



“十二五”普通高等教育本科国家级规划教材



普通高等教育“十一五”国家级规划教材

# 化学信息学

第二版

李梦龙 文志宁 主编

李益洲 郭延芝 战玉华 副主编

译外  
学外  
译外

HUAXUE XINXUE



化学工业出版社



“十二五”普通高等教育本科国家级规划教材



普通高等教育“十一五”国家级规划教材

# 化学信息学

第二版

# 化 子 口 忌 子

欣喜的是，现在越来越多的化学家开始对秦和这个团队所做量子化学数据的分析中，

李益洲 郭延革 战玉华 副主编

《数据科学与机器学习》



化 学 工 业 出 版 社

· 北京 ·

《化学信息学》(第二版)是“十二五”普通高等教育本科国家级规划教材、普通高等教育“十一五”国家级规划教材。本书在第一版的基础上，着重针对中、外文期刊网络数据库的检索过程进行了介绍，并扩充了专利文献与技术标准数据库的相关内容。全书主要分为四大部分，共12章，其中第1章概述了化学信息学的产生及特点；第2~8章讲述了化学信息的来源，包括纸质版的手册，书籍，搜索引擎，目前广为使用的中、外文期刊文献数据库以及专利文献与技术标准数据库；第9~11章介绍了化学信息的处理工具(即化学软件)、处理方法(相关化学计量学算法)以及定量构效关系(QSAR)的原理及应用，并新增了部分例子；在第12章中，对生物信息学领域的研究进行了概述。

《化学信息学》(第二版)可作为高等院校化学、化工及相关专业本科“化学信息学”课程的入门教材，另外，书中提供了大量与化学信息学相关的网址，亦可作为研究生的参考书籍。

# 学信息学

## 第二版

主编 宁志文 李楚李

图书在版编目(CIP)数据

化学信息学/李梦龙，文志宁主编. —2版. —北京：化学工业出版社，2018.7

“十二五”普通高等教育本科国家级规划教材 普通高等教育“十一五”国家级规划教材

ISBN 978-7-122-32235-7

I. ①化… II. ①李… ②文… III. ①计算机应用-化学-信息检索-高等学校-教材 IV. ①G252.97

中国版本图书馆CIP数据核字(2018)第112648号

责任编辑：杜进祥

文字编辑：向东

责任校对：边涛

装帧设计：关飞

出版发行：化学工业出版社（北京市东城区青年湖南街13号 邮政编码100011）

印 装：三河市双峰印刷装订有限公司

787mm×1092mm 1/16 印张16 3/4 字数440千字 2018年9月北京第2版第1次印刷

购书咨询：010-64518888（传真：010-64519686）售后服务：010-64518899

网 址：<http://www.cip.com.cn>

凡购买本书，如有缺损质量问题，本社销售中心负责调换。

定 价：40.00元

版权所有 违者必究

# 前言

HUAXUE XINXIXUE  
化学信息学

化学是一门以实验为基础的古老学科，科学家们致力于探索新物质的各类化学属性。随着实验技术的进步，人们获取数据的能力已有了很大的提高。时至今日，面对呈指数级增长的化学数据，化学家们发现在本学科的探索中，最大的瓶颈已不再是如何获取未知物质的性质，而是如何从已有的实验数据中提取更多的化学信息，以总结规律。1995年，美国著名化学家 Brown 指出，化学家习惯于将99%的精力和资源用于数据的收集上，其余1%用于数据的分析和处理，将数据转化为信息。

庆幸的是，现在越来越多的化学工作者们开始注意到这个问题并投身于化学数据的分析中，化学信息学(cheminformatics)在这种情况下应运而生，它是汇集化学、数学、信息科学等交叉学科知识的研究领域，主要是通过对化学信息的检索、整理、分析以及可视化，最终完成将数据转化为信息的过程。德国 Johann Gasteiger 出版的《化学信息学教程》一书亦指出，化学信息学的任务就是运用信息学的方法来解决化学的问题。

目前，化学信息学主要涵盖了化学信息的获取、化学信息的表达以及化学信息的处理三个方面的内容。作为一门新的基础课程，如何尽快让化学专业的本科生了解并掌握其中涉及的概念、方法以及网络资源是该学科建设亟待解决的主要问题。

本书第一版是“十二五”普通高等教育本科国家级规划教材、普通高等教育“十一五”国家级规划教材。本次修订，着重对如下内容进行了更新：(1)着重针对中文及外文期刊网络数据库检索过程的介绍进行了扩充，单独在第4章、第5章对目前常用的数据库进行了详细的介绍；(2)在第6章中，对第一版未曾提到的专利文献与技术标准进行了介绍；(3)在生物信息学相关章节中，对目前热门的领域，如精准医疗等进行了介绍；(4)对第一版各章节中的网络资源进行了更新，并新增了一些例子，便于学生们在学习过程中对内容进行理解。

本书由李梦龙、文志宁任主编，李益洲、郭延芝、战玉华任副主编。感谢清华大学图书馆的战玉华老师撰写了本书第4章以及在后期修改过程中提出的宝贵意见。在本书编写过程中，笔者实验室的博士研究生柳媛以及硕士研究生梁羽、谢凡凡、覃柳、吉越、胡文、郝颖异、黄雨瑶、郭佳丽等同学亦参与了资料的收集与整理工作，另外，化学工业出版社在此过程中提供了大力的支持，在此一并表示诚挚的感谢。

由于化学信息学涉及面广，编者的水平和时间有限，书中不足之处在所难免，恳请广大读者批评指正。

## 2.4 图书编者

2.4.1 生命科学图书馆

2.4.2 中国科学院大连化学物理研究所图书馆

2.4.3 中国科学院国家科学图书馆

编者

2018年3月

# 目 录

哈工大出版社

HUAXUE XINXIXUE

化学信息学

## 第1章 概述

1.1 什么是化学信息学	1
1.2 化学信息学的诞生背景	1
1.3 信息科学在化学领域的应用	2
1.4 化学信息采集接口	2
1.5 化学信息的结构和特点	3
1.6 化学信息的工作方式	4
1.7 化学信息学的应用	4
1.7.1 化合物结构绘制	4
1.7.2 化学数据库设计与开发	4
1.7.3 化学反应体系模拟	5
1.7.4 计算机辅助波谱解析	5
1.7.5 化合物结构与活性关系预测	5
1.7.6 实验室信息管理系统	5
1.8 展望	5

## 第2章 化学信息来源

2.1 辞典	7
2.2 手册	7
2.3 化学期刊	9
2.3.1 综合类期刊	10
2.3.2 有机化学期刊	10
2.3.3 分析化学期刊	11
2.3.4 无机化学期刊	12
2.3.5 物理化学期刊	12
2.4 图书馆资源	13
2.4.1 生命科学图书馆	13
2.4.2 中国科学院大连化学物理研究所图书馆	13
2.4.3 中国科学院国家科学图书馆	14

2.4.4 国家科技图书文献中心	16
2.4.5 清华大学图书馆	17
2.4.6 中国国家图书馆	18
2.4.7 哈佛大学图书馆	19
2.4.8 斯坦福大学图书馆	19
<b>2.5 化学化工资源信息平台</b>	20
2.5.1 化学信息网	20
2.5.2 Computer Aided Chemistry Tutorial	21
2.5.3 Wilton High School Chemistry	22
<b>扩展阅读：信息</b>	22

## 第3章 信息搜索引擎 ..... 23

<b>3.1 概述</b>	23
3.1.1 搜索引擎的原理	23
3.1.2 搜索引擎的历史及发展趋势	24
<b>3.2 搜索引擎的定义及分类</b>	27
3.2.1 全文搜索引擎	27
3.2.2 目录索引类搜索引擎	27
3.2.3 元搜索引擎	27
3.2.4 垂直搜索引擎	28
<b>3.3 搜索引擎查询方法</b>	28
3.3.1 模糊查询	28
3.3.2 精确查询	29
3.3.3 逻辑查询	29
3.3.4 查询范围限制	29
<b>3.4 常用的全文搜索引擎</b>	30
3.4.1 百度	30
3.4.2 Google 中国	31
3.4.3 维基百科	31
<b>3.5 常用的学术搜索引擎</b>	33
3.5.1 BASE	33
3.5.2 SciTech Connect	34
3.5.3 CiteSeerX	35
<b>3.6 常用的元搜索引擎</b>	35
3.6.1 Dogpile	35
3.6.2 Excite	36
3.6.3 Ixquick	36
3.6.4 Mamma	37
3.6.5 Metacrawler	38
3.6.6 ProFusion	38

3.6.7	Savvysearch	39
<b>3.7 专业搜索引擎</b>		39
3.7.1	专业搜索引擎的优势	39
3.7.2	著名的专业搜索引擎	39
<b>扩展阅读：百度搜索技巧</b>		41

## 第4章 文献检索概述

<b>4.1 文献检索的意义</b>		42
4.1.1	文献的定义	42
4.1.2	信息检索的含义	42
4.1.3	文献检索的意义	42
<b>4.2 文献检索的基础知识</b>		42
4.2.1	逻辑运算符	42
4.2.2	通配符	43
4.2.3	位置算符	43
4.2.4	主要检索字段	44
4.2.5	文献资源分类	45
4.2.6	主要文献类型	45
<b>4.3 文献检索基本流程</b>		46
<b>4.4 文献检索途径及案例</b>		47
4.4.1	研究类文献检索	47
4.4.2	数据事实检索	50

## 第5章 中文学术论文资源数据库

<b>5.1 中国知网</b>		65
5.1.1	中国知网概况	65
5.1.2	中国学术期刊（网络版）	65
5.1.3	中国学术辑刊全文数据库	65
5.1.4	中国博士学位论文全文数据库	67
5.1.5	中国优秀硕士学位论文全文数据库	71
5.1.6	中国重要会议论文全文数据库	71
5.1.7	CNKI 国家科技成果数据库	72
<b>5.2 万方数据知识服务平台</b>		73
5.2.1	万方数据知识服务平台概况	73
5.2.2	中国学术期刊数据库	74
5.2.3	中国学位论文全文数据库	74
5.2.4	中国学术会议文献数据库	75
5.2.5	中文科技报告数据库	76
<b>5.3 维普网</b>		77
5.3.1	维普期刊资源整合服务平台概况	77

5.3.2 中文科技期刊数据库	77
5.3.3 中国科学指标数据库	879
<b>第6章 外文学术论文资源数据库</b>	81
<b>6.1 Web of Science</b>	81
6.1.1 Web of Science 概况	81
6.1.2 Web of Science Core Collection	83
6.1.3 Derwent Innovations Index	91
6.1.4 INSPEC	93
6.1.5 Chinese Science Citation Database	93
6.1.6 Journal Citation Reports	94
<b>6.2 美国《化学文摘》</b>	94
6.2.1 美国化学文摘概况	94
6.2.2 美国《化学文摘》(网络版)	95
6.2.3 SciFinder 检索方式及检索结果	95
<b>6.3 《工程索引》</b>	104
6.3.1 Engineering Village 概况	104
6.3.2 Ei Compendex 检索方式及检索结果	105
<b>6.4 ScienceDirect</b>	109
6.4.1 ScienceDirect 数据库简介	109
6.4.2 ScienceDirect 检索方式及结果	109
<b>6.5 Wiley Online Library</b>	113
6.5.1 Wiley Online Library 数据库简介	113
6.5.2 Wiley Online Library 检索方式及检索结果	114
<b>6.6 SpringerLink</b>	118
6.6.1 SpringerLink 数据库简介	118
6.6.2 SpringerLink 检索方式及检索结果	118
<b>6.7 其他外文文献数据库</b>	120
6.7.1 Online Computer Library Center	120
6.7.2 ProQuest 平台	121
<b>第7章 专利文献与技术标准</b>	123
<b>7.1 专利文献</b>	123
7.1.1 专利概述	123
7.1.2 美国专利文献概述	125
7.1.3 中国专利文献概述	128
7.1.4 中国专利全文数据库	129
7.1.5 万方中外专利数据库	129
7.1.6 SIPO 中国及多国专利查询系统	130
<b>7.2 技术标准</b>	132

7.2.1	技术标准概述	132
7.2.2	国际标准化组织(ISO)简介	132
7.2.3	ISO9000族标准及ISO14000族标准简介	132
7.2.4	中国标准的分类	133
7.2.5	国家标准全文数据库	134
7.2.6	中国行业标准全文数据库	134
7.2.7	中外标准数据库	135

## 第8章 化学信息数据库资源 ..... 136

8.1	数据库简介	136
8.1.1	数据	136
8.1.2	数据库	136
8.1.3	数据库管理系统	136
8.1.4	数据库系统(database system)	137
8.2	数据库历史及分类	137
8.2.1	数据库历史	137
8.2.2	数据库的模型分类	139
8.3	三类化学信息数据库	139
8.3.1	文献数据库	139
8.3.2	事实数据库	140
8.3.3	结构数据库	140
8.4	Internet上的化学化工数据库	140
8.4.1	Spectral Database for Organic Compounds	140
8.4.2	NIST Chemistry WebBook	141
8.4.3	NIST Atomic Spectra Database	142
8.4.4	Reaxys数据库	143
8.4.5	剑桥晶体数据库	144

## 第9章 化学软件 ..... 147

9.1	概述	147
9.2	化学软件的分类	148
9.3	语言软件和依托算法的化学计算软件	149
9.3.1	MATLAB	149
9.3.2	R语言	161
9.4	绘图软件	170
9.4.1	ChemBioOffice	170
9.4.2	ACD/ChemSketch5.0	171
9.4.3	Symyx Draw	173
9.5	化学分析仪器数据处理软件	174
9.5.1	GRAMS	174

9.5.2 MestReNova .....	177
9.5.3 Origin .....	178
<b>9.6 分子模拟软件 .....</b>	<b>179</b>
9.6.1 分子对接软件 AutoDock .....	179
9.6.2 量子化学计算软件——Gaussian 程序 .....	180
9.6.3 分子动力学模拟软件 .....	182
<b>第10章 信息处理与数据挖掘 .....</b>	<b>185</b>
10.1 概述 .....	185
<b>10.2 数据的标准化 .....</b>	<b>186</b>
<b>10.3 特征提取与优化 .....</b>	<b>186</b>
10.3.1 主成分分析 .....	186
10.3.2 偏最小二乘法 .....	189
10.3.3 逐步回归分析 .....	189
10.3.4 遗传算法 .....	191
<b>10.4 信号处理方法 .....</b>	<b>192</b>
10.4.1 协方差与相关系数 .....	193
10.4.2 自、互相关分析 .....	193
10.4.3 功率谱密度 .....	194
10.4.4 傅里叶变换 .....	194
10.4.5 小波变换 .....	195
<b>10.5 机器学习方法 .....</b>	<b>198</b>
10.5.1 K 最近邻法 .....	198
10.5.2 概率神经网络 .....	198
10.5.3 分类回归树 .....	199
10.5.4 助推法 .....	200
10.5.5 人工神经网络 .....	201
10.5.6 支持向量机 .....	204
<b>10.6 数据库挖掘技术 .....</b>	<b>206</b>
10.6.1 聚类算法 .....	206
10.6.2 决策树算法 .....	207
<b>10.7 Web 数据挖掘技术 .....</b>	<b>207</b>
10.7.1 Web 内容挖掘 .....	207
10.7.2 Web 结构挖掘 .....	207
10.7.3 Web 日志挖掘 .....	208
<b>扩展阅读：化学计量学 .....</b>	<b>208</b>
<b>第11章 QSAR及药物设计 .....</b>	<b>209</b>
11.1 概述 .....	209
11.2 QSAR 模型的分类 .....	210

11.2.1	二维定量构效关系 (2D-QSAR)	210
11.2.2	三维定量构效关系 (3D-QSAR)	212
11.2.3	多维定量构效关系	214
11.2.4	方法评价	215
<b>11.3 定量构效关系研究中常用的回归分析法</b>		<b>216</b>
11.3.1	多元线性回归	216
11.3.2	主成分回归	217
11.3.3	偏最小二乘回归	217
11.3.4	投影寻踪回归	218
11.3.5	非线性方法	219
<b>11.4 药物设计</b>		<b>219</b>
11.4.1	类药性	220
11.4.2	脂水分布系数 $\lg P$	220
11.4.3	脑血分配系数	220
11.4.4	肠穿透性	221
11.4.5	水溶性	221
11.4.6	毒性	221
<b>11.5 精准医疗</b>		<b>221</b>
<b>11.6 高通量筛选</b>		<b>222</b>
<b>11.7 QSAR 方法的应用示例</b>		<b>222</b>
<b>扩展阅读：定量构效关系</b>		<b>226</b>
<b>第12章 生物信息学</b>		<b>227</b>
<b>12.1 什么是生物信息学</b>		<b>227</b>
<b>12.2 生物信息学的发展历程</b>		<b>228</b>
<b>12.3 生物信息学的研究内容</b>		<b>230</b>
12.3.1	生物信息挖掘	230
12.3.2	药物设计	231
12.3.3	基因组学	231
12.3.4	蛋白质组学	232
12.3.5	代谢组学	234
<b>12.4 生物信息学的研究方法</b>		<b>235</b>
12.4.1	数学统计方法	235
12.4.2	动态规划方法	235
12.4.3	模式识别技术	235
12.4.4	数据库技术	235
12.4.5	分子模型化技术	235
12.4.6	分子力学和量子力学计算	236
12.4.7	分子动力学模拟	236
<b>12.5 生物信息学的应用</b>		<b>236</b>

12.6 生物信息学的研究趋势及未来挑战	237
12.7 蛋白质功能研究	238
12.8 蛋白质数据库简介	239
12.8.1 综合性蛋白质数据库	240
12.8.2 专用性蛋白质数据库	241
12.9 蛋白质序列的特征提取方法	242
12.9.1 基于氨基酸组成和位置的特征提取方法	243
12.9.2 基于氨基酸物理化学特性的特征提取方法	244
12.9.3 基于数据库信息挖掘的特征提取方法	245
12.10 蛋白质相互作用	247
12.11 蛋白质网络	252
12.11.1 蛋白质-蛋白质相互作用网络	252
12.11.2 氨基酸残基网络	254
扩展阅读：高通量基因表达检测技术	255

## 参考文献

# 第1章

## 概 述

材料、能源和信息是构成物质世界的三个基本要素。随着社会发展的需要，人们逐渐认识到信息的重要性，并创立了信息论与信息科学。20世纪90年代初，随着美国“信息高速公路”计划的提出，信息科学和信息产业出现了前所未有的飞速增长，成为这一时代的重要标志。同时信息科学加快了向传统科学渗透，化学中的信息学理论基础不断成熟。正是在这一背景下，结合其使用的计算机和因特网（Internet）工具，化学工作者在科研实践中促成了化学信息这一新兴化学分支的出现。化学信息学（chemoinformatics 或 cheminformatics）是化学领域中飞速发展起来的一个新的分支，是建立在多学科基础上的交叉学科，它利用计算机技术和计算机网络技术，对化学信息进行表示、管理、分析、模拟和传播，以实现化学信息的提取、转化与共享，揭示化学信息的实质与内在联系，促进化学学科的知识创新。

### 1.1 什么是化学信息学

诺贝尔化学奖获得者法国化学家 J. M. Lehn 在获奖报告中首次提出化学信息的概念，对化学的发展而言具有深远的影响，具有深刻的时代意义。虽然在此之前众多的化学工作者没有对化学信息展开实质性的工作，但是传统有机化学、无机化学、生物化学、材料化学以及在受体设计、超分子形成过程的结构化学等方面所积累的大量实验数据，为构建化学信息提供了基础。

化学信息学是个广义的概念，它包含对化学相关信息的设计、创造、组织、存储、处理、恢复、分析、再开发、可视化及应用。另一种关于化学信息学的定义是从药物研发的角度提出的，认为化学信息学是各种信息资源的混合体，目的是将数据转化成信息，再把信息转化成知识，以期更快、更准确地进行药物筛选和设计。

### 1.2 化学信息学的诞生背景

近十年来，由于计算机及网络技术向智能化、网络化方向发展，应用计算机技术能解决的化学问题也愈来愈多，化学工作者不仅获得了大量物质结构的信息，而且这些信息较从前也更为精确，计算机技术与化学的相互渗透已成为化学和计算机科学工作者的研究热点。由于计算机主要是通过数值计算来解决问题，其特点是能快速进行大量复杂、烦琐的数学运算，而化学是对化学物质进行认识、分析、合成及利用，从而使化学工作者能够对物质化学结构进行解析、表征、模拟与设计；能够处理复杂体系的电子结构、几何结构与其性能关系；完成微观分子工程设计与化学模拟；开展功能材料的研究；进行生物活性分子和药物分



子的相互作用机制及定量构效关系 (quantitative structure-activity relationship, QSAR) 研究; 探讨固体表面结构、固体表面轨道相互作用规律; 实现分子以上层次聚集体 (超分子体系、界面体系等) 结构和性能的模拟等。

然而, 纵观早期这一领域的工作, 仅仅涉及计算机技术的一些应用层次, 要想将计算机技术深入应用到化学中就必须解决化学与计算机的结合问题, 从化学工作者的角度应用和设计计算机软、硬件, 满足化学工作者处理化学信息的要求。该领域的研究包括计算机与分析仪器的接口、化学类应用软件程序包的开发、化学物质结构数据库的开发和查询。1973 年在荷兰举办了 “Computer Representation and Manipulation of Chemical Information” 会议, 1975 年美国化学会将期刊 “Journal of Chemical Documentation” 更名为 “Journal of Chemical Information and Computer Sciences”, 继而又于 2005 年更名为 “Journal of chemical information and modeling”。由此可见, 化学信息学已逐渐发展壮大, 它将给 21 世纪的化学带来全新的面貌。

### 1.3 信息科学在化学领域的应用

20 世纪中叶, Shannon 在 1948 年发表了关于信息论的著名文章, 提出了信息熵计算公式:

$$H(X) = - \sum_{i=1}^n \frac{1}{p(x_i)} \log_2 [p(x_i)] \quad (1-1)$$

式中,  $H(X)$  为事件  $X$  的信息熵, 它可由该事件当中所有可能出现的情况  $x_i$  的概率  $p(x_i)$  计算得到。此后信息理论开始了它的发展, 这一理论最早是与通信技术相关联, 但在其诞生后十年左右, 即从纯粹数学研究渗透到无线电、电视、雷达、心理学、语义学、经济学、生物学等领域。Wiener 认为信息的实质是负熵, 并强调信息这种负熵是在调节过程中相互交换而产生的。

化学科学中的分析化学从其诞生起就具有信息科学的特征, Kateman 等从三个方面阐述了分析化学的任务: 利用已有的分析方法, 提供关于物质化学成分的信息 (日常例行分析工作); 研究利用不同学科的原理、方法, 取得有关物质系统的相关信息的过程 (分析化学的科学工作); 研究利用现有分析方法取得关于物质系统的策略 (分析实验室的组织工作)。Kowalski 更是明确提出 “分析化学作为信息科学”, 他认为分析化学不仅在过去是一门信息科学, 现在仍然是一门信息科学, 在化学的各个分支学科中, 分析化学负担的任务与其他分支学科的不同之处在于分析化学的研究对象, 它并非某种具体的实物, 例如无机或有机材料, 而是与这些化学组成或结构相关的信息以及研究获取这些信息的最优方法与策略。

此外, 由于化学中熵的概念与 Shannon、Wiener 等提出的信息理论中的熵有着共同的基础, 这两门学科之间存在着深刻联系, 分析化学以发展分析信息理论作为其基础理论的组成部分, 获得了向前发展的动力。随后众多的化学家根据其从事的分析化学工作, 发表了多篇用于分析化学的信息理论系列论文, 其中捷克学者 Eckschlager 完成了在此领域的第一部专著。

### 1.4 化学信息采集接口

从图 1-1 中我们可以发现, 化学信息和化学实践之间是通过一个信息采集接口相连的,

这与其他化学分支学科明显不同。信息采集接口也是化学信息学科一个极为重要的组成部分，它是现实化学世界通往化学信息的桥梁，也是化学信息的生命源泉。信息采集接口和化学信息的三个层次构成了整个化学信息学科。

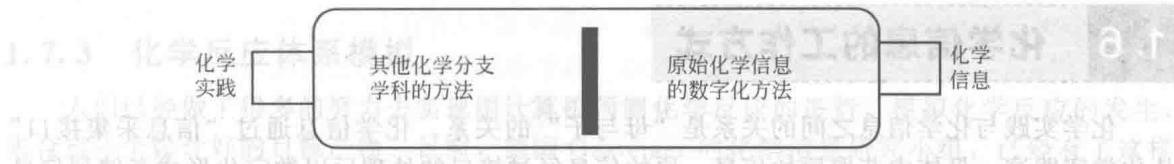


图 1-1 信息采集接口的内部结构

如图 1-1 所示，信息采集接口不是一套单一而又具体的软件或硬件，它实际上是一个综合了其他化学分支学科的某些方法以及原始化学信息数字化方法的方法集合。可以用作信息采集接口后端的方法众多、范围非常广泛，一般的各种化学图形处理软件、计算机化学应用软件、文字处理软件、数码照相机、具有 OCR (optical character recognition) 光学文字识别功能的扫描仪的录入系统等都可以作为信息采集接口的后端，因为它们都具有一个共同的特点，即能够把现实世界中的化学信息以数字化形式存储起来。对于接口的前端，只要是能从化学实践中获取化学信息的研究方法或仪器设备，如化学分析测量仪器、量子化学计算方法、分析化学及物理化学实验方法等，都可以用作信息采集接口的前端。

## 1.5 化学信息的结构和特点

仔细分析，我们可以发现化学信息的主体部分实际上是由三个层次构成的，即信息核心层、信息处理层和信息表示层，如图 1-2 所示，化学信息的这种分层结构本质上是计算机技术分层结构的反映。

在化学实践中产生出来并被计算机处理过的原始化学信息，比如在一个化学实验中发生的各种实验现象、记录的实验数据以及与化学实验相关的外界条件等，它们组成了化学信息的信息核心层。信息处理层由化学计量学方法、药物分子设计方法、QSAR 方法等能够对信息核心层中的数字化化学信息



图 1-2 化学信息的结构示意图

进行二次开发利用的计算方法组成。信息表示层处于化学信息学科的最外层，它根据信息核心层的特定要求在计算机信息科学中寻找适合表达化学信息的技术，从多个角度将化学信息以某种直观的形式，如基于计算机的图形、音频、视频等多媒体表示手段向化学工作者展示出来。信息处理层和信息表示层统称为化学信息学科的外层。

三个层次中最重要的层次为信息核心层，它在化学信息学科中处于基础核心地位，并决定了其他两个层次的构成。信息核心层对外层起着决定性作用，外层对信息核心层也能够产生一定的影响。由于受到计算机信息科学技术和仪器设备本身存在某些特点的影响，要求信息核心层中的化学信息必须要按照一定的要求进行组织和编排。例如，连续吸光度曲线在计算机中只能以离散数据点的形式存储，海量的光谱数据、化学化工测量数据或分子结构参数，唯有按一定的数据结构规范化并形成数据库甚至是专家系统，才能方便日后的使用信息处



理层对这些数据进行开发利用。信息处理层在对来自信息核心层的化学信息进行处理之后，所获得的结果一方面将交由信息表示层处理；另一方面，信息处理层将把某些处理结果和原始数据存储在信息核心层，使该层信息量甚至是一些局部结构发生变化。

## 1.6 化学信息的工作方式

化学实践与化学信息之间的关系是“母与子”的关系。化学信息通过“信息采集接口”从化学实践这一母体中获取原始信息，原始信息经过接口的处理后以数字化形式存储到信息核心层中，并通过外层将其重现出来给化学工作者。利用这些被数字化技术处理过的化学信息，化学工作者可以进行更深层次的化学研究实践，从而生产出新的数字化的原始化学信息，这就是化学信息的工作方式，如图 1-3 所示。

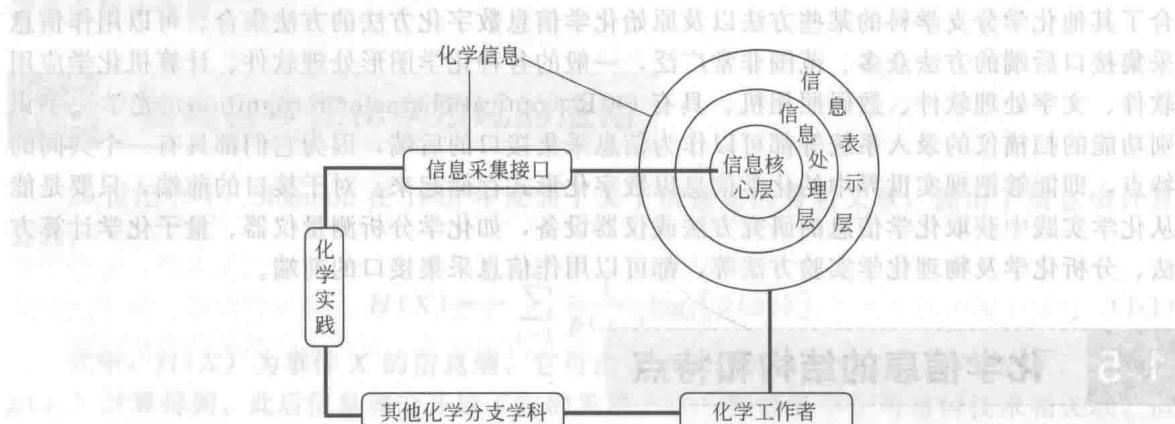


图 1-3 化学信息的工作方式

化学信息的这种工作方式与其他化学分支学科的工作方式基本相同，最大的区别在于化学工作者手中的研究工具是以计算机程序表达的化学计量学、QSAR 等可以对大量化学数据进行二次利用的计算化学方法，其研究的直接对象和得到的产品都是数字化了的化学信息。化学工作者的任务，一是利用现有的计算机软、硬件工具研究大量存储在信息核心层的原始化学数据，找出不同化学变量之间的关系，发现有实际意义的化学规律；二是改进现有的研究方法、开发新的研究手段以更新和完善化学信息外层以及不断丰富和修正信息核心层中的化学数据，为今后开展更深层次的研究工作奠定基础。

## 1.7 化学信息学的应用

### 1.7.1 化合物结构绘制

化学信息普遍存在于化学和计算机的结合之中，几乎每一个化学家都是一个绘制结构者，都会使用到 IsisDraw、ChemBioDraw、JchemPaint 等相关软件去绘制一个分子的二维或三维结构。然而，更深入的问题则需要化学信息学的方法来解决，例如如何将化学的结构有效地储存在计算机里；采用何种格式可以用来在不同类型的化学软件中交换数据等。

### 1.7.2 化学数据库设计与开发

化学数据库的设计、开发、维护和更新是化学信息的重要领域。化学结构和一些相关的

信息主要被存储在化学数据库中，如 Beilstein 数据库。自 1771 年起，化学数据库中存储了超过七百万个有机化合物的信息，可以由 CAS 登录号或分子式查询物质的物理和化学性质，包括光谱数据以及热力学参数。另外在构建这些数据库时，基于物质化学数据信息的数据库查询功能同样至关重要。

### 1.7.3 化学反应体系模拟

人们已经做了很多的努力去实现用计算机预测化学反应的进行，模拟化学反应的发生，来合成一个设计好的目标产物。目前，德国 Gasteiger 的化学信息研究小组，已经有了这样的一个系统，名叫 EROS (elaboration of reactions for organic synthesis)，能够进行包括两种反应物之间的反应结果的预测，或提供采用何种反应物的建议。

### 1.7.4 计算机辅助波谱解析

现代分析化学对于仪器的依赖越来越多，谱学手段发展很快，各种光谱、色谱、质谱以及多种仪器联用已成为了分析化学研究的重要工具，因而，计算机辅助谱图解析也成为了研究人员关注的焦点。一般而言，计算机辅助谱图解析方法可以大致分为两类：一是直接对谱图库进行相似性检索 (library search)；二是采取波谱模拟、模式识别及人工智能的手段进行间接识别。由此涌现出了一大批谱图解析专家系统，如由斯坦福大学组建的用于质谱、<sup>13</sup>C NMR 解析的 DENDRAL 系统；由日本左木慎一等研制的用于质谱、<sup>13</sup>C NMR、<sup>1</sup>H NMR、红外光谱解析的 CHEMICS 系统，等等。

### 1.7.5 化合物结构与活性关系预测

QSAR (quantitative structure-activity relationship) 定量构效关系方法尝试通过对一系列结构相似的药物分子进行分析，找出分子性质参数和生物活性之间的关系，并以此为依据去预测具有药效的新型分子的结构与性质。目前发展到三维的 3D-QSAR 实际上是 QSAR 与计算机分子图形学相结合的研究方法，是研究药物与受体间的相互作用，推测受体的图像及进行药物设计的有力工具。3D-QSAR 研究可分为受体结构已知及受体结构未知两种情况。受体结构已知（目前仅限于酶作为受体）时，可以根据 QSAR 的结果及计算机图形显示受体的三维结构，并随之进行有如“量体裁衣”式的设计。在受体结构未知的情况下（这是绝大多数情况），则可以根据激动剂或（和）拮抗剂的构效关系及计算机图形显示的化合物优势构象，推测受体的结构，然后进行药物设计，亦可以起到“量体裁衣”的作用。

### 1.7.6 实验室信息管理系统

实验室信息管理系统 (laboratory information management system, LIMS) 是将以数据库为核心的信息化技术与实验室管理需求相结合的信息化管理工具。该系统的建立是根据 ISO/IEC 17025: 2005-5-15《检测和校准实验室能力的通用要求》规范，结合网络化技术，将实验室的业务流程和一切资源以及行政管理等以合理方式进行管理。经过 LIMS 管理，可有效提高实验室样品检测效率、提高分析结果可靠性、提高对复杂分析问题的处理能力、协调实验室各类资源以及实现量化管理。经过多年的发展，国内很多大型实验室均认识到了信息化管理的重要性，引入了 LIMS 系统，并进行了不断的改进和升级。

## 1.8 展望

随着信息时代的到来，各类化学信息的相关数据不断涌现，这些数据在使我们获得更多