

найдена основа, то окончание принимается нулевым и строится парадигма в соответствии с пометами основы. Далее исходная словоформа сопоставляется со словоформами парадигмы. Если найдено совпадение, это означает, что основа выбрана правильно и построенная парадигма печатается в файл. Если первичный поиск не дал результатов, то окончание считается ненулевым, от исходной словоформы отрезается последняя буква, которая сохраняется в символьном массиве, а что остается от словоформы объявляется гипотетической основой. Гипотетическому окончанию приписывается набор грамматических значений в соответствии с приведенной выше табл. 1. Далее гипотетическая основа вновь сопоставляется с основами из базы данных, имеющими пометы на присоединение грамматических показателей с определенными значениями. Если эти пометы и признаки гипотетического окончания совпадают, разбиение словоформы на основу и окончание считается верным. Длина гипотетического окончания составляет не более трех символов. Если в результате поиска не обнаружена основа, то в выходной файл выводится сообщение: «Основа не найдена».

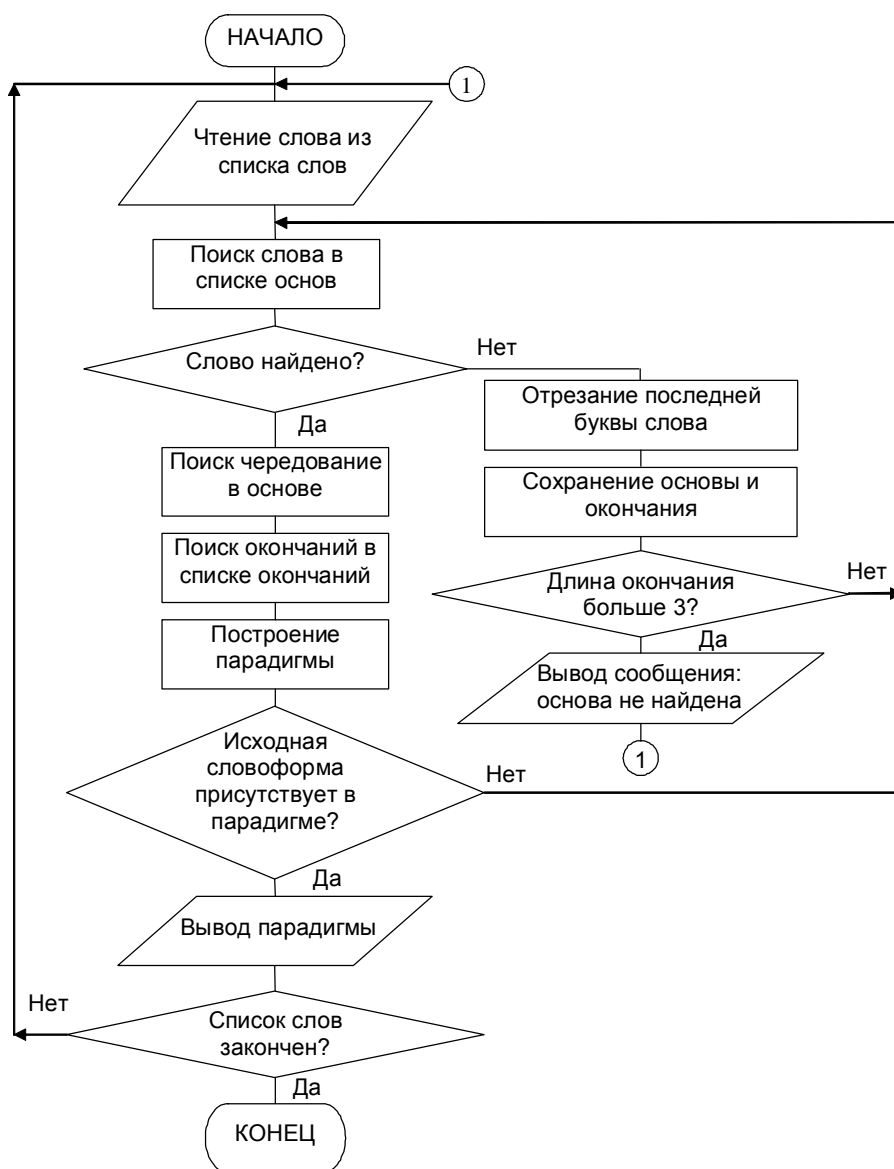


Рис. 1. Блок-схема алгоритма морфемного анализа.

Построение парадигмы происходит следующим образом. Когда основа найдена, анализируются соответствующие ей пометы, указывающие на часть речи, тип и номер склонения, и другие признаки, определяющие правила построения парадигмы. У некоторых словоформ присутствует чередование в основе, на которое указывает помета «*», а для прилагательных — еще и пометы «1» или «2». Различают несколько типов чередования, в зависимости от которых при построении парадигмы учитываются одна или две дополнительные основы. После проверки на чередование в основе осуществляется поиск возможных окончаний в соответствии с таблицей признаков. Далее строится парадигма исходного слова, причем вид парадигмы определяется пометами при найденной основе. Парадигма сохраняется в файле в удобном для пользователя виде. Если не возможно полное или частичное построение парадигмы, то выводятся соответствующие сообщения.

Разработанный программный модуль был использован для обработки словаря А. А. Зализняка. Первоначальный объем словаря составляет 97194 слова. Была сгенерирована расширенная база данных, содержащая исходные основы и дополнительные основы, получаемые в результате чередования; ее объем составил 117351 основу. В качестве проверки был обработан корпус словаря А. А. Зализняка и получена соответствующая статистика. Например, число парадигм для существительных составило 621682 слова, а для прилагательных — 882102 слова. Если при распознавании речи использовать подход «фонемы–словоформы», то только для существительных объем базы данных составил бы более 600 тысяч слов, в то время как при морфологическом анализе количество основ существительных составляет около 46 тысяч. Таким образом, очевидно, что использование морфологического модуля позволяет существенно сократить объем словаря. На рис. 2 приведена диаграмма распределения словоформ по частям речи, на рис. 3 отображено распределение слов с различными видами типами чередования.

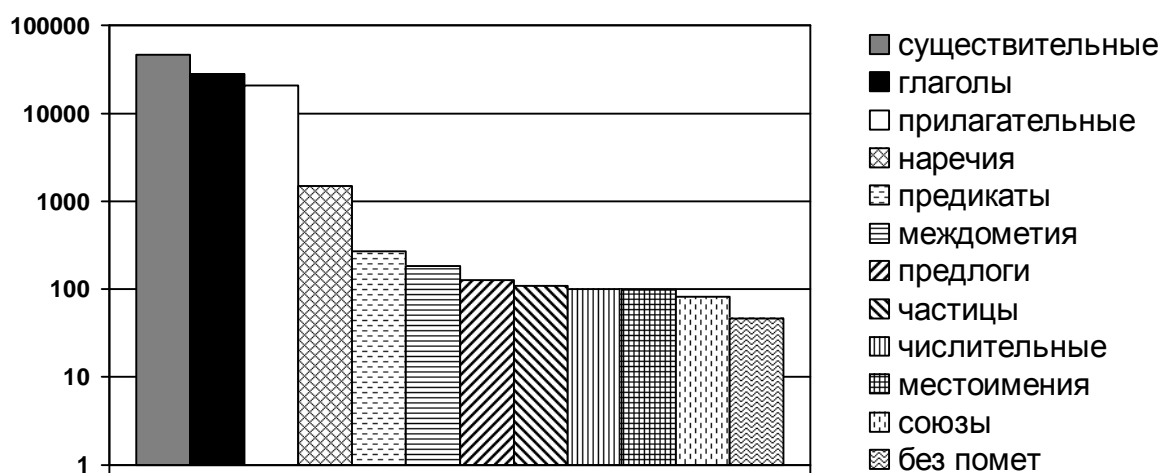


Рис. 2. Распределение словоформ по частям речи.

