

# PALI: A Language Identification Benchmark for Perso-Arabic Scripts

Sina Ahmadi, Milind Agarwal, Antonios Anastasopoulos

{sahmad, magarwa, antonis} @ gmu.edu

## Research Question

Can hierarchical approaches:

1. Identify languages that use unconventional writing systems in low-resource settings?
2. Resolve confusion among languages in custom-trained or off-the-shelf models?

## Data

We focus on bilingual language communities that use the Urdu, Persian, or Arabic scripts

### Data Sources:

- Wikipedia dumps
- Local news websites
- Small pre-published corpora

### URDU

Brahui, Punjabi, Kashmiri, Sindhi, Saraiki, Torwali

Balochi

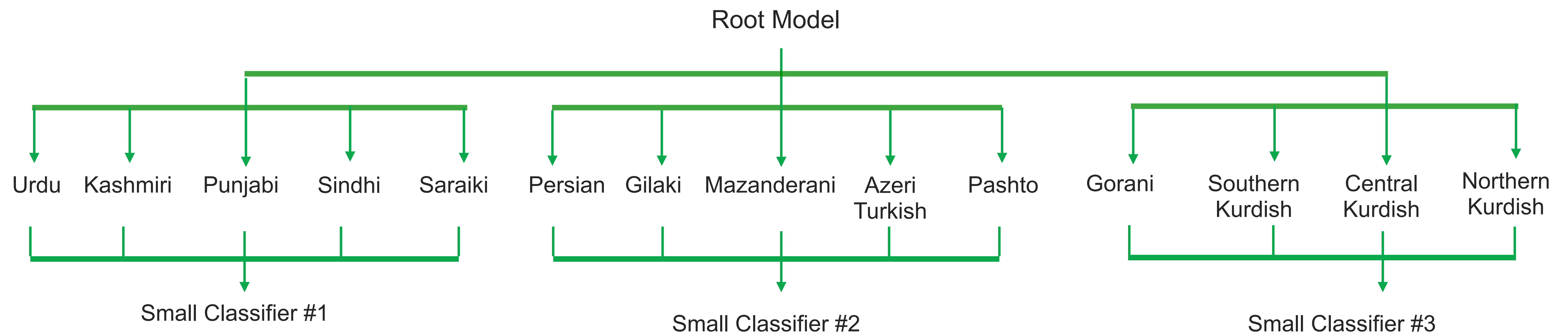
### PERSIAN

Pashto, Gilaki, Azeri Turkish, Mazanderani

### ARABIC

Northern/Central/Southern Kurdish, Gorani

Arabic



Urdu	15902	4	78	24	32
Kashmiri	3	15889	28	17	21
Punjabi	33	33	15782	26	95
Sindhi	10	5	12	15800	13
Saraiki	32	37	62	34	15818
	Urdu	Kashmiri	Punjabi	Sindhi	Saraiki

Noise %	Sentence
Clean	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا Second Kurdish photographers' exhibition in Belgium
20	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا
40	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا
60	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا
80	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا
100	دووه مین پیشانگه ها فوتوگرافه رین کورد ل به لجیکا

## Result Highlights

- Off-the-shelf and custom trained systems generally perform worse in noisy settings
- A confusion-based hierarchical classification can help build on top of a reasonable root system
- Such modeling can help identify languages that use unconventional writing systems

## Methodology

- Manually create script maps and generate data at different noise levels
- Inspect classification systems to identify highly-confused language clusters
- Off of the best root system, train specialized hierarchical classifiers to resolve confusion

### OURS (HIER)

### CUSTOM-TRAINED

### OFF-THE-SHELF

	F1
Hier	0.95
Root	0.94
MNB	0.10
MLP	0.10
fastText	0.27
CLD3	0.09
langid.py	0.13
Franc	0.13

## Check out our GitHub!

- Clean, noisy, merged data
- Compressed trained models for 19 languages
- All preprocessing and training scripts

