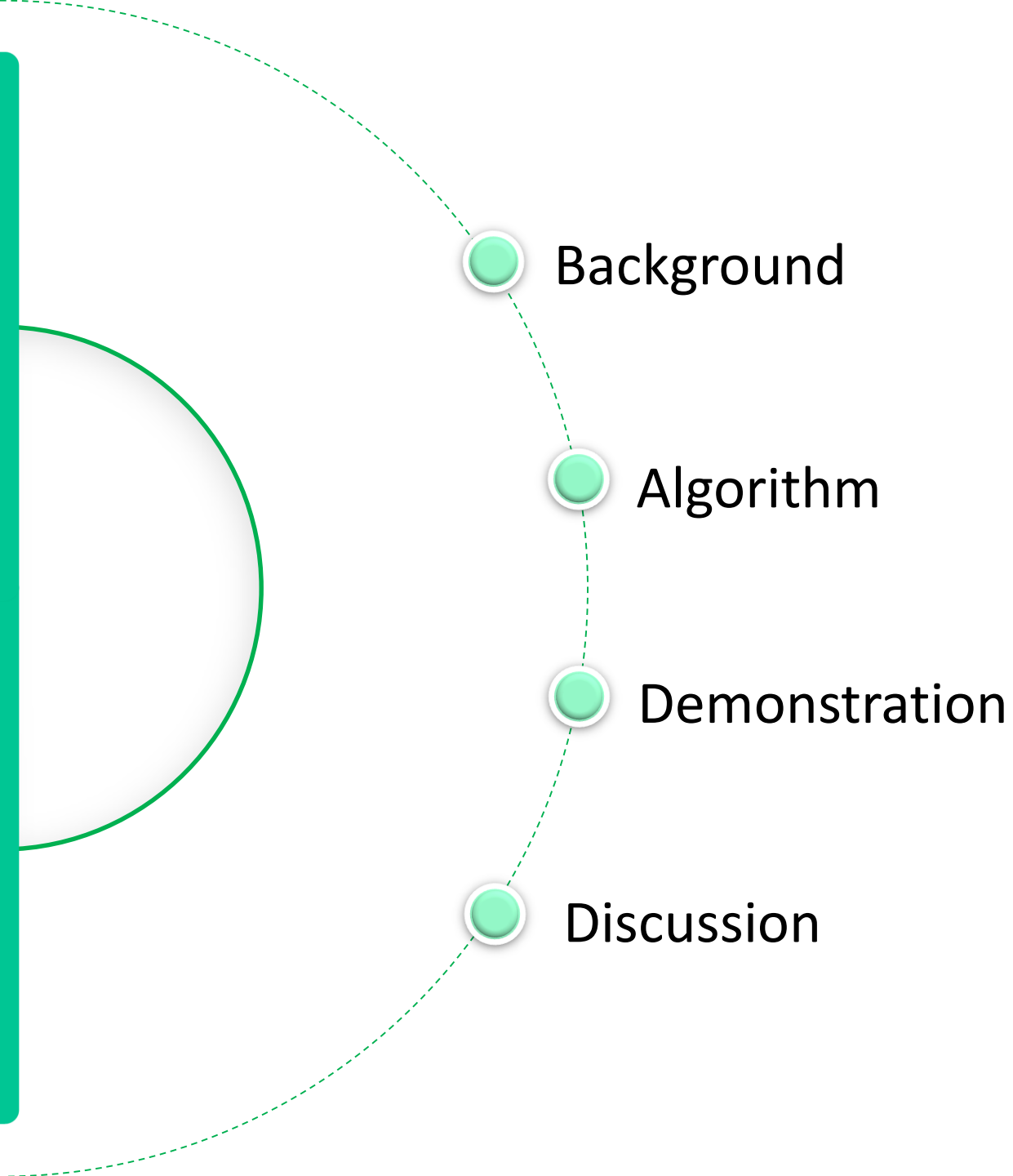




Chinese word segmentation

Presenter: Mingxin Chen

Group: Mingxin Chen, Mingwei Wei, Bin Jiang, Ange Tong



What is Chinese word segmentation?

Background

Algorithm

Demonstration

Discussion



like ... this ...

- | | |
|------------|---------------------|
| • 语言信息处理 | • 语言 / 信息 / 处理 |
| • 鉴萍老师美丽大方 | • 鉴萍 / 老师 / 美丽 / 大方 |
| • 诚实是一种美德 | • 诚实 / 是 / 一种 / 美德 |

Why segment?

Great
Significance!

Background

valuations

The basic of information retrieval, information extraction, information classification and so on.

Demonstration

applications

- Polyphone recognition (多音字识别)
- Text proofreading (文本校对) eg. [于预 >> 干预]

and more ...

.....

Discussion

Challenge

Background

What is challenge of Chinese word segmentation?

Algorithm

To solve the ambiguity, of course!!!

Demonstration

So .. What do you mean ... ambiguity?

Discussion

Like this ..

这个学生会打篮球

- 这个 / 学生 / 会 / 打 / 篮球
- 这个 / 学生会 / 打 / 篮球

OK, I got it.

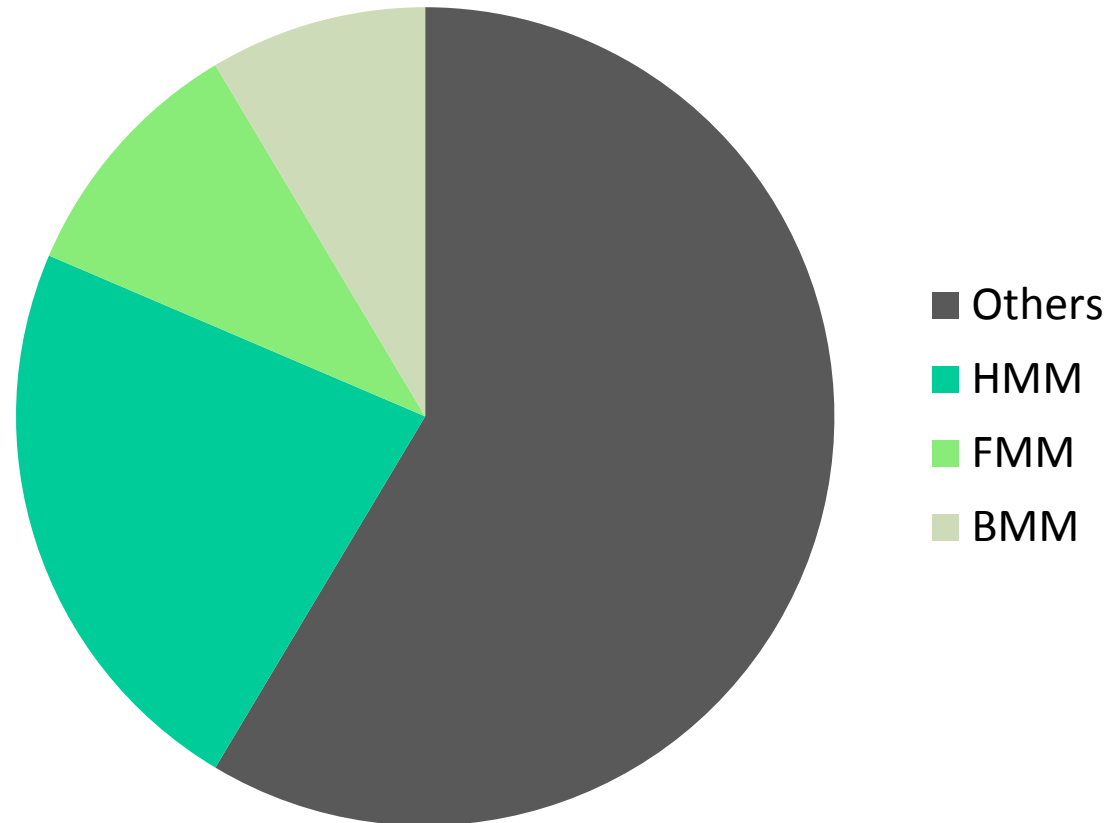
Algorithm classification

Background

Algorithm

Demonstration

Discussion



Algorithms in our solution

Hidden Markov Model Algorithm

Background

We define $C = \{B, M, E, S\}$, $O = O_1 O_2 \dots O_i$

Algorithm

- | | |
|-----|------------------|
| • B | • Begin of word |
| • M | • Middle of word |
| • E | • End of word |
| • S | • Single word |

O1O2O3O4O5O6O7O8O9O10

你/现在/应该/去/幼儿园/了

B BE BE S BME S

Demonstration

Discussion

Give us O , we should calculate C

That means we calculate $\operatorname{argmax}_C P((C_1, C_2 \dots C_i | O_1, O_2 \dots O_i)) = P(C|O)$ (1)

Simplify (1), we get $\operatorname{argmax}_C P(O|C)P(C)$ (2)

Analysis (2), we get final formula:

$\operatorname{argmax}_C P(O_1|C_1)P(O_2|C_2) \dots P(O_i|C_i) * P(C_1)P(C_2|C_1)P(C_3|C_2) \dots P(C_i|C_{i-1})$

A Pure-HMM Segmentation

Background

Algorithm

Demonstration

Discussion

We calculate three matrix

1

initial probability matrix

$$P_i = \text{Count}(C_i) / \sum \text{Count}(C_j)$$

2

state-transition matrix

$$A_{ij} = P(C_j | C_i) = \frac{P(C_i, C_j)}{P(C_i)} = \text{Count}(C_i, C_j) / \text{Count}(C_i)$$

3

emitter matrix

$$B_{ij} = P(O_j | C_i) = \frac{P(O_j, C_i)}{P(C_i)} = \text{Count}(O_j, C_i) / \text{Count}(C_i)$$

Add 1 smooth

$$B_{ij} = P(O_j | C_i) = (\text{Count}(O_j, C_i) + 1) / (\text{Count}(C_i) + N)$$

Viterbi Algorithm

Background

Algorithm

Demonstration

Discussion



Formula derivation

$$\begin{aligned} & \operatorname{argmax}_C P((C_1, C_2 \dots C_i | O_1, O_2 \dots O_i)) \\ &= \operatorname{argmax}_C P(O_1 | C_1) P(O_2 | C_2) \dots P(O_i | C_i) * P(C_1) P(C_2 | C_1) P(C_3 | C_2) \dots P(C_i | C_{i-1}) \\ &= \operatorname{argmax}_C P(O_1 | C_1) P(O_2 | C_2) \dots P(O_{i-1} | C_{i-1}) * P(C_1) P(C_2 | C_1) P(C_3 | C_2) \dots P(C_i | C_{i-1}) * P(O_i | C_i) \\ &= \operatorname{argmax}_C P((C_1, C_2, \dots, C_{i-1} | O_1, O_2, \dots O_{i-1})) * P(C_i | C_{i-1}) * P(O_i | C_i) \end{aligned}$$



$$\begin{aligned} & \operatorname{argmax}_C P((C_1, C_2 \dots C_i | O_1, O_2 \dots O_i)) \\ &= \operatorname{argmax}_C P((C_1, C_2, \dots, C_{i-1} | O_1, O_2, \dots O_{i-1})) * P(C_i | C_{i-1}) * P(O_i | C_i) \end{aligned}$$

Viterbi Algorithm

Background

Algorithm

Demonstration

Discussion

Process

Calculate the initial probability: $P(C1|O1) = P(C1) * P(O1|C1)$

Calculate $P((C1, C2.. Ct|O1, O2.. Ot)) =$
 $\text{argmax}_C P((C1, C2, \dots, C_{t-1}) | O1, O2, \dots O_{t-1})) * P(C_t|C_{t-1}) * P(O_t|C_t),$
using an array path to record the value of C_t when taking the maximum

Get the maximum sequence: $\text{argmax}_C P((C1, C2.. Ci | O1, O2.. Oi)),$
we can get the value of C_i when P taking the maximum.

backtracking:

Using the path array to backtrack and get the hidden sequence.
If the value is E or S then can be divided into a word.

A Pure-HMM Segmentation

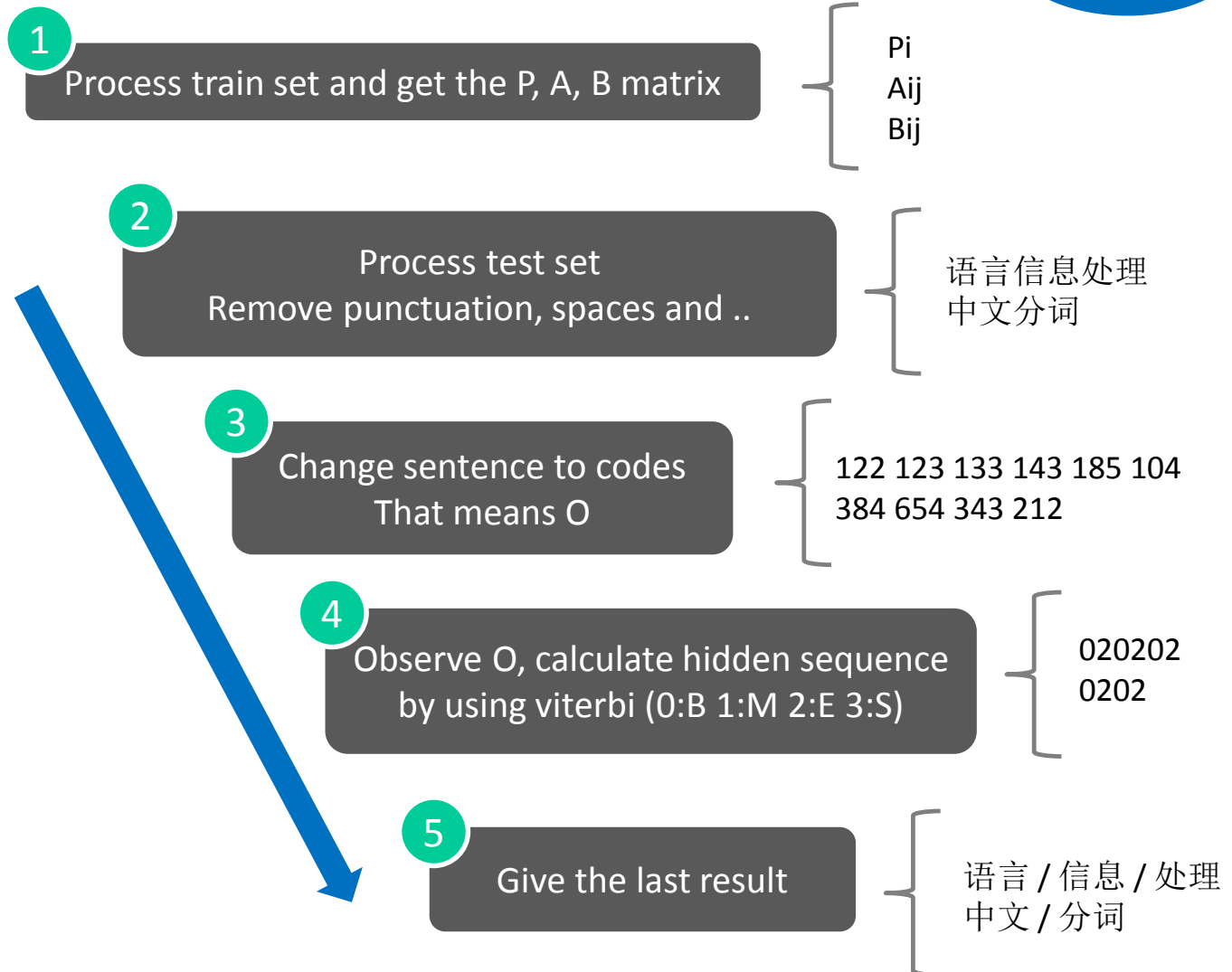
Example

Background

Algorithm

Demonstration

Discussion



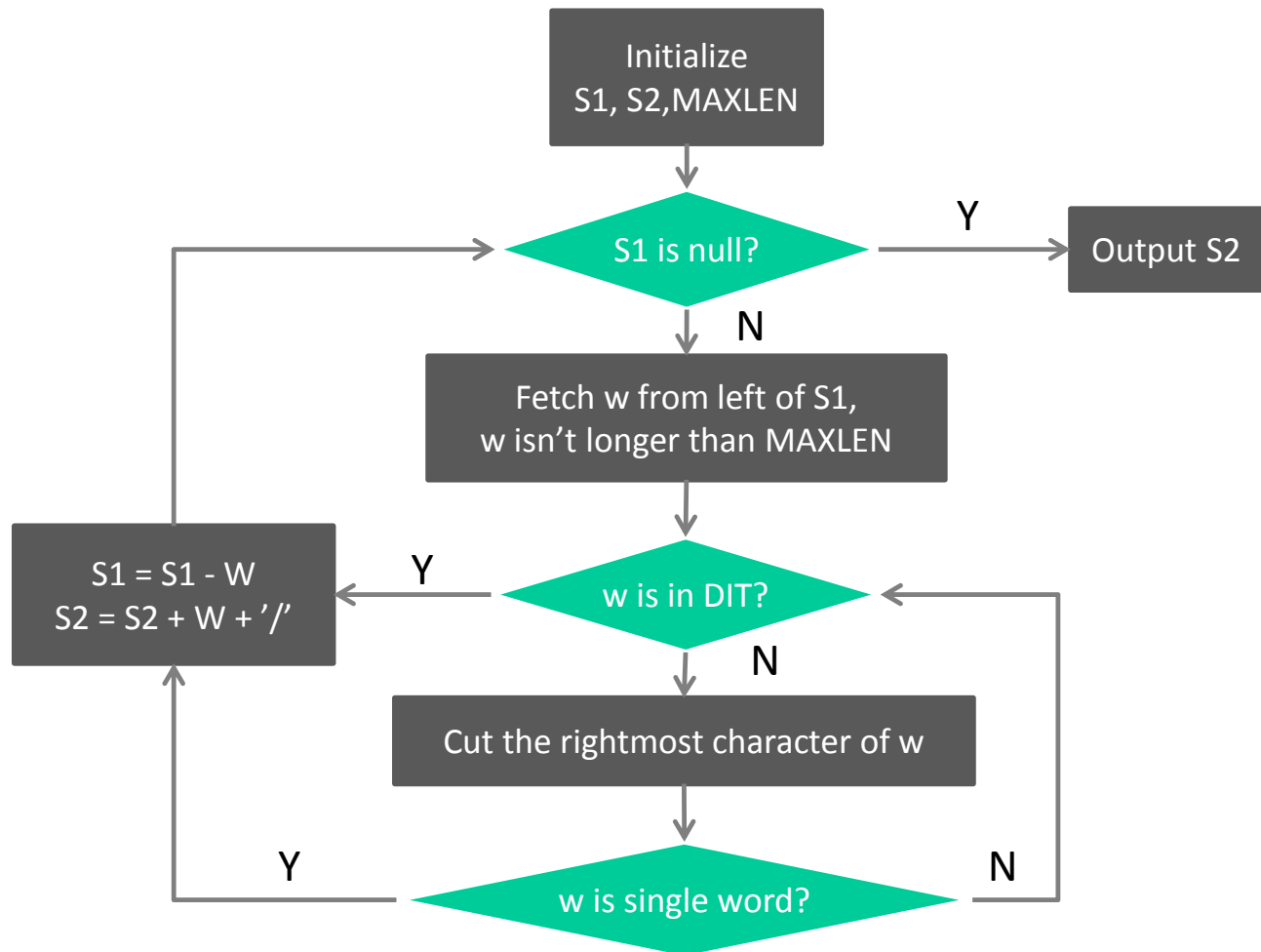
Forward Maximum Matching Algorithm

Background

Algorithm

Demonstration

Discussion



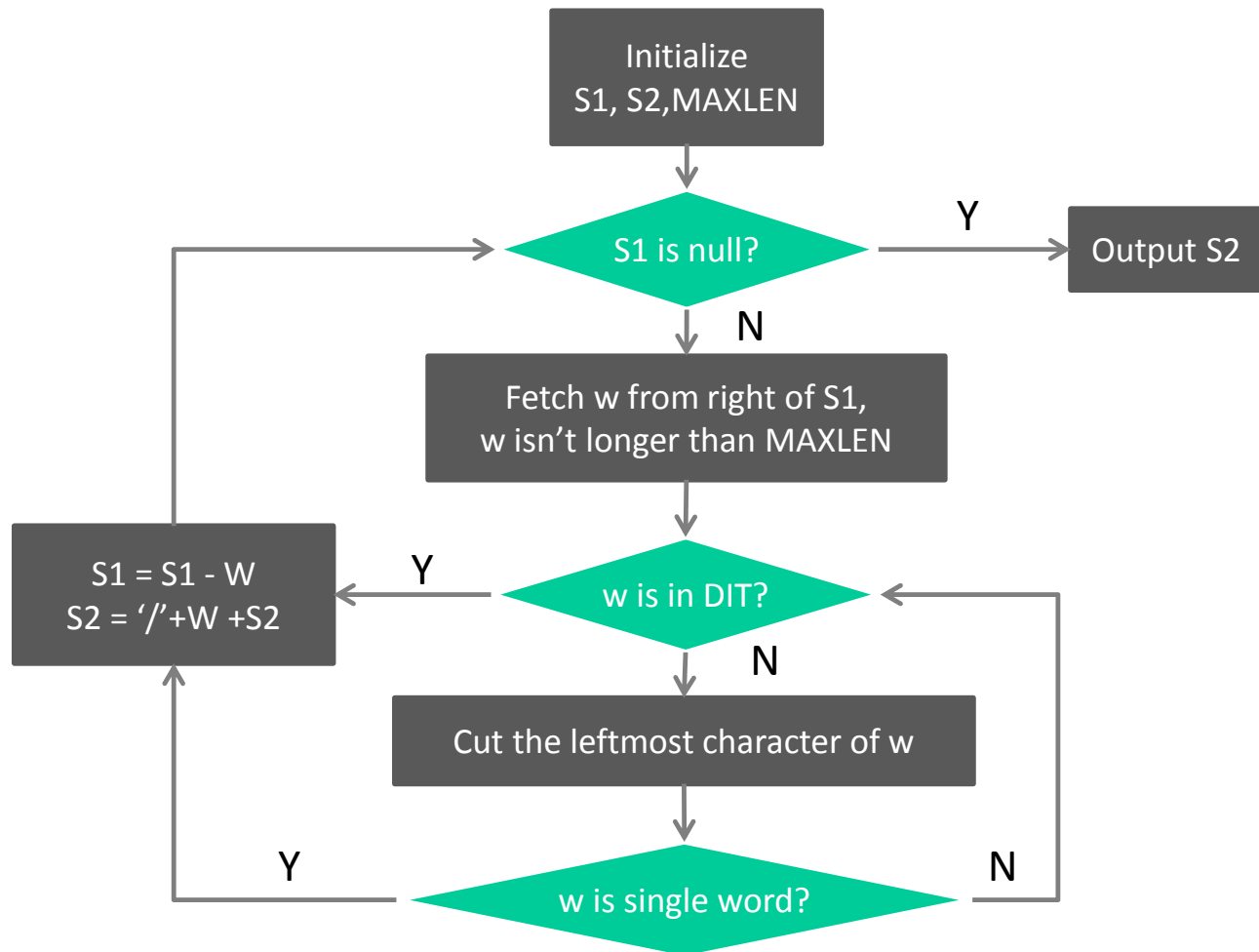
Backward Maximum Matching Algorithm

Background

Algorithm

Demonstration

Discussion



Demo

Background

Algorithm

Demonstration

Discussion



Demo

Background

Algorithm

Demonstration

Discussion



扬 帆 远 东 做 与 中 国 合 作 的 先 行
希 腊 的 经 济 结 构 较 特 殊 。
海 运 业 雄 踞 全 球 之 首 ， 按 吨 位 计 占 世 界 总 数 的 1 7 % 。
另 外 旅 游 、 侨 汇 也 是 经 济 收 入 的 重 要 组 成 部 分 ， 制 造 业 规 模 相 对 较 小 。
多 年 来 ， 中 希 贸 易 始 终 处 于 较 低 的 水 平 ， 希 腊 几 乎 没 有 在 中 国 投 资 。
十 几 年 来 ， 改 革 开 放 的 中 国 经 济 高 速 发 展 ， 远 东 在 崛 起 。
瓦 西 里 斯 的 船 只 中 有 4 0 % 驶 向 远 东 ， 每 个 月 几 乎 都 有 两 三 条 船 停 靠 中 国 港 口 。
他 感 受 到 了 中 国 经 济 发 展 的 大 潮 。
他 要 与 中 国 人 合 作 。
他 来 到 中 国 ， 成 为 第 一 个 访 华 的 大 船 主 。

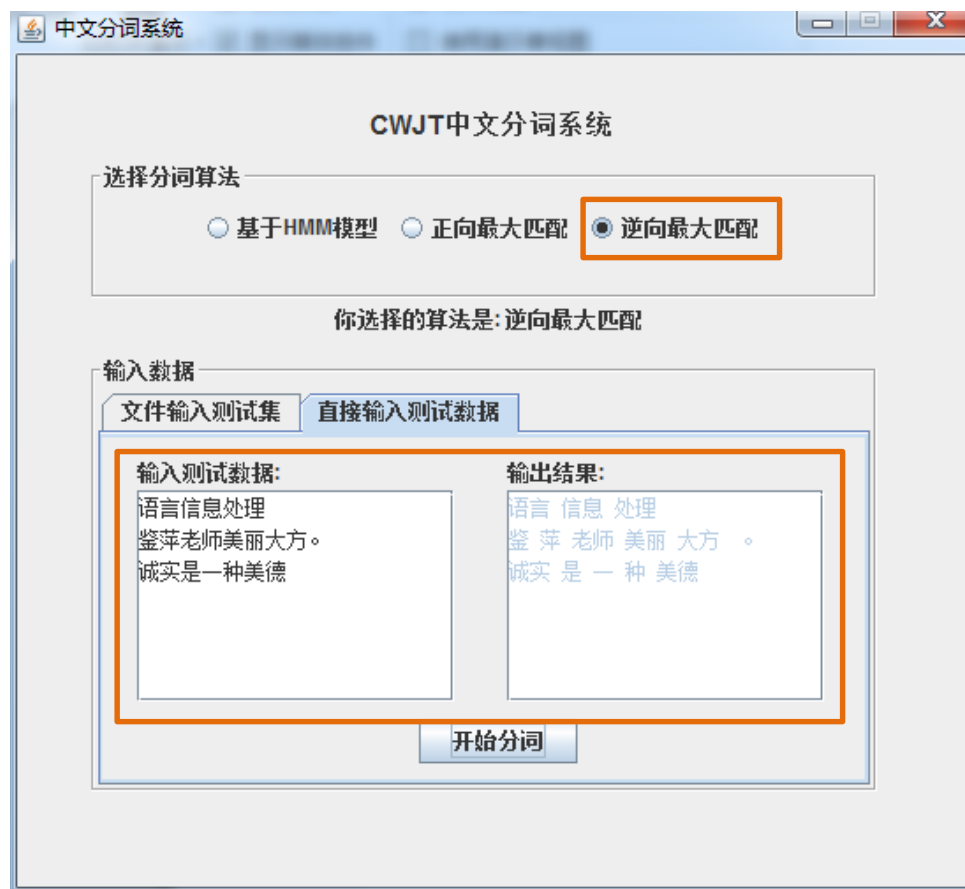
Demo

Background

Algorithm

Demonstration

Discussion



Score of HMM

Background

Algorithm

Demonstration

Discussion

```
151227 INSERTIONS: 2
151228 DELETIONS: 3
151229 SUBSTITUTIONS: 5
151230 NCHANGE: 10
151231 NTRUTH: 45
151232 NTEST: 44
151233 TRUE WORDS RECALL: 0.822
151234 TEST WORDS PRECISION: 0.841
151235 === SUMMARY:
151236 === TOTAL INSERTIONS: 8496
151237 === TOTAL DELETIONS: 4193
151238 === TOTAL SUBSTITUTIONS: 16237
151239 === TOTAL NCHANGE: 28926
151240 === TOTAL TRUE WORD COUNT: 106873
151241 === TOTAL TEST WORD COUNT: 111176
151242 === TOTAL TRUE WORDS RECALL: 0.809
151243 === TOTAL TEST WORDS PRECISION: 0.778
151244 === F MEASURE: 0.793
151245 === OOV Rate: 0.026
151246 === OOV Recall Rate: 0.431
151247 === IV Recall Rate: 0.819
```

The result is
not very
satisfactory

Score of FMM

Background

Algorithm

Demonstration

Discussion

```
148865 INSERTIONS: 2
148866 DELETIONS: 5
148867 SUBSTITUTIONS: 7
148868 NCHANGE: 14
148869 NTRUTH: 45
148870 NTEST: 42
148871 TRUE WORDS RECALL: 0.733
148872 TEST WORDS PRECISION: 0.786
148873 === SUMMARY:
148874 === TOTAL INSERTIONS: 6134
148875 === TOTAL DELETIONS: 4075
148876 === TOTAL SUBSTITUTIONS: 11818
148877 === TOTAL NCHANGE: 22027
148878 === TOTAL TRUE WORD COUNT: 106873
148879 === TOTAL TEST WORD COUNT: 108932
148880 === TOTAL TRUE WORDS RECALL: 0.851
148881 === TOTAL TEST WORDS PRECISION: 0.835
148882 === F MEASURE: 0.843
148883 === OOV Rate: 0.026
148884 === OOV Recall Rate: 0.305
148885 === IV Recall Rate: 0.866
```

not bad

Score of BMM

Background

Algorithm

Demonstration

Discussion

```
152233 INSERTIONS: 4
152234 DELETIONS: 3
152235 SUBSTITUTIONS: 7
152236 NCHANGE: 14
152237 NTRUTH: 45
152238 NTEST: 46
152239 TRUE WORDS RECALL: 0.778
152240 TEST WORDS PRECISION: 0.761
152241 === SUMMARY:
152242 === TOTAL INSERTIONS: 9502
152243 === TOTAL DELETIONS: 1908
152244 === TOTAL SUBSTITUTIONS: 10211
152245 === TOTAL NCHANGE: 21621
152246 === TOTAL TRUE WORD COUNT: 106873
152247 === TOTAL TEST WORD COUNT: 114467
152248 === TOTAL TRUE WORDS RECALL: 0.887
152249 === TOTAL TEST WORDS PRECISION: 0.828
152250 === F MEASURE: 0.856
152251 === OOV Rate: 0.026
152252 === OOV Recall Rate: 0.197
152253 === IV Recall Rate: 0.905
```

Recall rate
is high.

Thanks !
