



# 第十一届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2020

## 架构革新 高效可控



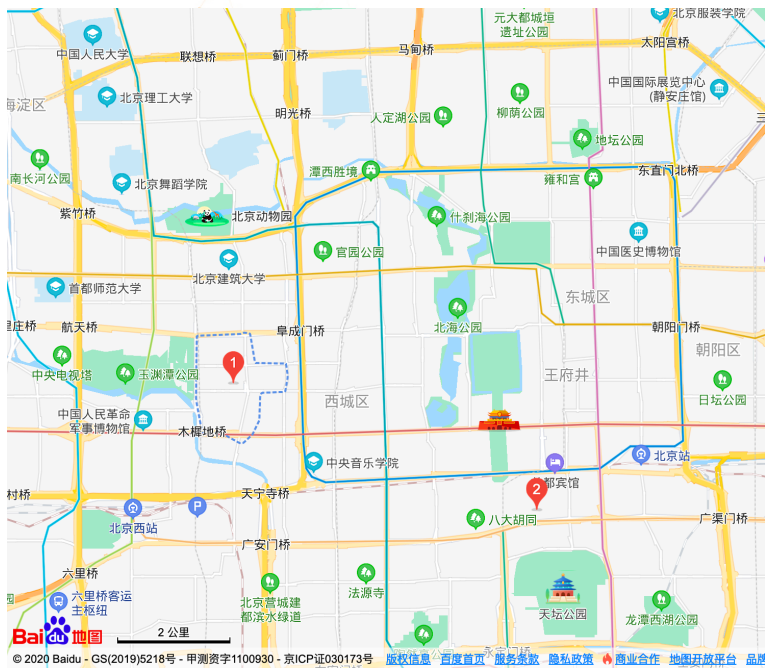
北京国际会议中心 | 2020/12/21-12/23

# 基于京东快递文本识别的自然语言优化之路

京东物流 王梓晨

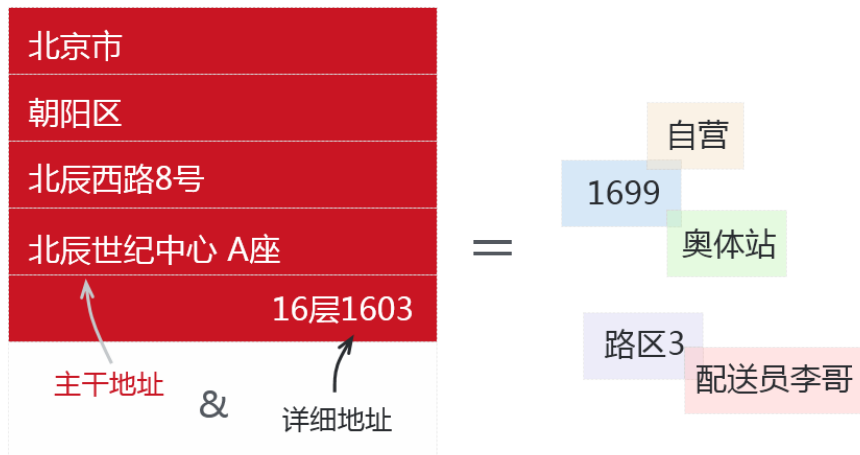


## 三里河的故事



## 地址有规范吗？

收件姓名是来自这个星球吗？



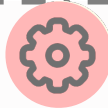
四级地址是什么？

北太平庄桥属于海淀还是朝阳？

## 业务意义



根据订单收货地址、  
商品属性、卖家类型  
等计算末端配送站信  
息



从ofc获得订单信  
息，下传订单信息  
到运单系统，控制  
下传速度

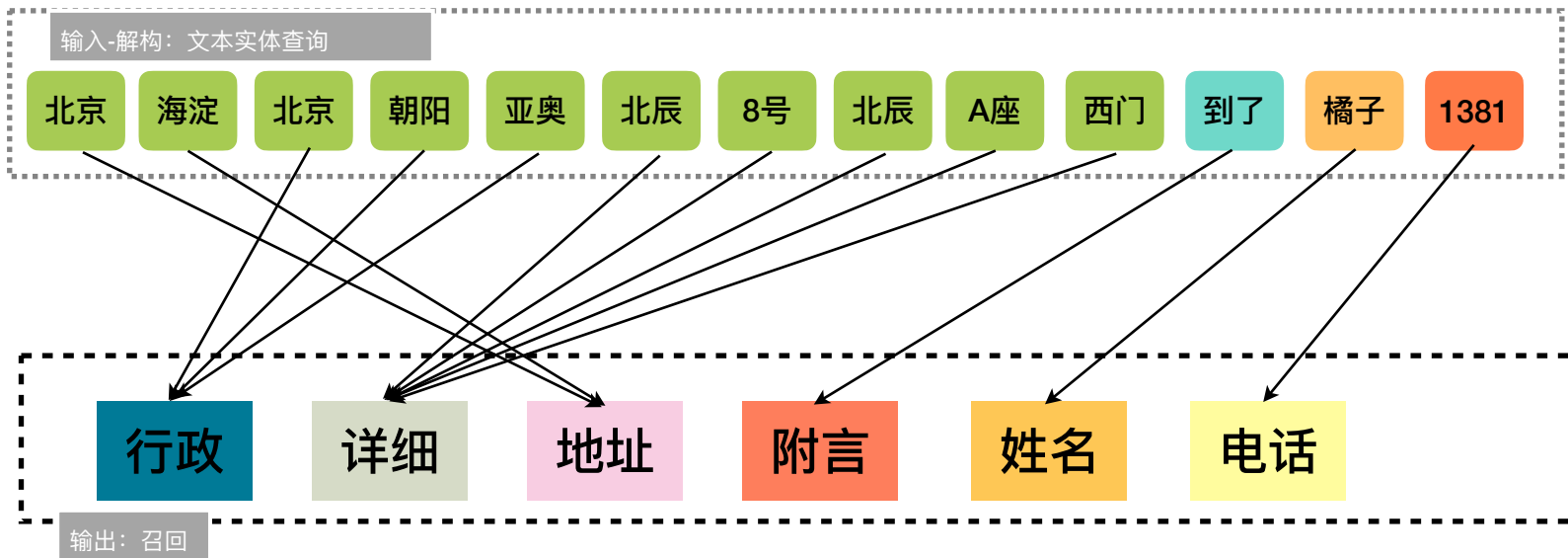


为仓库、分拣提供打  
印时的包裹标签信  
息：站点，路区，京  
鲜达等信息

### 分单的意义

1. 作为订单流程的一部分，下传订单信息到运单系统
2. 实现分拣中心快速作业
3. 指导站点提前安排配送资源
4. 让KA商家提前知道包裹配送的路由信息

## 目标：结构化



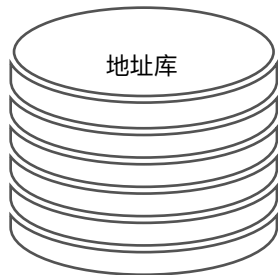
## 怎么做分词?

语料库  
地址库

到货请联系我前电话  
先电话联系  
送前电话联系以便留人签收电话联系

姓名库  
吴\*\*  
汪\*\*

电话库  
1392300\*\*\*\*  
1527015\*\*\*



北京大兴亦庄朝林广场A座  
小熊收18661861888



分词模型

label名称	词语
城市	北京
区县	大兴
poi	朝林大厦
人名	张三
电话	18661861888
附言	到了电话

## 地址分词

词典

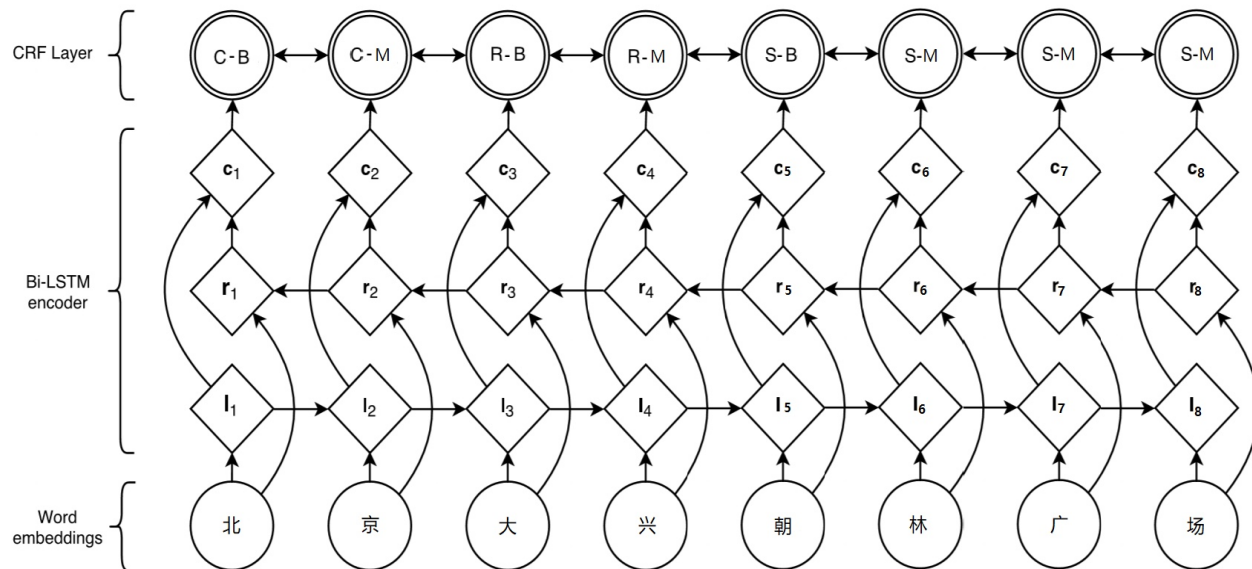


神经网络

嘉铭园-二区	c13332	下一级
武警一师家属楼	c13514	下一级
飘亮阳光广场	c13728	下一级
北京市朝阳区外国语小学二部	c13729	下一级
名人 大厦	c13731	下一级
慧忠北里 ( 扩建 )	c14070	下一级



## 文本解析网络结构



+词库 ... 遗留问题?

# 语义相似度

Word2Vec+TextCnn构建地址模型 “行政区划”

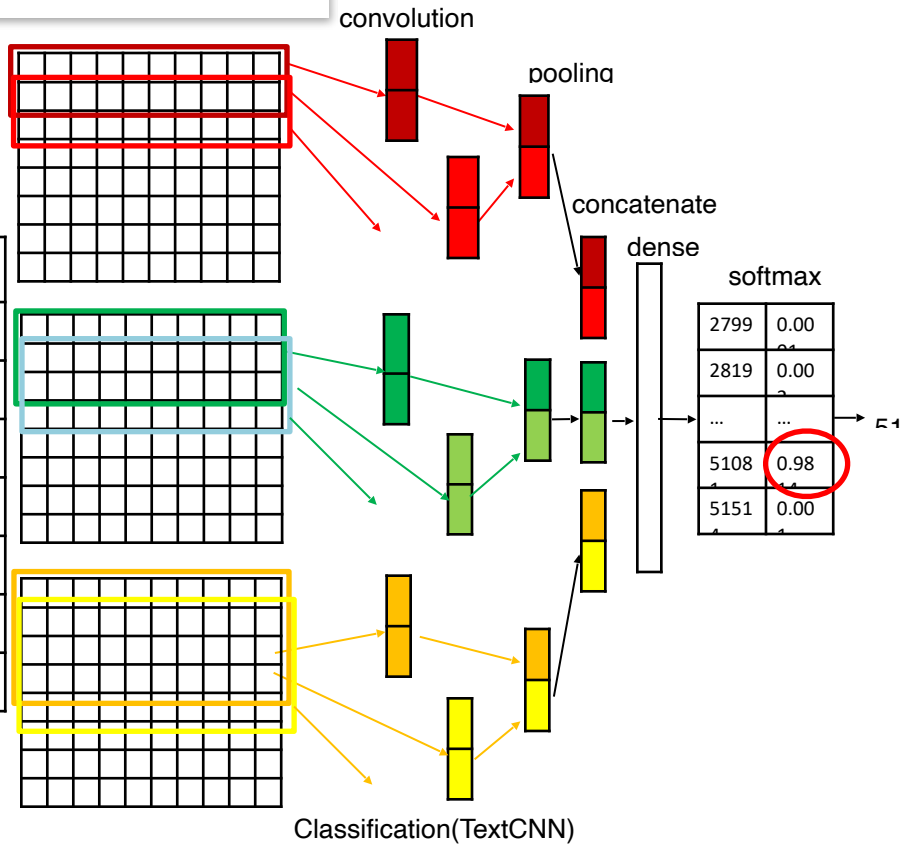
北京  
 大兴区  
 荣华  
 中路  
 19号  
 朝林  
 广场  
 A座

0.	0.	0.	0.2	...	0.	0.	0.	0.	0.
07	5	22	0		18	15	31	01	33
0.	-0.	-0.	0.5	...	0.	0.	-0.	-0.	0.
06	1	2	9		13	34	2	1	41
0.	0.	0.	0.0	...	-0.	0.	-0.	0.	0.
05	1	34	1		2	14	4	78	09
0.	0.	-0.	0.0	...	0.	0.	-0.	0.	0.
1	45	3	4		45	01	5	02	87
-0.	0.	0.	0.1	...	-0.	0.	0.	-0.	0.
8	02	23	7		6	08	19	2	12
0.	0.	0.	-0.	...	0.	-0.	-0.	0.	0.
07	1	37	1		02	3	6	37	22
0.	0.	0.	0.2	...	0.	-0.	0.	0.	0.
45	34	01	5		63	3	15	14	29
-0.	0.	0.	0.2	...	-0.	0.	0.	0.	0.
1	56	03	4		5	31	41	02	18

500dim

Address

Embedding(Word2Vec)



## 结构化基础数据建设架构

未登录词

地址挖掘

序列

语料

无监督



监督



其他?

地址结构化

地名消歧

自动补全

语义相似

词性分析

同X词

词性分析

地址库

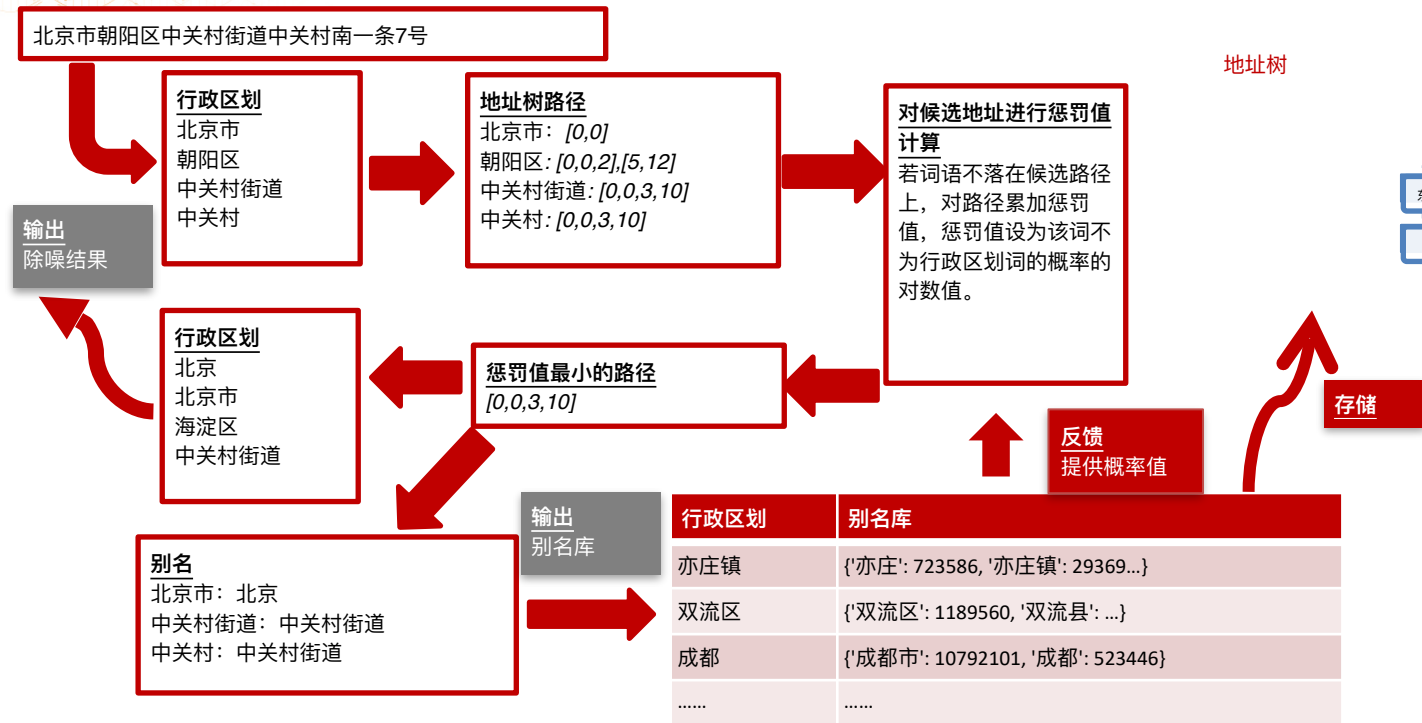
+

标注  
与检查

=

语料库

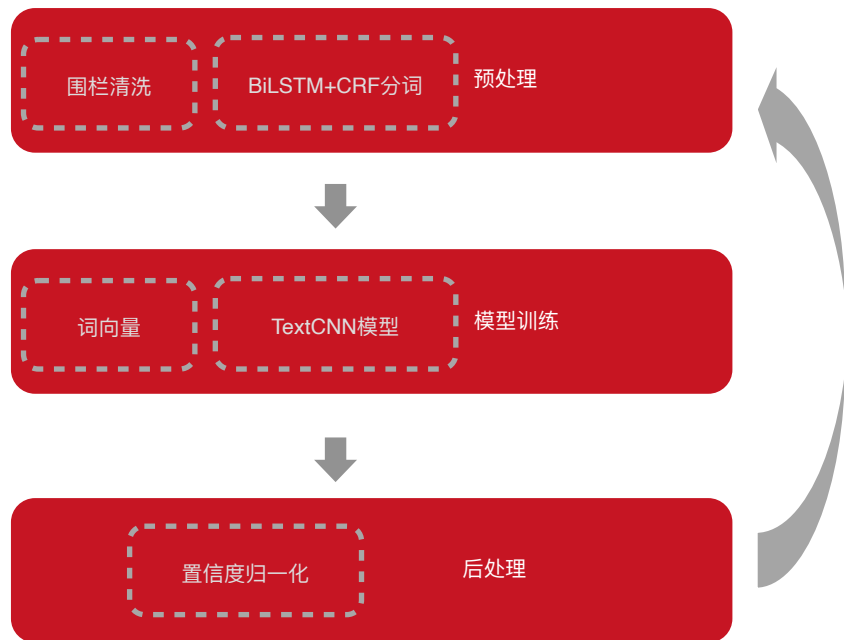
## 未登录词-编辑距离



## 相似度匹配

输出相似度匹配				
	建模	前处理	检索过程	后处理
	倒排表	重复率干扰	<div> <div>term匹配</div> <div> <div>优点</div> <div>缺点</div> </div> </div>	打分机制
	切词   不切词	召回限制	<div> <div>tf/idf</div> <div> <div>用户书写习惯不一致</div> <div>大词？小词？</div> <div>查找倒排索引表</div> </div> <div>topN</div> </div>	

## 相似度匹配



## 相似度匹配-打分

$$Score(Q, d) = \sum_i^N W_i * R(q_i, d)$$

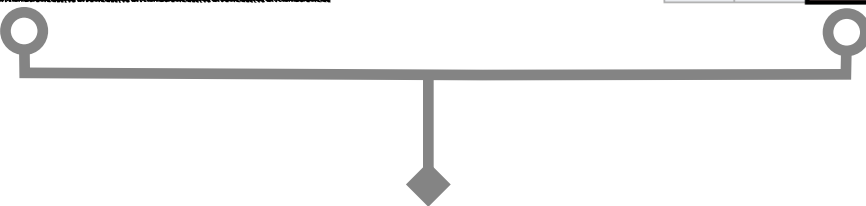
$W$  每个词在全文档的权重值

$R(q_i, d)$  每个查询词在当前文档中的得分

+

jaccard index( $A \cap B / A \cup B$ )

		$A$	
		0	1
$B$	0	$M_{00}$	$M_{10}$
	1	$M_{01}$	$M_{11}$



线性回归并以及 sigmoid 归一化来构建评估搜索返回结果置信度的模型

$$h_{\theta}(x) = \theta_0 + \theta_1 x_{tfidf\ score} + \theta_2 x_{site\ occupation} + \theta_3 x_{jaccard}$$



## 集腋成裘 · 聚沙成塔





# THANKS

