

汉语中介语动态追踪有声数据库建设的基本设想

袁丹 吴勇毅^①

[摘要] 本文介绍了汉语中介语动态追踪有声数据库不同于以往中介语语料库的特点,并讨论了数据库的建设流程,包括字表与词表的设计、信息库的建立、语音数据收集、语音标注、检索条件,重点在字表与词表构建中,强调语料库建设要服务于留学生汉语水平测试,并为开发留学生语音测试软件用于语音诊断及语音矫正提供研究基础。

[关键词] 汉语中介语;有声数据库;语音标注

The Innovation of Chinese Inter-language Dynamic Tracing Audio Database

Yuan Dan Wu Yongyi

[Abstract] This paper discusses the construction process of the Chinese inter-language dynamic tracing audio database, which with different characteristics from the previous inter-language database, includes word list designing, information database building, sound data collecting, sound labeling, search conditions, etc. The emphasis is on the construction of the word list and the word list, emphasizing that the construction of the corpus should serve the purpose of the Chinese proficiency test for international students, and provide a research basis for developing pronunciation test software for speech diagnosis and speech correction.

pronunciation test software for speech diagnosis and speech correction.

[Key words] Chinese inter-language; dynamic tracing; audio database

1 引言

国内对语音数据库建立最为关注的要数民族语言学界,赵尔琨等(1992)就已提出了建立藏语拉萨话语音声学参数数据库的设想,而后蒙古话(冯晓等,1997)、安多藏语(于洪志等,2007)、哈萨克语(孙致吉等,2008)等少数民族语言的声学数据库项目也相继

^① 略作简介,袁为吉林延边大学对外汉语学院教师,研究方向为语音学、方言学、社会语言学。吴勇毅为延边大学对外汉语学院教授、副院长,研究方向为语言学理论与应用语言学、第二语言习得、对外汉语教学。

本量的采集上,我们将采取“多多益善”的原则,即对不同母语背景的样本量的均衡性不做过

学习汉语(或入学)之初,就给他建档,动态采集其在不同学习阶段的语料,以达到动态追踪的目的。实验研究表明,二语习得者在语言学习之初要掌握母语音系中陌生的音位对立会比较困难,但随着学习时间的增加,这种困难会得到改进,甚至完全克服。如Liu和Jongman(2012)对美国学生习得汉语[ts]和[tsh]的研究结果表明:初等一级的美国学生能较好地掌握[ts]和[tsh]的对立,而初等二级的学生则两者都掌握得不好,不能掌握它们(Gambler vs grammy)的对立,而初等三级的学生则两者都能掌握。因此,对学生的语音习得过程进行动态追踪,可以发现哪些是学习者在习得过程中易于纠正的错误,哪些是学习者容易“化石化”难以改变的错误。另外为每个学习者建立一个语音数据库,可以对每个学习者进行追踪诊断,每隔一段时间即为学生进行测试,出具诊断报告,并给出建设性的纠正方案。

2.3 有声性

不同于以往的中介语口语语料库,本数据库十分注重录音样本的采集方式。以往的中介语口语语料库虽然也强调了有声性,但是对录音质量的控制却并不十分严格,大多采用录音笔录音。当然这和语料库的建库目的相关。以往的语料库一般用来做词汇、句法、语篇的

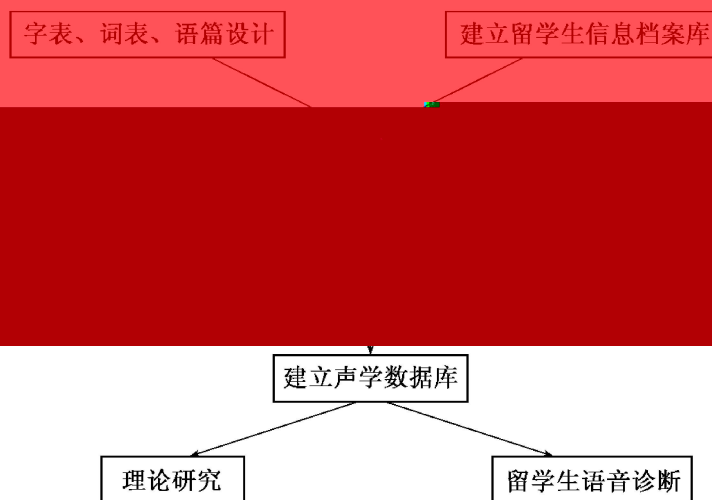
使用头戴式指向性话筒 AKG C520。

2.4 开放性和共享性

一个庞大的数据库,需要有庞大的数据样本量作为支撑,而数据库的开放性则是数据库得以不断更新和保证学术数据库质量和学术价值的关键。

“教育事业的进步,依赖于教育资源的共享,而实现最充分的资源共享是达此目的的首选”。实现共享性的前提是开放,如果经过充分论证,学术数据库,特别是教育数据库,建设者应该主动开放,开放有人可以在线直接浏览下载语音数据,真正实现数据共享(当然开放也可以无条件或有条件为本数据库提供符合要求的语音资料)。

来的录音材料都需要进行标注,然后进行声学分析,建立声学参数库,最后将这些数据用于理论研究以及应用于留学生语音诊断(具体流程参见下图)。以下我们将分步骤详细阐述建库流程。



3.1 字表、词表、语篇的设计

字表理论上要包括普通话声、韵、调配合的所有音节,但需要排除某些音节可能只有生僻字的音节,扣除后再经过挑选,确定普通话 1 000 个音节为数据库的字表。考虑到初级学生在辨识汉字上还有很大的困难,1 000 个音节只列出拼音不列汉字。

能根据经验来为那些三拼音节首字母设计两个字母的网址,如工工文调、这飞空音和个这飞空音的网址为 gggw 或 ggw 或 ggww 或 ggwww 或 ggwwww 或 ggwwwww。

初级学生只给拼音不给汉字。新编部分选录常用词要尽量做到声、韵、调搭配的全面性,可考虑设计多个较短小的语篇让学生复述。拼读部分可确定五个语篇,如,个人爱好、个人经历、家庭情况等让学生自由发挥讲述。时长为 10~15 分钟。朗读部分的语篇可以参考学

分的语篇可以看出每个学生个体的语音问题。

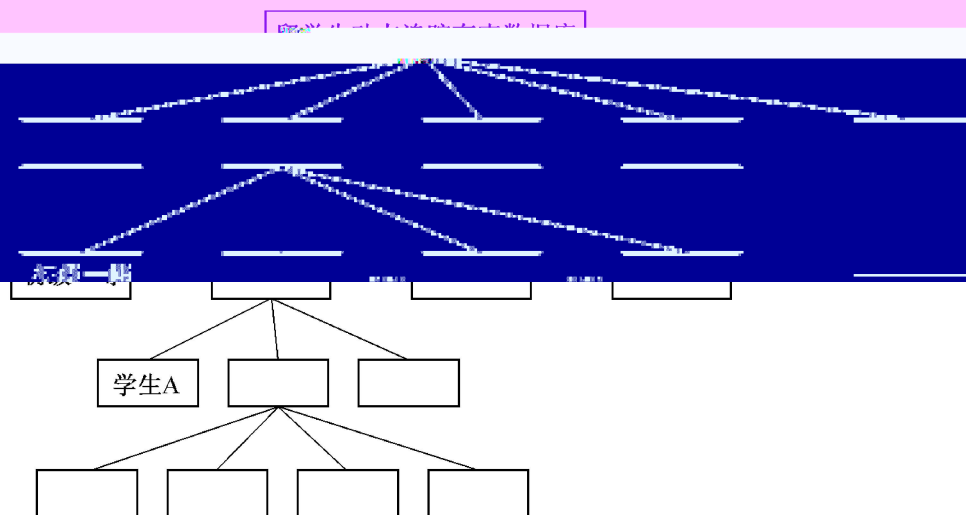
3.2 建立留学生信息档案库

要对留学生的语言习得情况进行研究或诊断,首先要对每一个留学生的档案建立信息库。档案库的建立,首先要按照不同母语背景建立分库,然后再按照不同等级建立次分库,最后再为每一个留学生建立一个信息库。

3.3 语音采集

为了达到动态追踪的目的,每学期分三次对留学生进行语音采集工作:入学考试一次、期中一次、期末一次。零起点的学生入学初不做采集,在入学一个月后做一次采集。事实

行采集外,我们还将进行母语者的语音采集,选择汉语普通话标准的15男、15女录音,作为和留学生语音偏误进行比较的参照库。数据库将首先按照不同母语背景建立分库,然后再按照不同等级建立分库,为每个留学生建立一个数据库,数据库里包括字库、词库、语篇库以及个人信息库,以便于日后的检索。具体如下图。



3.4 语音标注

采集到的所有录音都需要进行初次标注和二次标注。初次标注包括两层:第一层是音节标注,也就是切分出音节来;第二层是声韵母标注,也即将所有音节都标注出声母段和韵

用人工来进行测算,花费的时间是不可想象的。在 Praat 中做完初次标注后,就可以运用其

据体测第一层标注是时点标注,也就是标注一些重要的时点,如元音的起始和终止时点等。第二层标注是声韵母标注,标注出了声韵母,如元音标注为[pa]和[p^ha],声母标注为[p]和[p^h]。第三层标注是时点标注,标注出了起点(start)、中点(mid)和终点(end)。

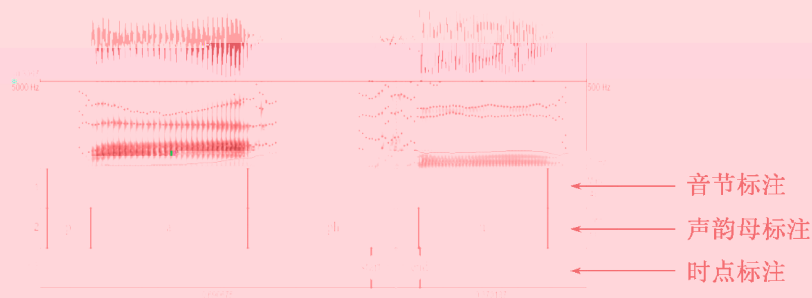


图 3.5 音节的初次标注和二次标注

3.5 语音分析

大规模的录音采样完成后,就可以对一些留学生汉语中介语的语音偏误项目进行声学测量参数的设定,建立留学生二语语音声学参数库。什么样的音素应提取什么样的声学参数,一般来说是有针对性的。元音一般提取第一共振峰(F1)和第二共振峰(F2)这两个声学参数,考察二语学习者的元音声学空间和母语者的元音声学空间的差异;送气声母可以提取送气声母的时长这个声学参数。如果发音偏误项目已经标注好,那么除了测量发音的

后,为语音偏误诊断工作的基础。

3.6 语音诊断

建立起不同母语背景者汉语中介语语音声学参数数据库后,就可以对不同母语背景学习者在习得汉语时的语音偏误进行研究和分析,列出偏误特点,开发留学生汉语语音测试软件,为留学生进行语音诊断,并且进行语音矫正。目前来看,还没有一个留学生汉语语音测试软件被应用于留学生汉语语音诊断。目前只在网络上开发了语音测试工具,开发了软件,可

语音识别并非用来判断发音正确与否,而是尽可能地排除个人口音、环境噪音等因素对识别结果的影响,所以对语音失误的宽容度比较高,因此只能作为判断语音的辅助工具,而不能

作为一个语音测试的软件。

3.7 理论探究

留学生汉语中介语有声数据库的建立可以为当前二语习得研究的理论探索提供大数据支持。本研究将在二语习得研究中对接以下理论。

(1981、1988、1991、1992)的SLM理论系统阐述了二语语音习得感知和产出中的母语迁移,建立留学生汉语中介语有声数据库,分不同母语背景对留学生汉语的语音习得偏误进行考察,可以为这两个理论提供汉语的例子,进一步完善这两个理论。虽然学者们普遍承认母语迁移在二语语音习得中的重要地位,但是也有一些学者指出,母语者和非母语者在发音上产生的差异并不完全是由于母语者的母语迁移,如Garcia and Herbert(1979)的论文《Some Phonological Errors in Second Language Learning》,在这以前已有其他的学者也提出过类似的观点(如Briere 1966,1968,Tarone 1978,Wode 1977,1978)^①。Bohn(1995)的研究更进一步指出,不仅在二语的言语产出中母语迁移不能解释所有的问题,在跨语言的言语感知中母语迁移也解释不能解释所有的问题。Bohn(1995)对德语学习者、中文学习者以及西班牙语学习者感知英语的[e]-[æ],[i]-[ɪ]进行了实验研究,英语的[e]-[æ],[i]-[ɪ]具有复杂的声学线索,既有元音素质的差异,也有时长的差异。实验前的假设为:德语只有[e],可以和英语的[e]-[æ]形成对立,且有时长对立,预测德语母语者会以元音音质差异和时长差异来区分英语的[e]和[æ];西班牙语有[i],而英语有[i]-[ɪ],没有时长差异,预测西班牙学习者仅以音质差异来进行区分,忽略时长的作用;汉语有[i],而英语有[i]-[ɪ],没有时长对立,但是时长是四个声调的辨别特征,预测母语为汉语的学习者和西班牙语学习者一样,以音质差异来区分[i]和[ɪ],时长的作用很小。实验结果表明,母语为母语者主要依靠元音音质的差异来区分[e]和[æ],时长的作用很小;母语为母语的学习者主要依靠时长来区分[i]和[ɪ],音质差异的作用很小;母语为汉语的学习者主要依靠时长来区分[i]和[ɪ],音质差异的作用很小。这些结果支持了“去母语”假说,即学习者音质的差异不能够满足区分母语对立的要求时,时长差异就会用来区分非母语的元音对立,不同母语的二语学习者会表现出共性。因此,通过对不同母语背景者发音音素增音效应的观察,除了证实不同母语背景者发音习得的差异性外,也可以考察产出和感知中的共性问题。

4 结 语

本研究探讨了建立汉语中介语动态追踪有声数据库建设的基本设想,归纳了大数据支持

的大数据支持,建立具有深刻的理论意义和实用价值,能为留学生汉语语音教学、语音测试、语音诊断以及语音矫正等提供大数据支撑。我们认为,通过将留学生的语音采集纳入到留学生入学测试、期中以及期末测试中这样的方式,大规模的数据采样完全可以做到持续而有序的进行。

① 转引自Bohn(1995)。

参考文献

- [1] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [2] 鲍怀翘, 徐昂, 陈嘉猷. 藏语拉萨话语音声学参数数据库[J]. 民族语文, 1992(5).
- [3] 呼和, 鲍怀翘, 确精扎布. 关于蒙古语语音声学参数数据库[J]. 内蒙古大学学报(人文社会科学版), 2007(7).
- [4] 娜孜古丽·吐斯甫那比. 哈萨克语语音声学参数数据库研编方法[J]. 民族翻译, 2015(2).
- [5] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [6] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [7] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [8] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [9] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [10] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [11] 曹志耘. 中国语言资源保护的理论与实践[C]. IACCL-23 会议论文, 2015.
- [12] Best, C. T. Nonnative and second-language speech perception: commonalities and complementarities