

中国语言生活绿皮书
国家语言文字工作委员会发布



中 国

Language Situation in China: 2009

语言生活状况报告

下 编

国家语言资源监测与研究中心 编

ZHONGGUO YUYAN SHENGHUO
ZHUANGKUANG BAOGAO (2009)

2009



商務印書館
THE COMMERCIAL PRESS

中国语言生活绿皮书
国家语言文字工作委员会发布



中国语言生活状况报告 (2009)

下 编

商務印書館
2010年·北京

图书在版编目 (CIP) 数据

中国语言生活状况报告.2009.下编/国家语言资源监测与研究中心编.一北京:商务印书馆,2010.10
(中国语言生活绿皮书)
ISBN 978 - 7 - 100 - 07281 - 6

I . ①中… II . ①国… III . ①社会语言学—研究报告—中国—2009 IV . ①H1

中国版本图书馆 CIP 数据核字(2010)第 129447 号

所有权利保留。

未经许可,不得以任何方式使用。

ZHONGGUO YUYAN SHENGHUO ZHUANGKUANG BAOGAO (2009)

中国语言生活状况报告 (2009)

下 编

国家语言资源监测与研究中心 编

商 务 印 书 馆 出 版

(北京王府井大街36号 邮政编码 100710)

商 务 印 书 馆 发 行

北京瑞古冠中印刷厂印刷

ISBN 978 - 7 - 100 - 07281 - 6

2010 年 10 月第 1 版 开本 787 × 1092 1/16

2010 年 10 月北京第 1 次印刷 印张 34 1/2

定价: 56.00 元



顾策主编 问划编 许嘉璐 赵沁平 郝平 李卫红
教育部语言文字信息管理司
李宇明
审订 陈章太 戴庆厦 陆俭明 邢福义

上 编

主编 周庆生
副主编 郭熙 周洪波
作者 (按音序排列)
白娟 蔡长虹 陈慧 陈章太 程祥徽 崔晓飞
邸宇维 高建平 高洋洋 郭济修 郭熙 平庆
何瑞 和丽峰 黄锦章 江修波 孔江 平阿
李晓华 李旭练 李阳 波传波 江平 道
骆峰 毛力群 倪彦斌 林晓先 魏鹏
汪锋 汪磊 翠王 姚培 刘畅
谢俊英 余桂林 战菊 张奇 魏丹
张映川 赵守辉 战军 张艳
周洪波 周庆生 祝晓宏 邹娟
邹玉华

下 编

主编 王铁琨
第一副主编 侯敏
副主编 杨尔弘 苏新春 何婷婷 赵小兵
作者 (按音序排列)
阿不都热依木·沙力 阿力木·木拉提 艾孜古丽·玉素甫
曹晖 陈敏 陈琪 陈雪 崔乐 高璐
郭曙纶 何婷婷 何伟 陈敏 崔江 格瓦尔·地汗
阚明刚 李安 陈永 宏 李艳欣 威斐
刘佳 刘俊 刘薇 刘林 刘齐 民嘎
祁坤钰 苏小康 苏新春 刘永 刘秋萍 卫新
王华英 王磊 王宁 刘奇 刘王 铁辉
王燕 王宇波 魏励 王继媛 王杨 曾青
于洪志 玉素甫·艾白都拉 王曾 青
张红春 张金爽 张蕾 张勇 曾小兵 張志平
庄晓云 邹煜

目 录

报纸、广播电视台、网络(新闻)用字用语调查	001
调查报告	001
用字总表	032
高频词语表	171
年度新词语调查	322
调查报告	322
新词语表	338
中国媒体年度流行语	415
中文博客专项调查	423
基础教育阶段小学语文教材汉字使用调查	435
调查报告	435
基础教育阶段小学语文教材生字位序表	459
现代维吾尔文网站用词调查	465
调查报告	465
维吾尔文网站高频词干表	484
小学藏语文新课标教材用词调查	498
调查报告	498
小学藏语文课本 500 高频词	511
语言资源监测与研究相关术语(2010 版)	525
图表目录	535
术语索引	540
后记	544

Contents

Survey of the Chinese Words and Expressions in Newspaper, Radio, Television and Internet (News)	001
A Research Report	001
A List of General Words	032
A List of High-frequency Words	171
Survey of New Words of the Year	322
A Research Report	322
A List of New Words	338
Popular Words and Phrases of the Year Used by Chinese Media	415
A Special Survey of Words and Expressions Used in Chinese Blogs	423
Survey of the Use of the Chinese Characters in the Chinese Textbooks (New Curriculum Standards) for Elementary School	435
A Research Report	435
A List of the New Characters in the Chinese Textbooks (New Curriculum Standards) for Elementary School by the Order	459
Survey of the Words Used in Modern Uighur Websites	465
A Research Report	465
A List of High-frequency Stems in Modern Uighur Websites	484
Survey of the Words Used in the Tibetan Textbooks (New Curriculum Standards) for Elementary School	498
A Research Report	498
A List of the 500 High-frequency Words in the Tibetan Textbooks (New Curriculum Standards) for Elementary School	500

Contents

Curriculum Standards) for Elementary School	511
Technical Terms Related to the Monitoring and Research on the National Language Resources (2010 Edition)	525
Index of Figures and Tables	535
Index of Technical Terms	540
Postscript	544

报纸、广播电视台、网络(新闻)用字用语调查

调 查 报 告

国家语言资源监测与研究中心利用国家语言资源监测语料库(包括平面媒体、有声媒体、网络媒体)的年度语料,已经连续发布了2005年度、2006年度、2007年度、2008年度的语言生活状况报告。其中,报纸、广播电视台、网络(新闻)的用字用语调查是每年都进行的调查项目。报纸、广播电视台、网络(新闻)的年度用字用语状况,可以反映媒体年度的语言使用实态,也可以透过这些字词语的使用状况看到年度的社会热点、重大事件等。2009年度的用字用语调查是在国家语言资源监测语料库2009年度的语料上进行的。

历时五年的语料积累,既可以用来反映五年共时的语言生活,也可以通过对比分析来反映五年历时的语言变化。在本报告的第四部分,对五年的字词语调查数据进行了比较,以期从共时、历时两个角度对五年的字词语使用状况进行分析。

一 调查使用的语料及调查内容

(一) 调查使用的语料

本年度的调查语料涵盖平面媒体、有声媒体、网络媒体三种,共计1 249 387个文本文件,1 237 492 014字符次(包括标点、符号及西文字符、数字等出现的次数),其中汉字出现1 007 019 960字次。

2009年度语料采集的依据及选择过程与往年一致。在进行此项调查时,为了使年度间的数据规模基本保持在同一个数量级上,总字符次保持在12亿左右,总汉字次保持在10亿左右。平面、有声、网络媒体的语料量仍按5:1:4的比例进行选取。

1. 报纸

平面媒体选择了 2009 年度 15 种报纸作为调查语料,选择时综合考虑了“发行量、发行地域、发行周期、媒体价值”等因素,同时考虑了语料的可获得性因素。发行量参考了 2008 年 6 月 2 日召开的第 61 届世界报业大会(瑞典·哥德堡)发布的“2008 年世界日报发行量前 100 名排行榜”(中国部分);媒体价值参考了由 2008 世界品牌大会(中国·北京)于 2008 年 6 月 2 日发布的“2008 年《中国 500 最具价值品牌》排行榜”。

选定的 15 种报纸是(按音序排列):《北京青年报》《北京日报》《北京晚报》《法制日报》《光明日报》《广州日报》《华西都市报》《今晚报》《南方周末》《钱江晚报》《人民日报》《深圳特区报》《羊城晚报》《扬子晚报》《中国青年报》。

报纸语料共计 795 046 个文本,633 705 267 字符次,其中汉字出现 511 924 302 字次。

2. 广播电视

广播、电视语料是根据播出的录音或录像转写的文本,选取的主要依据是节目收视率,并且综合考虑了传播媒介(广播、电视)、媒体级别(中央、地方)、传播广度(是否上星)、播出时间(是否黄金时段)、节目样态(独白、对话、综合)、文本现存(是否有转写好的文本)等因素。

2009 年度选取的广播、电视语料包括:中央电视台、北京电视台、上海文广新闻传媒集团(电视)、天津电视台、重庆电视台、广州电视台、山东电视台、山西电视台、安徽电视台、河南电视台、哈尔滨电视台、深圳电视台、石家庄电视台、洛阳电视台等 14 家电视台,以及中央人民广播电台、北京人民广播电台、上海文广新闻传媒集团(广播)、天津人民广播电台、重庆人民广播电台、山东人民广播电台、深圳人民广播电台、石家庄人民广播电台、洛阳人民广播电台、包头人民广播电台等 10 家广播电台,总计 169 个栏目,20 219 个文本。

广播、电视语料量为 124 884 433 字符次,其中汉字出现 102 707 631 字次。

3. 网络(新闻)

根据年度调查所确定的三种媒体语料量的比例,网络媒体语料只选取了新浪、腾讯两个网站的部分新闻语料。选取的方式是,对已采集的两个网站的全部语料分别以两个网站每一天的所有语料为单位,在一天的语料中按既定比例随机抽取文本,然后将以天为单位抽取得到的所有文本集合在一起,形成网络(新闻)的样本语料。



由此获得的语料共计 434 122 个文本, 478 902 314 字符次, 其中汉字出现 392 388 027 字次。

4. 语料说明

报纸语料是网络版的。广播、电视语料是由广播、电视节目转写的文本, 与原始有声语料之间存在些许差异。网络(新闻)语料来自新浪、腾讯 2009 年度的新闻页面, 这些语料全部从网络下载而得。用于调查的语料是纯文本的格式, 是由相应的页面文件转换而成的, 转换过程中去除 HTML 标签信息、广告信息, 报纸语料中去除转版信息、图片标题等短文本。

(二) 调查内容

本次调查的对象是汉字和词语。词语是计算机自动分词产生的分词单位, 既包括语文词, 也包括专名(人名、地名、组织机构名、其他专名)、时间表达式(如“2009 年 11 月、5 月、15 点 30 分”)以及结合紧密、使用稳定的短语(如“经济危机、世博会”)。调查项目主要有频次、频率、累加频率、出现文本数、使用率、累加使用率等, 调查内容既包括 2009 年度内三种媒体语言的共时考察, 也含有 2005—2009 五个年度间的历时考察。其中, 频次、频率、累加频率、出现文本数和使用率的含义及计算方法与 2005—2008 年度的《报纸、广播、电视、网络(新闻)用字用语调查》^①相同, 也可参考《中国语言生活状况报告(2008)》下编的“语言资源监测与研究相关术语”。^②

二 汉字使用情况

(一) 说明

1. 本次统计没有甄别文本中的别字。
2. 本次统计不包括汉字部件、乱码、无法显示的字符等。

^① 见国家语言资源监测与研究中心编《中国语言生活状况报告(2006)》下编第 2 页, 商务印书馆 2007 年版。

^② 见国家语言资源监测与研究中心编《中国语言生活状况报告(2008)》下编第 493—501 页, 商务印书馆 2009 年版。

(二) 基本情况

1. 字符总数：指全部语料中汉字、标点、符号等的总量，计 1 237 492 014 字符次。
2. 汉字总数：指全部语料中汉字出现的总字次，计 1 007 019 960 字次。
3. 字种数：指字形不同的汉字种数，计 10 204 个。
4. 共用字种数：指报纸、广播、电视、网络三种媒体中都出现的汉字，计 6 238 个。
5. 独用字种数：指只在报纸、广播、电视、网络某一种媒体中出现的汉字，计 2 359 个。

汉字使用情况的具体数据见表 1-1。三种媒体用字之间的关系除了表中列出的共用、独用字种之外，还有部分共用字种的情况。部分共用字种是指只出现在任意两种媒体中，而没有出现在第三种媒体中的汉字种数，比如报纸与广播电视台的部分共用字种计 165 个，报纸与网络（新闻）的部分共用字种计 1 357 个，这两部分的并集形成了报纸的部分共用字种，共计 1 522 个。为了使数据更加简洁、直观，表中没有列出部分共用字种数。

本次调查的全部语料用字（共计 10 204 个汉字）形成 2009 年度用字总表^①（简称“用字总表”），表中包括了汉字在语料中出现的频次和文本数两个最基础的数据，所有汉字根据使用频率由大到小排序。

表 1-1 2009 年度汉字使用情况

媒体	总字次	字种数	共用字种数	独用字种数
报纸	511 924 302	9 317	6 238	1 557
广播、电视	102 707 631	6 597		109
网络（新闻）	392 388 027	8 373		693
全部语料	1 007 019 960	10 204		2 359

(三) 汉字的覆盖率

汉字的覆盖率是调查的重要项目之一。它是反映汉字常用与否的重要指标，同时反映了在整个调查语料中汉字使用的分布情况。其统计结果见表 1-2。

^① 见国家语言资源监测与研究中心编《中国语言生活状况报告（2009）》下编第 32—170 页，商务印书馆 2010 年版。



表 1-2 2009 年度汉字对语料的覆盖情况

语料	达到 80%		达到 90%		达到 99%		达到 100%
	字种数	字种比例(%)	字种数	字种比例(%)	字种数	字种比例(%)	字种数
报纸	609	6.54	982	10.54	2 460	26.40	9 317
广播、电视	546	8.28	893	13.54	2 298	34.83	6 597
网络(新闻)	584	6.97	943	11.26	2 294	27.40	8 373
全部语料	602	5.90	970	9.51	2 400	23.52	10 204

(四) 与现行规范字表的比较

1. 用字总表前 2 500 字与一级常用字比较

用字总表前 2 500 字与《现代汉语常用字表》^①一级常用字(2 500 字)比较,用字总表中有 342 字是一级常用字中所没有的。将用字总表按照前 500 字、501 至 1 500 字、1 501 至 2 500 字分为三段,每一段中没有出现在一级常用字中的汉字列于表 1-3。

表 1-3 用字总表前 2 500 字与一级常用字比较

范围	一级常用字之外的字
前 500 字	尔(1个)
501—1 500 字	媒 钊 频 韩 伊 诺 措 迪 辑 综 俄 萨 澳 伦 曼 菲 姆 莱 洛 郭 聘 署 杭 蒂 姚 账 谓 咨 凌 拓 卢 屏 娜 埃 拟 邓 沪 弗 艾 鹏 琳 蔡 邦 贾 冯 浦 翔 斌 谐 逊(50个)
1 501—2 500 字	兹 帕 曹 穆 敦 胎 赫 潘 秦 肖 募 玲 曝 蒋 徽 聊 莞 霍 彭 粤 颁 颇 氚 韦 墅 兑 戈 奈 涵 许 莉 吕 袁 枚 颖 鸿 卦 履 玛 魏 讼 崔 涯 谭 癌 硕 卓 侯 廷 铭 郁 吁 肇 憾 匪 寓 耶 砸 怡 旭 娅 淮 茨 厢 蒋 晖 碳 峻 汶 磊 辖 庞 抑 妮 妆 莎 契 鑑 烳 瓷 魅 舜 柯 苑 蕾 岳 晰 雯 铝 湘 娟 琼 遷 逸 嘛 鼎 逾 汶 薇 歧 坤 吴 骏 邱 婷 滞 芯 赋 淑 杉 绛 邵 蓉 彰 弥 罕 雇 尹 幽 琪 薛 勃 坠 彦 弘 轴 淀 衷 拦 仲 巢 俟 奢 肇 挫 瘪 溢 萧 楠 鮑 浏 货 擅 尬 尷 瞬 飘 豫 莲 亨 摆 彬 坪 宠 渝 卿 殷 廖 莫 肇 澄 爵 澜 馨 珊 佐 瘤 坎 邹 勘 擎 祭 裸 轩 熙 玮 啭 犀 晔 淦 霆 奎 啥 呵 阴 蔓 靓 槛 莢 遂 哮 撼 腺 馈 贱 倩 坡 凸 逛 利 俞 眇 龚 焉 媛 奕 殿 蕴 攻 蔚 吻 喻 咖 窦 瑶 莹 侷 伽 郝 芭 缅 谍 衍 辐 骚 鹤 腕 啡 昔 瑰 扳 瞳 甸 瑟 仕 靖 邢 隧 麟 迄 饮 崩 函 贬 倪 玄 惟 埔 噪 晤 溃 闫 遏 娅 崩 粹 垦 秉 媳 吾 励 寂 懈 碟 睹 眇 腻 缪 枢 湛 藤 聚 蟹 炫 皓 夭 犀 潘 讶 咯 硅 奢 炜 帷 荔 儒(291个)

① 国家语言文字工作委员会、国家教育委员会 1988 年联合发布。

3. 用字总表前 7 000 字与《现代汉语通用字表》^①比较

用字总表前 7 000 字与《现代汉语通用字表》(7 000 字)比较,用字总表中有 725 字是《现代汉语通用字表》中所没有的。具体情况见表 1-5。

表 1-5 用字总表前 7 000 字与《现代汉语通用字表》比较

范围	《现代汉语通用字表》之外的汉字
前 3 000 字	𠂔 咎(2 个)
3 001—5 000 字	森 基 魏 琥 囤 漢 婦 玖 塑 犁 瘴 舛 睹 眇 珑 吋 嵩 钜 昇 岚 焰 硏 疋 祐 噫 全 孽 鮑 珪 喇 瞳 吮 弩 握 後 坻 那 劍 岬 朮 塚 豔 樑 嗜 曝 戎 斧(54 个)
5 001—7 000 字	鍾 嘶 琪 德 嫌 迳 璞 眇 紊 炅 狃 峯 帕 谈 駁 逃 甦 這 犀 鍼 啓 墉 啖 墾 有 積 炙 俟 佐 個 芬 露 寞 湿 還 會 勢 崑 旣 呀 喃 點 炕 篓 姮 沒 撐 魏 莎 枢 媚 妖 賚 瑶 素 塑 蹤 陞 埃 黃 驳 內 負 畔 查 栢 滴 漏 裸 郎 線 誌 間 琦 詠 激 洨 酒 艄 艨 壓 娛 現 檻 辻 謝 別 時 紮 棋 褚 係 順 増 嫩 罢 蘭 高 亾 紛 虬 遷 嶺 佢 鶴 始 尔 唉 嬉 脖 洩 為 呪 搨 涂 牀 盤 福 國 锺 踏 關 調 呴 間 芮 價 訴 依 嘴 繢 臺 師 媒 話 鮑 呼 炣 啟 啟 鵠 將 迴 嶺 望 𩫱 過 祇 買 咎 紿 跪 驪 𩫱 桀 趟 繫 鵠 坐 遁 鑊 觀 繢 蠱 刈 侷 板 菴 瑰 徹 珪 進 萬 鑑 眇 煙 煙 業 沢 啓 麼 淬 煥 鵬 丘 漲 埠 敏 艸 萌 漾 長 穗 說 機 鮑 捏 菴 題 煙 穀 場 愛 葵 媚 醜 經 夭 宮 性 抛 酪 可 啟 乘 覺 磡 佺 兒 莖 開 岡 裡 鮑 委 偶 吖 嘘 強 乘 脣 艸 眇 資 舜 樸 請 兩 態 峯 宦 幾 蘭 诉 彦 該 肩 較 𠂔 燥 勳 岌 讀 珍 聲 菴 菴 邶 變 發 滯 傢 對 見 厢 嘴 初 漢 肅 稔 𠂔 保 墉 罷 晓 虹 崑 梆 區 倘 濟 牝 彪 嘉 媛 韶 鵠 製 佢 沙 繩 並 鑄 嵩 穎 奥 動 久 柯 鏡 裝 熙 涅 藥 焱 佈 哲 混 鶲 圖 蛙 格 預 黑 球 錄 滾 芮 頭 濡 桃 明 芮 漏 皚 頭 漏 賦 煙 滅 峴 晴 嵩 東 啟 蒼 勸 嵩 咨 韶 達 楷 漢 鄭 覽 數 汶 鍾 橫 鶲 曾 脩 說 那 嫣 嶠 涠 姬 閔 睽 滯 塊 奴 猥 龍 覩 翩 布 鈴 汝 奕 篓 遊 結 雲 沖 呀 嵩 莠 格 姆 昊 莎 沖 助 婉 呒 篓 繢 貲 廵 遷 莊 兮 眇 錄 虞 蔽 繢 勦 驟 塙 鬱 謂 體 崑 兒 岑 莩 紅 口 紅 天 火 柿 坪 出 篓 繢 扇 鴛 妹 議 𠂔 植 忒 鑿 梵 鰐 痘 損 讀 罩 露 物 楷 詞 險 歷 陌 𠂔 瑞 語 銀 沢 閻 蟲 達 鑿 無 億 繢 馬 邸 蘆 記 傑 妻 車 止 從 国 球 錄 𠂔 犹 柏 墉 認 書 豐 產 鰐 穀 苟 濡 紅 紅 削 壞 堙 堙 堙 許 許 許 許 許 許 許 跨 餘 眇 錄 彪 樂 獎 俱 楷 甸 鴨 樂 學 瞭 確 鰐 肅 週 眇 篮 箱 稅 𠂔 蔻 髮 參 豐 岡 嶠 華 卉 實 增 勝 驅 準 塘 銀 乾 佛 楷 勝 𠂔 許 許 許 許 連 鏡 潤 篓 呕 廾 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 檻 貞 簾 印 級 柏 始 抹 瀑 眇 廪 陸 儒 凱 利 崙 豊 慶 傕 體 陽 窓 晶 淬 決 採 纔 勇 賣 增 隻 抠 疣 電 崴 貌 嘴 兩 隨 寶 壞 優 繢 聽 韶 丟 鈿 呕 倉 煙 門 總 潢 滴 鏰 難 戶 穂 絶 銘 傕 鈴 汐 張 淩 純 禅 淬 鎔(669 个)

① 国家语言文字工作委员会、中华人民共和国新闻出版署 1988 年联合发布。

在历年的年度用字总表前 7 000 字与《现代汉语通用字表》的比较中,2009 年度用字与《现代汉语通用字表》的差别是最大的。^① 探究其中的原因,可归为如下几个方面:

(1) 年度语言生活由于汉字的“繁简之争”“《通用规范汉字表》(征求意见稿)公开征求意见”而受到社会的普遍关注。众多语言文字工作者和广大民众对这些问题提出了自己的建议,在媒体中有许多讨论,其中不乏对文字规范工作计深虑远的建言献策,这使 2009 年度的用字总数增加,繁体字、生僻字增加。

(2) 网络语言中汉字的使用标新立异,且传播速度快。2008年度“囧、囧、囧”等汉字在网络上的流行,2009年度众多生僻字也在网络上被发掘,尤其一些字形比较奇特的合体汉字,如“熯、鼴、亞、王、𠂇、曇、𦵹、𡿆、𡷁、𡷃、𡷄、𡷅、𡷈、𡷉”。在网络世界里,这些汉字的本义已经不重要了,字形是其在现代社会表义的重要依据。如“夭”不再表示“天”的本义,而是“王八”义;“𠂇”不再表示“光”,而是“火化”义。但有些仍保持原义,如“𡷃”表示“不要”义,字音则由“不”的声母与“要”的韵母组成。此外,由于汉字“囧”在网上的流行,似乎在追随这样一种汉字结构带来的效果,一些由两个或三个相同的汉字成字部件构成的汉字在网络中也悄然流行开来,如“甡、囍、孖、屾、皕、垚、骉、犇、囍、皛、燚”等。这些现象值得进一步关注与探究。

这 725 个汉字包括：不规范的简化字（“礶、钖、㚱、鉅、桠、鍾、迳”）7 个、繁体字（如“噃、後、礮、襲、這、會、點”等）211 个、异体字（如“陞、迺、埶、詰”等）113 个、旧印刷字形（如“暨、苟、弑、茲、稅、沒、別”等）22 个、日本汉字（如“禪、畊、沢、麼、𠙴”等）21 个、旧计量字（“浬、吋、呎”）3 个。

4 《现代汉语通用字表》与用字总表比较

用字总表中未出现的《现代汉语通用字表》中的字有 162 个，具体如下：

锿 摧 鞍 鞍 璧 璋 鮕 汗 泠 潺 涣 嵴 疎 疏 疏 疏 疏 疏 疏 疏 疏 疏 疏
剗 爰 銚
橫 戈 蓋 淦 賦 滬 翳 到 倦 脍 溥 酝 銛 羦 敘 溢 肪 蔡 莺 穎 紹 振 憂 蔽
瞓 蔽 級 縱 僭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭 脭
酸 鮨 僇 蠻 蟻 翫 漵 漵 肿 鴻 蛸 噴 踰 踰 踰 踰 踰 踰 踰 踰 踰 踰 踰
阡 依 绦 塘 婁 淞 痴 芮 梅 煩 煩 煩 煩 煩 煩 煩 煩 煩 煩 煩 煩

^① 见国家语言资源监测与研究中心编《中国语言生活状况报告(2009)》下编第30页,商务印书馆2010年版。

繁 踏 庠 淚 镊 鹅 痰 豉 脍 脍 蹤 槛

(五) 汉字使用的其他情况

用字总表中的繁体字、异体字、不规范的类推简化字、旧计量单位用字、日本汉字等,出现在用字总表的 2 874 位以后,即覆盖率达到 99.51%之后,与去年的 3 118 位相比有所提前。2009 年度汉字的使用情况相对来说比较活跃,与该年度大众积极关注并参与语言文字工作有较大的关系。具体情况见表 1-6。

表 1-6 汉字使用的其他情况统计

类型	全部	报纸	广播、电视	网络(新闻)
旧印刷字形	86	65	9	65
繁体字	1 217	769	234	762
不规范的简化字	17	16	6	10
旧计量用字	4	3	1	3
异体字	396	335	64	241
日本汉字	84	61	13	60

(六) 报纸、广播、电视、网络(新闻)语料与全部语料的汉字频率比值

计算报纸、广播、电视、网络(新闻)各媒体前 2 500 个年度高频字与全部语料中相应汉字的频率比值^①,并按照频率比值的降序排列比较结果,可以观察三种媒体语料的用字特点。表 1-7 分别列出了报纸、广播、电视、网络(新闻)三种媒体语料中频率比值排在前 20 位的汉字。

表 1-7 报纸、广播、电视、网络(新闻)汉字频率比值分析

媒体	前 2 500 字中频率比值在前的 20 个汉字
报纸	墅 龄 粤 腺 阜 荔 践 禽 塘 废 瘤 绣 穗 尿 晴 肠 瞳 厨 胃 奴
广播、电视	伽 窝 咱 嘛 呢 扁 您 殖 饲 啊 邢 呀 它 么 旱 猪 吗 怎 你 徽
网络(新闻)	卦 浪 罢 描 页 寸 蒂 娱 踢 篮 弗 浏 球 芯 杆 霆 扳 兹 玮 棋

1. 从表 1-7 可以看出,广播、电视语料的口语化特征最为显著,记录语气词、称谓词的汉字居多;报纸语料中“粤、圳、禺”等显示了地域性报纸的特点;网络媒

^① 见国家语言资源监测与研究中心编《中国语言生活状况报告(2008)》下编第 498 页,商务印书馆 2009 年版。

体中的“浪、页、浏”等表明了网络媒体语料采集来源及形式的特点。

2. 媒体的用字特色也体现了内容上的不同侧重点, 报纸语料的“墅、幢、奴”等与社会生活中的房地产有关; 网络媒体上“娱、踢、篮、球、兹、棋”等则反映了网络媒体关注文体娱乐的特质。这些汉字在用字总表以及各媒体用字表中的频率、排序情况见附表1。

三 词语使用情况

本次调查仍然使用中国科学院自动化研究所研制的分词标注系统。

(一) 基本情况

1. 分词单位总数: 指由分词软件对语料切分得到的字符串的总数, 计717 321 946次。其中标点、符号等出现123 795 923次, 其他分词单位出现593 526 023次。

2. 总词语数: 即不包含标点、符号、纯西文、纯阿拉伯数字、数字与西文混等形式、网址等的分词单位, 共计592 414 821词次。

3. 词种数: 2 348 100个。

4. 共用词种数: 193 416个。共用词种是指报纸、广播、电视、网络(新闻)三种媒体语料库都用到的词语。基本数据见表1-8。与汉字的调查相同, 这里只列出了媒体间的独用、共用情况, 未列部分共用的情况。

表1-8 2009年度词语使用情况

媒体	总词语数	词种数	共用		独用	
			词种数	比例 (%)	词种数	比例 (%)
报纸	298 944 198	1 579 054	193 416	12.25	1 049 055	66.44
广播	60 964 722	432 014		44.77	165 638	38.34
电视	232 505 901	1 084 828		17.83	579 032	53.38
网络(新闻)	592 414 821	2 348 100		8.24	1 793 725	76.39

(二) 词语的覆盖率

表1-9列出了覆盖率为10%到90%和91%到100%各段的词种情况。从表中可以看出, 词种数的明显上升是在覆盖率为90%之后, 覆盖率99%到100%

