



Developing LGR for the Root Zone

Sarmad Hussain | Director IDN Programs, ICANN | 28 May 2017

Agenda

1

Background
of TLD
Program

2

Repertoire

3

Variants

4

WLE Rules

5

XML
Formulation
and LGR Tool

6

Conclusions
and Next
Steps

A world map where the continents are defined by a complex network of white nodes and connecting lines, set against a solid teal background. The nodes vary in size and are densely packed in some areas, creating a digital or network-like appearance of the globe.

Background of TLD Program

ASCII Domain Name Label

www.cafe123.com



Third Level
Domain

Second Level
Domain

Top Level
Domain (TLD)

Top Level Domains (TLDs)

- ⦿ Country Code TLDs (ccTLDs)
 - ⦿ .sg, .cn, .kh, .la, .mm, .th, .ca, ...
 - ⦿ Two letter [a..z] codes, reserved for countries and territories by ISO 3166 standard
- ⦿ Generic TLDs (gTLDs)
 - ⦿ .com, .org, .net, .edu, ... - organizations
 - ⦿ New gTLDs – 1930 applications in 2012

Domain Stakeholders

- ⦿ ICANN
- ⦿ Registry
- ⦿ Registrar
- ⦿ Reseller
- ⦿ Registrant
- ⦿ End-User

ASCII Domain Name Label

www.cafe123.com



Third Level
Domain

Second Level
Domain

Top Level
Domain (TLD)



Forming ASCII Labels

Use LDH

- Letters [a-z]
- Digits [0-9]
- Hyphen (LDH)

Label length = 63

Other constraints (e.g. on hyphen)

Forming ASCII Labels

Use only Letters

- Letters [a-z]

Label length = 63

Internationalized Domain Name (IDN) Labels



Syntax of IDN Labels

Valid U-Label: Unicode code points as constrained by IDNA 2008

Valid A-Label - “xn--” followed by punycode of U-Label of length 59

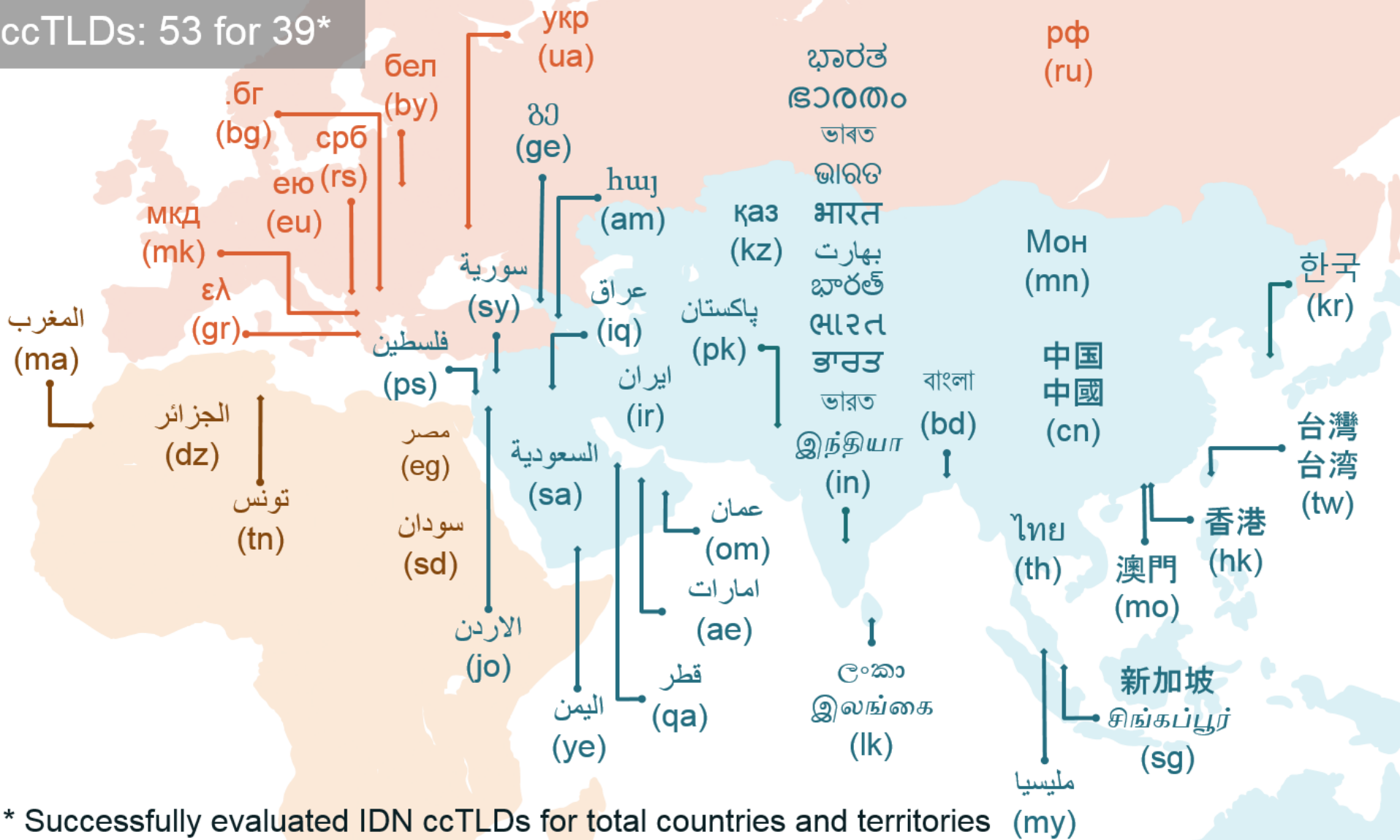
Syntax of IDN Labels

Valid U-Label, further constrained by the “letter” principle for TLDs

Valid A-Label

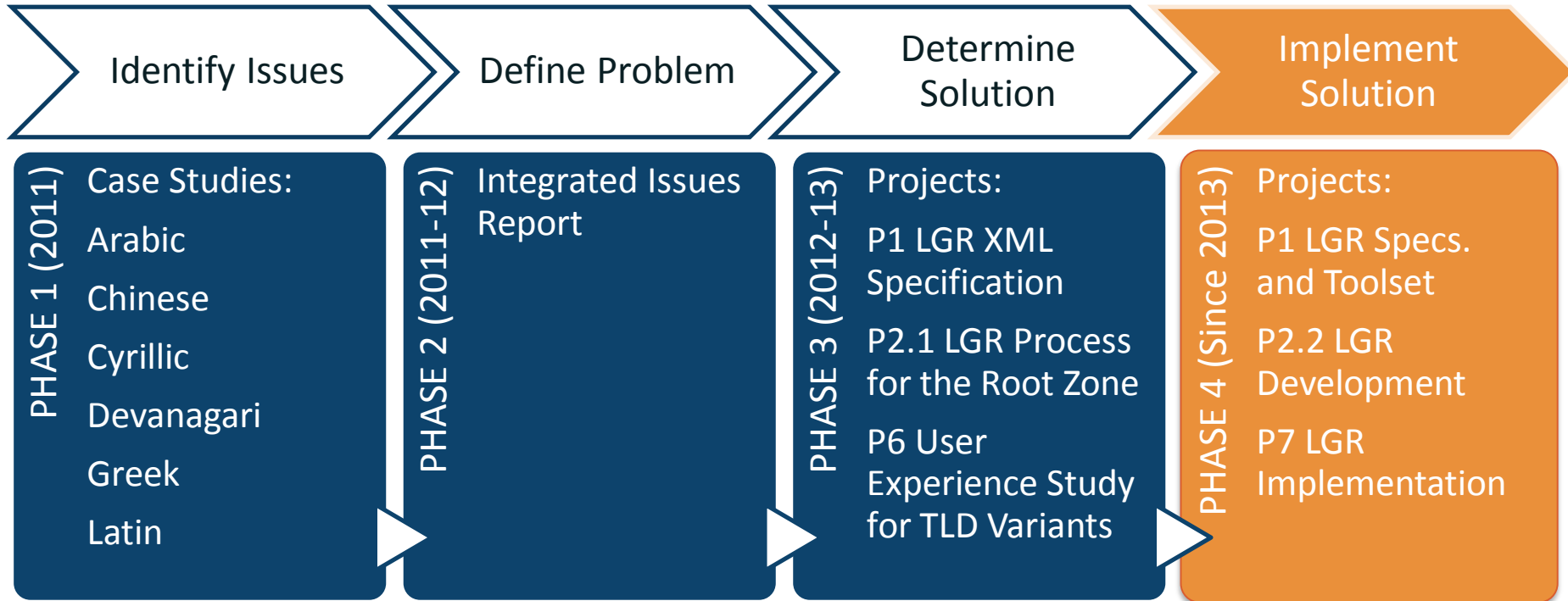
IDN Country Code Top-Level Domains

ccTLDs: 53 for 39*



* Successfully evaluated IDN ccTLDs for total countries and territories

IDN TLD Program



Community agreed to define a Label Generation Rules (LGR)

Reports and documentation of all completed projects available at:
<https://www.icann.org/resources/pages/reports-2013-04-03-en>

Label Generation Rules for the Root Zone

- ⦿ For the Root Zone, single “table” containing data for all scripts
 - ⦿ As it is a shared resource, must be conservative
 - ⦿ Must be stable and secure
 - ⦿ Must be based on inclusion based analysis
- ⦿ For each script or writing system:
 - ⦿ Which code points are valid for use?
 - ⦿ Are any of these code points variants of each other?
 - ⦿ Are there any additional constraints on the labels?

LGR for the Root Zone

Unicode

	000	001	002	003	004	005	006	007
0	NUL	DLE	SP	0	@	P	`	p
1	SOH	DC1	!	1	A	Q	a	q
2	STX	DC2	"	2	B	R	b	r
3	ETX	DC3	#	3	C	S	c	s
4	EOT	DC4	\$	4	D	T	d	t
5	ENQ	NAK	%	5	E	U	e	u
6	ACK	SYN	&	6	F	V	f	v
7	BEL	ETB	'	7	G	W	g	w
8	BS	CAN	(8	H	X	h	x
9	HT	EM)	9	I	Y	i	y
A	LF	SUB	*	:	J	Z	j	z
B	VT	ESC	+	;	K	[k	{
C	FF	FS	,	<	L	\	l	
D	CR	GS	-	=	M]	m	}
E	SO	RS	.	>	N	^	n	~
F	SI	US	/	?	O	_	o	DEL

• • •

	120	121	122	123	124	125	126	127	128	129	12A	12B
0	ሀ	ሐ	ሠ	ሰ	ቀ	ቐ	በ	ተ	ጎ	ነ	አ	ኮ
1	ሁ	ሑ	ሡ	ሱ	ቁ	ቑ	ቡ	ቱ	ጊ	ጊ	ኡ	
2	ሂ	ሐ	ሢ	ሲ	ቂ	ቊ	ቢ	ቲ	ጋ	ጋ	ኢ	ኩ
3	ሃ	ሐ	ሣ	ሳ	ቃ	ቍ	ባ	ታ	ጌ	ጌ	ኣ	ኣ
4	ሄ	ሐ	ሤ	ሴ	ቄ	቎	ቤ	ቴ	ግ	ግ	ኤ	ኮ
5	ህ	ሐ	ሥ	ስ	ቅ	቏	ቦ	ት	ገ	ገ	እ	ሱ
6	ሆ	ሐ	ሦ	ሶ	ቆ	ቐ	ቦ	ቶ	ገ	ገ	አ	
7	ሀ	ሐ	ሧ	ሰ	ቇ	ቑ	ቦ	ቶ	ገ	ገ	አ	
8	ለ	መ	ረ	ሸ	ቈ	ቐ	ሸ	ቐ	ገ	ገ	ከ	ከ
9	ሉ	መ	ሩ	ሹ			ሸ	ቐ	ገ	ገ	ከ	ከ
A	ሊ	ሚ	ሪ	ሺ	ቀ	ቀ	ሸ	ቐ	ገ	ገ	ከ	ከ
B	ላ	ማ	ሪ	ሻ	ቁ	ቁ	ሸ	ቐ	ገ	ገ	ከ	ከ
C	ሌ	ሚ	ሪ	ሽ	ቁ	ቁ	ሸ	ቐ	ገ	ገ	ከ	ከ
D	ል	ም	ር	ሽ	ቀ	ቀ	ሸ	ቐ	ገ	ገ	ከ	ከ
E	ሎ	ም	ሮ	ሾ			ሸ	ቐ	ገ	ገ	ከ	ከ
F	ሏ	ሚ	ሪ	ሽ			ሸ	ቐ	ገ	ገ	ከ	ከ

• • •

LGR for the Root Zone

Unicode

IDNA2008 – by IETF

LGR for the Root Zone

Unicode

IDNA2008

Maximal Starting Repertoire – by Integration Panel of ICANN

LGR for the Root Zone

LGR Proposal – by **Generation Panel** of Script Community

Unicode

IDNA2008

Maximal Starting Repertoire (MSR)

X		X		X
X				
X		X		
X				
X			X	X
X				X
X				X

	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	ا	ب	ي	ذ	-	ر	و	پ	ت	ث	غ	گ	ة	ي	و	و	ب	ب	ش	ب	ک				و
1	ا	ب	ء	ر	ف	ا	ا	ا	خ	ز	ف	گ	ه	ي	و	ا	ث	پ	ز	ب	ا				و
2	م	و	آ	ز	ق	ا	ا	ا	خ	ز	ب	گ	ا	ا	و	ا	پ	ک	خ	ج	ز				و
3	ب	و	ا	س	ك	ا	ا	ا	ج	ر	ف	گ	ا	ا	و	ا	ت	ک	ا	ظ					و
4	ا	ب	و	ش	ل	ا	ا	ا	ج	ر	ف	گ	و	-	و	ا	ن	ک	ا	ق				و	و
5	ا	ب	ا	ص	م	و	ا	ا	خ	ر	پ	ل	و	ه	ر	ا	ب	م	ئ	ق				و	و
6	ا	ب	ا	ض	ن	ا	ا	ا	ج	ر	ق	ل	و	ا	ا	ا	ن	م	ئ	ث				و	و
7	ا	ب	ا	ط	ه	ا	ا	ا	ج	ر	ف	ل	و	ا	ا	ا	خ	ن	ي	م				و	و
8	و	و	ب	ظ	و	ا	ا	ا	ث	ر	ق	ل	و	و	و	ا	ج	ن	و	ن				و	و
9	ا	و	ا	ع	ي	ا	ا	ا	د	ر	ک	ن	و	ا	ا	ا	ڈ	ن	و	ن				و	و
A	ا	و	ا	ت	ي	ا	ا	ا	د	ب	ب	ک	و	ا	و	ا	د	ل	ا	ا				و	و
B	ا	ب	ا	ث	ا	ا	ا	ا	ب	ا	ک	ن	و	ا	و	ا	ر	ز	ا	و				و	و
C	ا	ALM	ج	ک	ا	ا	ا	ا	ت	ا	پ	ن	ي	و	و	ا	ش	ز	ج	ب				و	و
D	ا		ح	ئ	ا	ا	ا	ا	د	ي	ص	ك	ئ	ي	ا	ا	ش	ش	ش	ا				و	و
E	ا	ا	خ	ئ	ا	ا	ا	ا	پ	ا	ض	ك	ه	ئ	ا	ا	ع	ج	ش	ا				و	و
F	ع	ا	د	ئ	ا	ا	ا	ا	ت	ا	ظ	ك	خ	و	ا	ا	غ	ج	ك	ص				و	و

Label Generation Rules (LGR)

0632	ك	0652	3	0672	ح	ر	فبا	0682	0692	06A2	06B2
0643	ك	0653	3	0673	ح	ر	فبا	0683	0693	06A3	06B3
0644	ل	0654	4	0674	ح	ر	فبا	0684	0694	06A4	06B4
0645	م	0655	5	0675	ح	ر	فبا	0685	0695	06A5	06B5
0646	ن	0656	6	0676	ح	ر	فبا	0686	0696	06A6	06B6
0647	ه	0657	7	0677	ح	ر	فبا	0687	0697	06A7	06B7
0648	و	0658	8	0678	ح	ر	فبا	0688	0698	06A8	06B8
0649	ي	0659	9	0679	ح	ر	فبا	0689	0699	06A9	06B9
064A	ي	065A	%	067A	ح	ر	فبا	068A	069A	06AA	06BA

- Valid code points
- Variants code points

کابل

کابل

- Label constraints
 - Cannot mix ک and ك in a label

کککته ✓

کککته ✓

کککته x

کککته x

MSR and LGR

	12C	12D	12E	12F	130	131	132	133	134	135	136	137
0	ሰ	ዐ	ዠ	ደ	ጀ	ጐ	ጠ	ጰ	ፀ	ፐ	✳	፳
	12C0	12D0	12E0	12F0	1300	1310	1320	1330	1340	1350	1360	1370
1		ዐ	ዠ	ደ	ጀ		ጠ	ጰ	ፀ	ፐ	፳	፳
		12D1	12E1	12F1	1301		1321	1331	1341	1351	1361	1371
2	ሰ	ዐ	ዠ	ደ	ጀ	ጐ	ጠ	ጰ	ፀ	ፐ	፳	፳
	12C2	12D2	12E2	12F2	1302	1312	1322	1332	1342	1352	1362	1372
3	ሰ	ዐ	ዠ	ደ	ጀ	ጐ	ጠ	ጰ	ፀ	ፐ	፳	፳
	12C3	12D3	12E3	12F3	1303	1313	1323	1333	1343	1353	1363	1373
4	ሰ	ዐ	ዠ	ደ	ጀ	ጐ	ጠ	ጰ	ፀ	ፐ	፳	፳
	12C4	12D4	12E4	12F4	1304	1314	1324	1334	1344	1354	1364	1374
5	ሰ	ዐ	ዠ	ደ	ጀ	ጐ	ጠ	ጰ	ፀ	ፐ	፳	፳
	12C5	12D5	12E5	12F5	1305	1315	1325	1335	1345	1355	1365	1375
6		ዐ	ዠ	ደ	ጀ		ጠ	ጰ	ፀ	ፐ	፳	፳
		12D6	12E6	12F6	1306		1326	1336	1346	1356	1366	1376
7			ዠ	ደ	ጀ		ጠ	ጰ	ፀ	ፐ	፳	፳
			12E7	12F7	1307		1327	1337	1347	1357	1367	1377
8	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12C8	12D8	12E8	12F8	1308	1318	1328	1338	1348	1358	1368	1378
9	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12C9	12D9	12E9	12F9	1309	1319	1329	1339	1349	1359	1369	1379
A	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CA	12DA	12EA	12FA	130A	131A	132A	133A	134A	135A	136A	137A
B	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CB	12DB	12EB	12FB	130B	131B	132B	133B	134B		136B	137B
C	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CC	12DC	12EC	12FC	130C	131C	132C	133C	134C		136C	137C
D	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CD	12DD	12ED	12FD	130D	131D	132D	133D	134D	135D	136D	137D
E	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CE	12DE	12EE	12FE	130E	131E	132E	133E	134E	135E	136E	137E
F	ዐ	ዐ	ዐ	ዐ	ጐ	ጐ	ጐ	ጐ	ፀ	ፐ	፳	፳
	12CF	12DF	12EF	12FF	130F	131F	132F	133F	134F	135F	136F	137F

Which **code points** must be included in the Root Zone

- ⊙ Are exclusions from MSR (pink) correct?
- ⊙ What must be included in LGR?

⊙ *“everyday, general purpose [use ...] in a stable and widespread manner”*

Are there any **variant code points**

- ⊙ Two code points when replaced produce labels considered confusingly similar by an end-user

Are there any **label-level constraints**

- ⊙ Well-formedness of a cluster?
- ⊙ Constraints on initial or final position in a label?
- ⊙ Other?

Root Zone LGR Procedure

Generation Panels

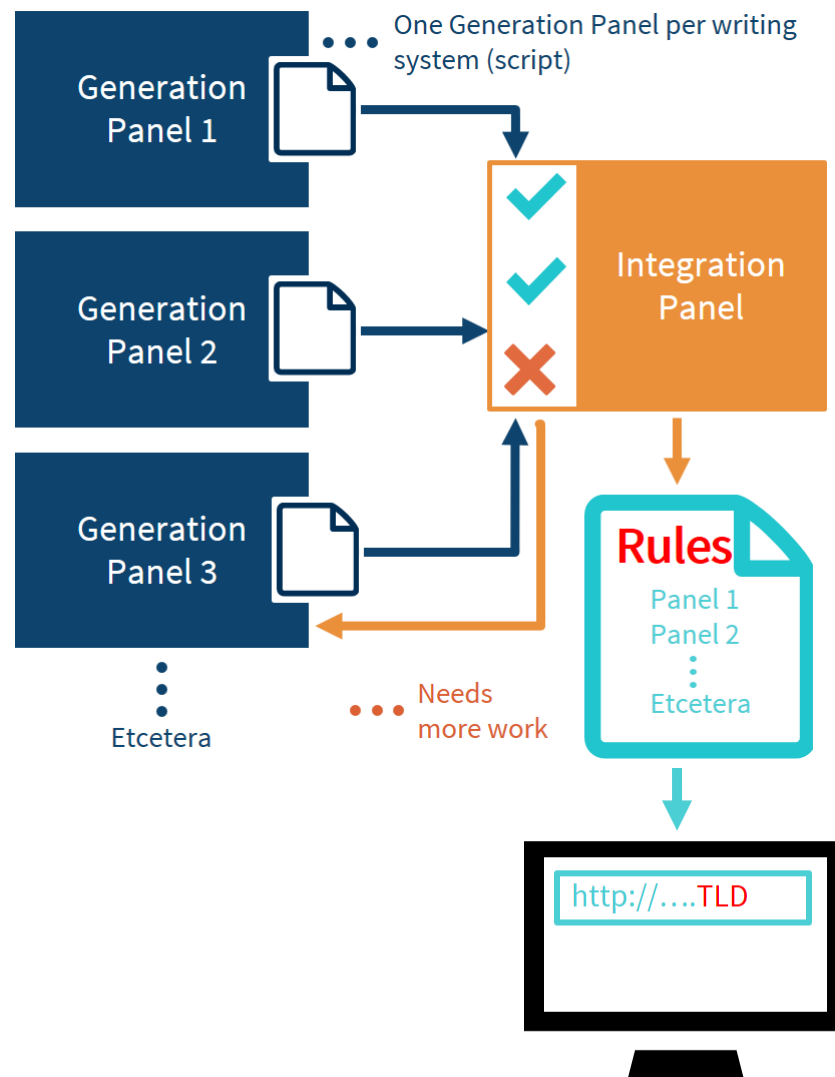
- Generate proposals for script specific LGRs, based on community expertise and requirements

Integration Panel

- Integrates them into common Root Zone LGR while minimizing the risk to Root Zone as shared resource

Label Generation Rules (LGR)

- Which labels are permissible
- Which variant labels exist
- Which variant labels may be allocated



LGR Specification

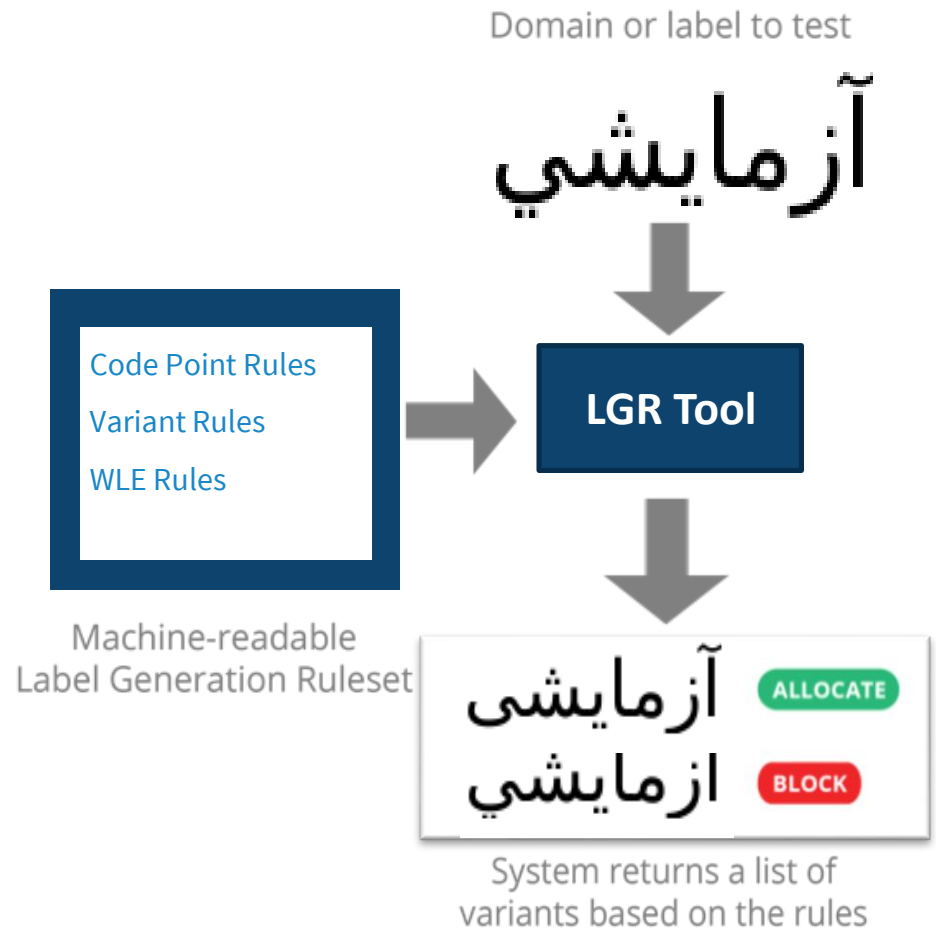
- ◉ Label Generation Rulesets (LGRs) used to generate domain name labels, per [RFC 7940](#)

- ◉ Example: excerpt from MSR-2 XML file

```
...  
<range first-cp="0780" last-cp="07B0" tag="sc:Thaa" ref="3"/>  
<char cp="07B1" tag="sc:Thaa" ref="5"/>  
<char cp="08A0" tag="sc:Arab" ref="12"/>  
<range first-cp="08A2" last-cp="08AC" tag="sc:Arab" ref="12"/>  
<range first-cp="08E4" last-cp="08EF" tag="sc:Arab" ref="12"/>  
<range first-cp="08F4" last-cp="08FE" tag="sc:Arab" ref="12"/>  
<range first-cp="0901" last-cp="0903" tag="sc:Deva" ref="0"/>  
<char cp="0904" tag="sc:Deva" ref="6"/>  
<range first-cp="0905" last-cp="0939" tag="sc:Deva" ref="0"/>  
<range first-cp="093A" last-cp="093B" tag="sc:Deva" ref="11"/>  
<char cp="093C" tag="sc:Deva" ref="0"/>  
<range first-cp="093E" last-cp="094D" tag="sc:Deva" ref="0"/>  
<char cp="094F" tag="sc:Deva" ref="11"/>  
<range first-cp="0956" last-cp="0957" tag="sc:Deva" ref="11"/>  
<char cp="0972" tag="sc:Deva" ref="9"/>  
<range first-cp="0973" last-cp="0977" tag="sc:Deva" ref="11"/>  
<range first-cp="0979" last-cp="097A" tag="sc:Deva" ref="10"/>  
<range first-cp="097B" last-cp="097C" tag="sc:Deva" ref="8"/>  
<range first-cp="097E" last-cp="097F" tag="sc:Deva" ref="8"/>  
<range first-cp="0981" last-cp="0983" tag="sc:Beng" ref="0"/>  
...
```

LGR Toolset (beta)

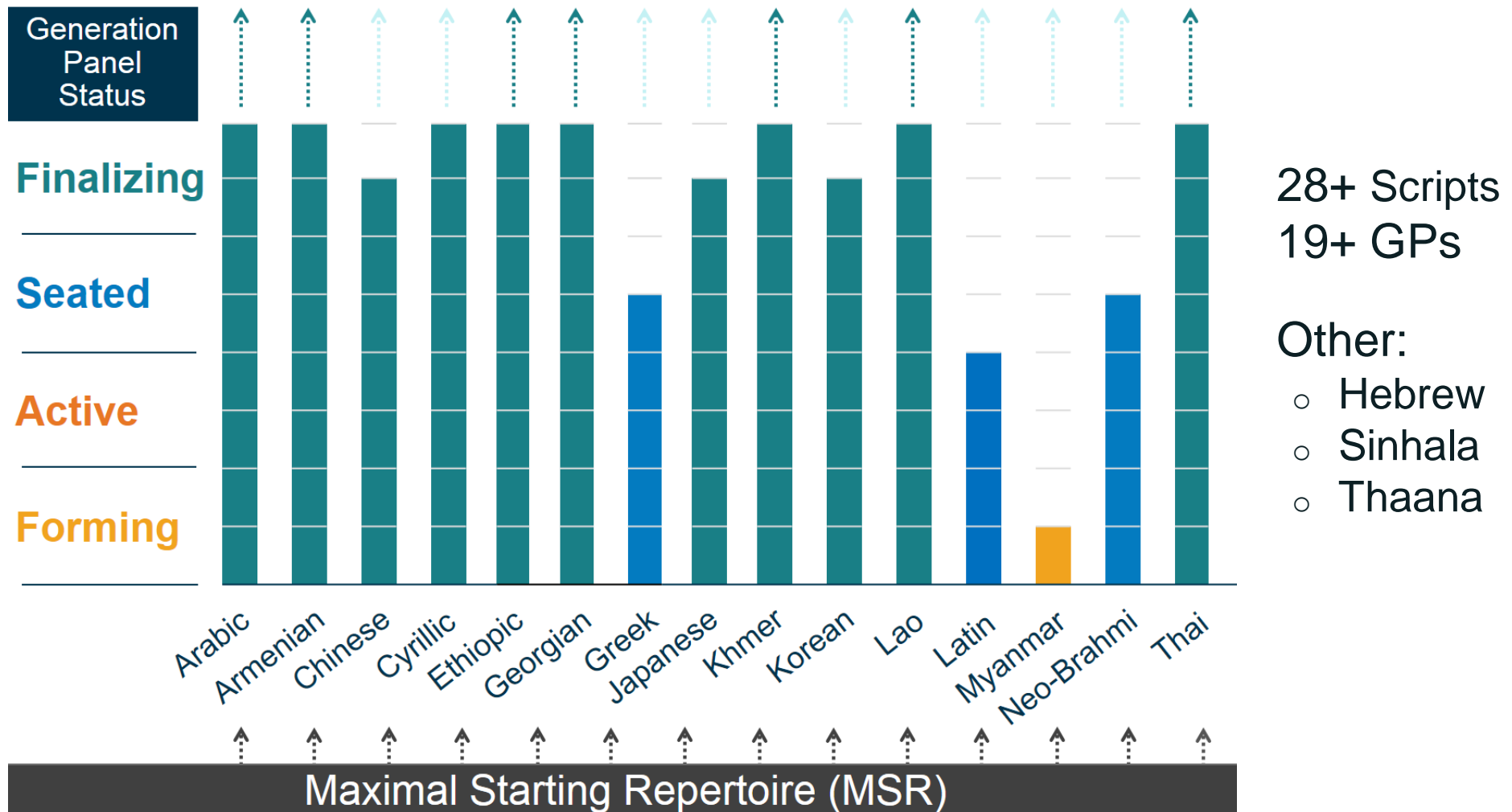
- LGR Toolset allows for the following:
 - Create a LGR
 - Use a LGR to validate label and variants
 - Manage LGRs
- Online beta deployment
 - Visit <https://lgrtool.icann.org/>
- Open source package(s) release
 - Released at github: [lgr-core](#), [lgr-django](#), [munidata](#)
- [User guide](#) available for further details



Root Zone LGR Development Status

Label Generation Rules (LGR)

May 2017



Linguistic Diversity in Africa

- More than 3000 languages in Africa

- Writing systems include:

- Ge'ez (Ethiopic)
- N'ko
- Tifinagh
- Vai
- ...
- Arabic (introduced)
- Latin (introduced)

የሰው ፡ ልጅ ፡ ሁሉ ፡ ሲወለድ ፡ ነጻና ፡ በክብርና ፡ በመብትም ፡
 እኩልነት ፡ ያለው ፡ ነው ፡ የተፈጥሮ ፡ የማስተዋልና ፡ ሕሊናው ፡
 ስላለው ፡ አንዱ ፡ ሌላውን ፡ በወንድማማችነት ፡ መንፍስ ፡
 መመልከት ፡ የገባዋል ።

ዳጋማ ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡
 ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡
 ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡
 ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡ ለሌሎች ስህጻን ስለሆነ ፡

የሰው ፡ ልጅ ፡ ሁሉ ፡ ሲወለድ ፡ ነጻና ፡ በክብርና ፡ በመብትም ፡
 እኩልነት ፡ ያለው ፡ ነው ፡ የተፈጥሮ ፡ የማስተዋልና ፡ ሕሊናው ፡
 ስላለው ፡ አንዱ ፡ ሌላውን ፡ በወንድማማችነት ፡ መንፍስ ፡
 መመልከት ፡ የገባዋል ።

The Characters of the N'Ko Script

consonants	ፈ	ፑ	ገ	ተ	ወ	ገ	ገ	ፑ	ፑ	ፑ
	m	l	k	r	d	c'	j	t	p	b
	ግ	ግ	ፑ	ገ	ገ	ፑ	ፑ	ፑ	ፑ	ፑ
	ŋ	y	w	h	n	ɲ	f	gb	s	rr

Alphabet Nationale du Tchad (ANT)

A - LETTRES MINUSCULES

a ' b b̄ c ch d d̄ dh e

ε ə f g gb h h̄ i ì j

k kh kp l m mb mv n ñ nd

ng nj η o ɔ p q r rh s

sl t th u v vb w y ȳ z zl



Repertoire

What is the Goal?

- ⦿ Goal is to create a mnemonic system for use in the Domain Name System (DNS)
 - ⦿ A mechanism to remember IP address
 - ⦿ Must remain secure and stable in use – if DNS is confusing to users, then the motivation is not met
 - ⦿ Not required to completely cover a language or a script
 - ⦿ May not form labels which are words in a language
 - ⦿ Not restricted to “correct” spellings
 - ⦿ May not carry a meaning in the “lexical” sense

Principles

1. Longevity – stable across Unicode versions
2. Least Astonishment– take into account the population using a code point
3. Contextual Safety – sensitive to ways in which code point may be used in malicious ways
4. Conservatism – any code point inclusion decision is as conservative as practicable

Principles

5. Inclusion – default is excluded, then add code point which is safe based on usability and confusability
6. Simplicity – rules determining use should be simple to understand
7. Predictability – rules determining whether a code point is included are predictable for others to reach the same conclusion
8. Stability – if permitted, taking it out very hard

Principles

9. Letter – Code point “will be alphabetic” in RFC 1123. Same principle so exclude code points not normally used to write words or used for purposes other than writing words

Questions to Ask

1. Is it contained in the Maximal Starting Repertoire?
2. Is it used with the script defined in the scope of the GP
3. Is it suitable in identifiers?
 - a. Is it in widespread modern use?
 - b. Is it not technical / religious / limited use only?
 - c. Is it not really a punctuation / symbol?
 - d. Is it really necessary for representing identifiers?
4. Is the Unicode encoding of the code point stable?
 - a. Are there any rendering issues?

Questions to Ask

5. What are the DNS security & stability concerns? rendering issue, homoglyph of non-PVALID code points?
6. How accessible would a TLD containing that code point be?
 - a. Are there input/keyboard concerns?
7. What are the risks if the code point is not included?
8. What are the risks if it is?
9. Is it in tension with any of the Principles in any way?
10. Does it always appear in a fixed sequence?

“everyday, general purpose [use ...]”

Level	Label	Description
0	International	The language is widely used between nations in trade, knowledge exchange, and international policy.
1	National	The language is used in education, work, mass media, and government at the national level.
2	Provincial	The language is used in education, work, mass media, and government within major administrative subdivisions of a nation.
3	Wider Communication	The language is used in work and mass media without official status to transcend language differences across a region.
4	Educational	The language is in vigorous use, with standardization and literature being sustained through a widespread system of institutionally supported education.
5	Developing	The language is in vigorous use, with literature in a standardized form being used by some though this is not yet widespread or sustainable.

<https://www.ethnologue.com/about/language-status>

“everyday, general purpose [use ...]”

6a	Vigorous	The language is used for face-to-face communication by all generations and the situation is sustainable.
6b	Threatened	The language is used for face-to-face communication within all generations, but it is losing users.
7	Shifting	The child-bearing generation can use the language among themselves, but it is not being transmitted to children.
8a	Moribund	The only remaining active users of the language are members of the grandparent generation and older.
8b	Nearly Extinct	The only remaining users of the language are members of the grandparent generation or older who have little opportunity to use the language.
9	Dormant	The language serves as a reminder of heritage identity for an ethnic community, but no one has more than symbolic proficiency.
10	Extinct	The language is no longer used and no one retains a sense of ethnic identity associated with the language.

<https://www.ethnologue.com/about/language-status>

How to Document the Repertoire

- ⦿ Document general but relevant information
 - a) History of script
 - b) Script characteristics
 - c) Languages using the script – standard name, ISO 639 code, name in local script, places language is spoken, other relevant information (e.g. EGIDS no.)
 - d) Criteria of language included in analysis (and excluded from analysis)
 - e) Types of code points – which types are included and which code points are excluded
 - f) Table of code points – with evidence/reference of use for each code point and any additional relevant information

Sources of Information - Languages

- ⊙ National governmental sources
- ⊙ www.ethnologue.com website
- ⊙ www.omniglot.com website
- ⊙ Published research and books
- ⊙ Field research
- ⊙ Others?

Sources of Information - Languages

COUNTRY

LANGUAGES

STATUS

MAPS

FEEDBACK

Expand All

Collapse All

1 (National)	Show Details »
2 (Provincial)	Show Details »
3 (Wider communication)	Show Details »
4 (Educational)	Show Details »
5 (Developing)	Show Details »
6a (Vigorous)	Show Details »
6b (Threatened)	Show Details »
7 (Shifting)	Show Details »
8a (Moribund)	Show Details »
8b (Nearly extinct)	Show Details »
9 (Dormant)	Show Details »
9 (Second language only)	Show Details »
10 (Extinct)	Show Details »

⦿ Ethnologue - <http://www.ethnologue.com/country/ET/status>

Sources of Information - Languages

2 (Provincial)

Hide Details 

Afar

[aar] 2 (Provincial). Statutory provincial working language in Afar Region (1994, Constitution, Art 47). 1,280,000 in Ethiopia (2010 UNSD). L2 users: 22,800 in Ethiopia. 906,000 monolinguals (1994 census). Total users in all countries: 1,828,800 (as L1: 1,806,000; as L2: 22,800).

Oromo, West Central

[gaz] 2 (Provincial). Statutory provincial working language in Oromia Region (1994, Constitution, Article 47(3)). West Central Oromo [gaz] is lingua franca of the area. 8,920,000 (1994 census). 25,500,000 all Oromo speakers in Ethiopia (2007 census). Ethnic population: 30,000,000.

Somali

[som] 2 (Provincial). Statutory provincial working language in Somali Region (1994, Constitution, Article 47(3)). 4,610,000 in Ethiopia (2010 UNSD). L2 users: 95,600 in Ethiopia. 2,880,000 monolinguals.

Tigrigna

[tir] 2 (Provincial). Statutory provincial working language in Tigray Region (1994, Constitution, Article 47(3)). 4,320,000 in Ethiopia (2010 UNSD). L2 users: 147,000 in Ethiopia. 2,820,000 monolinguals. Total users in all countries: 7,899,400 (as L1: 7,747,400; as L2: 152,000).

Sources of Information - Languages

Tigrigna



LANGUAGE

FEEDBACK

A language of Ethiopia

ISO 639-3

tir

Alternate Names

Beta Israel, Tigray, Tigrinya

Population

4,320,000 in Ethiopia (2010 UNSD). L2 users: 147,000 in Ethiopia. 2,820,000 monolinguals. Total users in all countries: 7,899,400 (as L1: 7,747,400; as L2: 152,000).

Location

Tigray region border areas; Amhara and Afar regions.

Language Maps

Djibouti, Eritrea and Ethiopia

Language Status

2 (Provincial). Statutory provincial working language in Tigray Region (1994, Constitution, Article 47(3)).

Classification

Afro-Asiatic, Semitic, South, Ethiopian, North

Typology

SOV; noun head final; gender (masculine/feminine); definite article; verb affixes mark person, number, gender of subject; passives; aspect; 33 consonant and 14 vowel phonemes.

Language Use

Also use Amharic [amh]. Used as L2 by Amharic [amh], Kunama [kun], Xamtanga [xan].

Language Development

Literacy rate in L2: 27%. Taught in primary schools. Fully developed. Bible: 1956.

Language Resources

OLAC resources in and about Tigrigna

Writing

Ethiopic script [Ethi], used since 13th or 14th century.

Other Comments

Christian.

Sources of Information - Repertoire

- ⊙ References which could be used to demonstrate “everyday, general purpose [use ...]”
 - a) National standard published by the government
 - b) Books published by Ministry of Education, e.g. for primary school
 - c) Common publications, e.g. newspapers
 - d) Other?

Strategies for Documenting the Repertoire

- ⦿ Strategy 1: Code Point Analysis
 - ⦿ For each code point in MSR
 - ⦿ Determine if it is used by one or more languages included
 - one example is sufficient
 - ⦿ Determine if the code point is required for the language(s)
 - ⦿ Document reference and reason for inclusion
- ⦿ Strategy 2: Language Analysis
 - ⦿ For each included language short-listed by GP
 - ⦿ Determine the required code points
 - ⦿ Document reference and reason for inclusion
 - ⦿ Review code points which are not analyzed

Example

Item #	Unicode Code Point	Glyph	Name and GC	Some languages using the character	Language, with EGIDS value	Reference
1	0621	ﺀ	ARABIC LETTER HAMZA;Lo	Arabic, Urdu, Punjabi, Sindhi	1 Arabic	[RFC 5564]
			...			
3	0623	ﺀ	ARABIC LETTER ALEF WITH HAMZA ABOVE;Lo	Arabic, Malay, Torwali	1 Arabic	[RFC 5564]
				
81	06AE	ﻙ⋮	ARABIC LETTER KAF WITH THREE DOTS BELOW;Lo	L'Alphabet National du Tchad (ANT)	1 ANT	[ANT]

Exercise

Item #	Unicode Code Point	Glyph	Name and GC	Some languages using the character	Language, with EGIDS value	Reference
1						
2						
3						
4						
5						

A world map where the continents are defined by a network of white dots and connecting lines, set against a solid orange background. The dots vary in size, and the lines are thin and white. The word "Variants" is written in white, bold, sans-serif font on the left side of the map.

Variants

What is the Goal?

- ⦿ Successfully defining variant rules for an LGR is not trivial
- ⦿ Code point or code point sequences causing two (or more) labels functionally “the same” in a script
- ⦿ Make the mnemonic system to minimize user confusion
- ⦿ Conservatism requires
 - ⦿ maximizing “blocked” variants
 - ⦿ minimize “allocatable” variants

Questions to Ask

1. Would a reasonable person with native knowledge of the script consider a pair of code points interchangeable?
2. Would such a person be unable to determine which of these interchangeable code points was used by appearance?
3. Is there an alternative representation?
4. What should the disposition of any defined variants be?
5. Should any of the variants of this code point be contingent on context?
6. Is each set of code point variant mappings symmetric?

Questions to Ask

7. Is each set of code point variant mappings transitive?
8. Are any variants contemplated that are in tension with any of the Principles?
9. Are the variants designed so that they lead to the minimal required number of allocatable variant labels?
10. Are the variants designed so that, in doubtful cases, they block potential variant labels?

Variant Relationships and Types

- ⦿ Variants are symmetric
 - ⦿ $A = B \Rightarrow B = A$
- ⦿ Variants are transitive
 - ⦿ $A = B \text{ and } B = C \Rightarrow A = C$
- ⦿ Variant code points can be of two types
 - ⦿ Allocatable
 - ⦿ Blocked
- ⦿ The types are directional
- ⦿ Label disposition calculated based on types of individual code points
 - ⦿ A single blocked type causes the whole label to be blocked

Example

0641	ف	فب	بف	بفب	1 (06A7)
06A2	فا	فبا	بفا	بفبا	5 (0641)
06A7	ف	فب	بف	بفب	5 (0642)
0642	ق	قب	بق	بقب	6 (06A7)

0641	06A2	a	Used interchangeably in Africa for languages using Western (African) orthography
0641	06A7	b	
0641	0642	b	
06A2	06A7	b	
06A2	0642	b	
06A7	0642	a	Used interchangeably in Africa for languages using Western (African) orthography

Whole Label Evaluation (WLE) Rules

Goal

- ⦿ Goal is to reduce label space
 - ⦿ Preventing labels which should not be possible for various reasons
 - ⦿ Not licensed by the script (but not spelling rules)
 - ⦿ Cause security issues
 - ⦿ Cause usability constraints
 - ⦿ Other?
 - ⦿ Reducing allocatable label by making them blocked in certain cases
 - ⦿ Put in contextual contexts for code points or their sequences

Examples

- ⦿ Cannot mix Persian Kaf and Arabic Kaf
- ⦿ Combining vowel mark must follow a consonant in Lao script
- ⦿ Subjoining consonant must follow a consonant in Khmer script
- ⦿ A label cannot start with a combining mark



LGR XML Specification

LGR XML Structure

- ⊙ “lgr” element has three sub-elements:
 - “meta”: all meta-data associated with the LGR, such as its authorship, what it is used for, implementation notes and references.
 - “data”: the substantive code point data
 - “rules” (optional): information on contextual and whole-label evaluation rules, if any
 - with any specific “action” elements providing the disposition of labels and their variants

LGR XML Structure

```
<?xml version="1.0"?>  
  <lgr xmlns="urn:ietf:params:xml:ns:lgr-1.0">  
  
    <meta> ... </meta>           //optional  
  
    <data> ... </data>  
  
    <rules> ... </rules>         //optional  
  
  </lgr>
```

Demo of online LGR Tool

Engage with ICANN



Thank You and Questions

Reach us at: IDNProgram@icann.org

Website: icann.org/idn



twitter.com/icann



[gplus.to/icann](https://plus.google.com/icann)



facebook.com/icannorg



weibo.com/ICANNorg



linkedin.com/company/icann



flickr.com/photos/icann



youtube.com/user/icannnews



slideshare.net/icannpresentations