

Maszyna stosowa działa na następującej zasadzie: włóż liczbę na stos, zdejmij argumenty operacji arytmetycznej z wierzchołka stosu i włóż jej wynik na stos. Na przykład, obliczenie wyrażenia $2 + 3 * 5$ wymaga stosu wielkości 3. Obliczenie to przebiega następująco: włóż 2 na stos, włóż 3 na stos, włóż 5 na stos, wykonaj mnożenie, wykonaj dodawanie. Obliczenie wyrażenia $3 * 5 + 2$ wymaga stosu wielkości 2. Obliczenie to przebiega następująco: włóż 3 na stos, włóż 5 na stos, wykonaj mnożenie, włóż 2 na stos, wykonaj dodawanie.

2. Napisz specyfikację (dla Lex'a i Bison'a) analizatora, który wczyta wyrażenie arytmetyczne i wyznaczy maksymalną głębokość wyrażenia. Na przykład dla $3 + 3 * (5 - 1)$ wynikiem powinno być 3, a dla $(2 - 1) + (1 - 2)$ wynikiem powinno być 2.
3. Rozszerz przykład kalkulatora z wykładu o możliwość definiowania nazwanych wyrażeń. Wartość takiego wyrażenia zależy od wartości zmiennych i jest obliczana za każdym razem, gdy takie wyrażenie jest używane.
4. Rozważamy następujący rodzaj kompresji słów nad alfabetem $\{0, 1\}$. Mając dane słowo s (o długości będącej potęgą dwójki, $|s| = 2^i$), jeśli:

- s składa się z samych znaków 0, $s = 0^{2^i}$, to jego kodem jest 0,
- s składa się z samych znaków 1, $s = 1^{2^i}$, to jego kodem jest 1,
- w przeciwnym przypadku napis s jest dzielony na równe części s_1 i s_2 długości 2^{i-1} , $s = s_1 s_2$, a jego kod jest postaci: $9kod(s_1)kod(s_2)$.

Na przykład:

- $kod(1) = 1$,
- $kod(0111) = 99011$,

- $kod(01001111) = 9990101$,
- $kod(00110101) = 99019901901$.

Napisz specyfikacje dla Lex'a i Bison'a dające program, który: wczyta zestaw zakodowanych słów, po jednym w wierszu i dla każdego kodu wypisze jego odkodowaną postać.

Wykład 15. Maszyny Turinga i obliczalność

15.1 Wstęp

W tym wykładzie zajmiemy się pojęciem obliczalności. Poznamy maszyny Turinga, model teoretyczny opisujący pojęcie obliczalności, równy pod względem siły wyrazu wszelkiego rodzaju językom programowania i komputerom.

Pojęcie obliczalności było badane zanim narodziła się informatyka, a nawet zanim powstały pierwsze komputery, bo było to już w latach 30-tych ubiegłego wieku. W ramach badań nad formalizacją podstaw matematyki, badano jakie zbiory formuł można uzyskać na drodze mechanicznego przepisywania napisów zgodnego z określonymi zasadami. Prowadziło to do sformalizowania intuicji dotyczącej tego co można „automatycznie obliczyć”. Pojawiło się wiele modeli formalizujących takie intuicyjne pojęcie obliczalności. Wśród nich można wymienić maszyny Turinga, rachunek λ , czy systemy Posta. W tym wykładzie zajmiemy się maszynami Turinga, gdyż spośród tych modeli są one najbliższe komputerom.

Jak się okazało, wszystkie te formalizmy są sobie równoważne — w każdym z nich można symulować każdy z pozostałych formalizmów. Zostało to ujęte w tzw. hipotezie Church’a, mówiącej, że wszystkie te modele opisują to samo intuicyjne pojęcie *obliczalności*. Oczywiście nie jest to formalne twierdzenie, które można by udowodnić, gdyż dotyczy ona intuicji, ale jak dotąd nic nie podważyło słuszności tej hipotezy.

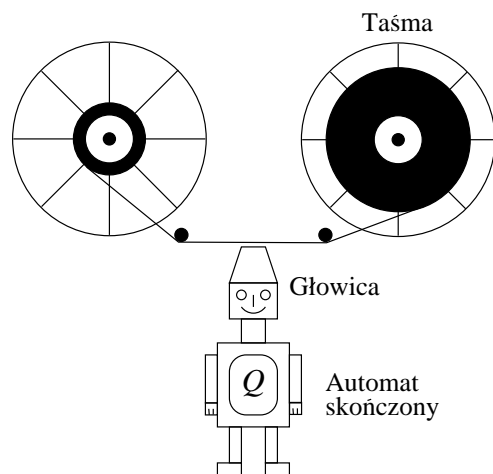
15.2 Uniwersalność

Formalizmy modelujące pojęcie obliczalności nie tylko są sobie nawzajem równoważne, ale też w każdym z nich można stworzyć model *uniwersalny* — taki, który potrafi symulować wszystkie inne modele danego rodzaju (na podstawie ich opisu). W szczególności istnieje uniwersalna maszyna Turinga, która potrafi symulować działanie dowolnej innej maszyny Turinga na podstawie jej opisu.

W pierwszym momencie, pojęcie uniwersalności może być niezrozumiałe lub zadziwiająca. Natomiast dzisiaj jest ono prawie oczywiste. Zamiast mówić o maszynach Turinga wybierzmy dowolny język programowania. Uniwersalny program to nic innego jak interpreter danego języka programowania napisany w nim samym. (Jeśli czytelnik zna język Scheme, to zapewne zetknął się z interpreterem Scheme’a napisanym w Scheme’ie.) Interpretery może nie są tak popularne jak kompilatory. Natomiast np. kompilator C++ napisany w C++ jest czymś naturalnym, a w pewnym sensie jest to uniwersalny program w C++. Za jego pomocą można uruchomić każdy inny program w C++.

15.3 Maszyny Turinga

Maszyna Turinga przypomina skrzyżowanie automatu skończonego z magnetofonem szpulowym. Tak jak automat stosowy to automat skończony wyposażony w dodatkową pamięć w postaci stosu, tak maszyna Turinga to automat skończony wyposażony w dodatkową pamięć w postaci nieskończonej taśmy, na której można zapisywać i odczytywać informacje.



Taśma ta jest podzielona na klatki. W każdej klatce można zapisać jeden symbol z ustalonego alfabetu. Nad taśmą przesuwa się głowica, którą steruje automat skończony. W jednym kroku automat odczytuje znak zapisany pod głowicą, może w tym miejscu zapisać inny znak, po czym może przesunąć głowicę w lewo lub prawo. Oczywiście automat sterujący głowicą w każdym kroku zmienia również swój stan.

Taśma ma początek, lecz nie ma końca — jest nieskończona. W pierwszej klatce taśmy jest zapisany specjalny znak, tzw. *lewy ogranicznik*. Jeżeli głowica znajduje się nad lewym ogranicznikiem, to nie może go zamazać ani przesunąć się na lewo od niego. Na początku, zaraz za lewym ogranicznikiem, zapisane jest słowo, które stanowi dane wejściowe dla maszyny Turinga. Oczywiście słowo to jest skończone. Za tym słowem taśma wypełniona jest w nieskończoność specjalnymi pustymi symbolami, tzw. *blank'ami*. Blank'i będziemy oznaczać przez \sqcup .

Będziemy tutaj rozważać jedynie deterministyczne maszyny Turinga, tzn. takie, w których automat sterujący głowicą jest deterministyczny. W momencie uruchomienia maszyny Turinga, głowica znajduje się nad pierwszą klatką taśmy (nad lewym ogranicznikiem) i jest w stanie początkowym. Maszyna ma dwa wyróżnione stany: akceptujący i odrzucający. Działa ona tak długo, aż znajdzie się w jednym z tych dwóch stanów. Jeśli jest to stan akceptujący, to słowo początkowo zapisane na taśmie zostało zaakceptowane, a w przeciwnym przypadku zostało odrzucone. Oczywiście jest możliwe, że maszyna Turinga nigdy nie znajdzie się w żadnym z tych stanów. Wówczas jej obliczenie trwa w nieskończoność i mówimy, że się *zapętliło*. Tak więc dla danego słowa maszyna Turinga może zrobić jedną z trzech rzeczy: zaakceptować to słowo, odrzucić je lub zapętlić się.

Nieznacznie zmieniając opisany model maszyny Turinga możemy go użyć do modelowania dowolnych obliczeń, a nie tylko akceptowania języków, czy problemów decyzyjnych. W tym celu wystarczy, że zamiast dwóch stanów: akceptującego i odrzucającego, mamy jeden stan końcowy. Wówczas to co jest zapisane na taśmie (z wyjątkiem lewego ogranicznika i blank'ów) w momencie gdy maszyna osiągnie stan końcowy, stanowi wynik obliczeń.

Def. 46. Maszyna Turinga to dowolna taka dziewiątka $M = \langle Q, \Sigma, \Gamma, \vdash, \sqcup, \delta, s, t, r \rangle$, że:

- Q , to skończony zbiór stanów,

- $\Sigma \subseteq \Gamma$, to alfabet wejściowy,
- Γ , to alfabet taśmowy,
- $\vdash \in \Gamma$, to lewy ogranicznik,
- $\sqcup \in \Gamma$, to symbol pusty, blank,
- $\delta : (Q \setminus \{t, r\}) \times \Gamma \rightarrow Q \times \Gamma \times \{L, P\}$, to funkcja przejścia; wymagamy, aby dla dowolnego $q \in Q$ istniał taki $q' \in Q$, że $\delta(q, \vdash) = (q', \vdash, P)$, czyli żeby \vdash nie był zamazywany, a głowica nie mogła przesunąć się na lewo od niego
- $s \in Q$, to stan początkowy,
- $t \in Q$, to stan akceptujący,
- $r \in Q, r \neq t$, to stan odrzucający.

□

Def. 47. Niech $M = \langle Q, \Sigma, \Gamma, \vdash, \sqcup, \delta, s, t, r \rangle$ będzie ustaloną maszyną Turinga. *Konfiguracja* maszyny M , to dowolna taka trójka:

$$\langle q, \vdash \alpha, k \rangle \in Q \times \Gamma^* \times \mathcal{N}$$

że $|\alpha| \geq k$. W konfiguracji $\langle q, \vdash \alpha, k \rangle$ q to stan maszyny, $\vdash \alpha$ reprezentuje zawartość taśmy (uzupełnioną w nieskończoność blank'ami), a k to pozycja głowicy.

Konfiguracja początkowa maszyny M dla słowa $x \in \Sigma^*$ to $\langle s, \vdash x, 0 \rangle$. □

Def. 48. Niech $M = \langle Q, \Sigma, \Gamma, \vdash, \sqcup, \delta, s, t, r \rangle$ będzie daną maszyną Turinga. Na konfiguracjach maszyny Turinga definiujemy relację $\rightarrow \subseteq (Q \times \Gamma^* \times \mathcal{N})^2$. Relacja ta opisuje przejścia między konfiguracjami odpowiadające pojedynczym krokom w obliczeniach M . Jest to najmniejsza relacja, która spełnia podane poniżej warunki.

Niech $\langle q, \vdash \alpha, k \rangle$ będzie konfiguracją i niech $q \neq t, q \neq r, \vdash \alpha = a_0 a_1 a_2 \dots a_n, n \geq k$. Niech $\delta(q, a_k) = (\langle \rangle q', a', x)$. Jeśli $x = L$ (czyli głowica przesuwa się w lewo), to mamy:

$$\langle q, \vdash \alpha, k \rangle \rightarrow \langle q', a_0 a_1 \dots a_{k-1} a' a_{k+1} \dots a_n, k - 1 \rangle$$

Jeśli $x = P$ i $k < n$ (czyli głowica przesuwa się w prawo, ale w obrębie $\vdash \alpha$), to mamy:

$$\langle q, \vdash \alpha, k \rangle \rightarrow \langle q', a_0 a_1 \dots a_{k-1} a' a_{k+1} \dots a_n, k + 1 \rangle$$

Jeśli zaś $x = P$ i $k = n$ (czyli głowica jest na prawym końcu $\vdash \alpha$ i przesuwa się w prawo), to mamy:

$$\langle q, \vdash \alpha, k \rangle \rightarrow \langle q', a_0 a_1 \dots a_{k-1} a' a_{k+1} \dots a_n \sqcup, k + 1 \rangle$$

Przez \rightarrow^* oznaczamy domknięcie zwrotnie przechodnie relacji \rightarrow . □

Relacja \rightarrow^* opisuje do jakich konfiguracji możemy dojść w wyniku obliczenia. Konfiguracje akceptujące, to te, które zawierają stan t , a odrzucające to te, które zawierają stan r .

Def. 49. Niech $M = \langle Q, \Sigma, \Gamma, \vdash, \sqcup, \delta, s, t, r \rangle$ będzie daną maszyną Turinga. *Język akceptowany* przez M składa się z tych słów, dla których obliczenie maszyny prowadzi do konfiguracji akceptującej:

$$L(M) = \{x \in \Sigma^* : \exists_{\alpha \in \Gamma^*, k \in \mathcal{N}} \langle s, \vdash x, 0 \rangle \rightarrow^* \langle t, \vdash \alpha, k \rangle\}$$

□

15.4 Języki częściowo obliczalne i obliczalne

Def. 50. Powiemy, że język $A \subseteq \Sigma^*$ jest *częściowo obliczalny*, jeżeli istnieje taka maszyna Turinga M , że $A = L(M)$. □

Wobec języków częściowo obliczalnych używa się też określeń częściowo rozstrzygalny i rekurencyjnie przeliczalny. Wszystkie te określenia oznaczają to samo pojęcie.

Zauważmy, że żeby maszyna Turinga M akceptowała język A , dla wszystkich słów z tego języka musi zatrzymywać się w stanie akceptującym, natomiast dla słów spoza języka A nie musi zatrzymywać się w stanie odrzucającym — może również się zapętlać. Intuicyjnie odpowiada to istnieniu programu komputerowego, który dla słów z języka A potrafi potwierdzić ich przynależność do języka, natomiast dla słów spoza tego języka wcale nie musi potrafić zaprzeczyć ich przynależności do języka, lecz może się zapętlić. Stąd słowo „częściowo” w nazwie.

Istnieje równoważna definicja języków częściowo obliczalnych, która mówi, że język A jest częściowo obliczalny, gdy istnieje taki program komputerowy, który wypisuje (w pewnej kolejności) wszystkie słowa należące do A . Jeżeli A jest nieskończony, to oznacza to, że każde słowo należące do A kiedyś zostanie wypisane.

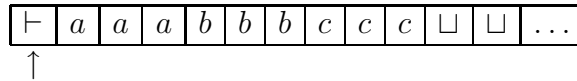
Def. 51. Powiemy, że język $A \subseteq \Sigma^*$ jest *obliczalny*, jeżeli istnieje taka maszyna Turinga M , że $A = L(M)$ i dla każdego $x \in \Sigma^*$ maszyna M albo akceptuje x , albo go odrzuca. □

Fakt 17. *Każdy język bezkontekstowy jest obliczalny.*

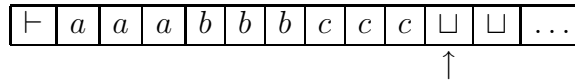
Fakt ten wynika stąd, że algorytm CYK (dla danej gramatyki bezkontekstowej) można zaimplementować w postaci maszyny Turinga. Ponieważ algorytm ten zawsze się zatrzymuje, więc i taka maszyna nie będzie się zapętlać.

Przykład: Pokazaliśmy poprzednio, że język $\{a^n b^n c^n : n \geq 0\}$ nie jest bezkontekstowy. Zapewne dla większości Czytelników napisanie programu, który wczytuje napis i sprawdza czy jest on postaci $a^n b^n c^n$ nie sprawiłoby kłopotu. Pokażemy, jak skonstruować maszynę

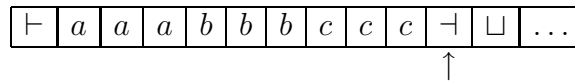
Turinga, która to robi. Nie podamy jej formalnej definicji (gdyż byłoby to nudne), ale opiszemy jej sposób działania. Oto przykładowa początkowa konfiguracja:



1. Nasza maszyna najpierw czyta słowo podane na wejściu i sprawdza czy pasuje ono do wzorca $a^*b^*c^*$, przesuwając głowicę w prawo, aż do napotkania pierwszego blanka. Jeżeli słowo na wejściu nie pasuje do wzorca, to odrzucamy je. Dodatkowo, jeżeli na wejściu jest dane słowo puste, to od razu je akceptujemy.



2. W miejsce pierwszego blank'a wstawiamy prawy ogranicznik.



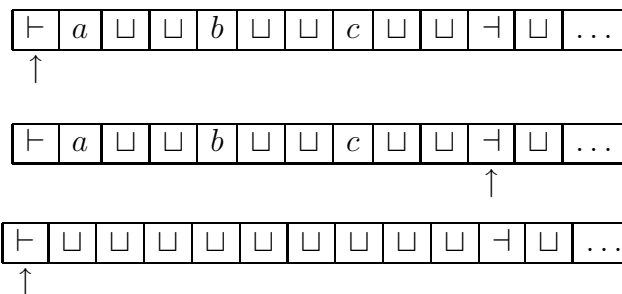
3. Następnie przejeżdżamy głowicą w lewo zamieniając pierwszy napotkany znak a , b i c na blanki. Jeżeli któregoś ze znaków zabrakło, to odrzucamy. Oznacza to, że znaków a , b i c nie było po równo.



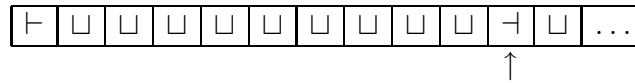
4. Przesuwamy głowicę w prawo, aż do prawego ogranicznika, sprawdzając czy między ogranicznikami są jakieś znaki inne niż blanki.



W takim przypadku skaczemy do kroku 3 i powtarzamy dwa ostatnie kroki tak długo, aż między ogranicznikami pozostaną same blank'i.



5. Gdy w końcu między ogranicznikami zostaną same blank'i, akceptujemy.



[Prezentacja multimedialna: animacja powyższej maszyny.]

Pokazaliśmy właśnie, że język $\{a^n b^n c^n : n \geq 0\}$ jest obliczalny, choć wiemy, że nie jest bezkontekstowy. Wynika stąd następujący fakt:

Fakt 18. *Klasa języków obliczalnych zawiera ściśle klasę języków bezkontekstowych.*

Czy każdy język jest częściowo obliczalny? Zdecydowanie nie! Wynika to stąd, że (dla ustalonego alfabetu wejściowego Σ) maszyn Turinga (z dokładnością do izomorfizmu) jest przeliczalnie wiele. Być może bardziej intuicyjne jest równoważne stwierdzenie, że w dowolnym języku programowania istnieje przeliczalnie wiele programów. Wszak każdy program to tylko słowo, a słów (nad ustalonym alfabetem Σ) jest przeliczalnie wiele. Natomiast języki to zbiory słów, czyli wszystkich możliwych języków (nad ustalonym alfabetem Σ) jest $2^{|\mathcal{M}|} > |\mathcal{M}|$. Z drugiej strony, każda maszyna Turinga czy program komputerowy akceptuje jeden określony język. Tak więc maszyn Turinga czy programów komputerowych jest **za mało**, żeby każdy język był częściowo obliczalny.

Powyższe rozumowanie jest poprawne, ale niekonstruktywne — nie dostarcza nam żadnego przykładu języka, który nie byłby częściowo obliczalny czy obliczalny. Dalej poznamy przykład takiego języka.

15.5 Wariacje na temat maszyn Turinga

W literaturze można spotkać różne warianty maszyn Turinga: maszyny z taśmą nieskończoną w obydwie strony, maszyny w taśmą wielościeżkową, maszyny z wieloma taśmami i głowicami, czy maszyny z dwuwymiarową płaszczyzną zamiast taśmy. Dalej w naszych rozważaniach czasem będzie nam wygodnie założyć, że maszyny Turinga posiadają któreś z tych rozszerzeń. Wszystkie te modle są równoważne przedstawionym tutaj maszynom Turinga. Nie będziemy tego formalnie dowodzić, ale krótko naszkicujemy odpowiednie konstrukcje.

Powiedzmy, że chcielibyśmy mieć maszynę wyposażoną w k ścieżek, przy czym na i -tej ścieżce są zapisane znaki z Γ_i . Głowica będąc w jednym położeniu może odczytać i zapisać równocześnie znaki na wszystkich ścieżkach. Tak naprawdę jest to równoważne klasycznej maszynie Turinga z alfabetem ścieżkowym $\Gamma_1 \times \Gamma_1 \times \dots \times \Gamma_k$.

Jeżeli chcemy, aby taśma była nieskończona w obie strony, to wystarczy, że złożymy ją (w miejscu rozpoczęcia wejścia) na pół i sklejmy w taśmę dwuścieżkową. Taka maszyna będzie w każdym kroku uwzględniać i modyfikować tylko jedną z dwóch ścieżek. Dodatkowo więc automat sterujący głowicą musi pamiętać nad którą ścieżką znajduje się głowica.

Jeśli chcielibyśmy mieć wiele taśm i głowic, to symulujemy to na maszynie z jedną głowicą i wieloma ścieżkami. Potrzebujemy taśmy o dwukrotnie większej liczbie ścieżek.

Połowa z nich będzie nam służyć do przechowywania zawartości taśmy symulowanej maszyny, a połowa będzie przeznaczona na znaczniki wskazujące położenia odpowiednich głowic. Oczywiście symulacja jednego kroku takiej maszyny może wymagać przejrzenia większego fragmentu taśmy i trwać odpowiednio długo. Niemniej jest to wykonalne.

15.6 Podsumowanie

W tym wykładzie poznaliśmy maszyny Turinga — model obliczeniowy równoważny programom komputerowym. Poznaliśmy też definicje klas języków obliczalnych i częściowo obliczalnych, które zbadamy bliżej w kolejnym wykładzie.

15.7 Skorowidz

- **Częściowa obliczalność** — język jest obliczalny, gdy istnieje maszyna Turinga akceptująca go. Nie musi ona zatrzymywać się dla słów spoza języka.
- **Maszyna Turinga** to teoretyczny model oddający pojęcie obliczalności.
- **Obliczalność** — język jest obliczalny, gdy istnieje taka maszyna Turinga akceptująca go, która zawsze się zatrzymuje.
- **Uniwersalna maszyna Turinga** to maszyna Turinga, która może symulować każdą inną maszynę Turinga na podstawie jej opisu oraz przetwarzanych przez nią danych.

15.8 Praca domowa

1. Opisz maszynę Turinga akceptującą język $\{ww : w \in \{a, b\}^*\}$. Maszyna ta nie może zapętlać się. Zasymuluj jej działanie dla słowa *abbabb*.

Wykład 16. Języki obliczalne, częściowo obliczalne i nie-obliczalne

W poprzednim wykładzie poznaliśmy definicje języków obliczalnych i częściowo obliczalnych. Jak dotąd, wszystkie przykłady języków, jakie rozważaliśmy były obliczalne. Poznamy teraz przykłady języków, dla których problem przynależności słów do języka nie jest obliczalny, a nawet nie jest częściowo obliczalny.

Poznamy też hierarchię Chomsky’ego, hierarchię klas języków, która podsumowuje większość poznanych przez nas klas języków i formalizmów.

16.1 Problem stopu

Mówiliśmy wcześniej o uniwersalnej maszynie Turinga. Sprecyzujmy jej postać. Ustalamy alfabet wejściowy $\Sigma = \{0, 1, \$\}$. Uniwersalna maszyna Turinga symuluje działanie danej innej maszyny Turinga na podstawie jej opisu oraz jej wejścia. Zarówno opis symulowanej maszyny, jak i jej wejście maszyna uniwersalna otrzymuje na swoim wejściu w postaci ciągów zer i jedynek. Wymaga to zakodowania opisu symulowanej maszyny Turinga jako (ciągu bitów) słowa nad alfabetem $\{0, 1\}$. Nie opisujemy tutaj takiego kodowania, tylko zakładamy, że takie kodowanie jest ustalone. (Jest to problem czysto techniczny.) Podobnie, symulowana maszyna może operować na słowach nad innym alfabetem niż $\{0, 1\}$. Zakładamy że słowa te są, w miarę potrzeby, niejawnie kodowane — to że znaki dowolnego alfabetu można zakodować jako ciągi bitów odpowiedniej długości jest oczywiste. Uniwersalną maszynę Turinga będziemy oznaczać przez U . Przyjmujemy, że akceptuje ona następujący język:

$$L(U) = \{\text{kod } M\$x : x \in L(M)\}$$

Wynik działania maszyny uniwersalnej jest dokładnie taki sam, jak wynik działania symulowanej maszyny M dla danego słowa x .

Problem stopu to problem decyzyjny polegający na stwierdzeniu, czy dana maszyna Turinga M zatrzyma się dla danego słowa x . Problem ten możemy zdefiniować w formie języka:

$$STOP = \{\text{kod } M\$x : M \text{ zatrzymuje się dla słowa } x\}$$

Fakt 19. *Język STOP jest częściowo obliczalny.*

Dowód: Maszyna Turinga M_{STOP} akceptująca język $STOP$ to prosta przeróbka maszyny uniwersalnej. M_{STOP} działa tak samo jak U , z tą różnicą, że gdy maszyna symulowana odrzuca (a więc zatrzymuje się), to M_{STOP} akceptuje. \square

Maszyna M_{STOP} nigdy nie odrzuca, natomiast zapętla się gdy symulowana maszyna się zapętla.

Tw. 6. *Język STOP nie jest obliczalny.*

Dowód: Dowód przebiega nie wprost. Załóżmy, że istnieje maszyna Turinga M' , która akceptuje język $L(M') = STOP$ i zawsze się zatrzymuje. Niech maszyna M'' będzie wynikiem przeróbki M' : M'' dostaje na wejściu kod maszyny M i działa tak jak M' dla słowa kod M kod M , przy czym jeśli M' odrzuca, to M'' akceptuje, a jeśli M' akceptuje, to M'' zapętla się. Tak więc

$$L(M'') = \{\text{kod } M : M \text{ zapętla się dla kodu } M\}$$

Ponieważ M'' nie odrzuca, więc zatrzymuje się dla kodu M , wtw., gdy M nie zatrzymuje się dla własnego kodu. Czy M'' zatrzyma się dla własnego kodu? Jeśli się zatrzymuje i akceptuje, to tym samym powinna się zapętlić. I na odwrót, jeżeli się zapętla, to powinna zaakceptować swój własny kod. Sprzeczność. Tak więc założenie, że język $STOP$ jest obliczalny jest fałszywe. \square

16.2 Własności klasy języków obliczalnych i częściowo obliczalnych

Poznaliśmy przykład języka, który nie jest obliczalny. Dalej poznamy przykład języka, który nie jest nawet częściowo obliczalny. Wcześniej jednak poznamy trochę własności języków obliczalnych i częściowo obliczalnych.

Fakt 20. *Jeśli A i \bar{A} są częściowo obliczalne, to A jest obliczalny.*

Dowód: Powiedzmy, że mamy dwie maszyny Turinga, M_1 i M_2 , akceptujące odpowiednio A i \bar{A} . Możemy skonstruować maszynę M (z dwoma ścieżkami i głowicami), która będzie symulować równocześnie działanie M_1 i M_2 . Maszyna ta patrzy, która z maszyn M_1 i M_2 zaakceptuje słowo i odpowiednio akceptuje lub odrzuca. Ponieważ któraś z maszyn M_1 , M_2 musi zaakceptować, więc M nie zapętla się. \square

Fakt 21. *Język \overline{STOP} nie jest częściowo obliczalny.*

Dowód: Wiemy, że $STOP$ jest częściowo obliczalny. Gdyby \overline{STOP} był częściowo obliczalny, to $STOP$ byłby obliczalny, a wiemy, że taki nie jest. \square

16.3 Gramatyki ogólne (typu 0)

Jak dotąd, dla każdej poznawanej klasy języków mieliśmy dwa równoważne modele — jeden nastawiony na opisywanie składni i jeden na modelowanie obliczeń. Czy istnieje formalizm nastawiony na opisywanie składni równoważny maszynom Turinga? Tak. Są to Gramatyki ogólne (tzw. typu 0). Ich definicja jest analogiczna do definicji gramatyk bezkontekstowych, ale po obu stronach produkcji mogą stać słowa.

Def. 52. *Gramatyka ogólna (typu 0), to dowolna taka czwórka $G = \langle N, \Sigma, P, S \rangle$, gdzie:*

- N to (skończony) alfabet symboli *nieterminalnych*,
- Σ to (skończony) alfabet symboli *terminalnych*, $\Sigma \cap N = \emptyset$,
- P to skończony zbiór *produkcji*, reguł postaci $\alpha \rightarrow \beta$ dla $\alpha, \beta \in (N \cup \Sigma)^*$,
- $S \in N$ to *aksjomat* wyróżniony symbol nieterminalny.

□

Produkcja $\alpha \rightarrow \beta$ pozwala w trakcie wyprowadzenia zamienić dowolne podśłowo α na β . Definicje wyprowadzenia i języka generowanego są analogiczne do tych dla języków bezkontekstowych. Jednak siła wyrazu gramatyk ogólnych jest dużo większa.

Tw. 7. *Niech $A \subseteq \Sigma^*$ będzie językiem. Następujące dwa stwierdzenia są równoważne:*

- $A = L(M)$ dla pewnej maszyny Turinga M , czyli A jest częściowo obliczalny,
- $A = L(G)$ dla pewnej gramatyki ogólnej (typu 0).

Szkic dowodu: Nie będziemy tutaj szczegółowo dowodzić. Dowód przebiega przez pokazanie implikacji w obie strony.

- Mając daną gramatykę ogólną możemy skonstruować maszynę Turinga, która będzie konstruować wszystkie możliwe wyprowadzenia, w kolejności wg. rosnącej ich długości. Jeśli słowo dane na wejściu ma wyprowadzenie, to zostanie ono w końcu skonstruowane i maszyna wówczas zaakceptuje to słowo. Jeżeli nie, a wszystkich możliwych wyprowadzeń jest nieskończenie wiele, maszyna zapętli się. Niemniej będzie ona akceptować język generowany przez daną gramatykę.
- Mając daną maszynę Turinga możemy skonstruować gramatykę, która będzie symulować działanie maszyny wstecz. Słowo pojawiające się w wyprowadzeniu będzie reprezentować zawartość taśmy, z wstawionym w odpowiednim miejscu znacznikiem reprezentującym położenie głowicy i stan maszyny. Dla każdego możliwego ruchu maszyny mamy produkcję cofającą go. W tworzonej gramatyce możemy wyprowadzić z aksjomatu wszystkie słowa reprezentujące konfiguracje akceptujące (jest to język regularny), następnie możemy stosować produkcje cofające nas w obliczeniach. Na koniec, jeśli możemy dojść do konfiguracji początkowej maszyny, to mamy produkcje usuwające znaczniki i pozostawiające słowo dane na wejściu.

Tak więc możemy wyprowadzić dokładnie te słowa, dla których istnieją obliczenia akceptujące je. □

Zauważmy jednak, że język generowany przez gramatykę ogólną może być nierozstrzygalny. Jest tak dlatego, że dla każdego języka częściowo obliczalnego mamy generującą go gramatykę ogólną, ale nie każdy język częściowo obliczalny jest obliczalny.

16.4 Gramatyki kontekstowe

Gramatyki kontekstowe to takie gramatyki ogólne, w których prawe strony produkcji nie są krótsze od lewych.

Def. 53. *Gramatyka kontekstowa*, to dowolna taka czwórka $G = \langle N, \Sigma, P, S \rangle$, gdzie:

- N to (skończony) alfabet symboli *nieterminalnych*,
- Σ to (skończony) alfabet symboli *terminalnych*, $\Sigma \cap N = \emptyset$,
- P to skończony zbiór *produkcji*, reguł postaci $\alpha \rightarrow \beta$ dla $\alpha, \beta \in (N \cup \Sigma)^*$, $|\beta| \geq |\alpha|$ i $\alpha \neq \varepsilon$,
- $S \in N$ to *aksjomat* wyróżniony symbol nieterminalny.

□

Def. 54. Powiemy, że język A jest kontekstowy, jeżeli istnieje generująca go gramatyka kontekstowa. □

Klasie języków kontekstowych odpowiadają maszyny Turinga ograniczone liniowo i gramatyki kontekstowe. Maszyny Turinga ograniczone liniowo to takie maszyny, których głowica nie wyjedzie dalej niż (o jeden znak za) słowo dane na wejściu. Inaczej mówiąc, są to maszyny, które dysponują pamięcią proporcjonalną do długości słowa danego na wejściu.

Fakt 22. *Dla każdego języka bezkontekstowego istnieje generująca go (z dokładnością do ε) gramatyka kontekstowa. Inaczej mówiąc, klasa języków kontekstowych zawiera w sobie klasę języków bezkontekstowych.*

Dowód: Zauważmy, że gramatyki bezkontekstowe w postaci normalnej Chomsky'ego są równocześnie gramatykami kontekstowymi. □

Przykład: Język $\{a^n b^n c^n : n \geq 1\}$ jest nie tylko obliczalny, ale i kontekstowy. Oto generująca go gramatyka:

$$\begin{aligned} S &\rightarrow abc \mid XaAbBC \\ Xa &\rightarrow aX \\ Xb &\rightarrow bX \\ Xc &\rightarrow cX \\ XA &\rightarrow aAX \mid aX \\ XB &\rightarrow bBX \mid bX \\ XC &\rightarrow cYC \mid cc \\ cY &\rightarrow Yc \\ BY &\rightarrow YB \\ bY &\rightarrow Yb \\ AY &\rightarrow XA \end{aligned}$$

Oto przykładowe wyprowadzenie słowa $aaabbbccc$:

$$\begin{array}{l}
 S \rightarrow XaAbBC \rightarrow aXAbBC \rightarrow aaAXbBC \rightarrow \\
 \rightarrow aaAbXBC \rightarrow aaAbbBXC \rightarrow aaAbbBXC \rightarrow \\
 \rightarrow aaAbbBcYC \rightarrow aaAbbBYcC \rightarrow aaAbbYBcC \rightarrow \\
 \rightarrow aaAbYbBcC \rightarrow aaAYbbBcC \rightarrow aaXAbbBcC \rightarrow \\
 \rightarrow aaaXbbBcC \rightarrow aaabXbBcC \rightarrow aaabbXBcC \rightarrow \\
 \rightarrow aaabbbXcC \rightarrow aaabbbcXC \rightarrow aaabbbccc
 \end{array}$$

Fakt 23. *Języki generowane przez gramatyki kontekstowe są obliczalne.*

Dowód: Dla gramatyki liniowej długość słowa w wyprowadzeniu nie może maleć. Można więc prześledzić wszystkie możliwe do wyprowadzenia słowa o długości nie przekraczającej długości danego na wejściu słowa. \square

16.5 Hierarchia Chomsky'ego

Hierarchia Chomsky'ego to hierarchia czterech poznanych przez nas klas języków. Została ona stworzona przez Noama Chomsky'ego w ramach jego badań lingwistycznych. Są to języki: regularne, bezkontekstowe, kontekstowe i częściowo obliczalne. Każdej z tych klas odpowiada jeden rodzaj gramatyki. Odpowiadają im gramatyki: liniowe, bezkontekstowe kontekstowe i ogólne. Równocześnie tym samym klasom odpowiadają cztery różne modele obliczeniowe: automaty skończone, automaty stosowe oraz maszyny Turinga (ograniczone liniowo i dowolne). To, że niezależnie stworzone modele do opisu języków odpowiadają tym samym klasom wskazuje na ich głębsze znaczenie.

Zwyczajowo klasy języków tworzące hierarchię Chomsky'ego są oznaczane liczbami od 0 do 3. Następująca tabelka podsumowuje hierarchię Chomsky'ego oraz pozostałe klasy języków, które poznaliśmy w trakcie tego kursu. Są one podane w kolejności od klasy najszerszej do najwęższej.

Nr	Gramatyki	Maszyny	Języki
0	ogólne	maszyny Turinga	częściowo obliczalne
—	—	maszyny Turinga z własnością stopu	obliczalne
1	kontekstowe	maszyny Turinga ograniczone liniowo	kontekstowe
2	bezkontekstowe	automaty stosowe	bezkontekstowe
3	liniowe oraz wzorce	automaty skończone	regularne

16.6 Podsumowanie

W tym wykładzie poznaliśmy przykłady języków, które nie są, odpowiednio, obliczalne i częściowo obliczalne. Jest to problem stopu i jego dopełnienie. Poznaliśmy też gramatyki ogólne i kontekstowe. Gramatyki ogólne są równoważne maszynom Turinga i opisują

dokładnie klasę języków częściowo obliczalnych. Poznane przez nas w trakcie tego kursu gramatyki i klasy języków (liniowe/regularne, bezkontekstowe, kontekstowe i częściowo obliczalne/ogólne) tworzą hierarchię Chomsky'ego. Hierarchia ta podsumowuje większość poznanych przez nas klas języków, z wyjątkiem klasy języków obliczalnych.

16.7 Skorowidz

- **Gramatyki kontekstowe** to takie gramatyki ogólne, w których prawe strony produkcji nie są krótsze od lewych.
- **Gramatyki ogólne (typu 0)** to gramatyki, w których po obu stronach produkcji mogą stać dowolne słowa.
- **Hierarchia Chomsky'ego** to hierarchia czterech klas języków: regularnych, bezkontekstowych, kontekstowych i częściowo obliczalnych, którym odpowiadają cztery rodzaje gramatyk: liniowe, bezkontekstowe, kontekstowe i ogólne. Zobacz także artykuł w Wikipedii.
- **Problem stopu** to język (lub równoważnie problem decyzyjny) złożony z tych kodów maszyn Turinga M i słów x , że M zatrzymuje się dla x . Język ten jest częściowo obliczalny, ale nie obliczalny. Jego dopełnienie nie jest nawet częściowo obliczalne.

16.8 Praca domowa

Napisz program, który wypisuje swój własny kod źródłowy. Uwaga:

- Program ma wypisywać swój kod źródłowy po skompilowaniu, usunięciu pliku źródłowego i uruchomieniu.
- Dopuszczalny jest dowolny, powszechnie dostępny, kompilowany język programowania (np. C, C++, Java, Pascal). Nie może to być jednak język, w którym kod skompilowany zawiera literalnie kod źródłowy.
- Pusty plik nie jest dobrym rozwiązaniem.