# SaudiNIC's Proposed Solution

## Registry-level Multilingual Arabic Script IDN Registration

**MENOG 8, Khobar, May 14-18, 2011**

Presented by:

AbdulRahman  Al-Ghadir

Developed by  SaudiNIC/CITC Staff:

AbdulAziz  Al-Zoman

Raed  Al-Fayez
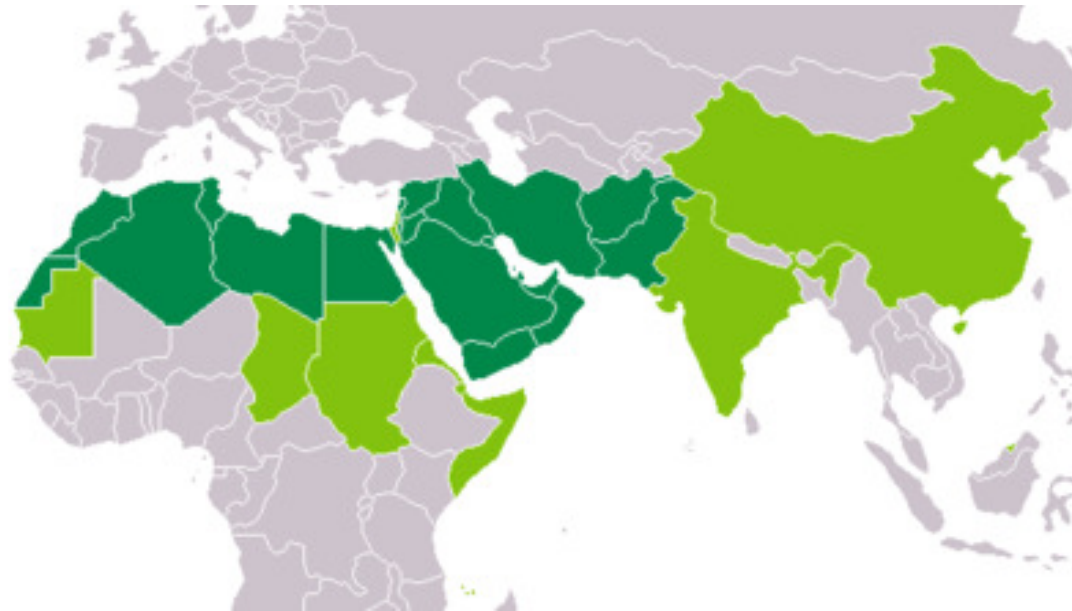
AbdulRahman  Al-Ghadir

# Content

- **Arabic Script Major Issues**

- **Confusing Similar Characters**

- **Proposed solution**
  - Characteristics
  - Language-level required tables
  - Language-level required tables
  - Language-level operation

- **Conclusion**

# Arabic Script

- The **2nd** most widely used alphabetic writing system in the world
- Used by **many languages** such as:
  - Persian, Urdu, Turkish, Kurdish, Pashto, Jawi, ...



Source: http://en.wikipedia.org/wiki/Arabic_script

# Arabic Script IDN - Major Issues

## IDNA Table
### Pvalid, Disallowed, ContextO



## Non-spacing Marks



## bidirectional



## Digit

| | | |
|---|---|---|
| 1.European digits | U+0030 .. U+0039 | (0123456789) |
| 2.Arabic-Indic digits | U+0660 .. U0669 | (٠١٢٣٤٥٦٧٨٩) |
| 3.Eastern Arabic-Indic digits | U+06F0 .. U+06F9 | (٠١٢٣۴۵۶٧٨٩) |

## ZWNJ/ZWJ

| Examples not using ZWNJ | Examples not using ZWNJ |
|---|---|
| طبل | طبل |
| input[0] = U+0637 | input[0] = U+0637 |
| input[1] = U+0628 | input[1] = U+200c |
| input[2] = U+0644 | input[2] = U+0628 |
| | input[3] = U+0644 |

## Combining Marks

| ى | + | ٔ | = | ئ ٔ ـٔ ئ | is confusing with | ئ ٔ ـٔ ئ |
|---|---|---|---|---|---|---|
| U+0649 | | U+0654 | | U+0649 U+0654 | | U+0626 |
| Description: Alef Maksura + Hamza Above <> Yeh With Hamza Above | | | | | | |
| Comments: This is a Unicode confusable! | | | | | | |

| ى | + | ٔ | = | ئ ٔ ـٔ ئ | is confusing with | ئ ٔ ـٔ ئ |
|---|---|---|---|---|---|---|
| U+06cc | | U+0654 | | U+06cc U+0654 | | U+0626 |
| Description: Farsi Yeh + Hamza Above <> Yeh With Hamza Above | | | | | | |
| Comments: This is Unicode confusable! | | | | | | |

- **There are a number of groups of characters that have the same shapes (Homoglyph).**
  - eg. Kaf, Heh, Yeh, Alef, … groups



**كلى.com**

Site # **1**

Arabic Label (U-Label) : كلى

Unicode (U+) : U+06CC , U+0644 , U+06A9

ASCII Label (A-Label) : xn--ghb5rwd.com

Now see where this domain will go : كلى.com

Why the domain name was registered :

1. To illustrate that an IDN script-based registration (particularly with Arabic Script) is a good ground for phishing.
2. To encourage the international communities in general and language communities ( that use Arabic) in particular to work together in finding solutions to these kind of problems.
3. To speed up the process for ccTLD IDN fast track where each IDN ccTLD will serve only one language.
4. If new IDN-gTLDs will be open, then they should focus on languages rather than scripts till a solution has been reached (see point 2)

For more information see this the Power Presentation that was delivered on the ICANN regional meeting, Dubai 1-3 April 2008 [will be available soon]

| input[0] = U+06a9 |
| input[1] = U+0644 |
| input[2] = U+06cc |

| input[0] = U+0643 |
| input[1] = U+0644 |
| input[2] = U+0649 |

**كلى.com**

Site # **2**

Arabic Label (U-Label) : كلى

Unicode (U+) : U+0643, U+0644, U+0649

ASCII Label (A-Label) : xn--fhbcp.com

Now see where this domain will go : كلى.com

Why the domain name was registered :

1. To illustrate that an IDN script-based registration (particularly with Arabic Script) is a good ground for phishing.
2. To encourage the international communities in general and language communities ( that use Arabic script) in particular to work together in finding solutions to these kind of problems.

| | 060 | 061 | 062 | 063 | 064 | 065 | 066 | 067 | 068 | 069 | 06A | 06B | 06C | 06D | 06E | 06F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0600 | 0610 | ▨ | ذ 0630 | ـ 0640 | 0650 | ٠ 0660 | 0670 | پ 0680 | 0690 | غ 06A0 | گ 06B0 | ۀ 06C0 | ی 06D0 | 06E0 | ۰ 06F0 |
| 1 | 0601 | 0611 | ء 0621 | ر 0631 | ف 0641 | 0651 | ١ 0661 | آ 0671 | خ 0681 | ڑ 0691 | ڡ 06A1 | ڱ 06B1 | ۃ 06C1 | ۑ 06D1 | 06E1 | ۱ 06F1 |
| 2 | 0602 | 0612 | آ 0622 | ز 0632 | ق 0642 | 0652 | ٢ 0662 | أ 0672 | خ 0682 | ڒ 0692 | ب 06A2 | ڲ 06B2 | ۂ 06C2 | 06D2 | 06E2 | ۲ 06F2 |
| 3 | 0603 | 0613 | أ 0623 | س 0633 | ك 0643 | 0653 | ٣ 0663 | إ 0673 | ج 0683 | ړ 0693 | ڣ 06A3 | ڳ 06B3 | ۃ 06C3 | 06D3 | 06E3 | ۳ 06F3 |
| 4 | | 0614 | ؤ 0624 | ش 0634 | ل 0644 | 0654 | ٤ 0664 | 0674 | چ 0684 | ڔ 0694 | ڤ 06A4 | گ 06B4 | و 06C4 | 06D4 | 06E4 | ۴ 06F4 |
| 5 | | ط 0615 | إ 0625 | ص 0635 | م 0645 | 0655 | ٥ 0665 | ڊ 0675 | ڇ 0685 | ڕ 0695 | پ 06A5 | ل 06B5 | ۅ 06C5 | �ە 06D5 | 06E5 | ۵ 06F5 |
| 6 | ؆ 0606 | 0616 | ئ 0626 | ض 0636 | ن 0646 | 0656 | ٦ 0666 | ؤ 0676 | چ 0686 | ږ 0696 | ڦ 06A6 | ڶ 06B6 | ۆ 06C6 | 06D6 | 06E6 | ۶ 06F6 |
| 7 | ؇ 0607 | 0617 | ا 0627 | ط 0637 | ه 0647 | 0657 | ٧ 0667 | ؤ 0677 | ڈ 0687 | ڗ 0697 | ڧ 06A7 | ڷ 06B7 | ۇ 06C7 | 06D7 | 06E7 | ۷ 06F7 |
| 8 | ؈ 0608 | 0618 | ب 0628 | ظ 0638 | و 0648 | 0658 | ٨ 0668 | ئى 0678 | ڈ 0688 | ژ 0698 | ڨ 06A8 | ڸ 06B8 | و 06C8 | 06D8 | 06E8 | ۸ 06F8 |
| 9 | ؉ 0609 | 0619 | ة 0629 | ع 0639 | ى 0649 | 0659 | ٩ 0669 | ڙ 0679 | ډ 0689 | ڙ 0699 | ک 06A9 | ڹ 06B9 | ۉ 06C9 | 06D9 | 𝑀 06E9 | ۹ 06F9 |
| A | ؊ 060A | 061A | ت 062A | غ 063A | ي 064A | 065A | ٪ 066A | ث 067A | د 068A | ښ 069A | ڪ 06AA | ں 06BA | ۊ 06CA | 06DA | 06EA | ۺ 06FA |
| B | ؋ 060B | ؛ 061B | ث 062B | ڱ 063B | ً 064B | 065B | ٫ 066B | پ 067B | ڋ 068B | ڛ 069B | ګ 06AB | ڻ 06BB | ۋ 06CB | 06DB | 06EB | ض 06FB |
| C | ، 060C | ▨ | ج 062C | کپ 063C | ٌ 064C | 065C | ٬ 066C | ت 067C | ڌ 068C | ڜ 069C | ڬ 06AC | ڼ 06BC | ى 06CC | 06DC | 06EC | غ 06FC |
| D | ؍ 060D | ▨ | ح 062D | ڍ 063D | ٍ 064D | 065D | ٭ 066D | ث 067D | ڍ 068D | ڝ 069D | ڭ 06AD | ڽ 06BD | ی 06CD | 06DD | 06ED | ۀ 06FD |
| E | ؎ 060E | ؞ 061E | خ 062E | ڎ 063E | َ 064E | 065E | ؍ 066E | ڀ 067E | ڎ 068E | ض 069E | لا 06AE | ھ 06BE | ی 06CE | 06DE | ۮ 06EE | ۾ 06FE |
| F | ؏ 060F | ؟ 061F | د 062F | ڏ 063F | ُ 064F | ▨ | و 066F | ث 067F | ڏ 068F | گ 069F | گ 06AF | چ 06BF | ف 06CF | 06DF | ۯ 06EF | ھ 06FF |

Confusing Similar Characters
Language-based Keyboards

06A9

0643

06A9

# Confusing Similar Characters
## Summary of Problems/Challenges

- **Security issues** (stability, trust,…) e.g. phishing
  - Some should be addressed at language level first
- **Input devices (keyboards) are based on languages**
- **Not all Arabic-script languages are ready:**
  - Not widely/commonly used
  - Language community are not ready
  - Hard to make decisions on behave of other language communities
  - Pressure to start with ready languages
- **Many problems have been escalated from the protocol to be handled by the registry (e.g. variants, bundling ..etc)**
- **… and yet has to provide a simple and transparent registration services**

isa

- **Assume there are 4 variants to letter (ﻫ)**

| | | | |
|---|---|---|---|
| هدهد | هدهد | ﮨدهد | ەدهد |
| هدهد | هدﮨد | هدەد | هدهد |
| هدﮨد | هدەد | ﮨدهد | ﮨدﮨد |
| ﮨدەد | ەدهد | ەدﮨد | ەدەد |

**16 possible ways to write "هدهد"**

**Only 4 are confusingly similar (25%)**

- **It is expected that some domains will have a large number of variants, e.g.:**
  - There are **16,384** possible variants to write the domain "هيئة-الاتصالات-وتقنية-المعلومات"

Q. How to know if a variant of a domain name has been registered?

Store all variants

**Store only a master key**

Waste of Time and Resources

# Characteristics of the Proposed Solution

- **Work for both ccTLDs and gTLDs**
- **Easy and fast** to be deploy by any registry
- **Extendable** to allow for adding new languages as they become ready
- **Simple** and **Transparent** for end users
  - Do not annoy/confuse end users with technical/special
  - Regular users should be able to register whatever they can type using available keyboards

- **Language Table (LT)**
  - A set of code points (Base characters) to be used by a registry for registering IDN domains in the corresponding language.
  - LT can have Alphabetical, Numbers and Separators (Hyphens, Dots)
- **Variant Table (VT)**
  - A table that records all relations of the LT characters with other characters across the script.
  - Each relation is defined depending on its similarity either:
    - Exact similarity: refers to identically look between base character and another character (e.g. exact match/mirror image).
    - Typo similarity: refers to almost look between base character other character (e.g. typo/style match).
  - Consists of a list of records, each record contains:
    - Base character (from LT),
    - List of other characters (variants) with:
      - A set of positions of similarity [**B**eginning , **M**edial, **F**inal, **I**solated],
      - Relation type (**E**xact, **T**ypo)

# Language-level Required Tables
## Building a Variant Table

List all possible shapes for the basic character

Search for all its variants from the rest of the Arabic script

Then compare the basic character with its variants in all possible positions.

Find all similarity position(s).

Record the similarity (type & position)

| | I | F | M | B |
|---|---|---|---|---|
| ف 0641 FEH | ف | ف | ـفـ | فـ |
| ف 06A7 QAF WITH DOT | ف | ڧ | ـڧـ | ڧـ |

```
          0641; 06A7(FI:T), 06A7(BM:E)
48   0636;
49   0637;
50   0638;
51   0639;
52   063A;
53   0641; 06A7(FI:T), 06A7(BM:E)
54   0642;
55   0643; 06A9(FI:T), 06A9(BM:E), 06AA(BMFI:T)
56   0644;
57   0645;
58   0646; 06BA(BM:E)
59   0647; 06BE(M:E), 06BE(BFI:T), 06C1(I:E), 06C1(MF:T), 06D5(FI:E)
60   0648;
61   0649; 06CD(FI:T), 06D2(FI:T)
62   064A; 067B(BMFI:T), 06D0(BMFI:T)
63   0660; 0030(BMFI:T)
```

- **Language Tables (one for every supported language)**
  - Users can only register domains using base characters from only one language table.

- **Group Variant Table (GVT):**
  - Generated from variant tables..
    - It combines all VTs into one table that group all base characters with all relations across script.
    - Each variant list will be assigned to a unique group key (master key) that identify that group and will be used for generating the Master Key.



$VT_n$

$VT_2$

$VT_1$

+

**GVT**

# Registry-level Required Tables
## Generating Group Variant Table

CVT:

0643; 06A9(FI:T), 06A9(BM:E), 06AA(BMFI:T)
0644;

GVT:

G1B; 0643(Q), 06A9(Q), 06AA(Q), 06A9&0643(E), 06AA&0643(T)
G1M;0643(Q), 06A9(Q), 06AA(Q), 06A9&0643(E), 06AA&0643(T)
G1F; 0643(Q), 06A9(Q), 06AA(Q), 06A9&0643(T), 06AA&0643(T)
G1I; 0643(Q), 06A9(Q), 06AA(Q), 06A9&0643(T), 06AA&0643(T)
G2B; 0644(Q)
G2M;0644(Q)
G2F; 0644(Q)
G2I; 0644(Q)

Language 1

VT | LT

Language 2

VT | LT

.
.
.

Language #

VT | LT

CVT

GVT

MK | Variant List

$$BaseChar(CVT) = \bigcup_{n}^{i=0} BaseChar(VTi)$$

VT: Variant Table
GVT: Group Variant Table

LT : Language Table
MK: Master Key
CVT: Combained Variant Table

<KEY>; [<char_hex> "("<pos>":Q)" ("," | ε)]* [<char_hex1>"&"<char_hex2>"&"  "("<pos>:<rel>")" ("," | ε)]*

**GVT keys**

**Relations for variant characters**

| 177 | G41M; | 063A(Q) | |
| 178 | G41F; | 063A(Q) | |
| 179 | G41I; | 063A(Q) | |
| 180 | G42B; | 0641(Q), 06A7(Q), | 06A7&0641(E) |
| 181 | G42M; | 0641(Q), 06A7(Q), | 06A7&0641(E) |
| 182 | G42F; | 0641(Q), 06A7(Q), | 06A7&0641(T) |
| 183 | G42I; | 0641(Q), 06A7(Q), | 06A7&0641(T) |
| 184 | G43B; | 0642(Q) | |
| 185 | G43M; | 0642(Q) | |

**Keys are used for Querying GVT**

**Check** if the input string follows certain language (using LT). **1**

**Generate** UNICODE code for that input. **2**

**Identify** the **position** for each character depending on language properties (UNICODE Standard). **3**

**Generate** Master key by taking every code from (step 3) and do simple lookup in GVT. **4**

كرة → check if characters are part of language — YES → 0643, 0631, 0629 → 0643(B:Q), 0631(M:Q), 0629(I:Q)

**1**

NO ↓ Fail

**2**

**3**

Generate Master Key using GVT ↓

GVT Key     G15B, G13M, G10I

**4**

.sa

- **Find all Exact strings using Master Key for activation purpose.**

**Master Key**

G15B, G13M, G13I

Use Master key to reversal lookup and find all possible characters (Query relation)

| **Arabic** | **Persian** | **Urdu** |
|---|---|---|
| U+0643, U+0631, U+0629 | U+06A9, U+0631, U+0629 | U+06A9, U+0631, U+06C3 |
| كرة | كرة | كرة |

# Registrant Interface

- **Lookup process (whois)**
  - Check domain syntax under any supported language using LTs.
  - Check if the same domain is available or not.
  - If it is found return the unavailable/whois-information; otherwise continue
  - Get the master- key for the domain (based on GVT)
  - Check if the master- key was registered before or not
  - If master- key is found return unavailable/whois-information; otherwise return domain is available

- **Registration process:**
  - Registrant should select one of supported languages and a domain (U-Label)
  - Registry should accept inputs based on the selected language table
  - If domain name can be registered (available based on Lookup process ) then register the domain

- **Activation process (enable exact variants)**
  - Original Registrant can activate any exact variant from the registered domain's Master Key.
  - List possible Exact variants that can be typed using one of the LT without intermingling between them
  - Activate one/many of Exact variants (if not activated before)

# Registry-level Operation
## Adding New Languages

**1** For every key with Q in the new GVT which also exists in the old GVT, merge the 2 variant lists together.

**2** Add the remaining keys of the new GVT at the end of old GVT keys.

**3** Check the resulted GVT : if the keys with Q appear in different GVT keys or not

**Failed Merge**
Keys appear in different GVT keys

**Successful Merge**
Keys don't appear in different GVT keys

**Cure:**
Regenerate old GVT using existing VTs including the new VTs.
Then regenerate all old Master keys using new GVT!

# Conclusions

- **We tried to have a prototype that fulfill the concepts of script based registry that is:**
  - Optimized, Simple, Transparent, Automated, and addresses many local issues

- **Next steps:**
  - Automating the process of finding variant characters.
  - Variant TLDs should be delegated:
    - E.g. Arabic => كويت => U+0643 U+0648 U+064A U+062A
      Persian => کویت => U+06A9 U+0648 U+06CC U+062A

**Thank you !**

**Thank you**

شكرا

شكرا

U+0634 U+06A9 U+0631 U+0627

U+0634 U+0643 U+0631 U+0627

**G35B G44M G32F G22I**