



# SaudiNIC:

## Securing Arabic Domain Names

Raed Alfayez, SaudiNIC  
APTLD74, Tashkent, Sep 2018

# Agenda

- About SaudiNIC
- About Arabic Script/Language
- Securing Arabic domains
- What's Missing?

# About SaudiNIC

- **Administering** the domain name space under:
  - (.sa) since 1995
  - (.السعودية) Experimental since 2004, official since 2010.
- Operated by a government organization:
  - **CITC** (Communication and Information Technology Commission)
- Leading the local and regional communities efforts towards supporting **Arabic language** in Domain Names since **2001** (more than **15** years of experience)
- Participated in several **TF & WG** related to Arabic domains regionally and globally
- Developed: **Tools, Pilot Projects, Reports & RFC** related to Arabic domains

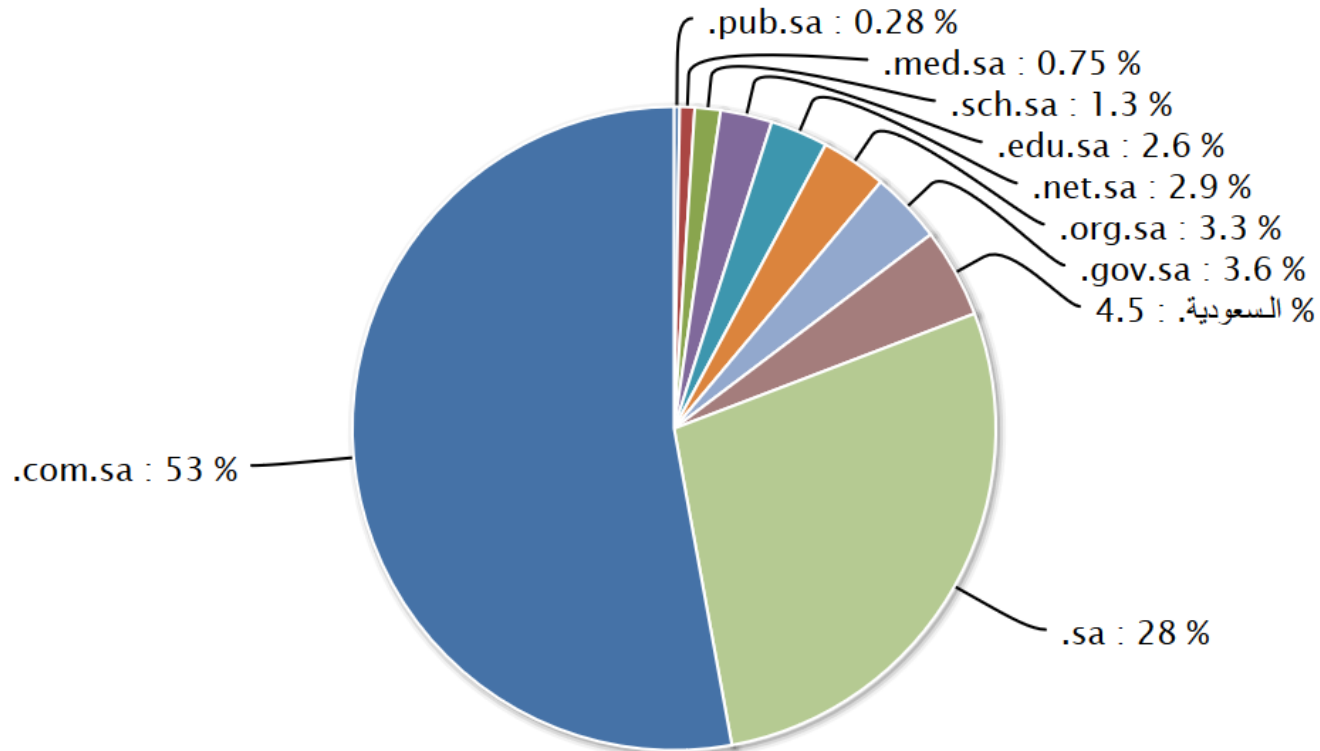
# About SaudiNIC

55,850 Saudi domain names

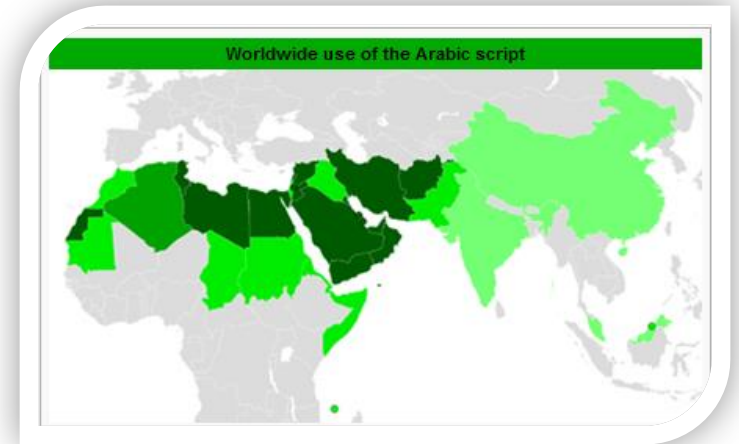
2,540 Arabic domains under (.السعودية)

466 Variant

## 2LD/3LD Domain Names Distribution %



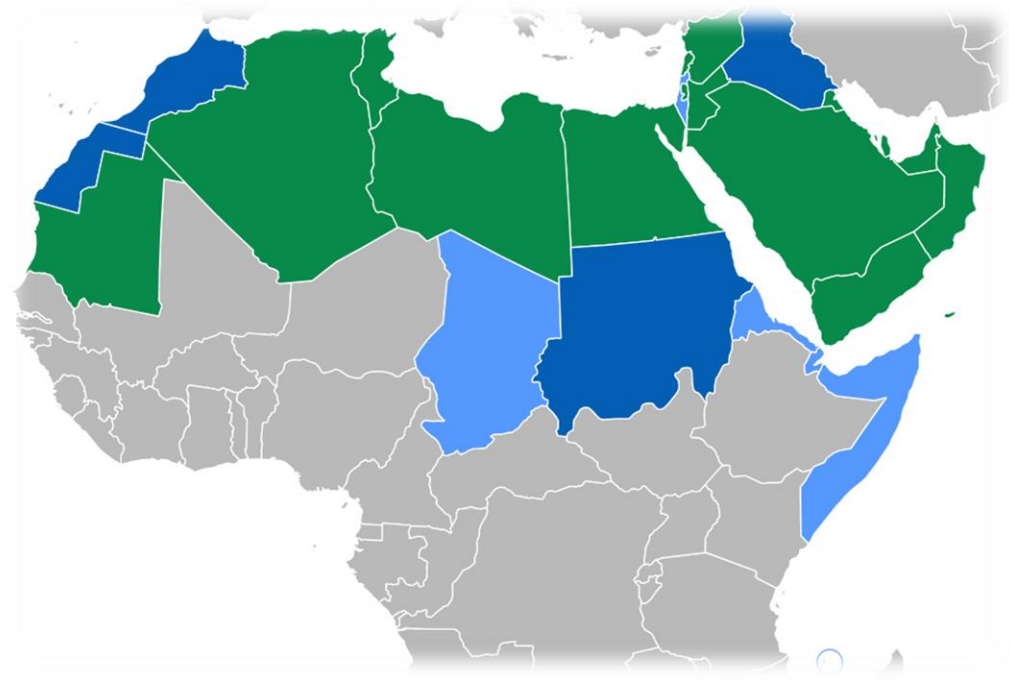
# Introduction: Arabic Script



- The **2<sup>nd</sup>** most widely used **alphabetic** writing system in the world
- Used by **many languages** such as:
  - Arabic, Urdu, Persian, Turkish, Kurdish, Pashto, ...etc
- It is widely used by more than **43 countries**
  - more than **one billion potential users** could be concerned in using Arabic script domain names.

# Introduction: Arabic Language

- Ranked as the **5<sup>th</sup>** language by native speakers in the world.
- Considered as Official/Co-official language in **25** country, ~ **295 million** native speakers.



# Introduction: Arabic Language /Script

- The Arabic **language** is a subset of the Arabic **script**

– **Script:** 0600..06FF, 0750..077F, 08A0..08FF **Language:** 0621..063A, 0641..064A, 0660..0669

	060	061	062	063	064	065	066	067	068	069	06A	06B	06C	06D	06E	06F	075	076	077	08A	08B	08C	08D	08E	08F
0	◌ْ	◌َ	ي	ذ	-	◌ِ	◌ِ	◌ِ	پ	ت	ع	گ	ة	ي	◌ِ	◌ِ	ي	يا	ث	ب	ك				◌ِ
1	◌ِ	◌ِ	ء	ر	ف	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
2	◌ِ	◌ِ	آ	ز	ق	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
3	◌ِ	◌ِ	أ	س	ك	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
4	◌ِ	◌ِ	و	ش	ل	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
5	◌ِ	◌ِ	ص	ض	م	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
6	◌ِ	◌ِ	ي	ض	ن	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
7	◌ِ	◌ِ	ا	ط	ه	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
8	◌ِ	◌ِ	ب	ظ	و	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
9	◌ِ	◌ِ	ة	ع	ي	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
A	◌ِ	◌ِ	ت	غ	ي	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
B	◌ِ	◌ِ	ث	ك	◌ِ	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
C	◌ِ	◌ِ	ج	ك	◌ِ	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
D	◌ِ	◌ِ	ح	ي	◌ِ	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
E	◌ِ	◌ِ	خ	ي	◌ِ	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ
F	◌ِ	◌ِ	د	ي	◌ِ	◌ِ	◌ِ	أ	خ	ز	ف	گ	ه	ي	◌ِ	◌ِ	ب	يا	ث	ب	ك				◌ِ

# Securing Arabic Domains

1

## Identifying allowed code points

- Identifying allowed code points at language level and/or script level
- Not allowing unsecure/unsafe code points that may effect the stability of the registry domain space (e.g. control code points, combining marks, diacritics, Tatweel ..etc).

2

## Identifying variants

- Variants within a language
- Variants across the whole script
- Reachability variants: So that a registered domain name can be accessed regardless of the input devices (language table) being used by the navigator users.

3

## No script mixing

- A label can be constructed only by Arabic code points and digits

4

## No language mixing

- (utilizing the **powerful** tools: Language tables)
- control input via the user interface
  - tremendously reduce the number of unnecessary/unrealistic variants
  - Registry choice to protect the TLD-space.

5

## No digits mixing

- Simplify the need for 3 digits sets in Arabic script
  - European digits
  - Arabic-Indic digits
  - Eastern Arabic-Indic digits

6

## WLE Rules (language/scripts)

- Safeguard to handle some special cases (e.g. similarity in certain positions)



# Securing Arabic Domains

## (1) Identifying code points at both language and script level

### Include only:

- ✓ Allowed code points at language level and/or script level

### Not included:

- Combining Marks
- Diacritics (Tashkeel)
- Special codes ( ZWNJ, ZWJ)
- Unused/Historical

### Non-spacing Marks

◌ْ	جمعية
064F	جمعية
◌َ	جمعية
0650	جمعية
◌ِ	جمعية
0651	جمعية

### 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	+	◌ْ	=	ئْ	is confusing with	ئْ
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	+	◌َ	=	ئِ	is confusing with	ئِ
U+06cc		U+0654		U+06cc U+0654		U+0626
Description: Farsi Yeh + Hamza Above <> Yeh With Hamza Above						
Comments: This is Unicode confusable!						

# Securing Arabic Domains

## (2) Identifying variants (language level)

خ

Hā

ح

Hā

ج

Ġim

ث

Ṭā

ث

Tā

ب

Bā

ا

Alif

ص

Ṣād

ش

Šīn

س

Sīn

ز

Zāī

ر

Rā

ذ

Ḍāl

د

Dāl

ق

Qāf

ف

Fā

غ

Ġain

ع

'Ain

ظ

Zā

ظ

Ṭā

ض

Ḍād

ي

Yā

و

Wāw

ه

Hā

ن

Nūn

م

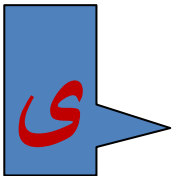
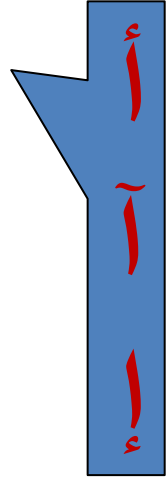
Mīm

ل

Lām

ك

Kāf



# Securing Arabic Domains

## (2) Identifying variants (script level)

- There are a number of **groups** of characters that have the **same shapes** (**Homoglyph**), eg.:

- Kaf group,
- Heh group,
- Yeh group,
- Alef group
- ...

	060	061	062	063	064	065	066	067	068	069	06A	06B	06C	06D	06E	06F
0	ا	ب	ت	ث	ج	ح	خ	د	ذ	ر	ز	س	ش	ص	ض	ط
1	آ	أ	إ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
2	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
3	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
4	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
5	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
6	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
7	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
8	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
9	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
A	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
B	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
C	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
D	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
E	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ
F	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ	أ

# Securing Arabic Domains

## (2) Identifying variants (international reachability)

- For reachability purposes, variants should be addressed to be activated by the registry, so that:
  - A registered domain name is accessed regardless of the input devices (language table) being used by the navigator users.
  - For example:
    - A user registered the domain “مكة” (all characters from the Arabic language)
    - if another user try to reach that domain name from an Internet café in Pakistan he/she will type “مكة” (all characters from the Urdu language)
    - If the “**activated**” variants were not allocated, delegated and hosted then the domain name will not be reachable!

Hence, reachability issue (based on input devices used by other language communities) should be carefully considered when defining variants (by language communities).



LANGUAGE	UNICODE	LABEL
Arabic	(U+0645) (U+0643) (U+0629)	مكة
Persian, Malay, Pashto	(U+0645) (U+06A9) (U+0629)	مكة
Urdu	(U+0645) (U+06A9) (U+06C3)	مكة



ك (0643)

ک (06A9)

# Securing Arabic Domains

## (3) No Script Mixing

### – Recommendation:

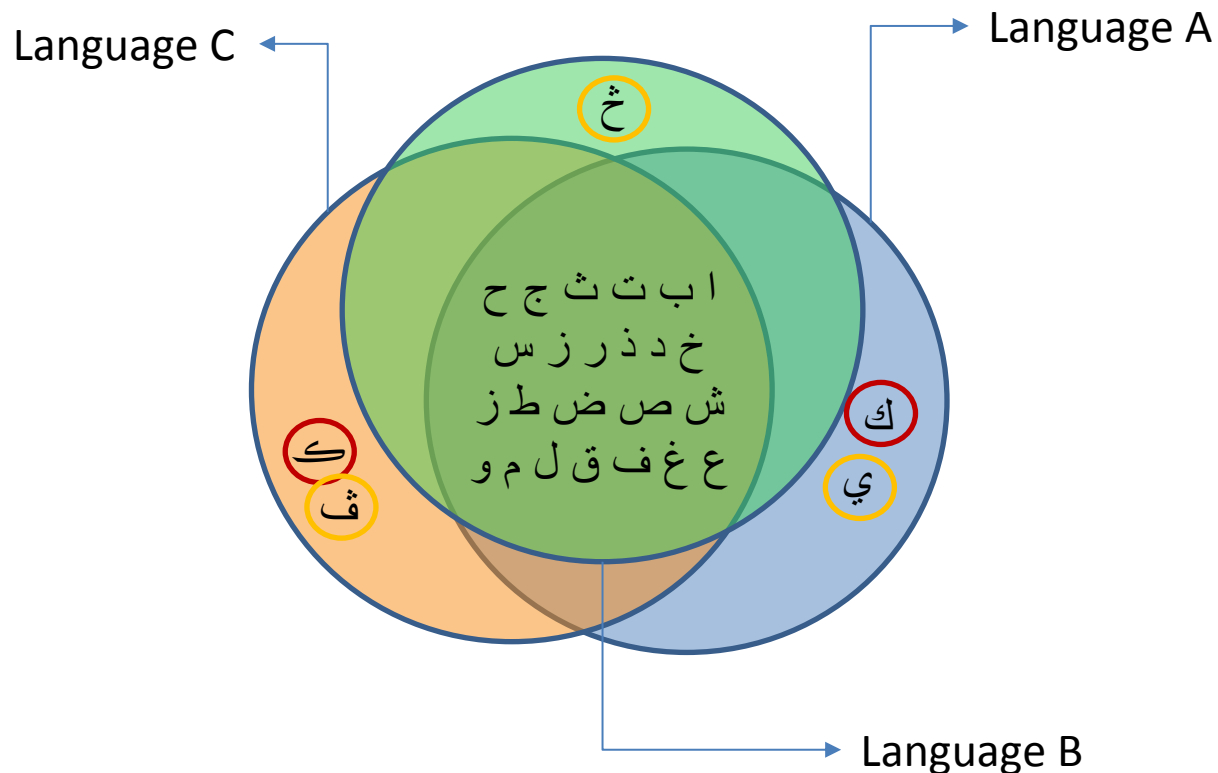
- No mixing within a label or cross labels in the same domain

Format	Example
<FirstName><Arabic-LastName>	Saleh الفلاني
<Arabic-LastName><FirstName>	الفلاني Saleh
<TLD>.<Arabic-Domain>	sa.رسيل
<Domain>.<Arabic-TLD>	raseel.السعودية

# Securing Arabic Domains

## (4) No language Mixing

- Not allowing code point that exists only in one language with another code points that exists only in another language in the same label
  - In the below example: ك & ك cant exist in the same label,
  - same thing for: ي & ي



# Securing Arabic Domains

## (4) No language Mixing

IDN	Total Variants	Allocatable	Blocked	Blocked due to Language Mixing
مكة-المكرمة	3239	34	3205	3181 (99.25%)
القرآن-الكريم	11999	111	11888	11836 (99.56%)
هيئة-الإعلام	47999	81	47918	47764 (99.68%)
كهف-الياسمين	28799	65	28734	28680 (99.81%)
كهف-اكيا	21599	47	21552	21534 (99.92%)

# Securing Arabic Domains

## (5) No Digit Mixing

### – Digits in the Arabic scripts:

- European digits U+0030 .. U+0039 (0123456789)
- Arabic-Indic digits U+0660 .. U0669 (٠١٢٣٤٥٦٧٨٩)
- Eastern Arabic-Indic digits U+06F0 .. U+06F9 (۰۱۲۳۴۵۶۷۸۹)

#	Input string	Display
1	مؤتمر2009 conference2009	<b>Acceptable:</b> Pure European digits
2	مؤتمر٢٠٠٩ conference٢٠٠٩	<b>Acceptable:</b> Pure Arabic-Indic digits
3	مؤتمر٢٠٠٩ conference٢٠٠٩	<b>Acceptable:</b> Pure Eastern Arabic-Indic digits
4	مؤتمر٩٢٠٠٩ conference٢٠٠٩	<b>Not-Acceptable:</b> Mix between European digits & Arabic-Indic digits
5	مؤتمر2٠٠٩ conference2٠٠٩	<b>Not-Acceptable:</b> Mix between European digits & Arabic-Indic digits
6	مؤتمر٢٠٠٩ conference٢٠٠٩	<b>Not-Acceptable:</b> Mix between Arabic-Indic digits & Eastern Arabic-Indic digits
7	مؤتمر2٠٠٩ conference2٠٠٩	<b>Not-Acceptable:</b> Mix between European digits & Eastern Arabic-Indic digits



# Securing Arabic Domains

## (6) WLE Rules

- HEH and TEH MARBUTA
  - They are variants only in the final form (at the end)
  - جده vs جدة
- ALEF MAKSURA:
  - Preventing connected ALEF MAKSURA
  - سند vs سد

# What's Missing?

Register and enable necessary variants:

مكة

مكة

مكة

Registry

Configure DNS & add need RRs (e.g. NS & A & CNAME) for:

xn--ogb5cf

xn--ogb9c4p

xn--hbb4rwc

DNS  
Hosting

Configure Email account and email aliases:

رائد@مكة

رائد@مكة

رائد@مكة

Email  
Hosting

Configure web-server and account and aliases:

```
<VirtualHost 10.10.10.10>  
  DocumentRoot "/makkah"  
  ServerName xn--ogb5cf  
  ServerAlias xn--ogb9c4p  
  ServerAlias xn--hbb4rwc  
</VirtualHost>
```

Web  
Hosting

**Automation for hosting IDN and Variants**

# Gift

- **SaudiNIC's Best Practices in Supporting and Managing Arabic Domain Names**

[http://www.nic.sa/docs/SaudiNIC\\_ADNBP.pdf](http://www.nic.sa/docs/SaudiNIC_ADNBP.pdf)

***Thank you***

**شكرًا**

للمزيد من المعلومات يمكنكم زيارة:

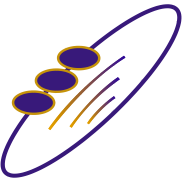
For more information you can visit:



سجل.السعودية

nic.sa

هيئة الاتصالات وتقنية المعلومات  
Communications and Information Technology Commission



هيئة-الاتصالات.السعودية

citc.gov.sa