# Ruling The Root ●

# *CJK Rules For The Root Zone*

*Kenny Huang, Ph.D.* 黃勝雄博士

*Member, CDNC / CGP*
*Co-author, RFC3743 IETF*
*Member, Executive Council, APNIC*
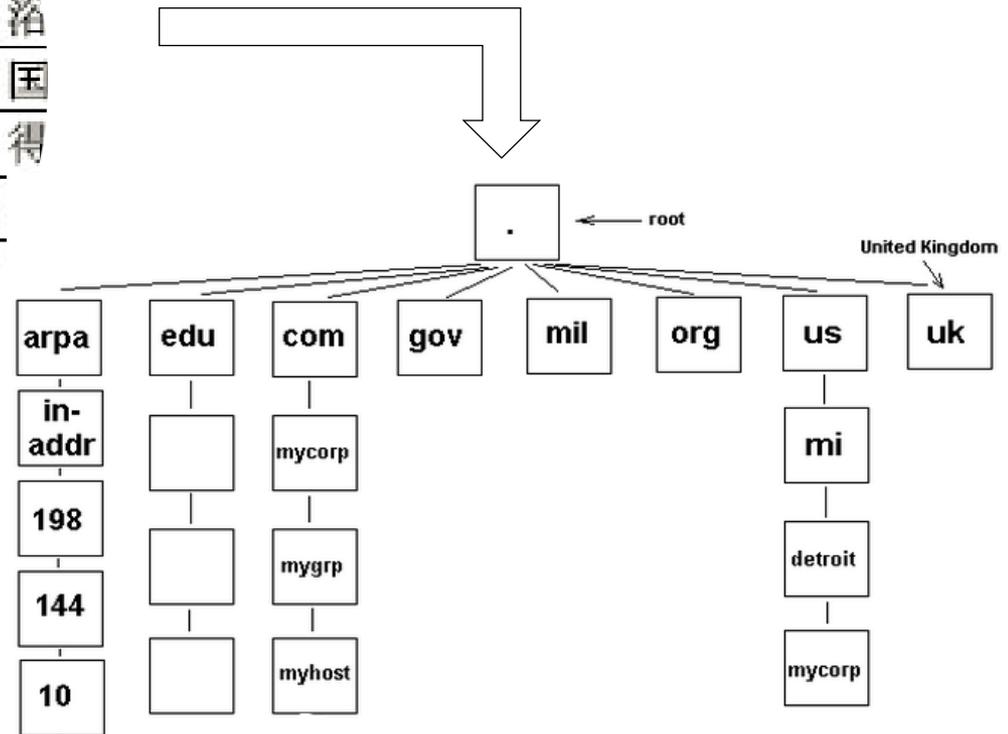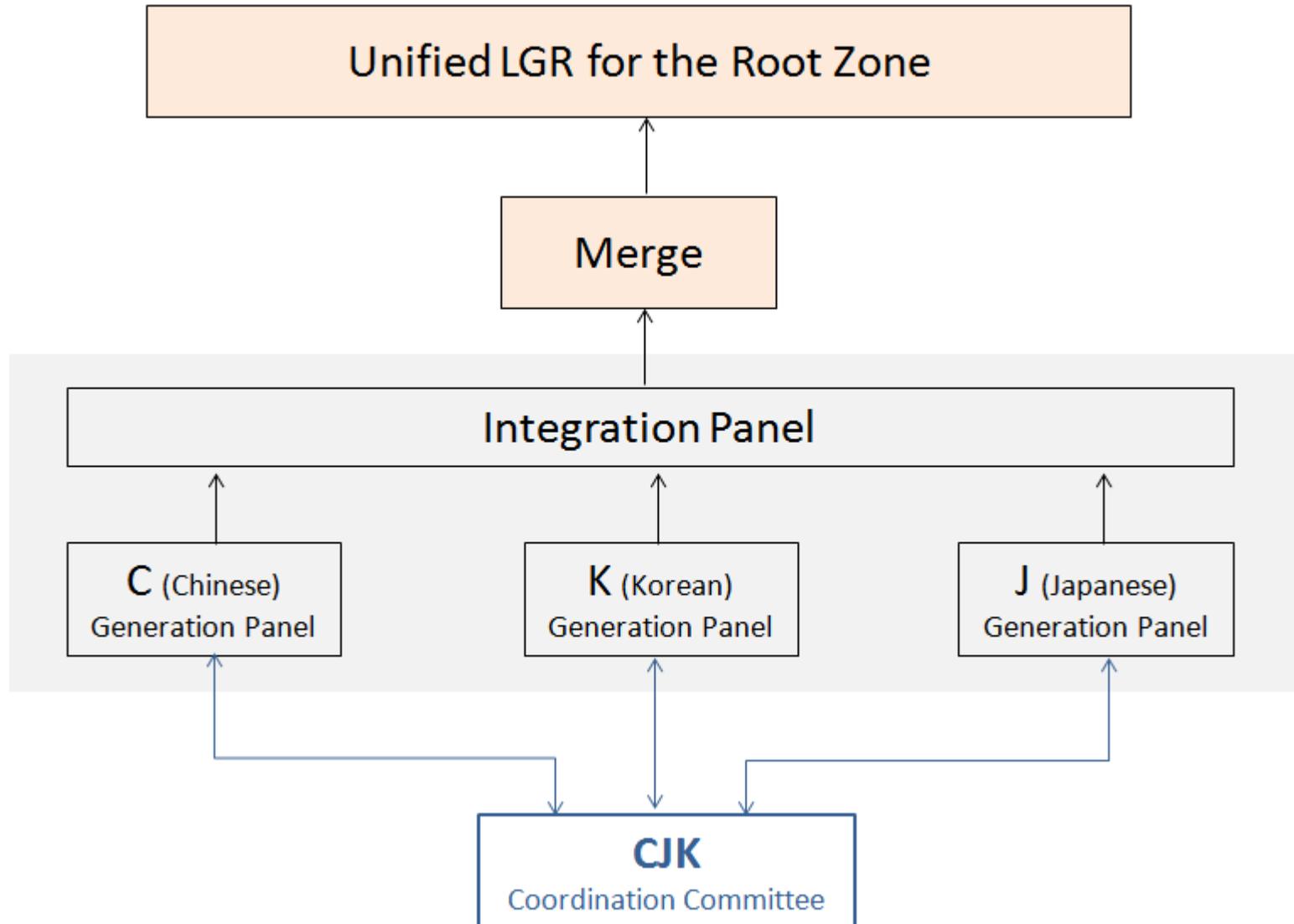*Member, Board of Directors, TWNIC*
*huangksh@gmail.com*
*2014.Jun*

ICANN FIFTY
22-26 JUNE 2014

# Problem : CJK Is Complicated

*Putting CJK labels in the root zone is even more complicated*

# Institutionalized Problem Solving : Structure
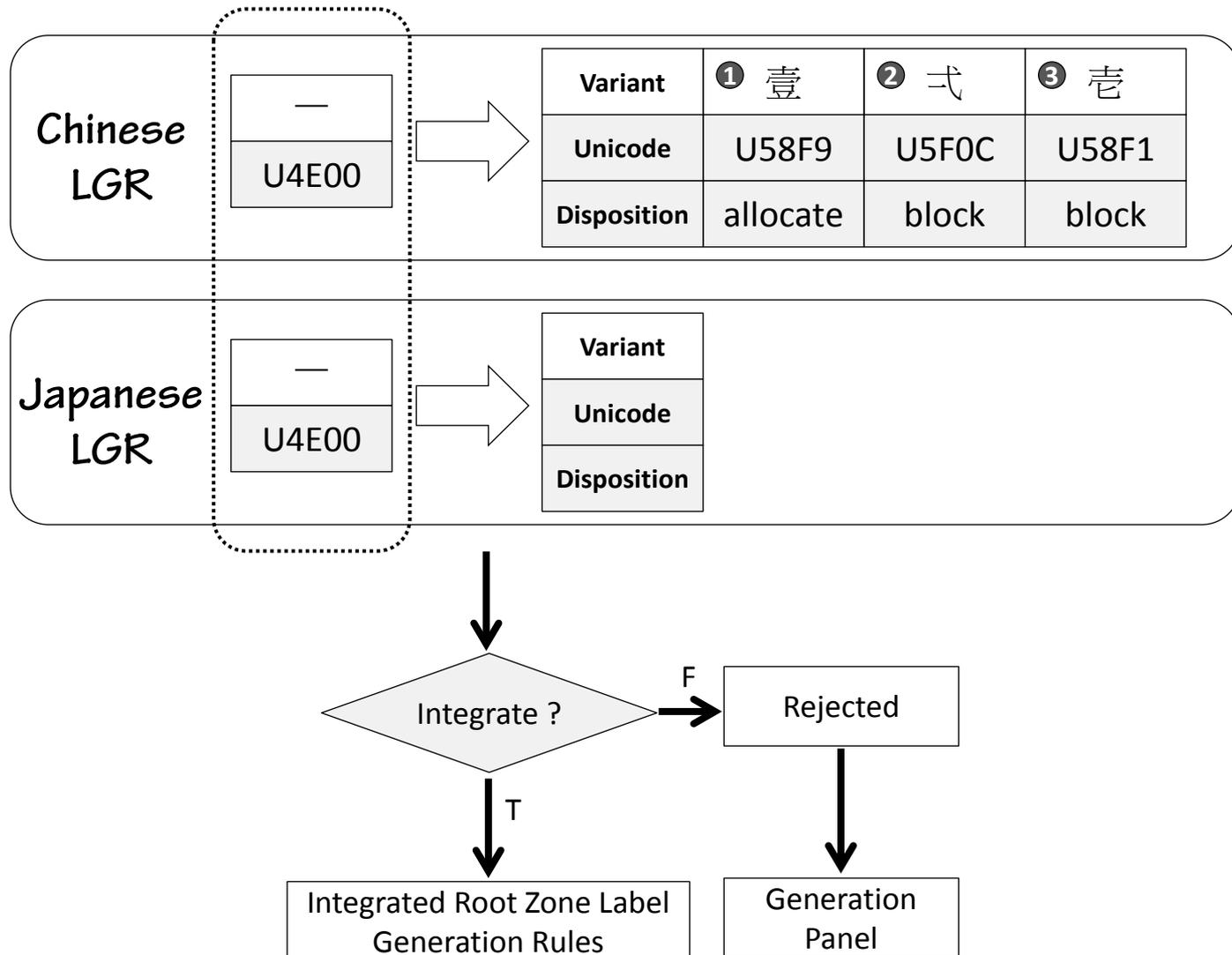
# Constraints for CJK LGR

| Independent Tasks | Coordination Tasks |
|---|---|
| ✓Each CJK Panel creates an LGR<br>✓Each LGR includes a repertoire and variants<br>  ➢Define labels permission<br>  ➢Define variants labels<br>  ➢Assign dispositions<br>    •Allocatable<br>    •Block | ✓If an LGR includes Han characters:<br>  ➢The variant *mappings* must agree for all the panels<br>  ➢The variant *types* may be different<br>  ➢The repertoires may be different |

*Presented by Lee Han Chuan  & IP,  Shanghai 2014 May 29

# Overlap Case Illustration

**Chinese LGR**
U4E00 一

| Variant | ❶ 壹 | ❷ 弍 | ❸ 壱 |
|---|---|---|---|
| Unicode | U58F9 | U5F0C | U58F1 |
| Disposition | allocate | block | block |

**Japanese LGR**
U4E00 一

| Variant |
|---|
| Unicode |
| Disposition |

Integrate ?

F → Rejected

T → Integrated Root Zone Label Generation Rules

Rejected → Generation Panel

5

# High Level Conflict Strategies

$$X$$
CJK overlap

C: rule $R_c$
J : rule $R_j$
K: rule $R_k$

## $R_{cjk}$ Conflict For Label X

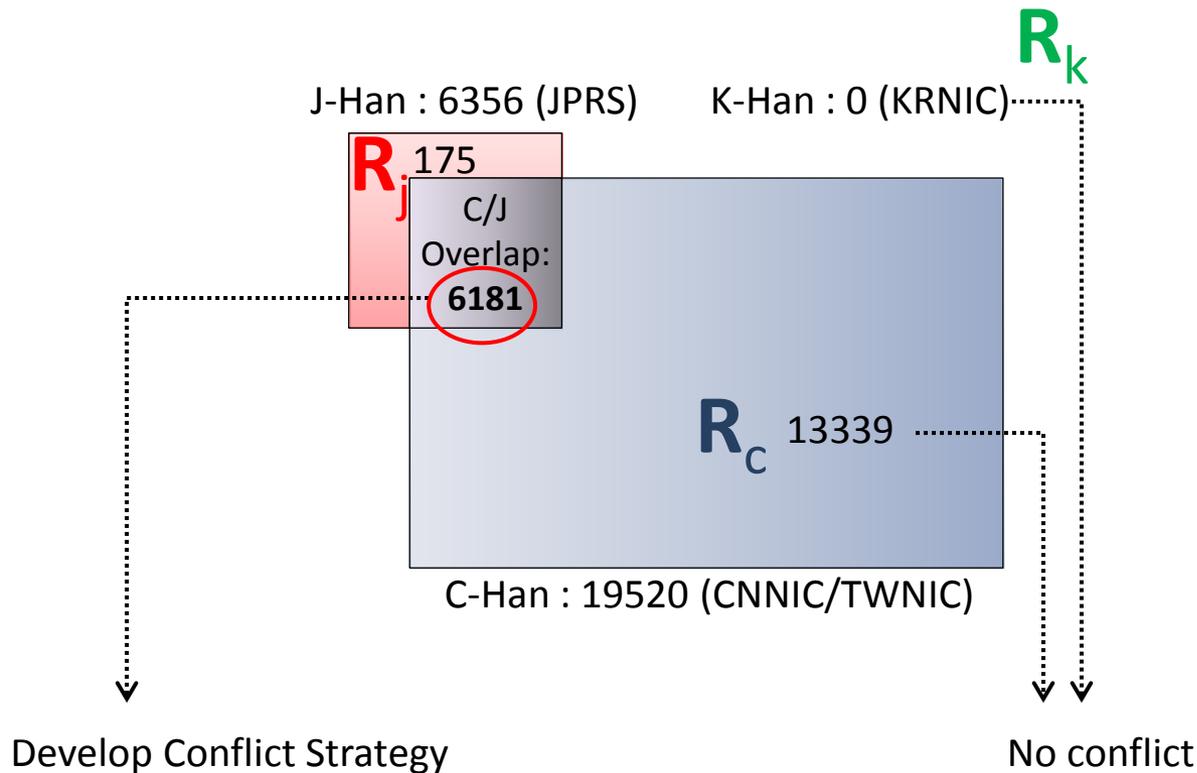| ID | Strategy | Pros | Cons | Rank |
|----|----------|------|------|------|
| 1 | Adopt X<br>Abandon $R_{cjk}$ | Permit X | No label rule | |
| 2 | Adopt X<br>Intersection $\cap$ ($R_{cjk}$) | Permit X<br>Permit $\cap$(variants/disp) | Rules changed | |
| 3 | Adopt X<br>Union $\cup$ ($R_{cjk}$) | Permit X<br>Permit $\cup$ (variants/disp) | Rules changed | |
| 4 | Abandon X and $R_{cjk}$ | No conflict | Label not available | |
| 5 | Adopt rules based on frequency of use | Fair & scientific approach | Rules changed; fairness doesn't mean appropriate | |

# Unified CJK LGR Illustration

## Chinese LGR

| 一 U4E00 → | Variant | ❶ 壹 | ❷ 弍 | ❸ 壱 |
|---|---|---|---|---|
| | Unicode | U58F9 | U5F0C | U58F1 |
| | Disposition | allocate | block | block |

## Japanese LGR

| 一 U4E00 → | Variant | | | |
|---|---|---|---|---|
| | Unicode | | | |
| | Disposition | | | |

## Union → Integrated LGR

| 一 U4E00 → | Variant | ❶ 壹 | ❷ 弍 | ❸ 壱 |
|---|---|---|---|---|
| | Unicode | U58F9 | U5F0C | U58F1 |
| | Disposition | allocate | block | block |

## Intersection → Integrated LGR

| 一 U4E00 → | Variant | | | |
|---|---|---|---|---|
| | Unicode | | | |
| | Disposition | | | |

# CJK Integration Methodology
## Divide & Conquer (D&C)



**Diversified CJK Demands**

C Demands

J Demands

K Demands

Plan and Define

**Strategic Direction**

**Unified CJK Rules**

Root Zone Admin

CJK Rules

CJK Repertoire

Requires

Requires

Evaluation Method

LGR Constrains

MSR

Minimal Viable Solution

Variant Dispositions

Repertoire

Split

JK Overlap

CJK Overlap

CJ Overlap

CK Overlap

CJ Usage Pattern

CK Usage Pattern

Repertoire

Repertoire

Merge

# Splitting Non-overlapping Code Points From Repertories

**CJK Han-overlap in IANA IDN Repository**

$R_k$

J-Han : 6356 (JPRS)     K-Han : 0 (KRNIC)

$R_j$  175

C/J Overlap: **6181**

$R_C$  13339

C-Han : 19520 (CNNIC/TWNIC)

$R_C$
Chinese LGR

$R_j$
Japanese LGR

$R_k$
Korean LGR

Develop Conflict Strategy     No conflict

| unified code points | |
|---|---|
| | 13339 |
| + | 175 |
| | 13514 |

*Problem Domain (Unsolved Overlap) : 6181*

# Engineering Design

**Computation for Word Usage and Frequency**



Matching → **C/J overlap code points**

TC : Apple News
SC : Sina News
JP : Mainichi News

*Sample size is statistical significant*

**usage**

**frequency of use**

Split unused code points

Split code points of low frequency of use

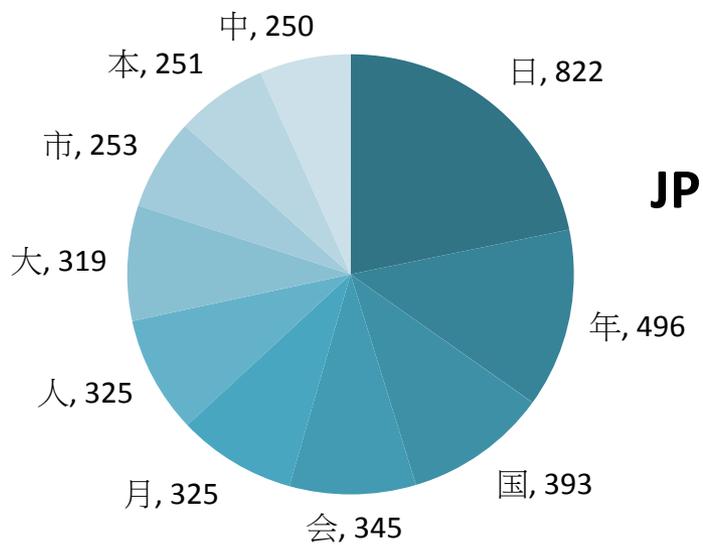# Splitting Unused Code Points from The Overlap

**C / J Overlap Data Set : 6181**
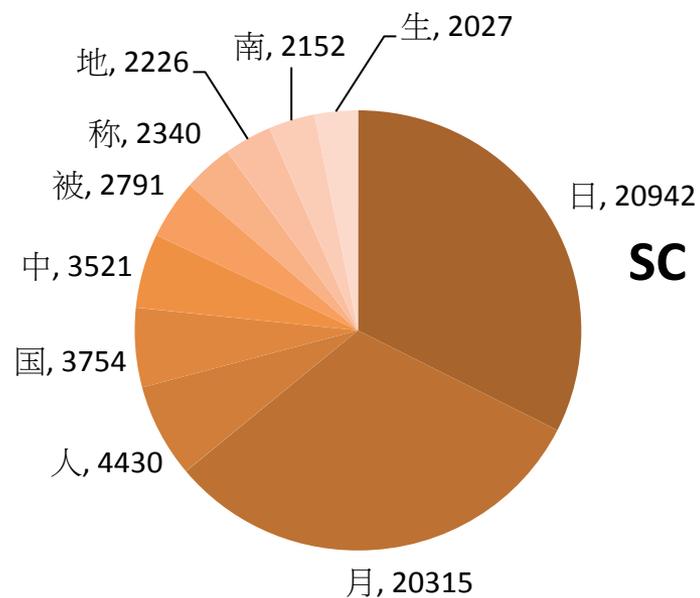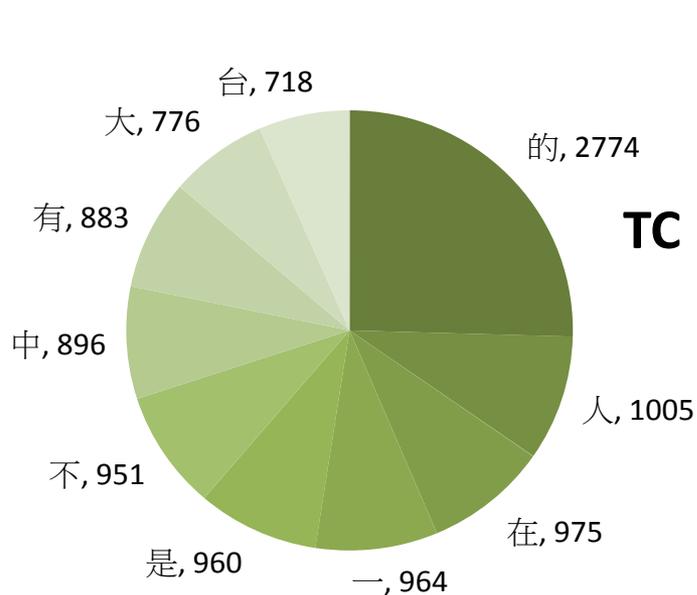
total unused : 2739

$R_j$ J only : 203

C / J usage overlap : *1312*

$R_c$ C only : 1927

total used : 3442

| unified code points |
|---|
| 2739 |
| 203 |
| +       1927 |
| 4869 |

*Problem Domain (Unsolved Overlap) : 1312*

# Computing Frequency of Use of Code Points

Initial Data Set : 1312

# Top 10 Most Popular Words

Top 20 : Chinese Frequency of Use > Japanese Frequency of Use

Generated Data Set : 939

**Top 20 : Chinese Frequency of Use < Japanese Frequency of Use**

Generated Data Set : 363

Frequency of Use %

Legend: C-Freq, J-Freq

| Character | C-Freq | J-Freq |
|-----------|--------|--------|
| 7530 | | 0.1988 |
| 672C | | 0.2788 |
| 5E74 | | 0.5512 |
| 6771 | | 0.1956 |
| 4F1A | | 0.3834 |
| 9023 | | 0.1644 |
| 540C | | 0.2366 |
| 6C0F | | 0.1056 |
| 7A0E | | 0.1212 |
| 91CE | | 0.1088 |
| 898B | | 0.1688 |
| 6C7A | | 0.1422 |
| 7C73 | | 0.1344 |
| 5316 | | 0.1622 |
| 5CF6 | | 0.0856 |
| 5E02 | | 0.2812 |
| 8B70 | | 0.1588 |
| 52DD | | 0.0912 |
| 53D6 | | 0.1134 |
| 4EE3 | | 0.1288 |

15

Chinese Frequency of Use = Japanese Frequency of Use

Generated Data Set : 10

# Frequency of Use Reassembly

C / J Usage Overlap Data Set : 1312

Freq J > C : 363

**R**$_j$

J = C
10

**R**$_c$    Freq C > J : 939

| unified code points |
|---|
| 363 |
| +    939 |
| 1302 |

*Problem Domain (Unsolved Overlap) : 10*

# Data Processing & Computation Recap



**Filtering Process**

**Filtering Process**

**>20K Han Code Points**
Splitting Non-overlapping

**6181 CJK Overlap**
Splitting Unused

**1312 Usage Overlap**
*Frequency of Use Computation*

LOGIC Design

10 Code Points

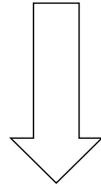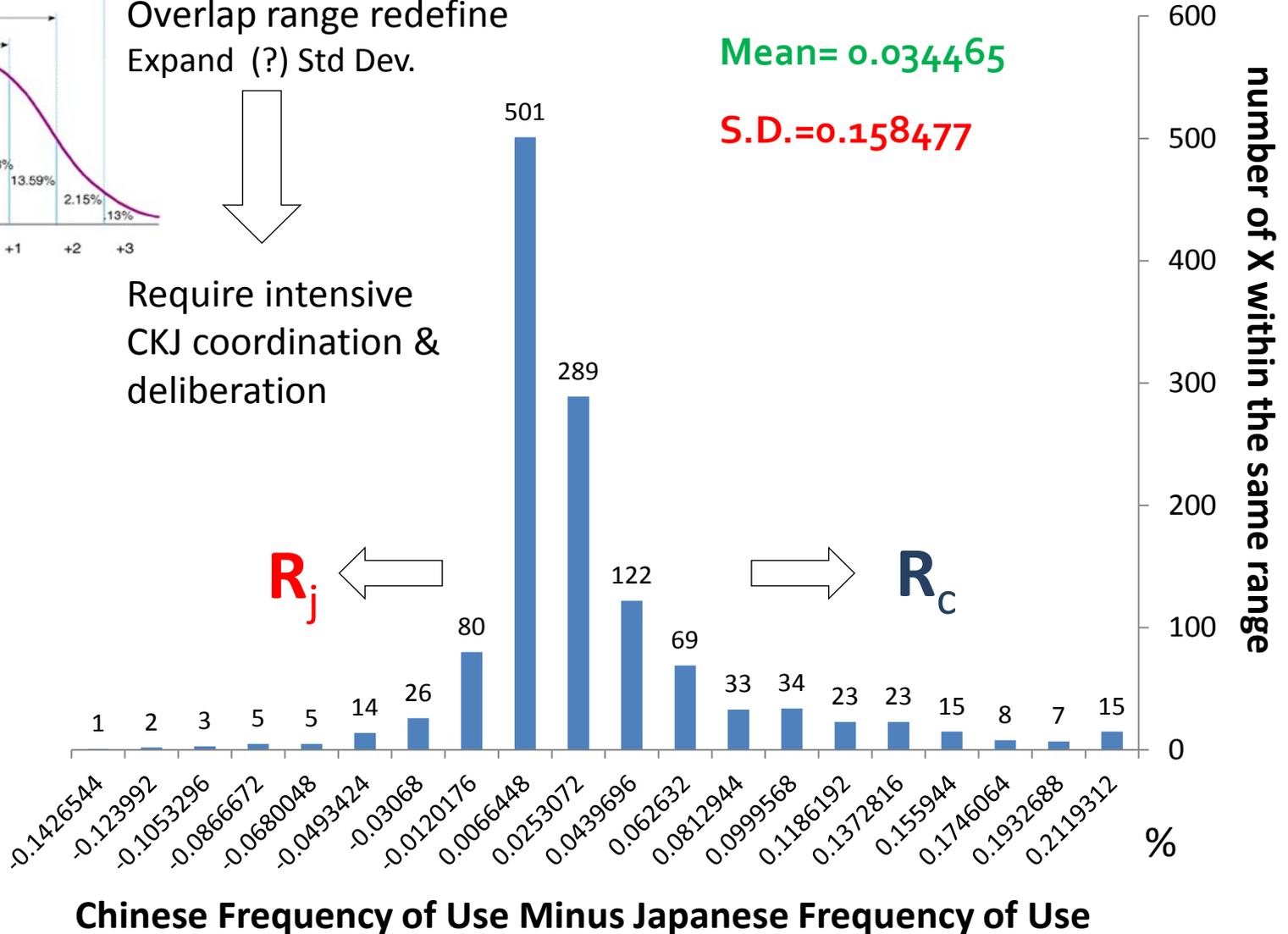| Methodology Review | CJK Coordination | Re-Sampling & Computation | Statistical Justification |

*Problem domain was effectively reduced*

# Future Work



Overlap range redefine
Expand (?) Std Dev.

Require intensive
CKJ coordination &
deliberation

Mean= 0.034465

S.D.=0.158477

$R_j$

$R_c$

number of X within the same range

%

**Chinese Frequency of Use Minus Japanese Frequency of Use**

# Re-consider Language Tag



Language tag support

Sources of Language Tag

- RFC 2860 : The name space of language tags is administered by IANA
- ISO Standard 639 :
  - when a language has both an IANA-registered tag and a tag derived from an ISO registered code, one MUST use the ISO tag.
  - Maintenance Agency : International Information Centre for Terminology (Austria)

## Perfection Syndrome

"Engineering isn't about perfect solutions; it's about doing the best you can with limited resources." Randy Pausch

Thank You
Question?

Kenny Huang, Ph.D.
huangksh@gmail.com