



Кафедра «Автоматизированные станочные системы»
Dept. of Automated Manufacturing Systems

Классификация грамматик и языков



4 типа грамматик по Хомскому:

V^+ — множество всех цепочек над алфавитом V без λ ;
 V^* — множество всех цепочек над алфавитом V ,
включая λ .

Ноам Хомский
(Noam Chomsky)

Тип 0: грамматики с фразовой структурой

На структуру их правил не накладывается никаких ограничений: для грамматики вида $G(VT, VN, P, S)$, $V = VN \cup VT$ правила имеют вид: $\alpha \rightarrow \beta$, где $\alpha \in V^+$, $\beta \in V^*$. Это самый общий тип грамматик. В него подпадают все без исключения формальные грамматики, но часть из них, к общей радости, может быть также отнесена и к другим классификационным типам. Дело в том, что грамматики, которые относятся только к типу 0 и не могут быть отнесены к другим типам, являются самыми сложными по структуре.

Практического применения грамматики, относящиеся только к типу 0, не имеют.

Тип 1: контекстно-зависимые (КЗ) и неукорачивающие грамматики

В этот тип входят два основных класса грамматик:

Контекстно-зависимые грамматики $G(VT, VN, P, S)$, $V = VN \cup VT$ имеют правила вида: $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$, где $\alpha_1, \alpha_2 \in V^*$, $A \in VN, \beta \in V^+$.

Неукорачивающие грамматики $G(VT, VN, P, S)$, $V = VN \cup VT$ имеют правила вида: $\alpha \rightarrow \beta$, где $\alpha, \beta \in V^+, |\beta| \geq |\alpha|$.

При построении предложений КЗ-грамматик один и тот же нетерминальный символ может быть заменен на ту или иную цепочку символов в зависимости от того контекста, в котором он встречается.

Цепочки α_1 и α_2 в правилах грамматики обозначают **контекст** (α_1 — левый контекст, а α_2 — правый контекст), в общем случае любая из них (или даже обе) может быть пустой. Говоря иными словами, **значение одного и того же символа может быть различным в зависимости от того, в каком контексте он встречается.**

При построении компиляторов такие грамматики не применяются

Неукорачивающие грамматики имеют такую структуру правил, что при построении предложений языка, заданного грамматикой, любая цепочка символов может быть заменена на цепочку символов не меньшей длины.

Тип 2: контекстно-свободные (КС) грамматики

Контекстно-свободные (КС) грамматики $G(VT, VN, P, S)$, $V = VN \cup VT$ имеют правила вида: $A \rightarrow \beta$, где $A \in VN, \beta \in V^+$. Такие грамматики также иногда называют неукорачивающими контекстно-свободными (НКС) грамматиками (видно, что в правой части правила у них должен всегда стоять как минимум один символ).

КС-грамматики широко используются при описании синтаксических конструкций языков программирования. Синтаксис большинства известных языков программирования основан именно на КС-грамматиках

Тип 3: регулярные грамматики

К типу регулярных относятся два эквивалентных класса грамматик: левостолбчатые и правостолбчатые.

Левостолбчатые грамматики $G(VT, VN, P, S)$, $V = VN \cup VT$ могут иметь правила двух видов: $A \rightarrow B\gamma$ или $A \rightarrow \gamma$, где $A, B \in VN, \gamma \in VT^*$.

В свою очередь, правостолбчатые грамматики $G(VT, VN, P, S)$, $V = VN \cup VT$ могут иметь правила тоже двух видов: $A \rightarrow \gamma B$ или $A \rightarrow \gamma$, где $A, B \in VN, \gamma \in VT^*$.

Регулярные грамматики используются при описании простейших конструкций языков программирования: идентификаторов, констант, строк, комментариев и т. д.

Для классификации грамматик всегда выбирают максимально возможный тип, к которому она может быть отнесена. Сложность грамматики обратно пропорциональна номеру типа, к которому относится грамматика. Грамматики, которые относятся только к типу 0, являются самыми сложными, а грамматики, которые можно отнести к типу 3 — самыми простыми.

Классификация языков

Тип 0: языки с фразовой структурой

Это самые сложные языки, которые могут быть заданы только грамматикой, относящейся к типу 0. Если язык относится к типу 0, то для него невозможно построить компилятор, который гарантированно выполнял бы разбор предложений языка за ограниченное время на основе ограниченных вычислительных ресурсов.

К сожалению, **все естественные языки относятся к фразовым**. Структура и значение фразы естественного языка может зависеть не только от контекста данной фразы, но и от содержания того текста, где эта фраза встречается. Одно и то же слово в естественном языке может не только иметь разный смысл, в зависимости от контекста, но и играть различные роли в предложении. Именно поэтому столь велики сложности в автоматизации перевода текстов, написанных на естественных языках

Тип 1: контекстно-зависимые (КЗ) языки

Тип 1 — второй по сложности тип языков. В общем случае время на распознавание предложений языка, относящегося к типу 1, экспоненциально зависит от длины исходной цепочки символов.

Языки и грамматики, относящиеся к типу 1, применяются в анализе и переводе текстов на естественных языках. Распознаватели, построенные на их основе, позволяют анализировать тексты с учетом контекстной зависимости в предложениях входного языка (но они **не учитывают содержание текста**, поэтому для точного перевода с естественного языка требуется вмешательство человека). На основе таких грамматик может выполняться автоматизированный перевод с одного естественного языка на другой, ими могут пользоваться сервисные функции проверки орфографии и правописания в языковых процессорах.

В компиляторах КЗ-языки не используются

Тип 2: контекстно-свободные (КС) языки

КС-языки лежат в основе синтаксических конструкций большинства современных языков программирования,

Тип 3: регулярные языки

Регулярные языки — самый простой тип языков. Поэтому они являются самым широко используемым типом языков в области вычислительных систем. Время на распознавание предложений регулярного языка линейно зависит от длины входной цепочки символов.

Регулярные языки лежат в основе простейших конструкций языков программирования (идентификаторов, констант и т. п.), кроме того, на их основе строятся языки ассемблеров, а также командные процессоры, символьные управляющие команды и другие подобные структуры.

Чем все это безобразиие распознавать

Для языков с *фразовой структурой* (тип 0) необходим распознаватель, имеющий неограниченную внешнюю память. Поэтому для языков данного типа нельзя гарантировать, что за ограниченное время на ограниченных вычислительных ресурсах распознаватель завершит работу и примет решение о том, принадлежит или не принадлежит входная цепочка заданному языку. Практического применения языки с фразовой структурой не имеют.

Для *контекстно-зависимых языков* (тип 1) распознавателями являются двусторонние недетерминированные автоматы с ограниченной памятью. Количество шагов, необходимых автомату для распознавания входной цепочки, экспоненциально зависит от длины этой цепочки.

Экспоненциальная зависимость времени разбора от длины цепочки существенно ограничивает применение распознавателей для контекстно-зависимых языков. Такие распознаватели применяются для автоматизированного перевода и анализа текстов на естественных языках, когда временные ограничения на разбор текста несущественны.

Для *контекстно-свободных языков (тип 2)* распознавателями являются односторонние недетерминированные автоматы с магазинной (стековой) внешней памятью — МП-автоматы. При простейшей реализации алгоритма работы такого автомата он имеет экспоненциальную сложность, однако путем некоторых усовершенствований алгоритма можно добиться полиномиальной (кубической) зависимости времени, необходимого на разбор входной цепочки, от длины этой цепочки. Следовательно, можно говорить о полиномиальной сложности распознавателя для КС-языков.

Пример: грамматика целых десятичных чисел

$G_1\{0,1,2,3,4,5,6,7,8,9,-,+ \},\{S, T, F\},P_1,S)$:

P_1 :

$S \rightarrow T \mid +T \mid -T$

$T \rightarrow F \mid TF$

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

По структуре своих правил данная грамматика G_1 относится к **контекстно-свободным грамматикам (тип 2)**. Ее можно отнести и к типу 0, и к типу 1, но максимально возможным является именно тип 2, поскольку к типу 3 эту грамматику отнести никак нельзя:

строка $T \rightarrow F \mid TF$ содержит правило $T \rightarrow TF$, которое недопустимо для типа 3, и хотя все остальные правила этому типу соответствуют, одного несоответствия достаточно.

Та же грамматика, но по-другому:

$G1'$ ($\{0,1,2,3,4,5,6,7,8,9,-,+ \}, \{S, T\}, P1', S$):

$P1'$:

$S \rightarrow T \mid +T \mid -T$

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

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

По структуре своих правил данная грамматика $G1$ является праволинейной и относится к **типу 3**.

Та же грамматика, но левوليнейная:

$G1''$ ($\{0,1,2,3,4,5,6,7,8,9,-,+ \}, \{S, T\}, P1'', S$):

$P1''$:

$T \rightarrow + \mid - \mid \lambda$

$S \rightarrow T0 \mid T1 \mid T2 \mid T3 \mid T4 \mid T5 \mid T6 \mid T7 \mid T8 \mid T9 \mid S0 \mid S1 \mid S2 \mid S3$
 $\mid S4 \mid S5 \mid S6 \mid S7 \mid S8 \mid S9$