

# **Technical and Standardisation Challenges for Nepalese Writing**

**Pat Hall**

**Language Technology Kendra**

**Nepal**

# the changing way we write



# the changing way we write



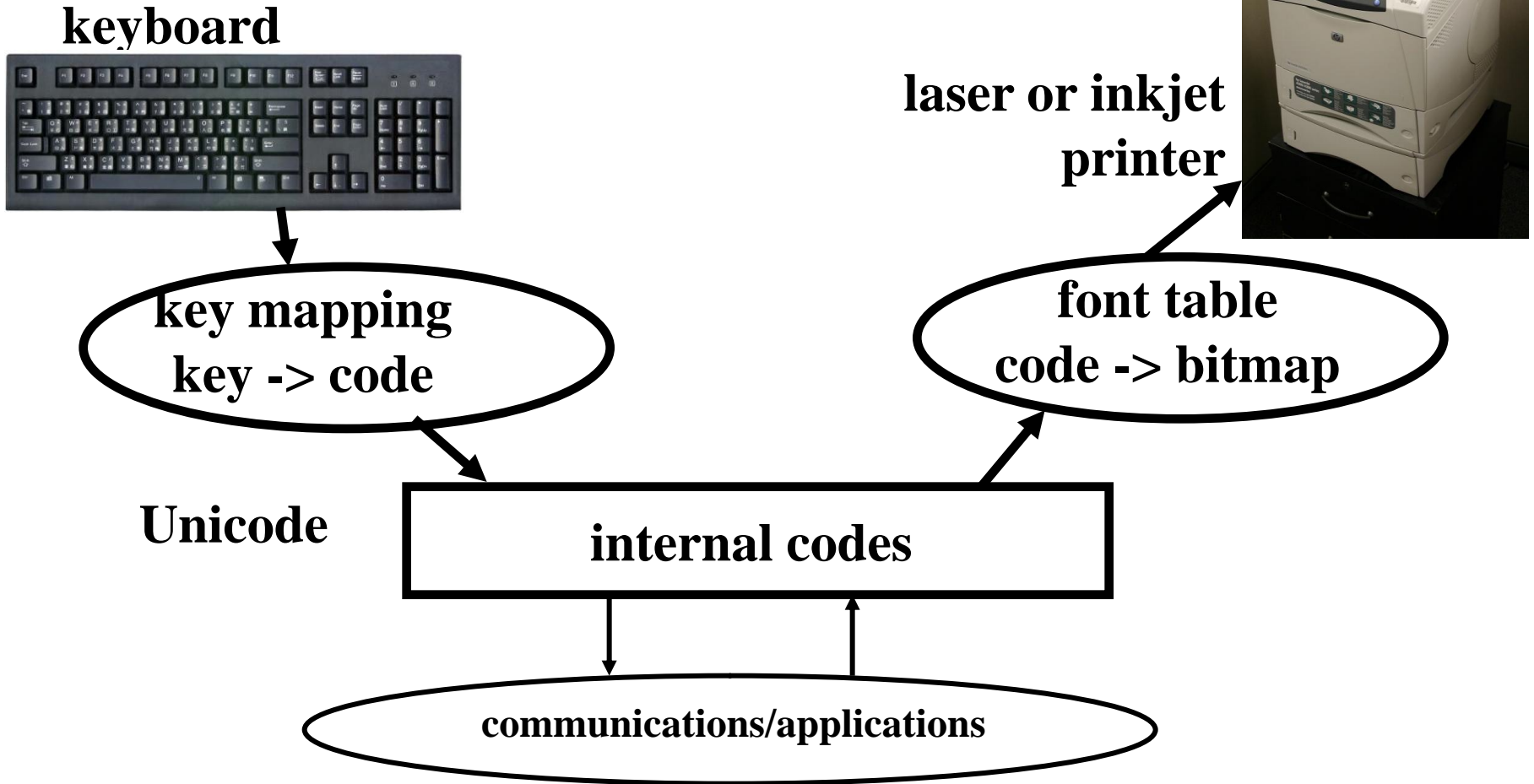
# the changing way we write



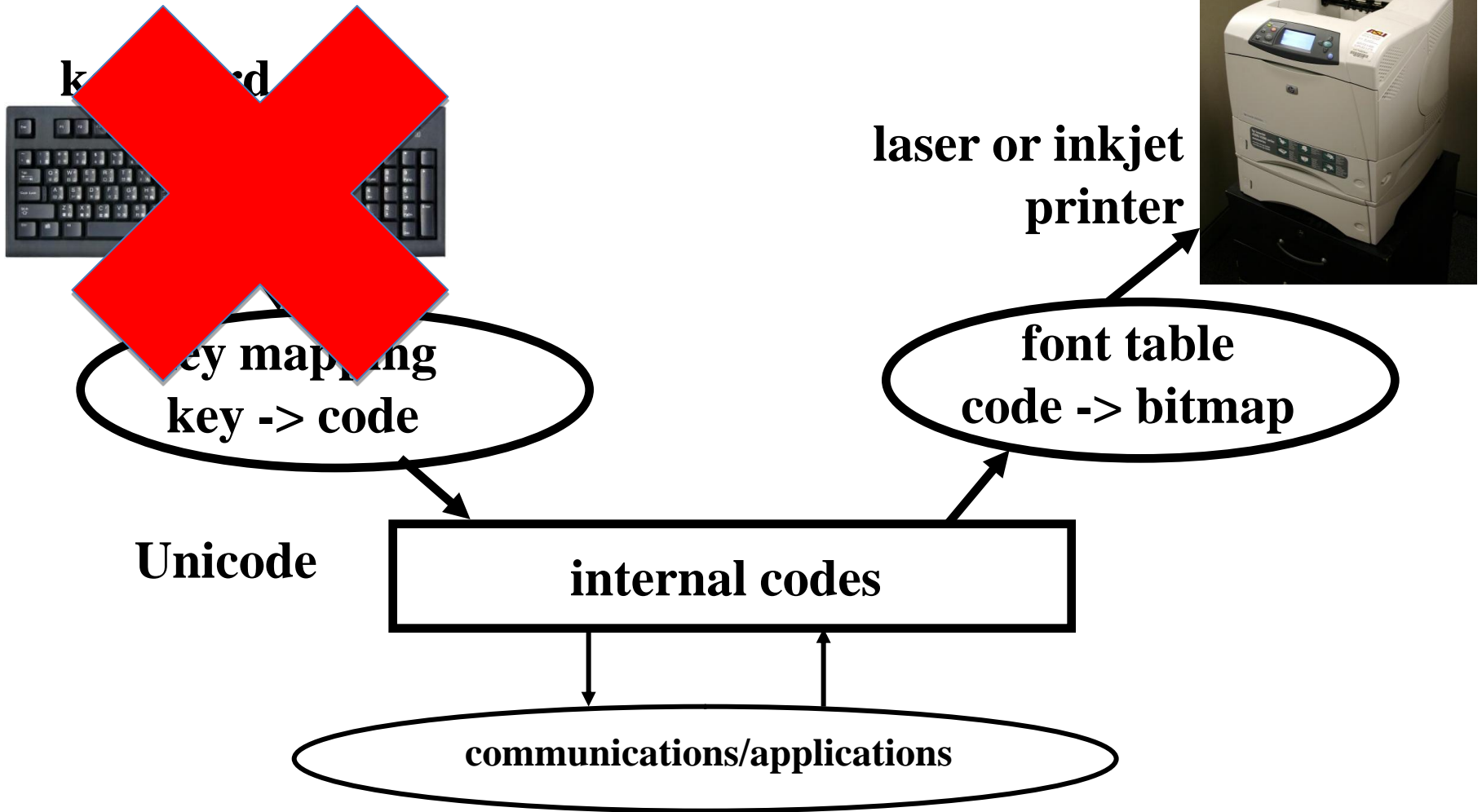
# the changing way we write



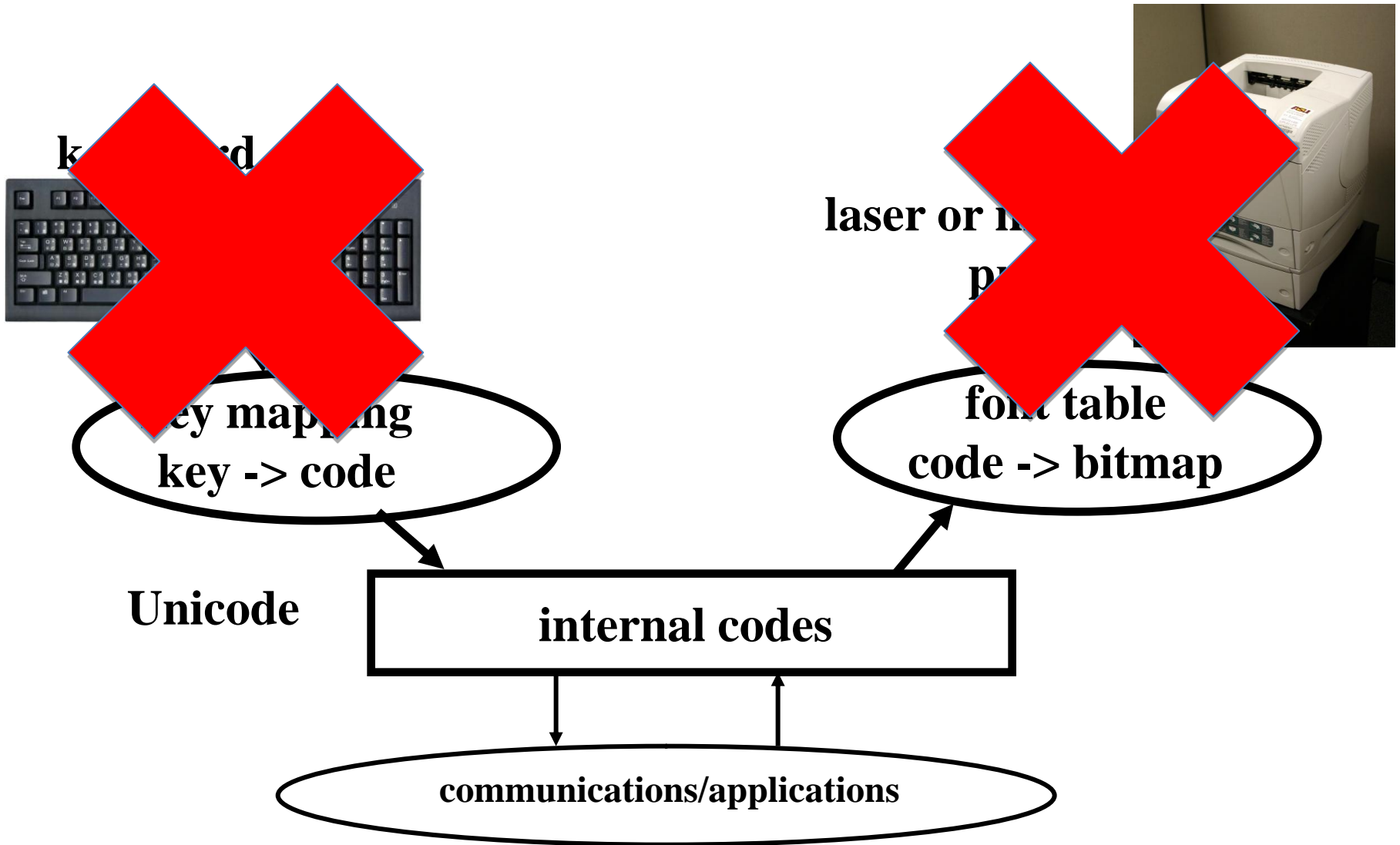
# but what matters is



# but what matters is

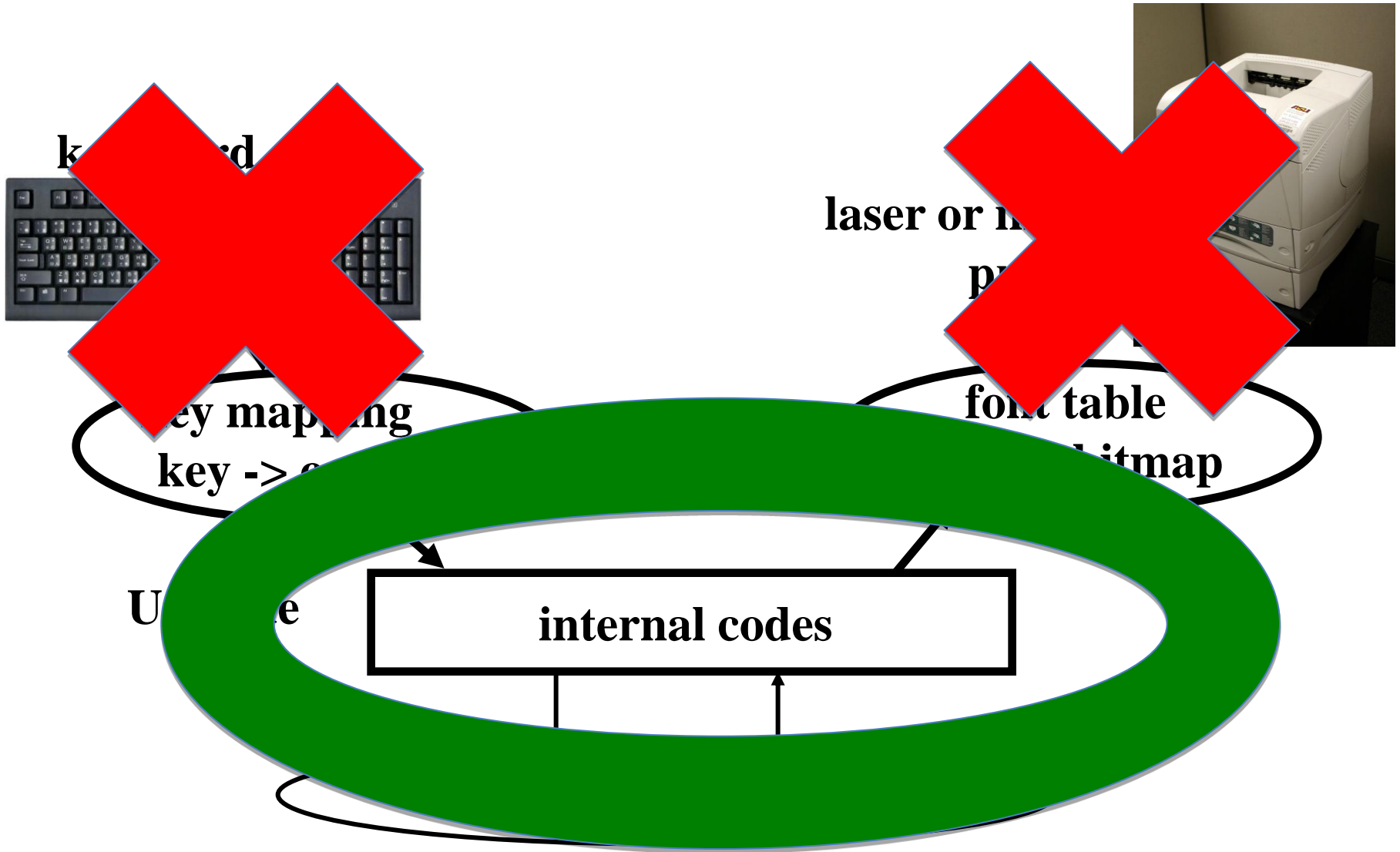


# but what matters is





# but what matters is



**deep structure characters**

# Example 1. Limbu

- written in Devanagari under one language policy
  - people know this
  - technology available
  - glottal stop added to Unicode Devanagari
- but had own writing system  
Sirijanga
  - created 18<sup>th</sup> Century, revived 1920s
  - Unicode table 2002 by Michailovsky and Everson

	190	191	192	193	194
0	𑄀	𑄁	𑄂	𑄃	𑄄
1	Z	𑄆	𑄇	𑄈	
2	𑄉	𑄊	𑄋	𑄌	
3	𑄍	𑄎	𑄏	𑄐	
4	𑄑	𑄒	𑄓	𑄔	𑄕
5	𑄖	𑄗	𑄘	𑄙	𑄚
6	𑄛	𑄜	𑄝	𑄞	𑄟
7	𑄠	𑄡	𑄢	𑄣	𑄤
8	𑄥	𑄦	𑄧	𑄨	𑄩
9	𑄪	𑄫	𑄬	𑄭	S
A	𑄮	𑄯	𑄰	𑄱	X
B	𑄲	𑄳	𑄴	𑄵	𑄶
C	𑄷	𑄸			𑄹
D	𑄺				𑄻
E	𑄼				V
F	Z				𑄿

# Limbu text in both writing systems

འཕགྱུམ། །། རྩོམ་ཅི།  
རྩུང་ལྷ། གྱིང་རྩུ རྩུ རྩུའི  
རྩུའི་ལྷ། གྱིང་, རྩུའི་ལྷ།  
གྱིང་, གྱིང་ ལྷ། རྩུའི་ལྷ།  
ལྷ།, རྩུའི་ལྷ། ལྷ། རྩུའི་ལྷ།  
ལྷ། རྩུའི་ལྷ། ལྷ། རྩུའི་ལྷ། །།

खेप्सुङ् ॥ आङ्गा?  
कुवाङ्भे?साबा नु नेहि  
तुकमाल्ले सिदा? केमेप्पा  
सिदा?साबा पोखाङ्बाल्ले  
खुने? कुहिम्मो मेम्बे:क्के  
आनिङ्वा?ए मेद्येनेन् ॥

**one-to-one correspondence**

**one language two writing systems - could we have one encoding of the language with two fonts?**

# one language multiple writing systems

## not unprecedented but unusual

- **Turkish**
  - Arabic to Roman
- **Central Asia**
  - Arabic to Cyrillic to Roman
- **Sindhi and Punjabi**
  - in Arabic in Pakistan
  - in Devanagari in India

encode the language  
the phonemes  
not the writing!  
re-map fonts

# Example 2. Newari / Nepal Bhasha

- **oldest written tradition in Nepal.**
- **1990s hack fonts produced**
- **2000 Unicode proposal drafted by Everson**
- **2008/10 Unicode drafts from Newars**
- **meetings of Nepal Lipi Guthi**
- **2010 Ranjana Unicode project driven by northern antiquarian interests.**

# 2010 Unicode Ranjana project

“Since Ranjana is visually and structurally similar to the Lantsa and Wartu scripts used for Buddhist Sanskrit documents in Tibet (China), Bhutan, Mongolia, Nepal, Sikkim and Ladakh (India) it has been considered would be practical to merge these two scripts (Lantsa and Wartu) with Ranjana for encoding purposes.”



the emphasis was on the writing, with much reference to antiquarian texts

the Newars pleaded for Prachalit, and after much fierce debate, the project ended within a year with no proposal

“Analysis of the chief features suggests that a number of these can be unified on structural grounds.”

# Analysis of Rabison Shakya's book

prachalit

seven more!

why?

ranjana

bhujimmola

क ka	ख kha	ग ga	घ gha	ङ ṅa	च ca
ख cha	ज ja	झ jha	ञ ṅa	ट ta	ठ tha
ड da	ढ dha	ण na	त ta	थ tha	द da
ध dha	न na	प pa	फ pha	ब ba	भ bha
म ma	य ya	र ra	ल la	व va	श ṣa
ष ṣa	स sa	ह ha	क्ष kṣa	त्र tra	ज्ञ jña

क ka	ख kha	ग ga	घ gha	ङ ṅa	झ ṅha
च ca	छ cha	ज ja	झ jha	ञ ṅa	ट ta
ठ tha	ड da	ढ dha	ण na	त ta	थ tha
द da	ध dha	न na	न्ह nha	प pa	फ pha
व ba	भ bha	म ma	म् mha	य ya	ह्य hya
र ra	ह rha	ल la	ल् lha	व va	व्ह vha
श ṣa	ष ṣa	स sa	ह ha	क्ष kṣa	त्र tra

क ka	ख kha	ग ga	घ gha	ङ ṅa	च ca
ख cha	ज ja	झ jha	ञ ṅa	ट ta	ठ tha
ड da	ढ dha	ण na	त ta	थ tha	द da
ध dha	न na	प pa	फ pha	ब ba	भ bha
म ma	य ya	र ra	ल la	व va	श ṣa
ष ṣa	स sa	ह ha	क्ष kṣa	त्र tra	ज्ञ jña

conjuncts?

conjuncts?

# prachalit rearranged in vargs

	unvoiced		voiced		nasalised	
1	क क ka	ख ख kha	ग ग ga	घ घ gha	ङ ङ ṅa	झ झ ṅha
2	च च ca	छ छ cha	ज ज ja	झ झ jha	ञ ञ ṅa	
3	ट ट ta	ठ ठ tha	ड ड da	ढ ढ dha	ण ण ṅa	
4	त त ta	थ थ tha	द द da	ध ध dha	न न na	न् न् ṅha
5	प प pa	फ फ pha	ब ब ba	भ भ bha	म म ma	म् म् ṅha

**Aspirated nasals!**  
**Not in Ranjana**  
**or Bhujimmola!**  
**Glyphs look like**  
**conjuncts!**  
*Are they*  
*conjuncts?*

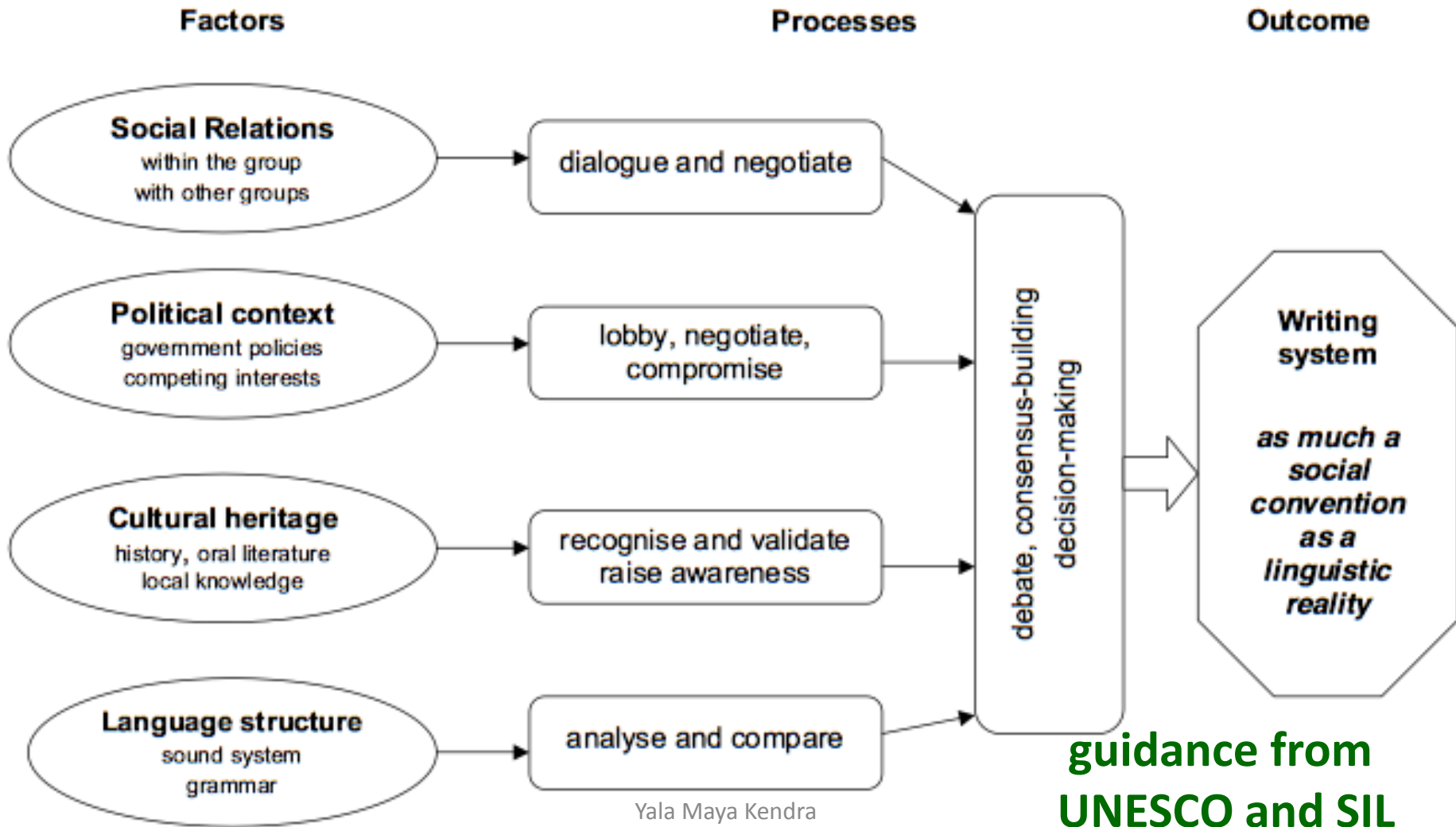
aspirated  
 unaspirated



# some possible conclusions

- **one language with many styles of writing**
  - **cannot map onto Devanagari**
    - **because more phonemes**
  - **but needs checking using contrastive pairs**
- **need a distinctive single unicode block**
  - **and many fonts, at least one per style**

# Example 3. Unwritten languages



# central stages

1. determine phonemic inventory

2. map to some writing system(s)

- new glyphs as necessary
- usually locally dominant system
- develop spelling system
- encode writing

In Nepal mostly  
based on Devanagari  
Michael Noonan survey  
Bhim Regmi's proposal

3. OR encode phonemes directly

- direct link to speech
  - TTS reads the text to people
  - ASR “writes” the speech
- can write in many scripts
- standarize codes and mapping to characters

Lohrung request  
for Roman

# but many want to write in Roman why not write in both?

नेपाली भाषा नेपाल बाहेक भारतमा पर्ना

nepālī bhāṣā nepāl bāhek bhārtmā pni

संवैधानिक हैसियत प्राप्त भाषा हो । नेपाल,

s̃vaidhānik haisiyt prāpt bhāṣā ho. nepāl,

भारत र संसारका अरू विभिन्न स्थानमा छरिएर

bhārt r s̃sārkā arū vibhinn sthānmā chrier

diacritic needs positioning

implicit short vowel missing

the Romanagari transliteration is rendered directly  
from the Devanagari codes following the  
Library of Congress ALA-LC romanization standards

# Example 4. Gurung and Magar

ᱠᱟᱨᱥᱟᱦᱟᱜ ᱠᱟᱨᱥᱟᱦᱟᱜ (गुरुङ्ग लिपि)

व्यन्जन (व्यन्जन)

ᱠ	ᱡ	ᱢ	ᱣ	ᱤ
क	ख	ग	घ	ङ
ᱥ	ᱦ	ᱧ	ᱨ	
च	छ	ज	झ	
ᱩ	ᱪ	ᱫ	ᱬ	
ट	ठ	ड	ढ	
ᱮ	ᱯ	ᱰ	ᱱ	ᱲ
त	थ	द	ध	न
ᱴ	ᱵ	ᱶ	ᱷ	ᱸ
प	फ	ब	भ	म
ᱺ	ᱻ	ᱼ	ᱽ	
य	र	ल	व	
᱿	᱾			
स	ह			

designs  
for new  
writing

Yala Maya Kendra

Magar (Akkha) Scripts

Vowels

ᱠ	ᱡ	ᱢ	ᱣ
A	a	I	ee
ᱥ	ᱦ	ᱧ	ᱨ
U	oo	Ey	ai

Consonants

[www.gorkhatimes.wordpress.com](http://www.gorkhatimes.wordpress.com)

ᱠ	ᱡ	ᱢ	ᱣ	ᱤ	
Ka	Kha	Ga	Gha	Nga	
ᱥ	ᱦ	ᱧ	ᱨ		
Cha	Chha	Ja	Jha		
ᱩ	ᱪ	ᱫ	ᱬ	ᱭ	
Ta	Tha	Da	Dha	Na	
ᱮ	ᱯ	ᱰ	ᱱ	ᱲ	ᱳ
Pa	Pha	Ba	Bha	Wa	ma
ᱴ	ᱵ	ᱶ	ᱷ	ᱸ	ᱹ
Sa	Ha	Ra	La	Ya	Nya

# what is going on?



Mark Turin

**marking identity by being visually distinct**

**could we mark language identity in font variants?**

# What should be standardized?

- **separate code sets for each language**
  - + **separate identity for each language**
  - **will require around 5,000 code points**
  - **extra cost per language**
- **one, or a few, unified code sets for Nepal**
  - + **fewer standards to write**
  - + **can share cost of technology**
  - **agreement required over many communities**
- **mapping from codes to characters**

# how should they be standardised?

- **Nepal Standards Organisation**
  - ✓ needs to be broadened to mirror ISO
  - ✓ must reach agreement within Nepal first
- **International Standards Organisation**
  - ✓ need champion member
    - India?
  - ✓ and champion expert
    - somebody who knows Nepal
  - ✓ attending meetings important but costly



# Thank You

---

## key issues for writing languages

- **the phonemes and their computer codes**
  - some way of writing these for debate
  - standardize mapping phonemes to characters
- **how to reach agreement**
  - enshrine that in standards