

医疗知识图谱的敏捷构建和实践

刘升平 博士
资深技术专家，云知声 AI Labs

云知声的业务



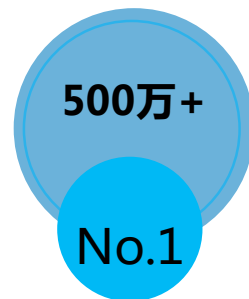
Athena 智慧大脑



产品技术架构

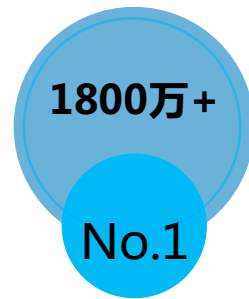


两大商业实践领域



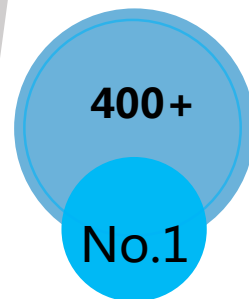
家居业务

产品出货量



车载业务

产品出货量



医疗业务

合作医院数



教育业务

口语评测日调用量

知性会话：基于语境知识图谱的人机对话系统



物理语境

时间/地点/场所

天气

情绪和
情感

设备显示

设备感知

言语语境

上下文

主题及
焦点

设备反
馈

知识语境

人类常
识

领域知
识

用户画
像

Agent画
像

设备信
息库

语境的生命周期：

请求级

会话级

长期

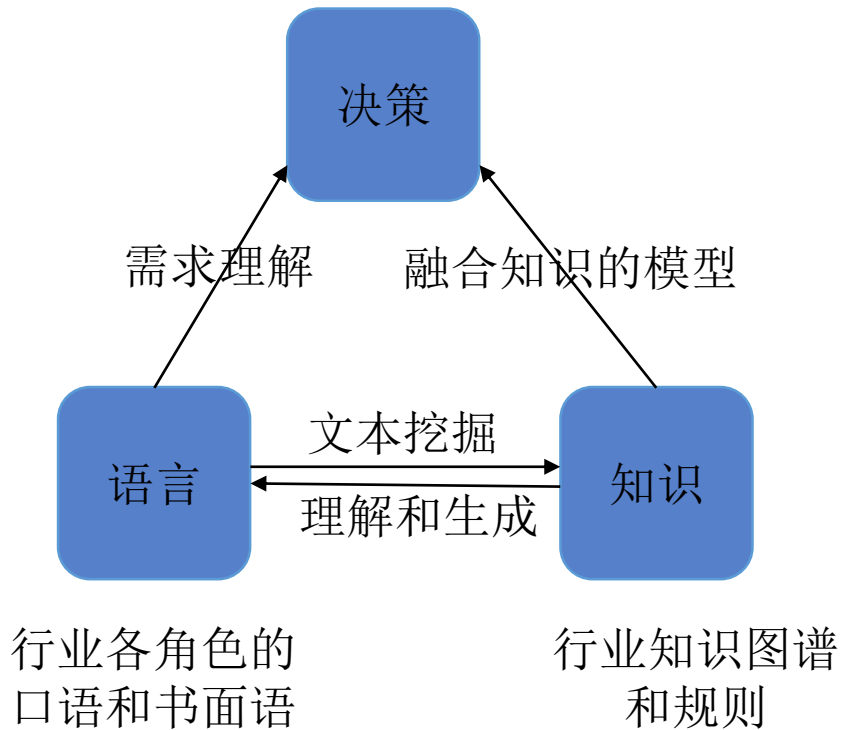
基于语境知识图谱的人机对话示例



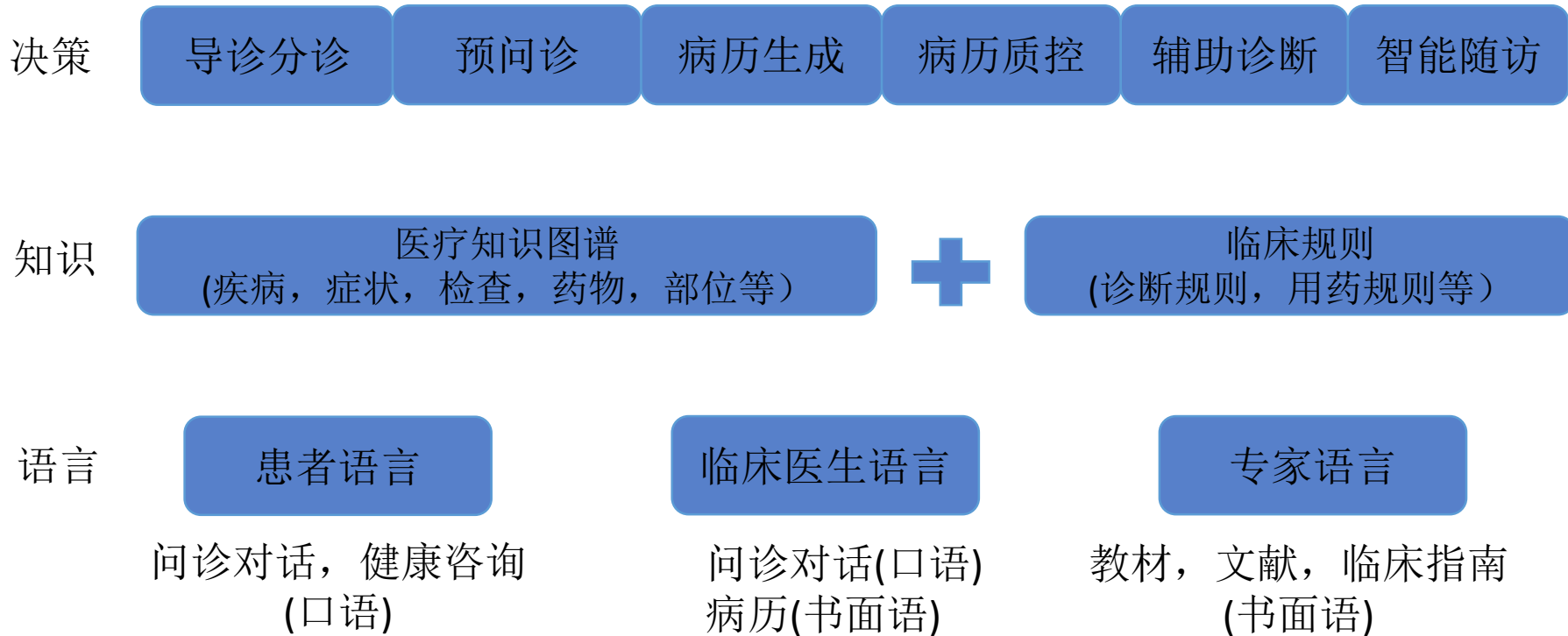
行业认知智能

行业认知智能三要素:

- 语言
- 知识
- 决策
- 知识引导自然语言理解和生成
- 知识辅助决策

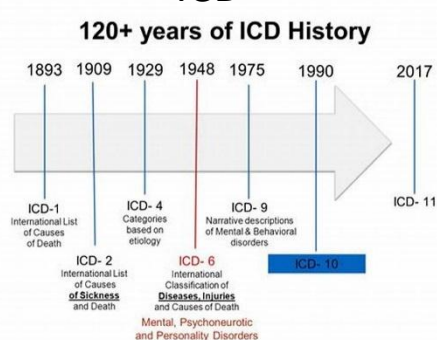


认知医疗-以医疗知识图谱为核心



医疗知识图谱的演进

ICD



MeSH

LOINC

SNOMED CT

UML-S

Gene Ontology

SNOMED CT OWL Ed.

GALEN

FMA

Google Medical KG

DrugBank

Bio2RDF

LinkedLifeData

医学术语系统

医学本体

医学知识图谱

1900

2000

2012

术语系统 (Terminology System)

Organization of Concepts

Concept
C0001621

Term
L0001621

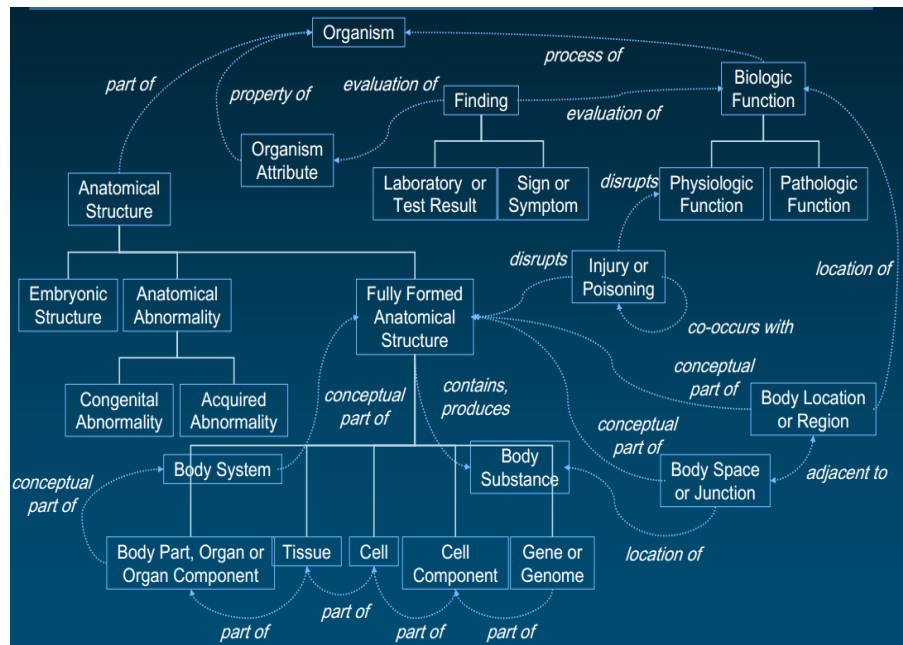
S0011231 Adrenal Gland Disease
A0020266 MeSH
A7568579 NCI Thesaurus
S0000441 Disease of adrenal gland
A0001264 SNOMED 1982
A6917004 SNOMED Clinical Terms
S0481705 Diseases of Adrenal Gland
A0014499 SNOMED 1982
S0220090 Diseases, adrenal gland
A0049924 MeSH

Term
L0181041

S0632950 Disorder of adrenal gland
A0688820 Read Codes
A4778687 SNOMED Clinical Terms
S0354509 Adrenal Gland Disorders
A6996540 MedlinePlus
A7576253 NCI Thesaurus
A7561794 Psychological Index Terms

Term
L1279026

S1520972 Nebennierenkrankheiten
A7500884



医学术语系统本体化：以SNOMED CT为例

Table 1 – Characteristics of description logic EL^{++}

Existential quantification	$\exists \text{FindingSite.AppendixStructure}$
Conjunction	$\exists \text{AssociatedMorphology.Inflammation} \sqcap \exists \text{FindingSite.AppendixStructure}$
Necessary condition	$\text{Acute Appendicitis} \sqsubseteq \text{Appendicitis}$
Necessary and sufficient conditions	$\text{Appendicitis} \equiv \exists \text{AssociatedMorphology.Inflammation} \sqcap \exists \text{FindingSite.AppendixStructure}$
General inclusions	$\text{Ulcer} \sqcap \exists \text{has-location.Stomach} \sqsubseteq \text{Ulcer} \sqcap \exists \text{has-location.}(\text{Lining} \sqcap \exists \text{part-of.Stomach})$
Class disjointness	$\text{BodyPart} \sqcap \text{Organism} \sqsubseteq \perp$
Domain restrictions	$\exists \text{has-location.} \top \sqsubseteq \text{Disease}$
Range restrictions	$\top \sqsubseteq \forall \text{has-location.BodyPart}$
Role hierarchy	$\text{Proper-part-of} \sqsubseteq \text{part-of}$
Role reflexivity	$\varepsilon \sqsubseteq \text{part-of}$
Role transitivity	$\text{part-of} \circ \text{part-of} \sqsubseteq \text{part-of}$
Right identity on roles	$\text{has-location} \circ \text{part-of} \sqsubseteq \text{has-location}$
Concrete domains	$\text{Minor} \equiv \text{Person} \sqcap <_{18\text{year}}(\text{age})$
Nominals	$\text{Kangaroo} \sqsubseteq \exists \text{has-origin.}\{\text{Australia}\}$
Class assertions	$\text{London} \in \text{GeographicLocation}$
Role assertions	$(\text{London, England}) \in \text{has-location}$

SNOMED CT
==
OWL EL^{++}

SNOMED CT中的身体部位: *SEP-Triplets*

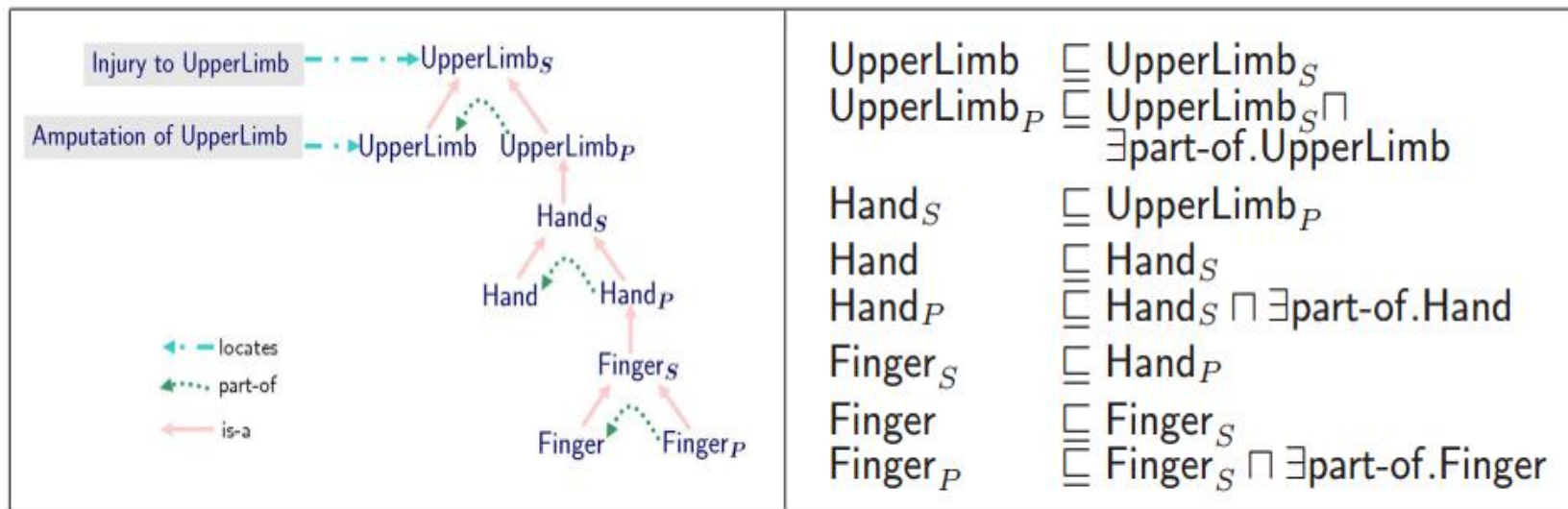


Fig. 1. Complete SEP-triplets in SNOMED CT.

SNOMED CT本体化: *EL++ for SEP-Triplets*

Finger \sqsubseteq BodyPart $\sqcap \exists \text{proper-part-of}.\text{Hand}$	(1)
Hand \sqsubseteq BodyPart $\sqcap \exists \text{proper-part-of}.\text{UpperLimb}$	(2)
UpperLimb \sqsubseteq BodyPart	(3)
AmputationOfFinger \equiv Amputation $\sqcap \exists \text{has-exact-location}.\text{Finger}$	(4)
AmputationOfHand \equiv Amputation $\sqcap \exists \text{has-exact-location}.\text{Hand}$	(5)
AmputationOfUpperLimb \equiv Amputation $\sqcap \exists \text{has-exact-location}.\text{UpperLimb}$	(6)
InjuryToFinger \equiv Injury $\sqcap \exists \text{has-location}.\text{Finger}$	(7)
InjuryToHand \equiv Injury $\sqcap \exists \text{has-location}.\text{Hand}$	(8)
InjuryToUpperLimb \equiv Injury $\sqcap \exists \text{has-location}.\text{UpperLimb}$	(9)
proper-part-of \circ proper-part-of \sqsubseteq proper-part-of	(10)
proper-part-of \sqsubseteq part-of	(11)
part-of \circ part-of \sqsubseteq part-of	(12)
ϵ \sqsubseteq part-of	(13)
has-exact-location \sqsubseteq has-location	(14)
has-location \circ proper-part-of \sqsubseteq has-location	(15)

Fig. 2. A re-engineered extract of SNOMED CT without SEP-triplets.

SNOMED CT to OWL

- SNOMED CT OWL Reference Set: 2018 July International Edition
 - Draft: 2018 Jan Edition
- [SNOMED CT to OWL toolkit](#)

3.1. Feature Overview

SNOMED CT supports the following types of class restrictions:

- existential quantification to a class expression (ObjectSomeValuesFrom)
- intersection of classes (ObjectIntersectionOf)
- existential quantification to a literal (DataHasValue)

SNOMED CT supports the following axioms, all of which are restricted to the allowed set of class expressions:

- class inclusion (SubClassOf)
- class equivalence (EquivalentClasses)
- object property inclusion (SubObjectPropertyOf) with or without property chains, and data property inclusion (SubDataPropertyOf)
- class disjointness (DisjointClasses)
- property equivalence (EquivalentObjectProperties and EquivalentDataProperties),
- transitive object properties (TransitiveObjectProperty)
- reflexive object properties (ReflexiveObjectProperty)



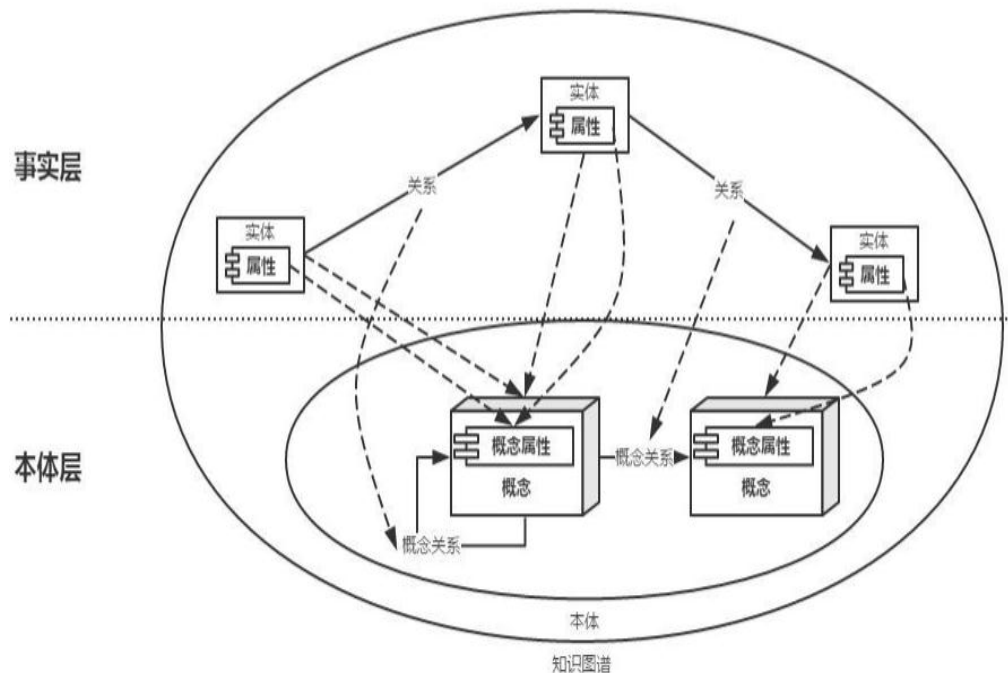
SNOMED CT OWL Guide

Linked with [Benevolent Medical Terminology](#)
SNOMED CT document sharing [http://www.snomed.org](#)
Publication Date: 2018-07-11

© Copyright 2018 International Health Terminology Standards Development Organisation

术语系统，本体，知识图谱，知识库

- 术语系统
 - 偏重term，不够形式化
 - Language-centered view
- 本体
 - 偏重概念模型，可以有规则(公理)和实例
 - is the study of *what there is*
- 知识图谱
 - 偏重实体及关系, 图模型
 - 常包含一个本体(schema)
- 知识库
 - 泛指知识的一种组织形式



医疗知识图谱应用(1): 问诊对话

APP



在候诊区使用

生成病历中的
病史信息

生成辅助诊断

机器人会话交互

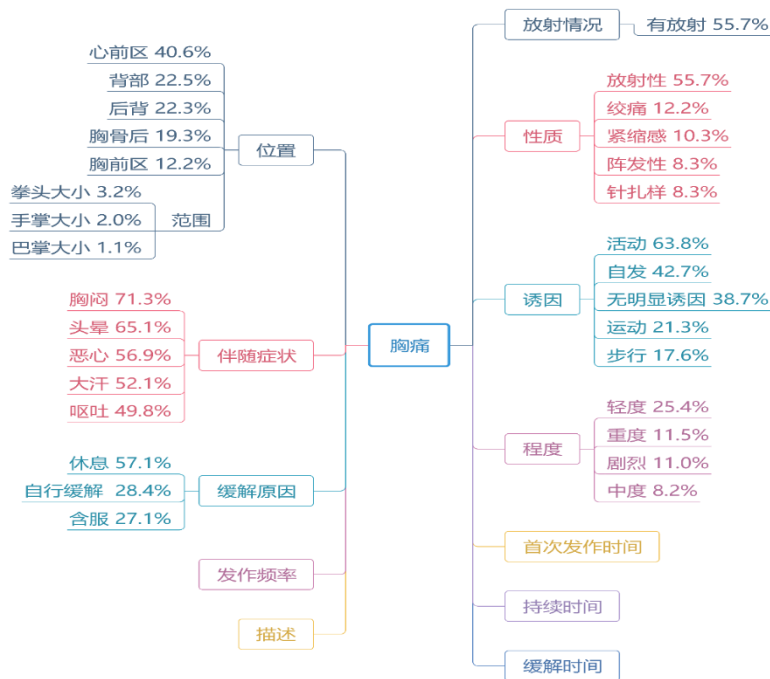


问诊对话依赖的知识

以症状作为入口
把症状的发生看成是事件

- 发生部位
- 发生时间和频率
- 缓解，加重因素
- 程度
- 性状
- 伴随症状

胸痛 is a Class



医疗知识图谱应用(2): 病历质控

主诉

形式质控

不能缺项
长度不超过20个字

内涵质控

包括主要症状和持续时间
能导出第一诊断
主诉和现病史相关相符

入院
诊断

不能缺项

诊断确切，依据充分
主次排列有序
无漏诊，误诊

病历质控依赖的知识

- 以疾病诊断为入口，判断诊断相关信息是否完整和一致
 - 疾病相关症状
 - 疾病相关个人史，既往史，家族史等
 - 疾病的诊断标准
 - 疾病是否特点人群相关
 - 疾病发生部位

胸痛 is an Instance

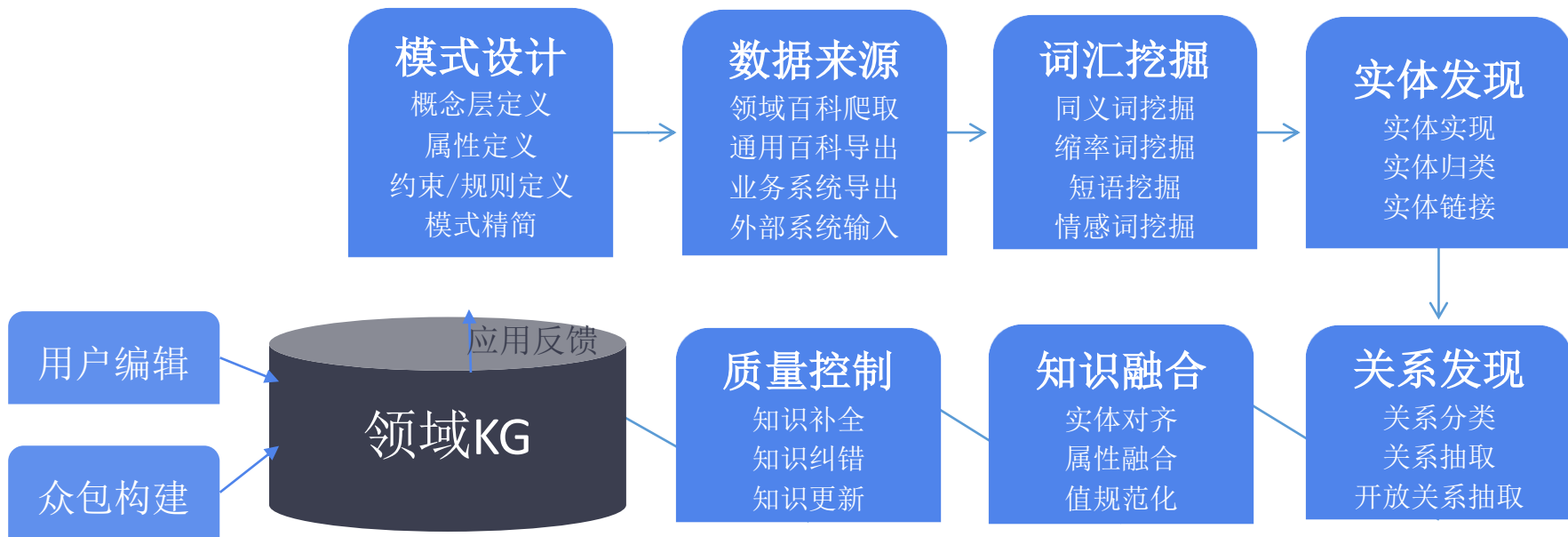
如何对知识图谱做测试

评估方法	方法说明	评级层次
基于黄金标准评估	将所构建的本体与黄金标准（一个公认的比较成熟的本体或是人工标注术语集）进行比较，罗列出其不足并进行改进。	词汇数据层，层级分类层，语义关系层
基于本体任务/应用的 本体评估	一个特定应用环境中，测试一组本体，看哪个本体最适合该应用，这些应用包括搜索、问答、推荐、决策等。	词汇数据层，层级分类层，语义关系层，应用层
数据驱动评估	通过衡量本体与领域语料的匹配度或本体的领域覆盖度来评估本体，或使用其他参考数据来辅助本体评估过程，这种方法常与文本分析、机器学习技术结合	词汇数据层，层级分类层，语义关系层
基于指标的评估（人工 评估）	基于一套预先定义好的原则、准则、标准等进行评估的方法，其多是从构建本体的原则来评估本体。	词汇数据层，层级分类层，语义关系层，应用层

我们的评估指标

评估方法	评估指标
基于指标的评估	一致性：是否存在一个term用在多个不同的地方； 精确性（人工）：是否存在多个实体表示同一个意思； 正确性（人工）：实体的属性，关系是否正确； 相关性（人工）：是否跟领域紧密相关
基于黄金标准的评估	Term覆盖率，关系的准确率和覆盖率 （如：以CCKS 2017和2018医疗实体评测为金标准）
基于应用的评估	基于知识图谱的NER效果和病历质控等应用效果

领域知识图谱的一般构建方法



医疗领域知识图谱构建的挑战

- 
- 冷启动
 - 敏捷构建
 - 缺少医疗专家

医疗知识图谱现状

国外

UML-S

FMA

DrugBank

SNOMED
CT

ICD-10

LinkedLife
Data

LOINC

GALEN

Bio2RDF

RxNorm

Gene
Ontology

More...

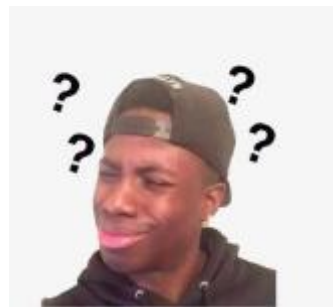
国内(中文版)

MESH

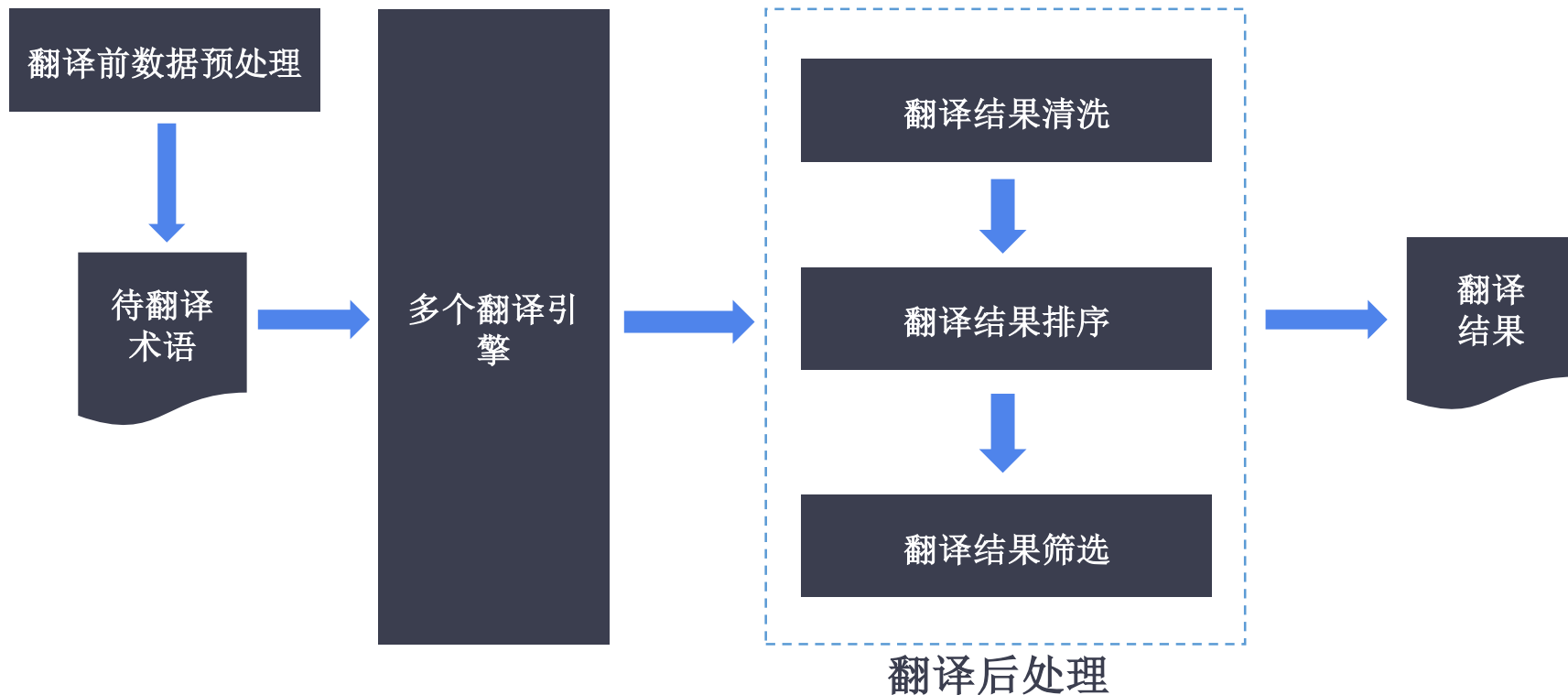
ICD-9-CM

ICD-10

症状知识图谱
@OpenKG



冷启动-基于UMLS的中文汉化



翻译结果排序和筛选

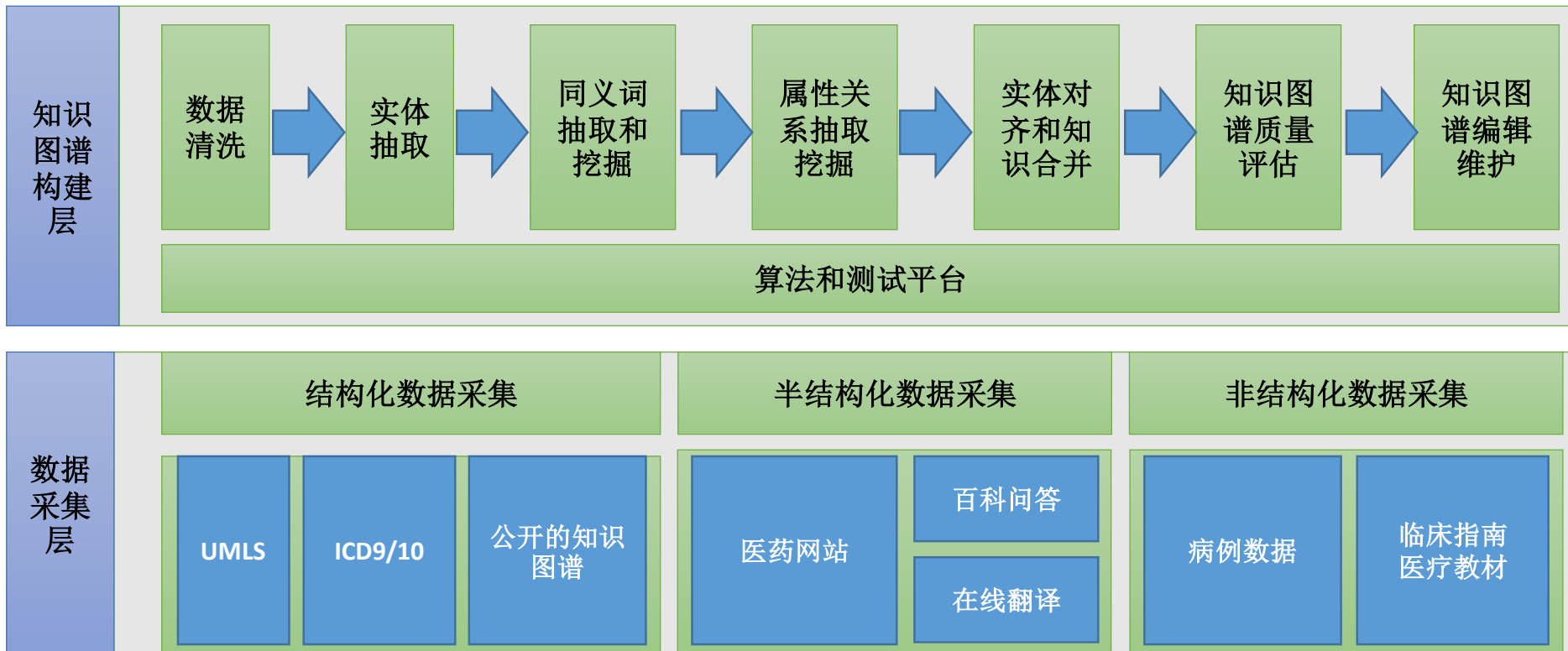
- 排序方法（按照以下依据依次比较）

- A. 是否标记为医疗专业术语 [医]
- B. 在病历文本中出现的频次
- C. 在海量医疗文本中出现的频次
- D. 在海量通用语料中出现的频次
- E. 候选翻译术语的最大长度

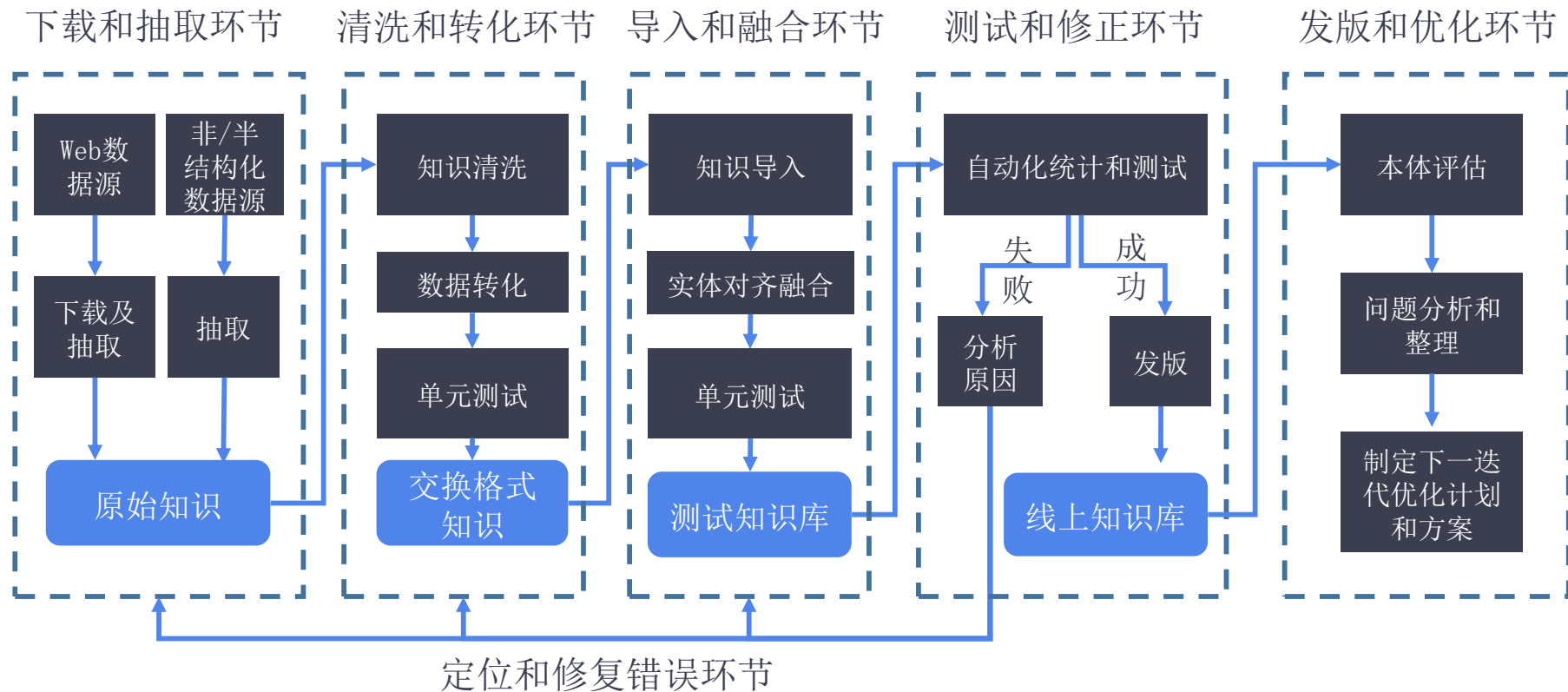
- 筛选依据

- A. 满足A，或 $B > k$ ，或 $C > k$ ，都认为是正确的翻译结果
- B. 其他需要人工处理

知识图谱构建总体架构



敏捷构建



演化和版本管理

- 版本描述

- 版本号及改进点描述
- 本版本与上一版本的差异比较结果文件：增加或减少了哪些概念，属性，实体，关系等

- 版本比较

- 根据不同版本的知识库文件，生成比较结果文件

- 版本恢复

- 按照操作日志回滚（轻量级）
- 根据版本差异比较文件恢复
- 根据版本备份文件恢复（回溯到某个发布版本）

- 版本发布

- 每一次开发迭代都要完成发版
- 发版要生成完整的版本描述和备份文件并存档

知识图谱的迭代

- **发现和梳理知识图谱存在的问题**

- 以应用为导向，根据知识图谱的应用效果，提出改进意见
- 领域专家通过可视化操作平台，分析知识库，并找出存在的问题
- 根据知识图谱构建各环节缺陷提出待优化的问题

- **版本规划**

- 列出所有需要改进的点
- 综合考虑实际应用需求、知识图谱质量要求以及开发成本等因素排优先级
- 确定下一版本发版计划和开发方案，并给出发版号

- **对知识图谱各个环节展开开发**

- 按照预先讨论的方案开发
- 中间知识图谱要快速生成，以确保其他环节有可用知识
- 做好单元测试，最大限度减少传递错误的次数

缺少医疗领域专家，怎么办

- 利用少数的医疗专家做人工抽查评测
- 利用网上公开的知识
- 从文本中挖掘知识：病历，教材，百科等

医疗知识图谱敏捷构建的体会总结

- 应用驱动
- 敏捷构建最关键的两个技术点
 - 自动化评测
 - 知识融合
- 医疗知识图谱的形式化是下一步的重心
 - 知识本体的表示与推理: OWL EL++ with Meta-Modelling support
 - 病历大数据结合下的海量数据查询与推理

谢谢！