

Ontology and Its Application

本体方法及其应用

甘健侯 姜 跃 夏幼明◎著 |



科学出版社

Ontology and
Its Application

本体方法及其应用

——甘健侯 姜 跃 夏幼明◎著

科学出版社

北京

研究基地

民族教育信息化教育部重点实验室
中国科学院计算机网络信息中心
云南省高校智能信息处理重点实验室
昆明理工大学冶金与能源工程学院

前言

·本·体·方·法·及·其·应·用·



本体论原本是一个哲学概念，是表达哲学理论的一个术语。20世纪90年代初，本体概念被广泛地引用到计算机领域特别是人工智能和知识工程研究领域。本体已经成为知识工程、自然语言处理、协同信息系统、智能信息集成、Internet智能信息获取、知识管理等各方面普遍研究的热点。

首先，本书介绍本体的概念、本体的组成、建立本体的原则和一般方法、本体的常用关系、本体常用开发工具以及典型本体等。其次，本书介绍本体描述语言，包括XML、RDF/RDFS、OWL等。

本书结合旅游信息资源本体、高校就业管理领域本体、常用软件本体构建介绍本体的开发步骤及开发过程中存在的问题，使读者对本体及其构建有一个较清晰的认识，为其深入的研究和应用打下一个良好的基础。

本体映射是解决本体异构问题的手段之一，概念语义相似度和相关度计算对实现本体集成和信息的语义检索起重要的作用。本书介绍本体映射相关概念、方法和典型系统，并对基于本体的概念语义相似度和相关度进行探讨和研究。

描述逻辑是基于对象知识表示的形式化，依据提供的构造器，在简单的概念和关系上构造出复杂的概念和关系。本书在基本描述逻辑ALC的基础上添加构造器——最大数量约束($\leq nR.C$)、最小数量约束($\geq nR.C$)、传递关系(R^+)、反关系(R^-)、关系并($R_1 \sqcup R_2$)、关系复合($R_1 \circ R_2$)、个体实例集($\{a\}$)，提出扩展描述逻辑 ALC^+ ，通过证明得到 ALC^+ 的一些性质。

在以上研究的基础上，本书以“软件开发”领域为例，设计并开发一个基于Lucene和本体的语义检索原型系统；基于科学家资源服务领域，对已有数据进行完善、更新和管理，构建科学家资源领域本体，并介绍基于本体的科学家资源服务平台框架。

从2003年开始，在云南师范大学林毓材教授的指导和帮助下，作者在语义Web领域进行系统研究。在本书出版之际，向林毓材教授表示衷心感谢。

感谢中国科学院计算机网络信息中心阎保平研究员、昆明有色冶金设计研究院谢刚教授，两位恩师在百忙中审阅本书初稿，并给出修改意见和建议，让作者受益匪浅。

云南师范大学徐天伟、文斌、李金绪、刘江涛、张姝老师，保山学院段寿建老师，中国人民银行昆明中心支行胡绍波等参加了课题的研究，并做出了很好的科研成果。

本书得到国家自然科学基金项目、中国科学院CNIC创新基金项目、云南省自然科学基金项目、云南省教育厅自然科学基金重点项目、云南师范大学学术著作出版基金的资助。在此，一并表示感谢。

由于作者水平有限，在完整性、准确性等方面难免存在问题。希望得到广大读者特别是专家和同行的指点，共同探讨有关问题，交流研究经验，使研究工作取得更大的进步。

—
王健侯

2011年5月10日于中国科学院

目 录

·本·体·方·法·及·其·应·用·



前言

第一篇 本体与本体描述语言 1

◆ 第1章 本体基础	3
● 1.1 本体概述	3
● 1.2 本体的组成	4
● 1.3 本体建立的原则	5
● 1.4 本体建立的一般方法	5
● 1.5 本体描述语言	7
● 1.6 本体中的常用关系	8
1.6.1 IS-A 关系	9
1.6.2 Instance-Of 关系	9
1.6.3 Member-Of 关系	10
1.6.4 Before 关系和 After 关系	10
● 1.7 常用的本体开发工具	11
1.7.1 Protégé	11
1.7.2 Apollo	12
1.7.3 OILED	12
1.7.4 OntoEdit	13
1.7.5 OntoSaurus	13
1.7.6 WebODE	13
● 1.8 其他工具	14
1.8.1 Jena 简介	14
1.8.2 Lucene 简介	15
● 1.9 典型本体介绍	18
1.9.1 CYC	18

1.9.2 WordNet	19
1.9.3 SUMO	19
1.9.4 知网	19
1.9.5 国家知识基础设施	20
1.9.6 CREAM	21
1.9.7 OntoWebber	21
1.9.8 其他模型	22
● 1.10 本体的研究和应用	23
◆ 第2章 语义 Web 与本体描述语言	24
● 2.1 语义 Web 概述	24
2.1.1 语义 Web 的概念、定义	24
2.1.2 语义 Web 的模型	25
2.1.3 Unicode 和 URI	26
2.1.4 本体层	26
2.1.5 逻辑、证明和信任	27
2.1.6 数字签名和加密	27
● 2.2 本体描述语言	27
2.2.1 XML	28
2.2.2 RDF	29
2.2.3 RDFS	34
2.2.4 OWL	36
● 2.3 OWL 本体语言的描述	37
2.3.1 命名空间定义	38
2.3.2 本体头定义	38
2.3.3 类定义	39
2.3.4 个体定义	43
2.3.5 属性定义	45
● 2.4 OWL 类构造器和原子解释	48
2.4.1 OWL 类构造器	48
2.4.2 OWL 原子解释	49
● 2.5 OWL 实例	50
● 2.6 语义 Web 的应用	55
2.6.1 智能信息检索	55
2.6.2 企业间数据交换及知识管理	55
2.6.3 Web 服务	55
2.6.4 基于代理的分布式计算	56

2.6.5 基于语义的数字图书馆	56
◎ 2.7 语义 Web 研究面临的问题和挑战	56
第二篇 本体技术	59
◆ 第3章 本体构建	61
◎ 3.1 构建旅游信息资源本体	61
3.1.1 构建旅游信息资源本体的目标	61
3.1.2 旅游信息资源本体构建过程	61
3.1.3 确定本体范围和术语	61
3.1.4 定义类和类的层次体系	62
3.1.5 定义类的属性	62
3.1.6 生成实例	63
◎ 3.2 高校就业管理领域本体构建	63
3.2.1 枚举领域本体的重要术语	63
3.2.2 复用现有的本体	63
3.2.3 定义类和类层次	63
3.2.4 定义类的属性	64
3.2.5 生成实例	64
◎ 3.3 常用软件本体构建	64
3.3.1 定义类和类的层次体系	64
3.3.2 定义常用软件的属性	66
3.3.3 创建常用软件实例	66
3.3.4 规则定义	66
3.3.5 常用软件领域知识推理系统的总体框架	67
◆ 第4章 本体映射	69
◎ 4.1 本体映射概述	69
4.1.1 本体异构及解决方案	69
4.1.2 本体映射概念及模型框架	70
◎ 4.2 常用的本体映射方法	72
4.2.1 基于语法的映射方法	72
4.2.2 基于概念实例的映射方法	73
4.2.3 基于概念定义的映射方法	74
4.2.4 基于概念结构的映射方法	74
4.2.5 基于规则的映射方法	75
4.2.6 基于统计学的映射方法	75
4.2.7 基于机器学习的映射方法	75

4.2.8 本体代数方法	76
4.2.9 本体聚类方法	76
※ 4.3 本体映射方法的分类	76
4.3.1 模式级与实例级	77
4.3.2 匹配粒度（元素级与结构级）	77
4.3.3 基于语言与基于约束	77
4.3.4 匹配基数	78
※ 4.4 本体映射典型系统介绍	78
4.4.1 Cupid	78
4.4.2 COMA	79
4.4.3 SF 方法	79
4.4.4 OntoMorph 系统	80
4.4.5 S-Match 动态多维概念映射算法	80
※ 4.5 目前本体映射存在的问题	80
◇ 第5章 基于本体的概念语义相似度和相关度计算	83
※ 5.1 概念语义相似度和相关度研究概述	84
5.1.1 语义相似度和相关度的概念及两者的关系	84
5.1.2 常用的语义相似度和相关度计算方法	85
5.1.3 语义相似度和相关度的评估方法	86
5.1.4 概念语义相似度和相关度的研究现状	86
※ 5.2 基于知网的词语语义相似度计算研究	87
5.2.1 知网简介	87
5.2.2 基于知网的词语语义相似度计算	89
5.2.3 基于知网的词语语义相似度计算的改进与实现	91
※ 5.3 基于领域本体的概念语义相似度和相关度的计算研究	94
5.3.1 基于领域本体的概念语义相似度计算	96
5.3.2 基于领域本体的概念语义相关度计算	101
5.3.3 结合领域本体的语义相似度和语义相关度的计算方法	102
※ 5.4 基于概念相似度和相关度的查询扩展	103
5.4.1 查询扩展技术概述	103
5.4.2 基于本体的查询扩展	104
5.4.3 基于领域本体概念间相似度和相关度的查询扩展	104
第三篇 本体推理方法——描述逻辑	109
◇ 第6章 基本描述逻辑 ALC	111
※ 6.1 描述逻辑及其发展	111

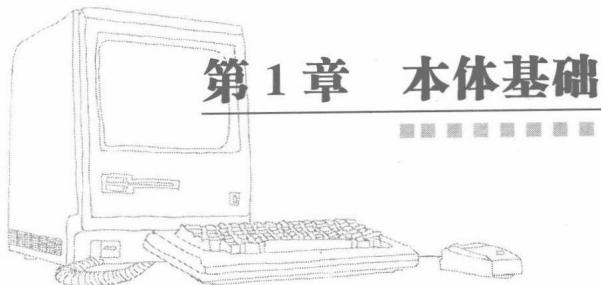
6.1.1 描述逻辑概述	111
6.1.2 描述逻辑的发展过程	113
6.1.3 描述逻辑的研究内容	114
6.2 基本描述逻辑 ALC 简介	114
6.2.1 基本描述逻辑 ALC 的语法与语义	114
6.2.2 基本描述逻辑 ALC 的知识库	115
6.2.3 基本描述逻辑 ALC 中的推理概述	116
6.2.4 基本描述逻辑 ALC 的推理过程	117
6.2.5 用 Tableau 算法进行推理的例子	120
◇ 第 7 章 扩展描述逻辑 ALC ⁺ 形式系统	124
7.1 扩展描述逻辑 ALC ⁺ 的形式化公理体系	124
7.1.1 扩展描述逻辑 ALC ⁺ 的语法	124
7.1.2 扩展描述逻辑 ALC ⁺ 的语义	125
7.1.3 扩展描述逻辑 ALC ⁺ 公理及其解释说明	127
7.2 扩展描述逻辑 ALC ⁺ 的基本性质	130
7.3 扩展描述逻辑 ALC ⁺ 系统的可靠性和完全性	149
7.4 扩展描述逻辑 ALC ⁺ 到谓词逻辑的转换	161
7.5 扩展描述逻辑 ALC ⁺ 与框架表示法的关系	162
7.5.1 框架及其组成	162
7.5.2 基于框架的知识库到扩展描述逻辑 ALC ⁺ 表示的知识库转换 过程	163
7.5.3 基于框架的知识库到扩展描述逻辑 ALC ⁺ 表示的知识库转换 示例	163
7.6 扩展描述逻辑 ALC ⁺ 与简单概念图的关系	164
7.6.1 简单概念图	164
7.6.2 简单概念图与扩展描述逻辑 ALC ⁺ 的关系对应	165
◇ 第 8 章 Web 本体语言 OWL 与扩展描述逻辑 ALC ⁺ 的关系	167
8.1 Web 本体语言 OWL 简介	167
8.1.1 Web 本体语言 OWL 的设计目标	167
8.1.2 Web 本体语言 OWL 的语法	168
8.2 扩展描述逻辑 ALC ⁺ 与 OWL 的对应	170
8.2.1 构造器的对应	170
8.2.2 扩展描述逻辑 ALC ⁺ 描述 OWL 中的部分公理	171
8.3 本体语言 OWL 描述的知识用扩展描述逻辑 ALC ⁺ 表示及推理的 示例	172

◆ 第9章 描述逻辑的应用	177
9.1 描述逻辑应用于概念建模	177
9.2 描述逻辑应用于软件工程领域	179
9.2.1 LaSSIE 系统和 CODEBASE 系统	180
9.2.2 CSIS 和 CBMS 系统	181
9.3 描述逻辑应用于语义 Web	183
第四篇 本体应用系统	185
◆ 第10章 基于本体的语义检索原型系统设计与实现研究	187
10.1 基于本体的语义检索模型设计	187
10.1.1 本体建立与管理模块	188
10.1.2 信息获取模块	189
10.1.3 Lucene 检索引擎模块	189
10.1.4 用户查询扩展和结果反馈模块	190
10.2 基于本体的语义检索原型系统设计与实现	190
10.2.1 系统开发平台及工具	191
10.2.2 基于本体的语义检索原型系统各模块的设计与实现	191
10.2.3 语义检索原型系统和传统信息检索系统的检索效果对比 ..	194
◆ 第11章 基于本体的科学家资源服务平台研究	199
11.1 概述	199
11.1.1 主要研究工作	199
11.1.2 技术方案	199
11.2 主要功能	200
11.2.1 科学家资源信息的获取	200
11.2.2 基于语义 Web 的科学家资源领域本体构建研究	200
11.2.3 基于本体的科学家信息资源网站自动生成技术研究	201
11.2.4 科学家资源个性化推荐技术研究	201
11.2.5 基于本体的科学家信息服务综合集成平台	201
11.3 科学家资源关系数据库数据模型构建	202
11.3.1 科学家基础数据	203
11.3.2 科学家科研项目数据	203
11.3.3 科学家科研成果数据	204
11.3.4 其他数据对象	205
11.3.5 科学家资源关系数据库中数据字段与本体推理的关系	205
11.4 科学家资源本体库构建	206
11.4.1 科学家资源概念层次树	206

11.4.2 科学家资源本体中的常用关系	206
11.4.3 谓词定义与扩展	207
11.4.4 操作符定义	208
11.4.5 IF-THEN 规则表示	209
11.4.6 科学家资源服务的基本知识推理	209
附录一 研究领域专业术语.....	211
附录二 重要的 Web 资源	212
附录三 RDF 类	213
附录四 RDF 属性	214
附录五 OWL 类	215

| 第一篇 |

本体与本体描述语言



第 1 章 本体基础



1.1 本体概述

20世纪90年代初，本体概念被广泛地引用到计算机领域特别是人工智能（AI）和知识工程研究中，因为AI和知识工程需要开发一个领域共享的、公共的概念，实现知识共享和重用。在AI领域，本体通常被称为领域模型（Domain Model）或概念模型（Conceptual Model），是关于特定知识领域内各种对象、对象特性以及对象之间可能存在关系的理论。通过对应用领域的概念和术语进行抽象，本体形成了应用领域中共享和公共的领域概念，可以描述应用领域的知识或建立一种关于知识的描述。本体的抽象可能是很高层次的抽象，也可能是针对特定领域的概念抽象。本体已经成为知识工程、自然语言处理、协同信息系统、智能信息集成、Internet智能信息获取、知识管理等各方面普遍研究的热点。因此，随着高度结构化的知识库在AI和面向对象系统中的出现，对于实际应用和理论研究，本体变得日益重要。

最近十年以来，各种研究机构和知识工程研究者提出了多种面向AI、具有细微差别的本体定义。

- (1) 一个本体是一个非形式的概念化系统；
- (2) 一个本体是由一个逻辑理论表示的概念系统；
- (3) 一个本体定义了组成主题领域词汇的基本术语和关系以及用于组合术语和关系以定义词汇的外延规则；

(4) 本体是概念模型的明确规范说明。

其中, Tom Gruber 的定义被引用最多, “本体是概念模型的明确规范说明”。Studer 等总结认为: “本体是共享概念模型明确的形式化规范说明。”

这包含 4 层含义: 概念模型 (Conceptualization)、明确 (Explicit)、形式化 (Formal) 和共享 (Share)。

“概念模型”指通过抽象出客观世界中一些现象的相关概念而得到的模型, 概念模型所表现的含义独立于具体的环境状态; “明确”是指所使用的概念及使用这些概念的约束都有明确的定义; “形式化”指本体是计算机可读的 (即能被计算机处理); “共享”指本体中体现的是共同认可的知识, 反映的是相关领域中公认的概念集, 即本体针对的是团体的共识。

根本上, 本体的作用是为了构建领域模型。例如, 在知识工程过程中, 一个本体提供了关于术语概念和关系的词汇集, 通过该词汇集可以对一个领域进行建模。虽然不同的本体之间存在一些差异, 但它们之间存在普遍的一致性。针对应用领域中一些特殊的任务, 知识表达可能还需要一种在很高的普遍性层次上的本体抽象概念。

当语义通过一定的形式添加到网络资源上之后, 下一步工作就是如何使得这些资源被理解和共享。对于 Web 上不同的数据资源, 它们对同一个概念可能采用不同的标识符。如对于“下载工具”这个概念, 可以使用〈DownloadTools〉, 也可以使用〈Download_tools〉。为了识别这些标记所代表的概念, 将本体论的方法引入到语义 Web 中来。

在语义 Web 中, 本体具有非常重要的地位, 它是解决语义层次上 Web 信息共享和交换的基础。就 Web 而言, 本体可以应用在如下方面: ①提高搜索引擎的精确度, 它只需根据元数据查找网页, 而不会像现在用语义含糊的关键词进行全文搜索; ②利用本体从 Web 页面到相应的知识结构和推理规则建立关系; ③本体还可用于电子商务网站, 使买卖双方可以基于机器进行交流等。

1.2 本体的组成

一个本体, 一般首先给出一组概念的层次性结构, 概念间的包含关系、组成关系、划分关系等。

(1) 分类层次结构: 如常用软件、旅游信息资源中的分类层次。

(2) 按 IS-A 和 Part-of 关系组织、组成概念结构。

(3) 概念的语义描述, 应该不局限于静态结构。其他如“先后关系”、“因果关系”以及语义复杂的“参照关系”, 具有丰富语义的关系往往无法清楚地表达出它们的语义来。

一般来说, 一个本体可由概念类、关系、函数、公理和实例等 5 种元素组成。

(1) 本体中的概念是广义上的概念，它既可以是一般意义上的概念，也可以是任务、功能、行为、策略、推理过程等。本体中的这些概念通常构成一个分类层次。

(2) 本体中的关系表示概念之间的一类关联，典型的二元关联如子类关系形成概念类的层次结构，一般情况下用 $R: C_1 \times C_2 \times \dots \times C_n$ 表示概念类 C_1, C_2, \dots, C_n 之间存在 n 元关系 R 。

(3) 函数是一种特殊的关系。其中，第 n 个元素相对于前 $n-1$ 个元素是唯一的，一般情况下，函数用 $F: C_1 \times C_2 \times \dots \times C_{n-1} \rightarrow C_n$ 表示。

(4) 公理用于表示一些永真式。更具体地，在许多领域中，函数之间或关联之间也存在着关联或约束。

(5) 实例是指属于某概念类的基本元素，即某概念类所指的具体实体，特定领域的所有实例构成的领域概念类在该领域中的指称域。

1.3 本体建立的原则

Tom Gruber 给出了 5 条设计本体的基本准则：

(1) 明确性和客观性：本体应该有效地传达所定义的术语内涵。

(2) 一致性：一个本体应该前后一致，即由它推断出来的概念定义应该与本体中的概念定义一致。

(3) 可扩展性：可扩展性是指一个本体提供一个共享的词汇，它应该在预期的任务范围内提供概念的基础，同时，它的表示应该使得人们能够单调地扩展和专门化说明这个词汇，即人们应该能够在不改变原有定义的前提下，以这组存在的词汇为基础定义新的术语。

(4) 最小编码偏差：本体应该处于知识的层次，而与特写的符号级编码无关。

(5) 最小本体承诺：一个本体应该在提供必需的共享知识条件下，要求有最小的本体承诺。

除了上述原则外，J. Arpirez 等提出本体设计应该坚持如下几点原则：

(1) 尽可能使用标准术语；

(2) 同层次概念之间保持最小的语义距离；

(3) 可以使用多种概念层次，采用多重继承机制来增加表达能力。

1.4 本体建立的一般方法

常用的本体开发方法有：

Uschold 和 King 的“骨架法”：由英国爱丁堡大学 AI 应用研究所基于开发