

# 全国馆藏管理与服务平台 改进

索昂  
2020-11-25

“联合目录建设带来的变化是革命性的。”

全国联合目录中的馆藏建设

索昂

2019-4

## 全国联合目录

- 2010年10月全国图书馆联合编目系统上线试运行，
- 2010年10月对分中心和主要的成员馆发布了馆藏收割的公函，呼吁建立全国联合目录
- 2011年联合编目中心对全国公益性图书馆实行书目数据免费政策，并把提交馆藏作为成员馆的对等义务写在了联编中心的协议书内；
- .....
- 截止6月底，联合编目中心揭示全国图书馆的馆藏量为**27376367**。

## 馆藏收割的难点：

- 目前国内的公共图书馆不到3000家。
- 截至2013年6月21日，全国普通高等学校（不含独立学院）共计2198所，其中独立设置民办普通高等学校424所；全国成人高等学校共计298所，其中民办成人高等学校1所。——百度百科
- 其他图书馆.....
- **国内海量的图书馆和图书室**

## 2.3 系统截图——图书馆管理

书号	书名	作者	来源	状态
00000001	计算机组成原理	唐朔飞	2002-07-07 12:00:00	正常
00000002	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000003	数据库系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000004	计算机网络原理	唐朔飞	2002-07-07 12:00:00	正常
00000005	数据结构与算法	唐朔飞	2002-07-07 12:00:00	正常
00000006	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000007	数据库系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000008	计算机网络原理	唐朔飞	2002-07-07 12:00:00	正常
00000009	数据结构与算法	唐朔飞	2002-07-07 12:00:00	正常
00000010	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常

## 2.3 系统截图——采集进度

已采集书数	书号总数	采集率	采集率	占比 (alpha/书号总数)
2300	14557	35613	14520	98.35%
4358	1448643	36156	113843	7.86%
725	14373	9475	6336	44.16%
6543	402118	982095	357656	36.54%
1709	278098	562132	278674	99.95%
768	18749	21135	18965	87.39%
17791	346511	344628	297253	85.79%
740	251573	20612	112358	44.66%
3853	696027	362084	547777	78.56%
5497	109288	18065	40969	37.46%
7274	326225	1759421	12836	0.4%

Know

## 实体资源馆藏工作的背景

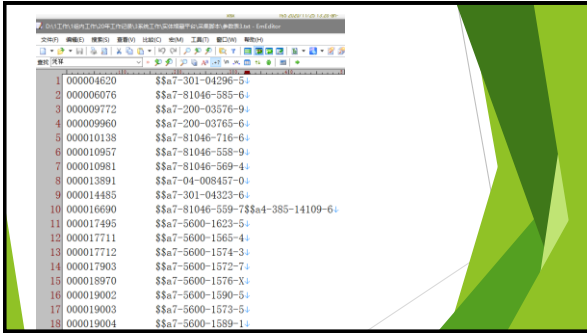
- 待处理馆藏量巨大，存量数据超过1.2亿条书目数据，每年新增约1000万条书目数据；
- 数据质量参差不齐，数据来源复杂。
- 基础查重算法：十四字段查重：010\$a。
- 改进方案：结合八字段查重，对不同的查重结果设置区间，减少查重者需要消耗的精力
- 查重数据量巨大导致传统查重方法无法达到实际效果；
- 查重工作人员的经验复用问题；
- 查重工作流程标准化问题。

## 查重工作界面

书号	书名	作者	来源	查重结果
00000001	计算机组成原理	唐朔飞	2002-07-07 12:00:00	正常
00000002	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000003	数据库系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000004	计算机网络原理	唐朔飞	2002-07-07 12:00:00	正常
00000005	数据结构与算法	唐朔飞	2002-07-07 12:00:00	正常
00000006	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000007	数据库系统原理	唐朔飞	2002-07-07 12:00:00	正常
00000008	计算机网络原理	唐朔飞	2002-07-07 12:00:00	正常
00000009	数据结构与算法	唐朔飞	2002-07-07 12:00:00	正常
00000010	操作系统原理	唐朔飞	2002-07-07 12:00:00	正常



[illegible][illegible]



1	000004620	\$\$a7-301-04296-5
2	000006076	\$\$a7-81046-585-6
3	000009772	\$\$a7-200-03576-9
4	000009960	\$\$a7-200-03765-6
5	000010138	\$\$a7-81046-716-6
6	000010957	\$\$a7-81046-558-9
7	000010981	\$\$a7-81046-569-4
8	000013891	\$\$a7-04-008457-0
9	000014485	\$\$a7-301-04323-6
10	000016690	\$\$a7-81046-559-7\$\$a4-385-14109-6
11	000017495	\$\$a7-5600-1623-5
12	000017711	\$\$a7-5600-1565-4
13	000017712	\$\$a7-5600-1574-3
14	000017903	\$\$a7-5600-1572-7
15	000018970	\$\$a7-5600-1576-X
16	000019002	\$\$a7-5600-1590-5
17	000019003	\$\$a7-5600-1573-5
18	000019004	\$\$a7-5600-1589-1

► 新的采集脚本，本质上是一种细化了数据范围的，提高采集针对性。