



# 医疗知识图谱的构建及应用

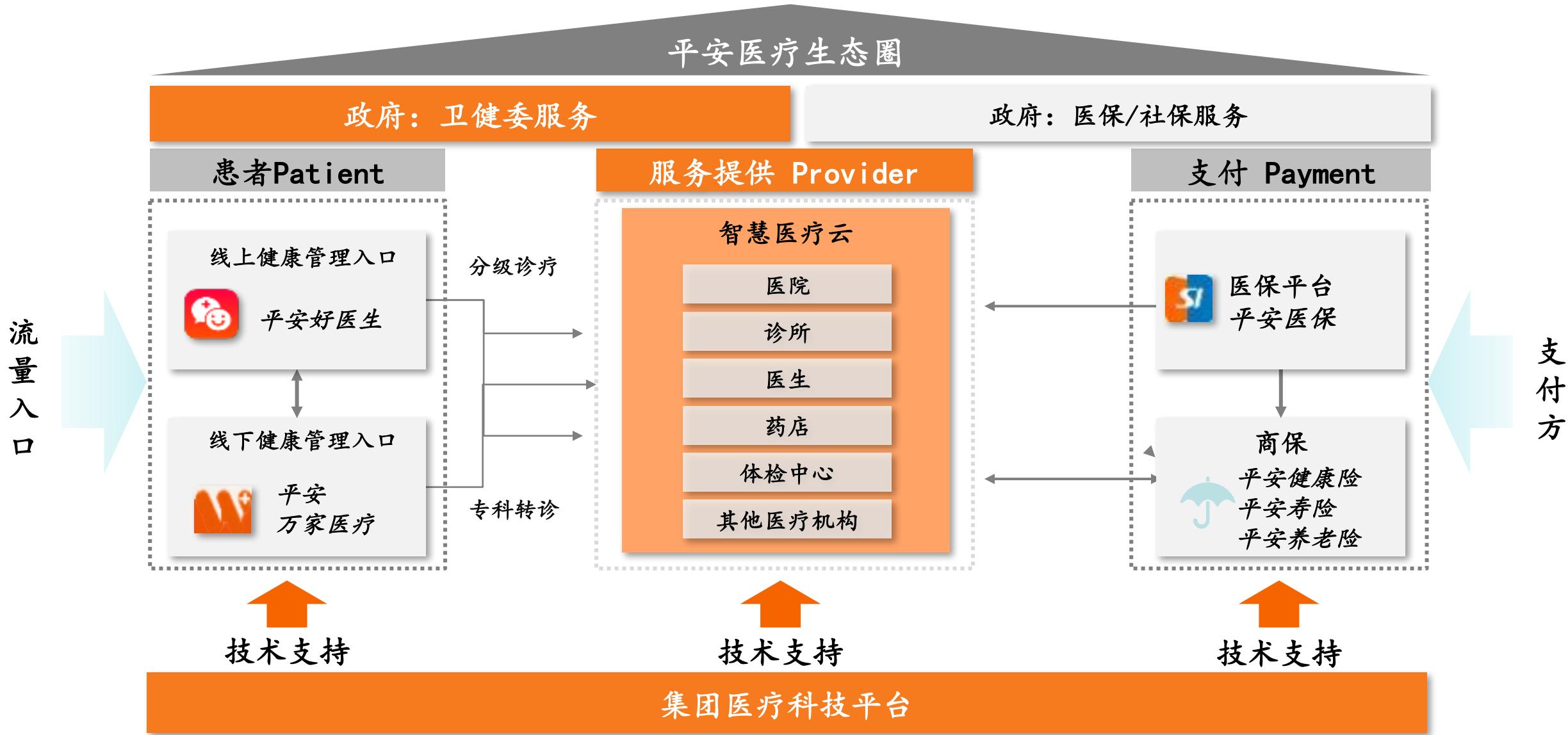


# 平安医疗科技

PING AN HEALTHCARE TECHNOLOGY

倪渊，平安医疗科技医疗文本处理负责人  
2018年12月8日

# 平安集团PPP模式医疗生态圈



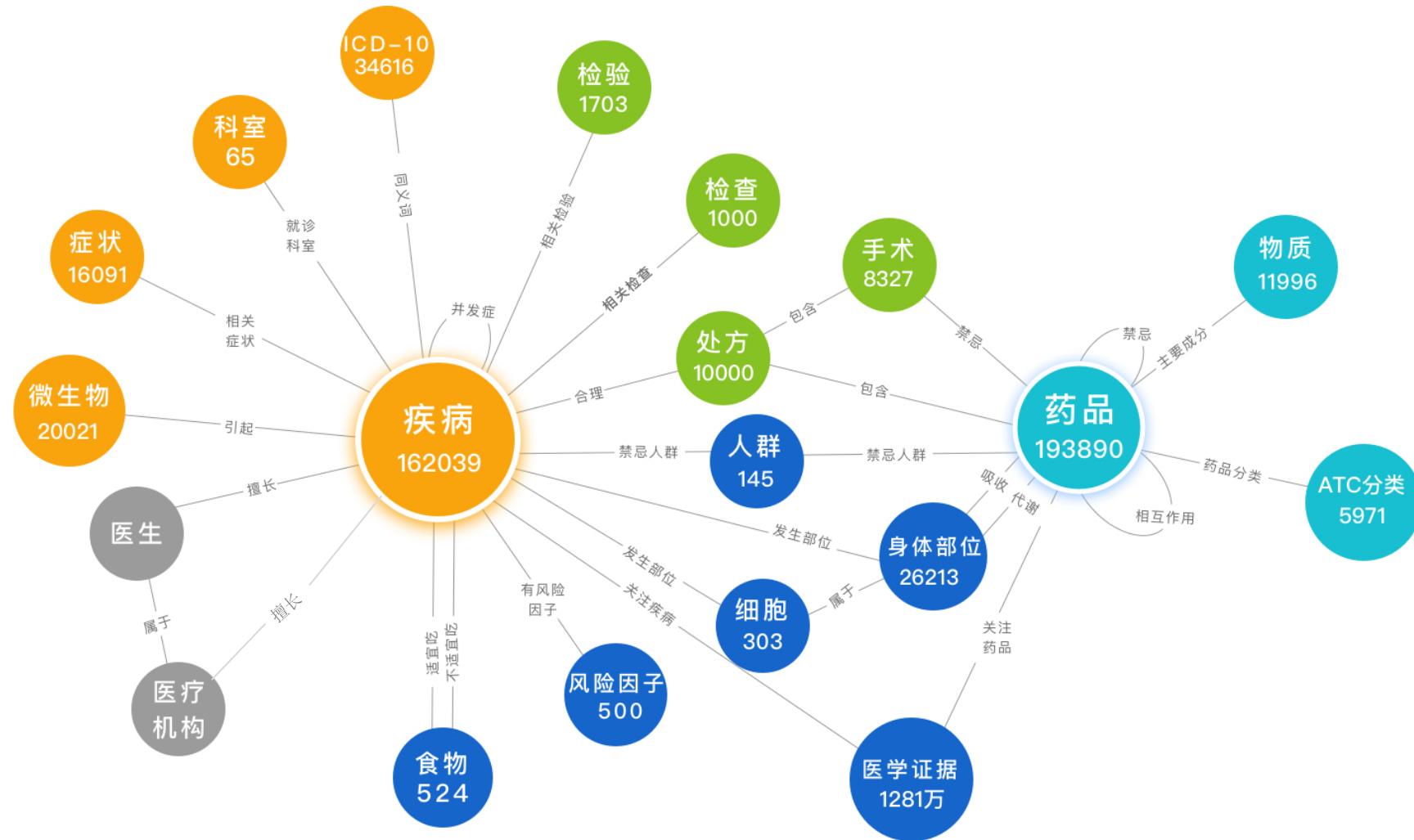
# 集团医疗科技平台



## A2 5大医疗信息库

<b>疾病库</b>	<b>治疗库</b>	<b>药品库</b>
<ul style="list-style-type: none"><li>• 疾病及疾病相关属性和关系</li><li>• 常见疾病知识库</li><li>• 多个城市地方病种库</li><li>• 不同版本ICD编码</li><li>• 疾病语义概念</li><li>• 疾病同义词库</li><li>• 症状同义词库</li></ul>	<ul style="list-style-type: none"><li>• 治疗模式</li><li>• 精准人群</li><li>• 检查检验信息</li><li>• 手术信息</li><li>• 临床路径</li><li>• 经典病例</li></ul>	<ul style="list-style-type: none"><li>• 药品信息</li><li>• 各省市药品目录</li><li>• 药品中标价格</li><li>• 药物不良反应</li><li>• 药物疗效</li><li>• 临床试验</li></ul>
<b>个人健康库</b>	<b>医疗机构&amp;医生库</b>	
<ul style="list-style-type: none"><li>• 疾病风险因子</li><li>• 患者常见问答</li><li>• 疾病知识问答库</li><li>• 科普文章</li><li>• 生活方式</li></ul>	<ul style="list-style-type: none"><li>• 医疗机构&amp;医生名称</li><li>• 医疗机构级别</li><li>• 主要科室</li><li>• 专家人数</li><li>• 住院床位</li><li>• 医生简介</li><li>• 医生资质</li><li>• 医生特长</li></ul>	

知识图谱：以结构化的形式描述客观世界中的概念及其关系



- 50类医学概念
  - 191种医学关系
  - 100种医学属性
  - 60万医学术语
  - 530万医学关系
  - 1000万医学证据

## A3 知识图谱：高血压知识图谱示例



# 平安医疗知识图谱-整体架构

融合医学核心概念以及医学临床证据

智能服务

医疗知识图谱查询  
导览及可视化

基于自然语言的  
知识图谱交互

基于知识图谱的  
医患教育

基于知识图谱的  
结构化和标准化

基于知识图谱的  
决策支持

数据平台

爬虫

MongoDB

图数据库

Cassandra

HBase

BerkeleyDB

ElasticSearch

自然语言处理平台

机器学习平台

图谱构建

Schema管理

数据结构化，槽填充

数据驱动

多图谱融合

图谱质量评价

图谱规则集

知识图谱 schema  
管理工具

实体识别

属性识别

模板  
定义  
和管  
理

基于数  
据挖掘  
的相  
关性

实体对齐

图谱冲突解决工具

规则校验工  
具

实体链接

关系抽取

对齐审核  
工具

图谱一致性校验

推理规则

数据层

医学核心概念

OMAHA  
开放医疗与健康联盟

中国医学科学院  
北京协和医学院  
医学信息研究所 图书馆  
Institute of Medical Information / Medical Library, CAMS & PUMC

疾病  
知识库

检查检验  
知识库

症状  
知识库

药品  
知识库

身体部位  
知识库

手术  
知识库

医学临床证据  
经典病例

RWE  
模型



随机对照  
的系统评价  
随机对照试验  
全或无病案研究  
队列研究的系统评价  
对列研究或较差随机对照研究  
结果研究；生态学研究  
病例对照研究的系统评价  
病例对照研究  
单个病例系列研究  
未经明确讨论或基于生理学，实验室研究或“第一  
原则”的专家意见

临床指南

医学  
科普



医学知识+证据图谱

临床指南

医学文献



检查

症状



疾病

药品



检验

膳食营养



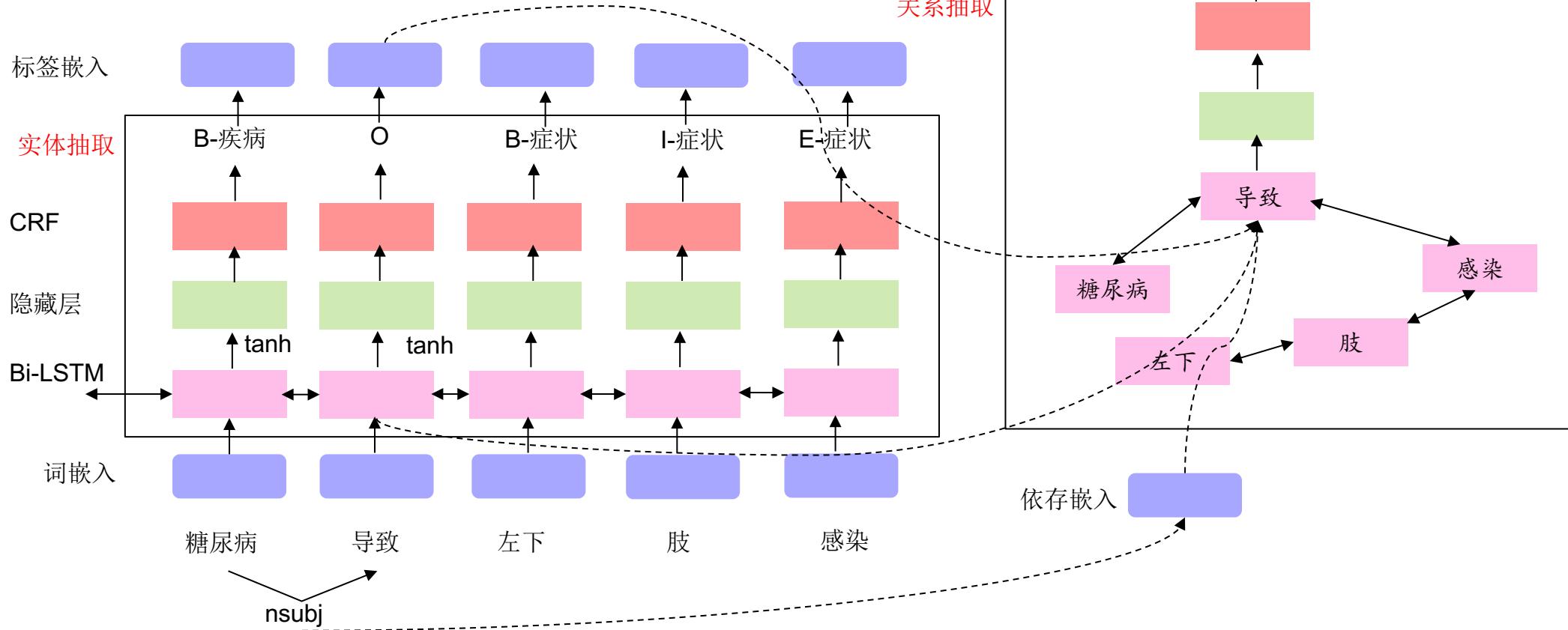
RWE模型

医学科普



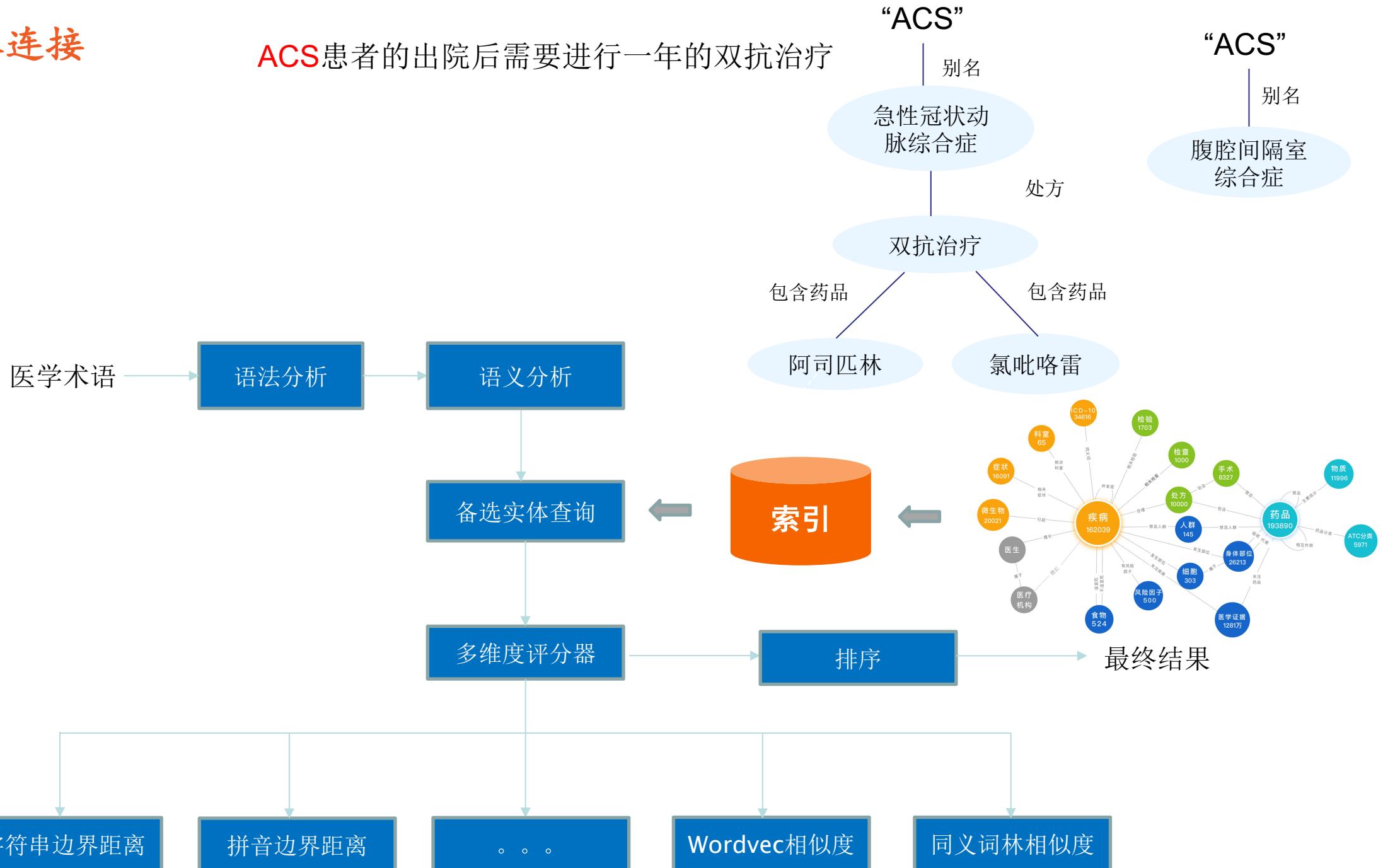
# 实体识别&关系抽取

- 基于深度学习的端到端的命名实体识别以及关系抽取，利用多任务方式同时提高命名实体识别以及关系抽取的精度
- 神经网络+知识图谱+人工知识
  - 神经网络: Bi-LSTM+CRF, Tree-LSTM
  - 知识图谱: 基于n-gram的匹配
  - 人工知识: 中文字形, 患者特征

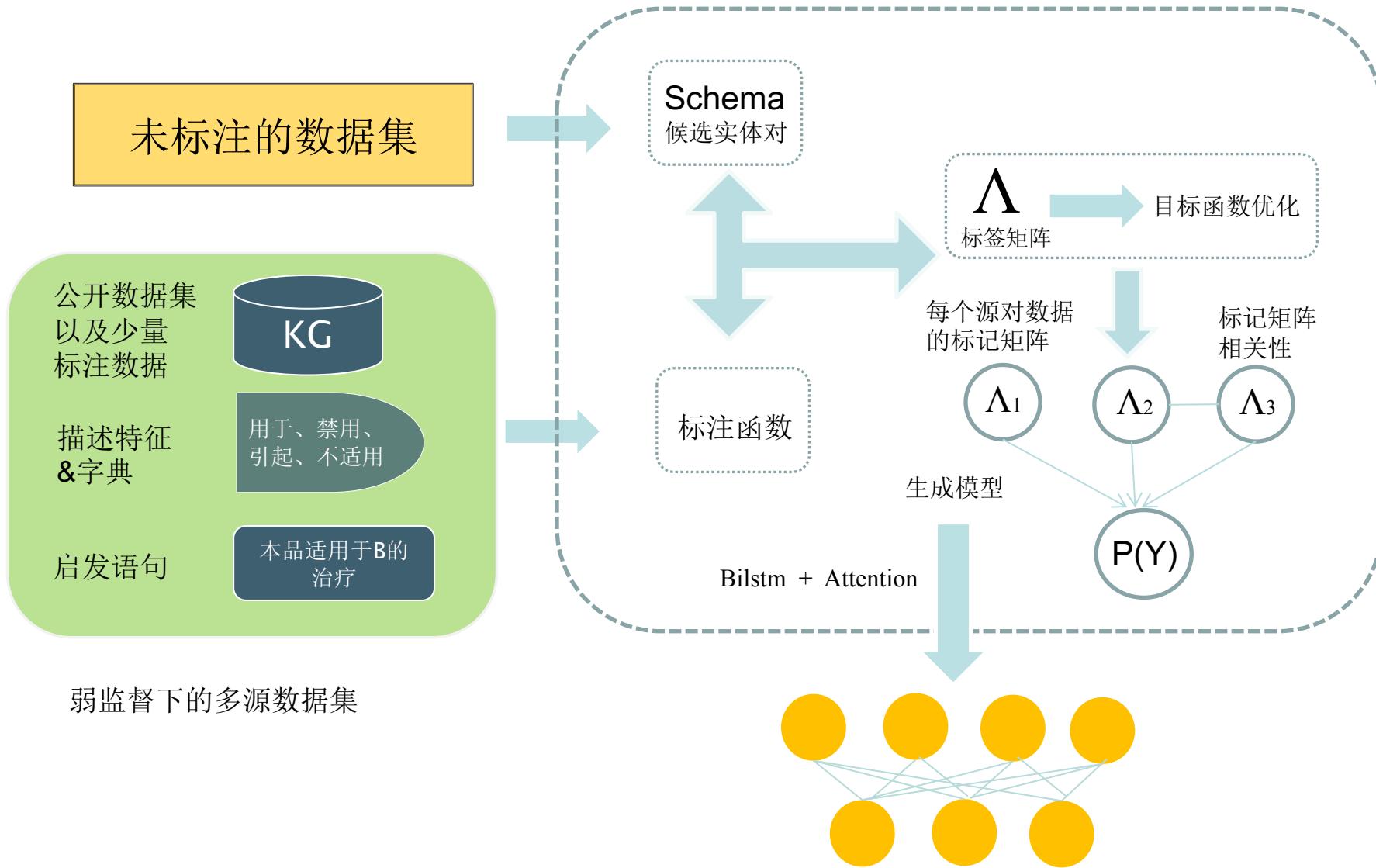


# 实体连接

ACS患者的出院后需要进行一年的双抗治疗



# 弱监督的信息抽取



# 图谱融合

## 医疗领域存在多种形式的同义词

### □ 不同翻译名称

- 阿司匹林
- 阿士匹灵
- 阿斯匹林

### □ 拼写错误

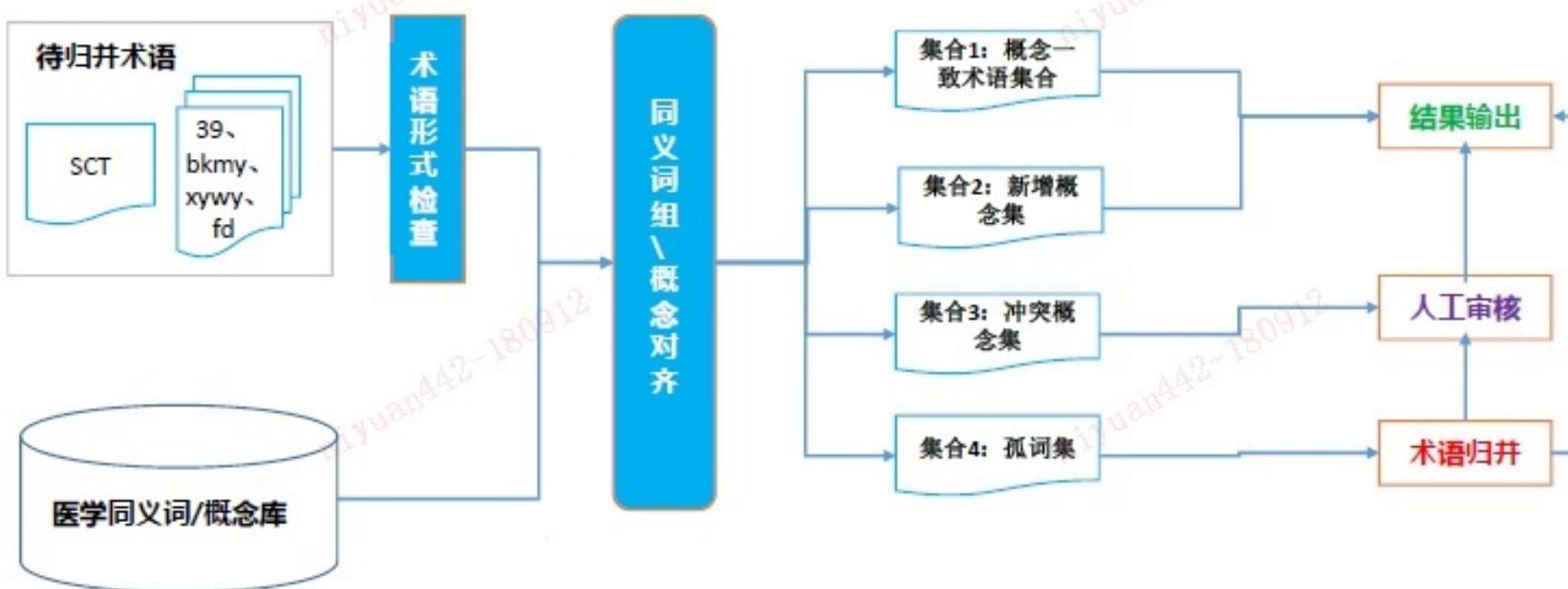
- 天竺癀
- 天竹黄

### □ 别名

- 小檗碱
- 黄连素片

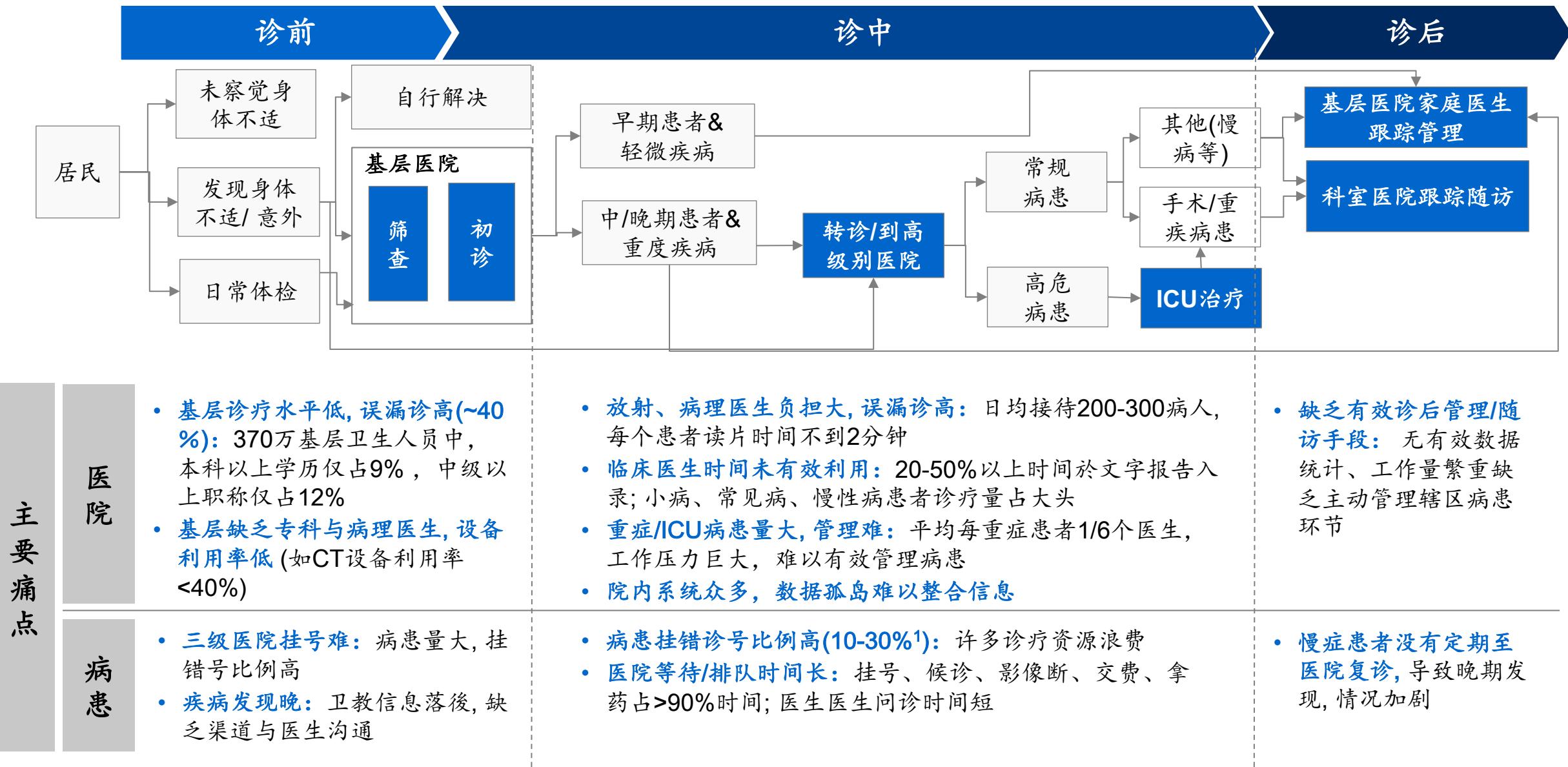
## 阿托伐他汀相关药物

阿托伐他汀分散片  
阿托伐他汀片  
山乐汀阿托伐他汀钙片(社区)  
阿托伐他汀片(山乐汀)  
氨氯地平/阿托伐他汀钙片  
氨氯地平阿托伐他汀片(多达一)  
阿托伐他汀片(立普妥)  
阿托伐他汀片(阿乐)  
阿托伐他汀钙分散片  
阿托伐他汀钙片  
阿托伐他汀钙片(医) 20mg/tab(\*\*)  
阿托伐他汀钙片(合资)  
阿托伐他汀钙片(山乐汀)  
阿托伐他汀钙片(山德士)  
阿托伐他汀钙片(立普妥)  
阿托伐他汀钙片(立普妥)H  
阿托伐他汀钙片(立普妥20mg)  
阿托伐他汀钙片(阿乐)  
阿托伐他汀钙片(限)  
阿托伐他汀钙片\_(立普妥)  
阿托伐他汀钙胶囊  
阿托伐他汀钙胶囊(国产)  
阿托伐他汀钙胶囊(限)



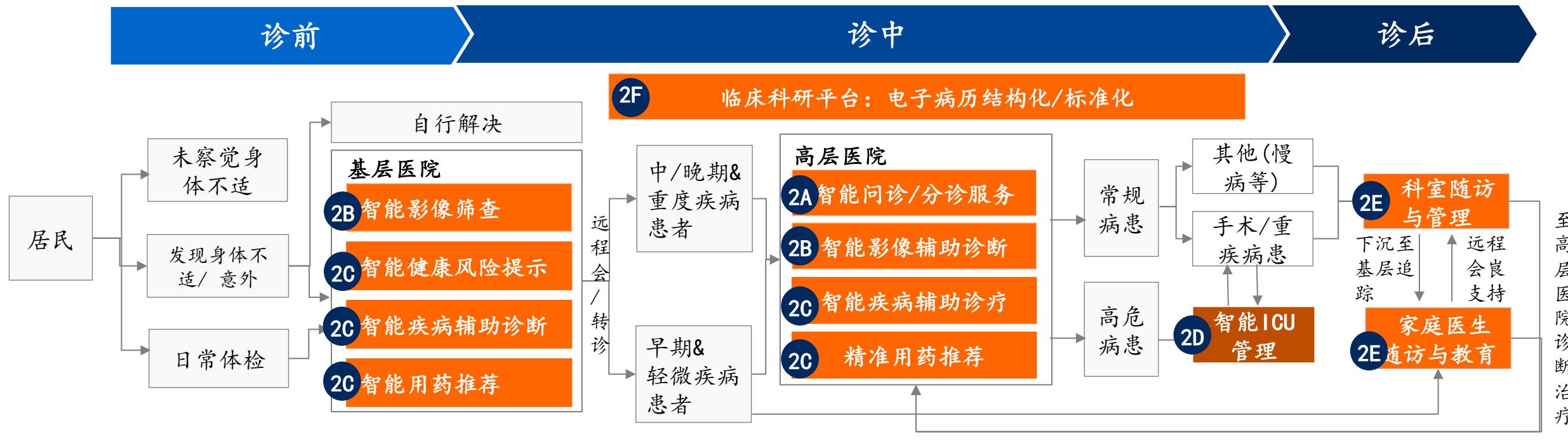
# 医联体/医院在诊前、诊中、诊后存在诸多痛点

关键环节



<sup>1</sup> 2015年首都医科大学附属北京朝阳医院急诊科主任郭树彬统计; 前几大误挂包括皮肤科误挂外科(~31%), 外科误挂内科(~30%), 口腔科误挂外科(~12%), 外科误挂妇产科(~11%), 中医误挂外科(~10%)

# 智慧医院CDSS解决方案全景图



## 2A 智能分诊 / 导诊

- 全病种(候诊大厅与药剂科)诊前数据收集&机器问答智能分诊
- 带来高效且精准的智能就医体验

## 2B 多模态智能影像辅助诊疗

- 提供~35种基于影像的辅助诊疗
- 包含智能诊报告&远程放射会诊等模块, 提升诊疗效率, 降低误诊率

## 2C 智能疾病辅助诊疗与治疗推荐

- ~30种常见疾病智风险提示、能辅助诊疗与用药建议, 规范临床诊疗和用药标准
- 精细管理医疗资源, 提升诊疗准确性

## 2D 智能ICU管理

- 死亡时间预测和诊疗资源预测、管理等智能模块
- 有效提高医生诊疗品质与效率

## 2E 智能患者教育/随访

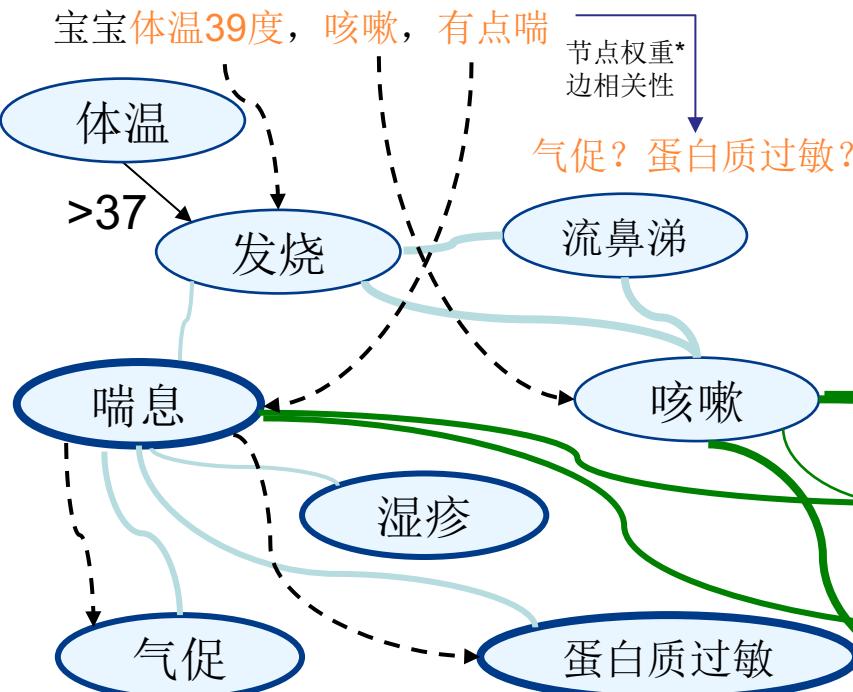
- 全病种基于自然语言理解的智能随访和诊疗
- 协助基层医院提升诊疗管理与随访效率

## 2F 智能临床科研平台

- 利用深度学习技术, 自动抽取电子病历中的关键信息
- 基于医疗知识图谱的医学概念标准化

## 诊前(分诊)

基于症状相关性以及重要度的推荐, 完善主诉

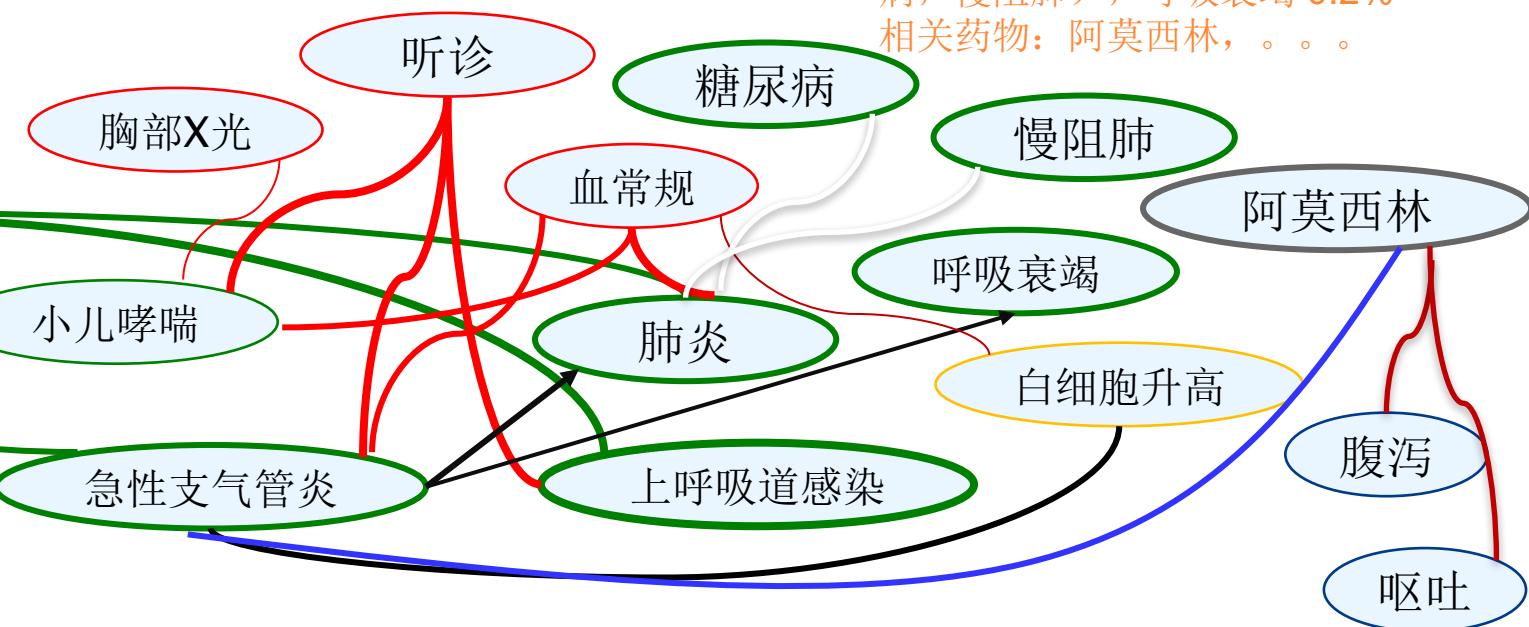


## 诊中(诊断辅助)

基于患者主诉, 选择相关疾病, 并按照节点权重及相关性进行排序。同时进一步给出排序的推荐检查

疾病: 急性支气管炎  
上呼吸道感染  
肺炎

检查: 听诊  
血常规



## 诊中(治疗辅助)

检查检验结果的解释。根据医生诊断, 提示可能得并发症及预后。提示可能得用药

急性支气管炎  
并发症: 肺炎 5% (风险因素: 糖尿病, 慢阻肺), 呼吸衰竭 0.2%  
相关药物: 阿莫西林, ...

## 数据准备



美国20年救护车使用  
调查公共数据集，包含210万次就诊记录

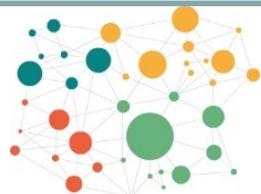


全球最大的医学  
文献数据库



斯坦福大学临床  
数据仓库

## 医学知识

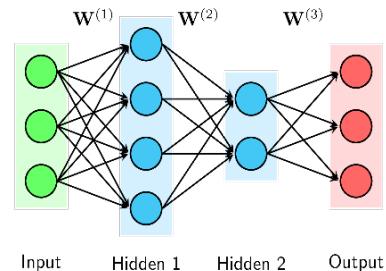


医学知识图谱和疾  
病的临床指南。包括1.2万个症状概念，  
14万个疾病概念，  
以及相互关系

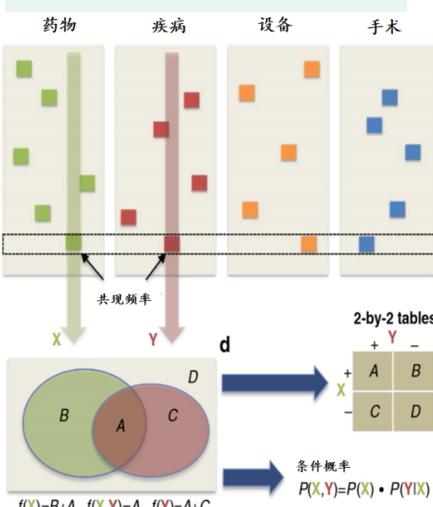
## 数据建模方法

## 构建深度学习模型

- 通过医学知识图谱及中的疾病和症状关系  
构建**基础网络**
- 利用公共数据进行增量学习，学习优化目标  
考虑先验知识 - 在损失函数中加入**知识  
对应的罚项**，提高模型精度。



## 构建贝叶斯网络



- 应用 PubMed 和 Stanford 数据构建**概率模型**，支持更多种类的疾病
- 通过**模型融合**技术生成最后的诊断列表

## 模型应用- 辅助诊断



多层肥胖

疑似诊断	
2型糖尿病	<a href="#">[详情]</a>
肥胖症	<a href="#">[详情]</a>
糖尿病视网膜病变	<a href="#">[详情]</a>
1型糖尿病	<a href="#">[详情]</a>
尿道炎	<a href="#">[详情]</a>
急性尿炎	<a href="#">[详情]</a>

主诉	多层肥胖
肥胖	多尿 多饮 多食 乏力 口渴 视物模糊 皮肤瘙痒 反复感染 尿频 尿急 尿痛
体征	身高 170 cm 体重 80 kg 收缩压 140 mmHg 舒张压 80 mmHg
糖化血红蛋白%	7.17 空腹血糖mmol/L 7.66
血压mmol/L	3.17 血清激素umol/L 120
高密度脂蛋白胆固醇mmol/L	1.17 低密度脂蛋白胆固醇mmol/L 4.17
筛查结果: 视网膜病变	中度 57.21% <a href="#">查看详情</a>
检查报告	

## 研究价值和进展

- 价值：利用患者症状、体征等信息为医生推荐疑似诊断，减少误诊、漏诊
- 进展：概率模型覆盖500+种疾病，全科常见30+种疾病的诊断模型的准确率为95%\*
- 合作：北京大学第一医院全科

## 数据准备

## 临床指南

- 收集整理了20+国家临床指南和专家共识
- 包括特定疾病的诊断、治疗和预防的方法和路径。

## 知识图谱

- 药品: 17万种
- 药品和疾病关系: 35万条
- 药物相互作用: 21万条



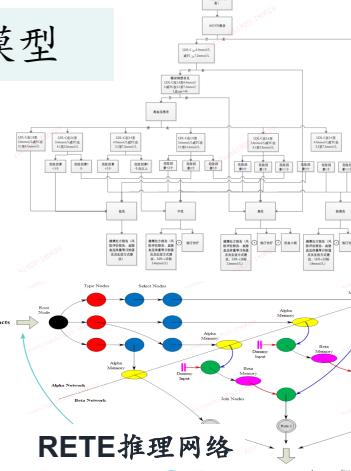
## 临床数据

- 重庆卫健委电子病历数据: 包括35种常见疾病的就诊115万次。
- 理赔数据: 包括35种常见病就诊1700万次。
- 每次就诊包括患者基本信息, 诊断, 检验检查信息, 疾病史, 用药史和处方信息。

## 建模方法

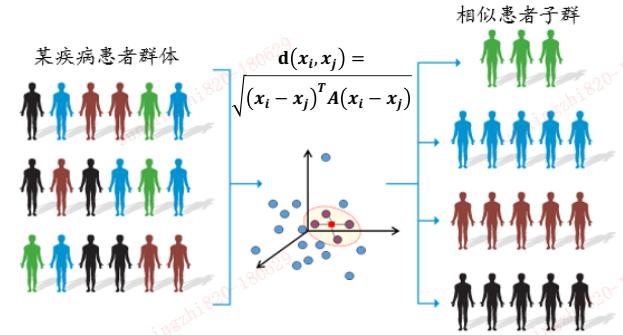
## 基于知识的治疗模型

- 知识表示: 将临床指南表示为决策树, 进而翻译成可执行的规则。
- 知识推理: 采用推理引擎将规则应用到患者数据, 生成治疗方案。



## 基于数据的治疗模型

- 通过关联规则分析发现常见治疗模式。
- 通过精准分群技术, 找出在临幊上相似患者的人群。
- 在相似人群中的进行治疗模式的个性化推荐。
- 采用模型融合技术把推荐结果和知识模型结果整合



## 模型应用- 药物推荐

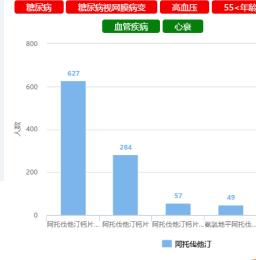
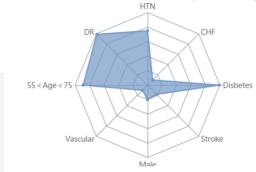
- 根据患者的信息给出处方推荐

用药推荐		药品通用名	药品商品名	规格	用法	剂量	频率	价格	备注
用药方案	血压管理	ACE/ARB	卡托普利	25mg	口服	2片	Q.d.	78元	
	血糖管理	二甲双胍	盐酸二甲双胍片	0.5g	口服	1片	Tid.	28元	
	血脂管理	阿托伐他汀	阿托伐他汀钙片(立普妥)	20mg	口服	2片	Q.d.	62元	
	血小板管理	阿司匹林	阿司匹林肠溶片	0.1g	口服	2片	Q.d.	16元	

- 提供推荐依据:

来自相似疾幊分析、临床指南、医学文献的证据		相关数据集中, 与当前患者相似的患者数	
根据《2016欧美高危人群糖尿病管理指南》和《2013年中国糖尿病指南》, 对此类病人需要使用高剂量他汀治疗。	提到他汀治疗的文献数: 2796篇	与当前患者相似的患者数: 26%	2329/8958人
相似患者中, 使用该处方的患者数: 48%	Pubmed Link	相关数据集中, 与当前患者相似的患者数: 26%	2329/8958人
4363654人		使用该处方后, 患者血脂达标率: 90%	114575756人
		使用该处方后, 患者依从性: 85%	15757655人

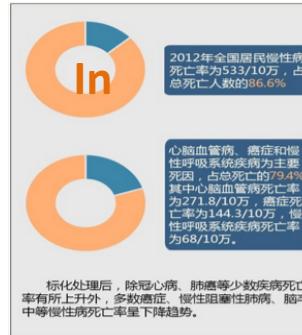
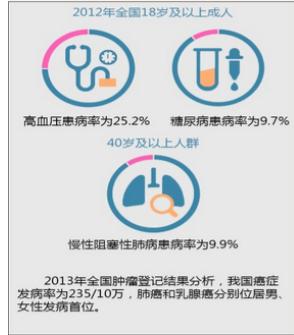
与当前患者相似的患者数: 26% (2329/8958)



## 研究价值和进展

- 价值: 规范医生治疗, 提升患者满意度
- 进展: 完成3种疾病(高血压, 糖尿病, 房颤)的依据知识的治疗模型, 正在开发基于数据的30+种疾病的治疗推荐模型
- 合作: 上海中山医院全科、赛诺菲

3亿高血压患者， 1亿糖尿病患者，  
5000万慢阻肺



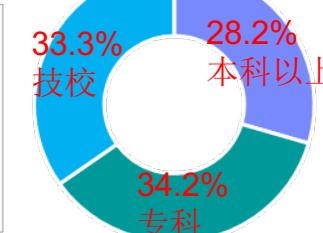
每1万人的社区医  
生数量



每1万人的药师数  
量



药师的教育水  
平



患者在慢病管理当中会有各种疑问



上亿注册用户

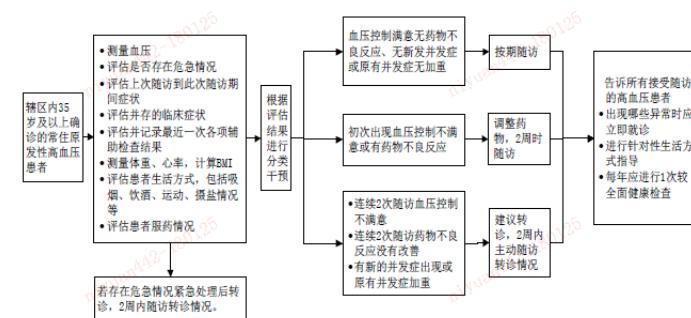


百万级日门诊量

大量健康教育相关问题，可以从医疗知识库当中找到答案

糖尿病患者应该如何饮食？  
立普妥和络活喜可以一起吃  
吗？  
空腹血糖8.6，有问题吗？

卫计委要求家庭医生要定期对慢性病随访



随访内容包括

症状：头痛，头晕，胸闷等

体征：血压，体重，BMI，心率等

生活方式：吸烟，饮酒，运动等

用药：种类，用法用量

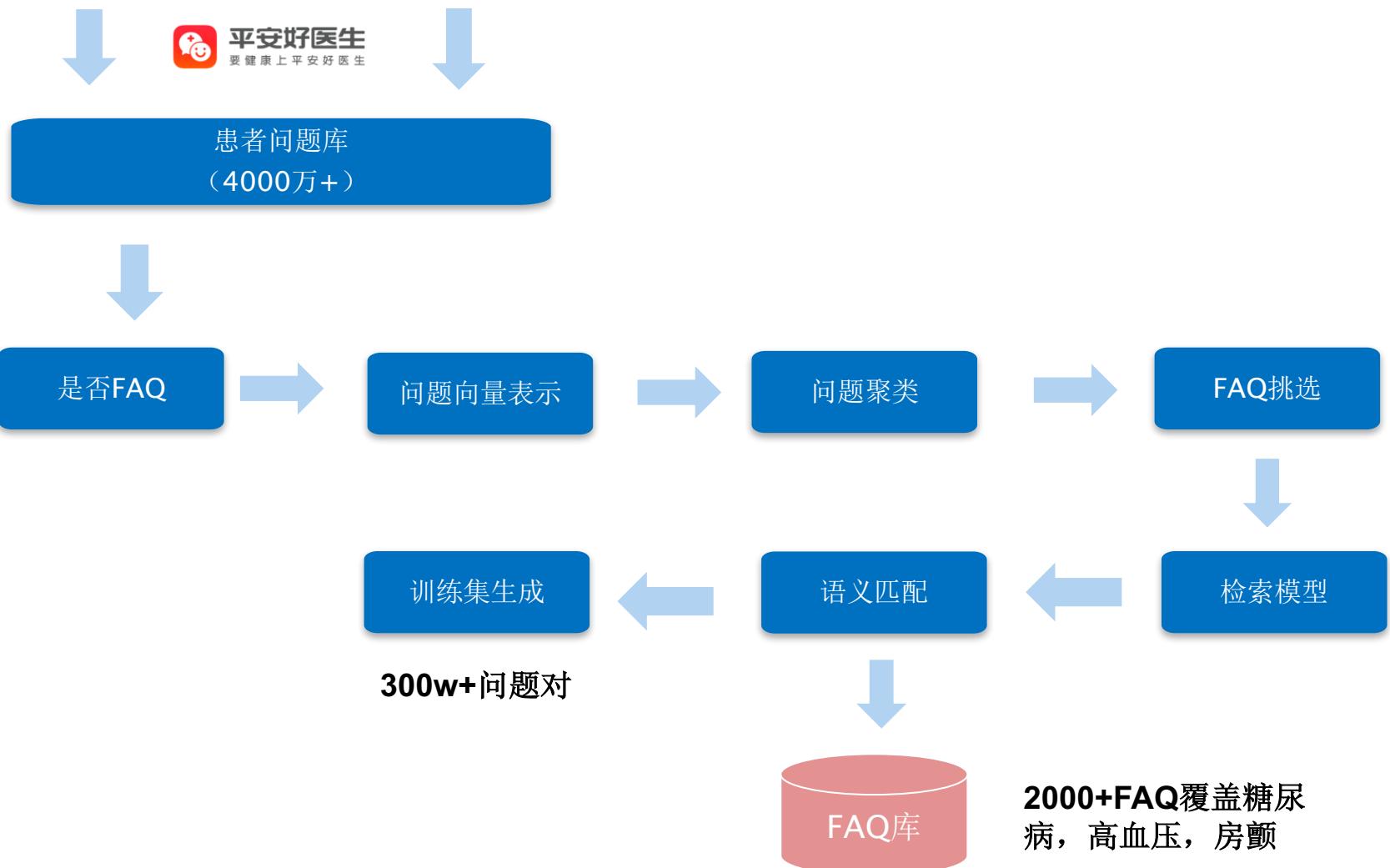
[“超过5亿人”的家庭医生去哪了](#) - 网易新闻

[数据丰满现实骨感,5亿人的家庭医生都在哪儿?-中青在线](#)

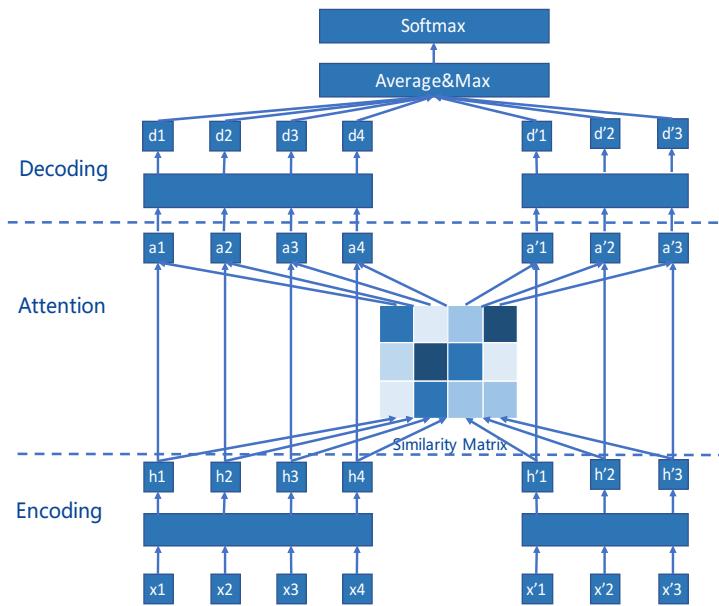
[5亿人有家庭医生？医生:连病人家门在哪都不知道!家庭医... 新浪新闻](#)

# 疾病健康问答-问题相似度匹配

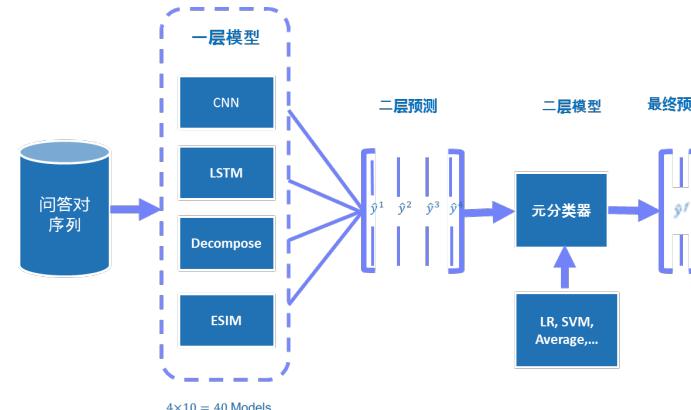
## 数据驱动的FAQ问题集准备



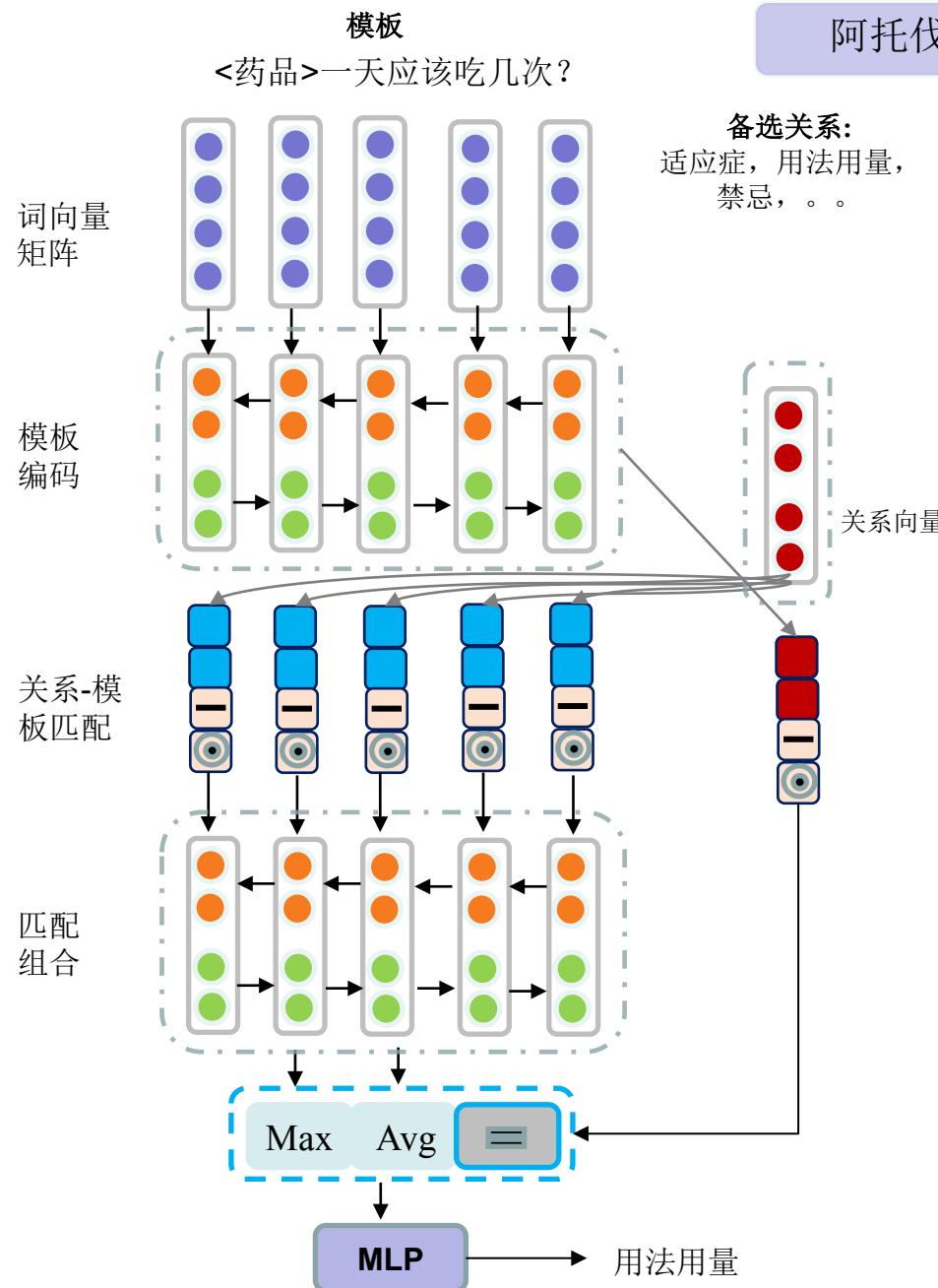
## 基于深度注意力网络的文本匹配



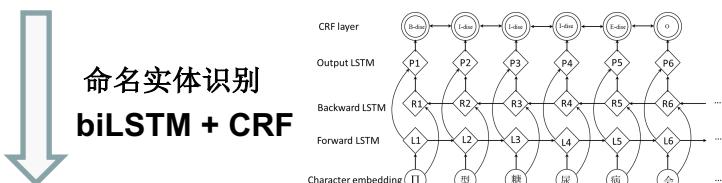
## 集成学习技术



病种	问题库	FAQ	问题组	训练集	精度
糖尿病	170,000	90,000	700	340万对	93%
高血压	110,000	80,000	707	110万对	95%
房颤	10,000	5,000	200	10万对	92%
卒中	60,000	20,000	350	35万对	92%
慢阻肺	20,000	10,000	437	14万对	91%



阿托伐他汀一天应该吃几次?



阿托伐他汀一天应该吃几次?

↓ 实体链接

阿托伐他汀钙片 阿托伐他汀钙分散片 阿托伐他汀钙胶囊 阿托伐他汀胶囊

↓ 多轮对话的实体消岐

片剂 胶囊

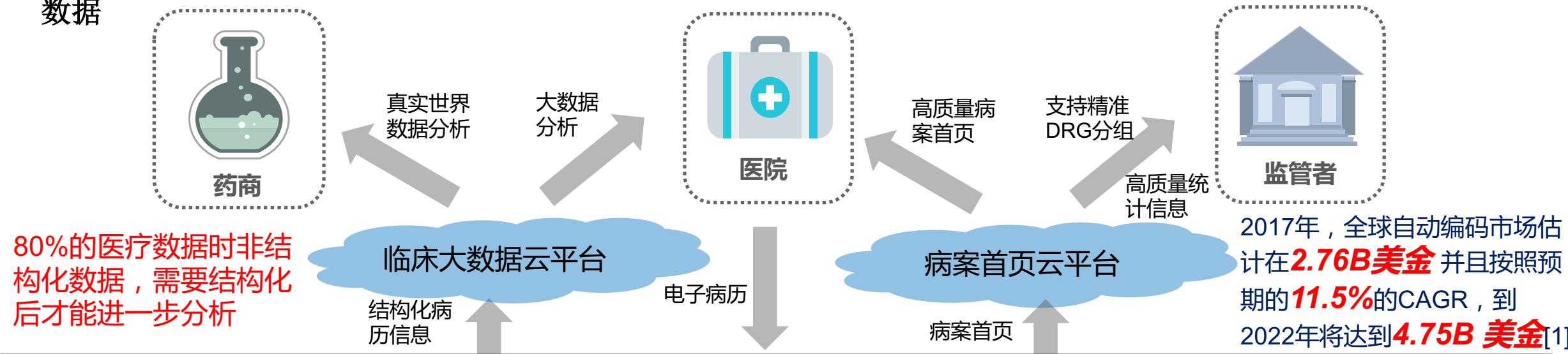
↓ 立普妥 阿乐

↓ 立普妥(10mg) 立普妥(20mg) 立普妥(40mg)

↓ 立普妥 (10mg) 常用的起始剂量为 10 mg 每日一次...

# 基于知识图谱的电子病历理解和智能编码

每年**70亿人次**就诊患者，产生**100亿+**电子病历数据；通过电子病历结构化以及自动编码服务，整合电子病历数据



患者诉**发冷**，出现**寒战**，**头昏**，**恶心**未吐。血常规：白细胞数 $17.83*10^9/L$ ，。。。查体：左肾区叩击痛阳性，未触及肿块。。。考虑患者**泌尿系感染**合并**SIRS**。制定治疗方案：1停止静点头孢替唑，给予**头孢哌酮舒巴坦钠** $2.03/日$ 静点

自然语言理解理解

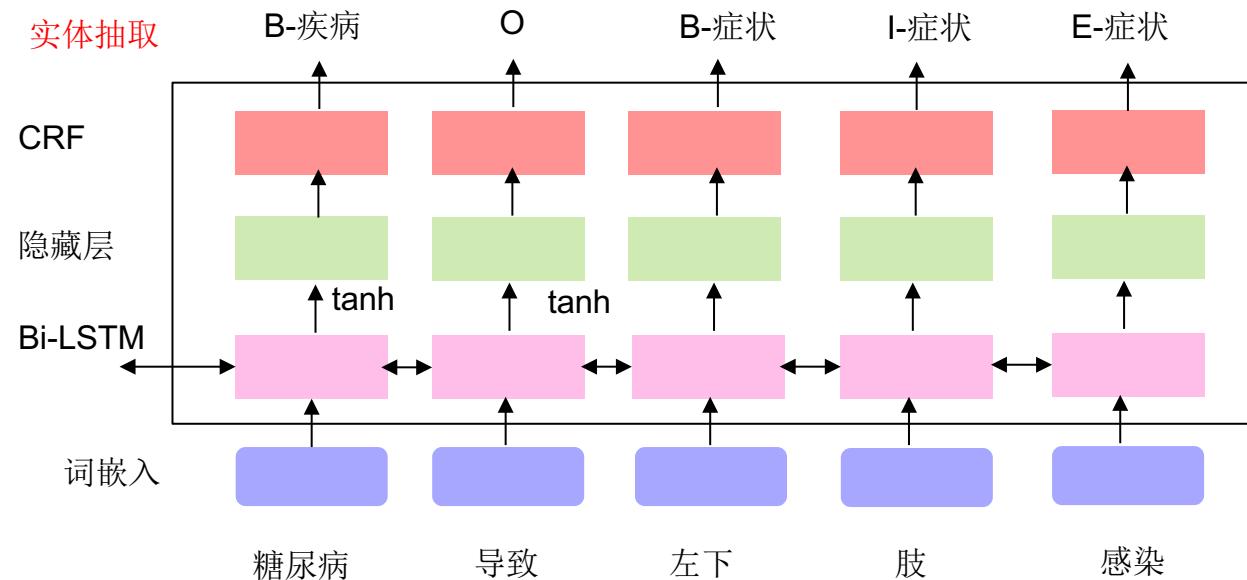
症状	有	寒战，头昏，恶心
	无	呕吐
检验	白细胞	$17.83*10^9/L$
体征	有	左肾区叩击痛阳性
	无	肿块
诊断		泌尿系感染合并SIRS
治疗	无	头孢替唑
	有	头孢哌酮舒巴坦钠

自动编码

ICD-10疾病编码

ICD-9-CM手术编码

## □ LSTM+CRF 模型



## □ 模型只能解决部分问题

- 无恶心，呕吐，无发烧，发冷
- 尿尿尿不出尿来3天
- 全身浅表淋巴结未触及肿大，头颅无，五官端正
- 今查房，症状同前，体查同前，治疗同前
- 巴林斯基征，巴氏征，巴宾斯基征，巴宾斯基，babinski征，巴彬斯基征，babinskisign。。

# 电子病历结构化和标准化产品

数据看板，快速浏览结构化电子病历的统计信息

电子病历结构化标准化



支持对于结构化电子病历的语义查询

电子病历结构化标准化

主要诊断名称: 面部肿物

性别: 男

添加筛选条件

确定

自动抽取电子病历文本中的信息（例如疾病，症状，症状属性等），并利用知识图谱将信息标准化。生成标准化的结构化数据



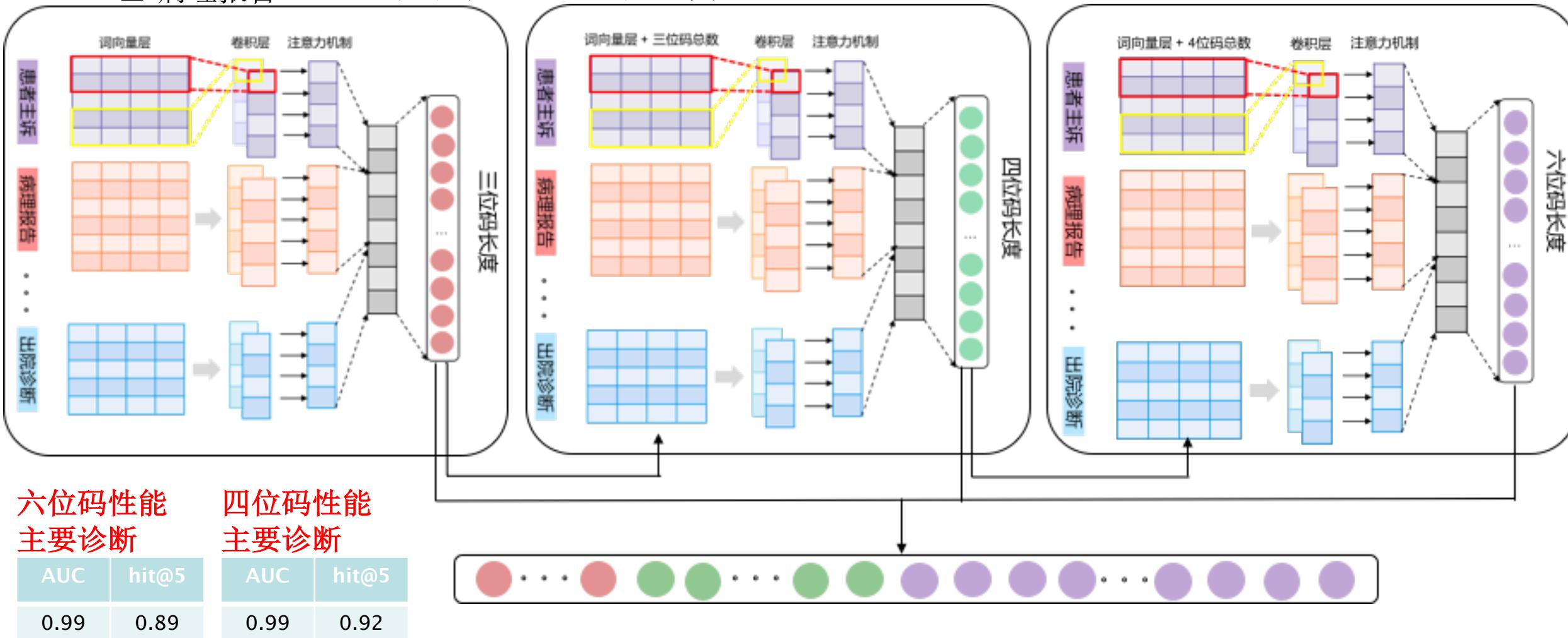
文件名	病历描述	上传时间
EMR-6.txt	中年男性, 49岁, 主因: 腹痛、腹泻2个月...	2018-8-23
EMR-71.txt	半年, 加重1周于2017年2月15日入院...	2018-8-23

## 基于机器学习的自动编码

- 考虑整体的电子病历来提取特征
  - 患者主诉
  - 病理报告
  - 就诊过程
  - 出院诊断

- 出院医嘱
- 患者基本信息

- 通过机器学习来自动的推荐编码，包括主要诊断，次要诊断以及手术编码
- 通过医学知识图谱来进一步提高性能



## 模型比较

训练数据: 22万+电子病历, 包含临床诊断描述, 科室, 患者信息等和icd编码

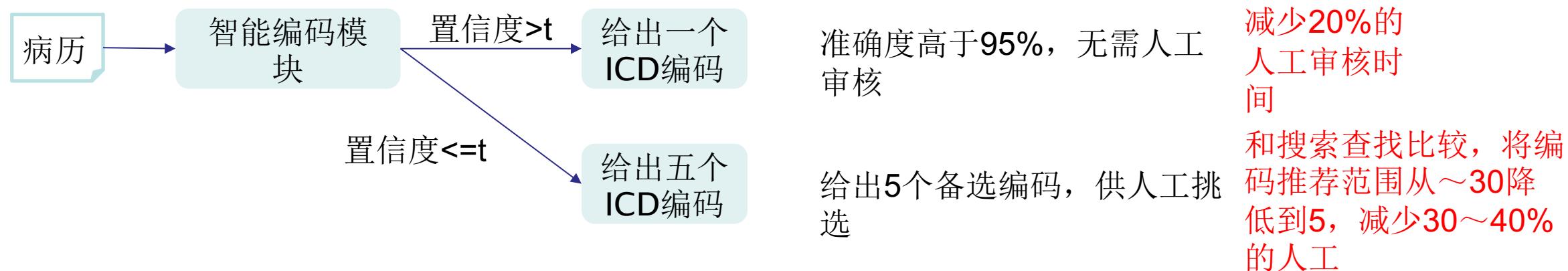
模型 (Label)	Hit@1	Hit@5	coverage hit@1=0.95	coverage hit@1=0.90
模型1(一层六位码模型) (六位码)	0.62	0.85	2%	5%
模型2(层级模型) (六位码)	0.66	0.85	17%	28%
模型3(层级多任务模型) (六位码)	0.73	0.89	20%	32%

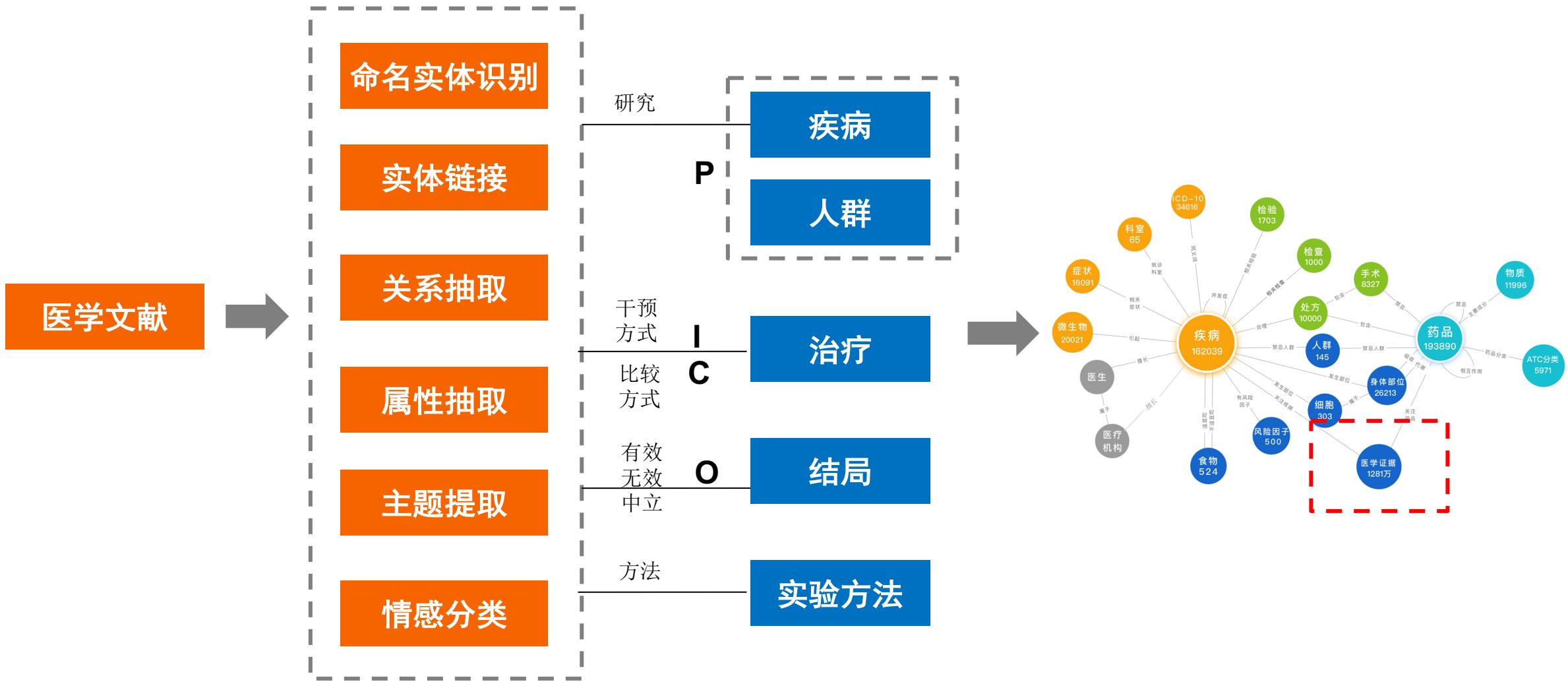
hit@1 = 模型只给出一个备选编码的时候就是正确编码的百分比

hit@5 = 模型给出5个备选编码的时候包含正确编码的百分比

coverage @ (hit@1=0.95) 在hit@1可以达到0.95的时候, 模型可以覆盖的病历百分比

coverage @ (hit@1=0.90) 在hit@1可以达到0.90的时候, 模型可以覆盖的病历百分比

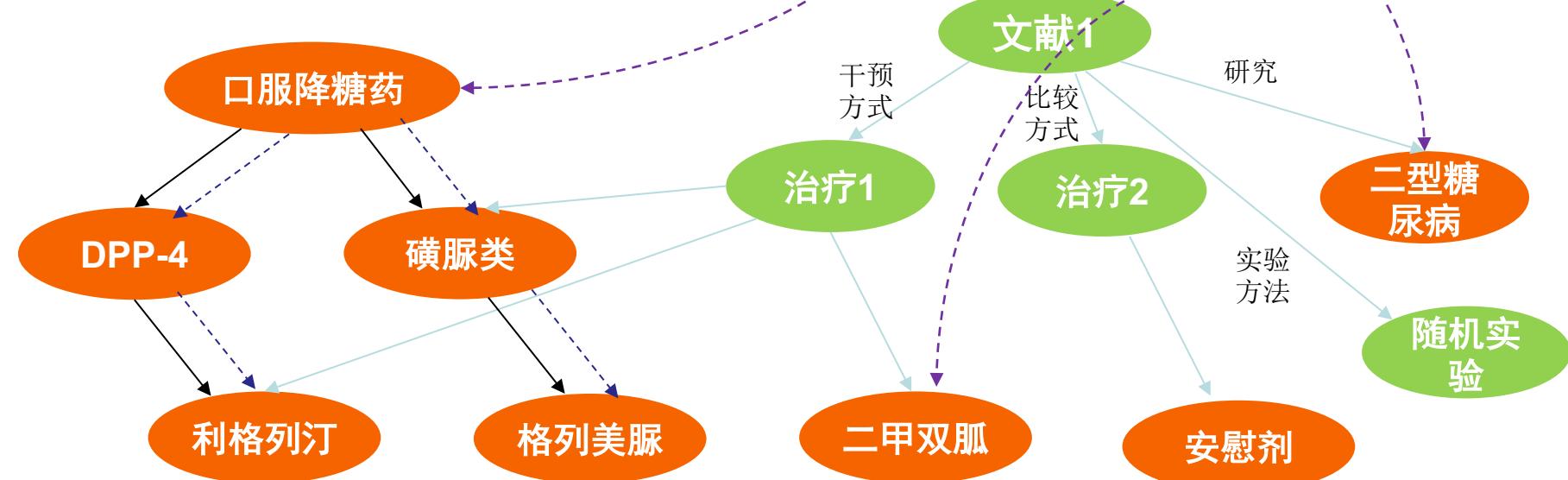




## 利格列汀联合二甲双胍及磺脲类药物在2型糖尿病患者治疗中有效性、安全性分析

**摘要：**目的：研究利格列汀联合二甲双胍及磺脲类降糖药物治疗（二甲双胍及磺脲类药物血糖控制不佳的）的2型糖尿病患者有效性及安全性的相关数据。方法：48例2型糖尿病患者按2：1比例随机分为利格列汀组与安慰剂组，两组均用药24周。结果：利格列汀与二甲双胍及磺脲类降糖药物联合应用24周可以显著降低糖化血红蛋白（HbA1c）的水平（与安慰剂组比较，HbA1c较基线时下降0.62%，P<0.0001）。结论：利格列汀与二甲双胍及磺脲类降糖药物联合应用安全有效。

2型糖尿病患者，已经在服用二甲双胍，血糖控制不满意，加入哪种口服降糖药合适？





Thanks