# Evaluation of Finite State Morphological Analyzers Based on Paradigm Extraction from Wiktionary

Ling Liu
Mans Hulden

FSMNLP
Sep 6, 2017

# Outline

‣ Motivation of the Study

‣ How the Morphological Analyzer works

‣ Data

‣ Evaluation and Result

# Outline

‣ **Motivation of the Study**

‣ How the Morphological Analyzer works

‣ Data

‣ Evaluation and Result

# Motivation of the Study

‣ Wiktionary: morphological inflection tables for many languages

| | indicative | |
|---|---|---|
| **present** | ich schreibe | wir schreiben |
| | du schreibst | ihr schreibt |
| | er schreibt | sie schreiben |
| **preterite** | ich schrieb | wir schrieben |
| | du schriebst | ihr schriebt |
| | er schrieb | sie schrieben |
| **imperative** | schreib (du) schreibe (du) | schreibt (ihr) |

WIKTIONARY
*the free dictionary*

‣ Wiktionary Morphological Database: 350 languages

# Motivation of the Study

‣ Forsberg and Hulden (2016): a method to convert morphological inflection tables into unweighted and weighted finite transducers for parsing and generation

- Evaluated on German, Spanish, Finnish

| Language | | Lemma | L+MSD | MSD |
|---|---|---|---|---|
| German | nouns | 77.06 | 69.44 | 79.50 |
| | verbs | 90.02 | 89.76 | 92.78 |
| Spanish | verbs | 96.92 | 96.92 | 97.43 |
| Finnish | nounadj | 70.29 | 69.68 | 91.59 |
| | verbs | 90.44 | 90.44 | 98.02 |

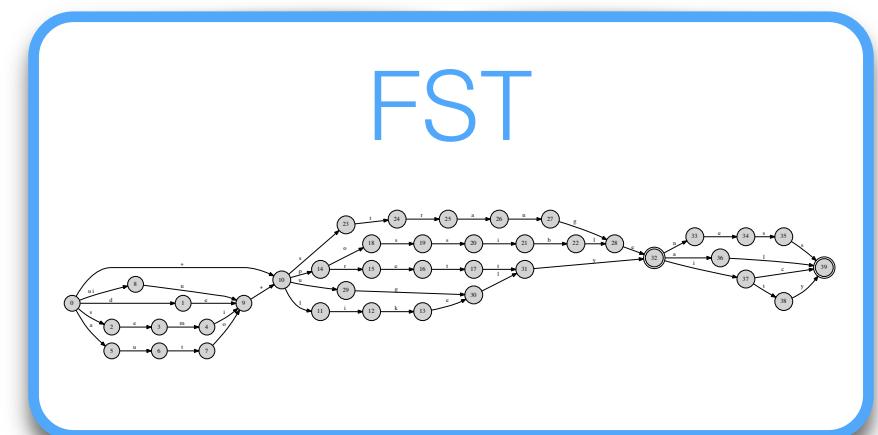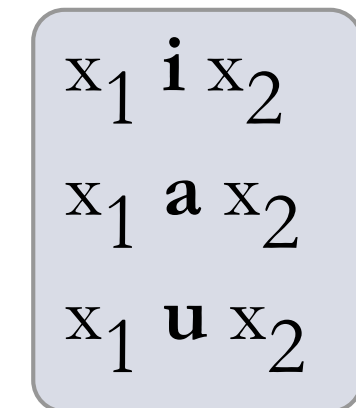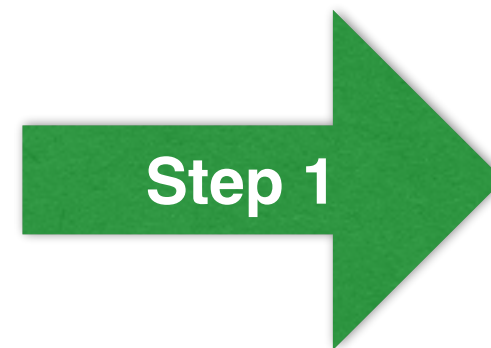**Evaluation of the weighted model (Forsberg and Hulden, 2016)**

# Outline

‣ Motivation of the Study

‣ **How the Morphological Analyzer works**

‣ Data

‣ Evaluation and Result

# Workflow of the Morphological Analyzer

## Inflection examples

| | indicative | |
|---|---|---|
| **present** | ich schreibe | wir **schreiben** |
| | du schreibst | ihr schreibt |
| | er schreibt | sie **schreiben** |
| **preterite** | ich schrieb | wir schrieben |
| | du schriebst | ihr schriebt |
| | er schrieb | sie schrieben |
| **imperative** | schreib (du) / schreibe (du) | schreibt (ihr) |

## Paradigmatic Models

$x_1 \mathbf{i} x_2$

$x_1 \mathbf{a} x_2$

$x_1 \mathbf{u} x_2$

**Step 1**

**Step 2**

WIKTIONARY
*the free dictionary*

## Ranking Analyses

blargashed
1. blargash[V;PST]
2. blargash[V;V.PTCP;PST]
3. blargashe[V;PST]
…

**Step 3**

FST

# Generalization from inflection tables

## Inflection examples

| | indicative | |
|---|---|---|
| present | ich schreibe | wir **schreiben** |
| | du schreibst | ihr schreibt |
| | er schreibt | sie **schreiben** |
| preterite | ich schrieb | wir schrieben |
| | du schriebst | ihr schriebt |
| | er schrieb | sie schrieben |
| imperative | schreib (du) schreibe (du) | schreibt (ihr) |

WIKTIONARY
*the free dictionary*

**Step 1**

## Paradigmatic Models

$$x_1 \; \mathbf{i} \; x_2$$
$$x_1 \; \mathbf{a} \; x_2$$
$$x_1 \; \mathbf{u} \; x_2$$

# Generalization

The **common parts (stem)** are calculated by extracting the **Longest Common Subsequence** from related forms

inflection table

**ring**
**rang**
**rung**
**rings**
**ringing**

*Ahlberg, Forsberg, Hulden (2014, 2015)

# Generalization

The **common parts (stem)** are calculated by extracting the **Longest Common Subsequence** from related forms

inflection table

**r** i **ng**

**r** a **ng**

**r** u **ng**

**r** i **ng** s

**r** i **ng** ing

*Ahlberg, Forsberg, Hulden (2014, 2015)

# Generalization

The **common parts (stem)** are calculated by extracting the **Longest Common Subsequence** from related forms

inflection table

**r i ng**

**r a ng**

**r u ng**

**r i ng s**

**r i ng ing**

LCS = rng

*Ahlberg, Forsberg, Hulden (2014, 2015)

# Generalization

The **common parts (stem)** are calculated by extracting the **Longest Common Subsequence** from related forms

inflection table

**r i ng**

**r a ng**

**r u ng**

**r i ng s**

**r i ng ing**

LCS = rng

$x_1 = \textbf{r}$

$x_2 = \textbf{ng}$

*Ahlberg, Forsberg, Hulden (2014, 2015)

# Generalization

Formal claim: the common parts (stem) are calculated by extracting the **Longest Common Subsequence** from related forms*

inflection table

$$\mathbf{r\ i\ ng}$$
$$\mathbf{r\ a\ ng}$$
$$\mathbf{r\ u\ ng}$$
$$\mathbf{r\ i\ ng\ s}$$
$$\mathbf{r\ i\ ng\ ing}$$

$$\underline{x_1} \quad \underline{x_2}$$

LCS = **rng**

$x_1 = \mathbf{r}$
$x_2 = \mathbf{ng}$

$\longrightarrow$

"paradigm"
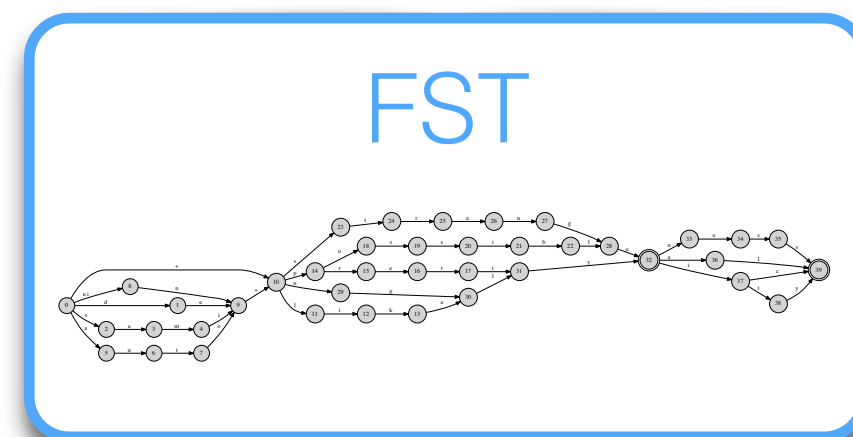
$x_1 + i + x_2$
$x_1 + a + x_2$
$x_1 + u + x_2$
$x_1 + i + x_2 + s$
$x_1 + i + x_2 + ing$

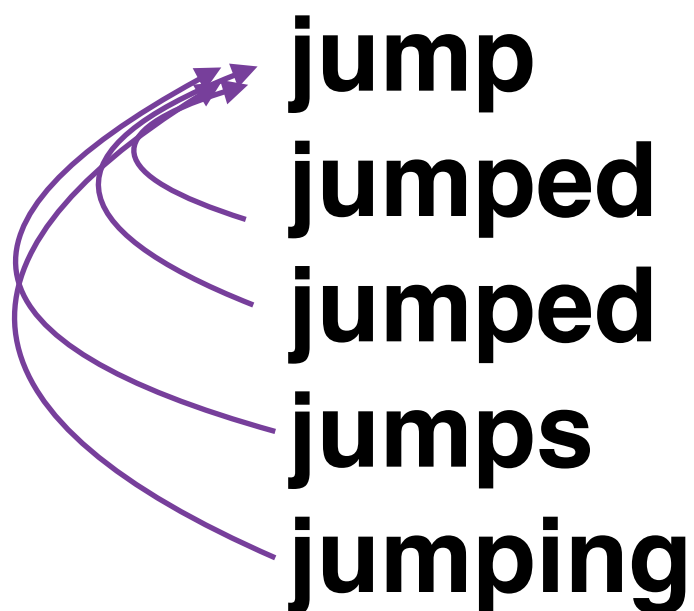*Ahlberg, Forsberg, Hulden (2014, 2015)

# Generalization

## Inflection tables

| |
|---|
| ring |
| rang |
| rung |
| rings |
| ringing |

| |
|---|
| jump |
| jumped |
| jumped |
| jumps |
| jumping |

| |
|---|
| drink |
| drank |
| drunk |
| drinks |
| drinking |

## Paradigms

| |
|---|
| $x_1+i+x_2$ |
| $x_1+a+x_2$ |
| $x_1+u+x_2$ |
| $x_1+i+x_2+s$ |
| $x_1+i+x_2+ing$ |

$x_1$
$x_1+ed$
$x_1+ed$
$x_1+s$
$x_1+ing$

| |
|---|
| $x_1+i+x_2$ |
| $x_1+a+x_2$ |
| $x_1+u+x_2$ |
| $x_1+i+x_2+s$ |
| $x_1+i+x_2+ing$ |

## Collapsed

| |
|---|
| $x_1+i+x_2$ |
| $x_1+a+x_2$ |
| $x_1+u+x_2$ |
| $x_1+i+x_2+s$ |
| $x_1+i+x_2+ing$ |

# From paradigm to FST

## Paradigmatic Models

$$x_1 \; \mathbf{i} \; x_2$$

$$x_1 \; \mathbf{a} \; x_2$$

$$x_1 \; \mathbf{u} \; x_2$$

**Step 2**

FST

# From paradigm to FST

## Lemmatization

| | | |
|---|---|---|
| **jump** | $x_1$ | **infinitive** |
| **jumped** | $x_1 + ed$ | **simp past** |
| **jumped** | $x_1 + ed$ | **past part** |
| **jumps** | $x_1 + s$ | **simp pres 3sg** |
| **jumping** | $x_1 + ing$ | **pre part** |

# From paradigm to FST

## Lemmatization

**jump**        $x_1$              **infinitive**

**jumped**      $x_1 + ed$         **simp past**

**jumped**      $x_1 + ed$         **past part**

**jumps**       $x_1 + s$          **simp pres 3sg**

**jumping**     $x_1 + ing$        **pre part**

# From paradigm to FST

## Lemmatization

**jump**     $x_1$     **infinitive**
**jumped**     $x_1 + ed$     **simp past**
**jumped**     $x_1 + ed$     **past part**
**jumps**     $x_1 + s$     **simp pres 3sg**
**jumping**     $x_1 + ing$     **pre part**



*$x_1$*            *ed:ε*

# From paradigm to FST

## Add inflection information

| | | |
|---|---|---|
| **jump** | $x_1$ | **infinitive** |
| **jumped** | $x_1 + ed$ | **simp past** |
| **jumped** | $x_1 + ed$ | **past part** |
| **jumps** | $x_1 + s$ | **simp pres 3sg** |
| **jumping** | $x_1 + ing$ | **pre part** |



jumped > jump[V;PST]

# From paradigm to FST

More lemmatization and analysis example

**rel**y          $x_1 + y$          **infinitive**
**rel**ied        $x_1 + ied$        **simp past**
**rel**ied        $x_1 + ied$        **past part**
**rel**ies        $x_1 + ies$        **simp pres 3sg**
**rel**ying       $x_1 + ying$       **pre part**

@

0 →@→ 1 →**<i:y>**→ 2 →**<e:ε>**→ 3 →**<d:ε>**→ 4 →ε:[→ 5 →**ε: V;PST**→ 6 →ε:]→ 7

relied > rely[V;PST]

# Building the analyzer

**Paradigm**

$x_1+i+x_2$

$x_1+a+x_2$

$x_1+u+x_2$

$x_1+i+x_2+s$

$x_1+i+x_2+ing$

<span style="color:purple">analyzers</span>

**Paradigm**

$x_1$

$x_1+ed$

$x_1+ed$

$x_1+s$

$x_1+ing$

<span style="color:purple">analyzers</span>

m transducers

$$\text{Analyzer} = \ f_1 \cup f_2 \cup \ldots \cup f_1 \cup \ldots \cup f_m$$

# From paradigm to FST



*jumped > jump[V;PST]*



*relied > rely[V;PST]*

verified                    tried                    died

verify[V;PST]              try[V;PST]               dy[V;PST]
verifi[V;PST]              tri[V;PST]               di[V;PST]
…                          …                        die[V;PST]

                                                    …

# Ranking Analyses

FST

**Step 3**

## Ranking Analyses

**blargashed**
1. **blargash**[V;PST]
2. **blargash**[V;V.PTCP;PST]
3. **blargashe**[V;PST]
…

# Language models over variables (WFSTs)

| | | |
|---|---|---|
| **jump** | $x_1$ | **infinitive** |
| **jumped** | $x_1 + ed$ | **simp past** |
| **jumped** | $x_1 + ed$ | **past part** |
| **jumps** | $x_1 + s$ | **simp pres 3sg** |
| **jumping** | $x_1 + ing$ | **pre part** |

*jump*
*watch*
*look*
*listen*
*work*
*ask*
*…*

Infer a language model!

# From paradigm to WFST

**jump**
**jumped**
**jumped**
**jumps**
**jumping**

$x_1$
$x_1 + ed$
$x_1 + ed$
$x_1 + s$
$x_1 + ing$

**infinitive**
**simp past**
**past part**
**simp pres 3sg**
**pre part**

$x_1$

```
LM_{X1}  --e:ε/0.0-->  ○  --d:ε/0.0-->  ○  --ε:[/0.0-->  ○  --ε: V;PST/0.0-->  ○  --ε:]/0.0-->  ◎
```

n-gram model                    jumped > jump[V;PST]

# From paradigm to WFST

**jump**
**jumped**
**jumped**
**jumps**
**jumping**

$x_1$
$x_1 + ed$
$x_1 + ed$
$x_1 + s$
$x_1 + ing$

**infinitive**
**simp past**
**past part**
**simp pres 3sg**
**pre part**

$x_1$

**LM$_{X1}$**

e:ε/0.0    d:ε/0.0    ε:[/0.0    ε: V;PST/0.0    ε:]/0.0

n-gram model

jumped > jump[V;PST]

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|---|---|---|---|---|---|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Example analysis (weighted)

**Both are correct**

| rank | log_prob | paradigm | variables | lemma | mst |
|------|----------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

**verified**

# Outline

‣ Motivation of the Study

‣ How the Morphological Analyzer works

‣ **Data**

‣ Evaluation and Result

# Data

‣ Wiktionary Morphological Database

‣ UniMorph project ([https://unimorph.github.io/index.html](https://unimorph.github.io/index.html))

‣ 55 Languages

‣ 19 Language groups

‣ 10 scripts

# Outline

‣ Motivation of the Study

‣ How the Morphological Analyzer works

‣ Data

‣ **Evaluation and Result**

# Evaluation Task

‣ Lemmatization and morphosyntactic information tagging

‣ 90% for training; and 10% for test (unless less than 50 inflection tables)

‣ The evaluation data is **disjoint** from the training data

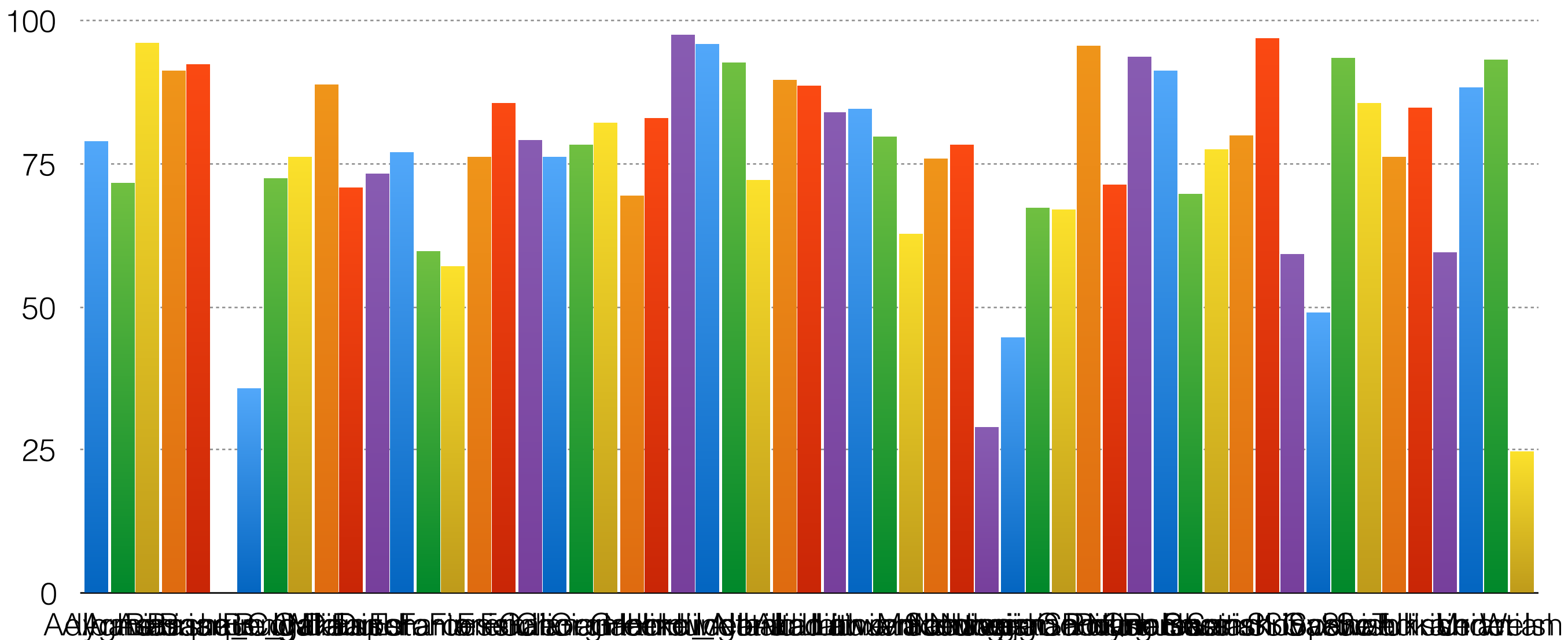‣ The first-ranked analyses

‣ Recall

-lemma

-lemma + POS

-lemma + MST

| rank | score | paradigm | variables | lemma | mst |
|------|-------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

# Evaluation Task

‣ Lemmatization and morphosyntactic information tagging

‣ 90% for training; and 10% for test (unless less than 50 inflection tables)

‣ The evaluation data is **disjoint** from the training data

‣ The first-ranked analyses

‣ Recall

  -lemma

  -lemma + POS

  -lemma + MST

| rank | score | paradigm | variables | lemma | mst |
|------|-------|----------|-----------|-------|-----|
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;PST]** |
| 1 | -11.44 | p4_unmarry | (1=verif) | **verify** | **[V;V.PTCP;PST]** |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;PST] |
| 2 | -18.36 | p1_dribble | (1=verifi) | verifie | [V;V.PTCP;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;PST] |
| 3 | -30.49 | p20_preempt | (1=verifi) | verifi | [V;V.PTCP;PST] |

# Result

‣ Paradigms are extracted successfully for all languages.

‣ **Lemmatization recall:**

- Low end: 0% (Basque)

- High end: 97.5% (Hindi)

‣ **Lemma-POS recall:**

- Low end: 0% (Basque)

- High end: 97.0% (Hindi)

‣ **Lemma-tag recall:**

-  Low end: 0% (Basque)
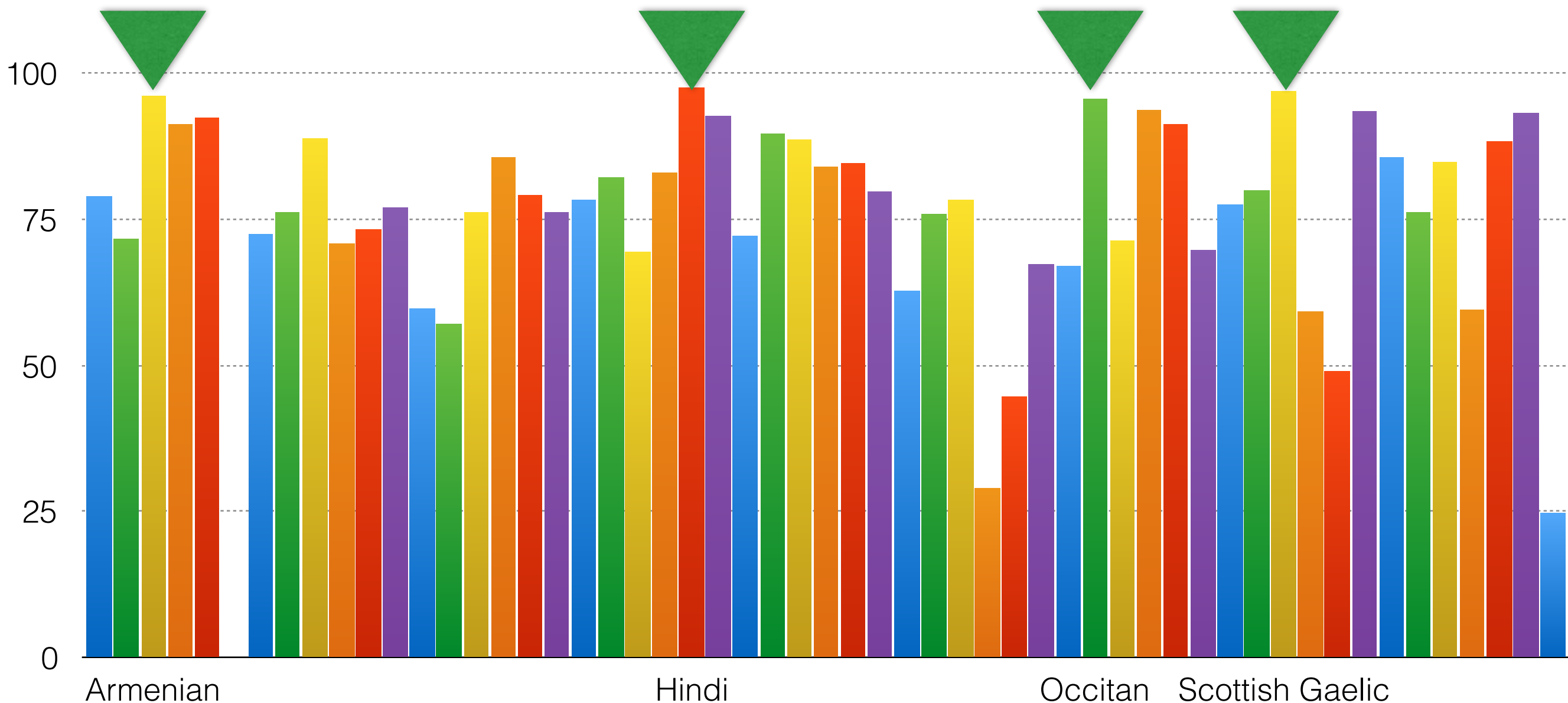
- High end: 96.9% (Hindi)

# Result Overview: Lemma Recall
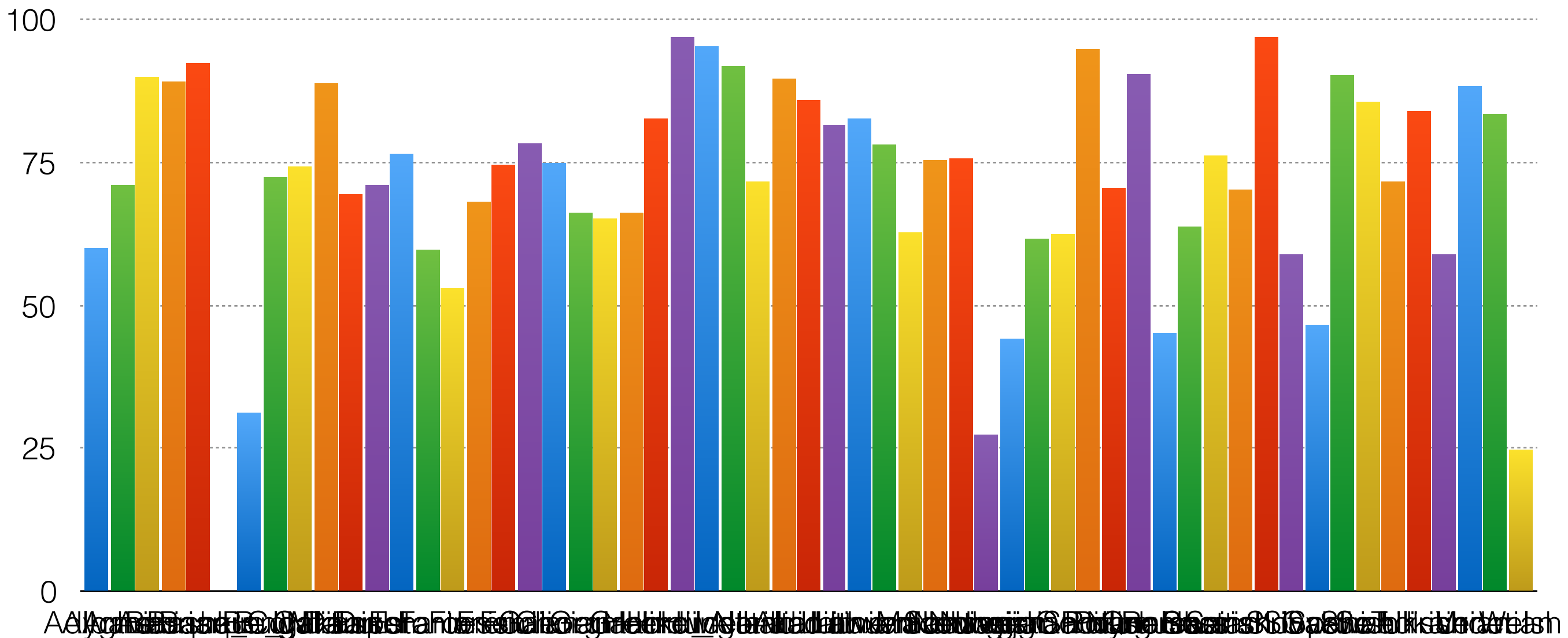
# Result Overview: Lemma Recall < 30%
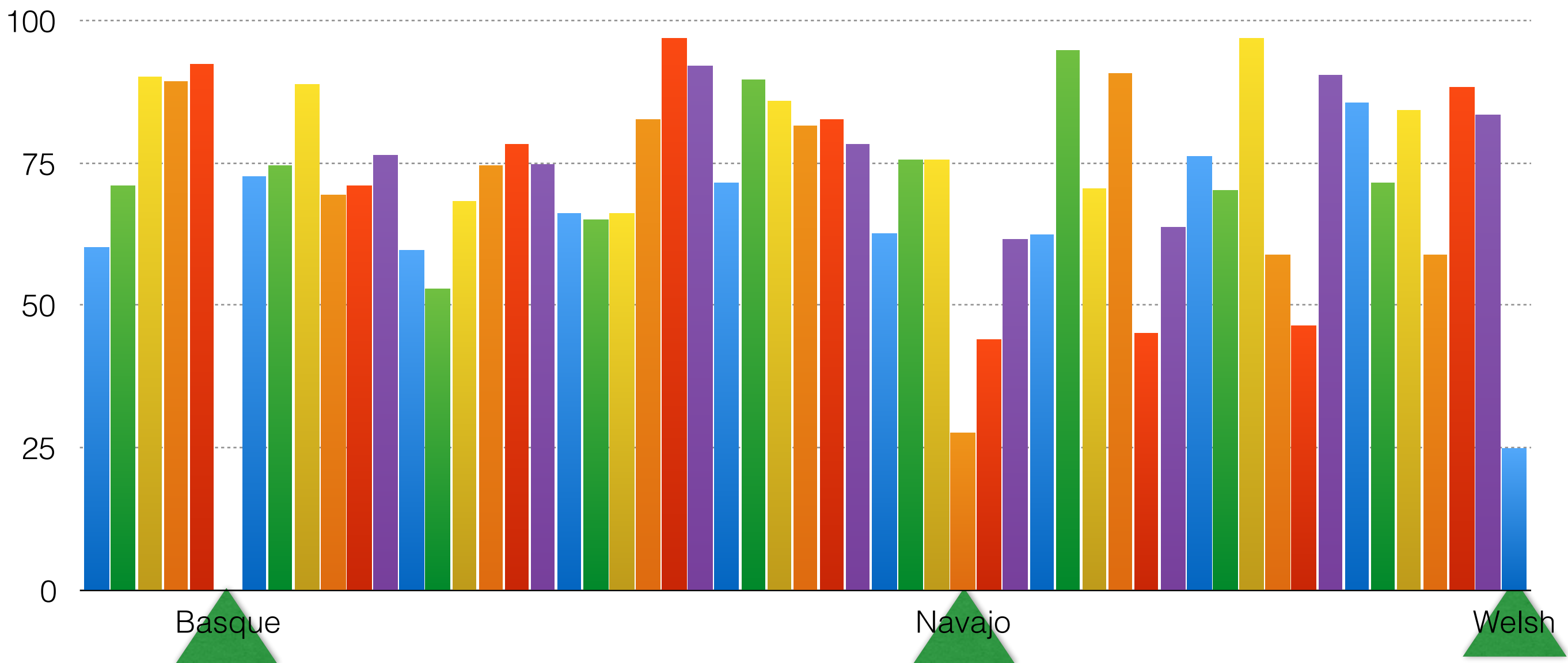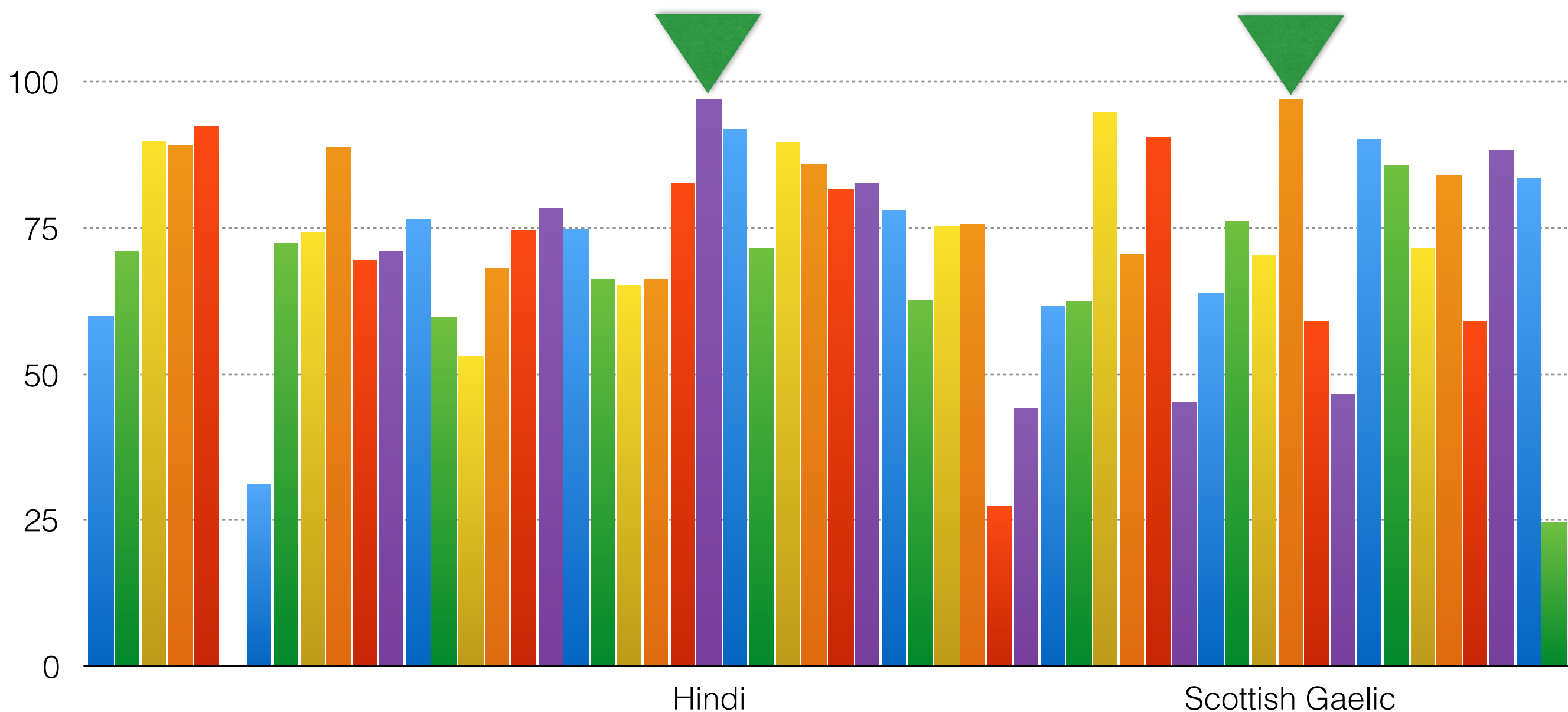
# Result Overview: Lemma Recall > 95%

# Result Overview: Lemma+POS Recall

Result Overview:
Lemma+POS Recall < 30%
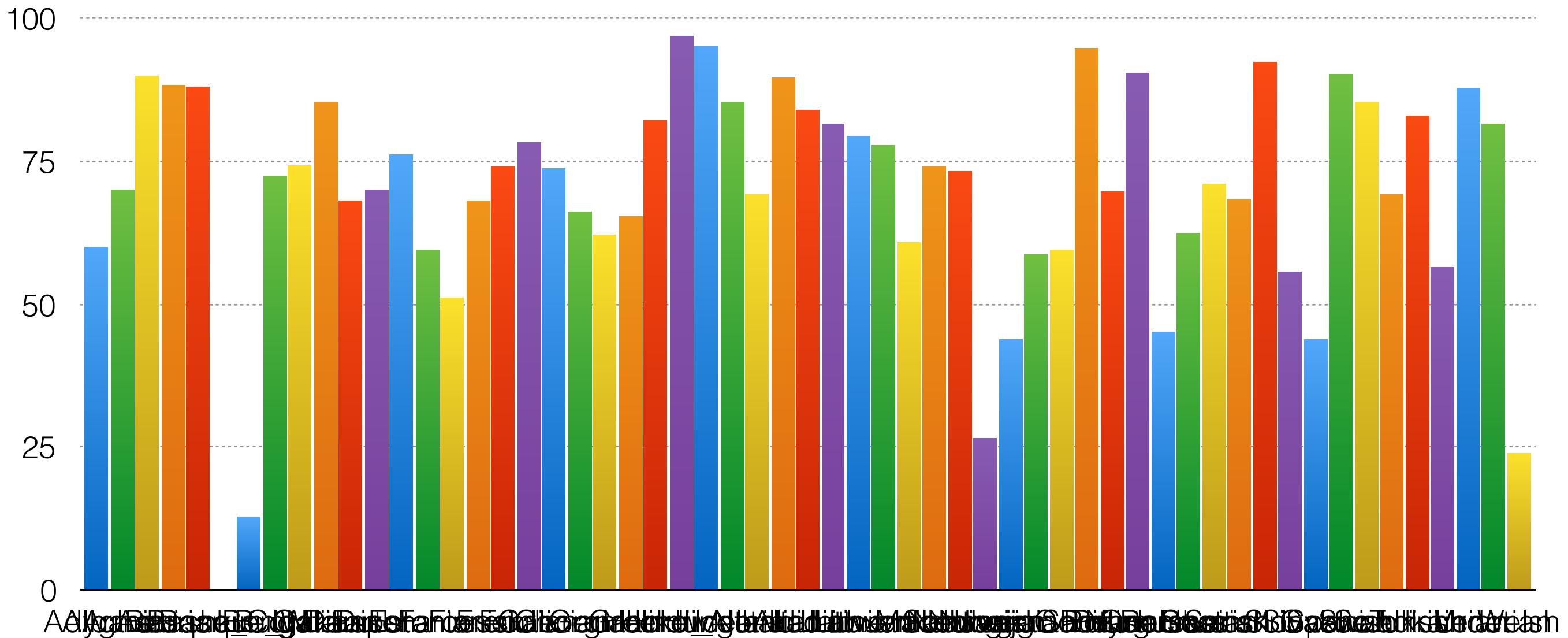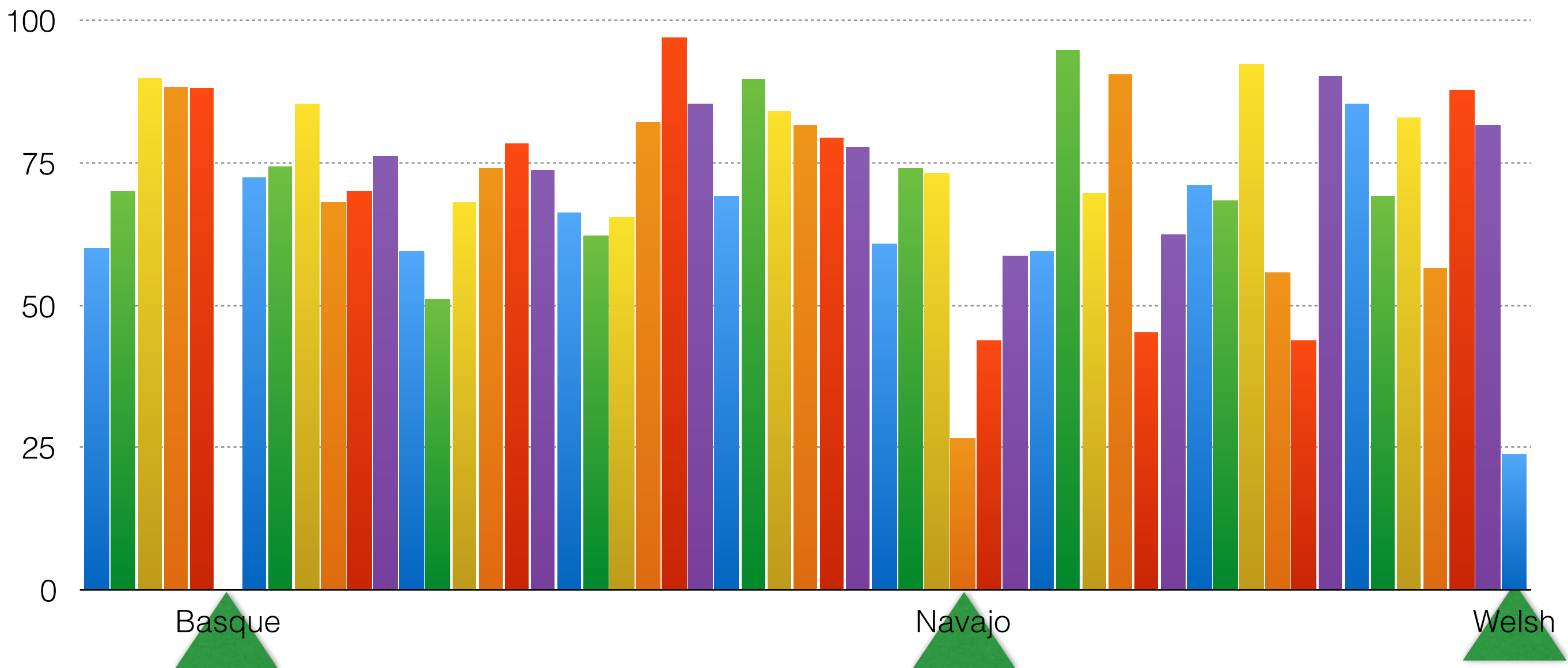
# Result Overview:
# Lemma+POS Recall > 95%

Hindi          Scottish Gaelic

# Result Overview: Lemma+Tags Recall



Evaluation of Finite State Morphological Analyzers

# Result Overview:
# Lemma+Tags Recall < 30%

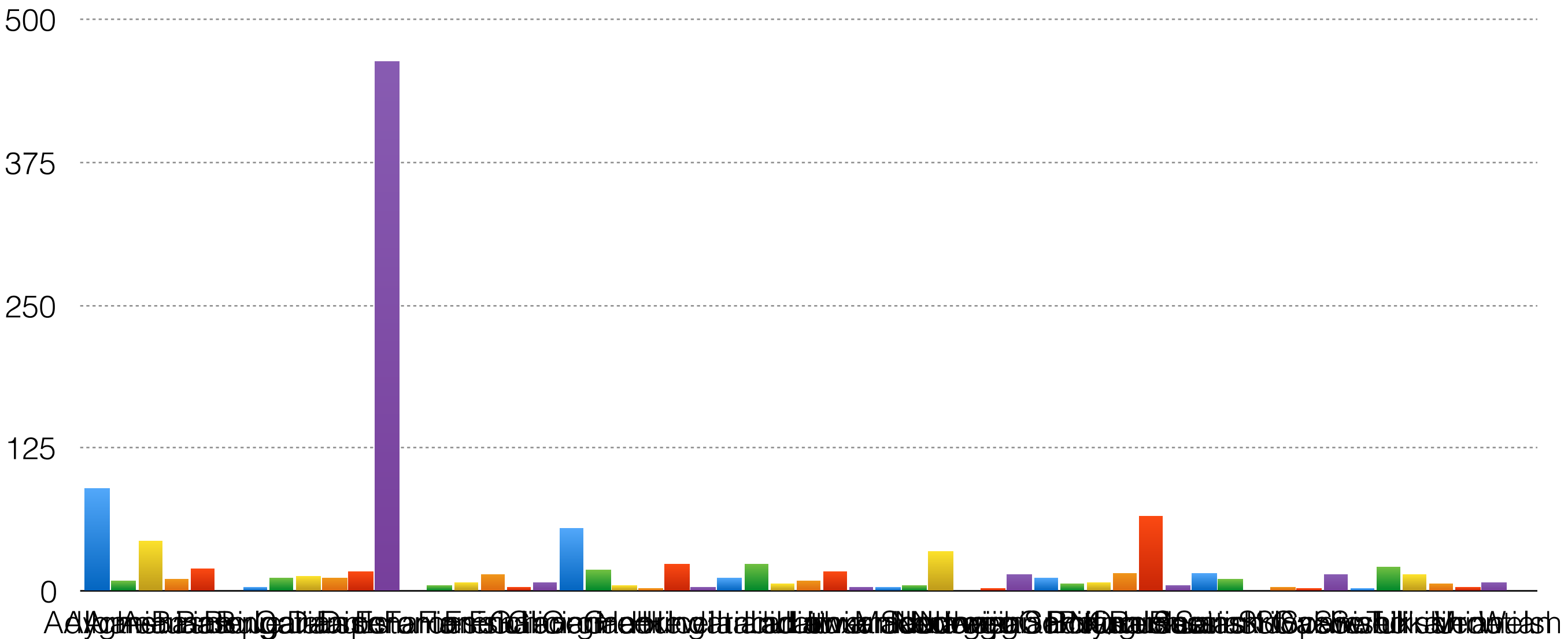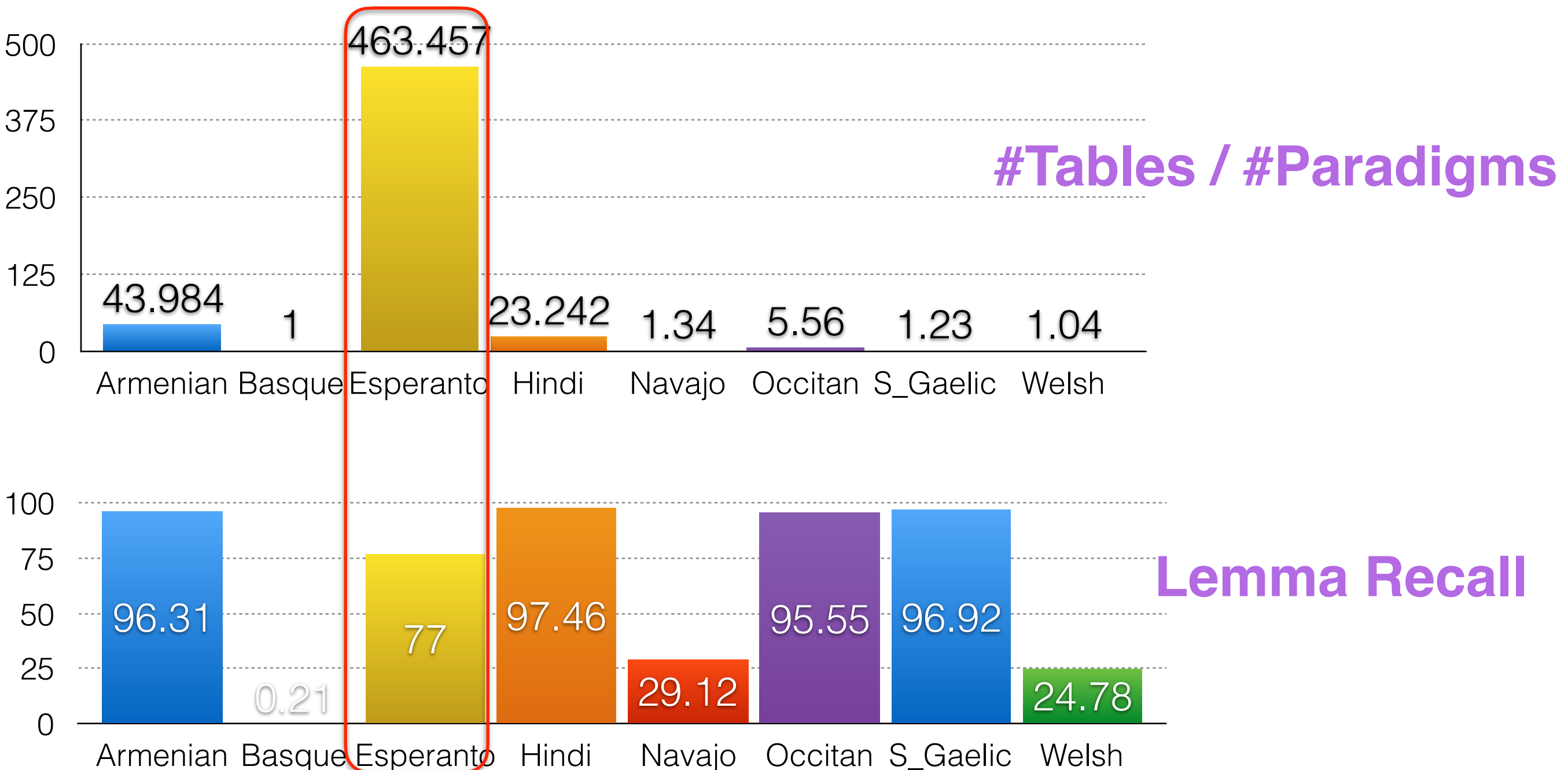# Result Overview: Lemma+Tags Recall > 95%

# Overview: #Tables/#Paradigms

# Result

‣ The results seem to correlate strongly with the amount and representativeness of available data.



**#Tables / #Paradigms**

Armenian: 43.984, Basque: 1, Esperanto: 463.457, Hindi: 23.242, Navajo: 1.34, Occitan: 5.56, S_Gaelic: 1.23, Welsh: 1.04

**Lemma Recall**

Armenian: 96.31, Basque: 0.21, Esperanto: 77, Hindi: 97.46, Navajo: 29.12, Occitan: 95.55, S_Gaelic: 96.92, Welsh: 24.78

# Result

‣ The results seem to correlate strongly with the amount and representativeness of available data.



**#Tables / #Paradigms**

**Lemma Recall**

# Result

‣ The results seem to correlate strongly with the amount and representativeness of available data.



**#Tables / #Paradigms**

**Lemma Recall**

# Result

‣ The results seem to correlate strongly with the amount and representativeness of available data.



**#Tables / #Paradigms**

| | Armenian | Basque | Esperanto | Hindi | Navajo | Occitan | S_Gaelic | Welsh |
|---|---|---|---|---|---|---|---|---|
| | 43.984 | 1 | 463.457 | 23.242 | 1.34 | 5.56 | 1.23 | 1.04 |

**Lemma Recall**

| | Armenian | Basque | Esperanto | Hindi | Navajo | Occitan | S_Gaelic | Welsh |
|---|---|---|---|---|---|---|---|---|
| | 96.31 | 0.21 | 77 | 97.46 | 29.12 | 95.55 | 96.92 | 24.78 |

# Wrap-up

‣ Simple method to construct weighted FST from labeled data

‣ Robust performance for inflectional morphology

‣ Large representative data is critical for the performance.

# Thank You