

Chapter 1

Notes

This project uses state of the art Chinese/Japanese handwriting recognition methods in order to provide an Kanji teaching application with an error correction.

Conceptually, the application is an e-learning environment for Japanese characters, intended for the foreign learner of the Japanese language. In order to provide more than a multiple choice method, like most other systems, the application contains a handwriting recognition engine that can be used preferably with a hand-held device like a PDA, but generally any stylus input device.

The handwriting recognition method used is similar to the one proposed by Chen, J.-W. and S.-Y. Lee (1996) in their article "A Hierarchical Representation for the Reference Database of On-Line Chinese Character Recognition" in "Advances in Structural and Syntactical Pattern Recognition", Volume 1121 of "Lecture Notes in Computer Science", pp. 351--360. Berlin/Heidelberg, Germany: Springer.

(Nakagawa, Tokuno, Zhu, Onuma, Oda, and Kitadai 2008) report their recent results of online Japanese handwriting recognition and its applications. Their article gives important insights into character modeling, which are employed in this application.

Japanese HWR

Steven Buraje Poggel
steven.poggel@gmail.com

November 17, 2009

DFKI

IUI

Prof. Wolfgang Wahlster

Department of Computational Linguistics

Saarland University

Contents

1	Notes	1
2	Introduction	1
2.1	Motivation	1
2.1.1	Integrating NLP and e-learning	2
2.1.2	Another subsection with a yet unknown title	2
2.2	A CJK environment	2
2.3	Running text	3
3	Japanese Script	4
3.1	A Short History of the Japanese Writing System	4
3.2	The Japanese Writing System Today	4
3.3	Writing Japanese - typical errors	4
4	On-Line Handwriting Recognition	5
4.1	Introduction	5
4.2	The Art of Handwriting	6
4.3	Automated Recognition of Handwriting	6
4.3.1	Short History of Handwriting recognition	6
4.3.2	Recognition vs Identification	6
4.3.3	Interpretation of Handwriting	7
4.3.4	Hardware requirements	7
4.3.5	On-Line vs. Off-Line recogniton	8
4.4	A Typical On-Line HWR application	9
4.4.1	Data capturing	9
4.4.2	Preprocessing	9
4.4.3	Character Recognition	9
4.4.4	Postprocessing	9
4.5	HWR of Hanzi and Kanji	9
4.5.1	The current State-of-the-Art in Japanese and Chinese Character Recognition	10
4.5.2	Overview of a typical OJCCR system	10

5	e-learning	11
5.1	General E-Learning methods	11
5.2	E-Learning of languages	11
5.3	E-Learning of Japanese	11
5.3.1	Conceptual issues	11
5.3.2	Japanese e-learning software	11
6	Conceptual design of Kanji-Coach	12
6.1	General requirements	12
6.2	Tackling the difficulties of the Japanese script	12
6.3	Integration of HWR into learning process	12
6.4	Use cases	12
7	Technical Design of the Application	13
7.1	System Architecture	13
7.2	Framework and Devices	13
7.2.1	Pen Input Device	13
7.2.2	Desktop Computer	13
7.3	GUI	13
7.4	Communication	13
8	Handwriting Recognition Engine	14
8.1	Capturing Data	14
8.2	Data Format	14
8.3	Database	14
8.4	Recognition Architecture	15
8.5	Stroke recognition process	15
8.5.1	From point list to vectors	15
8.5.2	Handling curves	15
8.5.3	Handling all that other stuff that requires some math	15
8.6	Radical recognition process	15
8.7	Character recognition process	15
8.8	Error recognition	15
8.8.1	How to deal with typical errors when writing Japanese	15
8.9	HWR applied to e-learning of Japanese Kanji	15
8.9.1	Integration of HWR into e-learning app	15
8.9.2	Error handling	15
9	Implementation and Evaluation	16
9.1	Implementation Details	16
9.2	Evaluation of the HWR	16
9.2.1	Evaluation Metrics	16
9.3	Evaluation of E-Learning Application with Integrated HWR	16
9.4	Evaluation of the Error Hints	16
10	Conclusions	18

Abstract

Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet. Lorem dolor sit amet. Lorem ipsum dolor sit amet.

Chapter 2

Introduction

2.1 Motivation

In the history of Computational Linguistics there have been a several attempts to integrate natural language processing techniques with existing technologies. This work is one just like that. Concretely, we will try to create a handwriting recognition for Japanese Kanji. That seems interesting, because Kanji is an iconographic writing system, thus handwriting recognition (HWR) can follow different patterns than in alphabetical writing systems like latin.

Studying Japanese language is a complex task, because a new learner has to get used to a new vocabulary that - coming from a European language - has very little in common with the vocabulary of his mother tongue, unlike in European languages where quite often there are several intersections. The learner also needs to learn a new grammar system. Broadly speaking, most of the central European languages follow a subject-verb-object (SVO) structure. Japanese follows a subject-object-verb (SOV) structure therefore creating additional difficulty, comparable with German subclause structures that are a source of error for learners of German. Yet, the most notable difference for a language learner with a central European mother tongue is of course the writing system. The Japanese writing system uses three different scripts. The Kana scripts Hiragana and Katakana are syllabic, each character represents a syllable. Each syllable consists of either a vowel, a consonant and a vowel, or a consonant cluster and a vowel. Hiragana and Katakana represent roughly the same set of syllables and both have around 40-50 characters that can be modified with diacritics and thus yield additional syllable representations. Therefore, these scripts are a hurdle, but relatively unproblematic, due to their limitation in number of characters. Besides, they look quite distinct,

so there is the problem of confusing one character with another, but this is limited to a relatively short period of learning those characters.

Kanji, however, is an iconographic writing system that has around 2000 characters, which are built up of around 200 subunits called 'radicals'. So one part of the complexity lies in the number of characters. The other part of the complexity lies in the general concept of representing an idea or concept with a character instead of representing the phonemes of the spoken language with graphemes in connection with some language specific pronunciation rules. Another difficulty lies in connecting the characters with their pronunciations. Most characters have multiple pronunciations and for a language learner, studying Japanese vocabulary is a double or triple task compared to languages using a Latin or at least an alphabetic writing system. Therefore, the two tasks of learning the Kanji and studying the vocabulary together can epitomise a very high learning curve. A subordinated issue connected to that is that quite often subjectively 'simple' vocabulary comes with complex Kanji. Some e-learning applications have taken on that issue by creating a learning environment in which a learner can connect learning vocabulary with studying the Kanji.

xxx: see santosh2009: the statement of need

2.1.1 Integrating NLP and e-learning

In this project, we would like to approach the issue of studying Kanji in an e-learning application. The novelty about it is a handwriting recognition that gives the learner the ability to actually practise writing the Kanji, instead of the rather limited multiple choice recognition that most other applications use.

2.1.2 Another subsection with a yet unknown title

2.2 A CJK environment

Rather than selecting a CJK font as the main document typeface, you might want to define a CJK environment for text fragments used in the midst of a document using a normal Roman font. This allows me to say `\begin{CJK}東光\end{CJK}` to generate 東光, without putting the whole paragraph into the Far Eastern font. Or I could define a command that takes the CJK text as an argument, so that `\cjkb{北京}` produces 北京. It's that easy!

2.3 Running text

コンピューターは、本質的には数字しか扱うことができません。コンピューターは、文字や記号などのそれぞれに番号を割り振ることによって扱えるようにします。ユニコードが出来るまでは、これらの番号を割り振る仕組みが何百種類も存在しました。どの一つをとっても、十分な文字を含んではいませんでした。例えば、欧州連合一つを見ても、そのすべての言語をカバーするためには、いくつかの異なる符号化の仕組みが必要でした。英語のような一つの言語に限っても、一つだけの符号化の仕組みでは、一般的に使われるすべての文字、句読点、技術的な記号などを扱うには不十分でした。

これらの符号化の仕組みは、相互に矛盾するものでもありました。二つの異なる符号化の仕組みが、二つの異なる文字に同一の番号を付けることもできるし、同じ文字に異なる番号を付けることもできるのです。どのようなコンピューターも（特にサーバーは）多くの異なった符号化の仕組みをサポートする必要があります。たとえデータが異なる符号化の仕組みやプラットフォームを通過しても、いつどこでデータが乱れるか分からない危険を冒すこととなるのです。

Chapter 3

Japanese Script

3.1 A Short History of the Japanese Writing System

3.2 The Japanese Writing System Today

kurze erwahnung der morphologie. hiragana an verben zur konjugation.
zusammenhang verben / nomen in kanji.
uppercase / lowercase nicht vorhanden. etc.

3.3 Writing Japanese - typical errors

Chapter 4

On-Line Handwriting Recognition

4.1 Introduction

Handwriting is a very personal skill to individuals. It consists of graphical marks on a surface, it can be used to identify a person, it has the main purpose of communication. This is achieved by drawing letters or other *graphemes*, which in turn represent parts of a language. The characters have a certain basic shape, which must be recognisable for a human in order for the communication process to function. There are rules for the combination of letters, which have the ability - if known to the reader - to help recognise a character or word.

Handwriting was developed as a means of communication and to expand one's own memory. With the advent of each new technology the question arose, if handwriting was going to survive. However, the opposite seems to be the truth: For example, the printing press increased the number of documents available and therefore increased the number of people who learned to read and write. Through the increased rate of alphabetisation, naturally there was an increased use of handwriting as a means of communication.

In various situations handwriting seems much more practical than typing on a keyboard. For instance children at school are using notepads and pencils or ink pens, which are regarded as a better tool to teach writing by German teachers. Therefore it can be concluded that there is little danger of the extinction of handwriting as a communication tool. In fact, as the length of handwritten messages decreases, the number of people using handwriting increases (Plamondon and Srihari 2000).

4.2 The Art of Handwriting

xxx: see tappert1990 IV handwriting properties and recognition problems
xxx: was ist das generelle problem? xxx: welche probleme treten dabei auf? xxx: related problems: see liujaegernakagawa2004 1.1 xxx: write something about the problems of different scripts here, too, as in tappert1990 IV: problems.

4.3 Automated Recognition of Handwriting

4.3.1 Short History of Handwriting recognition

Handwriting recognition (HWR) as a technological discipline performed by machines has been around for many years. The quality of the systems recognising handwriting has improved over the decades. It is the key technology to pen-based computer systems. The first research papers concerned with *pattern recognition* on computers were published in the late 1950ies, *Handwriting recognition* as an individual subject in the early 1960ies. (Goldberg 1915) describes in a US Patent a machine that can recognise alphanumeric characters as early as 1915. However, despite the surprise how early such a device was invented, it should be taken into consideration that that was before the times of modern computers, therefore the methods he employs are quite different from the algorithms used after the advent of computers, more concretely, computers with screens.

(Tappert et al. 1990) describe in their review the development of handwriting recognition, which was a popular research topic in the early 1970ies and then again in the 1980ies, due to the increased availability of pen-input devices. Generally speaking, handwriting recognition (HWR) involves automatic conversion of handwritten text into a machine readable character encoding like ASCII or UTF-8. Typical HWR-environments include a pen or stylus that is used for the handwriting, a touch-sensitive surface, which the user writes on and an application that interprets the strokes of the stylus on the surface and converts them into digital text. Usually, the writing surface captures the x-y coordinates of the stylus movement.

4.3.2 Recognition vs Identification

Handwriting recognition is the task of transforming a spatial language representation into a symbolic representation. In the English language (and many others) the symbolic representation is typically 8-bit ASCII. However, with *Unicode* being around for more than a decade now, storage space on harddisks not being as much of an issue any more and *RAM*

being readily available to the Gigabytes, it has become more common to use a *UTF-8* encoding, which is a variable-length character encoding for Unicode (The Unicode Consortium 2000). Akin disciplines to handwriting recognition are *handwriting identification*, which is the task of identifying the author of a handwritten text sample from a set of writers, assuming that each handwriting style can be seen as individual to the person who wrote it. The task of *signature verification* is to determine if a given signature stems from the person who's name is given in the signature. Thus, handwriting identification and verification can be used for analysis in the field of jurisdiction. They determine the individual features of a handwritten sample of a specific writer and compare those to samples given by a different or the same writer. By analysing those features one can find out if a piece of handwritten text is authentic or not.

4.3.3 Interpretation of Handwriting

Handwriting recognition and interpretation are trying to filter out the writer-specific variations and extract the text message only. This conversion process can be a hard task, even for a human. Humans use context knowledge in order to determine the likeliness of a certain message in a certain context. For instance, a handwritten message on a shopping list that could be read as *bread* or *broad* due to the similarities of the characters for 'e' and 'o' in some cursive handwriting styles, will be interpreted as *bread*, since it is a much more likely interpretation in the shopping list domain. However, if the next word on the shopping list is *beans*, the likelihood for the interpretation of the first word as *broad* rises, because the collocation *broad beans* is a sequence that is likely on a shopping list, at least more likely than having the interpretation *bread* and then *beans* without a clear separation between the two. Even with non-handwritten, but printed characters, the human mind can be tricked because of the brain's ability to perform these interpretations within milliseconds without conscious thinking. An example of that are modern T-Shirt inscriptions that state things like *Pozilei* in a white font on a green ground (the German police colours in most federal states are green and white), which German native speakers usually read as *Polizei* (police), because that is the most likely interpretation.

4.3.4 Hardware requirements

xxx: see santosh2009 basic tools / techniques / digitizer technology,

Several different hardware commercial products are available in order to capture the x-y coordinates of a stylus or pen. Graphics tablet like

the products of the Wacom Co., Ltd.¹ are popular input devices for hand motions and hand gestures. The use of pen-like input devices has also been recommended, since 42% of mouse users report feelings of weakness, stiffness and general discomfort in the wrist and hand when using the mouse for long periods (Woods et al. 2002). Moreover there are PDAs and Tablet PCs, where the writing surface serves as an output device, i.e. an display at the same time. New generation mobile phones also contain touch-displays, but for those it is more common to be operated without a stylus. Those devices interpret user gestures, however the input is given directly with the users fingers. Another rather new development are real-ink digital pens. With those, a user can write on paper with real ink, and the pen stores the movements of the pen-tip on the paper. The movements are transferred to a computer later. It can be expected that with technologies like Bluetooth it may be possible to transfer those data in real-time, not delayed.

4.3.5 On-Line vs. Off-Line recognition

xxx: see plamondon2000 1.5. blend systems! xxx: see santosh2009: off-line / on-line chapter

On-line HWR means that the input is converted in *real-time*, *dynamically*, while the user is writing. This recognition can lag behind the user's writing speed. (Tappert et al. 1990) report average writing rates of 1.5-2.5 characters/s for English alphanumerics or 0.2-2.5 characters/s for Chinese characters. In online systems, the data usually comes in as a sequence of coordinate points.

Off-line HWR is the application of a HWR algorithm after the writing. It can be performed at any time after the writing has been completed. That includes recognition of data transferred from the real-ink pens (see 4.3.4) to a computing device after the writing has been completed. The standard case of off-line HWR, however, is a subset of optical character recognition (OCR). A scanner transfers the physical image on paper into a bitmap, the character recognition is performed on the bitmap. An OCR system can recognise several hundred characters per second.

On-line devices have the dynamic information of the writing, since each point coordinate is captured at a specific point of time. Also, the system knows the input stroke sequence, their direction and speed of writing. All these information can be an advantage for an on-line system, however, off-line systems have used algorithms of line-thinning, such that the data consists of point coordinates, similar to the input of online systems (Tappert et al. 1990).

¹www.wacom.com

4.4 A Typical On-Line HWR application

A typical HWR application has several parts that follow up on each other in a procedural fashion.

- **Data capturing:** The data is captured through an input device like a writing surface and a stylus.
- **Preprocessing:** The data is segmented, noise reduction like smoothing and filtering are applied.
- **Character Recognition:** Feature analysis, stroke matching, time, direction and curve matching.

4.4.1 Data capturing

how is the data captured? what format? hardware? xxx: see plamondon2000 1.4. xxx: see santosh2009 sampling

4.4.2 Preprocessing

xxx: see santosh2009 pre-processing xxx: see tappert1990 preprocessing: segmentation, noise reduction. xxx: see santosh2009 noise elimination xxx: see santosh2009 normalization xxx: see santosh2009 repetition removal

4.4.3 Character Recognition

xxx: see tappert1990 VI shape recognition. xxx: see plamondon2000: 3.1.1 different models xxx: see all the substroke stuff, santosh, shimodaira2003, nakai2003: very short, properly in OLCCR

4.4.4 Postprocessing

xxx: what happens after the recognition process? xxx: see tappert1990 again in postprocessing chapter.

4.5 HWR of Hanzi and Kanji

- **Warum:** Um einen Ueberblick ueber HWR-Techniken fuer Japanische Schriftzeichen und verschiedene Herangehensweisen zu verschaffen.
- **Nutzen:** Leser kann sich ein Bild darueber verschaffen, in welchem Kontext sich die Applikation bewegt.

- Was: research different approaches, see what the focus on, what their specialty is and report about them. Take different specialist papers and compare them.
- Wie: Wiss. Report. / Zusammenfassung. Vergleich.

4.5.1 The current State-of-the-Art in Japanese and Chinese Character Recognition

From the 1990s onwards, On-Line Japanese and Chinese Character Recognition (OJCCR) systems have been aiming at loosening the restrictions imposed on the writer when using an OJCCR system. Their focus shifted from recognition of block style script ('regular' script) to fluent style script, which is also called 'cursive' style. Accuracies of up to about 95% are achieved in the different systems.

xxx: bla. says the opposite. (Chen and Lee 1996) oder auch xxx: (Nakagawa, Tokuno, Zhu, Onuma, Oda, and Kitadai 2008) und (Nakai, Shimodaira, and Sagayama 2003) xxx: zu guter letzt: (Santosh and Nattee 2009)

4.5.2 Overview of a typical OJCCR system

xxx: (Liu, Jaeger, and Nakagawa 2004) have said:

xxx: graphic: handwritten -> character segmentation -> ... -> character codes see fig. 3 of liujaegernakagawa2004

Broadly speaking, from an abstract viewpoint, typical handwriting recognition systems for Chinese and Japanese characters have the same structure like the systems for latin-based alphabets. The process begins with *Character segmentation*, goes on with *Preprocessing*, *Pattern description*, *Pattern recognition* and ends with *Contextual processing*, if applicable. However, there are differences to the standard process, due to the nature of the Chinese characters (see ??). Especially the pattern representation is divers in the different OJCCR systems, whereas it is naturally more alike in the systems focussing on latin characters. This is due to the fact that the latin alphabet is rather small, but has more variation concerning writing style, whereas the Chinese alphabet has a larger inventory of characters, but less variation in how to write a character - at least - it is widely agreed upon a 'proper' stroke sequence for a character, even across country borders.

xxx: liujaeger2004: 4.1. structural representation. statistical representation.

xxx: liujaeger2004: character classification xxx: liujaeger2004: very short: contextual processing.

Chapter 5

e-learning

5.1 General E-Learning methods

5.2 E-Learning of languages

in section

5.3 E-Learning of Japanese

5.3.1 Conceptual issues

5.3.2 Japanese e-learning software

put all your bashing and criticism here

Chapter 6

Conceptual design of Kanji-Coach

6.1 General requirements

6.2 Tackling the difficulties of the Japanese script

6.3 Integration of HWR into learning process

6.4 Use cases

Chapter 7

Technical Design of the Application

7.1 System Architecture

7.2 Framework and Devices

7.2.1 Pen Input Device

7.2.2 Desktop Computer

7.3 GUI

7.4 Communication

Chapter 8

Handwriting Recognition Engine

The sections of this chapter are more the result of a brainstorming than a proper thought-through chapter design. xxx: see santosh2009 for mathematical stuff: nice description of what I'm doing

8.1 Capturing Data

this should deal with how the data are captured during the process mouse coordinates and stuff

8.2 Data Format

how is the data structured? radicals, strokes, characters, xml-format

8.3 Database

where did I get it from? how many chars are in there? how are they accessible? what format?

8.4 Recognition Architecture

8.5 Stroke recognition process

8.5.1 From point list to vectors

8.5.2 Handling curves

8.5.3 Handling all that other stuff that requires some math

8.6 Radical recognition process

8.7 Character recognition process

8.8 Error recognition

8.8.1 How to deal with typical errors when writing Japanese

Error recognition

focus on technical aspects

Error handling

focus on technical aspects

8.9 HWR applied to e-learning of Japanese Kanji

8.9.1 Integration of HWR into e-learning app

educational aspects / the e-learning view

8.9.2 Error handling

educational aspects / the e-learning view

Chapter 9

Implementation and Evaluation

9.1 Implementation Details

Pointer auf CD und auf Appendix mit Beispielinteraktionen (diese mit Foto). Screenshots. Zahlen zur Erkennung - z.B. wie lange dauert es, ein zeichen zu erkennen?

9.2 Evaluation of the HWR

9.2.1 Evaluation Metrics

evaluation method: counting precision and recall section about precision and recall - the odd numbers. how can that be done honest and useful? how can I get meaningful numbers at all?

9.3 Evaluation of E-Learning Application with Integrated HWR

qualitative auswertung, keine zahlen, sondern fragebogen. fuehlt der lerner sich unterstuetzt? glaubt er, dass es schneller geht als ohne HWR? besser als auf papier?

9.4 Evaluation of the Error Hints

use cases, inwieweit helfen die fehlerhinweise? geht das lernen dann wirklich schneller? wie laesst sich der mehrwert bewerten? system kann sagen:

wo liegt die verwechslung? warum war das falsch?

Chapter 10

Conclusions

Chapter 11

Outlook

Was kann man hiermit fuer chinesische / koreanische machen? HWR als Service fuer andere Applicationen.

References

- Chen, J.-W. and S.-Y. Lee (1996). A Hierarchical Representation for the Reference Database of On-Line Chinese Character Recognition. In *Advances in Structural and Syntactical Pattern Recognition*, Volume 1121 of *Lecture Notes in Computer Science*, pp. 351--360. Berlin/Heidelberg, Germany: Springer.
- Goldberg, H. E. (1915, December). Controller. *United States Patent 1,116,663*.
- Liu, C.-L., S. Jaeger, and M. Nakagawa (2004). Online Recognition of Chinese Characters: The State-of-the-Art. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 26, 198--213.
- Nakagawa, M., J. Tokuno, B. Zhu, M. Onuma, H. Oda, and A. Kitadai (2008). Recent Results of Online Japanese Handwriting Recognition and Its Applications. In D. Doermann and S. Jaeger (Eds.), *Arabic and Chinese Handwriting Recognition*, Volume 4768 of *Lecture Notes in Computer Science*, pp. 170--195. Berlin/Heidelberg, Germany: Springer.
- Nakai, M., H. Shimodaira, and S. Sagayama (2003). Generation of Hierarchical Dictionary for Stroke-Order Free Kanji Handwriting Recognition Based on Substroke HMM. In *Proc. Seventh Int'l Conf. Document Analysis and Recognition*, pp. 514--518.
- Plamondon, R. and S. N. Srihari (2000). On-line and off-line handwriting recognition: A comprehensive survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22(1), 63--84.
- Santosh, K. and C. Nattee (2009). A Comprehensive Survey on On-Line Handwriting Recognition Technology and its Real Application to the Nepalese Natural Handwriting. *Kathmandu University Journal of Science, Engineering and Technology* 6(I), 30--54.
- Tappert, C. C., C. Y. Suen, and T. Wakahara (1990). The State of the Art in Online Handwriting Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12(8), 787--808.

- The Unicode Consortium (Ed.) (2000). *The Unicode Standard. Version 3.0*. Addison-Wesley.
- Woods, V., S. Hastings, P. Buckle, and R. Haslam (2002). *Ergonomics of using a mouse or other non-keyboard input device*, Chapter 3, pp. 23. Number 045 in HSE research report. London: Health and Safety Executive.