# String Similarity Small Group Outcome

# Agenda

❖ Overview

❖ Task 1

❖ Task 2

❖ Task 3

❖ Conclusion

# Overview

# Background

**Charter Questions**

EPDP-IDN Charter asks to consider any adjustment to the string similarity review due to the variant implementation: (Charter Question E3)
- What role, if any, do the "withheld same entity" variants play? (Charter Question E1)
- What are the potential consequences for the other allocatable variant labels in the same set of a requested variant label, which is rejected as a result of the string similarity review? (Charter Question E3a)

**Staff Paper Recommendation**

String similarity review should compare strings under consideration not just against all allocated or applied-for strings, but also all variants of those strings (i.e., allocatable, withheld-same-entity, and blocked).

**EPDP Team Discussion**

The EPDP Team discussed three (3) possible levels of comparison among visually confusable strings, as well as analyzed the impact and potential consequences:

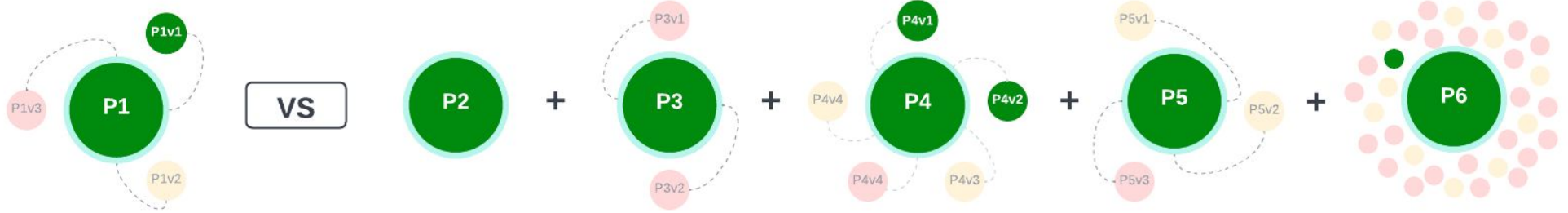Level 1: Primary + only requested allocatable variants

Level 2: Primary + all allocatable variants

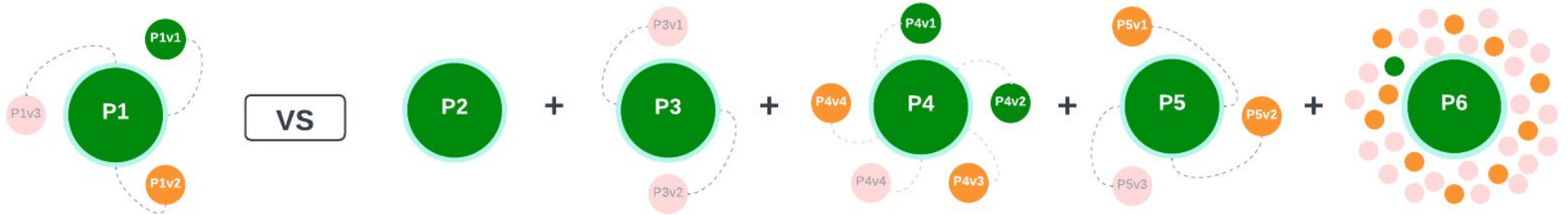Level 3: Source gTLD + all valid variants (blocked + allocatable)

***Primary**: The applied-for gTLD that also serves as the source gTLD for calculating its allocatable and blocked variants during the application process; the applicant may request to activate none, one, or more allocatable variants of such an applied-for gTLD.*

# Three Levels of Comparison



**Level 1** — Primary + ONLY Requested Allocatable Variants

P1 **VS** P2 + P3 + P4 + P5 + P6

**Level 2** — Primary + ALL Allocatable Variants

P1 **VS** P2 + P3 + P4 + P5 + P6

**Level 3** — Primary + ALL Allocatable and Blocked Variants

P1 **VS** P2 + P3 + P4 + P5 + P6

🟢 Requested Allocatable Label    🟠 Non-Requested Allocatable Label    🔴 Blocked Label

# Problem Statement

**String Similarity Small Group has been set up to tackle the following problems:**

*Problem 1: There is a divergence of opinions regarding which level is the most appropriate*

*Problem 2: The discussion has been largely academic based on abstract concepts*

# Small Group Tasks

Facilitate a comprehensible discussion by ***developing concrete examples of variants that are visually confusable***

**Task 1**: Develop ***concrete examples of strings*** that have blocked and/or allocatable variant labels and may be visually confusable with other strings in the same scripts or across scripts

- Develop practical examples – limit to visual similarity – that could happen in reality & indicate how feasible/possible such cases could happen
- Discuss whether any existing mechanisms that could help prevent such confusingly similar strings being delegated

**Task 2**: Demonstrate ***how these examples would be compared against each other in the string similarity review according to the three levels***, showcasing the impact on the review and the potential consequences

- Propose a String Similarity Review model with the view of minimizing security, stability, and user confusability risks

**Task 3**: Demonstrate ***how these examples would undergo the objection process according to the three levels***, showcasing the impact on the objection process and the potential consequences

- Identify which type of variants should be subject to the objection process

**Exclusion:** Complexity in implementation for Tasks 2 and 3 is out of scope – defer deliberation to EPDP Team

# Small Group Composition

| Member | Affiliation | Language Proficiency |
|--------|-------------|----------------------|
| Edmon Chung | Board Liaison | Chinese (Mandarin, Cantonese) |
| Hadia El miniawi | ALAC | Arabic |
| Imran Hossen | Independent | Bangla |
| Jerry Sen | RySG | Chinese (Mandarin) |
| Justine Chew (Small Group Lead) | ALAC | Malay |
| Michael Bauland | RrSG | German |
| Wael Nasr | Independent | Arabic |

**Note:**

- Between 18 May 2022 and 3 August 2022, the Small Group held a total of 10 meetings
- Small Group agreed to the 3 tasks stated in the assignment form during its first meeting on 18 May 2022
- Supported by ICANN staff with additional language proficiency
- Wael Nasr joined toward the end of small group deliberation

# Task 1

*Develop concrete examples of strings that have blocked and/or allocatable variant labels and may be visually confusable with other strings in the same scripts or across scripts*

# Example Strings

The group developed **eight (8) examples**, as contributed by both members and staff, and discussed their <span style="color:blue">primary</span>**,** <span style="color:green">allocatable</span> and, <span style="color:red">**blocked**</span> variants calculated by RZ-LGR

| No. | Label A | Label B | Label C | Practicality Consideration |
|-----|---------|---------|---------|---------------------------|
| 1 | Latin **bıß** | Cyrillic **вiss** | | Valid strings per RZ-LGR |
| 2 | Traditional Chinese 滙豐 | Simplified Chinese 汇丰 | | Real Chinese words with same meanings and variant relationship |
| 3 | Arabic بنئ | Arabic بنى | | Valid strings per RZ-LGR with at least one string that's meaningful in a language |
| 4 | Simplified Chinese 华鸟 | Traditional Chinese 华島 | | Real Chinese words with different meanings |
| 5 | Latin **rıch** | Latin **ṅch** | | Valid strings per RZ-LGR |
| 6 | Arabic ركى | Arabic رکﮯ | | Valid strings per RZ-LGR with at least one string that's meaningful in a language |
| 7 | Simplified Chinese 华为 | Simplified Chinese 华鸟 | Simplified Chinese 华岛 | Real Chinese words with different meanings |
| 8 | Japanese Kanji 一休 | Traditional Chinese 一體 | | Real Japanese and Chinese words with different meanings |

# Task 2

*Demonstrate how these examples would be compared against each other in the string similarity review according to the three levels, showcasing the impact on the review and the potential consequences*
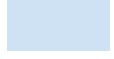
# Selected Examples for Comparison

| No. | Label A | Label B | Label C |
|---|---|---|---|
| 1 | Latin bıß | Cyrillic вiss | |
| 2 | Traditional Chinese 滙豐 | Simplified Chinese 汇丰 | |
| 3 | Arabic بنئ | Arabic بنى | |
| 4 | Simplified Chinese 华鸟 | Traditional Chinese 华島 | |
| 5 | Latin rıch | Latin ṅch | |
| **6** | **Arabic رکی** | **Arabic رگے** | |
| **7** | **Simplified Chinese 华为** | **Simplified Chinese 华鸟** | **Simplified Chinese 华岛** |
| 8 | Japanese Kanji 一休 | Traditional Chinese 一體 | |

Demonstrate why hybrid model is recommended

Demonstrate
- Applied-for gTLD vs. Existing gTLD
- Comparison among three strings
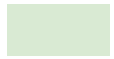
# Example 6: Two Applied-for Arabic TLDs

**Applied-for Primary Strings:**

رکی (A1)

رگے (B1)

**Allocatable Variants of Primary Strings:**

رکی (A2)

رگی (A3)

None

**Blocked Variants of Primary Strings:**

| | | | |
|---|---|---|---|
| رکی (A15) | رکئ (A4) | رگّب (B13) | رگی (B2) | رگؠ (B24) |
| رکی (A16) | رکي (A5) | رگی (B14) | رگی (B3) | رگے (B25) |
| رکے (A17) | رکّب (A6) | رگی (B15) | رگي (B4) | رگئ (B26) |
| رگئ (A18) | رکی (A7) | رگی (B16) | رگّب (B5) | رگی (B27) |
| رگي (A19) | رکی (A8) | رگے (B17) | رگی (B6) | رگي (B28) |
| رکّب (A20) | رکي (A9) | رگئ (B18) | رگّب (B7) | رگی (B29) |
| رگی (A21) | رکے (A10) | رگی (B19) | رگي (B8) | رگی (B30) |
| رگی (A22) | رکئ (A11) | رگي (B20) | رگی (B9) | رگی (B31) |
| رگي (A23) | رکي (A12) | رگّب (B21) | رگئ (B10) | رگي (B32) |
| رگے (A24) | رکّب (A13) | رگی (B22) | رگی (B11) | |
| | رکی (A14) | رگي (B23) | رگی (B12) | |

# Example 6: String Similarity Review

① → (B1) رگے

(A1) رکی

② (red)

(B2) رگئ   (B13) رگَب   (B24) رگِي
(B3) رگی   (B14) رگَے   (B25) رگی
(B4) رگِي   (B15) رگُئ   (B26) رگئ
(B5) رگَب   (B16) رگُي   (B27) رگی
(B6) رگی   (B17) رگَے   (B28) رگِي
(B7) رگی   (B18) رگُئ   (B29) رگَب
(B8) رگِي   (B19) رگی   (B30) رگی
(B9) رگَے   (B20) رگِي   (B31) رگی
(B10) رگُئ   (B21) رگُب   (B32) رگِي
(B11) رگی   (B22) رگی
(B12) رگُي   (B23) رگی

③ → (B1) رگے

(A2) رکی
(A3) رکی

④ (red)

(B2) رگئ   (B13) رگُب   (B24) رگِي
(B3) رگی   (B14) رگَے   (B25) رگی
(B4) رگِي   (B15) رگُئ   (B26) رگئ
(B5) رگَب   (B16) رگُي   (B27) رگی
(B6) رگی   (B17) رگَے   (B28) رگِي
(B7) رگی   (B18) رگُئ   (B29) رگَب
(B8) رگِي   (B19) رگی   (B30) رگی
(B9) رگَے   (B20) رگُي   (B31) رگی
(B10) رگُئ   (B21) رگُب   (B32) رگِي
(B11) رگی   (B22) رگی
(B12) رگُي   (B23) رگی

⑤ (yellow) → (B1) رگے

(A15) رکی   (A4) رکئ
(A16) رکِي   (A5) رکي
(A17) رکَے   (A6) رکُب
(A18) رکُئ   (A7) رکی
(A19) رکِي   (A8) رکی
(A20) رکُب   (A9) رکِي
(A21) رکی   (A10) رکَے
(A22) رکی   (A11) رکُئ
(A23) رکَي   (A12) رکِي
(A24) رکے
(A13) رکُب
(A14) رکی

# Example 6: String Similarity Review (Cont.)



**String Similarity Review may find the following confusingly similar strings**

(2) رکی (A1) & رگی (B3) & رگی (B6)

(4) رکی (A2) & رگی (B3) & رگی (B6)

(4) رکیے (A3) & رگی (B3) & رگی (B6)

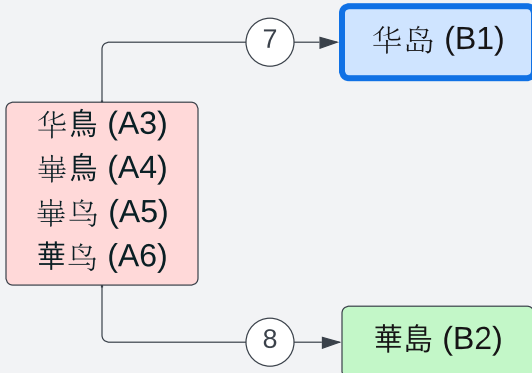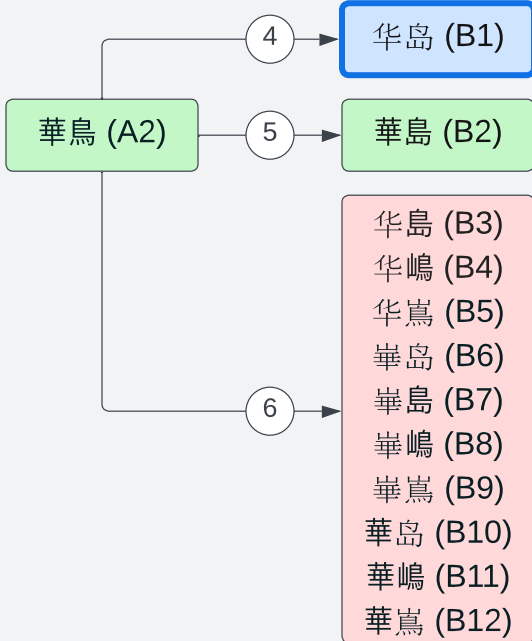(5) رگے (B1) & رکیے (A10) & رکیے (A17) & رکیے (A24)

**Potential Outcome of the String Similarity Review**

رکی (A1) & its variants A2-A24 AND رگے (B1) & its variants B2-B32 get processed in a contention set

**If the hybrid model were not used and blocked variants were not taken into account in String Similarity Review**

رکی (A1) and رگے (B1) would have been both delegated with the misconnection risk. E.g., a user may mistake رکی (A1) as رگی (B3), a blocked variant of رگے (B1), but arrive at site controlled by a registrant different to رگے (B1).

## Scenario 1: String Similarity Review of Applied-for String A1 & Existing String B1
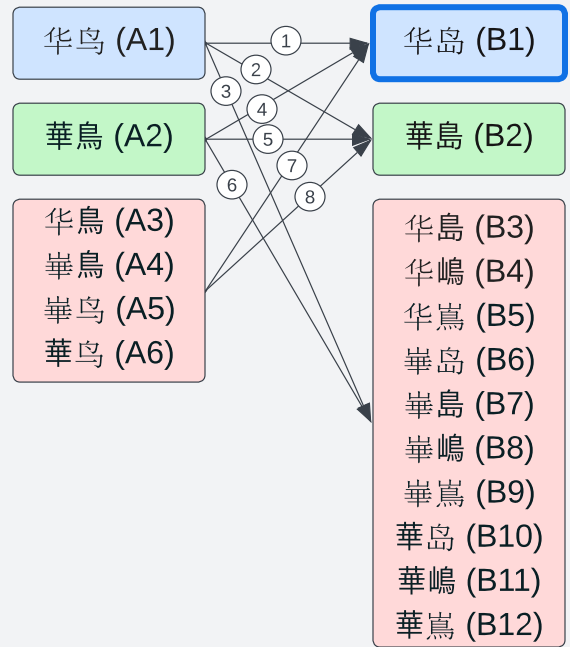
华鸟 (A1) —①→ 华岛 (B1)

华鸟 (A1) —②→ 華島 (B2)

华鸟 (A1) —③→
华島 (B3)
华嶋 (B4)
华嵨 (B5)
崋島 (B6)
崋嶋 (B7)
崋嶋 (B8)
崋嵨 (B9)
華島 (B10)
華嶋 (B11)
華嵨 (B12)

華鳥 (A2) —④→ 华岛 (B1)

華鳥 (A2) —⑤→ 華島 (B2)

華鳥 (A2) —⑥→
华島 (B3)
华嶋 (B4)
华嵨 (B5)
崋島 (B6)
崋嶋 (B7)
崋嶋 (B8)
崋嵨 (B9)
華島 (B10)
華嶋 (B11)
華嵨 (B12)

华鳥 (A3)
崋鳥 (A4)
崋鸟 (A5)
華鸟 (A6) —⑦→ 华岛 (B1)

华鳥 (A3)
崋鳥 (A4)
崋鸟 (A5)
華鸟 (A6) —⑧→ 華島 (B2)

## Scenario 1: Consolidated View

华鸟 (A1)
華鳥 (A2)
华鳥 (A3)
崋鳥 (A4)
崋鸟 (A5)
華鸟 (A6)

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

华岛 (B1)
華島 (B2)
华島 (B3)
华嶋 (B4)
华嵨 (B5)
崋島 (B6)
崋嶋 (B7)
崋嶋 (B8)
崋嵨 (B9)
華島 (B10)
華嶋 (B11)
華嵨 (B12)

### String Similarity Review may find the following confusingly similar pairs

① 华鸟 (A1) & 华岛 (B1)

⑤ 華鳥 (A2) & 華島 (B2)

⑥ 華鳥 (A2) & 崋島 (B7)
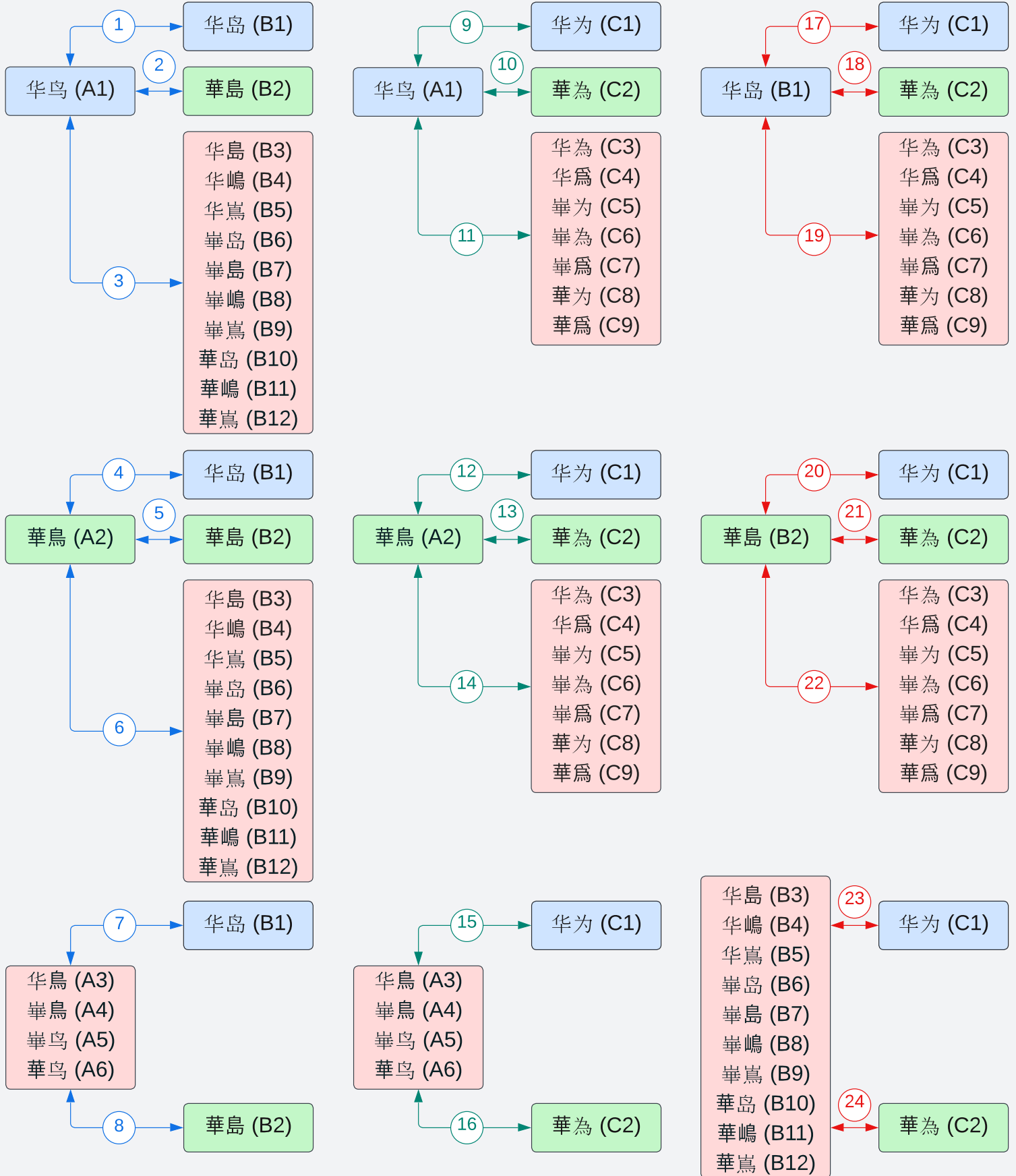
⑧ 崋鳥 (A4) & 華島 (B2)

### Potential Outcome of String Similarity Review

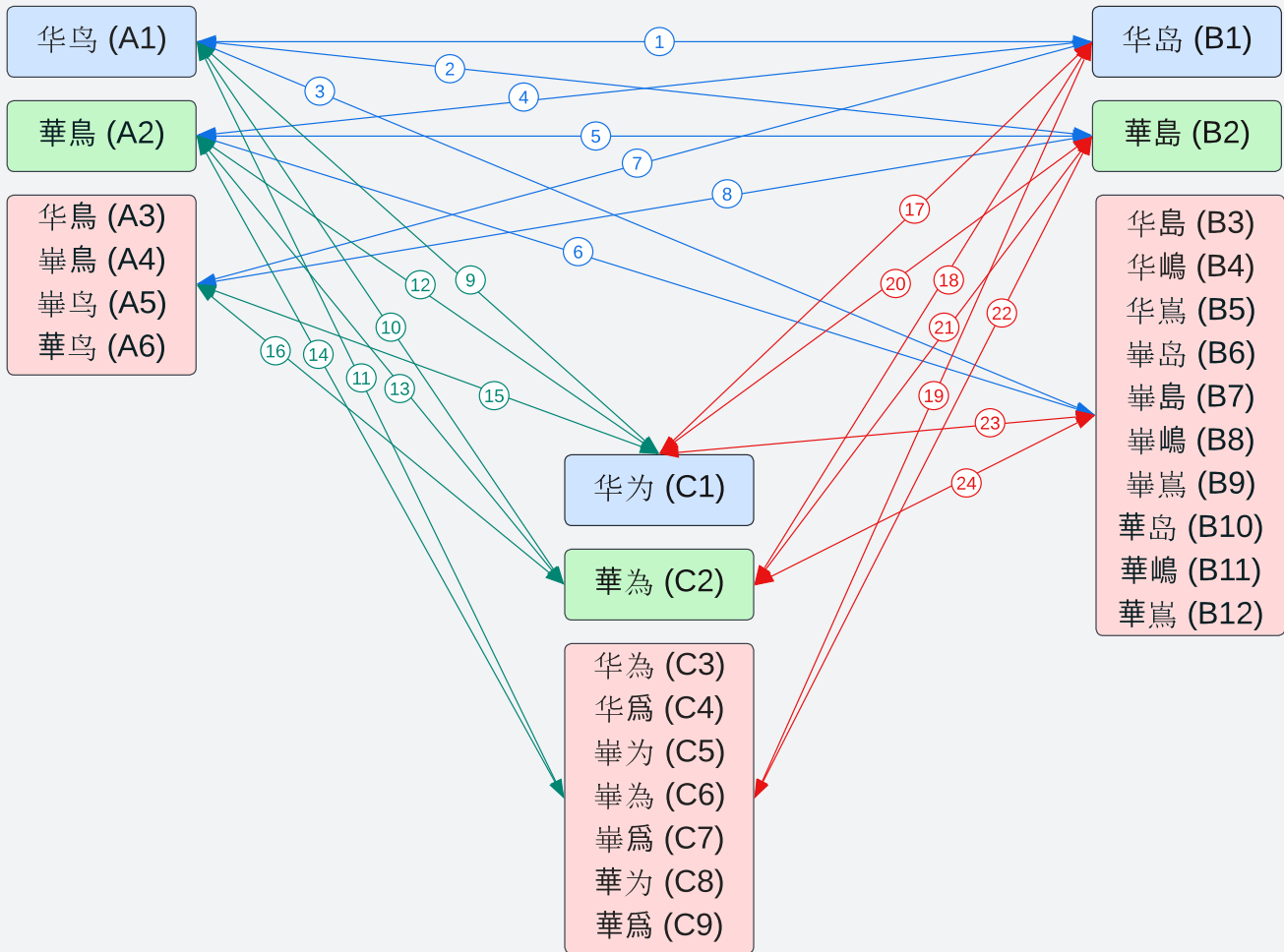华鸟 (A1) may be rejected due to its confusing similarity to the already delegated 华岛 (B1)

### Legend

Primary applied-for string

Primary existing string

Allocatable variant of primary string

Blocked variant of primary string

# Scenario 2: String Similarity Review of Applied-for Strings A1, B1 & C1

华鸟 (A1)
① → 华岛 (B1)
② ↔ 華島 (B2)
③ → 华島 (B3) / 华嶋 (B4) / 华嵨 (B5) / 崋島 (B6) / 崋島 (B7) / 崋嶋 (B8) / 崋嵨 (B9) / 華岛 (B10) / 華嶋 (B11) / 華嵨 (B12)

华鸟 (A1)
⑨ → 华为 (C1)
⑩ ↔ 華為 (C2)
⑪ → 华为 (C3) / 华爲 (C4) / 崋为 (C5) / 崋為 (C6) / 崋爲 (C7) / 華为 (C8) / 華爲 (C9)

华岛 (B1)
⑰ → 华为 (C1)
⑱ ↔ 華為 (C2)
⑲ → 华为 (C3) / 华爲 (C4) / 崋为 (C5) / 崋為 (C6) / 崋爲 (C7) / 華为 (C8) / 華爲 (C9)

華鳥 (A2)
④ → 华岛 (B1)
⑤ ↔ 華島 (B2)
⑥ → 华島 (B3) / 华嶋 (B4) / 华嵨 (B5) / 崋島 (B6) / 崋島 (B7) / 崋嶋 (B8) / 崋嵨 (B9) / 華岛 (B10) / 華嶋 (B11) / 華嵨 (B12)

華鳥 (A2)
⑫ → 华为 (C1)
⑬ ↔ 華為 (C2)
⑭ → 华为 (C3) / 华爲 (C4) / 崋为 (C5) / 崋為 (C6) / 崋爲 (C7) / 華为 (C8) / 華爲 (C9)

華島 (B2)
⑳ → 华为 (C1)
㉑ ↔ 華為 (C2)
㉒ → 华为 (C3) / 华爲 (C4) / 崋为 (C5) / 崋為 (C6) / 崋爲 (C7) / 華为 (C8) / 華爲 (C9)

华鳥 (A3) / 崋鳥 (A4) / 崋鸟 (A5) / 華鸟 (A6)
⑦ → 华岛 (B1)
⑧ → 華島 (B2)

华鳥 (A3) / 崋鳥 (A4) / 崋鸟 (A5) / 華鸟 (A6)
⑮ → 华为 (C1)
⑯ → 華為 (C2)

华島 (B3) / 华嶋 (B4) / 华嵨 (B5) / 崋島 (B6) / 崋島 (B7) / 崋嶋 (B8) / 崋嵨 (B9) / 華岛 (B10) / 華嶋 (B11) / 華嵨 (B12)
㉓ ↔ 华为 (C1)
㉔ ↔ 華為 (C2)

**Scenario 2**: Consolidated View

华鸟 (A1)

華鳥 (A2)

华鳥 (A3)
峀鳥 (A4)
峀鸟 (A5)
華鸟 (A6)

华为 (C1)

華為 (C2)

华為 (C3)
华爲 (C4)
峀为 (C5)
峀為 (C6)
峀爲 (C7)
華为 (C8)
華爲 (C9)

华岛 (B1)

華島 (B2)

华島 (B3)
华嶋 (B4)
华嵓 (B5)
峀島 (B6)
峀島 (B7)
峀嶋 (B8)
峀嵓 (B9)
華島 (B10)
華嶋 (B11)
華嵓 (B12)

**String Similarity Review may find the following confusingly similar pairs**

① 华鸟 (A1) & 华岛 (B1)

⑤ 華鳥 (A2) & 華島 (B2)

⑥ 華鳥 (A2) & 峀島 (B7)

⑧ 峀鳥 (A4) & 華島 (B2)

⑬ 華鳥 (A2) & 華為 (C2)

⑭ 華鳥 (A2) & 峀為 (C6)

⑯ 峀鳥 (A4) & 華為 (C2)

㉑ 華島 (B2) & 華為 (C2)

㉔ 峀島 (B7) & 華為 (C2)

**Potential Outcome of String Similarity Review**

华鸟 (A1) & its variants A2-A6 AND 华岛 (B1) & its variants B2-B12 AND 华为 (C1) & its variants C2-C9 get processed in a contention set

**Legend**

Primary applied-for string

Allocatable variant of primary string

Blocked variant of primary string

# Recommendation: Hybrid Model

**Summary:** *The small group recommends the **hybrid model**, which is a **mixed-level approach between level 2 and level 3***

**Goal:** *Mitigate any possibility of confusing similarity between one IDN TLD and another IDN TLD or any of its valid variant(s), vice versa*

*In practice, the string similarity review must be modified to compare:*

- **An applied-for primary IDN gTLD and <u>all of its allocatable variant label(s)</u>**

*Against:*

- **Existing TLDs and <u>all of their allocatable and blocked variant labels</u>;**

- **Strings requested as IDN ccTLDs and <u>all of their allocatable and blocked variant labels</u>;**

- **Other applied-for gTLDs in the same round and <u>all of their allocatable and blocked variant labels</u>;**

- **Reserved Names; and**

- **Any other two-character ASCII strings and <u>all of their allocatable and blocked variant labels</u> (*if the applied-for primary IDN gTLD is a two-character string*)**

*In addition, compare:*

- **<u>All of the blocked variant label(s)</u> of an applied-for primary IDN gTLD**

*Against:*

- **Existing TLDs and <u>all of their allocatable variant labels</u>**

**Note:** *Blocked variants of one IDN TLD should NOT be compared against blocked variants of another IDN TLD*

# Rationale for Hybrid Model

Considering the limited scope of security, stability and user confusability, the small group believes the <u>hybrid model</u> would:

- **Be sufficiently conservative** and can **help mitigate two types of failure modes** – denial of service and misconnection, which may have a higher likelihood to affect non-native speakers of certain scripts or languages

- **Help detect many more pairs of visually confusable strings** and **reduce the risks of failure modes**

- **Reduce computational complexity by not requiring comparison of blocked variant labels** of a primary applied-for IDN gTLD string against blocked variant labels of other existing and applied-for TLD strings

The small group also believes that:

- Level 1 and 2 may fail to detect some visually confusable strings and increase the risks of failure modes
- Level 3 unnecessarily compares blocked variants against each other with exponential increase of computational complexity

**Additional Considerations**

- While the pool of strings that needs to be considered will be large, **language experts in the String Similarity Review panel can evaluate the strings on a case-by-case basis**

- After the evaluation completes, there are **other mechanisms in the New gTLD Program** – e.g., limited appeal mechanism and objection processes – to review the string similarity panel's decision

# Conclusion

# String Similarity Review Recommendation

**Summary:** *The small group recommends the **hybrid model**, a **mixed-level approach between level 2 and level 3***

*The string similarity review must be modified to compare:*

- **An applied-for primary IDN gTLD and <u>all of its allocatable variant label(s)</u>**

*Against:*

- **Existing TLDs and <u>all of their allocatable and blocked variant labels</u>;**

- **Strings requested as IDN ccTLDs and <u>all of their allocatable and blocked variant labels</u>;**

- **Other applied-for gTLDs in the same round and <u>all of their allocatable and blocked variant labels</u>;**

- **Reserved Names; and**

- **Any other two-character ASCII strings and <u>all of their allocatable and blocked variant labels</u> (*if the applied-for primary IDN gTLD is a two-character string*)**

*In addition, the string similarity review must be modified to compare:*

- **<u>All of the blocked variant label(s)</u> of an applied-for primary IDN gTLD**

*Against:*

- **Existing TLDs and <u>all of their allocatable variant labels</u>**