# An Accessible Coded Input Method for Japanese Extensive Writing

**Takeshi Okadome**  **Junya Nakajima**  **Sho Ito**  **Koh Kakusho**

Kwansei Gakuin University/ 2-1 Gakuen, Sanda-shi, Hyogo-pref. 669-1337 Japan.

tokadome@acm.org  cmj88416@kwansei.ac.jp  baj886106@kwansei.ac.jp  kakusho@kwansei.ac.jp

## Abstract

The orthogonal code, *O-code*, proposed here is an asscessible coded input method for Japanese text typing. To a kana, it assigns a romaji (literally "Roman letter") sequence corresponding to the kana and, to each of the 458 most-frequently-used kanji characters, it also assigns a successive two key stroke code which differs from any romaji representation for a kana. With the standard English keyboard or virtual one on a tabletop, the features of the method enable a novice user to input Japanese text using the kana-to-kanji conversion by inputting the corresponding romaji for a kana and a well-trained user to type text fast with less fatigue.

## 1 Introduction

Widespread mobile phones and tabletops have brought us a chance of using various methods for text input not through a physical keyboard. In particular, smart phones use interactive and/or stroke-based soft "keyboard" on a touch screen. Imagine, however, when you must input a vast amount of text on a small touch screen with one or two of the fingers. Then you will soon need a physical standard keyboard or virtual one on a tabletop. (For example, see Lewis, LaLomia, and Kennedy (1999) and Rick (2010) for virtual keyboards.)

On the other hand, some of the methods adopt even different modes such as speech, handwriting, and gesture recognition (Kölsch and Turk, 2002). They have, however, not achieved good performance in practice.

The situations hold for Japanese text input. This paper discusses input methods for a great amount of Japanese text such as academic and research articles, novels, and long blogs. It introduces a Japanese input method that, with the standard English keyboard or virtual one, enables a novice user to input Japanese text and a well-trained user to type text fast with less fatigue.

## 2 Japanese Input Methods: Brief Review

### 2.1 Japanese writing system

Almost all normal Japanese writing is a mixture of kanjis, kanas, and others. Kanjis are Chinese characters. (For detail of Japanese writing system and discussion on a variety of Japanese typing methods, see Yamada (1983).) The largest kanji dictionary today lists more than 50,000 distinct kanjis. They are used for nouns and the verb roots, adjectives, adverbs, and the like. The set of kanas consists of two subsets of syllabaries, hiraganas and katakanas, each of which in turn is made up of about 80 letters. Hiraganas are used to write inflections and other grammatical parts of sentences and katakana are used for the transcription of foreign words. Many different kanjis that have different meanings may have the same reading. That is, they are represented by the same kana strings. This is called *homophone problem.*

### 2.2 Kana-to-kanji conversion

The complexity of Japanese writing system has led to the development of a variety of Japanese typing methods. Among the typing methods, kana-to-kanji conversion with the romaji (literally "Roman letters") input for the kanas is the most popular. If all kanjis are coded exactly as they are pronounced, and if all homophone problems are solved by syntactic and/or semantic analyses, then kana-to-kanji conversion might to be best solution for the input problem.

### 2.3 Coded input method

Unfortunately, too many obstacles prevent kana-to-kanji conversion from being an ideal input method. The homophone problem is the most severe one and, in kana-to-kanji conversion systems, all homophonic kanjis are displayed at once

31

and the user must select the correct one by an appropriate designation such as keying or pointing. This interactive feature of kana-to-kanji conversion prevents us from touch typing and makes us irritated in Japanese typing task.

Although kanji dictionaries lists many distinct kanjis, for ordinary use, we Japanese deal with only about 1,000 and the size of the necessary kanji set increases, up to perhaps 1,500 or so for an office of 10 people (Yamada, 1983). Focussing on the ordinary usage of kanjis, some researchers have developed coded input methods in which each of the kanas and frequently-used kanjis is coded in successive two (or three) key strokes on the standard English keyboard. Unlike kana-to-kanji conversion, a coded input method such as the *T-code* (Hiraga, Ono, and Yamada, 1980) and the *TUT-code* (Ohiwa, Takashima, and Mitsui, 1983) permit users to input Japanese texts in touch typing after users practice and acquire the codes. The problem in coded input methods is that it takes us 400 to 1,000 hours to acquire them.

An attepmt to combine a coded input method with ordinary kana-to-kanji conversion enables a novice user to input Japanese text if he/she at least learns the codes for the kanas. In particular, a coded input method with an enhancement of kana-to-kanji conversion called *kanji-kana mixture conversion* partly solves the homophone problem because, in the method, input strings may be partially represented in kanji, which reduces the number of homophones. (For kanji-kana mixture conversion, see Ono (1990).)

This paper introduces a new coded input method named the orthogonal code, *O-code*, in which, to each kana, a romaji sequence corresponding to the kana is assigned and, to each of the 458 most-frequently-used kanji characters, a successive two key stroke code which differs from any romaji representation for a kana is assigned. With the standard English keyboard, the features of the *O-code* enable a novice user to input Japanese text using the kana-to-kanji conversion by inputting the corresponding romaji (Roman letter) for a kana and a well-trained user to type text fast with less fatigue.

Today people ordinarily use kana-kanji conversion systems for inputting Japanese text. Because many of them can touch-type the romajis for the kanas, the *O-code* with kana-kanji conversion is more accessible than the other coded input methods in which even the kanas are coded in meaningless and non-associative two (or three) key-strokes.

## 3 Design Principles

### 3.1 Requirements

The design requirements of the *O-code* are its "accessibility" and "efficiency." That is, a code system are required to embody the following features:

1. Novice users who have learned English typing can easily use for Japanese text inputs together with kana-to-kanji conversion.

2. A high level of input speed is attainable.

3. Users suffer from less "fatigue." The fatigueness is indirectly measured by counting interactive choice of kanjis that are displayed on the screen, for example.

4. A high rate of accuracy may be maintained. That is, the codes should not be susceptible to erroneous finger motions.

The last three of them coincide with those in the *T-code* (Hiraga, Ono, Yamada, 1980).

These factors are closely correlated to each other. Of them, barrier-free (that is, less interactive) and speed are in a way the most decisive, and also the most appealing factors although they are somewhat contradictory.

A system design would be satisfied with a multi-stroke code system on the standard English keyboard, where we type a kana by the romaji input method and each of the most-frequently-used kanjis by two successive key strokes which differ from any romaji representation for a kana. This immediately implies that the system is usable by an untrained user and, for a trained user, he or she can type fast with less fatigue because of less interactive choices of kanjis.

### 3.2 Design

The *O-code* basically adopts the same strategy of the kanji-set selection and its assignment to two-stroke codes as the *T-code* does, because (1) the design policy of the *T-code* satisfies requirements 2-4 and (2) there are several multi-stroke systems currently implemented and put to use, and among them, the most famous one is the *T-code* that has the following specifications (Yamada, 1983; Hiraga, Ono, and Yamada, 1980).

1. Only 40 out of 48 printing key are used.

2. The shift key and space bar are not involved in character codes.

3. Out of 40 × 40 = 1,600 possible two-stroke pairs, some of those that involve the top-row keys are not used. In particular, the same-hand pairs involving the top-row keys are almost all unused. This gave us about 1,200 usable pairs, which should cover more than 95 percent of kanjis, or more than 98 percent of the total characters, in ordinary text. The coverage is almost 100 percent if the code set is tuned to the task of a specific group or individual.

4. There seem to be no scientific data indicating which is the better hand for starting alternate-hand stroking. The *T-code* chose to start with the right (and stop with the left); that is, the most frequently used characters are coded in RL pairs, then in RR pairs, and finally LR pairs, which makes the expected length of the alternating-hand sequences largest.

5. The actual assignment of codes to characters is made to optimize various parameters.

### 3.3 Design restrictions

Unlike the *T-code*, the *O-code* uses only 30 printing keys that are not the top-row keys. Out of 30 × 30 = 900 two-stroke pairs, only 458 pairs are used for the kanji coding. This is because,

1. to kanjis, the *O-code* cannot assign the two-stroke codes that match the first two letters in a romaji representing a kana or

2. it also cannot use those in which the first stroke is for a vowel: 'a,' 'e,' 'i,' 'o,' or 'u,' and a punctuation.

## 4 The Determination of the Kanji Set

The kanji set selection is based upon the frequencies of the usage of kanjis in texts on the web. For this, we use the 1-gram data in the Web Japanese N-gram data in the set of about 200,000,000 sentences that Google collected from the web and analysed (Kudo and Kazawa, 2007). The top ten kanjis are

1. 

2. 

3. 

4. 

5. 

6. 

7. 

8. 

9. 

10.   .

Considering a statistical fact that about 500 kanji and 150 kana characters are used by an average person for his daily use (Hiraga, Ono, and Yamada, 1980), although the set may change gradually, our selection of the kanji set seems reasonable.

## 5 Coding of the Kanji Set

The kanji coding is based on the efficiency of finger movements just like *T-code*. The coding method maps kanjis, arranged in the order of the frequency of usage, to the key pairs arranged in a suitable ordering as defined below. Second order adjustments will be made afterwards.

The method obtains the ordering of key pairs by assigning certain weights to certain characteristics of hand motions and using their linear sum for each key pair. It gives a larger weight to characteristics that seems to have a greater importance. The assignment will bring us the ordering of 458 key pairs on a 30-key keyboard.

This ordering, however, is not immediately usable to fix the assignment of characters directly, because the typing process is not a collection of isolated key pairs, but their continuous sequence. If, for example, key pair $(h, f)$ is with a high score, then $(f, h)$, its reverse, would also be with a high score, but the frequent appearances of these two key pairs would result in the frequent tapping motion of key pairs $(h, h)$ and $(f, f)$ in the interval of consecutive $(h, h)$s and $(f, f)$s, or vice versa, which are known to be less preferred. This would also be adverse to alternate hand stroking as well.

Through such considerations, Hiraga, Ono, and Yamada (1980) concludes that the desirable keyboard characteristics may be itemized as follows:

1. The whole typing procedure is to keep as much keying rhythm as possible. Fluent rhythm, as well as high average of typing speed, is best realized by alternate stroking by both hands. Thus, it would be our principal objective to let the code system be such that it would allow alternate hand stroking as much as possible.

2. Hands should not be moving up and down incessantly on key rows, but stay in the same row as much. Thus, strokes on the home row should be used as much as possible and excursions to other rows should be held minimum. Comparing between the upper and the bottom rows, all evidences point out that hands are more fluent on the upper row, so the ranking of rows should be in the preference order of the home, the upper, and the bottom.

3. Fingers should be loaded in proportion to their dexterity. In typing motions, fingers are divided into the stronger ones (index and middle fingers) and the weaker ones (ring and little fingers). Index and middle fingers are not so much different in their capacity and functions. However, we must keep in mind that each index finger must cover two inner columns. The difference between ring and little fingers is also not so obvious. Although a ring finger is superior in its stroking force in typing motions, a little finger may have the advantage of the twisting motion of the wrist. (In the *T-code*, little fingers will be given more emphasis than the ring.)

4. The number of awkward keying sequences must be decreased as much as possible. Almost all awkward key pairs are of one-handed stroking, again attesting to the importance of alternate hand stroking. The major awkward key pair sequences, in the order of their disadvantages are:

   (a) Hurdling: the stroking from the upper to the bottom row vice versa, jumping over the home row.
   (b) Reaching: the stroking of different keys with the same finger.
   (c) Tapping: the stroking of the same key.
   (d) Rocking: stroking with adjacent fingers, especially from an inner to an outer one.

The *T-code* weighting process starts by accommodating for condition (1) above. The key pairs are divided into 4 blocks, namely RL, RR, LL, and LR blocks, where symbols L and R stand for the hands that stroke the keys of the pair. The blocks are given preference in the order above, and key pairs in each of the blocks are then ordered by taking further conditions into account. Within individual blocks, conditions (2), (3), …, are evaluated and weighed accordingly, and the whole ordering is decided. Awkward sequences are deliberately given negative weights in order to bring down their ranking, thus decreasing their occurrences when the codes are used.

In the *O-code*, we assgin a key pair in the order to a kanji in the occurrence-frequency order in the selected kanji set described in the previous section unless the pair matches the first two letters in a romaji representation for a kana. The entire code table for the kanji set is given in the appendix.

## 6  Evaluation

The code system test would be to actually measure its typing speed and error rate on a real system. The test, however, requires much time and cost.

Hence, we show some statistical figures about *O-code* derived, again, from the Web Japanese N-gram data in the set of about 200,000,000 sentences that Google collected from the web and analysed (Kudo and Kazawa, 2007).

Figures 1 and 2 illustrate the finger loading and row distributions of the *O-code*. From the figures, we see for our code that:

1. Hands are evenly loaded, slightly lighter for the weaker left hand.

2. About 47 and 35 percents of strokes fall on the upper and home rows.

3. The loading of fingers is slightly different from the conjectured strengths of the fingers. In particular, the little finger loading of the left hand is relatively higher; Also the ring finger loading of the right hand is relatively higher.

In the *T-code*, about 24 and 56 percents of strokes fall on the upper and home rows and, in the *TUT-code*, about 34 and 56 percents of strokes do. Furthermore, in both the *T-code* and *TUT-code*, the loading of fingers is in a qualitative agreement with the conjectured strengths of the fingers. The
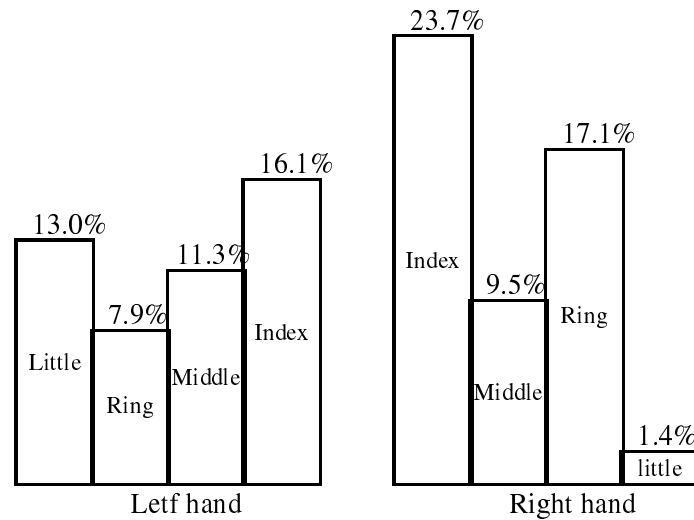
Figure 1: Finger loading distributions

results of these analyses may lead us to a conclusion that the *O-code* is less efficient than the *T-code* and *TUT-code*. This is because four out of the five vowels are on the upper row and vowel 'a' is typed by the little finger in the standard English keyboard.

## 7 Concluding Remarks

This paper proposes a new coded input method named the orthogonal code, *O-code*, for the input of Japanese texts. In the *O-code*, to each kana, a romaji sequence corresponding to the kana is assigned and, to each of the 458 most-frequently-used kanji characters, a successive two key stroke code which differs from any romaji representation for a kana is assigned. With the standard English keyboard, the features of the *O-code* enable a novice user to input Japanese text using the kana-to-kanji conversion by inputting the corresponding romaji for a kana and a well-trained user to type text fast with less fatigue. The principal goal has been to realize an input system that would allow a high degree of touch typing. The *O-code* is semi-optimal for experts, but it provides their easier accessibility by beginners.

We have implemented a *O-code* system using the Kanchoku Win system (Kanchoku Win, 2006) that enables us to input Japanese texts by a coded input method on the Windows. We have also developed the *O-code$_{DSK}$* for the Dvorak simplified keyboard (Dvorak, 1943). The *O-code$_{DSK}$* is more efficient that the *O-code*, because the five vowels are on the left-hand home row and

the high-frequency consonants on the right-hand home row. Furthemore, the restriction of the romaji system into the kunrei system that is an optimal romaji system permits us to encode more kanjis in both the *O-code* and the *O-code$_{DSK}$*.

Furthermore, we have constructed about 1,500 sentences for training the *O-code* using the eelll/JS (eelll/JS, 2005) that is a system for practicing a coded input method on a web browser.

### Acknowledgments

### References

1. Dvorak, A. (1943). There is a better typewriter keyboard. *National Business Education Quarterly*, 12, 51-58;66.
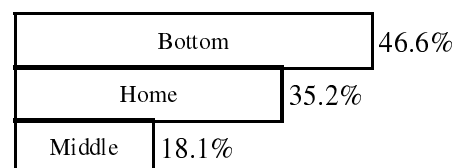
2. eeel/JS (2005). http://www.sato.kuis.kyoto-u.ac.jp/~yuse/tcode/eljs/ (accessed on 13 May 2010).



Figure 2: Row distributions

3. Kanchoku Win (2006). `http://www.sato.kuis.kyoto-u.ac.jp/\~yuse/tcode/kw/` (accessed on 13 May 2010).

4. Kölsch, M. and M. Turk (2002). Keyboards without keyboards: a survey of virtual keyboards. *Technical Report 2002-21*, UCSB, Santa Barbara, CA.

5. Kudo, T. and H. Kazawa (2007). *Web Japanese N-gram Version 1*, Gengo Shigen Kyokai.

6. Hiraga, Y., Y. Ono, and H. Yamada (1980). An assignment of key-codes for a Japanese character keyboard. *Proceedings of the 8th International Conference on Computational Linguistics*, 249-256.

7. MacKenzie, I. S., S. X. Zhang, and R. W. Soukoreff (1999). Text entry using soft keyboards. *Behaviour and Information Technology*, 18(4), 235-244.

8. Rick, J. (2010). Performance optimizations of virtual keyboards for stroke-based text entry on a touch-based tabletop. *Proceedings of the Twenty-Third Annual ACM Symposium on User Interface Software and Technology (UIST'10)* 77-86.

9. Ohiwa, H., T. Takashima, O. Mitsui (1983). A Method of touch-typing Japanese text. *IPSJ Journal*, 24(6), 772-779. (In Japanese)

10. Ono, Y. (1990). Auxiliary input methods for T-Code system : a kanzi-form combination and a kanzi-mixed conversion. *IPSJ Journal*, 31(3), 404-414. (In Japanese)

11. Yamada, H. (1983). Certain problems associated with the design of input keyboards for Japanese writing. In: *Cognitive Aspects of Skilled Typewriting*, edited by W. E. Cooper, 305-407, Springer-Verlag.

**Appendix:** *O-code* **table**

Total 458 Kanjis

| | | | | |
|---|---|---|---|---|
| kd | jd | kf | jf | hd |
| kg | hf | jg | hg | ks |
| ld | js | lf | hs | lg |
| ls | kr | jr | hr | jt |
| ht | kq | jq | hq | kw |
| jw | lr | hw | lt | lq |
| lw | yd | yf | yg | pd |
| pf | pg | ys | ps | yr |
| yt | pr | yq | pt | pq |
| yw | pw | kc | jc | kv |
| jv | hc | kb | hv | jb |
| hb | kz | jz | hz | kx |
| lc | jx | lv | hx | lb |
| lz | lx | md | mf | mg |
| mc | mv | mb | mz | mx |
| yc | yv | yb | pc | pv |
| yz | pb | pz | yx | px |
| mr | mt | mq | mw | jk |
| kj | hk | kh | hj | jh |
| k; | j; | h; | kl | lk |
| jl | lj | hl | lh | l; |
| kp | jp | hp | lp | yk |
| yj | yh | pk | pj | y; |
| ph | p; | yl | pl | yp |
| k, | j, | km | jm | h, |
| kn | hm | jn | hn | k/ |
| j/ | h/ | k. | l, | j. |
| lm | h. | ln | l/ | l. |
| mk | mj | mh | m; | ml |
| m, | m/ | m. | y, | ym |
| yn | p, | pm | y/ | pn |
| p/ | y. | p. | mp | df |
| fd | dg | gd | fg | gf |
| sd | sf | fs | sg | gs |
| dr | fr | dt | ft | gr |
| gt | dq | fq | gq | dw |
| fw | sq | sw | rd | rf |
| td | rg | tf | tg | qd |
| qf | qg | qa | wd | wf |
| rs | wg | qs | ws | rt |
| tr | qe | qr | rq | qt |
| tq | wr | rw | wt | tw |
| qw | wq | dc | dv | fc |
| fv | db | gc | fb | gv |
| gb | dz | fz | gz | sc |
| dx | sv | fx | sb | gx |
| sz | sx | cd | cf | vd |
| vf | cg | bd | vg | bf |
| bg | zd | zf | zg | xd |
| cs | xf | vs | xg | bs |

| | | | | |
|---|---|---|---|---|
| zs | xs | cv | vc | cb |
| bc | vb | bv | zc | cz |
| zv | vz | zb | bz | xc |
| cx | xv | vx | xb | bx |
| zx | xz | rc | rv | tc |
| rb | tv | tb | qc | qv |
| rz | qb | tz | qz | wc |
| wv | rx | wb | tx | qx |
| wz | wx | ce | cr | vr |
| ct | vt | br | bt | cq |
| zr | vq | zt | bq | zq |
| cw | xr | vw | bw | zw |
| xq | dk | dj | fk | fj |
| gk | fh | gj | gh | d; |
| f; | g; | sk | dl | sj |
| fl | gl | s; | sl | fy |
| dp | fp | gp | sp | rk |
| rj | tk | rh | tj | qk |
| qj | r; | qh | t; | q; |
| wk | wj | rl | wh | tl |
| ql | w; | wl | qi | qu |
| rp | qy | tp | qp | wy |
| qo | wp | d, | dm | f, |
| fm | dn | g, | fn | gm |
| gn | d/ | f/ | g/ | s, |
| d. | sm | f. | sn | g. |
| s/ | s. | ck | cj | vk |
| vj | bk | vh | bj | bh |
| zk | c; | zj | v; | zh |
| b; | z; | xk | cl | xj |
| vl | xh | bl | zl | x; |
| xl | c, | cm | v, | vm |
| cn | b, | vn | bm | bn |
| z, | c/ | zm | v/ | zn |
| b/ | z/ | x, | c. | xm |
| v. | xn | b. | z. | x/ |
| x. | r, | rm | t, | rn |
| tm | tn | q, | qm | r/ |
| qn | t/ | q/ | w, | wm |
| r. | wn | t. | q. | w/ |
| w. | ci | vy | cp | vp |
| bp | zp | xp | | |