

## Автоматы и формальные языки

Карпов Юрий Глебович профессор, д.т.н., зав.кафедрой "Распределенные вычисления и компьютерные сети" Санкт-Петербургского Политехнического университета

karpov@dcn.infos.ru

#### Курс по автоматам и формальным языкам

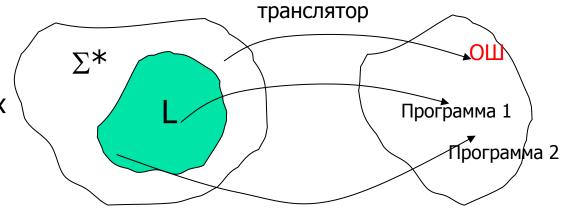
 Цель курса – изучение теоретических основ и принципов построения алгоритмов трансляции



- Входные последовательности символов (цепочки) конечные, они состоят из элементов конечных множеств
- Эти множества называются словарями, а элементы множеств символами
- Конечные последовательности символов называются цепочками (словами, предложениями, текстами)
- Основами трансляции являются теория и методы формального синтаксического и семантического анализа языков (в том числе языков программирования), а также алгоритмы анализа, основанные на этих теориях

#### Трансляторы, формальные языки, автоматы

- $\Sigma$  конечный словарь
- $\Sigma^*$  множество всех конечных цепочек из символов  $\Sigma$
- L подмножество  $\Sigma^*$ , язык



Выход

## программа

\_\_\_х Входной текст∈L x := ab + 23\*(x-b)

Транслятор

алгоритм обработки цепочек символов

0: LOAD 15

1: PUSH '23'

2: LOAD 21

3: LOAD 25

4: SUB

5: MULT

6: ADD

7: STORE 21

Поскольку абстрактные автоматы – это модели, использующиеся для обработки цепочек символов, то курс связан и с абстрактными автоматами



### **Цепочка** (конечная последовательность) символов

- программа на ЯВУ
- цепочка команд языка управления ОС
- фраза естественного языка
- описание ситуации
- описание действий по вязания кофты

# Выход

#### В общем случае - что угодно

- команды стековой машины
- цепочка (предложение) другого языка;
- последовательность запрашиваемых действий;
- образ, картинка, граф-схема программы;
- команды вязальной машины для вязания кофточки;
- команды на исполнительные механизмы робота

Входы разные, выходы разные, но ИДЕИ и методы, ЛЕЖАЩИЕ В ОСНОВЕ АЛГОРИТМОВ ТРАНСЛЯЦИИ – ОДНИ И ТЕ ЖЕ Наша задача – изучить эти идеи и методы

#### Зачем изучать теорию формальных языков

- Трансляторы языков программирования высокого уровня пишутся не очень часто, но знание основ построения трансляторов необходимо каждому
- В любой области знаний нужен свой специализированный язык. Например, нужно проверить правильность теоремы: "Даны три утверждения, истинных или ложных. Если из истинности первого следует истинность третьего, то верно, что если из истинности второго следует истинность третьего, то истинность первого или второго утверждения влечет истинность третьего". Как провести ее доказательство?

Эта теорема на формальном ЯЗЫКЕ логики:

$$(p \Rightarrow r) \Rightarrow ((q \Rightarrow r) \Rightarrow (p \lor q \Rightarrow r))$$

доказательство того, что эта формула – тавтология – тривиально

Любой интерфейс человека с машиной и межмашинный интерфейс - не аморфная последовательность 0 и 1. Последовательность имеет структуру, следовательно интерфейсы, входные языки пакетов программ и т.п. – все основаны на специализированных языках

С необходимостью разработки специализированных языков и построения транслятора для них сталкивается любой непримитивный программист



#### Важный, интересный материал

Специализированные формальные языки используются во всех областях деятельности человека, их изобретают для описания предметных областей, задания алгоритмов, ...

"*Границы моего языка есть границы моего мира*" Людвиг Витгенштейн

"Теория формальных языков — это цветок информатики, лепестками которого являются теория языков программирования и компиляция"

Альфред Ахо

# Мотивация

- Мы будем изучать теорию
  - Теорию абстрактных автоматов и формальных языков но с определенной практической целью: изучение методов построения трансляторов
- Изучение теории дает концепции и принципы, которые позволяют строить трансляторы и преобразователи любых языков (текстов)
- Теория формальных языков, формальных грамматик и синтаксический анализ одна из наиболее полно исследованных областей Computer Science. Ее теоретические результаты решают практические задачи: дают базу для построения трансляторов любых языков
- Семантический анализ и оптимизация кода два направления, в которых исследования продолжаются. В области синтаксиса ПОЧТИ все вопросы глубоко исследованы и хорошо поняты



#### "Деятельная" компонента

Курс по программной теории и инженерии не может быть чисто теоретическим: он должен включать обязательную часть "learning by doing"

Каждый студент должен построить свой транслятор

#### Как это сделать?

На лекциях подробно разбирается компилятор, построенный для базового языка Milan (Mini Language)

Каждому студенту предоставляется его собственное расширение возможностей этого языка, например, введение дополнительных управляющих конструкций или типов

Наряду с расширением базового компилятора реализацией этих конструкций, от студента требуется полное, доскональное понимание работы базового компилятора



Вход – программа на базовом языке MiLan

#### ЕСТЬ

будем вместе разрабатывать

Базовый компилятор для языка MiLan

Выход – программа для виртуальной стековой машины

Вход — программа на модифицированном языке MiLan++

У каждого студента свое расширение языка MiLan Разрабатывает студент

Компилятор для языка MiLan++ Выход — программа для виртуальной стековой машины (возможно, с доп. командами)

#### Организационные вопросы: отличия от предыдущего курса!!

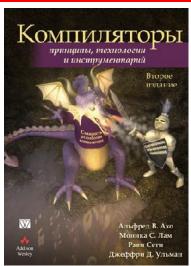
- Лекции
  - Каждая посещенная лекция дает +1 балл; посещение всех лекций дает +5 баллов
- Практическая работа
  - Практические занятия: посещение, активность и выполнение 10 самостоятельных работ от -10 до +15 баллов
  - Курсовая работа: построение своего расширения компилятора рекурсивного спуска базового языка Милан от - 10 до +5 баллов
  - Доп: за освоение Lex, Jacc (Bison), разработку компилятора сложных конструкций языков до +10 баллов
  - Доп: за освоение библиотеки обработки регулярных выражений до +10 баллов
- Экзамен письменный. Задачи на сайте //dcn.ftk.spbstu.ru/~karpov
  - за решение всех задач на экзамене до 80 баллов
- Всего можно набрать 140 баллов (15лек + 5лек + 15практ + 25курс + 80экз)
- Общая оценка выводится по сумме набранных баллов:

#### Автоматы и языки - в рамках одной теории

 Теория формальных грамматик и формальных языков тесно связана с теорией автоматов











М.В.Мозговой. Классика программирования: Алгоритмы, Языки, Автоматы, Компиляторы

Д. Хопкрофт, Р. Мотвани, Д. Ульман. Введение в теорию автоматов, языков и вычислений

А.Ахо, Р.Сети, Д.Ульман. Компиляторы: принципы, технологии и инструменты

Ю.Г.Карпов. Теория автоматов. Питер, СПб, 2003

Ю.Г.Карпов. Теория и технология программирования. Основы построения трансляторов. БХВ, 2006 (есть в библиотеке и на кафедре)

Т.Пратт, М.Зелковиц. Языки программирования: разработка и реализация.: Питер, 2002

Автоматы и формальные языки



|   | Автоматы и языки                               | _            | 4 л   |
|---|--|--------------|-------|
|   | Синтаксис: порождающие грамматики Хомского     | _            | 3 л   |
| • | Семантика: атрибутные трансляции и двусмысленн |              | 2 -   |
|   | КС-грамматики                                  | _            | 2 л   |
|   | Трансляция контекстно свободных (КС) языков    | _            | 6 л   |
|   | Доп. лекции                                    | -            | 2 л   |
|   |  | ====         | ====  |
|   |  | :            | 17 л  |
|   | Что должен знать и уметь студент               | 1 л (если ус | пеем) |
|   | (разбор задач по курсу)                        | ====         |       |

18 л



| ADTOMOTIL IA GOLUGIA  |   |   | 4 - |
|---|---|---|-----|
| Автоматы и языки  | • | _ | 4 л |
| • Лекция 1. Формальные языки. Примеры языков. Грамматики. КА              |   |   |     |
| <ul> <li>Лекция 2. Теория конечных автоматов-распознавателей</li> </ul>   |   |   |     |
| <ul> <li>Лекция 3. Трансляция автоматных языков</li> </ul>                |   |   |     |
| <ul> <li>Лекция 4. Регулярные множества и регулярные выражения</li> </ul> |   |   |     |
| Порождающие грамматики Хомского   | _ |   | 3 л |
| Атрибутные трансляции и двусмысленные КС-грамматики                       |   | _ | 2л  |
| Распознаватели КС-языков и трансляция                                     | _ |   | 6 л |
| Дополнительные лекции -   |   |   | 2 л |

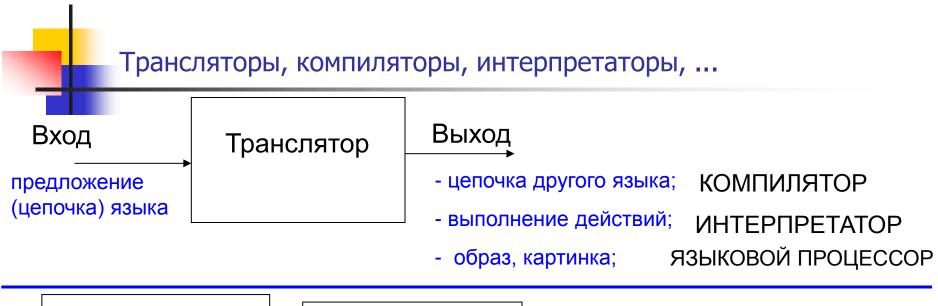


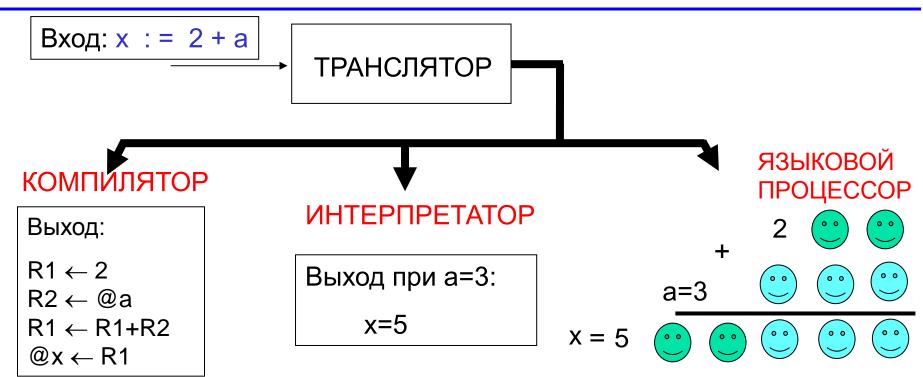
Формальные языки Примеры языков Грамматики Конечные автоматы как грамматики



#### Компиляторы, интерпретаторы, препроцессоры

- Входной язык называется исходным (source language), а входная программа на этом языке называется исходным кодом (source code)
- Машинный язык, на который осуществляется трансляция, называется объектным языком, результат трансляции называется объектным кодом
- Компьютер, на котором оттранслированная программа будет запущена, называется *целевой машиной* (target machine), он не обязательно совпадает с типом машины, на которой работает транслятор
- Если задачей является трансляция в язык низкого уровня, например, машинный язык или язык ассемблера, то такой транслятор называется компилятором
- Другая возможная задача при трансляции программы на алгоритмическом языке – непосредственное выполнение транслятором операторов исходного кода. В этом случае транслятор называется интерпретатором
- Языковой процессор программа, переводящая текст во что-нибудь, в программу на другом языке или, например, в картинку





# Отличия интерпретатора от компилятора Вход предложение (цепочка) языка х:=a+b\*c Транслятор Выход - цепочка другого языка; КОМПИЛЯТОР - выполнение действий; ИНТЕРПРЕТАТОР - образ, картинка; ЯЗЫКОВОЙ ПРОЦЕССОР

- Компилятор при обработке оператора присваивания x:=a+b\*c построит и выдаст последовательность команд объектного языка, осуществляющую вычисление значения выражения a+b\*c и засылку его в область памяти, отведенную для x
- Интерпретатор непосредственно выполнит действия, предписываемые этим оператором присваивания, т.е. выполнит a+b\*c для текущих значений правой части и перешлет полученное значение в область памяти для х. Классическим языком, для которого разработаны интерпретаторы, является Бэйсик
- Языковой процессор может цепочку х:=a+b\*с перевести
  - в картинку
  - в оператор присваивания на С: пару символов ':=' заменит на '=='
  - ...



#### Пример транслятора - интерпретатора

Вход

Транслятор

Выход

Запись решения шахматного этюда в стандартной шахматной нотации

1. Kpd2 Kb2

2. Kpc2 Kpa5

3. Kpb3 Kpa5

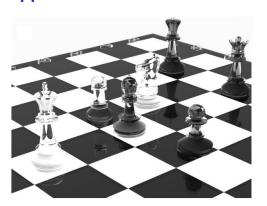
4. Kd6! c4+

5. K:c4+ Kpb5

6. ... ...

Здесь транслятор — интерпретатор, для каждого хода выдает: с какого поля на какое поле нужно переместить фигуру

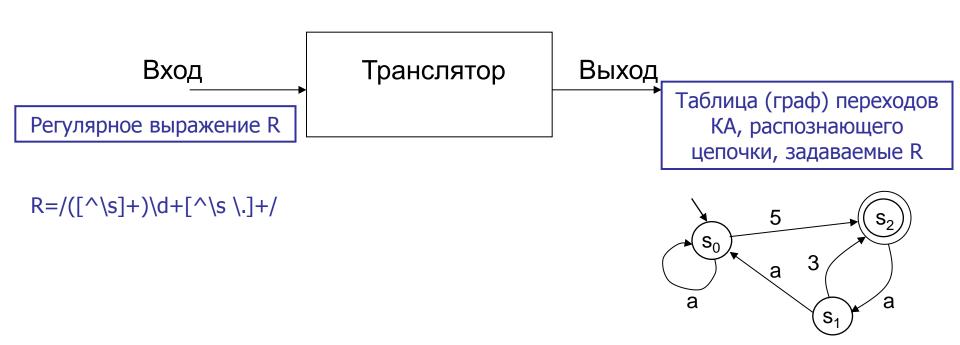
Движение фигур на шахматной доске



- Транслятор это преобразователь информации, который по входной цепочке символов выдает нужный выход
- Предложение (слово, цепочка) языка всегда конечно: предложение есть соединение слов, выражающее законченную мысль

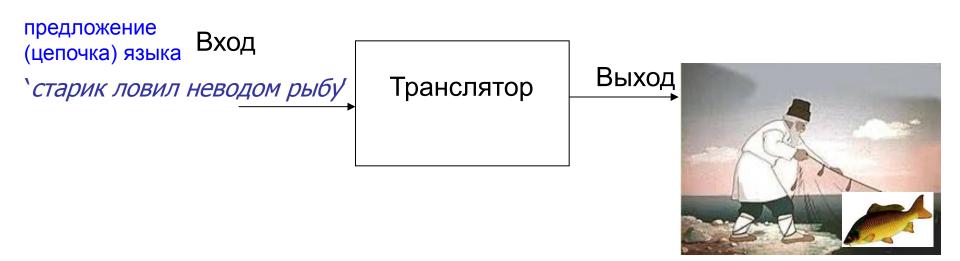


#### Пример языкового процессора





#### Пример языкового процессора



• Языковый процессор может транслировать текст в образы



# Основная цель курса: научить создавать трансляторы

#### Вопрос:

Как создать свой транслятор – компилятор или интерпретатор, и вообще какуюнибудь программу для обработки некоего текста, представляющего любую информацию:

- программу на некотором новом языке программирования,
- входной файл языка управления работами,
- просто какие-нибудь данные, вводимые с клавиатуры?

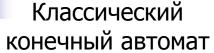
#### Ответ:

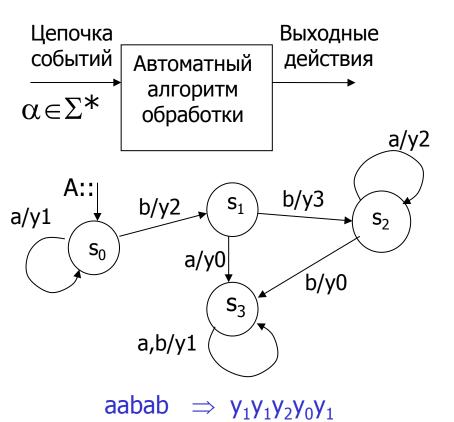
Теория формальных языков – это наука о том, как разобрать входную информацию, вычисляя ее "смысл".

Создать свой анализатор, свой парсер и "генератор" выходных данных можно на основе теории, представленной в курсе лекций

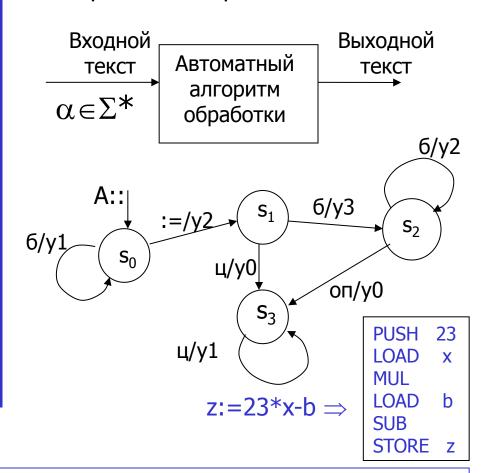
# 4

#### Связь теории автоматов и теории языков





Конечный автомат как "обработчик" предложений языков



Определение и обработка языков - это одна из самых неожиданных областей применения автоматов (не обязательно только конечных)

#### Языки – что это? Как определить язык формально?

- Неформально
  - Язык система знаков и правила их соединения для передачи информации последовательностями знаков
- Формальное определение языка должно быть очень широким и абстрактным, оно должно покрывать:
  - языки естественные "Мама мыла раму"
  - языки программирования "begin int k = 5; print (2\*k+8) end"
  - язык Римских чисел "СМХСІІ" (992), "ММХІ" (2011)
  - язык запросов к БД
  - язык формул логики высказываний "a $\checkmark$ (b $\oplus$ c), b $\Rightarrow$ (a  $\downarrow$  c  $\lor$  d)"
  - язык азбуки Морзе " · · · · · " − буква э, " — " − буква м,
  - язык для станков с ЧПУ (операторы перемещения инструмента, резать по дуге, поднять инструмент, подача в точку, резать от и до, ...)
  - язык управления ОС
  - язык для описания алгоритма вязания кофточки
  - ИТ.Д.

#### Основные определения и понятия

- Конечное множество объектов называется словарем.
- Элементы словаря называются символами.
- Цепочкой (символов) над словарем называют произвольную конечную последовательность символов словаря.
- **Пример.** Пусть  $\Sigma$  словарь, содержащий три символа, *a,b,c.* Над словарем  $\Sigma$  можно построить различные цепочки, например,  $\alpha = abbca$ ,  $\beta = cca$ ,  $\gamma = b$ .
- Важна пустая цепочка, для ее обозначения введен специальный символ ε.

$$\Sigma = \{a, b, c\}$$

$$\Sigma^* = \{\epsilon, a, b, c, aa, ab, abbca, ...\}$$

$$L \subseteq \Sigma^*$$

$$L = \{abb, abcbb, bbcc, aa, ab, ...\}$$

 Самое общее определение языка: Формальный язык – это некоторое множество предложений (цепочек), построенных из символов некоторого конечного словаря

#### Операции над цепочками

- Цепочки строятся из символов конечного словаря
- Пусть  $\Sigma$  словарь, содержащий три символа: a, b, c.  $\Sigma$ ={a, b, c}
  - над  $\Sigma$  можно построить различные цепочки, например,  $\alpha$ =abbca,  $\beta$ =cca
- Операция конкатенации двух цепочек:  $\alpha\beta$ =abbcacca;  $\beta\alpha$ =ccaabbca
- Пустая цепочка  $\varepsilon$  единица конкатенации:  $\varepsilon \alpha = \alpha \varepsilon = \alpha$
- Пусть  $\alpha$ β $\gamma$  цепочка над каким-нибудь словарем.  $\alpha$ ,  $\beta$ ,  $\gamma$ -тоже цепочки
  - $\alpha$  префикс цепочки  $\alpha\beta\gamma$
  - β подцепочка цепочки αβγ
  - $\gamma$  суффикс цепочки  $\alpha\beta\gamma$

В любой цепочке можно выделить (неоднозначно!) и префикс, и суффикс

- Операция подстановки в цепочку вместо одной подцепочки другой цепочки
  - Пусть правило замены: cac o b Цепочку abbcacca по этому правилу замены можно преобразовать в цепочку abbbca.

Это можно записать так:

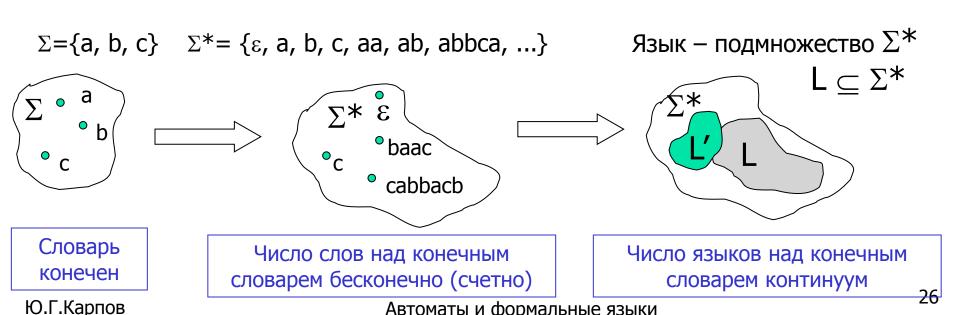
 $abbcacca \Rightarrow abbbca$ 



#### Язык: формальное определение

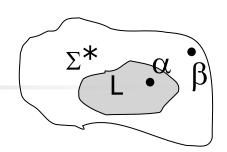
Языком над конечным словарем  $\Sigma$  называется произвольное множество конечных цепочек над этим словарем

- Цепочки языка называются словами (иногда предложениями)
- Над конечным непустым словарем можно определить бесконечное количество слов конечной длины (счетное множество)
- Очевидно, что над конечным непустым словарем можно определить бесконечное количество языков, т.е. подмножеств множества всех возможных слов (континуум, т.к. число подмножеств счетного множества имеет мощность континуум)



Автоматы и формальные языки

#### Примеры языков



- $\Sigma_1$ ={a,b,c}  $L_1$ ={abc, cc} конечный язык  $cc \in L_1$   $cbc \notin L_1$
- $\Sigma_2$ ={a,b,c}  $L_2$ =  $\varnothing$  пустой язык  $\mathsf{cc} \not\in \mathsf{L}_2$
- $\Sigma_3$ ={a,b,c}  $L_3$ =  $\Sigma^*$  все возможные цепочки из a, b и c принадлежат  $L_3$  асс∈ $L_3$ , так же, как и любая цепочка из символов словаря  $\Sigma_3$
- $\Sigma_4 = \{a,b,c\}$   $L_4 = \{a^nbc^m | n, m \ge 0\}$  язык  $L_4$  не из всех цепочек! ааааbcc  $\in$   $L_4$  cbaa  $\notin$   $L_4$
- $\Sigma_5 = \{a,b,c\}$   $L_5 = \{a^nbc^n | n \ge 0 \}$ •  $aabcc \in L_5$   $cbaa \notin L_5$
- $\Sigma_6$ ={a,b,c}  $L_6$ ={ $\alpha \in \Sigma_6^*$  | в  $\alpha$  количества вхождений a, b и c равны}  $\mathsf{ccbaba} \in \mathsf{L}_6$ ,  $\mathsf{ccb} \not\in \mathsf{L}_6$
- $\Sigma_7 = \{0,1\}$   $L_7 =$  множество четных двоичных чисел  $100 \in L_7$   $0111101 \notin L_7$

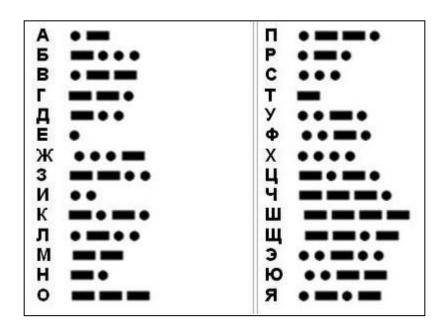


#### Пример языка: азбука Морзе

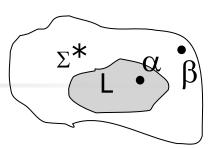
Словарь 
$$\Sigma = \{ \bullet, - \}$$

Язык азбуки Морзе – все кодировки букв (их конечное число)

Слова – последовательности букв – отделяются друг от друга разделителем (пробелом)



#### Примеры языков

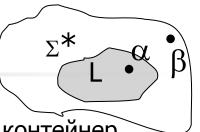


- $\Sigma_8 = \{ '+', '-', '0', ..., '9' \};$   $L_8 =$  множество целых констант  $-67 \in L_8 +234 \in L_8 -3 \in L_8 +2+3...4 \notin L_8$
- $\Sigma_{10} = \{ '+', '-', '0', ..., '9', '.' \}$   $L_{10} =$  множество вещественных констант  $+6.7 \in L_{10}$   $23+.+4..6 \notin L_{10}$ .
- $\Sigma_{11} = \{a\}$   $L_{11} = \{\alpha \in \Sigma^* \mid |\alpha| = 2k,$ где  $k = 0, 1, ...\}$  аааа $\in L_{11}$  а $\notin L_{11}$  (цепочки из четного числа а)
- $\Sigma_{12}$ ={a}  $L_{12}$ = цепочки из а длиной  $n^2$   $\{a^{n^2} | n=0,1,...,\}$  aaaa $\in$   $L_{12}$  aaa $\notin$   $L_{12}$

Языки могут быть конечными и бесконечными. Словарь всегда конечен

Примеры: Конечные языки: Азбука Морзе, L={abc, cc}, пустой язык, ... Бесконечные языки:  $L_8$ ,  $L_9$ ,  $L_{10}$ , ...

#### Примеры языков (3)



- $\Sigma_{13} = \{$ большой, мощный, тяжелый, ..., кран, локомотив, контейнер, ..., ставит, ведет, кладет, везет,... $\}$ ;  $L_{13} =$  русский язык большой локомотив везет тяжелый мощный кран  $\in L_{13}$  молодая красивая студентка мчалась галопом в карете по пыльной дороге  $\in L_{13}$  большой большой студентка контейнер красивый  $\notin L_{13}$
- $\Sigma_{14}$ ={`+',`-', ..., `0',..., `9', ..., begin, end, if, ..., };  $L_{14}$ =язык Паскаль end begin if :: begin 0:== ;  $\not\in L_{14}$  begin int a end  $\in L_{14}$
- $\Sigma_{15} = \{ \text{'+', '-', ... , '0', ..., '9', ..., begin, end, ...} \};$   $L_{15} = язык MiLan (Mini Language)$  end begin if :: begin  $0 :== ; \notin L_{15}$  begin a := 5; язык Милан основной язык

```
a := 5;

input (x);

a := x + a*3;

while x>0 do a:=a*2; x:=x-1

end \in L_{14}
```

Язык Милан – основной язык для наших примеров, для него будет строиться компилятор в рамках курсовой работы

#### Операции над языками

 Поскольку языки – это множества цепочек, то все операции над множествами применимы и к языкам:

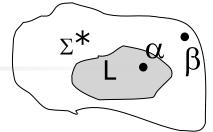
$$L_1=\{\varepsilon, acbac, b, ac\}; L_2=\{b, bca\};$$

■ Пример: теоретико-множественные операции  $L_1 \cup L_2 = \{\epsilon, \text{ acbac, b, ac, bca}\}; L_1 \cap L_2 = \{b\}; L_2 \setminus L_1 = \{\text{bca}\}$ 

#### Специальные операции над языками

- Конкатенация языков:  $L_1L_2 = \{\alpha\beta \mid \alpha \in L_1, \beta \in L_2\}; L^{n+1} = L^nL; L^0 = \{\epsilon\}$
- Итерация (звезда Клини): L\* =  $\bigcup_{k=0,1,...} L^k = \{\epsilon\} \cup L \cup LL \cup LLL \cup LLL \cup ...$  L+ =  $\bigcup_{k=1,2,...} L^k = L \cup LL \cup LLL \cup LLLL \cup ...$ 
  - Примеры:  $L_2^* = \{\varepsilon, b, bca, bb, bbca, bcab, bcabca, bbbca, bbbca, ... \}$   $\Sigma = \{a, b, c\}; \Sigma^* = \{\varepsilon, a, b, c, ab, bb, abbca, ... \}$  все цепочки из a, b, c
- Префиксное замыкание языка:  $Pref(L) = \{\alpha \mid (\exists \beta \in \Sigma^*): \alpha \beta \in L\}$ 
  - Примеры: L<sub>1</sub>={ε, acbac, b, ac}; Pref(L<sub>1</sub>) = {ε, b, a, ac, acb, acba, acbac} дополнение языка префиксами
     Σ = {a, b, c}; Σ\*= все возможные цепочки. Pref(Σ\*) = Σ\*
     L={ε, b, a, ac, bc, bca}, Pref(L) = L
     префиксно-замкнутый язык

#### Как задать язык? Грамматики



- Язык в общем случае бесконечное множество цепочек
   Но как его задать конечным образом, не перечисляя все цепочки?
- Что такое L<sub>9</sub> "множество правильных скобочных выражений"?
   Что такое "скобочное выражение"? Что такое "правильное"?
   Как задать новый язык программирования?
- Любая конечная формальная модель для задания языка называется грамматикой



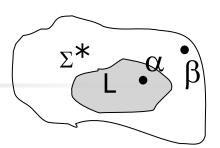
Формальная грамматика — способ выделения некоторого подмножества из множества всех возможных предложений.

Различают порождающие и распознающие грамматики; первые задают правила, с помощью которых можно построить любое предложение языка,

вторые позволяют по данному предложению определить, входит ли оно в язык



#### Пример порождающей грамматики



 Правила порождения бесконечного множества цепочек, не перечисляя их все



Например, оракул (диктатор):

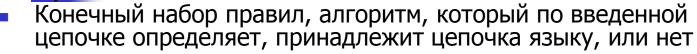
все, что он изрекает, правильно!

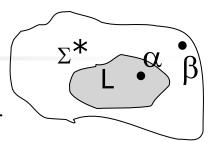
Например, партийная дисциплина:

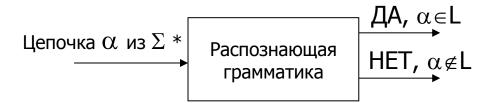
все, что наша партия заявляет, является правильным!

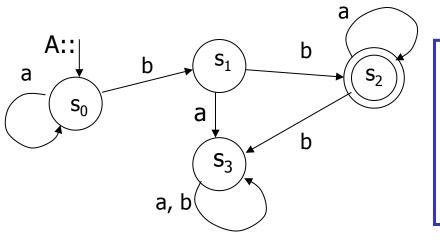
Мы вернемся к порождающим грамматикам позже

#### Распознающие грамматики









Конечный автомат без выхода можно считать распознающим механизмом, распознающей грамматикой,

КА может задавать язык!

Выходом автомата является тип состояния (Допускающее, НЕ Допускающее)

 $aaabba \in L_A$   $bbaa \in L_A$   $abb \in L_A$ 

aaa  $\not\in L_A$ abab  $\not\in L_A$ baaaba  $\not\in L_A$ bbba  $\not\in L_A$ 

#### Конечный автомат-распознаватель

- Конечный автомат-распознаватель  $A=(S, \Sigma, s_0, \delta, F)$ , где:
  - S конечное множество состояний
  - ullet  $\Sigma$  конечное множество входных символов
  - $s_0 \in S$  начальное состояние
  - $\delta$ : S ×  $\Sigma$  → S функция переходов

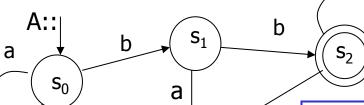
b

F ⊆S – множество финальных (допускающих) состояний



• Автомат А допускает язык L, если он допускает все цепочки этого языка, и только их





 $S_3$ 

a, b

$$L_{A} = \{ \alpha \in \Sigma^{*} \mid \delta^{*}(s_{0}, \alpha) \in F \}$$

Формально:  $\delta^*$ :  $S \times \Sigma^* \rightarrow S$ 

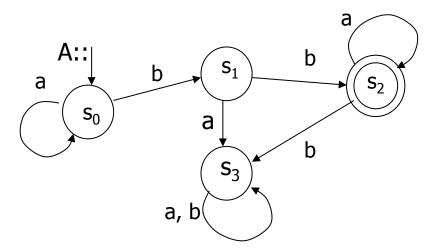
$$(\forall s \in S) \delta^*(s, \varepsilon) = s;$$

$$(\forall s \in S) (\forall \alpha \in \Sigma^*)(\forall a \in \Sigma) \delta^*(s, \alpha a) = \delta(\delta^*(s, \alpha), a)$$

# Пример конечного автоматараспознавателя

- Конечный автомат-распознаватель  $A=(S, \Sigma, s_0, \delta, F)$ , где:
  - S конечное множество состояний
  - $\Sigma$  конечное множество входных символов
  - $s_0 \in S$  начальное состояние
  - $\delta$ : S ×  $\Sigma \to$  S функция переходов
  - F ⊆S множество финальных (допускающих) состояний





Формально: 
$$A=(S, \Sigma, s_0, \delta, F)$$

S=(
$$s_0$$
,  $s_1$ ,  $s_2$ ,  $s_3$ );  $\Sigma = \{a, b\}$ ;

$$\delta$$
:  $S \times \Sigma \rightarrow S$ 

$$\delta(s_0, a) = s_0$$
;  $\delta(s_0, b) = s_1$ ;  $\delta(s_1, a) = s_3$ ;  $\delta(s_1, b) = s_2$ ;  $\delta(s_2, a) = s_2$ ;  $\delta(s_2, b) = s_3$ ;

$$F = \{s_2\}.$$



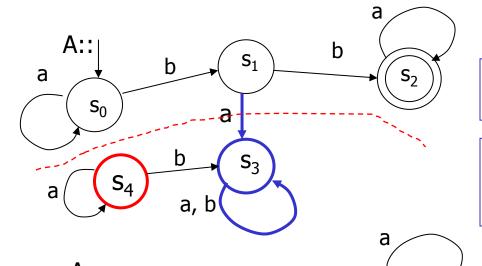
#### Конечный автомат-распознаватель

В автомате из состояния могут быть переходы не для всех входных символов  $L_{A} = \{a^{n} b b a^{m} \mid n, m \geq 0 \}$ 

b

 $S_2$ 





h

Состояние s4: автомат в него никогда не попадет

Состояние s3: если автомат в него попал, то никогда больше не попадет в допускающее состояние

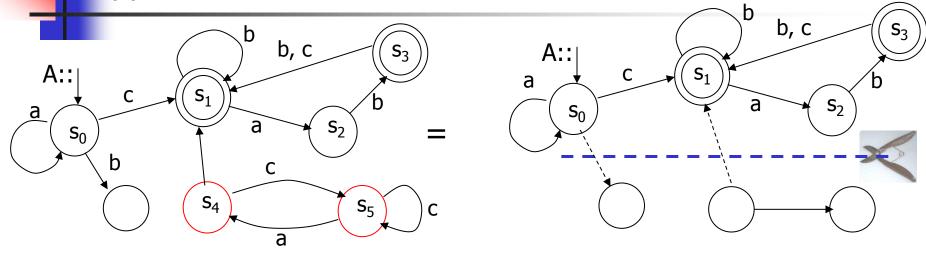
Состояние s<sub>4</sub> недостижимое Состояние s<sub>3</sub> некодостижимое

Такие состояния и соответствующие пути часто на рисунке выбрасывают. Это упрощает визуальное восприятие и облегчает понимание функционирования КА. Но если они реально есть, то они означают наличие дедлока

 $S_0$ 

a

#### Достижимые и недостижимые состояния



Достижимо(s) 
$$\Leftrightarrow$$
 ( $\exists \alpha \in \Sigma^*$ ):  $\delta(s_0, \alpha) = s$ 

$$\Sigma_A = \{ a, b, c \}$$

Недостижимо(s)  $\Leftrightarrow$  ( $\forall \alpha \in \Sigma^*$ ):  $\delta$ (s<sub>0</sub>, $\alpha$ )≠s

Общепринятая условность: рассматриваются только "приведенные" автоматы, у которых все состояния достижимы и кодостижимы. Выбрасывание недостижимых и некодостижимых состояний не изменяет языка, распознаваемого автоматом. Такая операция называется trimming (подстригание)

Задача 1: определите формально, какие состояния являются кодостижимыми

Задача 2: постройте алгоритм тримминга

# Приме

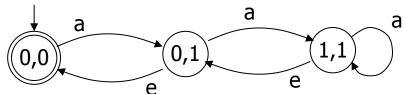
#### Пример КА для дискретной событийной системы

 Процессор обрабатывает поток работ, имея на входе буфер, который может хранить не более одной работы.
 Если пришедшая работа нашла систему занятой, эта работа отбрасывается

а: событие прихода работы

е: событие завершения обработки

состояние: <буфер В, процессор>

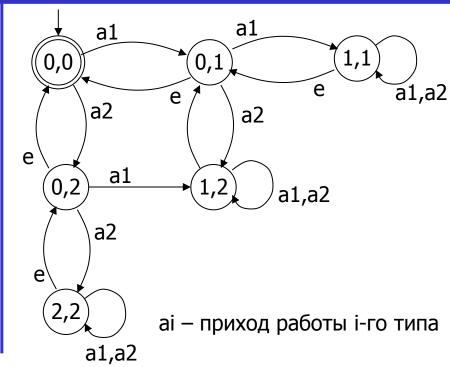


**Усложнение.** Работы во входном потоке могут быть двух приоритетов: 1 и 2. Приоритет 2 более высокий.

Если работа нашла систему занятой (буфер, и процессор заняты), то работа отбрасывается.

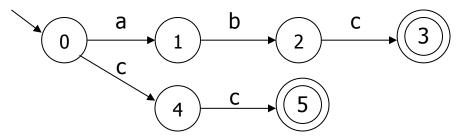
Если пришла работа приоритета 2, буфер пуст, а обрабатывается работа приоритета 1, то пришедшая работа выталкивает менее приоритетную, которая помещается в буфер

Система переходов задает язык



#### Построение автоматов для распознавания языков

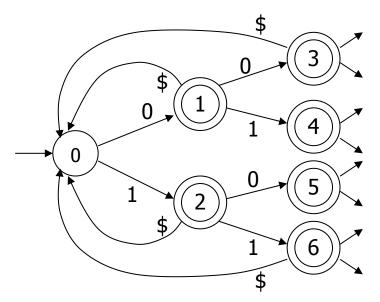
•  $\Sigma_1 = \{a,b,c\}$   $L_1 = \{abc, cc\}$  $cc \in L_1$   $cbc \notin L_1$ 

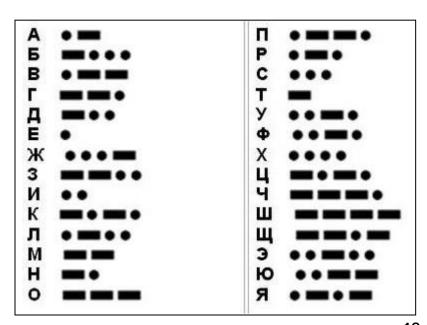


Любой конечный язык может быть распознан конечным автоматом

•  $\Sigma = \{0, 1, \$\}$  L = язык азбуки Морзе (\$ - разделитель между словами)

"··—·" – буква э, "——" – буква м, ...

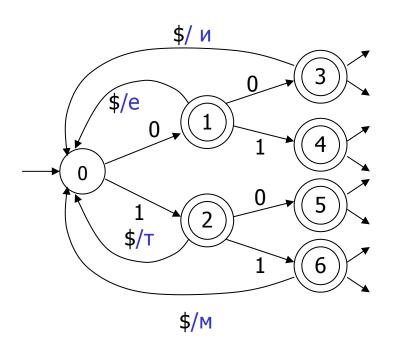


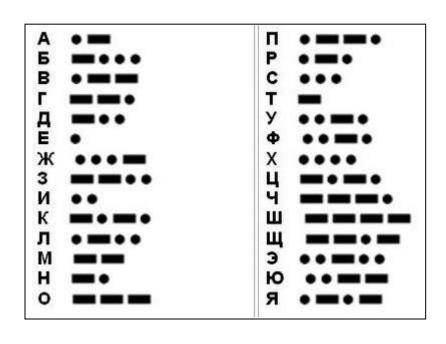




#### Автомат с выходом для трансляции языков

 Добавив выходы (семантические действия) к некоторым переходам, можем получить не только распознавание, но и трансляцию





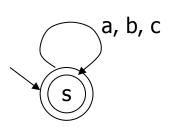
Мы потом покажем, что трансляция – перевод входной цепочки в некоторый выход – полностью основана на распознавании этой цепочки



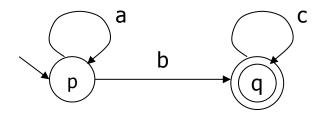
#### Построение автоматов для распознавания языков



■  $\Sigma_3$ ={a,b,c}  $L_3$ =  $\Sigma_3^*$  acc∈ $L_3$ , так же, как и любая цепочка из символов a, b, c



■  $\Sigma_4$ ={a,b,c}  $L_5$ ={a<sup>n</sup>bc<sup>m</sup>| n, m≥0} aaabcc∈L<sub>5</sub> cbaa∉L<sub>5</sub>



Если один язык "больше" другого. то это вовсе не значит, что автомат, его распознающий, сложнее. Например, язык L3 - самый "большой", он включает все возможные цепочки, но он распознается очень простым автоматом

Язык  $L_4 \subset L_3$ , но распознающий его автомат сложнее

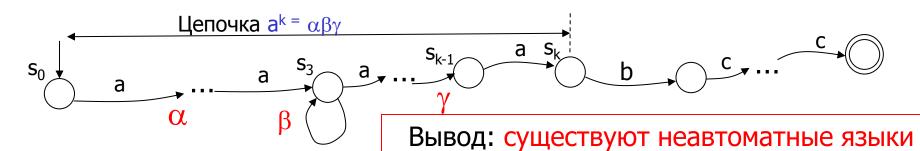
А язык  $L_5 = \{a^nbc^n\} \subset L_4$ , для него вообще нет конечного автомата, его распознающего. Все зависит от структуры цепочек языка!

#### Лемма о накачке. Доказательство неавтоматности языка

- $\Sigma_5 = \{a,b,c\} \quad L_5 = \{a^nbc^n | n \ge 0\}$   $aabcc \in L_5 \quad cbaa \notin L_5$
- Определение. Язык L называется автоматным, если существует конечный автомат, допускающий L
- Теорема. Язык {a<sup>n</sup>bc<sup>n</sup>} не автоматный

"Лемма о накачке"

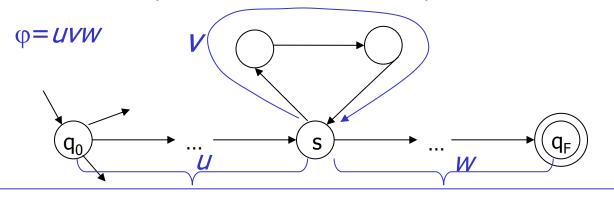
■ Доказательство. Предположим противное, т.е. что язык  $\{a^nbc^n\}$  автоматный. Тогда существует автомат A, распознающий этот язык. A имеет какое-то число состояний. Пусть A имеет k состояний. Возьмем цепочку  $\{a^kbc^k\}$ . Поскольку длина цепочки из a =числу состояний, то в автомате A есть состояние, которое цепочка, состоящая только из a, на пути в заключительное состояние проходит дважды, т.е.  $\alpha\beta\gamma bc^k \in L_5$ . Но если  $\alpha\beta\gamma bc^k \in L_5$ , то  $\alpha\beta^k\gamma bc^k \in L_5$ , причем  $\beta$  состоит из a и не пуста. Но уже цепочка  $\alpha\beta\beta\gamma bc^k \not\in L_5$ 



43

#### Лемма о накачке (pumping lemma, лемма о разрастании )

- Пусть хотим доказать, что L, содержащий бесконечное число цепочек, НЕ автоматный
- Предположим противное, т.е. пусть L является автоматным. Тогда существует какой-то конечный автомат A, распознающий L. У него сколько-то состояний. Сколько мы не знаем. Но знаем точно, что это число конечное. Пусть число его состояний равно k. Выберем в L цепочку φ с подцепочкой V длиной ≥ k. Эти цепочки мы можем выбрать ПРОИЗВОЛЬНО, так, чтобы доказательство было максимально простым
- Ясно, что на пути, который проходит автомат А при приеме подцепочки φ, какое-то состояние встретится не менее, чем два раза.



#### Пример:

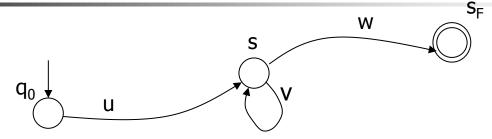
число состояний k, и цепочка длиной ≥ k

Пусть  $\phi = UVW$ , где U - цепочка от  $q_0$  до s, V - цепочка из s в s, и W - из s до  $q_F$  (финального состояния). Длина цепочки  $\phi$ :  $|\phi| \ge k$ . Цепочка  $\phi$  обязательно содержит подцепочку, которая переводит A из состояния s в то же состояние s, назовем ее v. Очевидно,  $\phi = UVW$ . Поэтому, наряду с цепочкой UVW, автомат A должен допускать

также и *UW*, UVVW, UVV ...VW, и т.д.



#### Лемма о накачке - формально



- Пусть L автоматный язык. Тогда существует такое число k (зависящее от L, k-длина "накачки" для L) со следующим свойством: для каждой строки φ языка L длиной по крайней мере k строка φ может быть записана как конкатенация трех подстрок u, v и w (φ = uvw), и при этом выполняются следующие условия: 1. |v|>0, 2. |uv|≤k, 3. для каждого целого i≥0, uviw∈L
- Лемма. Если L автоматный язык, то:

 $(\exists k \in \mathbb{N})(\forall \phi \in \mathbb{L}: |\phi| \geq k)(\exists u, v, w \in \Sigma^*)[\phi = uvw \land |uv| \leq k \land |v| \geq 1 \land (\forall i \in \mathbb{N} \cup \{0\})uv^i w \in \mathbb{L}]$ 

Для доказательства неавтоматности L мы используем контрапозицию этой импликации: Если:

(∀k∈N)(∃φ∈L:|φ|≥k)(∀u,v,w∈∑\*: φ=uvw ∧|uv|≤k ∧ |v|≥1) [(∃*i* $∈N∪{0}) uv 'w ∉L], то язык L неавтоматный.$ 

Проблема: Докажите, что вторая формула – контрапозиция Леммы!

В учебниках почитайте примеры использования Леммы для доказательства неавтоматности языков

#### Пример доказательства неавтоматности языка

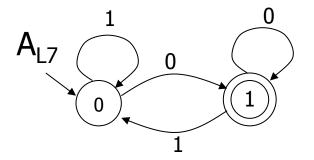
- Термин «накачка» отражает возможность многократного повторения (0, 1, 2, ...) некоторой подстроки в какой-нибудь подходящей строке подходящей длины автоматного языка
- Лемма используется для того, чтобы доказать неавтоматность некоторого языка, имеющего бесконечное число цепочек
- Пример доказательства. Докажем, что язык L={wcw<sup>R</sup> | w∈{a, b}} неавтоматный, w<sup>R</sup> обратная цепочка к w (например, для w= abaa w<sup>R</sup>=aaba)
- Доказательство от противного. Предположим противное, т.е. что L автоматный. Тогда существует конечный автомат A, распознающий этот язык. Пусть A имеет k состояний. Возьмем цепочку длиной 2k+1, первая часть ее w имеет длину k. Следовательно, w может быть представлена, как w=  $\gamma\mu\nu$ , так. что подцепочка  $\mu$  возвращается в то же состояние, из которого она началась, и  $|\mu| \ge 1$ .
- В соответствии с Леммой о накачке, если автомат А допускает wcw<sup>R</sup> =  $\gamma \mu \nu cw^R$ , то он допускает и  $\gamma \mu \mu \nu cw^R$ , и  $\gamma \mu \mu \mu \nu cw^R$ , и т.д. Но все эти цепочки не принадлежат языку L, поскольку все цепочки в L должны быть симметричными.
- Наше предположение о том, что язык L={wcw<sup>R</sup> | w∈{a,b}} автоматный, приводит к противоречию. Следовательно, этот язык – не автоматный

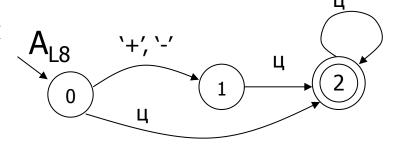


#### Построение автоматов для распознавания языков

- $\Sigma_6$ ={a,b,c}  $L_6$ ={ $\alpha \in \Sigma_6^*$  | в цепочке  $\alpha$  количества вхождений a, b и c равны}  $\mathsf{ccbaba} \in \mathsf{L}_6$ ,  $\mathsf{ccb} \notin \mathsf{L}_6$
- $\Sigma_7 = \{0,1\}$   $L_7 =$  множество четных двоичных чисел  $100 \in L_7$   $0111101 \notin L_7$
- $\Sigma_8 = \{ \text{'+', '-', '0',..., '9'} \};$   $L_8 =$  множество целых констант:  $-67 \in L_8$   $234 \in L_8$   $-3 \in L_8$   $+2+3...4 \notin L_8$

Язык  $L_6$  не автоматный

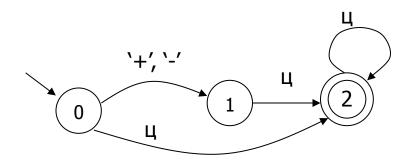




Язык L9 не автоматный



## Конечные автоматы: переход под воздействием символа из конечного множества символов



#### Как понимать переход, помеченный ц?

- Здесь символом ц обозначена любая цифра, фактически, конечное множество символов:  $\mu = \{0, 1, 2, ..., 9\}$
- Очевидно, что при реализации КА проверку того, принадлежит ли входной символ конечному множеству символов, выполнить всегда можно
- Например, поскольку кодировки всех цифр идут подряд, то для того, чтобы проверить, принадлежит ли входной символ s множеству ц, достаточно выполнить проверку  $0' \le s \le 9'$

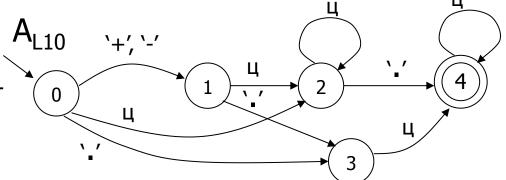


#### Построение автоматов для распознавания языков

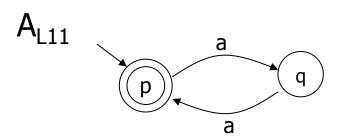
■ 
$$\Sigma_{10}$$
={'+', '-', '0',...,'9', '.'}  $L_{10}$ =

множество вещественных констант

6.7 $\in$ L<sub>10</sub> 23+.+4.. 6 $\notin$ L<sub>10</sub>



•  $\Sigma_{11} = \{a\}$   $L_{11} = \{\alpha \in \Sigma_{11}^* \mid$ длина  $\alpha$  равна 2k, где  $k = 0, 1, ...\}$  аааа $\in L_{11}$  а $\notin L_{11}$ 



• 
$$\Sigma_{12} = \{a\}$$
  $L_{12} = \{a^{n^2} | n = 0, 1, ..., \}$  aaaa $\in L_{12}$  a $\in L_{12}$ , aaa $\notin L_{12}$ 

Язык L<sub>12</sub> не автоматный



•  $\Sigma_{13}$ ={большой, мощный, тяжелый, ..., кран, локомотив, контейнер, ..., ставит, ведет, кладет, везет,...}  $L_{13}$ = русский язык

Большой тяжелый локомотив везет маленький красный контейнер  $\in L_{13}$ 

Большой кран большой поезд везет тяжелый везет кладет $otin L_{13}$ 

- Словарем естественного языка являются словоформы слова с их изменениями по родам, числам, падежам и проч. Всех словоформ конечное число, а предложений языка бесконечное число. Примеры словоформ: большой, ставит, кран, кладет везет, везли, везут, ..... Это все различные символы конечного словаря для русского языка. Т.о. цепочками естественного языка можно считать его предложения
- Русский язык (и другие естественные языки) НЕ являются автоматными
- Небольшие фрагменты естественных языков можно задать конечными автоматами





#### Языки программирования высокого уровня

- $\Sigma_{14} = \{ '+', '-', ..., '0', ..., '9', ..., begin, end, if, ..., \};$   $L_{14} =$  язык Паскаль end begin if :: begin  $0 := ; \notin L_{14}$  begin int a; end  $\in L_{14}$ 
  - Языки программирования высокого уровня не являются автоматными: они содержат вложенности конструкций произвольной глубины
  - Небольшие фрагменты языков программирования являются автоматными
  - Это:

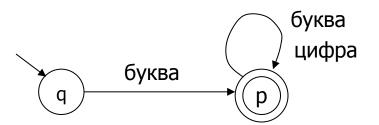
```
комментарии
идентификаторы (имена)
служебные слова (begin, while, for, end, ...)
константы
```

- Обычно обработка этих фрагментов выполняется на первых проходах трансляторов, где выбрасываются комментарии, пробелы, выделяются лексемы
- Лексемы -минимальные единицы языка, имеющие смысл
- Первые проходы трансляторов строятся на основе теории конечных автоматов распознавателей

# Язык Милан

```
\Sigma_{15} = \{ \text{'+', '-', ..., '0',..., '9', ..., begin, end, if, ..., } \}; L_{15} = язык MiLan (Mini Language) "end begin if :: begin 0 :== \text{;"} \not\in L_{15} "begin a :=5; input (x); a :=x+a*3; while x>0 do a:=a*2; x:=x-1 end" \in L_{15}
```

- Язык Milan не автоматный: в частности, он содержат вложенности операторов произвольной глубины
- Идентификатор это любая последовательность букв и цифр, начинающаяся с буквы



Это – автоматный подъязык языка Милан



# Дополнительный материал: Теорема Майхилла-Нероуда об автоматных языках

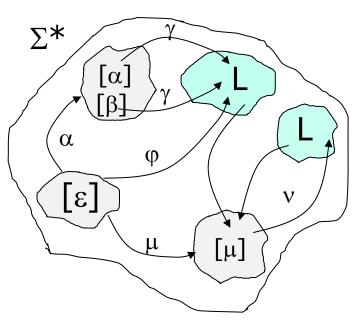
# Когда множество цепочек будет автоматным языком? Отношение эквивалентности на множестве цепочек

Правым инвариантом называется такое отношение эквивалентности R на  $\Sigma^*$ , что:

$$(\forall \alpha, \beta \in \Sigma^*) \alpha R\beta \Rightarrow (\forall \gamma \in \Sigma^*) \alpha \gamma R\beta \gamma$$

Пусть L — язык над словарем  $\Sigma$ . Правый инвариант  $R_L$ , ассоциированный с L:

$$(\forall \alpha, \beta \in \Sigma^*) \alpha R_L \beta \Leftrightarrow (\forall \gamma \in \Sigma^*) (\alpha \gamma \in L \Leftrightarrow \beta \gamma \in L)$$



Классы эквивалентности  $R_L$   $\delta([\alpha], a) = [\alpha a]$ 

- Классы эквивалентности  $R_L$  это множество таких предысторий (цепочек), что  $\alpha$  и  $\beta$  попадают в один класс ([ $\alpha$ ] = [ $\beta$ ]), если для любой цепочки  $\gamma$  либо обе цепочки  $\alpha\gamma$  и  $\beta\gamma$  принадлежат L, либо обе не принадлежат L
- На рисунке: [ε] класс эквивалентности, содержащий пустую цепочку;
- L класс эквивалентности, содержащий цепочки языка L;
- φ цепочка принадлежит языку L;
- [μ] класс эквивалентности цепочек, не принадлежащих L (любое продолжение не приведет в L). Соответствующее состояние A<sub>L</sub> некодостижимо

#### Теорема Майхилла-Нероуда

- **Теорема Майхилла-Нероуда**. (Джон Майхилл, Энил Нероуд, 1958 г.) Пусть L язык над словарем  $\Sigma$ . Следующие три утверждения эквивалентны:
  - 1. Язык L- автоматный, т.е. существует КА допускающий L.
  - 2. Существует правый инвариант R с конечным индексом (конечным числом классов эквивалентности) на  $\Sigma^*$ , такой, что L =  $\cup_{[x] \in \Pi} [x]_R$ , где  $\Pi$  некоторое множество классов эквивалентности R
  - 3. Существует отношение  $R_L$  конечного индекса, ассоциированное с L:

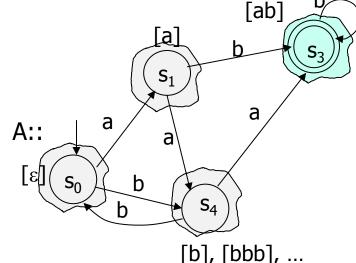
 $(\forall \alpha, \beta \in \Sigma^*) \ \alpha \mathsf{R}_\mathsf{L} \beta \Leftrightarrow (\forall \gamma \in \Sigma^*) \ (\alpha \gamma \in \mathsf{L} \Leftrightarrow \beta \gamma \in \mathsf{L})$ 

#### Доказательство (1 и 3).

Пусть  $R_L$  имеет классы эквивалентности  $[\alpha_1]$ , ...,  $[\alpha_n]$ . Определим конечный автомат  $A_L$  так:

- множество состояний S = {  $[\alpha_1]$ ,  $[\alpha_2]$ , ...,  $[\alpha_n]$  }.
- начальное состояние: [ε] (класс, включающий ε),
- функция переходов  $\delta([\alpha], a) = [\alpha a],$
- множество финальных состояний  $F=\{[\alpha] \mid \alpha \in L\}$

Очевидно, что  $L_A = \{ \alpha \mid \delta([\epsilon], \alpha) \in F \}$ 



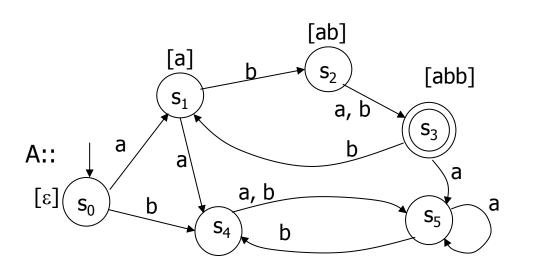
Состояния – классы эквивалентных предысторий.

Эквивалентность – относительно дальнейшего поведения: допущения цепочек языка L

#### Теорема Майхилла-Нероуда (2)

Пример отношения R – правого инварианта с конечным числом классов эквивалентности (конечного индекса) и соответствующего ему конечного автомата.

Классы эквивалентности R:  $\{\{\epsilon\}, \{b,aa,bbb, ababa, ...\}, \{a,abbb, ...\}, ... - соответствуют состояниям$ 



aRabbb, bRaa, abRabbbb, ...

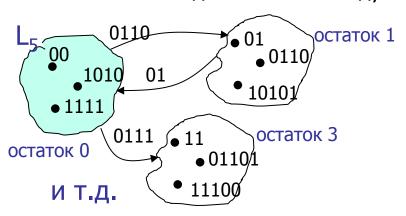
```
egin{array}{lll} s_0: & \{\epsilon\} & = [\epsilon] \\ s_1: & \{a, abbb, ...\} & = [a] \\ s_2: & \{ab, abbbb, ...\} & = [ab] \\ s_3: & \{aba, abab, ...\} & = [aba] \\ s_4: & \{b, aa, bab, aabb, ...\} & = [aa] \\ s_5: & \{bb, aab, bba, ...\} & = [bb] \\ \hline \end{array}
```

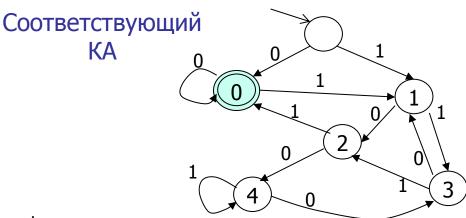
Теорема Майхилла Нероуда определяет необходимое и достаточное условия автоматности языка.

Она также позволяет доказать, что данный язык не автоматный.

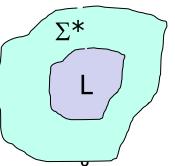
# Пример эквивалентности, ассоциированной с языком над $\Sigma$ , и соответствующего конечного автомата

- Одно из простейших отношений эквивалентности эквивалентность по модулю, (например, по модулю 5) на множестве неотрицательных целых чисел. Назовем его  $\pi_5$
- Числа A и B находятся в отношении  $\pi_5$ , если остатки от их деления на 5 равны. Отношение  $\pi_5$  имеет 5 классов эквивалентности:  $\{0, 5, 10, 15, 20, ...\}$  остаток от деления на 5 = 0,  $\{1, 6, 11, 16, 21, ...\}$  остаток от деления на 5 = 1 и т.д.
- Пусть [а] $\pi$  обозначает такой класс эквивалентности  $\pi_5$ , которому принадлежит а. Поэтому в десятичной системе счисления  $[0]_{\pi 5} = [10]_{\pi 5} = [325]_{\pi 5}$ ,  $[3]_{\pi 5} = [333]_{\pi 5} = [103]_{\pi 5}$
- Пусть  $\Sigma$ ={0, 1} и пусть отношение  $\rho_5$  на  $\Sigma^*$  объединяет цепочки, представляющие те числа в двоичной системе счисления, которые при делении на 5 дают один и тот же остаток. Ясно, что  $\rho_5$  эквивалентность на  $\Sigma^*$  с конечным индексом. Пусть язык L<sub>5</sub> это множество непустых цепочек над  $\Sigma$ , которые представляют двоичные числа, делящиеся без остатка на 5. Ясно, что  $\rho_5$  имеет 6 классов эквивалентности. Одному из них принадлежит только пустая цепочка  $\varepsilon$ , остальные объединяют числа с одним и тем же отношением делимости по модулю 5









- Зачем изучать теорию формальных языков (мотивация)?
  - Теория формальных языков это "цветок информатики", интереснейшая область математики, имеющая приложения в построении трансляторов
  - Каждый специалист в области информатики будет сталкиваться с необходимостью разработки языков интерфейсов, языков представления данных, формул, ... . Новые языки разрабатываются постоянно во всех инженерных областях при их автоматизации, для них нужно строить свои трансляторы
- Что такое формальный язык? Это просто множество цепочек
  - Формальный язык это множество (обычно бесконечное) конечных цепочек (слов, предложений) над конечным словарем
  - Для задания бесконечных объектов математика строит конечные модели
- Что такое грамматика? Это конечный формализм для задания языка
  - Существует два способа задания языка: распознавание и порождение.
     Трансляция входной цепочки основана на распознавании этой цепочки,
     т.е. на алгоритме проверки того, что эта цепочка принадлежит языку

#### Заключение (2)

- конечные автоматы как формализм задания языков
  - КА имеют множество применений. Одно из самых интересных применений.
     в области формальных языков. КА используются здесь как распознаватели "правильных" входных цепочек над конечным словарем. Такие автоматы не имеют выхода: выходом в них является тип состояния, в котором находится автомат (допускающее либо недопускающее)
  - Входная цепочка распознается автоматом, если она переводит его из начального в одно из финальных (допускающих) состояний
  - Каждый конечный автомат-распознаватель задает какой-нибудь язык.
     Язык, для которого существует распознающий его автомат, называется автоматным. Конечных автоматов распознавателей бесконечное число, автоматных языков тоже бесконечное число
- Не все языки являются автоматными
  - Не все формальные языки можно задать с помощью КА. Автоматные языки

     очень слабый подкласс формальных языков. Язык Милан, для которого
     мы будем строить транслятор, не автоматный. Неавтоматным является
     любой язык с вложенными конструкциями
  - Лемма о накачке удобный метод определения того, что заданный язык не является автоматным



### Спасибо за внимание