中文分词。

概念:由于中文的诗为基础的节转起,阅答之间从有明显的压力形化.图如零度的闪光海回标的.

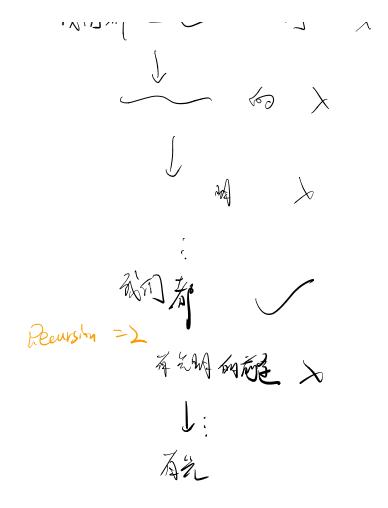
楚子洞图为词。

机械的词络

的最大正面征配名(14M. Maximum Matching Medal). Step 1: 取组切为设施到力 m个字符为正通字 程,m为调图中最长阅读人类。

安文:重成词界这些正规。正视成功,必谓, 为、海底市一个问切分出来。 老么成功,则将实在是几个 字子好,重并正知,更如切为出际有识

Q.9. 我们都有气服的预道。 X



Precursion ...

Result: AND TO THE ME

Distributed (Reserve Maximum Matching Method).

RMM.

Sill to Sind 5 Min Ale.

从文档湖面面扫描.每次取最末2:1字符.

(155=).

Tambula, 到提前加一片之

C. .. / " [MNIU , 11 5.

的闪南的 逐步闪电、每199季四的逐步在162.

实际处理时,实践之超过的树桃处堤。经过至了丘挡、河际状境,逐为河南巡游(MM).

(RMM 强和性益5 成分3MM).

①双向匹斯说。

MMS RAM A Riza.

将可干风时用以此与PMM 进约母插切多。 如应就从早极风,则分别之确。 若则投影从最近望。

00)

n-greens. (bi-greens
tri greens
4-greens

分車をみる

的国人