# RootLGR における JGP の具体的課題定義(案)

1 2

### 3 1 背景

4 5 1.1 歴史

c

- 6 国際化ドメイン名(Internationalized Domain Name、以降 IDN と記載)は、英語以外の言
- 7 語を母語とする人々にとって利用しやすいアドレスを提供するものとして、1998年にアジ
- 8 アを起点として検討が開始された。2000年にインターネットプロトコルの標準化を行って
- 9 いる IETF(Internet Engineering Task Force)で WG が設立され、2003 年に標準が RFC
- 10 3490、3491、3492 として制定された1。
- 11 IDN 標準化の際に、中国語圏から IDN の利用には異体字2問題の解決が必要だとの主張が
- 12 行われた。中国語における異体字問題とは、中国本土で使われている簡体字と台湾や香港な
- 13 どで使われている繁体字は基本的に 1 対 1 に対応付くものであり、簡体字の地域と繁体字
- 14 の地域間でのコミュニケーションのためには簡体字のアドレスと繁体字のアドレスは同一
- 15 のものとしてアクセスできる必要がある、というものである。しかし、この異体字のあるべ
- 16 き姿は、言語ごとに異なり、たとえば、共通的に漢字を使用する中国語と日本語においても、
- 17 異体字のとらえ方は異なっている。このような背景から、IETFにおける標準は世界中どこ
- 18 でも共通して使えるプロトコルである必要があるため、中国語の異体字問題は地域化
- 19 (Localization)であるとして、IDN プロトコルの検討対象外とされた。
- 20 中国語の異体字問題は、漢字を使う言語文化圏の中国・日本・韓国・台湾の NIC(Network
- 21 Information Center)が中心となって結成した JET(Joint Engineering Team)で検討が行わ
- 22 れ、ドメイン名登録時に異体字を含む文字列のバリエーションを同一登録者に紐付けて、そ
- 23 のまとまり(パッケージ)を不可分とするという運用方式が合意された。この方式は JET
- 24 Guidelines としてまとめられ、2004年に RFC 3743 として発行された。また、この考え方
- 25 は ICANN の IDN Implementation Guidelines にも取り入れられ、IDN 登録を行うレジス
- 26 トリは登録可能な文字のコードポイントとその異体字リストのテーブルを IANA に登録す
- 27 ることが推奨されている。

2829

#### 1.2 JET Guidelines の目的

30 31

JET Guidelines は、異体字を持つ漢字が含まれるドメイン名の登録時に、登録申請された

<sup>1 2010</sup> 年に、RFC 5890、5891、5892、5893、5894 で更新された。これらの標準を区別するため、2003 年版を IDNA2003、2010 年版を(方式が決定した年から)IDNA2008 と呼んでいる。

<sup>&</sup>lt;sup>2</sup> 一般には、同音同義の漢字であって、字体だけが異なり、どの文脈でも交換可能である 漢字と定義される。本書では、ある言語において同音同義の漢字であって、Unicode のコードポイントは異なるが同一の文字として取り扱われるべき漢字と定義する。

- 32 ドメイン名(文字列)と異体字で置き換えた文字列(群)を一つのまとまり(パッケージ)として、
- 33 同一登録申請者に紐付ける運用方式を提案するものである。JET Guidelines では、登録可
- 34 能な文字とその文字の異体字を結び付ける IDN テーブルのフォーマットと、パッケージを
- 35 生成するアルゴリズムを規定している。JET Guidelines で規定する IDN テーブルの構成
- 36 を以下に示す。

- 38 ▶ 構成要素
- 39 **✓** その IDN テーブルが適用される言語
- 40 ✓ その IDN テーブルに含まれる文字の出典のリスト(参照番号で区別)
- 41 ✓ その IDN テーブルのバージョン番号
- 42 ✓ 3カラムからなる IDN テーブル本体

43

- 44 ▶ IDN テーブル本体の各カラムの意味
- 45 第1カラム:登録可能文字(Unicode のコードポイント)と、出典の参照番号
- 46 第2カラム:ドメイン名での使用が推奨される異体字(群)と、その出典の参照番号
- 47 第3カラム:ドメイン名では使用されず予約される異体字(群)と、その出典の参照番号

- 49 各カラムはセミコロン(;)で区切られる。異体字が複数ある場合は、各異体字はカンマ(,)
- 50 で区切られる。以下に、例を示す(JET Guidelines から引用)。

Language Variant Table for zh-tw # 言語は台湾の中国語

Reference 1 CP950 (commonly known as BIG5) # 出典の参照番号 1 Reference 2 zVariant, zTradVariant, zSimpVariant in Unihan.txt Reference 3 List of Simplified Character Table (Traditional column)

Reference 4 zTradVariant in Unihan.txt

Version 1 20020701 # July 2002 # バージョン番号

#以下、IDN テーブル本体

#登録可能文字;推奨される異体字;予約される異体字

5718(1);5718(4);56E2(2),56E3(2) # 團;團;团,団

60F3(1);60F3(1); # 想;想;

6559(1);6559(1);654E(2) # 教;教;教;

6E05(1);6E05(1);6DF8(2) # 清;清;清

771F(1);771F(1);771E(2) # 真;真;眞

806F(1);806F(3);8054(2),8068(2) # 聯;聯;联,聯

96C6(1);96C6(1); # 集;集;

5253

この例で、「清(U+6E053) 真(U+771F) 教(U+6559)」という文字列が登録申請されたとする

と、登録者(ここでは登録者 A とする)に紐付けられる IDN のパッケージは以下となる。

5455

登録文字列:清(U+6E05) 真(U+771F) 教(U+6559)

予約文字列:清(U+6E05) 眞(U+771E) 教(U+6559)

清(U+6E05) 眞(U+771E) 教(U+654E)

清(U+6E05) 真(U+771F) 教(U+654E)

清(U+6DF8) 真(U+771F) 教(U+6559)

清(U+6DF8) 眞(U+771E) 教(U+6559)

清(U+6DF8) 真(U+771E) 教(U+654E)

清(U+6DF8) 真(U+771F) 教(U+654E)

56

57 登録者 A とは異なる別の登録者(ここでは登録者 B とする)が後から「清(U+6DF8) 眞

<sup>&</sup>lt;sup>3</sup> U+16 進数は、Unicode のコードポイントを示す。

58 (U+771E) 教(U+654E)」を登録申請しても、当該文字列中の文字はいずれも登録可能文字 59 に含まれていないため、登録できない。もし、いずれかの文字が登録可能文字であったとし 60 ても、その文字列は既に登録されている登録者 A のパッケージに含まれているため、登録 61 はできない。

6263

## 1.3 各 TLD における異体字の定義

64

- 65 いくつかのレジストリは、そのレジストリが運用する TLD 配下(主には SLD:第二レベルド
- 66 メイン名)で IDN を使用可能としている。.CN、.TW を中心とする中国語ドメイン名を登録
- 67 可能としているレジストリでは、それらレジストリが協力して共通の IDN テーブルを定義
- 68 し、そこから各レジストリが自レジストリで使う部分を切り出して JET Guidelines にした
- 69 がって定義し、それをルールとしている。
- 70 .JPでは、2001年に導入された汎用 JPドメイン名において日本語 JPドメイン名を使用可
- 71 能とすることをJPNICが決定し、専門家チームが日本語JPドメイン名のあり方を検討し、
- 72 外部の方々の意見も取り入れつつ、ルールを決定し、JPRS がそのルールに基づき運用して
- 73 いる。そのルールの特徴は、使用可能文字を JIS 第 1 水準と第 2 水準とし、異体字は設け
- 74 ないというものである。異体字を設けない理由は、日本では異なる文字コードを持つ文字は
- 75 異なる文字として扱うことが適切であるという判断をしたためである。

76

77 78

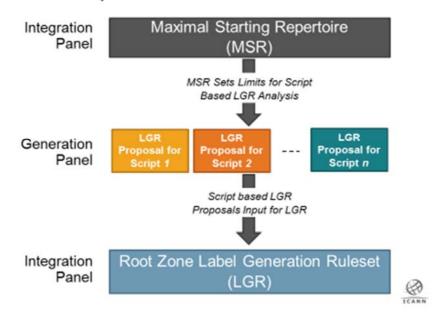
## 2 RootLGR と言語ルール

2.1 RootLGR とは

79 80

81 [RootLGR プロジェクトの背景、RootLGR の目的、RootLGR プロジェクトを構成する統合 82 パネルと各言語生成パネルの関係、現在の状況などをここで説明する予定。]

# LGR Development Process



8485

### 2.2 言語ルール

8687

RootLGR は Root ゾーン(TLD 文字列)用の統合ルールのことであり、各言語パネルが作成 した言語ルールを統合パネルが統合して作成する。各言語パネルが作成する言語ルールは JET Guidelines の IDN テーブルとほぼ同等であり、以下の情報から構成される。

8990

91

- ▶ 申請可能文字
- 92
   MSR<sup>4</sup>の中から申請可能な文字の範囲を言語パネルが選択する。選択できる文字は、言

   93
   語パネルがその言語で使用すると宣言した Script<sup>5</sup>の範囲に限られる。
- 94 ➤ 異体字リスト
- 95 個々の申請可能文字について、その文字が持つ異体字を一覧(リスト)にしたもの。異体 96 字リストに含まれる文字は、申請可能文字でなければならない。
- 97 ▶ 各異体字の属性
- 98 割り当て可能(その文字を含む文字列を Root ゾーンに登録することが可能)もしくはブ
- 99 ロック(その文字を含む文字列はRootゾーンに登録することが不可能)の2種類であり、
- 100 異体字リストに含まれる個々の文字について定義される。

<sup>4</sup> 統合パネルが作成した、ルートゾーン(TLD 文字列)として使用可能な文字の最大集合を定義したもの。Maximal Starting Repertoire の略称。

<sup>5</sup> Script は平仮名、片仮名、漢字などある言語を書き表すために使用される同種の文字(用字)のこと。

```
101 ▶ 文字列評価ルール
```

102 申請された文字列全体に対して適用されるルール。簡体字と繁体字の混在は許容しな 103 い(ブロックする)などの定義を行う。

104105

「1.2 JET Guidelines」で例として示した IDN テーブルを言語ルールで記述6すると以下の

106 ようになる。

```
<language>und-Hani</language> <!-- 言語は中国語 -->
<char cp="6E05" tag="sc:Hani"> <!-- 申請可能文字は「清」 -->-
  <var cp="6DF8" type="block" /> <var cp="6E05" type="both" />
  <!-- 異体字「清」はブロック、異体字「清」は簡体字と繁体字が同じ -->
</char>
<char cp="6DF8" tag="sc:Hani"> <!-- 清 -->
  <var cp="6DF8" type="block" /> <var cp="6E05" type="both" />
</char>
<char cp="771F" tag="sc:Hani"> <!-- 真 -->
  <var cp="771E" type="block" /> <var cp="771F" type="both" />
</char>
<char cp="771E" tag="sc:Hani"> <!-- 眞 -->
  <var cp="771E" type="block" /> <var cp="771F" type="both" />
</char>
<char cp="6559" tag="sc:Hani"> <!-- 教 -->
  <var cp="654E" type="block" /> <var cp="6559" type="both" />
</char>
<char cp="654E" tag="sc:Hani"> <!-- 教 -->
  <var cp="654E" type="block" /> <var cp="6559" type="both" />
</char>
<rules>
  <action disp="block" any-variant="block" />
  <action disp="allocate" only-variants="simp both" />
  <action disp="allocate" only-variants="trad both" />
  <action disp="block" any-variant="simp trad" />
</rules>
</data>
```

<sup>6</sup> IETF に提案されている Internet Draft "Representing Label Generation Rulesets using XML"(draft-davies-idntables)で形式が定義されている。

110		
110		JET Guidelines では、予約される異体字は登録可能文字の一覧に現れなくともよいが、
111		言語ルールでは異体字リストに含まれるすべての文字は申請可能文字でなければなら
112		ない。
113	>	JET Guidelines ではドメイン名登録を行うレジストリが独自の IDN テーブルを定義
114		することとしているが、言語ルールは言語コミュニティ(当該言語を使用している国・
115		地域などの言語文化圏)の合意に基づいて言語パネルが定義することとしている。
116	>	JET Guidelines の IDN テーブルはそれ自身が単独で適用されるルールであるが、言
117		語ルールは RootLGR を構成する要素であり単独で適用されることはない。
118		
119		2.2 RootLGR
120		
121	ID	NTLD として申請された文字列に対して適用されるルールは、言語パネルが作成した各
122	言言	語個別の言語ルールではなく、すべての言語ルールを統合パネルが統合した統合ルール
123	(Ro	ootLGR)である。これは、TLD 文字列はさまざまな言語の文字列が混在するため、特定
124	の	言語ルールを適用できないためである。
125	Ro	otLGR は、原則として各言語の言語ルールの和集合である。すなわち、異なる言語であ
126	つ	ても同じ Script を共有するもの(中国語と日本語の漢字など)であれば、同じ文字(Unicode
127	Ø:	コードポイント)に対して定義される異体字リストは、それぞれの言語で定義された異体
128	字	リストの和集合となる。
129	例:	えば、中国語ルールが「愛(U+611B)」と「爱(U+7231)」を異体字として

JET Guidelines の IDN テーブルと言語ルールの違いは以下の通りである。

と定義し、日本語ルールが「愛(U+611B)」を異体字なしの

135 と定義した場合、RootLGR は以下のようになる。

```
<data>
<char cp="611B" tag="sc:Hani"> <!-- 愛 -->
 <var cp="611B" type="simp" when="und-Hani" />
 <var cp="611B" type="allocate" when="und-Jpan" />
 <var cp="7231" type="trad" when="und-Hani" />
 <var cp="7231" type="block" when="und-Jpan" />
</char>
<char cp="7231" tag="sc:Hani"> <!-- 爱 -->
 <var cp="611B " type="simp" when="und-Hani" />
 <var cp="7231" type="trad" when="und-Hani" />
</char>
<rules>
 <action disp="block" any-variant="block" />
 <action disp="allocate" only-variants="simp both" />
 <action disp="allocate" only-variants="trad both" />
 <action disp="block" any-variant="simp trad" />
 <action disp="allocate" only-variants="allocate" />
</rules>
</data>
```

138 この例の RootLGR で、「愛(U+611B)」という文字列が日本語で登録申請されたとすると、

- 139 登録者(ここでは登録者 C とする)に割り当て可能な文字列は「愛(U+611B)」となり、「爱
- 140 (U+7231)」は登録者に紐付けられるがブロックされる文字列となる。登録者 C とは異なる
- 141 別の登録者(ここでは登録者 D とする)が後から「爱(U+7231)」を登録申請しても、登録者
- 142 Cに紐付けられているため、登録できない。

143144

## 3 JGP の具体的課題定義(案)

3.1 登録可能文字の範囲の決定

145146

147 JGP で登録可能とする文字の範囲を MSR の中から選ぶ。その範囲を選ぶ合理的な理由付

148 けが課題である。例えば、以下のような範囲が考えられる。

- 150 ▶ 常用漢字の範囲
- 151 ▶ 常用漢字+人名用漢字の範囲
- 152 ➤ JIS X 0208:2012 の第一水準・第二水準漢字の範囲

153	> など	
154		
155	3.2 中国語ルールの異体字が RootLGR に入った際の日本語への影響の理解	
156		
157	RootLGR に含む漢字の異体字として、JGP が CGP(中国語生成パネル)の異体字を容認	忍する
158	か個別の文字について協議するかは課題である。	
159	中国語ルールで定義された異体字を容認するとした場合は、TLD 文字列の登録申請を	日本
160	語だと宣言して行っても、中国語の異体字も区別なく異体字として扱われる。そのたる	め、次
161	のような状況が発生し得る。	
162		
163	▶ 登録者が意図していなかった文字を含む文字列が登録文字列に紐付けられブロッ	ックさ
164	れる。	
165	例:「航(U+822A)空(U+7A7A)」と「桁(U+6841)空(U+7A7A)」	
166		
167	日本語ルールが異体字を定義する・しないに関わらず、中国語ルールの異体字定義は日本語ルールが異体字を定義する・しないに関わらず、中国語ルールの異体字定義は日本語の表現を表現している。	3本語
168	に影響を与えることを一般利用者が理解する必要があり、その情報を周知することも	課題
169	である。	
170		
171	3.3 日本語ルールで異体字を定義するか	
172		
173	日本語ルールで異体字を定義する場合でもしない場合でも、その合理的な理由付けに	は課題
174	である。	
175	また、どのような観点からその理由付けを行うかも課題である。例えば、以下のような	よ観点
176	がある。	
177		
178	▶ 日本語 TLD 登録申請者の利益を最大にすること	
179	▶ 日本の一般利用者の混乱を最小にすること	
180	> など	
181		
182	3.4 日本語ルールが異体字を定義する場合、個々の異体字の根拠	
183		
184	日本語ルールで異体字を定義する場合、個々の異体字が何を根拠とするかは課題であ	る。こ
185	れは、上記「3.3日本語ルールで異体字を定義するか」で述べたとおり、どのような観点による。	見点に
186	より異体字を定義するかによって異なる。また、JGP が独自に異体字リストを作成で	上る場
187	合、それをどのように権威付けるかも課題である。	

189	3.5 日本語ルールで異体字を定義した場合の中国語への影響の把握
190	
191	日本語ルールで異体字を定義した場合で、かつ、その定義が中国語ルールが定義した異体字
192	と異なる場合、上記「3.2 中国語ルールの異体字が RootLGR に入った際の日本語への影響
193	の理解」で述べたことと同様の状況が発生し得る。これを中国語パネルとどのように調整す
194	るかは課題である。
195	
196	3.6 日本語ルール作成の迅速性
197	
198	新 gTLD の次のラウンドは早ければ 2016 年に開始されることが見込まれている。RootLGR
199	はそれまでに完成していることが必要であり、事前に CGP や KGP(韓国語生成パネル)との
200	調整を行うことを考慮すると、遅くとも2015年上半期中の日本語ルール作成が必要である。
201	迅速に日本語ルールを作成することは課題である。
202	
203	[EOF]