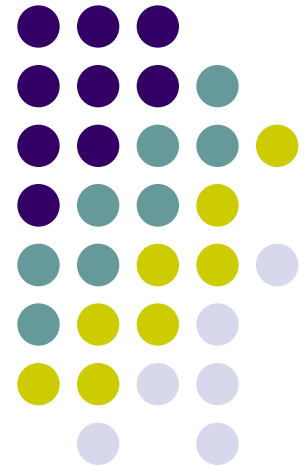


Символьные массивы. Поиск подстроки в строке. БМ-поиск, КМП - поиск





Дано. Текст **s** длины **n** и некоторая строка **p** длины **m** - образец.
Найти вхождение (первое) строки **p** в текст **s**. Определить позицию вхождения.
Рассмотрим 3 метода поиска.



1. Прямой поиск

s = “во дворе трава на траве дрова”

p = “дрова”

При прямом поиске последовательно идем по тексту **s**, пока не встретится символ, совпадающий с первым символом образца **p**. Сравниваем последующие символы по длине образца.

Если какой-то символ текста и образца не совпал, сдвигаемся в тексте **s** на одну позицию от начала сравнения с первым символом образца, и все повторяем.



Пример

→
во дворе трава на траве дрова
дрова ≠

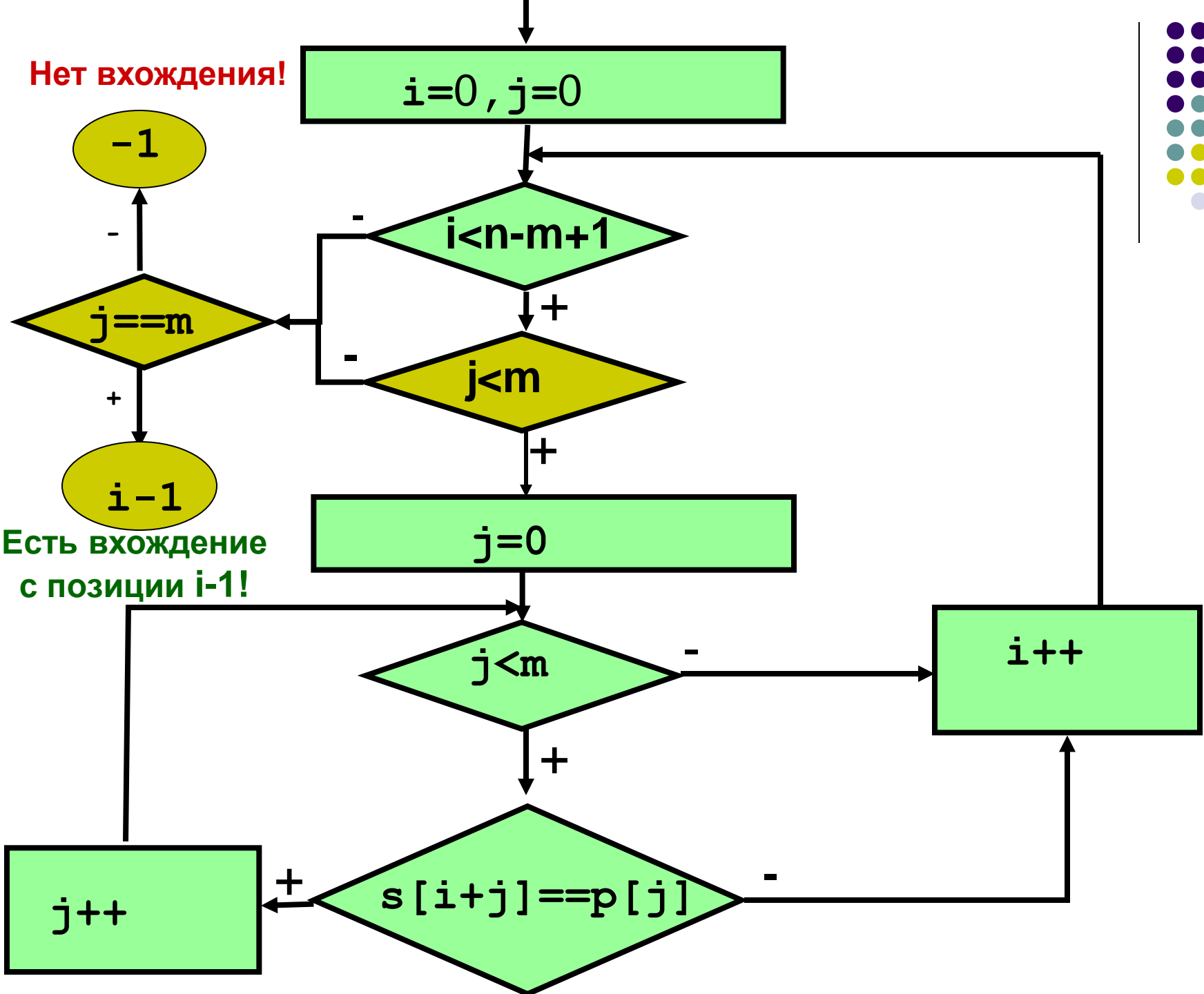
→
во дворе трава на траве дрова
дрова ...

Для сравнения с другими методами
важно заметить:

при неудачном сравнении – поиск
продолжаем всегда с начала образца.



Нет вхождения!



Есть вхождение
с позиции $i-1$!



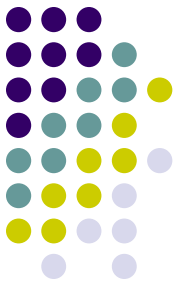
Обсуждение

В прямом методе поиска **не** анализируется структура образца.

В последующих 2-х методах предварительный анализ образца (**предтрансляция**) позволяет сформировать некоторую полезную информацию о его структуре.

В процессе поиска эта информация используется для того, чтобы некоторые операции сравнения **исключить**.

2. БМ-поиск (Бауэр и Мур)



S="МАМА МЫЛА РАМУ"

P="РАМУ"

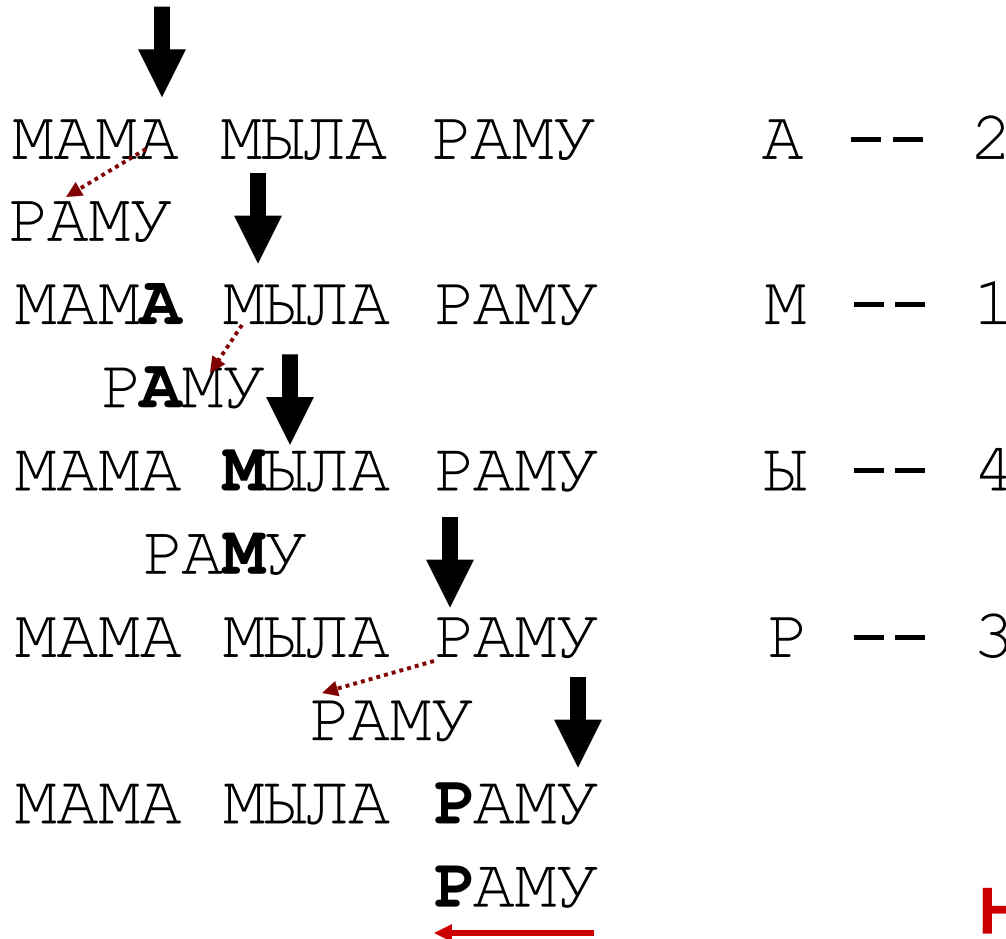
Сравнение с конца образца

A ≠ У, но в образце есть А,
совместим их

M ≠ У, но в образце есть М,
совместим их

Ы ≠ У, и в образце нет Ы,
сместимся на всю длину
образца

P ≠ У, но в образце есть Р,
совместим их



Нашли!

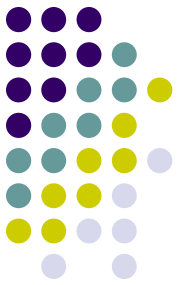
Идея алгоритма

Числа 2, 1, 3 – это расстояния символов образца до его конца.

Таким образом, заранее готовится **таблица чисел из 256 элементов**, в которой формируются расстояния символов образца (кроме последнего) до конца образца. Если каких-то символов **нет** в образце, то соответствующий элемент таблицы полагается **равным длине образца** (в нашем примере 4 – на столько происходит смещение в этом случае).

Для этого **все элементы** таблицы сначала полагаются равными длине образца **m**.

При поиске смещение добавляется к позиции той буквы текста **s**, которая соответствует последней букве образца на текущем шаге.

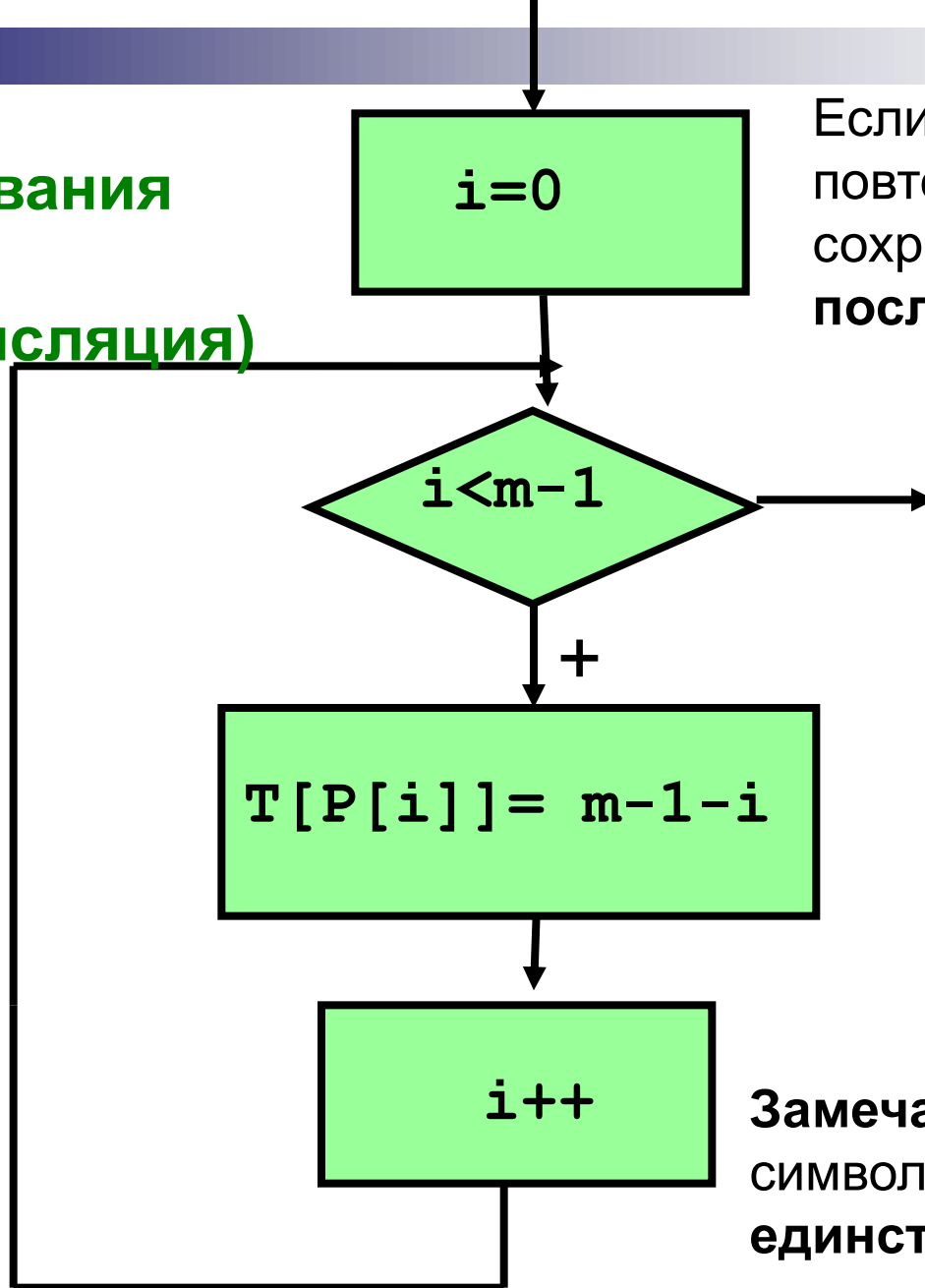




Таблица

**Символу образца соответствует
элемент таблицы с номером, равным
коду символа!**

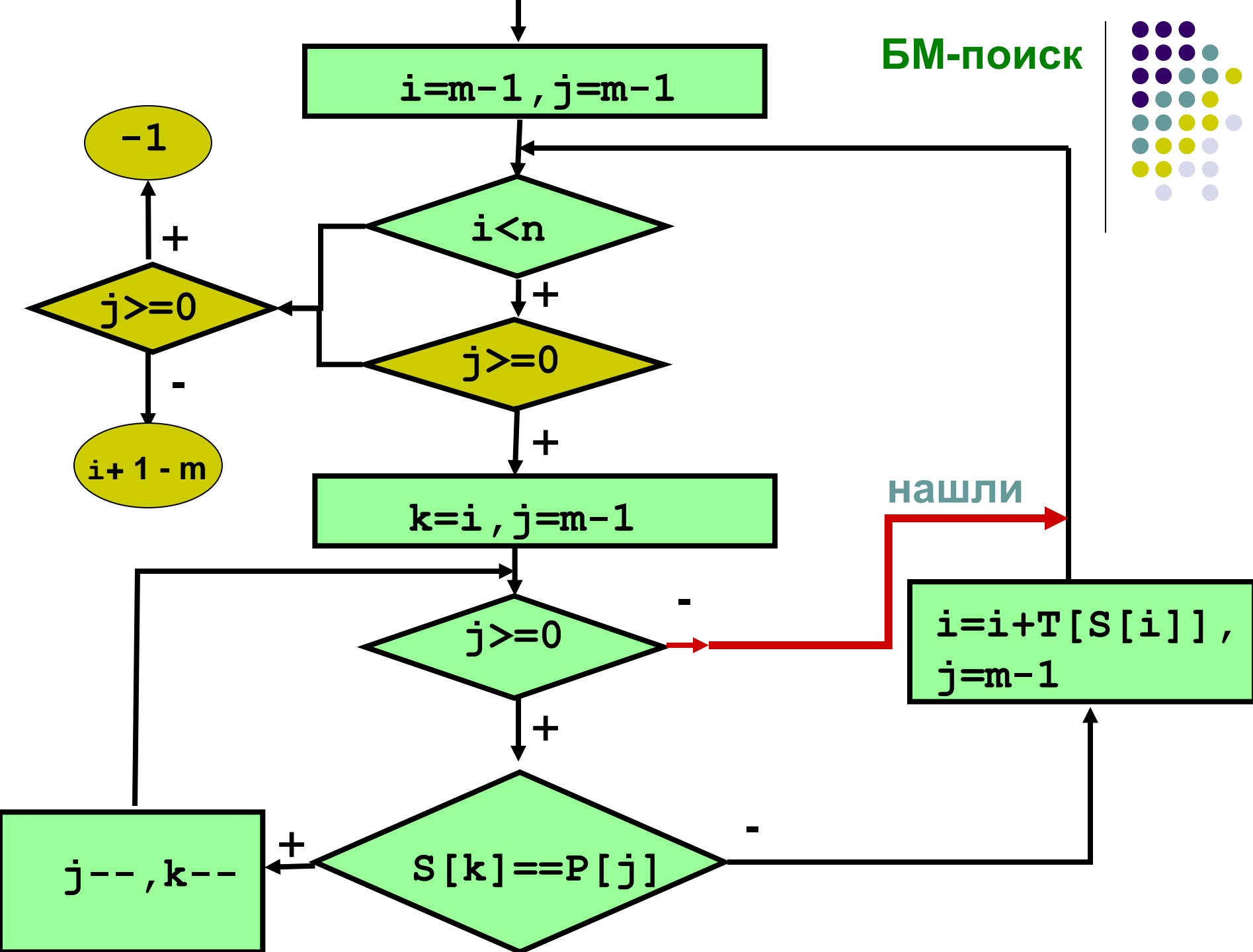
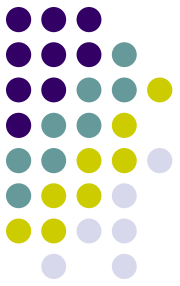
Алгоритм формирования таблицы (предтрансляция)



Если символ в образце повторяется, то в таблице сохраняется расстояние **последнего** из них.

Замечания. Если последний символ образца в нем **единственный**, то соответствующий элемент таблицы остается равным **m**.

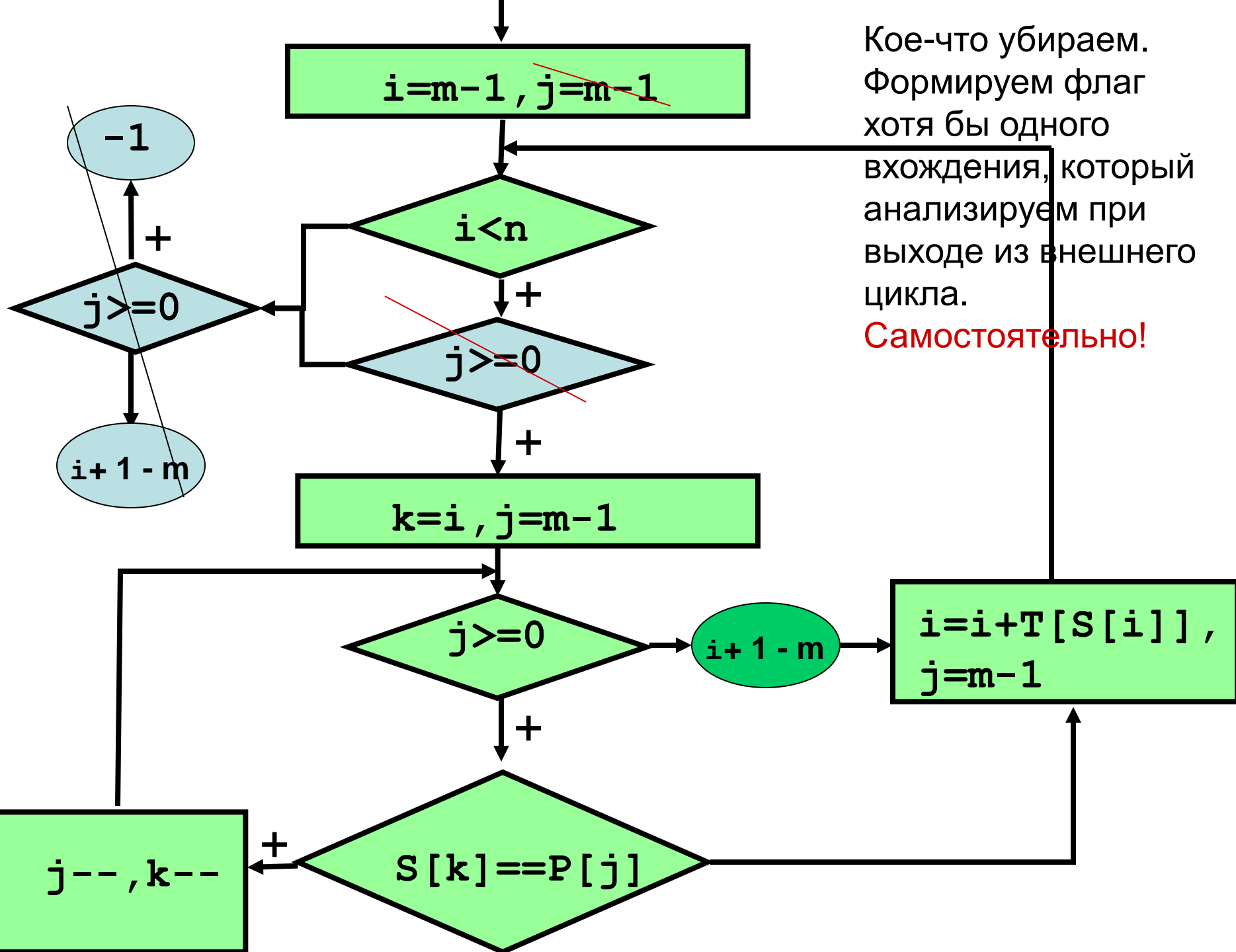
БМ-поиск





Замечание.

Для нахождения всех вхождений образца алгоритм модифицируется немного: после нахождения образца выводится позиция его начала (в нашей блок-схеме после внутреннего цикла) и поиск продолжается аналогично.



Пример из 1.



s = “во дворе трава на траве дрова”

Таблица

p = “дрова”

| — | ... | а | б | в | г | д | ... | о | р |
|---|-----|---|---|---|---|---|-----|---|---|
| 5 | ... | 5 | 5 | 1 | 5 | 4 | ... | 2 | 3 |

s = во дв⁺¹оре трава на траве дрова

p = др^ова

s = во дв⁺²оре тра^е на траве дрова

p = др^ова

е нет в образце!

s = во дв⁺⁵оре трава на траве дрова

p = др^ова

s = во дворе⁺¹ тра^ва на траве дрова

p = др^ова



s = во дворе трава на траве ⁺⁵
r = дрова

т нет в образце!

s = во дворе трава на траве ⁺⁵ дрова
r = дрова

пробела нет в образце!

s = во дворе трава на траве ⁺⁵ _ дрова
r = дрова

s = во дворе трава на траве _ дрова
r = дрова



| | | | | | | | | | |
|---|-----|---|---|---|---|---|-----|---|---|
| _ | ... | а | б | в | г | д | ... | о | р |
| 5 | ... | 5 | 5 | 1 | 5 | 4 | ... | 2 | 3 |



8 новых позиций индекса **i** вместо 25 в
прямом поиске!