

Building an IDN Domain Name Registry

Chris Wright

CTO - AusRegistry International

ICANN ^{no.} 35, Sydney, Australia

22nd June 2009

AusRegistry International

- Located in Melbourne, Australia
 - Involved in Domain Name Industry since 1999
 - ICANN Accredited Registrar since 2000
 - .au Registry Operator since 2002
- Domain Name Registry Services
 - Registry Systems and Software Provider
 - Consultancy Services
 - Our software and consultancy services have been used by several other TLDs including some IDN enabled ccTLDs

Overview

- Some of you will be applying for the IDN representation of your ccTLD
- Others want to use IDNs under your ASCII ccTLD
- Most of you, at some point soon, will need to begin supporting IDNs in your Registry Solution
 - Some notes from our experiences
 - Share what we have learnt
 - High-Level overview of part of our solution
 - Some interesting points to think about

So why did we implement IDNs?

- We supply Registry software & services to other TLDs
- We need to remain innovative and up-to-date
- We need to provide what our customers want

Our goals

- Implement IDNs to the current IDNA drafts
- Do so generically and flexibly
- Ensure implementation is easily maintainable
- Ensure implementation may be customised if required by customers
- Configurable to suit various local policies without sacrificing performance, security or stability

Registry Implementation

IDNA – Internationalised Domain Names in Applications

- Whilst it is a protocol in the dictionary definition of the term
- It is NOT a protocol in the sense that DNS, HTTP or EPP are protocols
- It's really three main things:
 - A way of converting a Unicode string into an ASCII string so that it can be used in the DNS protocol
 - A sequence of steps that a Registry must follow before accepting a name for registration
 - A sequence of steps that an Application must follow when looking up a name in the DNS

Why must we understand all of IDNA?

- IDNA assumes any required pre-processing has been performed by Registrars including:
 - ensuring the name is in Unicode NFC form
 - any other local processing that may be required (but is not defined in the IDNA specification)
- To maintain the integrity of the Registry it is important to check that all rules have been followed.

Steps a Registry Must Follow

- Verify that the name is in NFC form, reject if not
- If domain provided in A-label format, generate U-label version using punycode
- If domain provided in U-label form it is strongly advised not to accept it to avoid any ambiguities
- Validate that both A-label form and U-label form are in fact related, reject if not
- Reject any name with leading combining marks
- Reject any name that contains consecutive hyphens in the 3rd and 4th positions (in the U-label)

Steps a Registry Must Follow (cont.)

- Verify that the domain contains only valid code points as defined by the IDNA standards, reject if it doesn't
- Apply the joiner rules (context j rules), reject if these rules fail
- Verify that for each context o code point, a rule exists in the standard and that when the rule is applied the domain name is still valid, reject if any of these rules fail
- If the domain contains any right-to-left characters apply the BIDI rules, reject if any fail

Basic Implementation Summary

- Implementing these steps is relatively simple as they are well defined in the protocol
- A simple implementation of these can be achieved very quickly
- However there are many methods that can be used to efficiently implement these global steps in an elegant manner

Now we have a valid IDN name

Zone specific processing

- What needs to be done, how and why it should be done, is not documented anywhere
- However there are some VERY important steps that should be followed:
 - Checking for duplicate names (including complex equivalencies such as those created by the use of combining marks etc. – think variants or bundles)
 - Apply local policies
 - Validating against our language rules
 - Checking reserved lists

Checking for Duplicate Names

- Duplicate domains are domains that are considered 'the same' as one another
 - ASCII 'the same' is simply a case insensitive compare
 - IDNs this is not the only case
- Lots of reasons for this, including how people have 'modified' the use of their language to map to the limited character set previously available for domain names

Duplicate Example

- ASCII John's Cafe (because of convention)
 - johnscafe.com ← Sacrificing the é
- IDN John's Café (because now we can)
 - Johnscafé.com
- Shouldn't the two be considered the same name? i.e. Duplicates?

Implementing duplicates – The variant generation method

- The idea that one character is a variant of another character e.g.
 - ‘e’ and ‘é’
- When a domain is created using one representation the other representation is also considered registered or ‘blocked’
 - cafe.com
 - café.com
- This is done by ‘calculating’ all of the variants

Implementing duplicates – The variant generation method (cont.)

- This can happen at time of registration in which case all the variants are then stored for later comparisons

or

- This can happen on input to all commands (obviously very inefficient)

Implementing duplicates – The variant generation method

- Calculating and storing duplicates introduces overhead
- Consider a name where there is only one variation of several of the characters in the name e.g.

e → é

cafeeeeeeeeeeeeeeee.com

cafeeeeeeeeeeeeeeeeé.com

cafeeeeeeeeeeeeeeeeée.com

cafeeeeeeeeeeeeeeeeéé.com

.

.

caféééééééééééééééééé.com

In this fictitious case there is 2^{16} combinations i.e. 65,536 variations

Implementing duplicates – The variant generation method –

- If we have a domain name with just 32 characters in it, each with one variant we would have over 4 billion variants
- There has to be a better way!
- And there is...

Implementing duplicates – The canonical method

- Canonical representation of domain names isn't new
- ASCII domain names use the concept, its built into the protocol
- The overall premise is that we assign each character a canonical form

What do we mean by character?

- A character, for the sake of this discussion, is a sequence of one or more code points that represents one particular component of a word.

‘a’ is a character

ا (single code point) is a character

اَ (multiple code points) is a character

Assigning canonical form

- Each character is assigned a canonical form
- You can think of it as the base form of the character
- In most cases it just be the character itself
- Sometimes just the lowercase version, sometimes another code point entirely
- The actual character chosen doesn't really matter – its just a concept

Using the canonical form

- Define all canonical mappings for your zone
- Perform a simple substitution of each character for its canonical equivalent
 - This generates the canonical form of the label being registered
- Use this canonical form of the label as the unique key for the domain registration representing ALL forms of the domain name (without each of those forms having to be generated and/or stored)

Using the canonical form

- In our zone we allow the following characters with the canonical mappings listed:

a → a

c → c

e → e

é → e

f → f

Using the canonical form – An example

- We register the name cafe'.com and compute the canonical form
 café.com → cafe.com
- The domain is café.com but the unique label is cafe. So when someone tries to register cafe.com we compute the canonical form
 cafe.com -> cafe.com
- But this will NOT be allowed as a domain with that canonical label is already registered

Using the canonical form – Another example

- The name cafeeeeeeeeeeeeeeeee.com maps to cafeeeeeeeeeeeeeeeee.com → cafeeeeeeeeeeeeeeeee.com
 as does
 cafééééééééééééééééé.com → cafeeeeeeeeeeeeeeeee.com
 as does
 caféééééééééééééééééééééééé.com → cafeeeeeeeeeeeeeeeee.com
- So by storing the canonical form and checking all new registration attempts against it we have blocked all other registrations without actually having to calculate them all!

More on canonical...

- Mapping names to a canonical form is nothing new
 - Exactly what happens in existing domain name registries when we lower case names
 - Implied canonical mapping between upper case and lower case (implemented by a function)
 - Just also happens to be enforced by the DNS protocol itself

Making canonical work for us

- Just as we lower case the domain name provided to Registry functions such as:
 - Search
 - Domain Check / Update
 - Reserved List Matching
 - WHOIS
 - Etc.
- If we apply the canonical mapping to IDN names passed to registry functions everything just works

Benefits of using canonical

- It just works
- Its linear time regardless of the size of the domain names and desired variant configuration
- It provides speed and efficiency benefits, especially when compared to variant generation methods
- It saves space and memory
- Its a simple algorithm that is easy to implement, less error prone and easier to optimise

Some examples of how canonical 'just Works'

- Just as we lower case the domain name provided to Registry functions such as:
 - Search
 - Domain Check / Update
 - Reserved List Matching
- If we apply the canonical mapping to IDN names passed to Registry functions everything just works
- It is important to note that the canonical mappings have to be consistent for all languages within a zone

Bundles

Why Bundles?

- Sometimes blocking is just not enough
- In some scenarios it make sense that a Registrant can make use of multiple versions of a name e.g.
 - cafe.com
 - café.com

In simple terms...

- Its the same as the generating variant model, so it has the same issues
 - If in our zone configuration we said that we wanted the following character variant ‘provisioned’ or used to create ‘bundles’

١ → 1

- And then we registered the name

١١١١١١١١.com

- We still end up with...

Example

- The following variants to be provisioned

11111111.com

111111111.com

1111111111.com

11111111111.com

.

.

.

1111111111.com

- Which in this case would be 256 variants to be calculated, stored and provisioned in the zone file
- Canonical mappings can't help us here

Bundles (cont.)

- Character variants for blocking of registrations make all combinations important
- ... But when considering bundling.. If we look at the reason people desire variants, another option is presented

Continuing our example...

- In this case it makes sense that someone may enter either of the following domains:

11111111.com
11111111.com

- But does it really make sense that someone would type the following domains names:

11111111.com
11111111.com

- All combinations need to be blocked (which canonical mappings will do) , yet only two out of the 256 variants provisioned in the DNS are required.

Introducing Mutual Exclusion

Mutual Exclusion

- Mutual exclusion is not a new concept, it is used everywhere in modern-day life
- If we apply it to domain name variants we can achieve the desired behaviour e.g.

Primary Grouping	Sub-Grouping
Numerals	English Numerals e.g. 1,2,3,4,5...
	Arabic Numerals e.g. ١٢٣٤٥...

So the rule is...

- If a domain name contains any characters that are in one sub-group, it is not allowed to contain any characters from other sub-groups of the same primary group to be provisioned in the DNS
- i.e. The characters in one sub-group are mutually exclusive to the characters in another subgroup

Returning to our example...

Primary Grouping	Sub-Grouping
Numerals	English Numerals e.g. 1,2,3,4,5...
	Arabic Numerals e.g. ١٢٣٤٥...

- These are allowed:

١١١١١١١.com

11111111.com

- But these are not:

١١١١١١١١.com

١١١١1111.com

Validating Local Language Rules

What are local language rules?

- In short, they are and can be anything
 - Which unicode code points make up the language
 - Handling of edge cases
 - ae → AE
 - ss → SS
 - Final form sigma
 - and so on
- Important that the business rule engine is flexible and customisable enough to handle these requirements

Putting it all together

How can we represent IDN configuration?

- In a generic way
- That reduces the management and configuration overhead
- That is easily understood by non-technical people

Our Solution

Language Set

- Name
- Description

Canonical Mappings

- List of ALL code points from ALL languages in the Language set with their canonical equivalents

Language

- Name
- Tag
- Description

Allowed Code Points

- List of the code points allowed in that language

Mutual Exclusion Groups

- Configuration of exclusion groups for the language

Just the tip of the iceberg!

Many other areas to consider

- IDNA – Internationalised Domain Names in Applications
 - Registry Systems are also Applications in fact they are a collection of many different applications
 - We have to implement the application and the Registry portion of IDNA
- Changing Language rules in an already established zone!
- Effects on EPP
 - Command Extensions
 - Protocol Extensions
 - Returning Variants
- Security Considerations
 - Puny Code Overload
 - Puny Code Reverse Engineering
 - Handling of supplementary characters

Many more areas to consider

- Internationalising your Registry
- Unicode versions understood by software in use
- Registrars, Registrants
- Effects on DNS
 - Increase in zonefile size
 - DNAME vs NS records
 - Increase in complexities
- Infrastructure Requirements
- IDNs, variants & DNSSEC

IDNs are hard (to do right)...

- There are many creative and innovative solutions to all of the issues I have mentioned
- Start experimenting, share knowledge and seek help
- I will be talking a little more about some of the other issues at a higher level on Wednesday (?)
- Please feel free to speak to myself or another AusRegistry International staff member. We'd love to chat

