

# IDNA 2003 & IDNA2008

Patrik Fältström  
Head of Research and Development  
Netnod  
[paf@netnod.se](mailto:paf@netnod.se)

# What is IDN?

- Internationalized domain names (IDN) is a new feature in the Domain Name System
- It gives the ability for users to express domain names using non-latin script
- What is used in the DNS protocol is still latin script
- IDN's are identified by the prefix “xn--”

# Where do we see IDN?

- In a standard (IDNA)
  - What codepoints are allowed, and how?
- In Applications (IDNA)
  - Encode/decode to/from “xn--” form
- In policies that registries set up
  - What domain names can be registered?

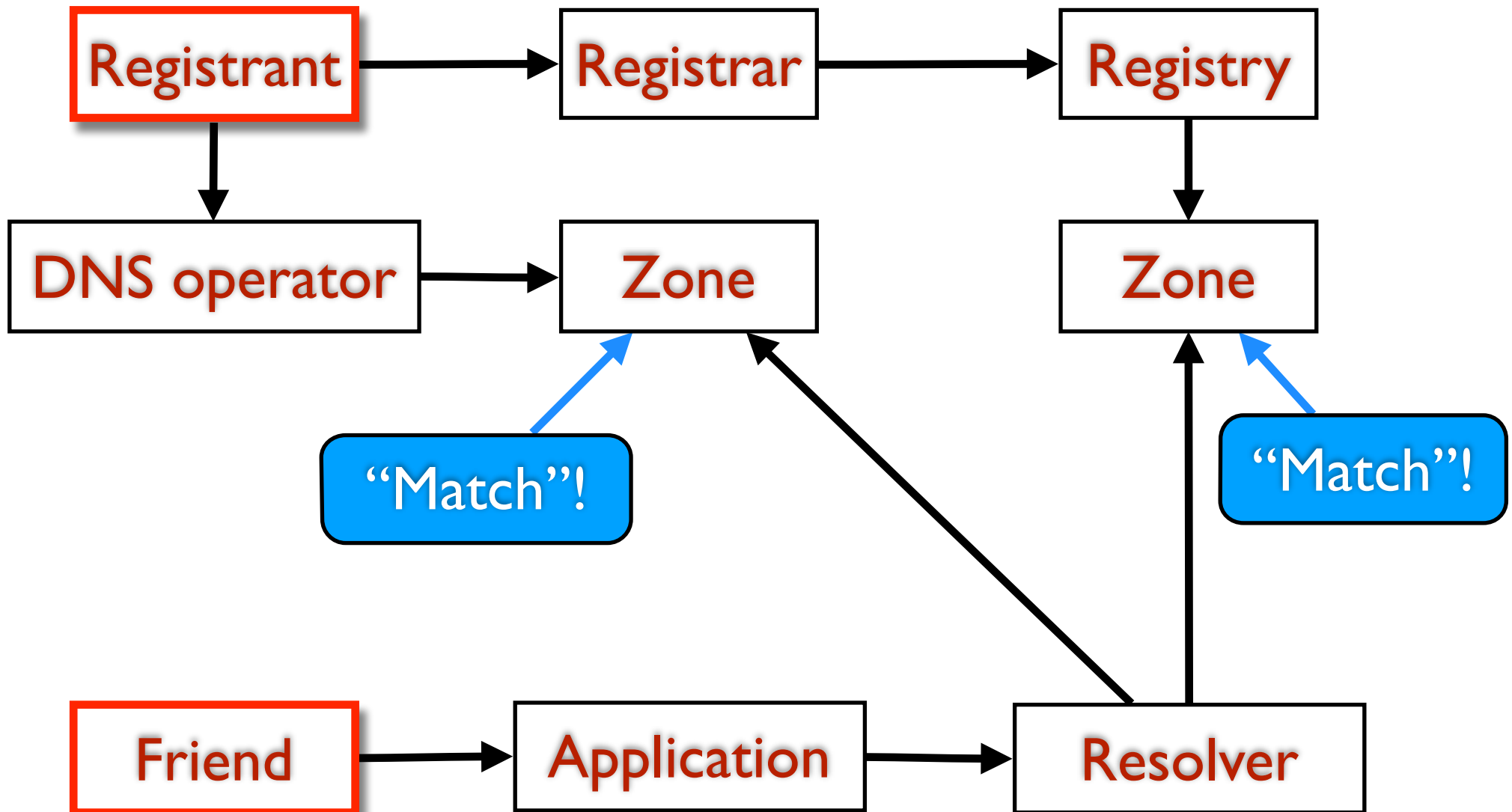
# Why not “just unicode”?

- Unicode include tons of characters, many of them look the same
- Unicode is updated now and then, and we need to agree on what version to use
- Many protocols and applications can only handle latin script

# More precisely

- I register a domain name:
  - fältström.se
- You type in a domain name:
  - FÄLTSTRÖM.SE
- Should you “find” what you want to find?
  - Solution, use: xn--fltstrm-5wa1o.se

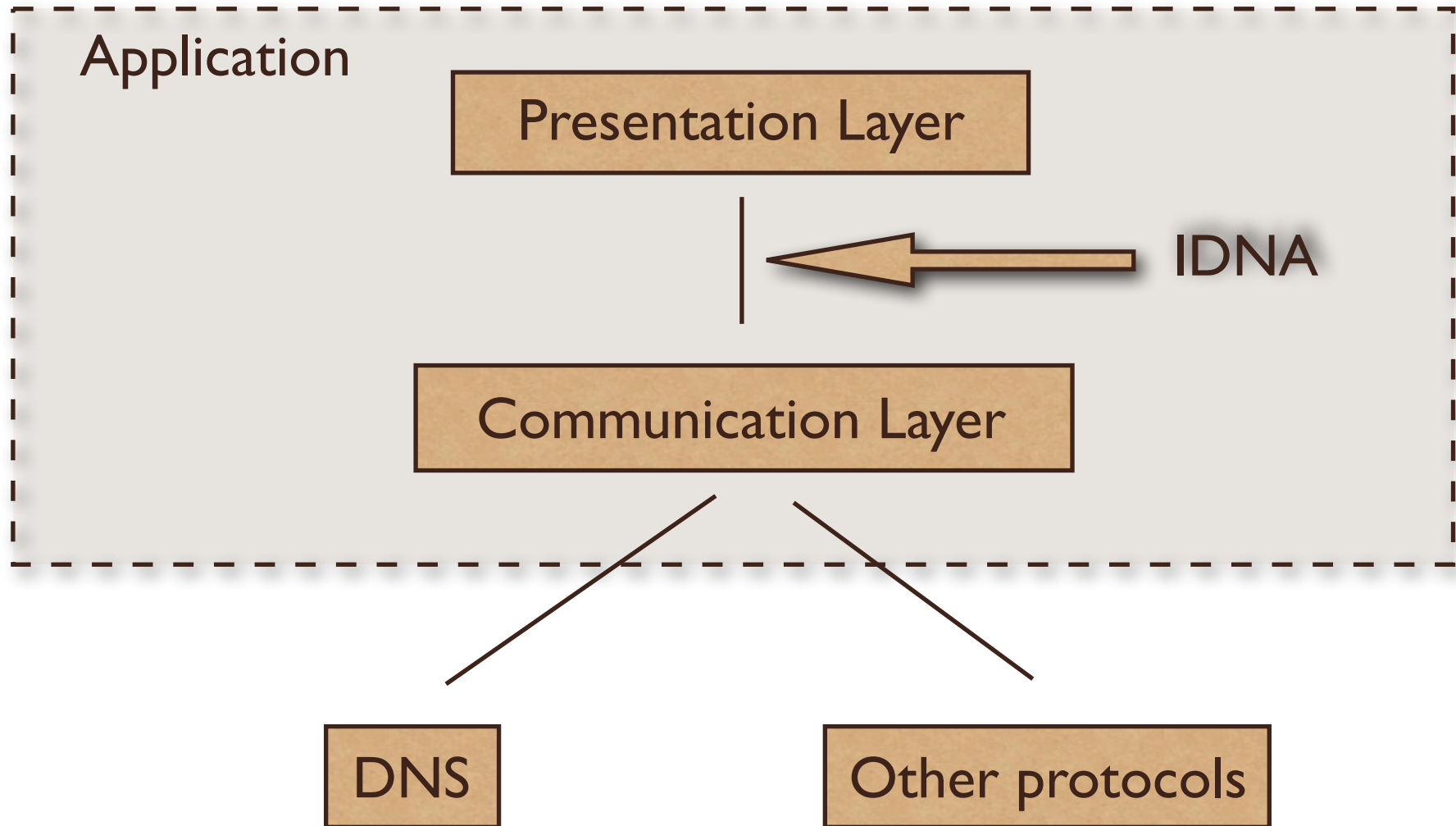
# In even more detail



# Conclusion

- For the implementation of IDN to work, what the registrant register and what other people use, must match when comparing the two domain names in various DNS servers in the world
- It is much easier to define what strings to send to the matching rule than change the matching rule in every DNS server

# Where is this applied?





# Spot The Difference...



≠



```
input[0] = U+0627  
input[1] = U+0654  
input[2] = U+066e  
input[3] = U+06ec  
input[4] = U+0627
```

```
input[0] = U+0623  
input[1] = U+0646  
input[2] = U+0627
```

# They Look The Same To Us ... But Not To A Computer



U+0623



=



+



U+0654

U+0627

# When 1 is not 1...

Arabic-Indic VS. Eastern Arabic-Indic digits

— ٩ ٨ ٧ ٦ ٥ ٤ ٣ ٢ ١ .  
— ٩ ٨ ٧ ٤ ٥ ٣ ٢ ١ .

١ ٢ ٣ ٧ ٨ ٩ .

input[0] = U+06f1  
input[1] = U+06f2  
input[2] = U+06f3  
input[3] = U+06f7  
input[4] = U+06f8  
input[5] = U+06f9  
input[6] = U+06f0

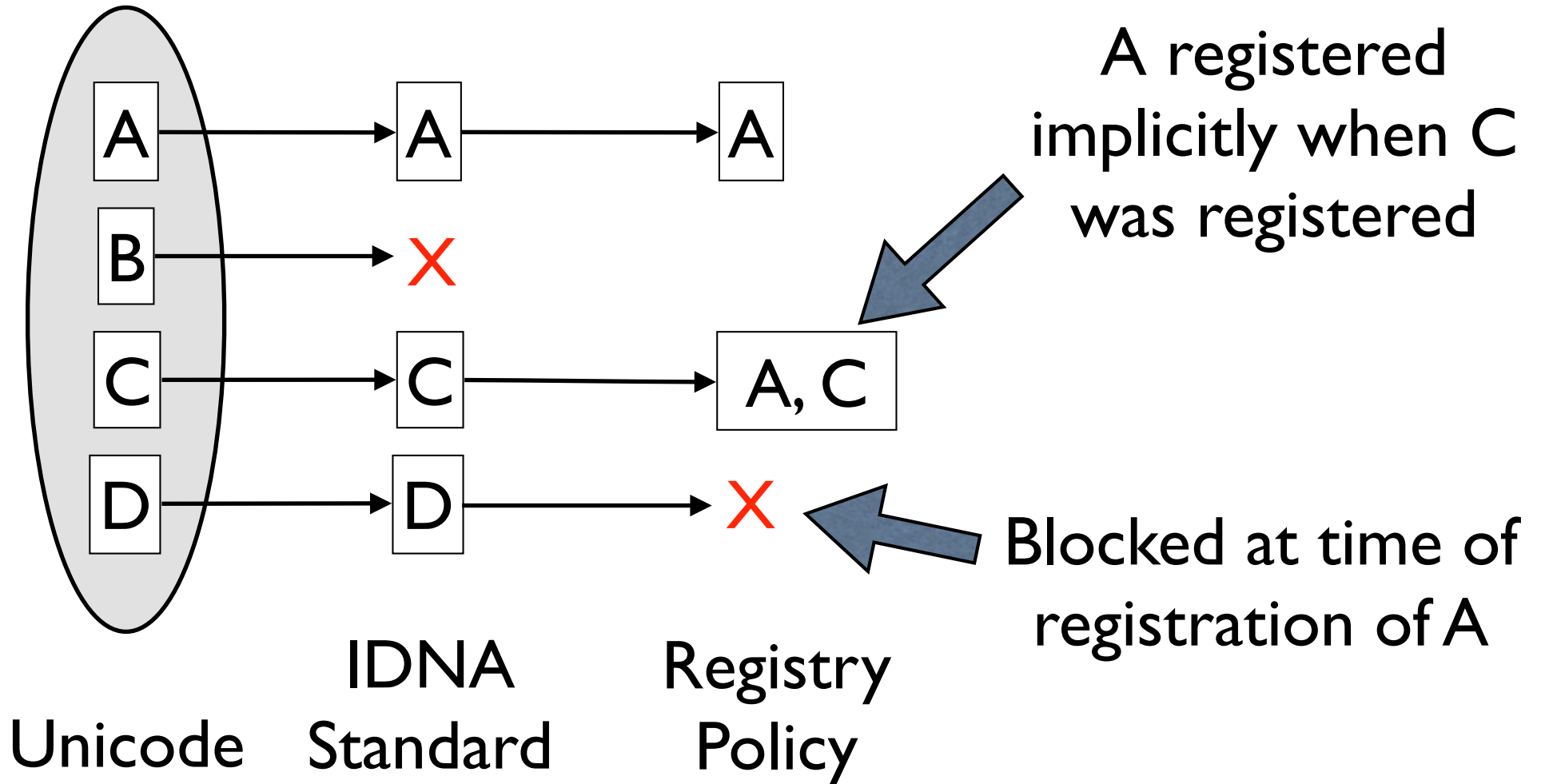
≠

١ ٢ ٣ ٧ ٨ ٩ .

input[0] = U+0661  
input[1] = U+0662  
input[2] = U+0663  
input[3] = U+0667  
input[4] = U+0668  
input[5] = U+0669  
input[6] = U+0660

# Variants, policies etc

“Similar”



# The Arabic Language use only a small part of the Arabic Script table

0600		Arabic														06FF	
	060	061	062	063	064	065	066	067	068	069	06A	06B	06C	06D	06E	06F	
0	0600	0610		0630 ذ	0640 -	0650	0660 ٠	0670 ١	0680 ٢	0690 ٣	06A0 ٤	06B0 ٥	06C0 ٦	06D0 ٧	06E0 ٨	06F0 ٩	
1	0601	0611	0621 ء	0631 ر	0641 ڤا	0651	0661 ١	0671 آ	0681 ٢	0691 ٣	06A1 ٤	06B1 ٥	06C1 ٦	06D1 ٧	06E1 ٨	06F1 ٩	
2	0602	0612	0622 آ	0632 ز	0642 ڤا	0652	0662 ٢	0672 ا	0682 ٢	0692 ٣	06A2 ٤	06B2 ٥	06C2 ٦	06D2 ٧	06E2 ٨	06F2 ٩	
3	0603	0613	0623 أ	0633 س	0643 ك	0653	0663 ٣	0673 ا	0683 ٢	0693 ٣	06A3 ٤	06B3 ٥	06C3 ٦	06D3 ٧	06E3 ٨	06F3 ٩	
4		0614	0624 و	0634 ش	0644 ل	0654	0664 ٤	0674 ا	0684 ٢	0694 ٣	06A4 ٤	06B4 ٥	06C4 ٦	06D4 ٧	06E4 ٨	06F4 ٩	
5		0615	0625 ا	0635 ڤا	0645 م	0655	0665 ٥	0675 ا	0685 ٢	0695 ٣	06A5 ٤	06B5 ٥	06C5 ٦	06D5 ٧	06E5 ٨	06F5 ٩	
6			0626 ١	0636 ٢	0646 ن	0656	0666 ٦	0676 و	0686 ٢	0696 ٣	06A6 ٤	06B6 ٥	06C6 ٦	06D6 ٧	06E6 ٨	06F6 ٩	
7			0627 ا	0637 ط	0647 ه	0657	0667 ٧	0677 و	0687 ٢	0697 ٣	06A7 ٤	06B7 ٥	06C7 ٦	06D7 ٧	06E7 ٨	06F7 ٩	
8			0628 ب	0638 ظ	0648 و	0658	0668 ٨	0678 ١	0688 ٢	0698 ٣	06A8 ٤	06B8 ٥	06C8 ٦	06D8 ٧	06E8 ٨	06F8 ٩	
9			0629 ة	0639 ٢	0649 ١	0659	0669 ٩	0679 ١	0689 ٢	0699 ٣	06A9 ٤	06B9 ٥	06C9 ٦	06D9 ٧	06E9 ٨	06F9 ٩	
A			062A ن	063A ٢	064A ١	065A	066A %	067A ١	068A ٢	069A ٣	06AA ٤	06BA ٥	06CA ٦	06DA ٧	06EA ٨	06FA ٩	
B	060B ڤا	061B ١	062B ١		064B ١	065B	066B ر	067B ١	068B ٢	069B ٣	06AB ٤	06BB ٥	06CB ٦	06DB ٧	06EB ٨	06FB ٩	
C	060C ١		062C ٢		064C ١	065C	066C ١	067C ١	068C ٢	069C ٣	06AC ٤	06BC ٥	06CC ٦	06DC ٧	06EC ٨	06FC ٩	
D	060D ١		062D ٢		064D ١	065D	066D *	067D ١	068D ٢	069D ٣	06AD ٤	06BD ٥	06CD ٦	06DD ٧	06ED ٨	06FD ٩	
E	060E ١	061E ١	062E ٢		064E ١	065E	066E ١	067E ١	068E ٢	069E ٣	06AE ٤	06BE ٥	06CE ٦	06DE ٧	06EE ٨	06FE ٩	
F	060F ١	061F ١	062F ٢		064F ١		066F ١	067F ١	068F ٢	069F ٣	06AF ٤	06BF ٥	06CF ٦	06DF ٧	06EF ٨	06FF ٩	



# Accepted characters for Arabic, Persian, Urdu, and Pashto

	060	061	062	063	064	065	066	067	068	069	06A	06B	06C	06D	06E	06F
0																
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

# Conclusions

- We have to agree on what codepoints can be used in domain names, specifically when multiple variations for expressing “the same” character exist
- Agreements have to be global, across language groups, as scripts are shared
- We need rules both for TLDs themselves and for 2nd level domain registrations

# In the beginning

- 3454 Preparation of Internationalized Strings ("stringprep"). P. Hoffman, M. Blanchet. December 2002. (Format: TXT=138684 bytes) (Status: PROPOSED STANDARD)
- 3490 Internationalizing Domain Names in Applications (IDNA). P. Faltstrom, P. Hoffman, A. Costello. March 2003. (Format: TXT=51943 bytes) (Status: PROPOSED STANDARD)
- 3491 Nameprep: A Stringprep Profile for Internationalized Domain Names (IDN). P. Hoffman, M. Blanchet. March 2003. (Format: TXT=10316 bytes) (Status: PROPOSED STANDARD)
- 3492 Punycode: A Bootstring encoding of Unicode for Internationalized Domain Names in Applications (IDNA). A. Costello. March 2003. (Format: TXT=67439 bytes) (Status: PROPOSED STANDARD)



# What was this - IDNA2003

- 3454 Specifies overall algorithm - stringprep
- 3490 Specifies IDN algorithm - IDNA
- 3491 Specifies Nameprep
- 3492 Specifies Punycode

# stringprep

- With profiles, any Unicode based string can be converted to another Unicode string so that they can be compared
  - Include illegal codepoints
  - Include mapping table
  - Give ability to create profiles
- Used for IDN, LDAP and other protocols

# idna

- Algorithm for how to convert a domain name with Unicode codepoints to ascii
- How to use the stringprep profile and unicode
- Includes specification on how to handle unallocated codepoints
- “core” to IDN standard

# nameprep

- Specific stringprep profile for unicode based domain names
- Convert a domain name with unicode codepoints to one of
  - Illegal domain name
  - Domain name with Unicode codepoints

# punycode

- Converts a label with unicode codepoints to a domain name in ascii
- Example:
  - fältström
  - xn--fltstrm-5wa1o

# What happened?

4690 Review and Recommendations for Internationalized Domain Names  
(IDNs). J. Klensin, P. Faltstrom, C. Karp, IAB. September 2006.  
(Format: TXT=100929 bytes) (Status: INFORMATIONAL)

# In short...

- Explains the problems in the earlier standards
  - Bidirectional scripts
  - Non-spacing codepoints
- Explains the problems with scripts not yet created when IDNA was written
- Explains problem with versioning of Unicode
  - Old standard based on Unicode 3.2

# Example

- If a label include a character that has right to left directionality, both first and last character of the string has to have right to left directionality
- Creates problem if for example the string ends with a codepoint with no directionality



יִירוּאָ

U+05D9 HEBREW LETTER YOD (R)  
U+05D9 HEBREW LETTER YOD (R)  
U+05B4 HEBREW POINT HIRIQ (NSM)  
U+05D5 HEBREW LETTER VAV (R)  
U+05D5 HEBREW LETTER VAV (R)  
U+05D0 HEBREW LETTER ALEF (R)  
U+05B8 HEBREW POINT QAMATS (NSM)

- Note that last codepoint has no directionality (Non Spacing Mark)

# New version, IDNA2008

- Also consists of a few documents
- Will not change punycode
- Backward compatible
- Does explicitly talk about what is possible “to register” in the DNS (old IDNA say what is possible “to use”, and that include mappings)
- **ONLY** define what is possible to register in DNS (U-label / A-label)

# New documents

- 5890 Internationalized Domain Names for Applications (IDNA): Definitions and Document Framework J. Klensin
- 5891 Internationalized Domain Names in Applications (IDNA): Protocol J. Klensin
- 5892 The Unicode Code Points and Internationalized Domain Names for Applications (IDNA) P. Faltstrom
- 5893 Right-to-Left Scripts for Internationalized Domain Names for Applications (IDNA) H. Alvestrand, C. Karp
- 5894 Internationalized Domain Names for Applications (IDNA): Background, Explanation, and Rationale J. Klensin

# RFC 5890

- Addresses the concerns in the IAB document RFC 4690
- Explain how the issues are resolved

# RFC 5891

- Replaces the IDNA specification
- Core specification of new IDN standard

# RFC 5892

- Defines algorithm to use to calculate whether a codepoint in Unicode is in one of the categories
  - PVALID (Protocol Valid)
  - CONTEXTO / CONTEXTJ
  - DISALLOWED
  - UNASSIGNED

# RFC 5893

- Gives specifics for bidirectional scripts

# But IDNA2003 had mappings

- Mappings are not part of IDNA2008
- Labels **MUST** be stable under NFC
- Codepoints in label **MUST** pass bidi requirements
- Codepoints **MUST** be ok according to algorithm specified in tables document (which might include contextual rules)
- We **MIGHT** separate documents on mapping, recommended behaviour for different applications etc



# Implications for registry

- Should not require any change for a registry that do have a policy for what subset of IDNA2003 codepoints one can register
- Changes might be required cases:
  - ➔ Registry that allow any IDNA2003 codepoint
  - ➔ Registry that did allow registration of the few codepoints that indeed have changed, ß for example
  - ➔ Registry that define codepoints based on possible input to IDNA2003 (and not what is actually registered in DNS)

# Why is this needed?

- IDNA standard **must** be independent of Unicode version
- IDNA standard **must** handle bidirectional scripts
- ...plus other things mentioned in RFC 4690

# What has happened?

- Unicode has moved from 5.2 to 7.0
- RFC6452 - The Unicode Code Points and Internationalized Domain Names for Applications (IDNA) - Unicode 6.0

# RFC6452

- Changes from Unicode 5.2 to 6.0
  - U+0CF1 KANNADA SIGN JIHVAMULIYA
  - U+0CF2 KANNADA SIGN UPADHMANIYA
  - U+19DA NEW TAI LUE THAM DIGIT ONE
- Conclusion: No update of IDNA2008 is needed

# Unicode TR#46

- Client software, such as browsers and emailers, faces a difficult transition from the version of international domain names approved in 2003 (IDNA2003), to the revision approved in 2010 (IDNA2008). The specification in this document provides a mechanism that minimizes the impact of this transition for client software, allowing client software to access domains that are valid under either system.
- The specification provides two main features: One is a comprehensive mapping to support current user expectations for casing and other variants of domain names. Such a mapping is allowed by IDNA2008. The second is a compatibility mechanism that supports the existing domain names that were allowed under IDNA2003. This second feature is intended to improve client behavior during the transitional period.

# Unicode TR#46

- Unicode Consortium TR#46 is incompatible with IDNA2008
- Includes good material for people that have implemented IDNA2003
- Transition must be made to IDNA2008

# draft-klensin-idna-5892upd-unicode70-00.txt

- Codepoint in Unicode 7.0:
  - U+08AI :ARABIC LETTER BEH WITH HAMZA ABOVE
- Existing allowed codepoints:
  - U+0628 :ARABIC LETTER BEH
  - U+0654 :ARABIC HAMZA ABOVE
- Should U+08AI be disallowed?

Patrik Fältström  
paf@netnod.se