



携程AI助力产品内容化实践

孙玉霞





孙玉霞

携程平台中心AI研发部，目前担任资深算法工程师，负责携程产品内容化中台相关算法的落地实践和优化。

5年+的自然语言处理经验，专注于文本挖掘/匹配/生成方向。2016年加入携程，负责小诗机，智能客服等项目。

目录

- 1 马可波罗平台介绍
- 2 主题内容挖掘
- 3 优质内容抽取和生成
- 4 总结

目录

1 马可波罗平台介绍

2 主题内容挖掘

3 优质内容抽取和生成

4 总结

内容化痛点



找主题难

选产品难

挖内容难



AI赋能

马可波罗平台

平台功能

特色发现

榜单对标

主题产品挖掘

主题图片挖掘

智能主题配置

主题优质文本抽取

主题文章挖掘

文章自动挂货

文章评级

微游记生成

短亮点

长推荐理由

首图优选

算法层

事件发现

实体链接

敏感内容发现

情感分析

内容聚合

主题匹配

文本生成

摘要抽取

词法分析

句法分析

NER

知识图谱

图文匹配

图片优美度

图片分类

图片去重

数据层

商品信息

用户图片

用户评论

游记

旅拍

问答

马可波罗平台

主题产品挖掘 - 产品主线的内容化

特色发现

榜单对标

主题产品挖掘

主题图片挖掘

主题文章挖掘

智能主题配置

主题优质文本抽取

文章自动挂货

文章评级

微游记生成

短亮点

长推荐理由

首图优选

标签: 131-油菜花 6-网红 46-雪景 361-民宿 362-奢华 33-购物 286-水上乐园 93-山水 29-温泉 159-古镇 27-情侣 78-亲子 4-亲子 18-滑雪 13-美食 39-樱花 1-海滨

酒店ID: 请输入酒店ID

批量输入

酒店名称: 请输入酒店名称

状态: 不限

评分: 最低评分 ~ 最高评分

标签得分: 最低得分 ~ 最高得分

星级: 不限

销量: 最低销量 ~ 最高销量

价格: 最低价格 ~ 最高价格

来源: 全部

评论数量: 最少评论量 ~ 最多评论量

90天夜间量: 最低近90天夜间量 ~ 最高近90天夜间量

目的地: 全部 全部 全部 全部 添加

请先在上方选择并添加目的地

重置 查询 编辑

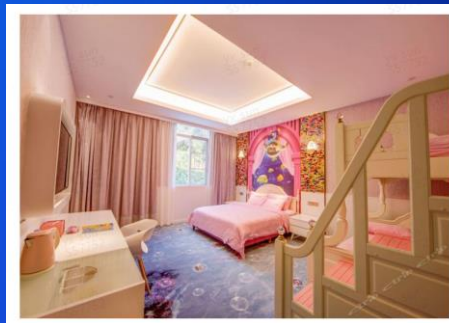
产品ID	产品名称	标签	标签得分	产品链接	图片	短亮点	推荐理由	设施服务	短视频	长视频	GIF	文章	优质评论
+	1700306	珠海长隆企鹅酒店	亲子	57.33	查看		高海洋王国很近, 对带...	高海洋王国很近, 对带...					高海洋王国很近, 家庭...
+	1378371	常州环球恐龙城恐龙主题度...	亲子	54.9	查看		亲子恐龙园游玩首选, ...	亲子恐龙园游玩首选, ...					亲子游首选, 孩子还想...
+	14996879	三亚海棠湾仁恒皇冠假日度...	亲子	54.83	查看		带孩子亲子首选, 海滩...	带孩子亲子首选, 海滩...					适合带孩子去, 儿童泳...
+	419382	杭州良渚君澜度假酒店	亲子	54.05	查看		亲子项目丰富, 还能泡...	亲子项目丰富, 还能泡...					依然不错, 亲子游首选...
-	710691	广州亚特兰酒店	亲子	46.67	查看		小朋友的天堂, 酒店有...	小朋友的天堂, 酒店有...					环境很适合亲子出行...

马可波罗平台

主题产品挖掘 - 产品主线的内容化

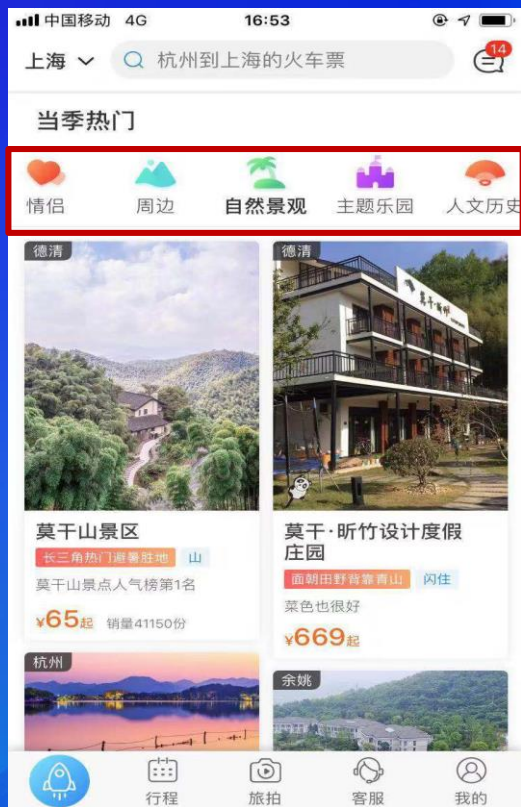
亲子

产品ID	产品名称	标签
1700306	珠海长隆企鹅酒店	亲子
1378371	常州环球恐龙城 恐龙主题度...	亲子
14996879	三亚海棠湾仁恒 皇冠假日度...	亲子
419382	杭州良渚君澜度 假酒店	亲子
710691	广州亚特兰酒店	亲子

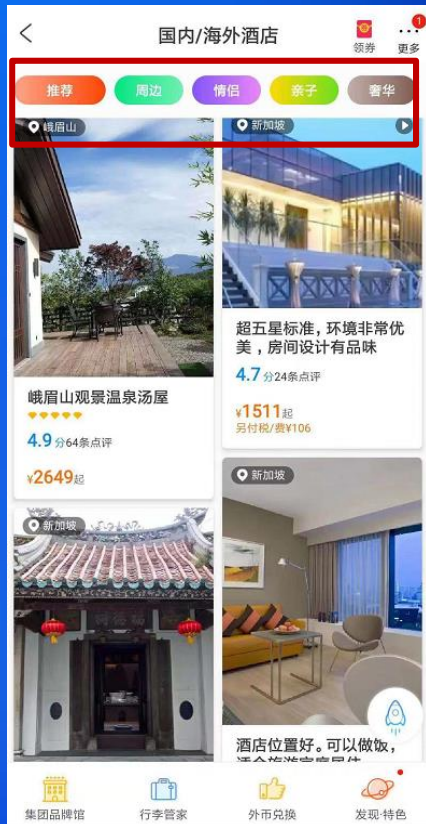


1. 这是一家性价比很高的亲子酒店，距离广州长隆动物园非常近，打车10分钟左右；
2. 亚特兰亲子酒店是小朋友的天堂，酒店里就有游乐园，能让小朋友乐不思蜀；
3. 房间走亲子风格，可爱卡通，小朋友喜欢；
4. 在长隆边上，去长隆玩比较适合，酒店内部装修风格适合亲子游，卫生服务很好；
5. 很适合亲子旅游，亲子房设计合理，床够大，餐厅也不错；

业务应用



首页瀑布流



酒店瀑布流

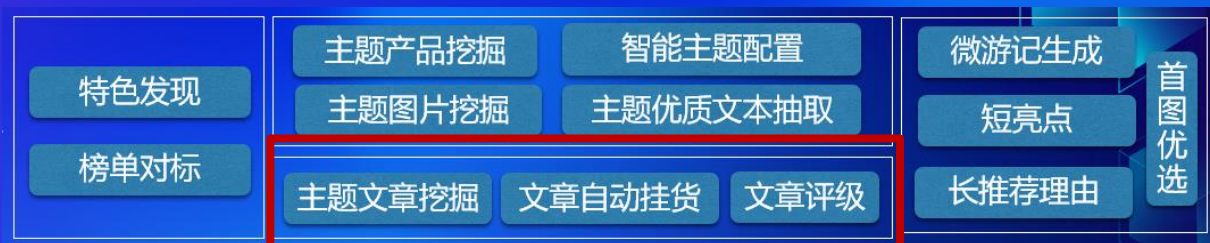


机票瀑布流

主题产品挖掘

马可波罗平台

主题文章挖掘 - 文章主线的内容化



1

- 根据当前话题自动获取最相关的文章。
- 话题文章聚合。

主题文章挖掘

2

- 对图片/文本等内容的量级以及质感进行综合评分，级别粗筛。
- 有效过滤低质文章，缩小候选文章的量。

文章评级

3

- 识别文章中提及到的城市/景点/酒店/餐厅等。
- 提高产品挂载的灵活性。

文章自动挂货

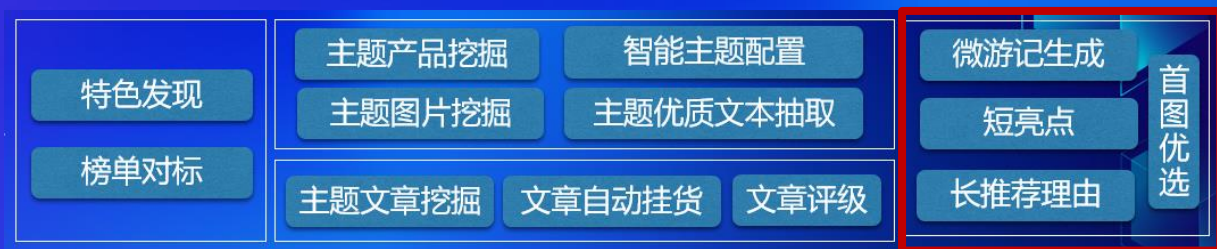
业务应用



攻略 - 话题圈子

马可波罗平台

内容多样化



素材库

- 优质长文本/短亮点抽取
- 首图优选
- 内容改写和生成

微游记生成

- 语义去重
- 标签化
- 图文匹配

文章框架

- 产品-多维度模式
- 主题模式
- 行程模式

业务应用



酒店 - 首页二屏



酒店 - 推荐理由



餐厅 - 推荐理由



IM+ 酒店推荐理由

目录

1 马可波罗平台介绍

2 主题内容挖掘

3 优质内容抽取和生成

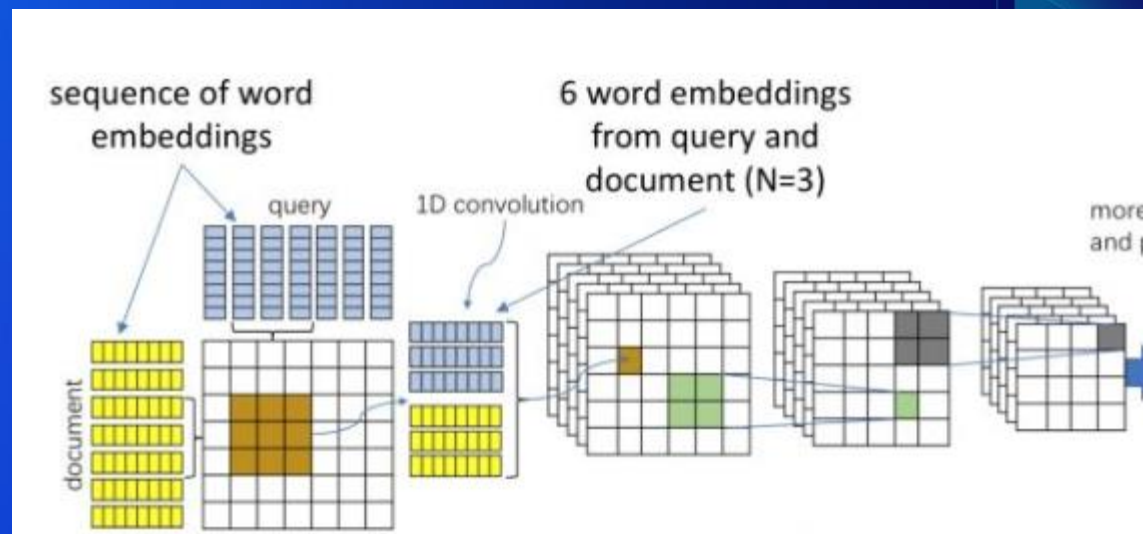
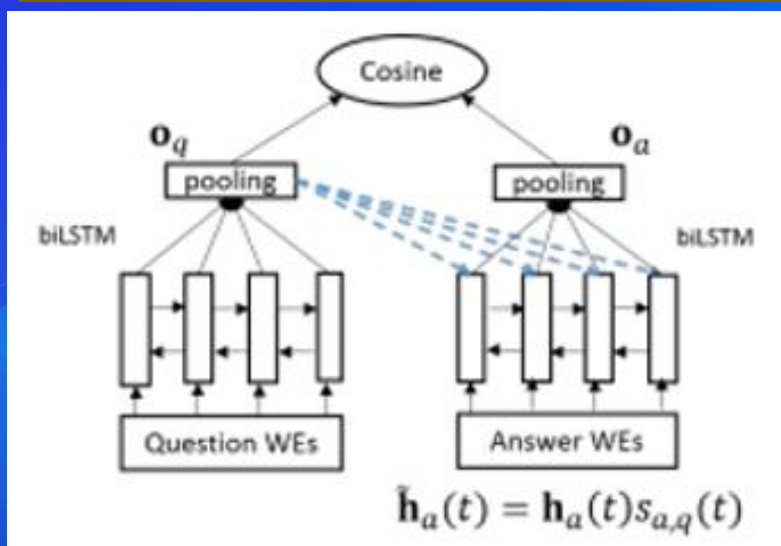
4 总结



2.1 主题产品/文章挖掘

匹配模型简单分类

- 交互式和非交互式
- Loss Pointwise (0, 1分类) 和 pairwise $\max\{0, \text{margin} - (S(Q, D+) - S(Q, D-))\}$



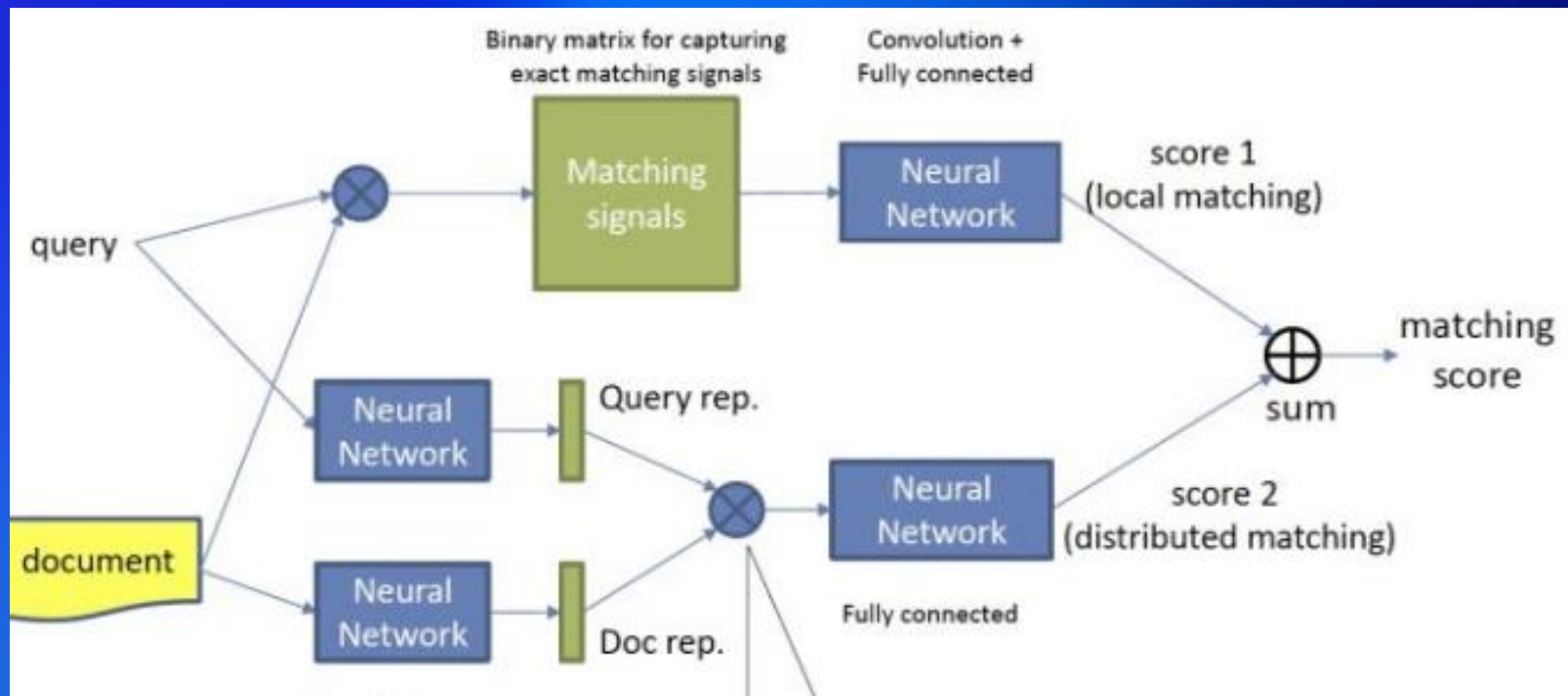
Part 1

- q,a使用同样的网络抽取特征
- 计算相似度前有交互，如 pooling、attention。

Part 2

- q,a对应的embedding后的word计算点积 matching signals
- CNN方法多次卷积
- 2分类模型

2.1 主题产品/文章挖掘

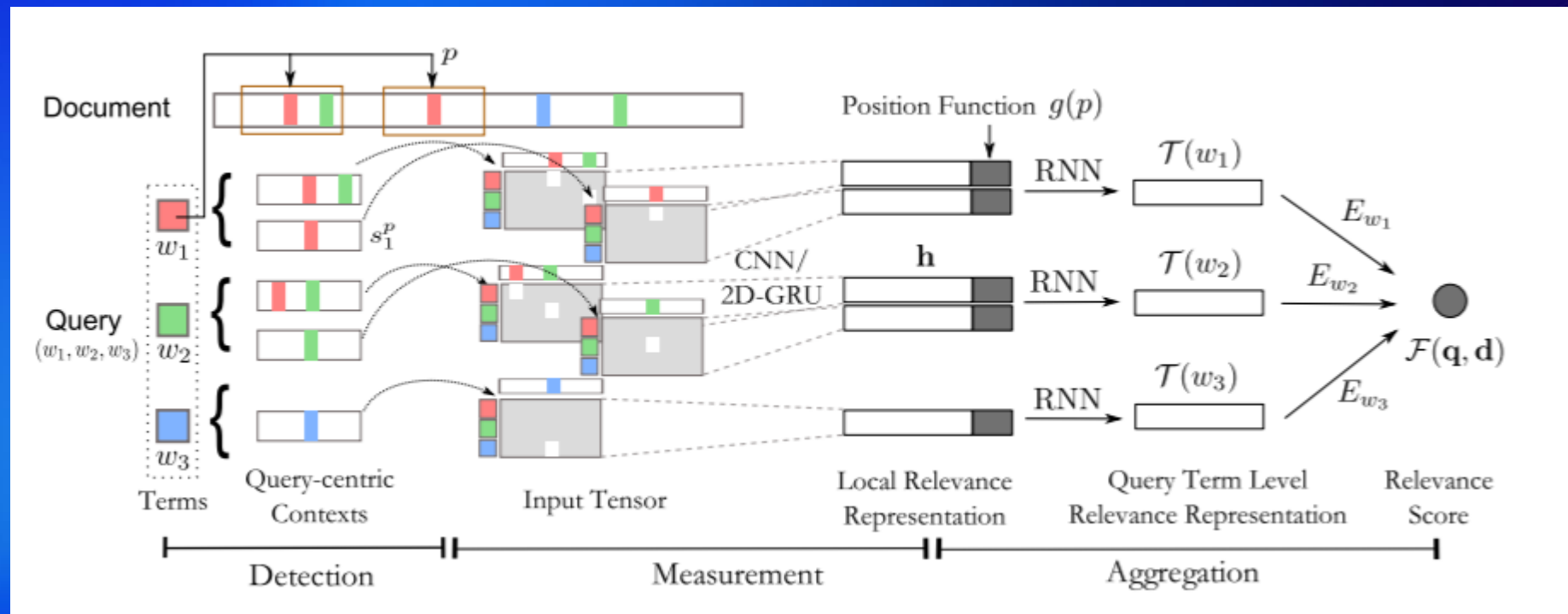


- 短文本之间的匹配
- 结合上文的两种抽取方法进行网络组合。
- 采用pairwise的训练方法，输入为 $(q, a+, a-)$, $a+$ 表示跟 q 匹配, $a-$ 不匹配。
- pointwise 每个输入为 $q, a, label$, $label$ 表示是否匹配。
- $Loss = \max\{0, \text{margin} - (S(Q, D+) - S(Q, D-))\} + \text{loss}$ (分类模型交叉熵损失)

Learning to Match using Local and Distributed Representations of Text for Web Search

2.1 主题产品/文章挖掘

- 短文本-长文本主题相关性
- 长文本主题不聚焦
- 长文本主题多样性
- 类似于query-document问题



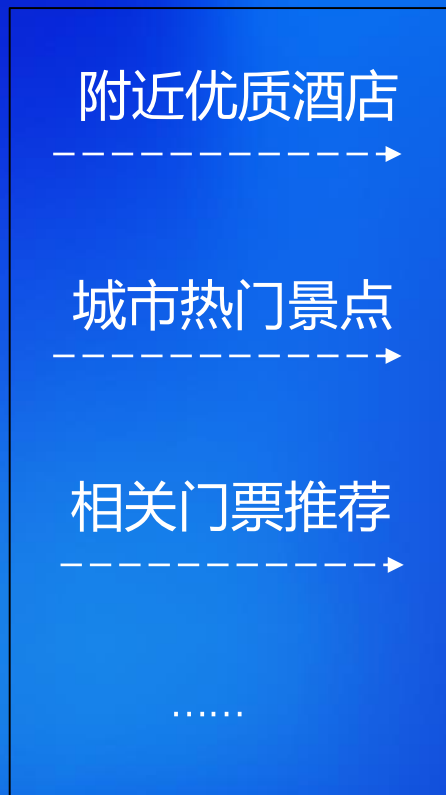
2.1 主题产品/文章挖掘



2.2 文章自动挂货



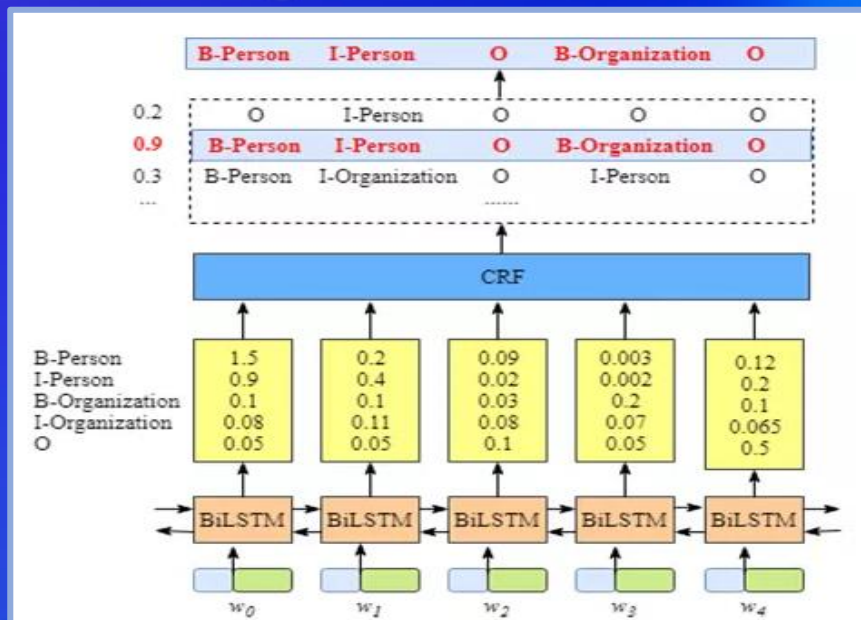
基础实体识别和链接



多场景灵活的挂货策略



2.2 文章自动挂货



1 识别不全

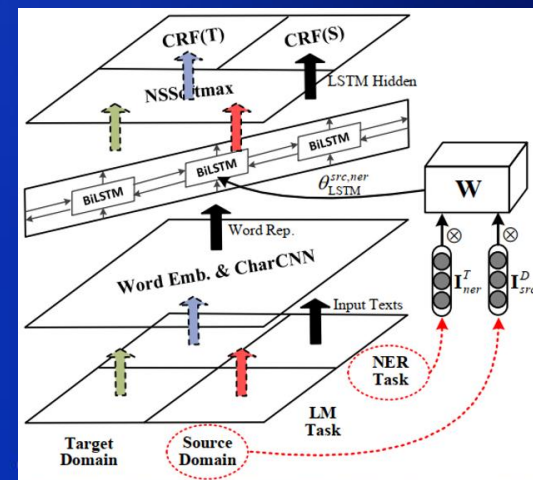
瘦西湖景区

实体补全+反查

2 识别有误

另外酒店还非常人性化地为我们把退房时间延迟.

句法分析+纠错



Cross-Domain NER using Cross-Domain Language Modeling

NER

人名/地名/机构名

类别细分

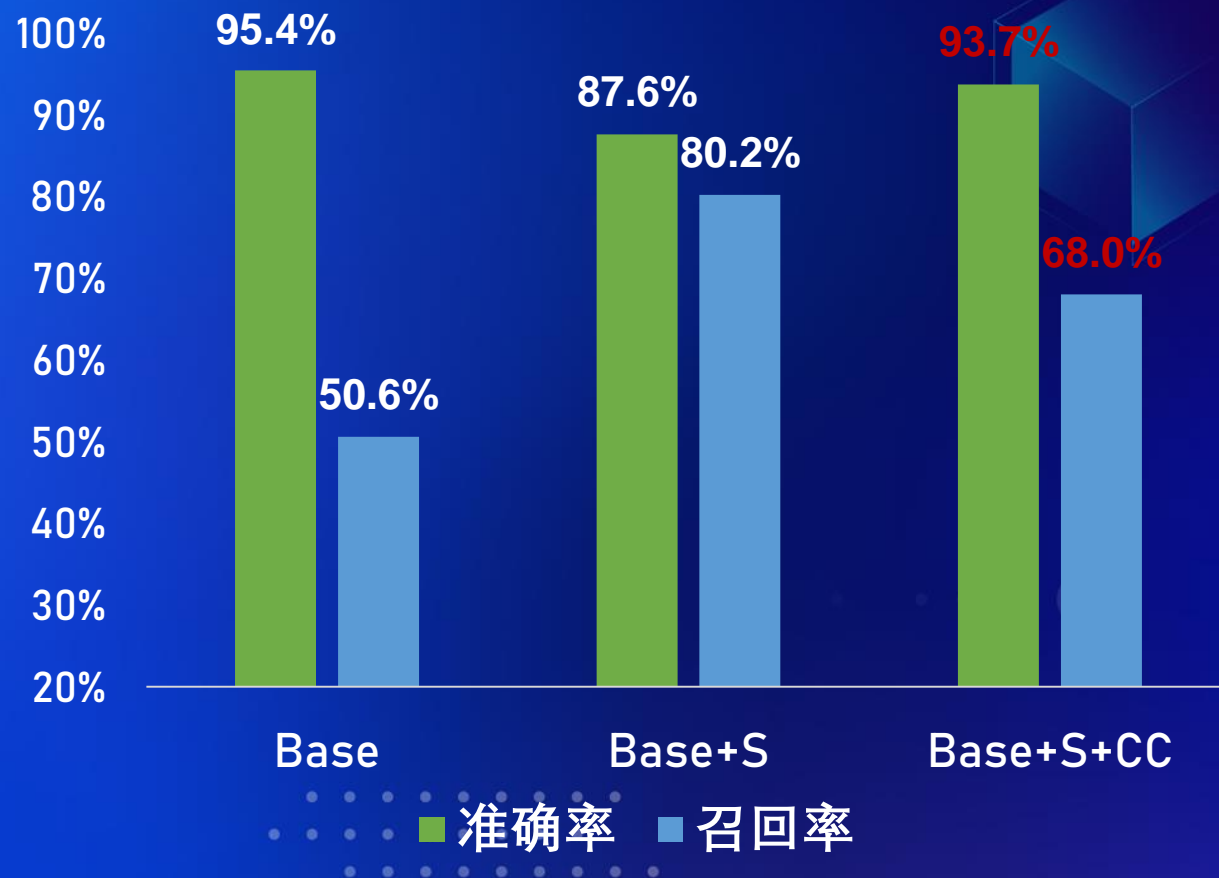
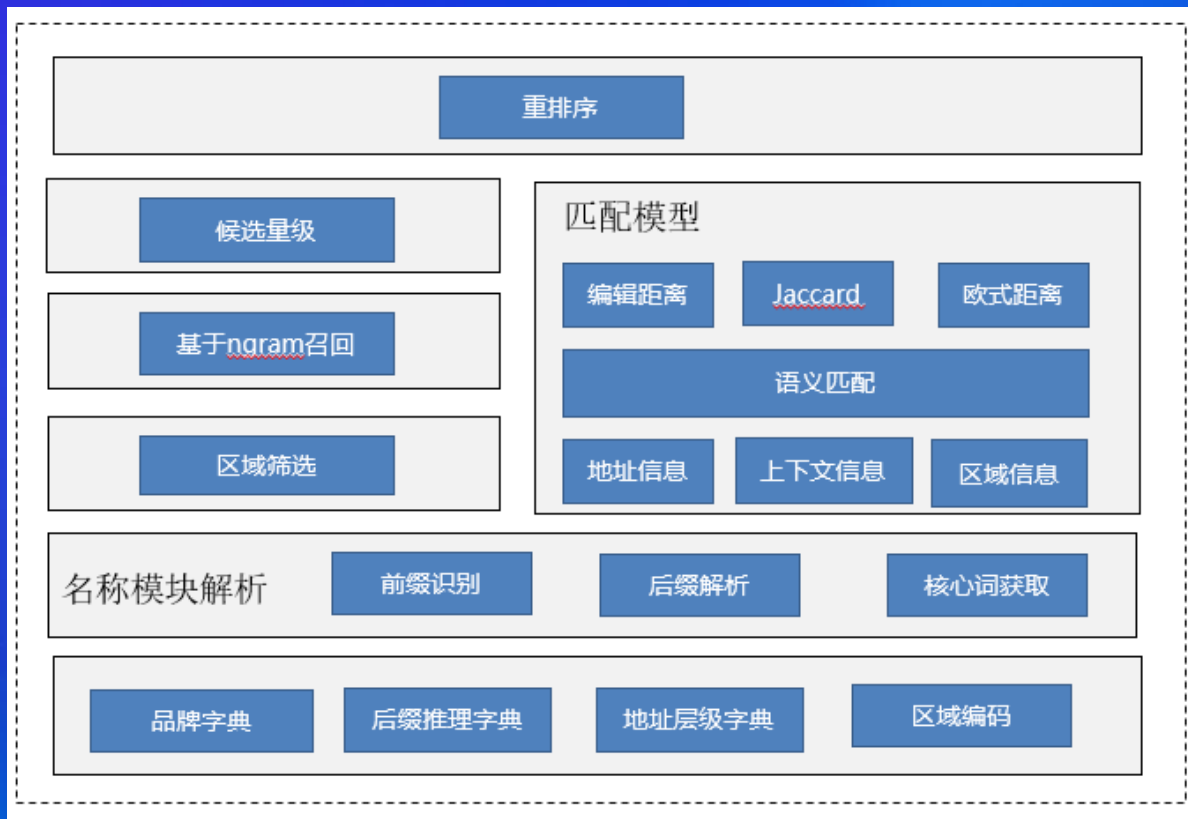
城市/景点/酒店

链接校验

反向校验

2.2 文章自动挂货

实体链接算法框架



2.3 主题图片挖掘

用户自定义主题

手动设置关联



153标签类别

层级类别

草原

梯田

湖泊

.....

动物园

徽派建筑

大象

企鹅

海豚

.....

目录

1 马可波罗平台介绍

2 主题内容挖掘

3 优质内容抽取和生成

4 总结

文章自
动评级

优质文
本自动
抽取

文本生
成

首图优
选

3.1 文章自动评级

选优

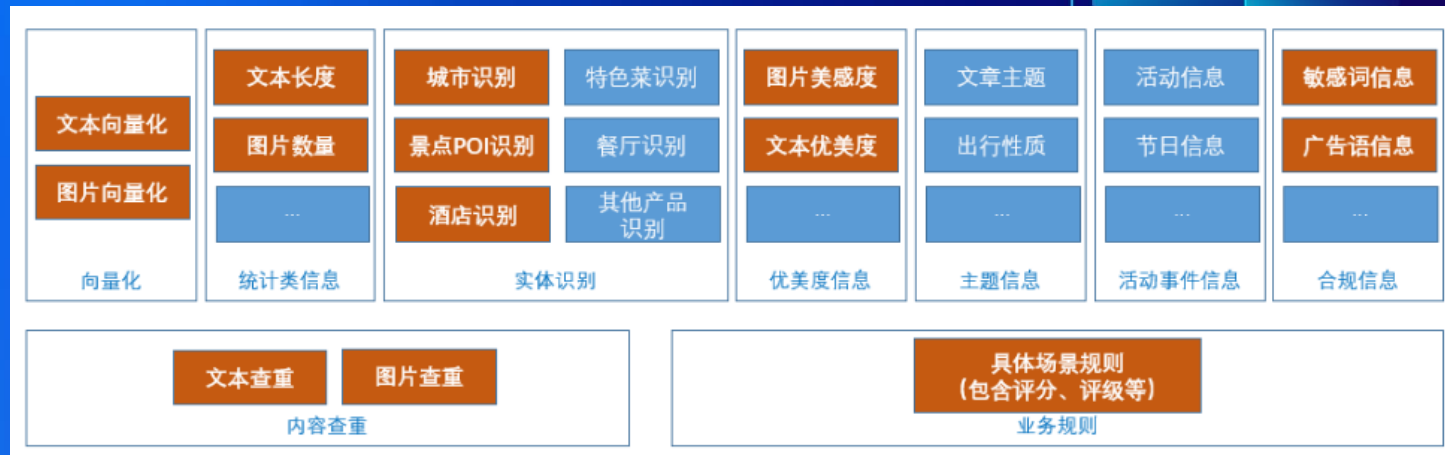
文章内容复杂

标准不一，难以定义

低质筛除

关键点的检测

标准较为清晰，单一



- 敏感内容识别
 - 广告识别
 - 负面情感识别
 - 旅游领域识别
-

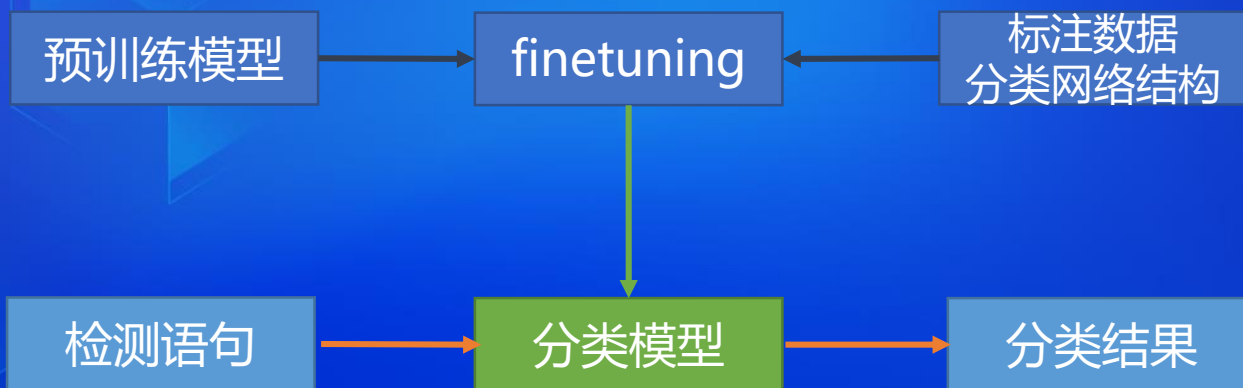
3.1 文章自动评级

论文：《BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding》

- 用Transformer方法进行特征提取，可同时利用左侧和右侧的词语信息
- Mask（遮挡%15）技术在语言模型上应用
- 能够输出word level 和sentence level
- 下游任务对接方便

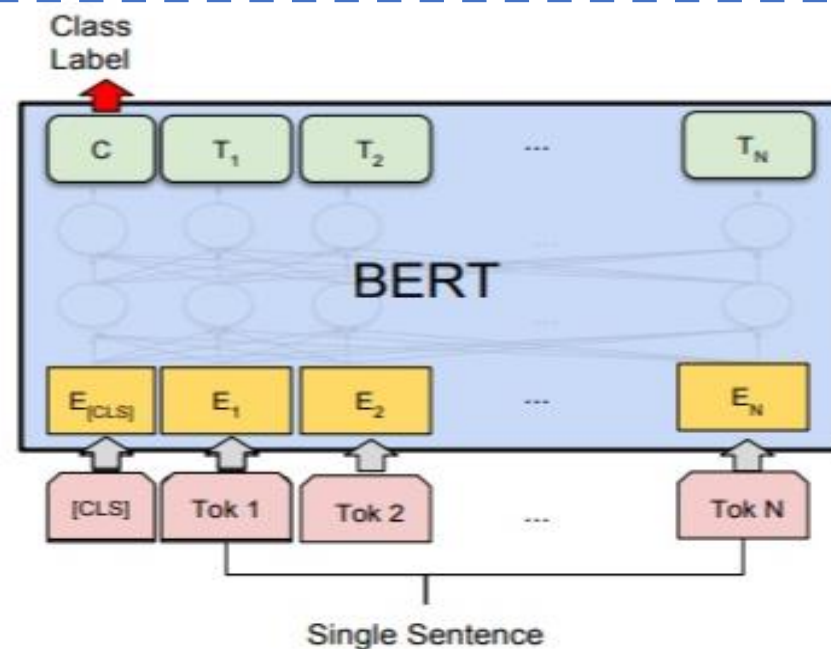
使用方法：

BERTBASE (L=12, H=768, A=12, Total Parameters =110M)



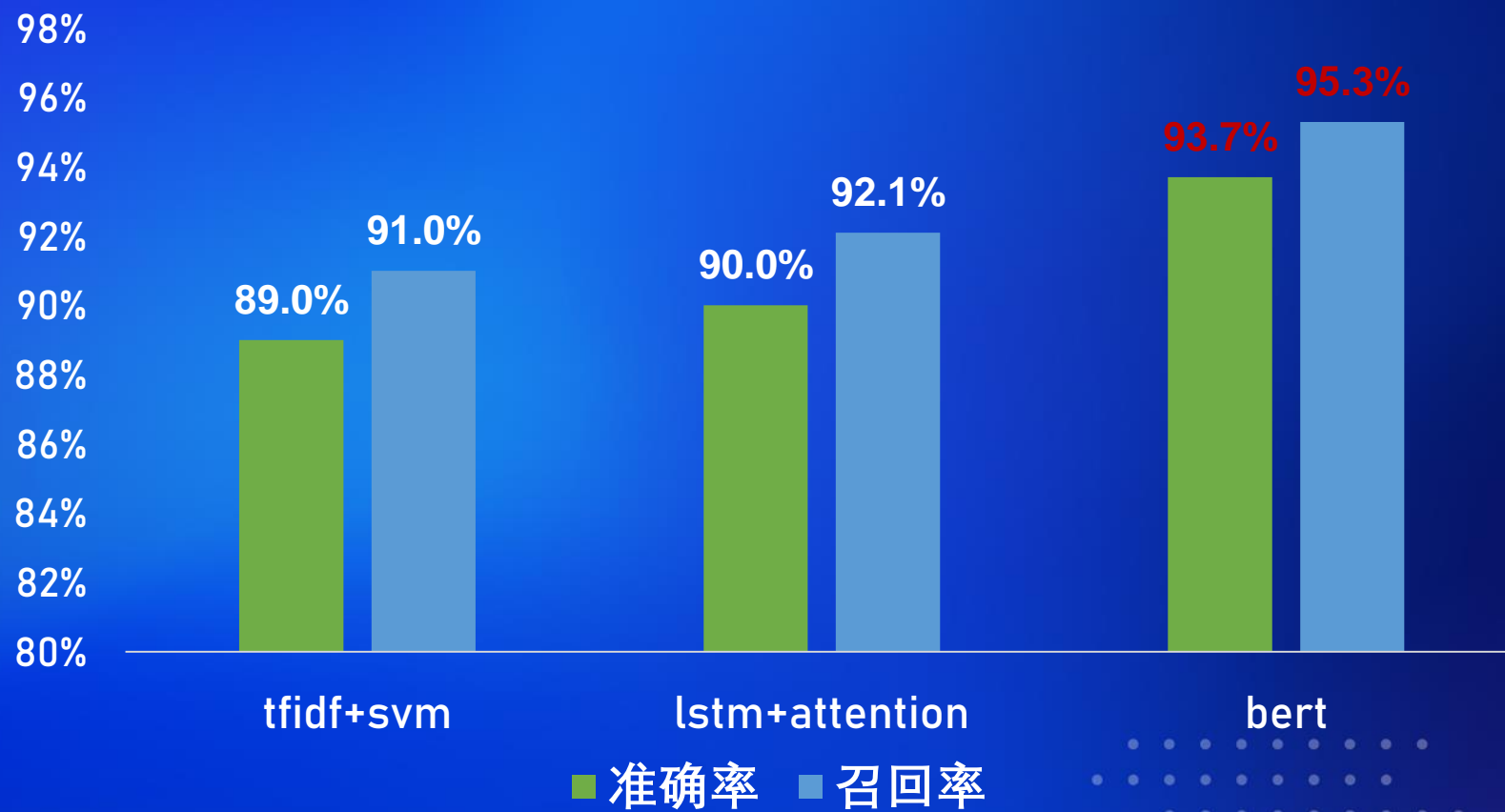
1. Token embedding 表示当前词的embedding
2. Segment Embedding 表示当前词所在句子的index embedding
3. Position Embedding 表示当前词所在位置的index embedding

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	E_{ring}	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}



3.1 文章自动评级

情感模型-不同模型效果对比图



3.2 优质文本内容抽取



3.2 优质文本内容抽取

语句基础信息

- 信息熵
- 词性占比
- 依存句法分析

语句层面

产品特征识别

- 实体识别
- 知识图谱

产品层面

场景维度评价

- 规定各场景（酒店，美食，景点）评价维度
- 建立维度的分类模型

场景层面

内容丰富度

3.2 优质文本内容抽取

词法句法
内容丰富度

- 信息熵:

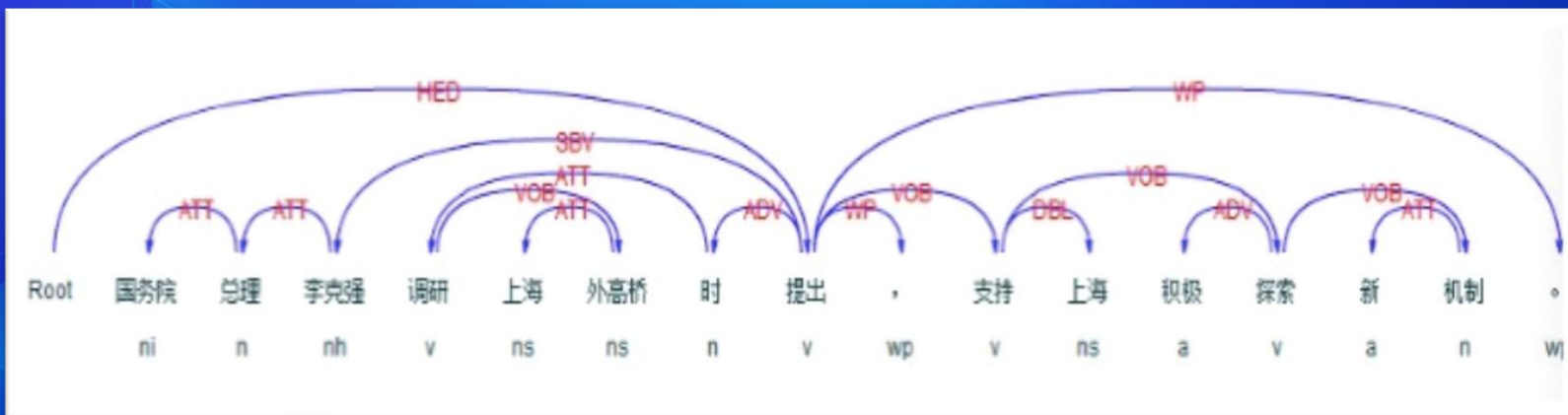
$$H_s = \sum_{i=1}^n p_i I_e = - \sum_{i=1}^n p_i \log_2 p_i$$

- 词性分析:

考虑到推荐理由的结果大概率的同时包含名词和形容词等词性占比。

- 依存句法分析:

重点分析某些关系类型的分布, SBV, ATT, VOB等等。



关系类型	Tag	Example
主谓关系	SBV	我送她一束花 (我 <-- 送)
动宾关系	VOB	我送她一束花 (送 --> 花)
间宾关系	IOB	我送她一束花 (送 --> 她)
定中关系	ATT	红苹果 (红 <-- 苹果)
状中结构	ADV	非常美丽 (非常 <-- 美丽)

3.2 优质文本内容抽取

信息点
内容丰富度

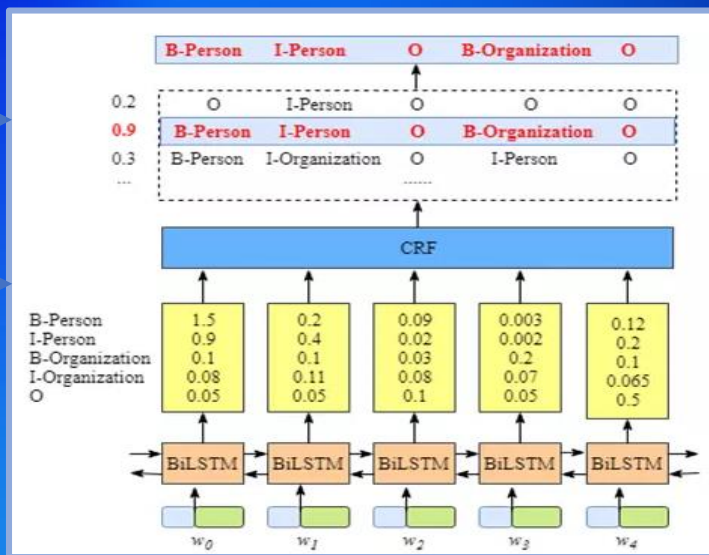
美食 (菜品, 口味等)

酒店 (设施, 交通等)

景点 (名称, 等)

动态提取
实体识别模型

BiLSTM-CRF



**小厨:

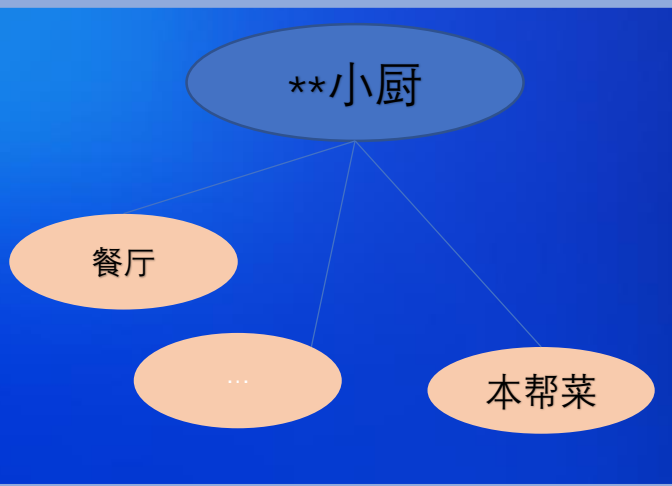
枣泥发糕, 糕有浓浓的荷叶香。
商户详情上说是本帮菜, 干锅
花菜强烈推荐。

识别结果:
干锅花菜, 枣泥发糕

识别结果:
本帮菜

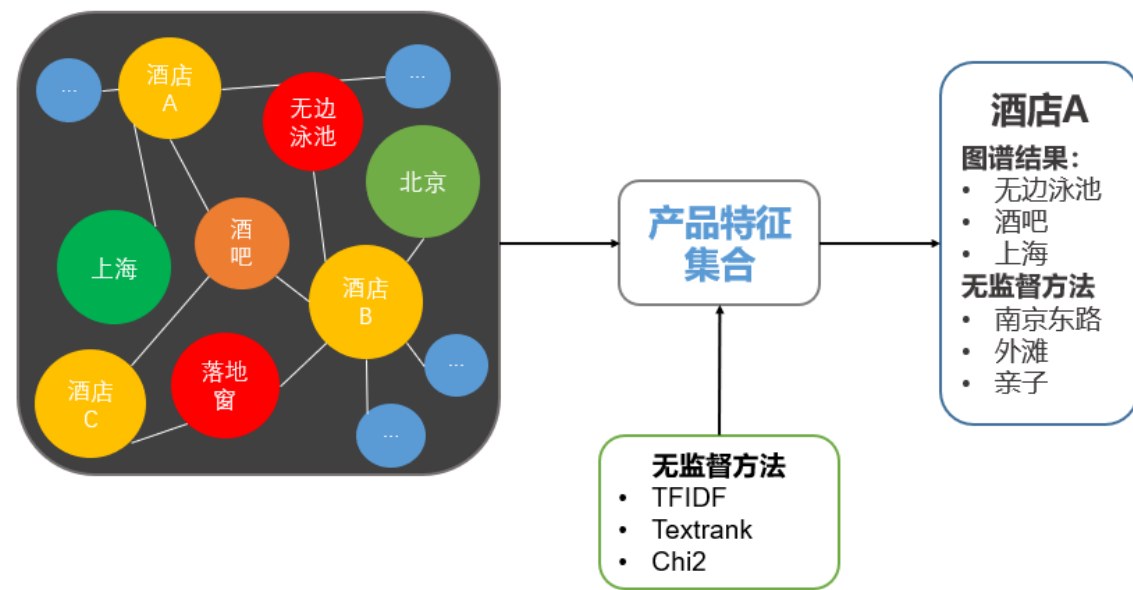
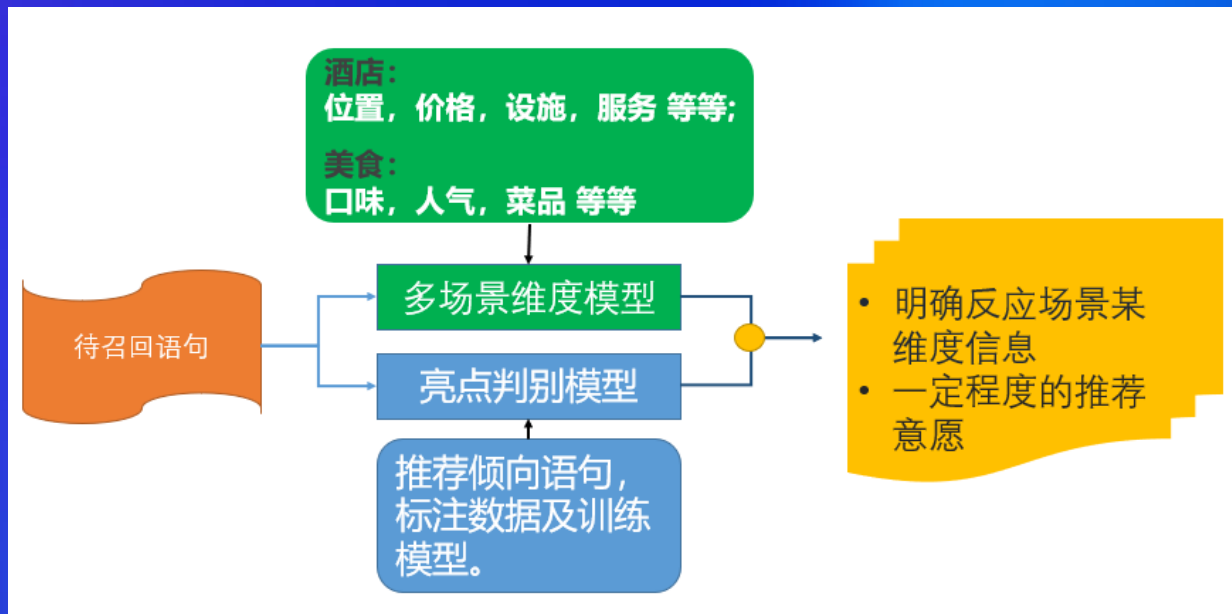
识别结果:
干锅花菜,
枣泥发糕,
本帮菜

静态提取
知识图谱



3.2 优质文本内容抽取

场景及特色
内容丰富度



3.2 优质文本内容抽取

早餐很好，菜品很好
房间隔音，喜欢沐浴用品
环境优美充满中西方底蕴
早餐丰富，位置很好
早餐品类丰富
地点很好历史悠久

重排前

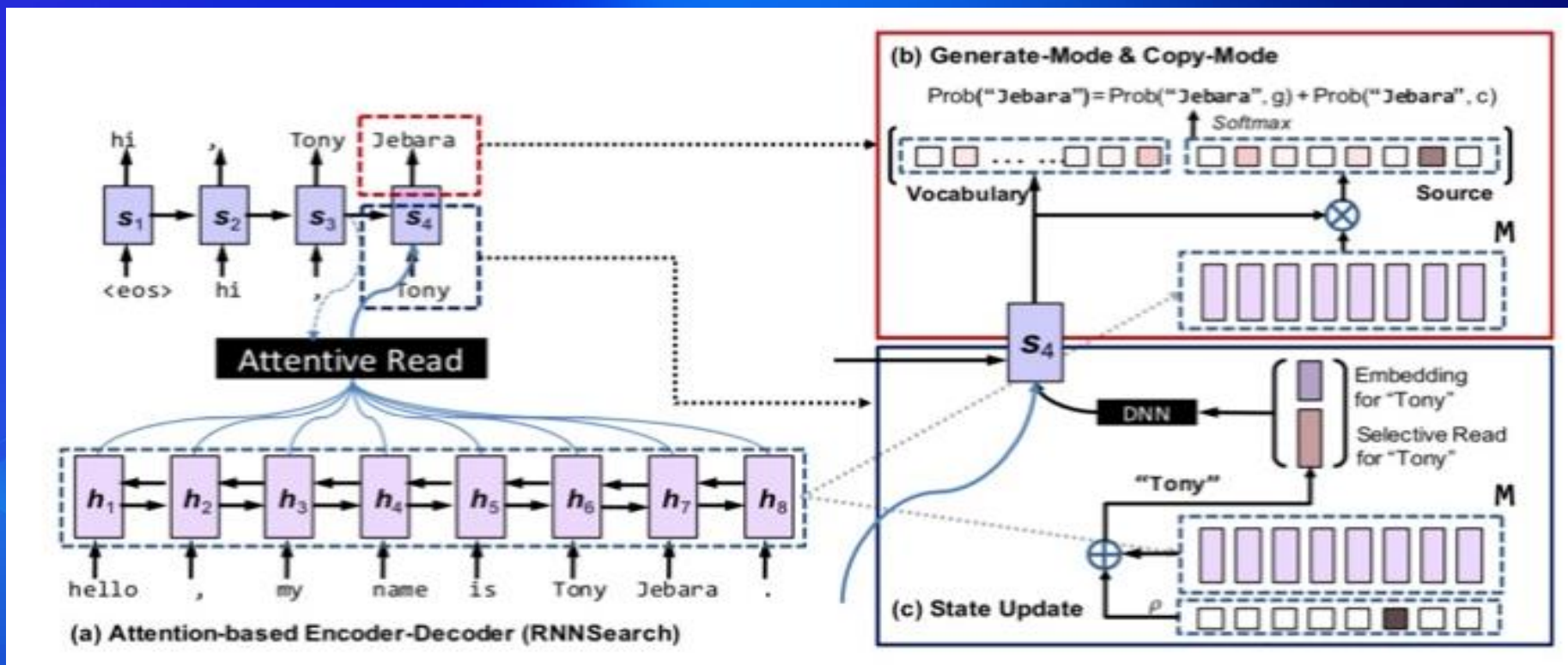
内容丰富度重排

9楼的行政酒廊安静，餐点精致
阳台酒吧是观赏外滩风景的绝佳位置
环境优美充满中西方文化底蕴
下午茶很好，爵士酒吧很喜欢
早餐不错，品类丰富，尤其是西点做得好
一楼的茉莉酒廊值得去

重排后

**饭店-效果展示

3.3 内容生成



优点:

- Copy机制, 较好处理 OOV 问题。
- 场景模式单一的化, 能够较好学习到。

缺点:

- 容易出现以偏概全的问题

独行在孔庙和国子监博物馆

面试结束后顺便去孔庙、国子监转了一圈放松心情走在古香古色的院子里感觉一切忧虑都云消雾散了#孔庙和国子监博物馆[地点]#

网红地洪崖洞

网红地打卡~洪崖洞洪崖洞重庆美食昨晚去了一趟洪崖洞，人真的超级多，不过景真的很美，还是很值得去的，一边是江景一边是建筑，在灯光的照耀下格外好看。

舌尖上吃日料

又来吃日料啦，这家寿司很棒呢，料足又好吃，那个薯你醉鳗是招牌哦

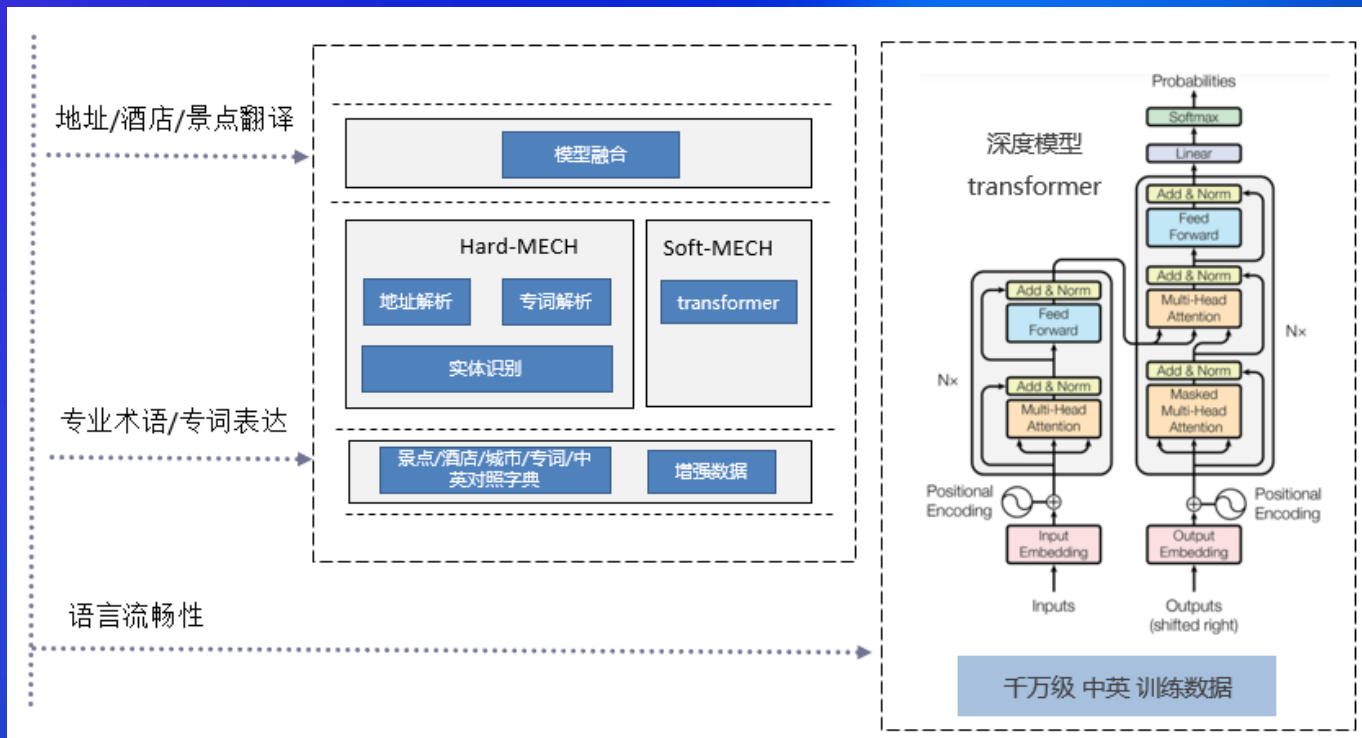
西安昭陵，你不可错过的地方

西安必打卡的地方：碑林+城墙。碑林博物馆馆藏文物十分丰富，可以说是中国书法艺术的百科全书，我最感兴趣的还是【昭陵六骏】，这是任何一部《中国美术史》必敲黑板的重点。

杭州|象山艺术公社

打卡杭州|象山艺术公社时间匆忙，很多线条，人景融合都没拍，建筑物简洁干练，线条框架感实足，关键是人好少鸭#象山艺术公社[地点]#

3.3 内容生成



优点:

- 借助千万级平行数据, 模型更加稳健
- 较为丰富的单向文本数据

缺点:

- 经过两次翻译, 信息双重丢失
- 对实体类型, 长尾词语的表达问题比较严重

去平江路和狮子林都很近, 走路三分钟就能到, 都很方便。
园林景观很美, 闹中取静, 有苏州特色
很安静, 离海信广场很近几分钟走路就可以到。
洗漱用品超级可爱, 精心配备了小朋友的拖鞋、牙刷。

改写

靠近平江路和狮子林, 步行3分钟, 很方便。
园景漂亮, 闹中取静, 有苏州特色。
非常安静, 距离海信广场仅几步之遥。
洗漱用品超级可爱, 配有儿童拖鞋和牙刷。

3.4 首图优选

香港如心海景酒店暨会议中心(L'hotel... (5星) 产品ID: 436824 城市: 香港 最后更新者: kk.he 最后更新时间: 2019-09-25 19:48 状态: 生效

当前图片 尺寸: 2268*1513(查看原图)

A精选 公共区域 外观 房间 其他 休闲 餐饮

☒ 公共区域 (2268*1513) ☒ 公共区域 (2000*1333) ☒ 外观 (1534*1024) ☒ 房间 (2268*1513)

☒ 房间 (2268*1513) ☐ 房间 (3264*1836) ☐ 其他 (2808*1872) ☐ 其他 (2000*1333)

☐ 其他 (2000*1333) ☐ 休闲 (2000*1332) ☐ 休闲 (2268*1511) ☒ 餐饮 (2268*1513)

已选投放图片 ①可删除及拖拽排序 清除全部

返回 失效, 下一个 上一个 生效, 下一个

3.4 首图优选

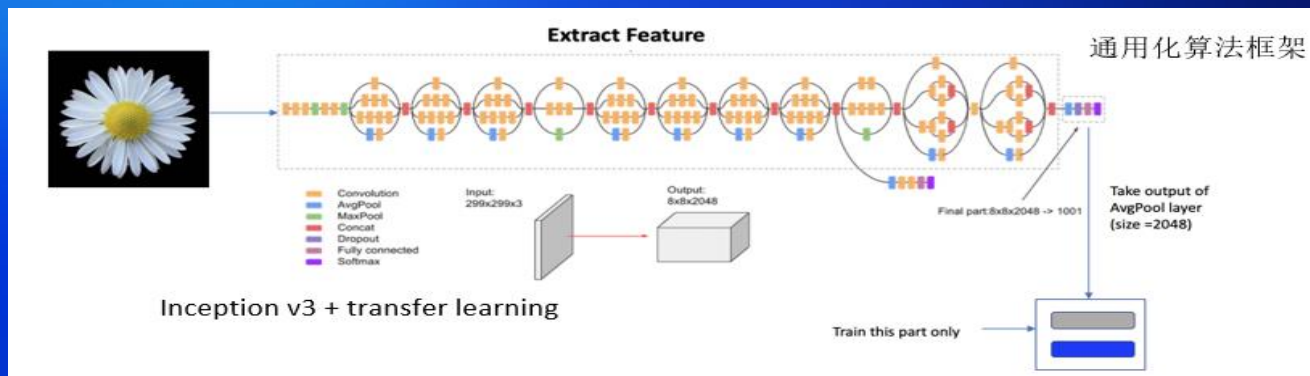


- 多分类：外观 游泳池 大床房 其他 等类别
- 首图场景化类别
- 特定视角，优质度优于覆盖度

1 首图分类

2 美感度

- 二分类 美/丑
- 概率值-图片美感度分值
- 场景样本歧视
- 训练样本的高区分度

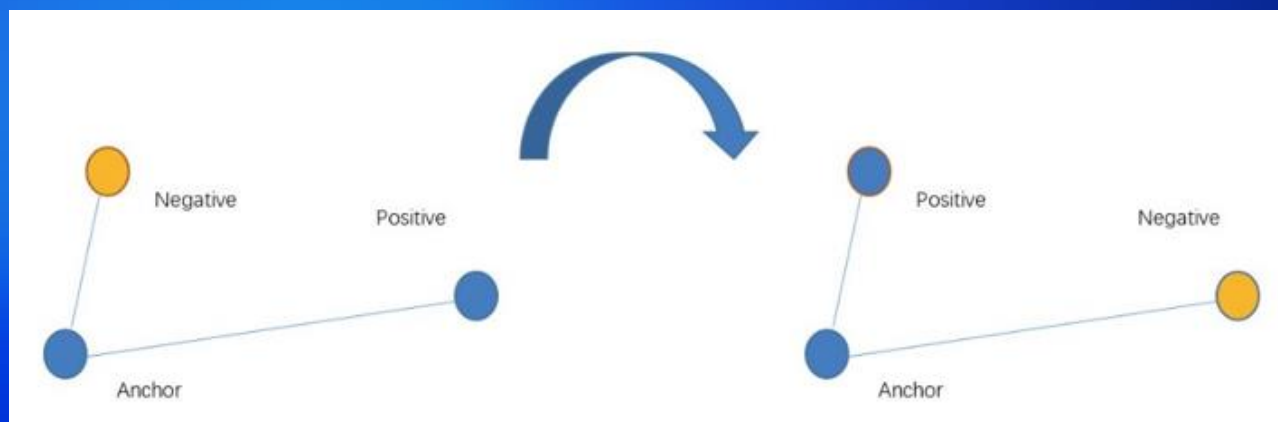
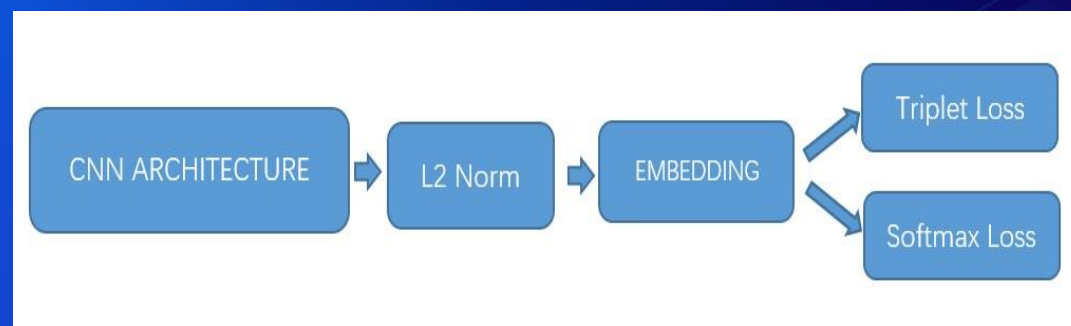


3 图片
去重

Triplet Loss: 最大化异类和同类间的距离差

$$L_t = \text{MAX}(|D_A - D_P| + \text{Margin} - |D_A - D_N|, 0)$$

1. 特征更为稳定
2. 特征维数可控, 特征冗余较小
3. 鲁棒性强



目录

- 1 马可波罗平台介绍
- 2 主题内容挖掘
- 3 优质内容抽取和生成
- 4 总结

1 结合用户反馈

通过用户反馈有效提高内容化的价值。

2 生成探索

图片自动增强/裁剪/
美化文案设计，文本生
成进一步探索。

3 内容聚合

内容多样性聚合，更加灵活。

Thanks For Watching