**Date**: 2008-03-21

---

**ISO/IEC JTC1/SC2/WG2**
**Coded Character Set**
**Secretariat: Japan (JISC)**

---

**Doc. Type:** Input to ISO/IEC 10646

**Title:** Representation of CJK Unified Ideographs in multi-column

**Source:** **Michel Suignard, project editor**
**Project:** JTC1 02.10646
**Status:** For review by WG2
**Date:** 2008-03-21
**Distribution:** WG2
**Reference:**
**Medium:**

---

The need to publish a new edition of ISO/IEC 10646 creates the requirement to republish all code charts containing CJK Unified Ideographs. However, sources have grown from eight to nine with the addition of the M source (Macao SAR) since the last edition: ISO/IEC 10646:2003. In addition, even in the last edition, three sources: KP (Democratic People Republic of Korea), H (Hong-Kong SAR), and U (UTC) were not adequately represented in the multi-column charts. However, is not practical to show all nine sources for all CJK characters because no character has all nine sources, and very few characters go beyond 6 sources. Furthermore, the CJK extensions (A, B, and C) have even less sources, calling for a more optimal format based on average number of sources for a given collection.

In fact, if one looks at the various CJK Unified collections, the number optimal of sources to be shown varies from 6 (CJK Unified Ideograph: 4E00-9FC5), 4 (CJK Unified Ideographs Extension A 3400-4DB5), and 2 (CJK Unified Ideographs Extension B: 20000-2A6D6, and proposed CJK Unified Ideographs extension C: 2A700-2B734).

The following pages show some examples of what could be done as a compromise between size and legibility. The main focus is to provide a more consistent amount of information per page, thus allowing a variable amount of columns per characters depending on the collections. Because the standard is typically not printed in its entirety anymore, paper preservation is not as critical; this in return allows some flexibility. In addition, some additional information is provided such as the radical image and the 'radical, stroke' count for each ideograph.

There are still some issues worth considering:

- For the main CJK Unified Ideographs collection (4E00-9FC5), should we preserve spatial positioning for the sources, therefore requiring frequently 2 rows when a source is not part of the first six sources, or should we show these additional sources highlighted in the same row? (The author prefers the second solution)
- Can we get away with the spatial positioning for all other collections, allowing a compression to 4 or even 2 sources per characters (with exceptions dealt through additional rows)? (The author is in favor)
- Some G sources (such as KangXi: G_KX) have very long labels. This can be accommodated by shortening to just the prefix (like 'G_KX') or allowing the first source (always G when existing) to bleed over the character

1

code position and pushing the radical.count info upward. Note that the source representation of the Japanese ARIB Kanji characters may have to be shortened to 9 characters instead of the current 11 to fit in the sub-column box. (The author prefers the second solution, that is preserving the information by tweaking the layout)



41D5
立 117.14
G_KX087321  T3-5A68  K3-2E29  KP1-63A4

- There are multiple 'radical,stroke' data sets. One has to be selected. The Unicode Unihan database contains such a data set. Ideally, the information should be added in the CJKU_SR.txt file.

The following pages show some samples:

- Format using large glyph: for reference only. Although the glyphs are large the amount of information per page is limited.
- Format using 2 columns per page with 6 sub-columns per characters allowing 6 sources per line. The first six sources are always represented in the same order if present (G, T, J, K, V, and KP) in the first line. If other sources exist (H, U, and M), they are represented in that order in a second line. The main issue with that format is that it is very wasteful in space when only 2 sources exits but on different lines (such as 4E04, 4E1A, and 4E1C in the example). This format allows about 32 characters per page (versus 48 for the current multi-column format for the main CJK Unified collection).
- Variant of the format above but representing the second line sources (H, U, and M) in holes in the first line starting from the right. These sources are highlighted to make clear that they are not related to the sub-column standard source. In the example, the H source is shown in the KP column for 4E04, 4E1A, 4E1C and the V column for 4E21. This format allows about 36 characters per page in the main CJK Unified collection.
- Format using 3 columns per page with 4 sub-columns per characters allowing 4 sources per line. The sources are represented in the order of appearance in the CJKU_SR.txt (G, T, J, K, V, H, KP, U, and M). This formats allows about 50-54 characters per page in the CJK Extension A collection (versus 48 for the current multi-column format for the same collection).
- Format using 5 columns per page with 2 sub-columns per characters allowing 4 sources per line. The sources are represented in the order of appearance in the CJKU_SR.txt (G, T, J, K, V, H, KP, U, and M). This formats allows about 80-90 characters per page in the CJK Extension B and Extension C collections (versus 128 for the current single-column format for Ext B and 32 for the multi-column format for Ext C).

Final note: the following examples use commercial fonts available to the author to represent source characters. Many sources were unavailable and are therefore not correctly represented. No claim is made about accuracy of these glyphs. It is assumed that appropriate fonts will be delivered to the contributing editors in charge of font production. The source info has been copied manually from CJKU_SR.txt, so the examples may contain mistakes. The real process will collect data directly from the file.

**Format using large glyphs (2 Cols)**

| | | | |
|---|---|---|---|
| **4E00** 一 | | | |
| 一 | | | |
| 1.0 | G0-523B | T1-4421 | J0-306C | K0-6C69 |
| 一 一 | | | |
| V1—4A21 | KP0-FCD6 | | |
| **4E01** 丁 | | | |
| 1.1 | G0-3621 | T1-4421 | J0-437A | K0-6F4B |
| 丁 丁 | | | |
| V1-4A22 | KP0-E8B9 | | |
| **4E02** 丂 | | | |
| 1.1 | G5-3021 | T4-2126 | J1-3021 |
| **4E03** 七 | | | |
| 1.1 | G0-465F | T1-4424 | J0-3C37 | K0-7652 |
| 七 七 | | | |
| V1—4A23 | KP0-EFA6 | | |
| **4E04** 上 | | | |
| 1.1 | GE-2121 | T3-2126 | J1-3022 | H-9EB3 |
| **4E05** 丁 | | | |
| 1.1 | GE-2122 | T3-2125 | J1-3023 |
| **4E06** 丆 | | | |
| 1.1 | G1-7D3D | K2-2121 |
| **4E07** 万 | | | |
| 1.2 | G0-4D72 | T2-2126 | J0-4B7C | K0-5832 |
| 万 万 | | | |
| V1—4A24 | KP0-DAB9 | | |
| **4E08** 丈 | | | |
| 1.2 | G0-5549 | T1-4437 | J0-3E66 | K0-6D5B |

丈 丈

V1-4A25   KP0-E6DD

| | | | |
|---|---|---|---|
| **4E09** 三 | | | |
| 1.2 | G0-487D | T1-4435 | J0-3B30 | K0-5F32 |
| 三 三 | | | |
| G0-487D | T1-4435 | | |
| **4E0A** 上 | | | |
| 1.2 | G0-494F | T1-4438 | J0-3E65 | K0-5F3E |
| 上 上 | | | |
| V1-4A27 | KP0-E1C2 | | |
| **4E0B** 下 | | | |
| 1.2 | G0-4F42 | T1-4436 | J0-323C | K0-793B |
| 下 下 | | | |
| V1-4A28 | KP0-F2BA | | |
| **4E0C** 丌 | | | |
| 1.2 | G0-5822 | T2-2127 | J1-3024 | K2-2122 |
| **4E0D** 不 | | | |
| 1.3 | G0-323B | T1-4462 | J0-4954 | K0-5C74 |
| 不 不 | | | |
| V1-4A29 | KP0-DFBE | | |
| **4E0E** 与 | | | |
| 1.3 | G0-536B | T2-212F | J0-4D3F | K2-2123 |
| 与 与 | | | |
| V1-4A21 | KP0-FCD6 | | |
| **4E0F** 丏 | | | |
| 1.3 | G3-3021 | T2-212D | K2-2124 |

| Code | Char | Ref | G | T | J | K | V | KP |
|---|---|---|---|---|---|---|---|---|
| 4E00 | 一 | 一 1.0 | G0-523B | T1-4421 | J0-306C | K0-6C69 | V1-4A21 | KP0-FCD6 |
| 4E01 | 丁 | 一 1.1 | G0-3621 | T1-4421 | J0-437A | K0-6F4B | V1-4A22 | KP0-E8B9 |
| 4E02 | 丂 | 一 1.1 | G5-3021 | T4-2126 | J1-3021 | | | |
| 4E03 | 七 | 一 1.1 | G0-465F | T1-4424 | J0-3C37 | K0-7652 | V1-4A23 | KP0-EFA6 |
| 4E04 | 丄 | 一 1.1 | G0-523B | T1-4421 | J0-306C | | | |
| | 上 | | H-9EB3 | | | | | |
| 4E05 | 丅 | 一 1.1 | GE-2122 | T3-2125 | J1-3023 | | | |
| 4E06 | 丆 | 一 1.1 | G1-7D3D | | | K2-2121 | | |
| 4E07 | 万 | 一 1.2 | G0-4D72 | T2-2126 | J0-4B7C | K0-5832 | V1-4A24 | KP0-DAB9 |
| 4E08 | 丈 | 一 1.2 | G0-5549 | T1-4437 | J0-3E66 | K0-6D5B | V1-4A25 | KP0-E6DD |
| 4E09 | 三 | 一 1.2 | G0-487D | T1-4435 | J0-3B30 | K0-5F32 | V1-4A26 | KP0-E1B5 |
| 4E0A | 上 | 一 1.2 | G0-494F | T1-4438 | J0-3E65 | K0-5F3E | V1-4A27 | KP0-E1C2 |
| 4E0B | 下 | 一 1.2 | G0-4F42 | T1-4436 | J0-323C | K0-793B | V1-4A28 | KP0-F2BA |
| 4E0C | 丌 | 一 1.2 | G0-5822 | T2-2127 | J1-3024 | K2-2122 | | |
| 4E0D | 不 | 一 1.3 | G0-323B | T1-4462 | J0-4954 | K0-5C74 | V1-4A29 | KP0-DFBE |
| 4E0E | 与 | 一 1.3 | G0-536B | T2-212F | J0-4D3F | K2-2123 | V1-4A2A | KP0-FCD6 |
| 4E0F | 丏 | 一 1.3 | G3-3021 | T2-212D | | K2-2124 | | |
| 4E10 | 丐 | 一 1.3 | G0-5824 | T1-4461 | J0-5022 | K2-2125 | V1-4A2B | KP1-3409 |
| 4E11 | 丑 | 一 1.3 | G0-3373 | T1-4460 | J0-312F | K0-7564 | V1-4A2C | KP0-EEC9 |
| 4E12 | 刃 | 一 1.3 | GE-2123 | T4-2139 | J1-3025 | | | KP1-340E |
| 4E13 | 专 | 一 1.3 | G0-5728 | | | | | |
| 4E14 | 且 | 一 1.4 | G0-4752 | T1-4562 | J0-336E | K0-7326 | V1-4A2D | KP0-ECA8 |
| 4E15 | 丕 | 一 1.4 | G0-5287 | T1-4561 | J0-5023 | K0-5D60 | V1-4A2E | KP0-DFC9 |
| 4E16 | 世 | 一 1.4 | G0-4A40 | T1-4560 | J0-4024 | K0-6126 | V1-4A2F | KP0-E5F9 |
| 4E17 | 丗 | 一 1.3 | GE-2124 | T4-2155 | J0-5242 | | | KP1-3413 |
| 4E18 | 丘 | 一 1.4 | G0-4770 | T1-4563 | J0-3556 | K0-4E78 | V1-4A30 | KP0-D0DF |
| 4E19 | 丙 | 一 1.4 | G0-317B | T1-455F | J0-4A3A | K0-5C30 | V1-4A31 | KP0-DDF9 |
| 4E1A | 业 | 一 1.4 | G0-5235 | | | | | |
| | 业 | | H-9EB2 | | | | | |
| 4E1B | 丛 | 一 1.4 | G0-3454 | | | | | |
| 4E1C | 东 | 一 1.4 | G0-362B | | | | | |
| | 东 | | H-9DD6 | | | | | |
| 4E1D | 丝 | 一 1.4 | G0-4B3F | | | | | |
| 4E1E | 丞 | 一 1.5 | G0-5829 | T1-4722 | J—3E67 | K0-632A | V1-4A32 | KP0-E4EF |
| 4E1F | 丢 | 一 1.5 | GE-2125 | T1-4723 | J1-3026 | K1-6D4A | | |

**Format using highlight to represent sources in non standard location (here: H source) (2 Cols)**

| | | | | | | |
|---|---|---|---|---|---|---|
| **4E00** 一 1.0 | 一 G0-523B | 一 T1-4421 | 一 J0-306C | 一 K0-6C69 | 一 V1-4A21 | KP0-FCD6 |
| **4E01** 一 1.1 | 丁 G0-3621 | 丁 T1-4421 | 丁 J0-437A | 丁 K0-6F4B | 丁 V1-4A22 | KP0-E8B9 |
| **4E02** 一 1.1 | 丂 G5-3021 | 丂 T4-2126 | 丂 J1-3021 | | | |
| **4E03** 一 1.1 | 七 G0-465F | 七 T1-4424 | 七 J0-3C37 | 七 K0-7652 | 七 V1-4A23 | KP0-EFA6 |
| **4E04** 一 1.1 | 丄 G0-523B | 丄 T1-4421 | 丄 J0-306C | | 丄 H-9EB3 | |
| **4E05** 一 1.1 | 丅 GE-2122 | 丅 T3-2125 | 丅 J1-3023 | | | |
| **4E06** 一 1.1 | 厂 G1-7D3D | | 厂 K2-2121 | | | |
| **4E07** 一 1.2 | 万 G0-4D72 | 万 T2-2126 | 万 J0-4B7C | 万 K0-5832 | 万 V1-4A24 | KP0-DAB9 |
| **4E08** 一 1.2 | 丈 G0-5549 | 丈 T1-4437 | 丈 J0-3E66 | 丈 K0-6D5B | 丈 V1-4A25 | KP0-E6DD |
| **4E09** 一 1.2 | 三 G0-487D | 三 T1-4435 | 三 J0-3B30 | 三 K0-5F32 | 三 V1-4A26 | KP0-E1B5 |
| **4E0A** 一 1.2 | 上 G0-494F | 上 T1-4438 | 上 J0-3E65 | 上 K0-5F3E | 上 V1-4A27 | KP0-E1C2 |
| **4E0B** 一 1.2 | 下 G0-4F42 | 下 T1-4436 | 下 J0-323C | 下 K0-793B | 下 V1-4A28 | KP0-F2BA |
| **4E0C** 一 1.2 | 丌 G0-5822 | 丌 T2-2127 | 丌 J1-3024 | 丌 K2-2122 | | |
| **4E0D** 一 1.3 | 不 G0-323B | 不 T1-4462 | 不 J0-4954 | 不 K0-5C74 | 不 V1-4A29 | KP0-DFBE |
| **4E0E** 一 1.3 | 与 G0-536B | 与 T2-212F | 与 J0-4D3F | 与 K2-2123 | 与 V14A2A | KP0-FCD6 |
| **4E0F** 一 1.3 | 丏 G3-3021 | 丏 T2-212D | | 丏 K2-2124 | | |
| **4E10** 一 1.3 | 丐 G0-5824 | 丐 T1-4461 | 丐 J0-5022 | 丐 K2-2125 | 丐 V1-4A2B | KP1-3409 |
| **4E11** 一 1.3 | 丑 G0-3373 | 丑 T1-4460 | 丑 J0-312F | 丑 K0-7564 | 丑 V1-4A2C | KP0-EEC9 |
| **4E12** 一 1.3 | 刃 GE-2123 | 刃 T4-2139 | 刃 J1-3025 | | | 刃 KP1-340E |
| **4E13** 一 1.3 | 专 G0-5728 | | | | | |
| **4E14** 一 1.4 | 且 G0-4752 | 且 T1-4562 | 且 J0—336E | 且 K0-7326 | 且 V1-4A2D | KP0-ECA8 |
| **4E15** 一 1.4 | 丕 G0-5287 | 丕 T1-4561 | 丕 J0-5023 | 丕 K0-5D60 | 丕 V1-4A2E | KP0-DFC9 |
| **4E16** 一 1.4 | 世 G0-4A40 | 世 T1-4560 | 世 J0-4024 | 世 K0-6126 | 世 V1-4A2F | KP0-E5F9 |
| **4E17** 一 1.3 | 丗 GE-2124 | 丗 T4-2155 | 丗 J0-5242 | | | 丗 KP1-3413 |
| **4E18** 一 1.4 | 丘 G0-4770 | 丘 T1-4563 | 丘 J0-3556 | 丘 K0-4E78 | 丘 V1-4A30 | KP0-D0DF |
| **4E19** 一 1.4 | 丙 G0-317B | 丙 T1-455F | 丙 J0-4A3A | 丙 K0-5C30 | 丙 V1-4A31 | KP0-DDF9 |
| **4E1A** 一 1.4 | 业 G0-5235 | | | | 业 H-9EB2 | |
| **4E1B** 一 1.4 | 丛 G0-3454 | | | | | |
| **4E1C** 一 1.4 | 东 G0-362B | | | | 东 H-9DD6 | |
| **4E1D** 一 1.4 | 丝 G0-4B3F | | | | | |
| **4E1E** 一 1.5 | 丞 G0-5829 | 丞 T1-4722 | 丞 J0-3E67 | 丞 K0-632A | 丞 V1-4A32 | KP0-E4EF |
| **4E1F** 一 1.5 | 丟 GE-2125 | 丟 T1-4723 | 丟 J1-3026 | 丟 K1-6D4A | | |
| **4E20** 一 1.5 | 北 G5-3023 | 北 T3-2262 | | | | |
| **4E21** 一 1.5 | 両 GE-2126 | 両 T3-22612 | 両 J0-4E3E | 両 K2-2126 | 両 H-994F | 両 KP1-3415 |
| **4E22** 一 1.5 | 丢 G0-362A | 丢 T3-2263 | | | 丢 V1-4A33 | KP1-3417 |
| **4E23** 一 1.6 | 乤 GE-2127 | 乤 T4-2335 | 乤 J1-3027 | | | 乤 KP1-3419 |

5

**Ext A format (3 Cols)**

| Code | Radical | G | T | K | KP/other |
|---|---|---|---|---|---|
| 41C0 | 穴 116.15 | G5-5E61 | T4-637C | K3-2E22 | |
| 41C1 | 穴 116.17 | G3-5F71 | T4-6922 | | KP1-6366 |
| 41C2 | 立 117.1 | G_KX0870? | T3-2434 | | |
| 41C3 | 立 117.3 | G_KX0870? | T3-2A46 | | |
| 41C4 | 立 117.3 | GS-237B | V0-3F51 | H-994B | |
| 41C5 | 立 117.4 | G3-5F4A | T4-2A76 | K3-2E22 | |
| 41C6 | 立 117.4 | G-HZ | T3-2E4A | | |
| 41C7 | 立 117.5 | GS-237A | T3-3325 | | |
| 41C8 | 立 117.5 | G3-5F4E | T4-2E6D | | |
| 41C9 | 立 117.5 | G_HZ | T3-3322 | | |
| 41CA | 立 117.5 | JA-2549 | H-8E55 | | |
| 41CB | 立 117.6 | G-KX0871? | T5-3446 | JA-254A | |
| 41CC | 立 117.7 | G_KX08711 | T3-3D6F | KP1-6386 | |
| 41CD | 立 117.7 | G3-5F4F | T4-396A | K3-2E23 | KP1-6386 |
| 41CE | 立 117.8 | G5-5E28 | T3-4348 | K3-2E24 | KP1-638D |
| 41CF | 立 117.8 | G3-5F50 | T4-3F54 | K3-2E26 | H-994E / KP1-6390 |
| 41D0 | 立 117.8 | G3-5F51 | T4-3F55 | | KP1-6389 |
| 41D1 | 立 117.8 | G5-5E2A | T4-3F5A | K3-2E27 | KP1-638A |
| 41D2 | 立 117.11 | G5-5E2D | T4-4563 | JA-254B | |
| 41D3 | 立 117.12 | G3-5F52 | T4-5752 | | KP1-639F |
| 41D4 | 立 117.13 | G3-5F54 | T4-5C3B | K3-2E28 | |
| 41D5 | 立 117.14 | G_KX08732 | T3-5A68 | K3-2E29 | KP1-63A4 |
| 41D6 | 竹 118.3 | G_KX08780 | T5-2B21 | K3-2E2A | |
| 41D7 | 竹 118.4 | G3-634E | T4-2E75 | KP1-63C3 | |
| 41D8 | 竹 118.4 | G5-6276 | T4-2E70 | K3-2E2B | KP1-63C8 |
| 41D9 | 竹 118.4 | G_KX08782 | T3-3328 | | |
| 41DA | 竹 118.4 | G5-6278 | T4-2E77 | K3-2E2C | |
| 41DB | 竹 118.4 | G_KX08791 | T3-3329 | V2-7F4B | H-8EFE |
| 41DC | 竹 118.4 | G3-634F | T4-2E73 | K3-2E2D | KP1-63B6 |
| 41DD | 竹 118.4 | G3-634A | T4-2E72 | | |
| 41DE | 竹 118.5 | G3-6355 | T4-3376 | K3-2E2E | KP1-63D0 |
| 41DF | 竹 118.5 | G5-6327 | T4-337D | K3-2E2F | |
| 41E0 | 竹 118.5 | G_KX08792 | T3-3774 | V0-3F71 | KP1-63EB |
| 41E1 | 竹 118.5 | G_KX08800 | T5-3448 | JA-254C | KP1-63E7 |
| 41E2 | 竹 118.5 | G3-6356 | T4-3377 | K3-2E30 | KP1-63D1 |
| 41E3 | 竹 118.5 | G5-6324 | T4-337B | K3-2E31 | KP1-63E3 |
| 41E4 | 竹 118.5 | G_KX0881 | T3-3775 | | |
| 41E5 | 竹 118.5 | G7-2368 | T6-4276 | | |
| 41E6 | 竹 118.5 | G3-6324 | T4-3378 | JA-7347 | |
| 41E7 | 竹 118.6 | G3-6363 | T4-396D | | |
| 41E8 | 竹 118.6 | G3-6371 | T4-3974 | KP1-63F9 | |
| 41E9 | 竹 118.6 | G3-6366 | T4-3971 | KP1-641B | |
| 41EA | 竹 118.6 | G_KX0882 | T3-3D74 | V3-3649 | |
| 41EB | 竹 118.6 | G3-637D | T4-3F58 | KP1-6414 | |
| 41EC | 竹 118.6 | G3-636A | T43972 | K3-2E32 | KP1-6409 |
| 41ED | 竹 118.6 | G5-632D | T3-3D7A | H-8D5F | KP1-6402 |
| 41EE | 竹 118.6 | G5-6334 | T4-3975 | JA-254D | V2-7F50 / KP1-641A |
| 41EF | 竹 118.6 | G5-632E | T3-3D73 | H-8E59 | KP1-641A |
| 41F0 | 竹 118.6 | G3-6379 | T4-396F | KP1-640A | |
| 41F1 | 竹 118.6 | G_KX0884 | T6-4C5B | JA-254E | |
| 41F2 | 竹 118.6 | GS-233A | T6-4C56 | | |
| 41F3 | 竹 118.6 | JA-254F | | | |

**Ext B and Ext C format (5 Cols)**

| Code | Radical.Stroke | Sources |
|---|---|---|
| 20000 | 一 1.1 | G_KX00750 T5-2515 |
| 20001 | 一 1.1 | G_HZ |
| 20002 | 一 1.1 | TF-2121 |
| 20003 | 一 1.2 | G_KX0076 T6-212F |
| 20004 | 一 1.2 | T6-212D |
| 20005 | 一 1.2 | G_HZ T6-212F |
| 20006 | 一 1.2 | K4-0002 |
| 20007 | 一 1.3 | G_KX00770 T6-2142 |
| 20008 | 一 1.3 | G_KX00770 T6-2143 |
| 20009 | 一 1.3 | T5-2133 KP1-3408 |
| 2000A | 一 1.3 | G_HZ |
| 2000B | 一 1.3 | G_HZ T3-2144 / J3-2E22 KP1-340C |
| 2000C | 一 1.3 | G_HZ |
| 2000D | 一 1.4 | G_KX00771 |
| 2000E | 一 1.4 | G_4K |
| 2000F | 一 1.4 | TF-213E |
| 20010 | 一 1.4 | TF-213F |
| 20011 | 一 1.4 | G_HZ |
| 20012 | 一 1.4 | G_HZ T6-222B |
| 20013 | 一 1.4 | G_HZ |
| 20014 | 一 1.4 | G_HZ T5-214D |
| 20015 | 一 1.4 | G_HZ |
| 20016 | 一 1.4 | V0-3F5F |
| 20017 | 一 1.4 | V0-3F60 |
| 20018 | 一 1.5 | G_KX007 T6-2340 |
| 20019 | 一 1.5 | G_KX007 T5-233E |
| 2001A | 一 1.5 | G_KX007 T6-233F |
| 2001B | 一 1.5 | G_HZ |
| 2001C | 一 1.5 | G_HZ |
| 2001D | 一 1.5 | G_HZ |
| 2001E | 一 1.5 | G_HZ |
| 2001F | 一 1.5 | G_HZ |
| 20020 | 一 1.5 | G_HZ T6-2467 |
| 20021 | 一 1.6 | G_KX007 T6-255F / H-9C71 |
| 20022 | 一 1.6 | G_KX0078 T5-232F |
| 20023 | 一 1.6 | TF-2274 |
| 20024 | 一 1.6 | G_HZ |
| 20025 | 二 7.5 | G_HZ T6-2567 |
| 20026 | 一 1.6 | G_HZ |
| 20027 | 一 1.6 | V0-354F |
| 20028 | 一 1.6 | V2-6E21 |
| 20029 | 刂 6.5 | G_HZ T6-2563 |
| 2002A | 一 1.6 | V0-456C |
| 2002B | 一 1.6 | V0-456D |
| 2002C | 一 1.7 | G_KX00781 T6-2937 |
| 2002D | 一 1.7 | G_KX0078 T6-293A |
| 2002E | 一 1.7 | G_KX0078 T6-2938 |
| 2002F | 口 31.5 | G_HZ |
| 20030 | 一 1.7 | G_HZ |
| 20031 | 一 1.7 | G_HZ |
| 20032 | 一 1.7 | V0-305F |
| 20033 | 一 1.7 | V2-6E25 |
| 20034 | 一 1.7 | V0-354A |
| 20035 | 一 1.8 | TF-2922 |
| 20036 | 一 1.8 | TF-2923 |
| 20037 | 口 30.6 | G_HZ |
| 20038 | 一 1.8 | G_HZ |
| 20039 | 口 30.6 | G_HZ |
| 2003A | 一 1.8 | G_HZ |
| 2003B | 一 1.8 | G_HZ T6-2E66 |
| 2003C | 夊 35.7 | G_HZ |
| 2003D | 一 1.9 | G_HZ |
| 2003E | 一 1.9 | H-9375 |
| 2003F | 一 1.9 | V2-6E27 |
| 20040 | 一 1.9 | V0-3F68 |
| 20041 | 一 1.10 | G_KX0078 T5-3072 / KP1-341E |
| 20042 | 一 1.10 | V2-6E26 |
| 20043 | 一 1.10 | G_HZ T5-3323 |
| 20044 | 一 1.10 | V0-354B |
| 20045 | 一 1.10 | G_HZ |
| 20046 | 一 1.11 | TF-3932 H-9376 |
| 20047 | 一 1.11 | TF-3933 |
| 20048 | 一 1.11 | G_HZ T6-472D |
| 20049 | 一 1.11 | G_HZ |
| 2004A | 一 1.11 | G_HZ |
| 2004B | 一 1.11 | TF-3B73 |
| 2004C | 一 1.12 | TF-4035 |
| 2004D | 一 1.12 | TF-4075 |
| 2004E | 一 1.13 | H-9548 |
| 2004F | 一 1.13 | G_HZ |
| 20050 | 一 1.13 | G_HZ |
| 20051 | 一 1.13 | V0-354C |
| 20052 | 一 1.13 | TF-4742 |
| 20053 | 一 1.14 | TF-4D56 |
| 20054 | 一 1.14 | V2-6F21 |
| 20055 | 一 1.14 | G_KX0078 T7-2121 / KP1-341F |

-------