

Об одном способе осуществления синтеза предложения  
при МД *на основе синтагматического анализа*.

Ч. А. Хоар.

my first article,  
written summer 1960, published  
in *Математика*

По обычным представлениям, в результате синтакси- *Перевод*  
тического анализа получается список синтагм, т.е.,  
список упорядоченных пар слов, первое из которых "под-  
чиняет" второе, причем для каждой синтагмы имеется при-  
знак, указывающий на то, какой тип связи имеет место в  
данной синтагме. При синтезировании предложения *другом* *На*  
*/или на том же/ языке*, из списка пар надо получить линей-  
ную последовательность слов, образующую синтаксически  
и морфологически правильное предложение. При этом  
надо определить, какой именно набор признаков синтагм  
оказывается наиболее удобным. В настоящей заметке я не  
буду рассматривать морфологическую сторону синтеза; я  
ограничусь вопросом *о расстановке слов в нужном порядке*  
*/что более существенно для синтеза предложения, например,*  
*на английском языке, и вообще представляется более сложной*  
*задачей/.*

Предлагаемый метод синтеза основывается на том,  
что во многих языках линии зависимости не пересекаются,  
т.е., если слова *a* и *b* входят в одну синтагму (*a, b*), и  
слова *c* и *d* входят в другую синтагму (*c, d*), то эти слова  
нельзя ставить в таком порядке, что линия, проведенная  
между *a* и *b* пересекается с линией, проведенной между  
*c* и *d*. *Это значит, следующие возможности исключаются:*

<i>a...c...b...d</i>	<i>a...d...b</i>
<i>a...b...c...d</i>	<i>b...a...c...d</i>
<i>a...c...b...d</i>	<i>c...b...d...a</i>
<i>b...d...a...c</i>	<i>d...b...c...a</i>

*4)* Это утверждение верно во всех случаях, когда легко  
по интуиции установить линии зависимости. В разборе  
более сложных случаев, надо обеспечить, чтобы осталось  
верным.

Мы можем использовать этот факт, если мы синтезируем фразу слова направо /разыскивая сперва первое слово, потом второе, и так далее/. Положим, что на данном шагу мы правильно расставили первые  $n$  слов /которые мы в дальнейшем обозначим расставленными словами/; мы разыскиваем теперь  $n+1$ -ое слово /разумеется, среди нерасставленных слов/. И здесь полезно определить понятие "косвенной зависимости". Слово  $x_n$  прямо зависит от слова  $b$ , если существует синтагма  $(b, x_n)$  в анализируемом предложении; слово  $x_n$  косвенно зависит от слова  $b$ , если существует в анализируемом предложении последовательность слов  $x_1, x_2, \dots, x_n$ , причем  $x_1$  прямо зависит от  $b$ ,  $x_{i+1}$  прямо зависит от  $x_i$  для каждого  $i$ , и  $x_n$  прямо зависит от  $x_n$ . Иными словами,  $x_n$  является в некотором смысле "предком" слова  $a$ . Теперь мы можем сформулировать следующие правила:

I. Если существует не расставленное слово, прямо зависящее от  $n$ -ого слова, то в качестве  $n+1$ -ого слова нельзя выбрать слово, которое не зависит /прямо или косвенно/ от  $n$ -ого слова /если бы мы выбрали такое слово, это неизбежно привело бы к пересечению линий зависимости/.

II. Если не существует слова, зависящего от  $n$ -ого, но есть слово, подчиняющее  $n$ -ое слово, то можно выбрать в качестве  $n+1$ -ого слова только это слово или слово, прямо или косвенно зависящее от него.  
Мы могли бы привести другие правила, ограничивающие возможности выбора следующего слова  $\neq$  но достаточно этих двух, чтобы показать каким именно образом выбор очередного слова может быть значительно сокращен.

Такие правила, однако, не определяют однозначно какое именно слово надо выбрать в случае если существует больше одного нерасставленного слова, зависящего от последнего / $n$ -ого/ и расставленного слова. Для таких случаев мы вводим особые признаки синтагм. Эти признаки должны указывать, стоит ли зависящее слово налево или направо от управляющего слова, и как "тесно" они связаны. В качестве таких признаков используются целые числа: отрицательные числа обозначают, что зависимое слово стоит налево, и положительные числа — направо; причем меньшая абсолютная величина обозначает более тесную связь, и .наоборот.

Если слово  $a$  более тесно связано со словом  $b$ , чем с  $c$ ,  
то обозначает, что мы должны стоять ближе к  $b$  чем  
т.е., порядок слов будет  $a, b, c$ .  
Т.е., /в случае левой за-  
висимости/, или  $a, c, b$  /в случае правой зависимости/.

Положим теперь, что два нерасставленных слова,  $b$  и  $c$ ,  
зависят от последнего / $a$  ого/ слова в синтагмах  $+1(a, b)$ ,  $\boxed{+3}$   
и  $+3(a, c)$ . В таком случае мы не можем выбрать слово  $c$ ,  
так как  $b$  более тесно связано с  $a$  чем  $c$  ( $|+1| < |+3|$ ).  
/Здесь мы предполагаем, что признаки приписываются с таким  
расчетом, чтобы никогда не встречались два слова, во-  
всех от данного слова с одним и тем же признаком/.

~~то слово  $b$ , которое зависит от  $a$  сама мышь~~  
По мы еще не можем сказать, что именно  $b$  есть следующее ~~следующее~~ помо-  
гательным признаком  
слово : если другое слово  $d$  "зависит налево" /т.е.,  
с отрицательным признаком/ от слова  $b$ , то  $d$  непременно  
должно стоять перед  $b$ . Если ~~два~~ слова "зависят на-  
лево" от  $b$ , то мы выбираем в качестве "кандидата" на  
следующее место то слово, которое ~~менее~~ <sup>тесно</sup> связано с  $b$ , то  
есть, которое должно стоять ~~далее~~ <sup>далее</sup> налево от него. Обо-  
значим это слово  $e$ . Если не существует не расставле-  
нного слова, зависящего от  $e$  налево, мы можем определенно  
ставить  $e$  на  $+1$ -ом месте, и потом весь процесс на-  
чинается с самого начала. Если, однако, существуют слова,  
 зависящие от  $e$  налево, мы выбираем на место "кандидата" то  
слово, которое должно стоять дальше налево : ~~бывший~~  
~~кандидат~~  $e$  отвергается. Мы продолжаем этот процесс,  
 пока не найдем "кандидат", от которого ничего не зависит  
налево : именно это слово будет  $+1$ -ым словом в син-  
тезируемом предложении. Можно грубо описать изложенный  
процесс следующим образом : искать "первое" слово, завися-  
щее направо от последнего расставленного слова  $a$  обозначен-  
ным его  $b$ . Потом искать "украинее" слово, право или кос-  
венко зависящее налево от  $b$  : ставить это слово на  
следующем месте после  $a$ .

Я еще не рассматривал такой случай, когда не существует не расставленного слова, зависящего направо от последнего расставленного слова. Этот случай легко разбирается. Мы переходим на  $n+1$ -ое слово, и обрабатываем его как будто оно является последним расставленным словом: если и у него нет зависящих направо, то мы переходим на  $n+2$ -ое слово, и т.д. Но когда мы нашли искомое слово наших, мы всё равно ставим его на  $n+1$ -ом месте: никаких вставок не может быть.

До сих пор я предполагал, что мы уже расставили первые  $n$  слов: теперь надо объяснить как выбирается первое слово /в случае  $n=0$ / . Здесь мы рассматриваем такой метод анализа, при котором остается одно и только одно слово, которое не зависит от другого слова. В самом начале синтеза удобно взять это независимое слово в качестве "первого" слова, и, идя слова направо, синтезировать вышеописанным образом ту часть фразы, которая стоит направо от него. Потом мы синтезируем левую часть фразы, начиная еще раз с независимого слова, но идя теперь справа налево. При этом, можно использовать тот же самый алгоритм, если только заменить поясоду слова "правое" на "левое", "за" на "перед", и "следующее" на "предыдущее".

Остается теперь лишь один вопрос: каким образом присвоить синтагмам признаки, т.е., как определить "тесноту" связей в синтезируемом предложении. Здесь я возьму в качестве иллюстрации английский язык. Можно перечислить для примера некоторые типы синтагм, в которых существительное фигурирует как управляющий член.

Тип синтагмы	Пример	Признак
(S, артикль)	the	-3
(S, прилагательное)	new	-2
(S, S как прилагательное)	<u>Leica</u> camera	-1
(S, S в антюзии)	(type III)	+1
(S, предлог)	<u>in</u> a leather case	+2
(S, причастие)	<u>bought</u> in Germany	+3
(S, соч. слово)	, <u>and</u> an erazurable	+4

8 Здесь мы считаем, что, когда существительное управляет несколько прилагательных, и только самое близкое связывается с ним, а другие связываются друг с другом в цепочке. Например X "современная вычислительная машина" анализируется в синтагмы ( машина, вычислительная ) ( вычислительная, современная )

Здесь я не претендую на полноту : вероятно, что понадобится больше 7 степеней тесности. Я надеюсь, однако, что из этого списка будет ясно, что в стандартных случаях можно легко установить степень тесности, если только знаем, какие части речи входят в данную синтагму. В нестандартных случаях, конечно, дело обстоит сложнее, и надо будет составить более или менее сложный алгоритм, чтобы разобрать такие случаи. Например, если причастие не входит целого оборота /т.е., от него ничего не зависит кроме, может быть, наречия/, то оно обычно стоит перед существительным на том же месте, где обычно стоят прилагательные : the savvite an exciting play, a play exciting much interest. С другой стороны, если от причастия зависит предложный оборот, оно обычно стоит за существительным : savvite a sacred book, a book sacred to the believers. Но это уже дело лингвистов, и не касается технического метода осуществления синтеза.

Для этого метода синтеза предполагается, что язык, на котором синтезируется предложение, является регулярным в следующем смысле. Положим, что хри слова a, b, c, d, и X зависят от третьего слова x, и что они должны стоять в порядке a...b...x...c...d. } Скажем, что в другом предложении } Если этот порядок остается именным, несмотря на возможное наличие ~~отсутствие~~ других слов, зависящих от X, то язык называется регулярным, и предлагаемый метод применяется без особых затруднений. В таком случае, можно приписать и признак каждой синтагме на основе свойств одной /данной/ синтагмы. В обратном случае, невыгодно было бы применить предлагаемый метод.

Достоинство предлагаемого метода состоит в следующем :

1. Лингвистическая проблема синтеза целого предложения /в случае регулярности/ и проблеме правильной расстановке только двух слов, зависящих от одного слова /так как правильный порядок всех слов, зависящих от данного слова однозначно определяется, если только знать правильный порядок всех пар слов/. Тогда можно составить алгоритм, который присваивает синтагмам их нужные признаки.

2. Алгоритм имеет циклический характер, так что программа будет очень проста, и для реализации процесса понадобится незначительное число ячеек машинной памяти.

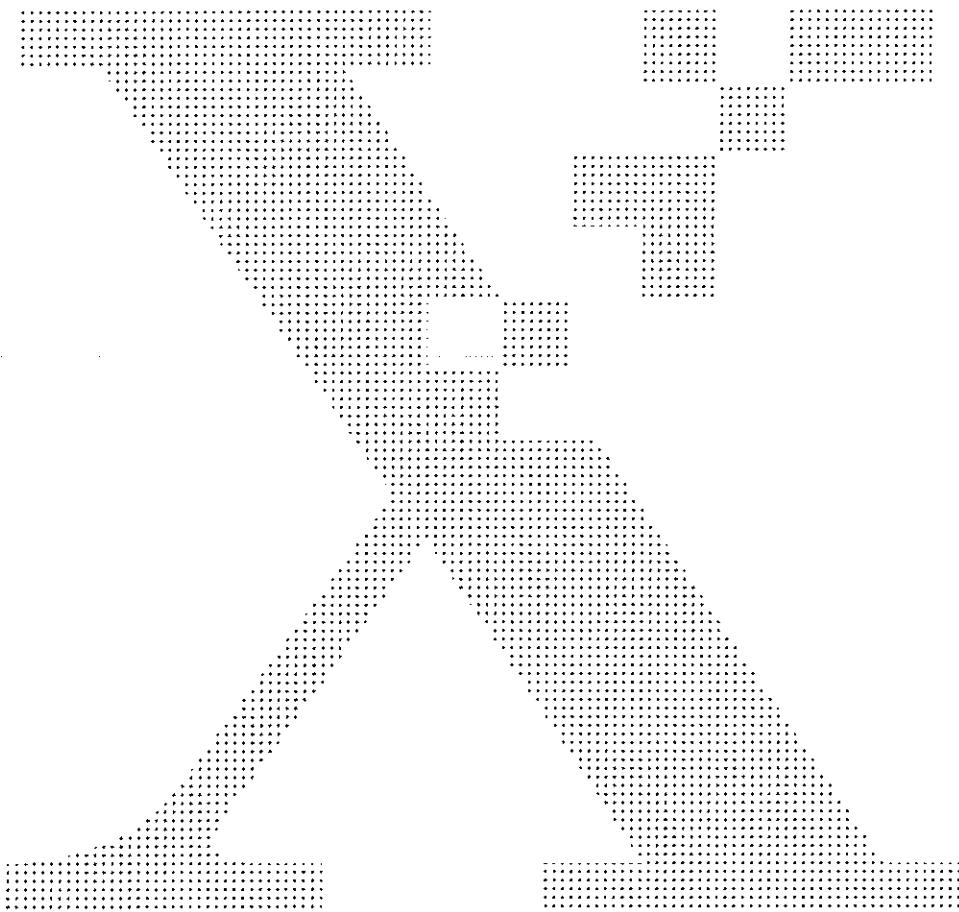
Недостатком предлагаемого метода является то обстоятельство, что его осуществление требует много машинных операций : чтобы найти на каждом шагу крайнее слово налево, надо будет пройти несколько раз вдоль списка синтагм. Этот недостаток может быть существенен, если мы переводим на языки, в котором вообще мало приходится переставлять.

Автор хочет выразить свою благодарность И.А.Мельчуку, который ~~каждых~~ поправил изложение во многих местах, и вели со мной разговоры, на которых эта статья основана.

thoare

# thoare

**Customer Order Form\_N1.pdf**  
**27-09-06 10:06**



thoare

# thoare

India.pdf  
27-09-06 10:05

