

Adapting the Tesseract Open Source OCR Engine for Multilingual OCR

Dar-Shyang Lee

Google Inc., 1600 Amphitheatre Pkwy, Mountain View, CA 94043, USA.

然而我们描述使Tesseract开源OCR引擎适应多种脚本和语言的工作，*effort* 导致对于启用通用的多语言操作，这样，*OCR engine for multiple scripts and languages*。Effort has been provided in addition to the new language-specific ones, *除了提供文本集之外，新语言所需的自定义要求也可以忽略，而 that*，尽管需要对于包括物理符号分析和消音后处理在内的答案，*beyond providing a corpus of text*。Although change was required to the module, *模块进行更改，但除几个限制外，无需对字符分类器进行任何更改*，*Tesseract* 分类器已经适应简体中文，*以随机的书籍样本中获得的英语（一种欧洲语言和俄语的混合）的测试*，*English, a mixture of European languages, and Russian, taken from a random sample of books*。Show a reasonable consistent word error rate of 2.73% and 5.78%, *而简体中文的字符错误率仅为 3.77%*，*and Simplified Chinese has a character error rate of only 3.77%*。

关键词 Tesseract 多语言OCR。

Tesseract, Multi-Lingual OCR.

1. Introduction

使用中文、日文和韩文输入代码 (CJQR) 随后是阿拉伯语 (34) 在 a decade ago, in favor of Chinese, Japanese and Korean (CJK) 然后是印地安人 [5] [6] 这些语言经合会 (OECD) 及 OCR 系统的开发 [7,8] 带来了更大的挑战。中文和日语共享汉字, 其中包括成对 classifiers, and also to the other components of OCR systems. 由于上万种不同的字符形状, 朝鲜语使用的汉字具有自己的教子 Chinese and Japanese share the Han script, which contains thousands of different character shapes. 阿拉伯文一到两个数量级。阿拉伯文字, which has several thousand more of its own, as well as using Han characters. the number of characters is one of two

单词中的位置发生变化。印地语将少量字母组合成代表音素的
with connected characters, and its characters change shape
数千种形状。当字母结合在一起时,它们形成曲线,其形状仅
according to the position in a word. Hindi combines a small
与原始字母模糊不清。Hindi then combines all symbols and single
represent syllables. As the letters combine they form ligatures
Oreka结合在一起,从而将CJK和阿拉伯语的问题结合在一起。
whose shape only vaguely resembles the original letters. Hindi
then combines the problems of CJK and Arabic, by joining all the
symbols in a word with a line called the shiro-reka.

研究方法使用了精定于语言的解决方法, 从而以某种方式避免了问题, 因为这个问题找到适用于所有语言的解决方案要简单得多。例如, 试图找到一个适用于所有语言的解决方案。例如, 韩文和印地文的字母符号集主要由数量少得多的组件组成。汉族中的部首字母, 韩文中的字母以及印地文中的字母。由于开发 radicals in Han, Jamo in Hangul, and letters in Hindi. Since it is much easier to develop a classifier for a small number of classes, 并从前部首组合中推断出实际字符。与韩文或印地文相比, 这种方法对韩语来说更容易, 因为

汉字中的部首变化不太,而在汉语中,部首经常被压缩以适合在 Han, the radicals often are squashed to fit in the character and mostly touch other radicals. Hindi takes this a step further by 同时通过改变辅音的形式来进一步采取这一步骤。阿拉伯语的 consonant ligature. Another example of a more language-specific work-around for Arabic, where it is difficult to determine the 另一种针对特定语言的变通方法,其中很难确定字符边界,以将连接的组件分割成字符。一种常用的方法是对单个画像素素 characters. A commonly used method is to classify individual vertical pixel strips, each of which is a partial character, and 符边界建模的隐马尔可夫模型将分类组合在一起 [3]. that models the character boundaries [3].

Google致力于提供尽可能多的语言服务(Voice)，因此我们也有兴趣将Tesseract开源OCR引擎[7, 8]改编为多种语言。本文讨论了我们为Tesseract Open Source OCR Engine [8, 9] 为 many languages, 也我们在使Tesseract全面国际化方面所做的努力, 以及其中一些可能令人惊讶的技巧。我们的方法是使用语言通用方法, 以尽量减少覆盖多种语言的人工工作。

2. Review Of Tesseract For Latin

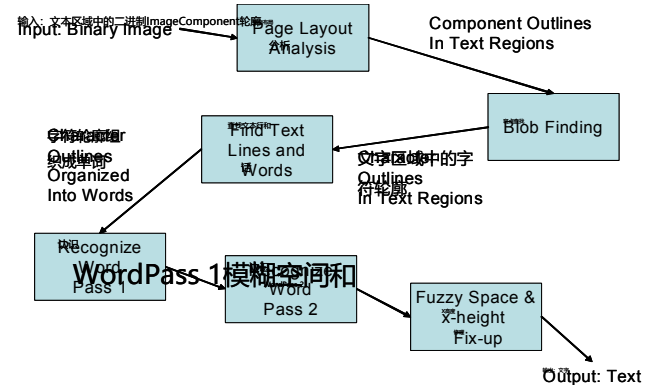
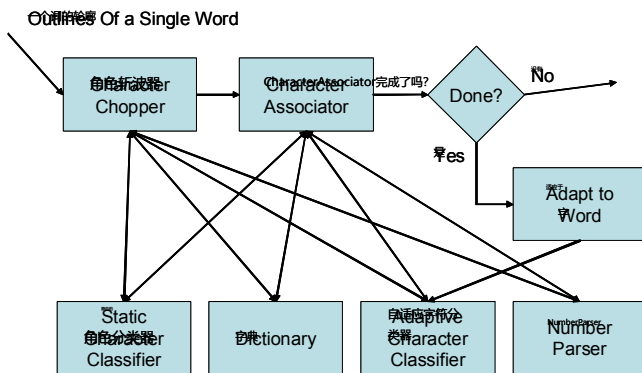


Figure 1. Top-level block diagram of Tesseract.

图 8 是 Tesseract 的基本组成部分的框图。在 Tesseract 1.0 的新网页面上，[The new page layout and analysis](#) 分析从开始就设计成与语言无关，但是该引擎的其余部分都是为从 the beginning to be language-independent, but the rest of the engine was developed for English. 白话文如何通用于其他语言，在英语开发的，没有太多思想，白话文如何通用于其他语言，在英语到当时的商业引擎严格用于黑白文本之后，Tesseract 的最初设计由 commercial engines at the time were strictly for black-on-white text, one of the original design goals of Tesseract was that it should recognize white-on-black. 可能是事实，证明白话文连接组件 (CC) 分析的方向进行，并控制组件的轮廓进行操作。CC 之活的第一步 and operating on outlines of the components. The first step after CC

图2是单词识别器的框图。在大多数情况下,一个blob对应于一个字母,一个blob corresponds to a character, so the word recognizer first classifies each blob, and presents the results to a dictionary search. 因此单词识别器首先对每个blob进行分类,然后将结果提供给出字典搜索。然后在单词中每个blob的分类器选择组合中,找到每个blob in the word. If the word result is not good enough, the next step is to chop poorly recognized characters, where this improves the character classifier's confidence. 如果单词结果不够好,则下一步是将无法识别的字符切掉,从而提高分类器的置信度。After using all the possible possibilities, the system first searches for the best match of the whole word, and then puts back together chopped character fragments, or parts of characters, that were broken into multiple pieces in the original image. 在利用所有可能的组合之后,系统首先搜索所得的最佳图,然后将切好的字符片段或部分字符分解成原始图像中的各个CC。在最佳优先搜索的每个步骤中,将对任何新的斑点组合进行分类,然后将分类器结果再次提供给字典。单词组合 combinations are classified, and the classifier results are given to the dictionary again. The output for a word is the character string listed in the dictionary. 列进行加权后,它具有最佳的基于距离的整体评级的字符串。对于 English version, 大多数标点符号规则都是硬编码的。The English version, most of these punctuation rules were hard-coded.



图像中的单词被处理两次。第一步是将成功的单词(即词典中的单词, successful words, being those that are in a dictionary and are not dangerously ambiguous, are passed to an adaptive classifier for final classification)一旦有足够的样本,就可以提供分类结果;即使是在第一次通过时。在第五遍中在第一遍的单词不够好,是

3.1 垂直文本和中文文本识别和检测文在不同程度上均水平或垂直提取所有文本行，并且经常在单个页面上混合方向。这个问题不是CC独有的，因为英语杂志的页面通常在照片或文章的English面使用竖排文字来表示摄影师或作者的信誉。通过页面布局分析，可以检测到竖排文字，如果制表位上的大多数CC的左侧都在垂直线被检测到。如果制表位上的大多数CC的左侧都在右侧选项卡上，而右侧都在右侧选项卡上，则制表位之间的所有内容都可能是行垂直文本。为了防止假阳性的不稳定性，进一步的限制要求垂直文本的CC之间的垂直中间距离应小于CC的垂直间距。如果页面上的大多数CC垂直对齐，则将页面旋转90度并再次运行页面布局分析，以减少在垂直文本中发现错误的可能性。然后少数原本水平的文本将在旋转后的页面中变为垂直文本，并且文本的正文将为水平。The minority originally horizontal text will then become vertical and the text of the text will be horizontal.

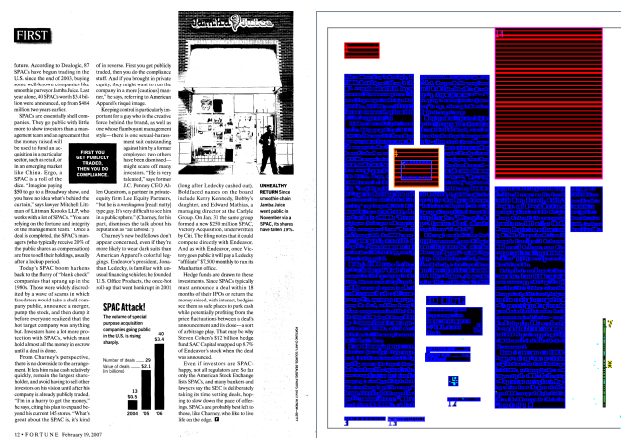


Figure 5 (a) A page containing a vertical text region.
(b) The detected regions with image in red, horizontal text in blue and yellow vertical text in yellow.

Figure 4. The vertical text is differentially rotated so it is oriented horizontally.

1. 2. 3. 4. 5. 6. 7. 8. 9. 10. 11. 12. 13. 14. 15. 16. 17. 18. 19. 20. 21. 22. 23. 24. 25. 26. 27. 28. 29. 30. 31. 32. 33. 34. 35. 36. 37. 38. 39. 40. 41. 42. 43. 44. 45. 46. 47. 48. 49. 50. 51. 52. 53. 54. 55. 56. 57. 58. 59. 60. 61. 62. 63. 64. 65. 66. 67. 68. 69. 70. 71. 72. 73. 74. 75. 76. 77. 78. 79. 80. 81. 82. 83. 84. 85. 86. 87. 88. 89. 90. 91. 92. 93. 94. 95. 96. 97. 98. 99. 100. 101. 102. 103. 104. 105. 106. 107. 108. 109. 110. 111. 112. 113. 114. 115. 116. 117. 118. 119. 120. 121. 122. 123. 124. 125. 126. 127. 128. 129. 130. 131. 132. 133. 134. 135. 136. 137. 138. 139. 140. 141. 142. 143. 144. 145. 146. 147. 148. 149. 150. 151. 152. 153. 154. 155. 156. 157. 158. 159. 160. 161. 162. 163. 164. 165. 166. 167. 168. 169. 170. 171. 172. 173. 174. 175. 176. 177. 178. 179. 180. 181. 182. 183. 184. 185. 186. 187. 188. 189. 190. 191. 192. 193. 194. 195. 196. 197. 198. 199. 200. 201. 202. 203. 204. 205. 206. 207. 208. 209. 210. 211. 212. 213. 214. 215. 216. 217. 218. 219. 220. 221. 222. 223. 224. 225. 226. 227. 228. 229. 230. 231. 232. 233. 234. 235. 236. 237. 238. 239. 240. 241. 242. 243. 244. 245. 246. 247. 248. 249. 250. 251. 252. 253. 254. 255. 256. 257. 258. 259. 260. 261. 262. 263. 264. 265. 266. 267. 268. 269. 270. 271. 272. 273. 274. 275. 276. 277. 278. 279. 280. 281. 282. 283. 284. 285. 286. 287. 288. 289. 290. 291. 292. 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 303. 304. 305. 306. 307. 308. 309. 310. 311. 312. 313. 314. 315. 316. 317. 318. 319. 320. 321. 322. 323. 324. 325. 326. 327. 328. 329. 330. 331. 332. 333. 334. 335. 336. 337. 338. 339. 340. 341. 342. 343. 344. 345. 346. 347. 348. 349. 350. 351. 352. 353. 354. 355. 356. 357. 358. 359. 360. 361. 362. 363. 364. 365. 366. 367. 368. 369. 370. 371. 372. 373. 374. 375. 376. 377. 378. 379. 380. 381. 382. 383. 384. 385. 386. 387. 388. 389. 390. 391. 392. 393. 394. 395. 396. 397. 398. 399. 400. 401. 402. 403. 404. 405. 406. 407. 408. 409. 410. 411. 412. 413. 414. 415. 416. 417. 418. 419. 420. 421. 422. 423. 424. 425. 426. 427. 428. 429. 430. 431. 432. 433. 434. 435. 436. 437. 438. 439. 440. 441. 442. 443. 444. 445. 446. 447. 448. 449. 450. 451. 452. 453. 454. 455. 456. 457. 458. 459. 460. 461. 462. 463. 464. 465. 466. 467. 468. 469. 470. 471. 472. 473. 474. 475. 476. 477. 478. 479. 480. 481. 482. 483. 484. 485. 486. 487. 488. 489. 490. 491. 492. 493. 494. 495. 496. 497. 498. 499. 500. 501. 502. 503. 504. 505. 506. 507. 508. 509. 510. 511. 512. 513. 514. 515. 516. 517. 518. 519. 520. 521. 522. 523. 524. 525. 526. 527. 528. 529. 530. 531. 532. 533. 534. 535. 536. 537. 538. 539. 540. 541. 542. 543. 544. 545. 546. 547. 548. 549. 550. 551. 552. 553. 554. 555. 556. 557. 558. 559. 560. 561. 562. 563. 564. 565. 566. 567. 568. 569. 570. 571. 572. 573. 574. 575. 576. 577. 578. 579. 580. 581. 582. 583. 584. 585. 586. 587. 588. 589. 590. 591. 592. 593. 594. 595. 596. 597. 598. 599. 600. 601. 602. 603. 604. 605. 606. 607. 608. 609. 610. 611. 612. 613. 614. 615. 616. 617. 618. 619. 620. 621. 622. 623. 624. 625. 626. 627. 628. 629. 630. 631. 632. 633. 634. 635. 636. 637. 638. 639. 640. 641. 642. 643. 644. 645. 646. 647. 648. 649. 650. 651. 652. 653. 654. 655. 656. 657. 658. 659. 660. 661. 662. 663. 664. 665. 666. 667. 668. 669. 670. 671. 672. 673. 674. 675. 676. 677. 678. 679. 680. 681. 682. 683. 684. 685. 686. 687. 688. 689. 690. 691. 692. 693. 694. 695. 696. 697. 698. 699. 700. 701. 702. 703. 704. 705. 706. 707. 708. 709. 710. 711. 712. 713. 714. 715. 716. 717. 718. 719. 720. 721. 722. 723. 724. 725. 726. 727. 728. 729. 730. 731. 732. 733. 734. 735. 736. 737. 738. 739. 740. 741. 742. 743. 744. 745. 746. 747. 748. 749. 750. 751. 752. 753. 754. 755. 756. 757. 758. 759. 760. 761. 762. 763. 764. 765. 766. 767. 768. 769. 770. 771. 772. 773. 774. 775. 776. 777. 778. 779. 780. 781. 782. 783. 784. 785. 786. 787. 788. 789. 790. 791. 792. 793. 794. 795. 796. 797. 798. 799. 800. 801. 802. 803. 804. 805. 806. 807. 808. 809. 810. 811. 812. 813. 814. 815. 816. 817. 818. 819. 820. 821. 822. 823. 824. 825. 826. 827. 828. 829. 830. 831. 832. 833. 834. 835. 836. 837. 838. 839. 840. 84

3.2 Text-line and Word Finding

修改后的改进行直找算法可从左起分析串为每个文本区域独立工作 each text region from layout analysis and begin by searching the neighborhood of small CCs (relative to the estimated text size) 并以搜索较小CC的邻域。它相对于估计的文本大小开始。以 neighborhood of small CCs (relative to the estimated text size) 找到最接近正文文本大小的CCs。如果邻域没有正文文本大小的CCs, 则将一个小的CC视为可能由噪音, 并得将其丢弃。 (通常性目被 discarded (An exception has to be made for dotted/dashed leaders, as typically found in a table of contents.) Otherwise, a larger neighbor is constructed and used in place of the bounding box of the small CC in the following projection.

确定切割线后, 将整个连接的零部件放置在文本线中, 使其垂直重叠 components are placed in the text line that they vertically overlap the most, (still) using the modified boxes) except when a 雷霸假定这些GG是来自多行触摸的字符, 因此需要在剪切行处进行剪别, 或者是在单字中 cut, 在这种情况下, tips 它们应位于重叠的, and so need cutting at the cut line, or drop caps, in which case 顶部符号, 即使位于五划拉伯语。该算法也能中常工作。 This algorithm works well, even for Arabic.

©ACM, 2009年. 这是作品的作者版本。经ACM许可将其张贴在此处供您个人使用, 不得进行重新分配。最终版本发布于2009年7月25日在西班牙巴塞罗那举行的2009年多语言OCR国际研讨会论文集上。 <http://doi.acm.org/10.1145/1577802.1577804> Workshop on Multilingual OCR 2009, Barcelona, Spain July 25, 2009. <http://doi.acm.org/10.1145/1577802.1577804>.

解决了很长的斑点序列问题。这是确定分割图时确定效率和质量的关键因素。这至少可以使得识别方法的长度减少到更易于管理的大小。[searching the segmentation graph, this would at least reduce the lengths of recognition units into more manageable sizes.]

如第2节中所述,黑白文本的检测基于轮廓的嵌套复杂度。同样,基于轮廓的过滤过程也可以拒绝非文本,包括半色调噪声、侧面的黑色区域或侧边栏或反向视频区域中的大型容器盒。过滤的一部分基于对斑点的拓扑复杂性的度量,并根据内部组件的数量、嵌套孔的层数、周长与面积之比等进行估算。但是,无论如何,繁体中文文字的复杂性通常都超过了包装在框中的英文单词的复杂性。因此,解决方案是对不同的语言应用不同的复杂度阈值,并依靠后续分析来恢复任何错误拒绝的斑点。

3.3 估算西里尔字母的x高度和上升高度,并将斑点块组织成行

组织成行后, Tesseract 估算每个文本行的x高度。x高度估算算法首先确定基于最大和最小可接受x高度的界限,根据为该块计算的初始行大小。然后,对于每像素,分别量化在像素上出现的斑点的边界框的高度,并将其聚合为直方图。从该直方图中, x高度查找算法查找两个最常见的高度模式。它们相距足够远,以达到潜在的x高度和上升高度。为了获得对某些噪声的鲁棒性,该算法确保了高度相对于该行上的斑点总数,选择为x高度和上升高度的模式具有足够的数目或出现次数。the x-height and ascender height have sufficient number or occurrences relative to the total number of blobs on the line.

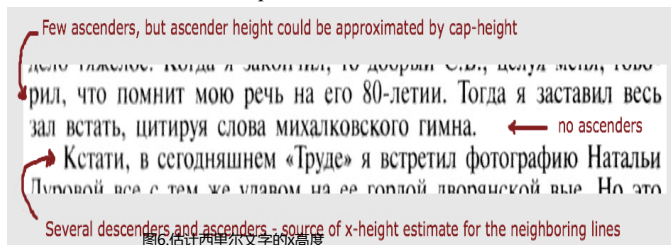
对于大多数拉丁字体,该算法效果很好。但是,当按原样应用于西里尔字母时, Tesseract 无法为大多数线找到正确的x高度。结果,在 Russian books 的数据集上, Tesseract 的单词错误率被证明是9%。如此高的错误率的原因是双重的。首先,西里尔字母的统计特性与拉丁字母显著不同。仅降低期望的每行上升数阈值并不是一个有效的解决方案,因为拉丁文本很少包含一个或不包含上升字母。因此,如此差的性能的第二原因是西里尔字母的大小写混淆不清。例如,在33种大写现代俄语字母中,只有6种小写字母的上半部分与大多数现代俄语字母明显不同。因此,当使用西里尔字母时, Tesseract 会很容易被错误的x高度信息误导,会轻易将小写字母识别为大写字母。

我们解决西里尔字母x高度问题的方法是调整行上的最小上升数。考虑下降统计信息,并更有效地使用同一文本块中相邻行的x高度信息。neighborhood是由分页版面分析确定的文本块,更有效地 (a block is a text region identified by the page layout analysis). 基于文本块点和行间距的大小,因此可能包含相似间距或相似字体大小的字母。

对于给定的文本块,改善的x高度查找算法使用行高度信息。首先,它尝试找到每行的x高度。基于on

计算结果的每一行都属于以下四类之一: (1) 找到x高度和上升模式的行; (2) 找到下降模式的行; (3) 可以使用的通用x高度; (4) 没有找到x高度和上升模式的行。如果在该行中找到第一类中具有可靠的x高度和上升模式,则将其高度估计用于第二类中具有类似x高度估计值的行 (存在降序的行)。对于第三类中的那些行 (没有可靠的x高度和上升模式), 则使用相同的x高度估计值。对于第四类中的那些行 (没有找到x高度和上升模式), 则使用相同的x高度估计值。如果逐行方法未能找到任何可靠的x高度和上升模式,则将汇总文本块中所有斑点的统计信息,并使用此累积信息重复搜索x高度和上升模式。

作为上述改进的副产品, Tesseract 改进了对x高度和上升模式的估计。



然而,即使是在人类读者也必须使用来自相邻文本块或有关书籍的上下文信息来确定的信息来确定的行是大写还是小写。the information from the context of the neighboring blocks or from knowledge about the common organization of the books to determine whether the given line is upper- or lower-case.

4. 字符识别在 Word Recognition 方面

克服的主要挑战之一是将最初为字母语言设计的内容扩展为处理诸如中国和日语之类的表意语言。这些语言的特点是符号大量,缺少清晰的单词边界。这对从小写字母很好界定的单词而设计的搜索策略和分类引擎进行了严格的测试。4.1 分段和搜索如第3.2节所述,对于下一节,描述所需的修改。

4.1 分段和搜索

As mentioned in section 3.2, for non-space delimited languages like Chinese, recognition units that form the equivalent of words in western languages now correspond to punctuation delimited phrases. Two problems need to be considered to deal with these phrases. Tesseract 在分割图上采用最佳优先搜索策略,该策略随Blob序列的长度呈指数增长。While this approach worked on shorter Latin words with fewer

当在字典中找到结果时, 一切分点和终止条件会给短语短语分类is found in the dictionary, it often exhausts available resources when classifying a Chinese phrase. To resolve this issue, we need to 经常耗尽可用资源。要解决此问题, 我们需要大幅度减少排列表中评估的分割点的数量, 并设计出更容易满足的终止条件。ated in the permutation and devise a termination condition that is easier to meet.

为了减少分割点的数量,我们将大致恒定的字符宽度的约束并 incorporate the constraint of roughly constant character widths in a mono-spaced language like Chinese and Japanese. 在这些语言中,字符多数具有相似的长宽比,并且在其位置上是全角还是半角 are either full-width or half-width in their positioning. Although the normalized width distribution would vary across fonts, and the 和数字或拉丁词的包含而变化,这并不罕见,但总的来说,这些约束为特定的分段点是否与另一个分段点兼容提供了强有力的指导。因此,使用与分割模型的偏差作为代价,我们可以消除许多令人难以置信的分割状态,并有效地减少了搜索空间。我们不使用此估计值基于最佳局部解决方案来修剪搜索空间 the search space based on the best partial solution, making it effectively a beam search. 当没有进一步的扩展可能会产生更好的解决方案时,这也提供了终止条件。Likely to produce a better solution.

另一个强大的约束是短语中字符脚本的一致性。当我们包含多个形状类 within a phrase. As we include shape classes from multiple 脚本之间, 跨不同脚本的字符之间的混淆将变得不可避免。景符我们可以 cross different scripts become inevitable. Although we can establish the 页面的主要脚本或语言, 我们也必须允许使用拉丁字符, 因为 in characters as well, since the occurrence of English words inside 外语书籍中英语单词的出现非常普遍。在一个识别单元内的字 foreign language books is so common. Under the assumption that 符具有相同脚本的假设下, 如果改善整个单元的总脚本 script, we would promote a character interpretation if it improves the 性, 我们将进步字符解释。但是, 如果单词或短语是真正的混 the overall script consistency of the whole unit. However, blindly 命脚本, 则基于先验盲目地提升脚本字符实际上可能会损害性 the performance of the word or phrase recognition task. 能。因此, 仅当 tope exploration 中超过一半的字符属于同脚本。 we apply the constraint only if over half the characters in the top interpreted as belonging to the same script. 本, 并且像其他排列一样, 对形状识别分数加权调整时, 本应 is weighted against the shape recognition score, like any other 用约束。 permutation.

4.2 形状Shape用于大量类别的分类器

4.2 形状分类用于大量类别的分类器仍然是一个研究问题,即使在今

特征是形状轮廓的多边形近似的分量。在训练中，从多边形重绘的每个元素导出 (x, y-位置, 方向, 长度) 的 4 维特征向量。然后将元素聚类以形成原型特征向量 (因此，名称为 Tesseral feature vectors)。 (Hence the name, Tesseral)。在识别中，多边形的元素被分解成等长的较短片段。因此，长度维从特征向量中消除了。多个短特征相互匹配 from the feature vector. Multiple short features are matched against each

来自训练的原型特征, 这使得分类过程对残缺字符的鲁棒性更高, making the classification process more robust against broken characters.

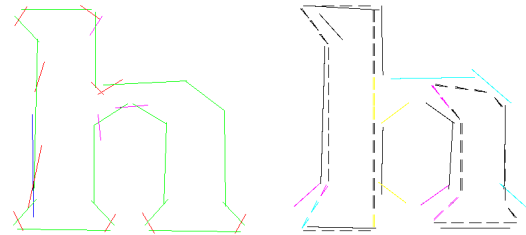


Figure 7.7.(a) Prototype for (b) missing h against the prototype.

图 7.7(a) 显示了字体 Times Roman 的字母 letter ‘t’ 的系列原型。Times Roman. The green line-segments represent cluster means of significant clusters that contain samples from almost every sample of ‘t’. 样本的样本。蓝色部分 segments 意味着与另一个簇合并 merged with another cluster to form a significant cluster. Magenta segments were not used, as they matched an existing significant 簇匹配。红色部分的样本数量不足以构成显著的, 并且无法与任何相邻的簇合并以形成显著的簇。

图 7(b) 显示了未知字符的较短特征如何与原型匹配, 以实现对抗 *h* against the prototype to achieve insensitivity to broken characters. 较薄的线是未知特征 (被打破的 *h*)。The short, thick lines are the features of the unknown, being a broken *h* and in original lines are the prototype features. 红色 *c* 与蓝色 *b* 匹配良好, 黄色 *y* 与黑色 *b* 匹配较差, 蓝色 *b* 与黄色 *y* 不匹配。垂直原型都很好匹配。cyan \rightarrow poor, and yellow \rightarrow no match. The vertical prototypes are all well matched, despite the fact that the *h* is broken.

形状分类器分两个阶段运行。第一阶段称为维数修剪，它使用与Locality Sensitive Hashing (LSH) [15] 密切相关的方法将字符集缩减为 short 个字符的简短 characters，using a method closely related to Locality Sensitive Hashing (LSH) [15]. The final stage computes the distance of the unknown from the prototypes of the characters in the short-list.

On-line algorithm designed as a simple and vital time-saving
通过优化 on 类修剪器通过分别考虑每个 3-D 特征来划分高维特
feature space, by considering each 3-D feature individually. In
place of the hash table of LSH, there is a simple look-up table,
范围为 [0, 255] 的整数向量 μ 每个字符集对应一个字符集, 该值
class in the character set, with the value representing the
表示该字符集的近似匹配度, 角色类型原型的特征。同量结果在末
期的所有特征之间求和。总分数在最高分数的 $\frac{1}{10}$ 部分之内的
类别将作为下一个阶段要归类的候选清单返回。余修总相对较
快, but its time scales linearly with the number of classes and also
with the number of features.

第二阶段对分类器计算每个特征与其最近原型的距离 d_i ，作为 i 维特征空间中从 μ_i 的nearest prototype-as the squared Euclidean distance d_i of the (x,y) feature coordinates from the prototype line 平方 d_i 再加上与该特征线成角度的加权 (w) 着原型 angle θ from the prototype:

$$2\alpha_{\text{eff}} + w\theta^2$$

从本质上讲,这是一个生成分类器。它可以计算与理想目标的距离,并使用以下公式将特征距离转换为特征证据。The feature distance is converted to feature evidence E_f using the following equation:

$$f_k d_{f, \square} \frac{1}{1 + k d_f^2}$$

常数 k 用于控制证据随距离衰减的速率 α 。当特征与原型匹配时, α decays with distance. As features are matched to prototypes, the feature evidence E_F is copied to the prototypes E_P . Since the 特征证据 E_F 被复制到原型 E_P 。由于原型期望多个特征与它们匹 配, 并且“最佳匹配”的收集是为了速度而独立完成的, 因此 the sums of feature and prototype evidence can be different. The sums 特征和原型证据的总和可能不同。总和由特征数量和原型长度 归一化。结果转换回距离。number of features and sum of prototype lengths L_p , and the result is converted back into a distance:

$$\alpha_{\text{决赛}} = 1 - \frac{\sum_f E_f^E + \sum_p E_p}{N_f + \sum_p L_p}$$

请注意, 实际的实现使用定点整数算术运算, 并且上述公式中省略了很多 ϵ 比例项或常数, 否则这些比例常数会混淆计算。obscure the calculations are omitted from the equations above.

第Ⅱ阶段分类器的部分优势在于, 每个候选类簇中都可以有多个理想值 (multiple ideals (known as *configs* in Tesseract), within each class label, thus allowing multi-modal distributions that may be caused by arbitrary character variations). 因此可以呈现出字体或版式上的任意差异引起的各模式分布, 上述匹配过程在计算最终距离时, 选择最佳匹配, 识别这个语义下, 该分类器因此实际上是跟最近的邻居 (the final distance). In this sense, the classifier is thus effectively a nearest neighbor classifier.

我们假设类修剪器和辅助分类器对于大量类来说效果很好, *work well for large numbers of classes*, 这是因为它们在很多个不同尺度的“弱分类器”之间使用投票, *because of their use of voting among multiple weak classifiers of small dimension*, 是依赖于一个高维度的分类器。这是提顿[15]背后的概念, *it is the very concept behind boosting* [15], 只是当前的权重分类器尚未加权, *except that currently the weak classifiers are not weighted*. 特征空间的尺寸被量化为256个级别, *The dimensions of feature space are quantized into 256 levels*, 这提供了足够的精度来存储CJK字符和印度音节的复杂形状, *This provides enough precision to store the complex shapes of CJK characters and Indic syllables*, 并且量纲计算避免了类似于类修剪器的尺寸诅咒, *and the d_{final} calculation avoids the curse of dimensionality in a similar fashion to the class pruner*.

5.上下文处理Contextual Post-Processing

单词列表生成新语言的字典的方法, 支持部分扩展语言模型 (the language model by providing a way to generate dictionaries for new languages). 为了紧凑和快速搜索, 这些字典由有向无环词图 (DAWG) 表示。在这些字典中, DAWG 数据结构用于顺序搜索 (fast search, these dictionaries are represented by directed acyclic word graphs (DAWGs)). In the original implementation, the dictionaries include the pre-generated system dictionary, the OCRred dictionaries (including the OCRred document and a user-provided word list).

最初,DAWG中的每个边都存储了一串0-63字符,代表用于DAWG中相应转换的字母。但是,letter这个表示是有限的,corresponding transition in the DAWG. This representation, however, was limiting, since 以这种形式处理多字符串索引和多字符Unicode字符很尴尬。修改了Dale Wong的数据结构,以存储字符分类器使用的unicharset。这大大简化了 the character classifier instead. This significantly simplified the

[illegible]

5.1 约束模式Constraint Patterns

没有泛化到拉丁文字之外。即使对于拉丁文字，Tesseract 也不是 hard-coded 并且 did not generalize beyond the Latin scripts. Even 接受很多 damage 有效的标点和数字模式。为了帮助 Tesseract 处理非拉丁文字，处理标点符号和数字，Tesseract 的培训过程扩展了代码，以收集和编码了组频繁出现的标点符号和数字，并 encode a set of frequently occurring punctuation and number patterns. The step for collecting these patterns was implemented 并行进行，以构造给定语言的字典。为了表示和匹配生成 的模式，使用了已经存在的用于生成和搜索单词 DAWG 的代 码。对确定给定分类器选择是否为给定语言的 Tesseract 中 的有效单词的算法进行了一些修改。对包含单词、标点和数 字模式的所有 DAWG 进行了同时搜索。通过此修改，可以删 除所有针对数字和标点的语言和脚本特定的硬编码规则。标 点和数字模式的生成和搜索过程被设计为完全由数据驱动，of generating and searching the punctuation and number patterns 到目前为止，对于任何语言和脚本，都不需要特殊的要求。no special casing for any language in particular.

5.2.2 Resolving Shape Ambiguities Tesseract's shape

Hyphenic 包 括 一 个 自 适 应 生 成 用 于 OCR 文 档 中 看 到 的 shape classifier，包 括 了 一 个 适 应 性 组 件 (adaptive component) 它 学 习 了 文 档 中 的 特 征 特 性 以 及 学 习 了 已 被 标 记 的 文 档 中 的 特 征 特 性 以 确 保 能 仅 适 应 明 确 的 字 典 词。如 果 满 足 两 个 约 束 条 件，则 将 OCR 的 classifier 只 适 应 于 一 个 无 歧 义 的 字 典 词。The 中 文 称 为“明确”第 一 个 是 形 状 无 歧 义 类 别 必 须 在 单 词 的 所 有 其 他 OCR'd word 是 “dubbed” unambiguous。If it satisfies two 选 择 中 识 别 出 明 显 的 赢 家 (即 最 佳 选 择 的 分 类 器 评 级 必 须 明 显 a clear winner among all the alternative choices for the word) (是 高 于 次 优 选 择 的 评 级)。第 二 个 约 束 是 不 能 存 在 字 典 词。其 形 状 对 于 同 词 的 最 佳 选 择 而 言 是 模 棱 两 可 的。此 要 求 对 识 别 速 度 也 很 重 要，因 为 一旦 两 个 字典 词 的 分 数 可 以 存 在 于 同 一 个 字 典 词 的 形 状 中，则 无 法 选 择 到 最 佳 选 择 的 字 典 词。This requirement is also important，它 就 可 以 接 受 识 别 结 果 并 停 止 对 该 单 词 的 进 一 步 处 理 (after score) once Tesseract finds such a word choice, it could accept the recognition result and stop further processing of the word.

对于拉丁文或 `ripless` 包含一个手工制作的数据文件，称为“冗长”（referred to as “longerous ambiguities” file）specifying which 候选文件，该文件指定大多数拉丁字体中哪些字母组合固有地模棱两可，一种使用其他脚本实现此功能的可汗屏解决方案，是开发一种自动生成任意给定 `n`-gram 对列表的自动方法，generating a list of ambiguous `n`-gram pairs for any given

语言וגא大型文本语料库中收集了这一组n-gram。在这种情况下，combined weight accounts for 95% of all n-grams in the language (为uni- bi-gram)。其组合权重占该语言中所有n-gram的95%。was collected from a large text corpus。The n-grams were rendered with n-gram用一组常用字体以几种降阶模式和曝光进行渲染，and exposures。Then, Tesseract's shape classifier was run over the rendered images上运行Tesseract的形状分类器，以获得每个n-gram的渲染图像上的一组得分。a set of top scoring classifications for egram的每组最高得分分类，resulting in n-gram的每组最终得分。was assigned to each n-gram。The n-gram的每个得分被分配给了每个n-gram。统计信息，并丢弃具有低分类得分的离群值。(在某些字体和降阶分类器模式下，某些n-gram的得分被丢弃，因为某些字体和降阶模式下的字符被渲染为无法识别，并且在某些情况下，只含污染数据)。然后，对于每个n-gram，计算其得分。The n-gram和相应的正确或不正确的OCR'd n-grams和的相应正确或不正确的n-gram对，计算歧义分数。模糊度分数被定义为形状分类器pair of ambiguity score was computed。The ambiguity score was 感知到的错误和正确的n-gram跨所有字体和降阶模式的一致性 between the wrong and correct n-grams (根据n-gram频率和分类与语言中正确的n-gram频率之间的相似度的函数。为了在less-gram的速度和准确性之间达到所需的平衡，有必要选择一阈值 between Tesseract's speed and accuracy, it was necessary to pick a threshold of the expected number of errors allowed to occur due 器错误统计信息计算) 而在发生的预期错误数上生成 n-gram frequency and classifier error statistics)。To generate the "危险歧义" 或作歧义歧义 n-gram歧义歧义分数的递增顺序进行排序。sorted in the non-increasing order of the ambiguity scores and the appropriate number of top-scoring ambiguities that ensured the desired expected error rate were included in the file。

使用这种自动方法生成的数据文件是与使用手工制作的文件同样危险的。可能的改进包括使用与使用拉丁字母 (FEIGS 数据集) 相比, 有可能在拉丁字母 (FEIGS 数据集) 数据集上使用的改进 (尽管在某些语言中结果稍差)。在俄罗斯数据集上使用自动生成的文件可以减少 10% 的单词错误率。对其它语言生成的文件进行检查也表明, 自动生成的文件包含相当多的常见形状混淆, 但是将相对应的数据集进行进一步测试需要量化改进。数据集将需要被用于量化改进。

5.3.3 Handling Highly Inflected Languages

和相关要素最大限度地增加这些语言同时最小化存储词典所消耗的空间始终是 个挑战, while minimizing the space consumed to store the dictionary. 从语料库以高度变化的语言生成字典是一项特别困难的任 务, because language is particularly difficult task. The frequency of words in highly inflected languages is more evenly distributed, and thus to 范围就需要使用更大的词典. 此外, 即使是更频繁使用的单词, Many single word forms of the word forms of even the more frequent words might not occur enough times in the training corpus to be 法包含在其中. 因此字典可能无法在训练语料库中概括. Figure 8 通过语料库的覆盖范围与选择形成词典的最常用单词的 problem on a collection of languages by graphing the coverage of the corpus against the number of most frequent words chosen to form the dictionary.

由于DAWG数据结构的头部和尾部压缩,添加DAWG中已经存在的单词结构, adding an inflected form of a word that already exists in the structure, 可能会导致插入单词大小上的微小增加。这是因为该单词的DAWG might result in a very small increase in the overall size of the dictionary. 这是因为该单词的beginnings and the endings of the word 则已在字典中添加单词 in the DAWG (for example it would be relatively cheap to add the word "talking" to the dictionary if "talk" and "making" have already been inserted).

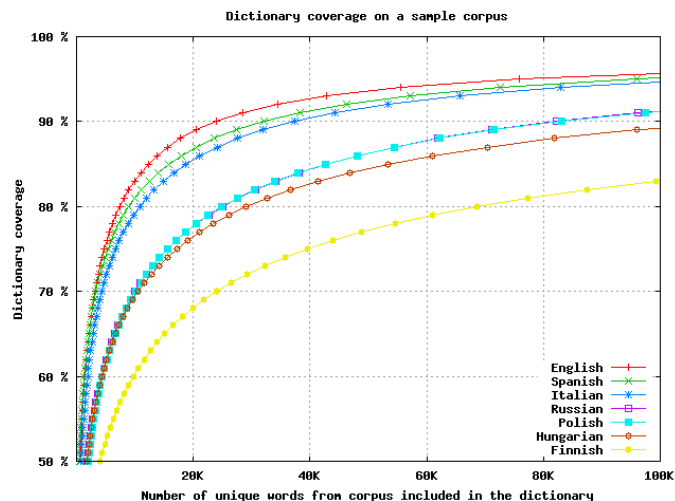


Figure 8. Corpus coverage of varying-size dictionaries in a collection of languages.

[illegible]

6. 状态实验结果

6. Status of Experimental Results

然后将数据集分为训练集、验证和测试集，其中训练和验证集用于开发过程中学习和基准化算法，测试集保留用于发布期间的最终评估。表1总结了几种语言的数据集大小和当前准确性。对于字母语言，我们报告了字符级和单词级的情误率。对于单词含义不明确的中文，我们仅报告字符替换率。EFGSD是英语、法语、意大利语、德语、西班牙语和荷兰语的组合。EFGSD is a combination of English, French, Italian, German, Spanish and Dutch.

对于简体中文中, 我们注意到不同书籍之间的错误率差异很大, 差异主要归因于字体的质量的变化。在页面质量和排版合理的地方, 错误主要是由于相似或接近相同的形状类别之间的混淆所致。我们计划增加形状匹配器中特征空间的容量, 这应该有助于区分相似的

那些在几乎相同的情况下，同种字体的类内字体可能由于类间字体而分离，它们的用法和先验差异如此之大，以至于当我们使用 CJK 语言模型时，很容易地对其进行纠正。

Table 1. Error rates over various languages

| 语言数量 | No. of chars
(millions) | 字数
of words
(millions) | 字符错误率
rate (%) | 单词错误率
error rate (%) |
|-------------------------|----------------------------|------------------------------|-------------------|-------------------------|
| 英语 English | 53.72 | 39 | 0.5 | 3.72 |
| 西班牙语 Spanish | 26.075 | 213 | 0.75 | 5.78 |
| 俄语 Russian | 35.48 | 38 | 1.35 | 5.48 |
| 简体中文 Simplified Chinese | 0.25 | 377 | NA | NA |
| 阿拉伯语 Arabic | 4.038 | 15.41 | 0.33 | 69.44 |

7. 结论和未来工作我们已经讨论了体工作适用于多种语言的实验，并且令人惊讶地发现这主要是工程问题。没有任何对分类器进行了重大更改，我们能够获得多种基于拉丁字母的语言，甚至简化的中文。但是我们的测试集包含新旧字体的混合，并且很大一部分错误是由于训练集不包含旧字体的事实造成的。这项工作尚未涵盖从右到左书写的语言，这主要是另一个工程问题，但是期望白语具有 Tesseract 可能无法解决的系列问题，即与另一个工程问题，即阿拉伯语具有其自己的 set of problems 有关。我们还没有讨论过的另一种语言是泰语，它带来了高度歧义的字符问题，并且与中文一样，单词之间没有空格。像 Thai，which poses problems of highly ambiguous characters, and like Chinese, does not have spaces between words.

一个重要的未来项目是改善训练过程，使其能够使用真实数据进行训练，而不仅仅是具有字符边界框的合成数据。这将大大提高识别的准确性。我们还需要测试阿拉伯语和泰语，因为我们目前还没有更多问题。对于中文、日语和泰语，我们需要允许使用语言模型来搜索任意串联的单词的空间，因为这些语言的单词之间没有空格。尽管德语复合功能会增加大小与更改和插入字母的复杂性，但相似的功能对德语也很有用。

8. 参考文献

[1] Nag, G., “汉字识别：过去十五年的前景”，第 9 卷，five-year perspective,” 9th Int. Conf. on Pattern Recognition, Nov 1988, pp163-167.
[2] 夏芳，基于知识的子模式分割与汉字分解，图像处理程序，1994，IEEE 国际会议，1994, Proc. ICIP-94, IEEE Int. Conf. vol.1, 13-16 Nov 1994, Conf. 第1卷, 1994年11月13日至16日, pp179-182.

[3] 陆志东, Ghiswara, S., Rana, P., Rajan, P., Bazzaz, M., Houl, J. “BBN-ADYLOS OCR 系统的先进技术”, Proc. 5th Int. Conf. on Document Analysis and Recognition, 1999, pp337-340.
[4] Karan, T., Marton, M., Bulboacă, O., Bulboacă, C., Onuț, S., 拉伯 OCR 引擎的开源实现, 2001, Proc. SPIE 3651, 7 Jan 1999, pp109-120.
[5] Bansal, V., R. M. K. Singh, “一个完整的 OCR 引擎用于印地语文本在 Devanagari 脚本”, Proc. 6th Int. Conf. on Document Analysis and Recognition, 2001, pp800-804.
[6] Govindaraj, V., et al., “A Template-based OCR Engine for Document Image Analysis for Libraries”, 2004, pp122-133.
[7] Google 官方博客, <http://googleblog.blogspot.com/2008/07/hitting-40-languages.html>.
[8] Smith, R., “A Tesseract OCR 引擎概述”, Proc. 6th Int. Conf. on Document Analysis and Recognition, 2007, pp629-633.
[9] Tesseract Open-Source OCR: <http://code.google.com/p/tesseract-ocr>.
[10] Smith, R., “Hybrid Page Layout Analysis via Tab-Stop Detection, Document Analysis and Recognition” Proc. 10th Int. Conf. on Document Analysis and Recognition, 2009.
[11] Smith, R., “A simple and efficient skew detection algorithm via text row accumulation”, Proc. 3rd Int. Conf. on Document Analysis and Recognition, 1995, pp1145-1148.
[12] Unnikrishnan, R., Smith, R., “Combined Script and Page Layout Analysis for Document Image Analysis”, Proc. 10th Int. Conf. on Document Analysis and Recognition, 2009.
[13] Gionis, A., Indyk, P., Motwani, R., “Similarity Search in High Dimensions via Hashing”, Proc. 25th Int. Conf. on Very Large Data Bases, 1999, pp518-529.
[14] Baluja, S., Covell, M., “Learning to hash: forging hash functions and applications”, Data Mining and Knowledge Discovery 17(3), Dec 2008, pp402-430.
[15] Schapire, R. E., “The Strength of Weak Learnability”, Machine Learning, 5, 1990, pp 197-227.