# Massively Multilingual Transfer for NER
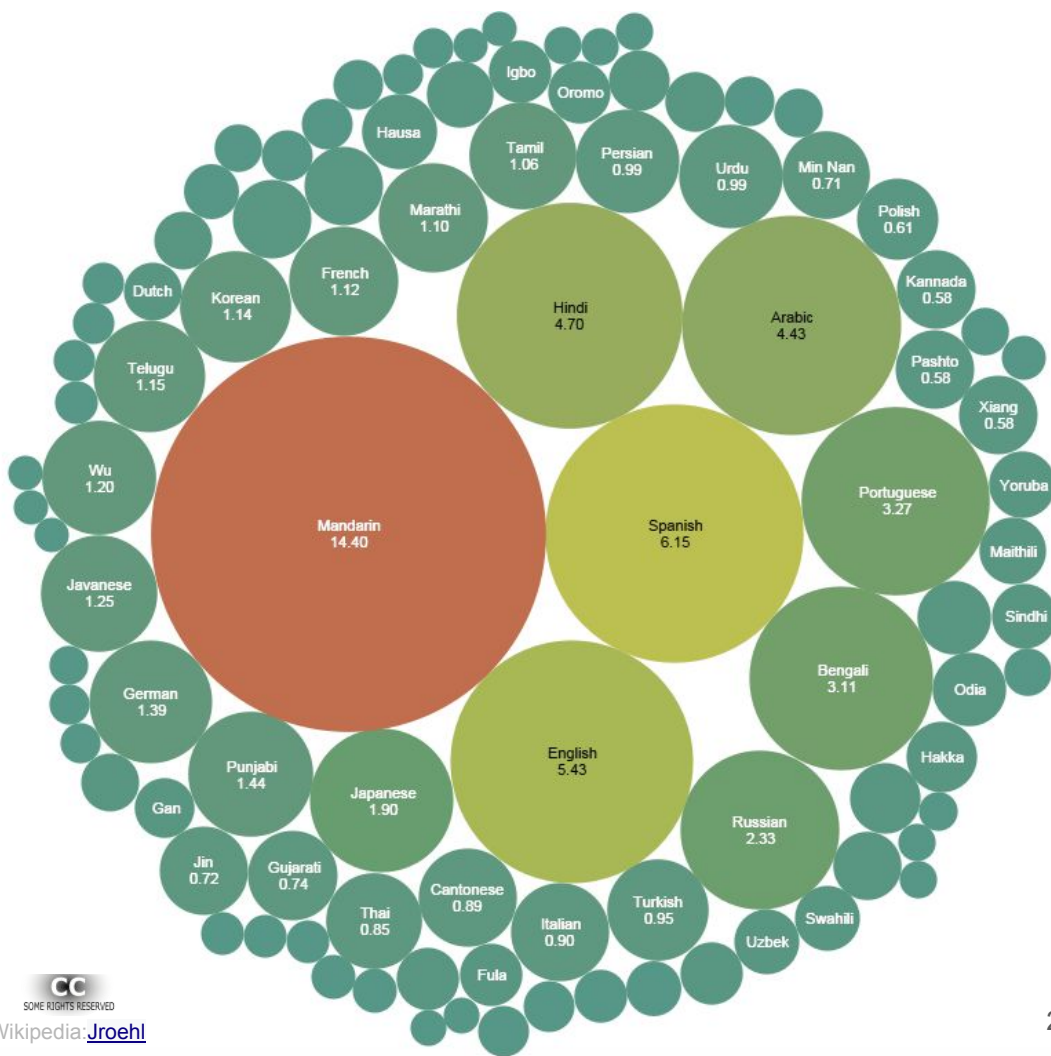
Afshin Rahimi, Yuan Li, and **Trevor Cohn**
University of Melbourne

Afshin Rahimi, Yuan Li, and **Trevor Cohn**
University of Melbourne

**6000+** languages

**≈ 1%** with annotation

Wikipedia:Jroehl

2

Emergency Response  Named Entity Recognition
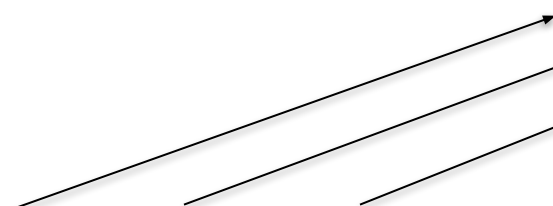
# Annotation Projection for Transfer

kailangan namin ng mas maraming dugo sa Pagasanjan .

**B-LOC**

......

we need more blood in Pgasanjan .
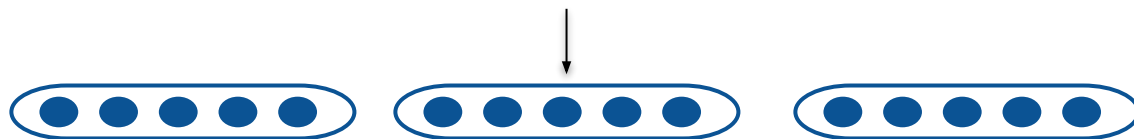
O   O   O      O   O   **B-LOC**   O

Yarowsky et al. (2001)

# Representation Projection for Transfer

kailangan namin ng mas maraming dugo sa Pagasanjan .

*language independent representation*
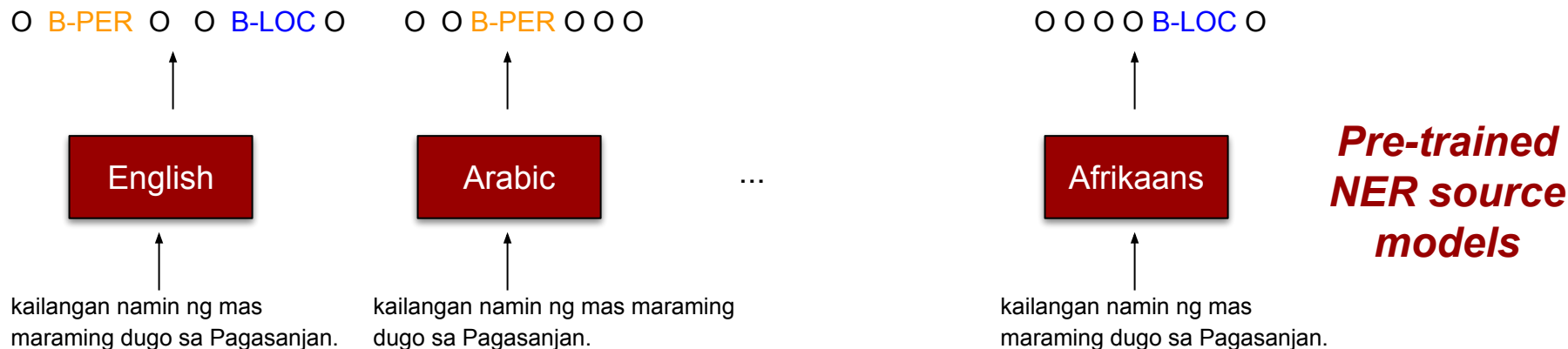
Cross-lingual word embeddings
(Lample et al., 2018)

**Mis-matched Model**

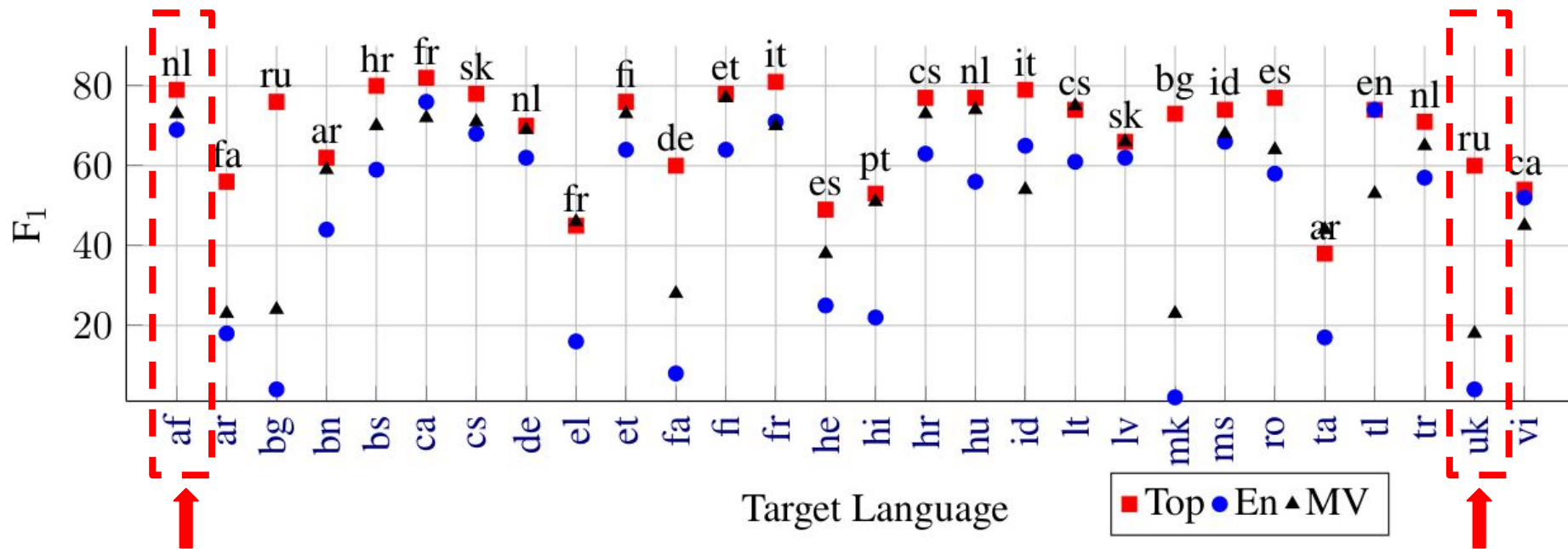**Ideal: source-target similar in word order, script, syntax**

O O O O O B-LOC O

# Direct Transfer for NER

Output: Labelled sentences in the target language

O  B-PER  O  O  B-LOC  O      O  O B-PER O O O                      O O O O B-LOC O

↑                            ↑                                          ↑

| English |        | Arabic |        ...        | Afrikaans |

↑                            ↑                                          ↑

kailangan namin ng mas        kailangan namin ng mas maraming        kailangan namin ng mas
maraming dugo sa Pagasanjan.  dugo sa Pagasanjan.                    maraming dugo sa Pagasanjan.
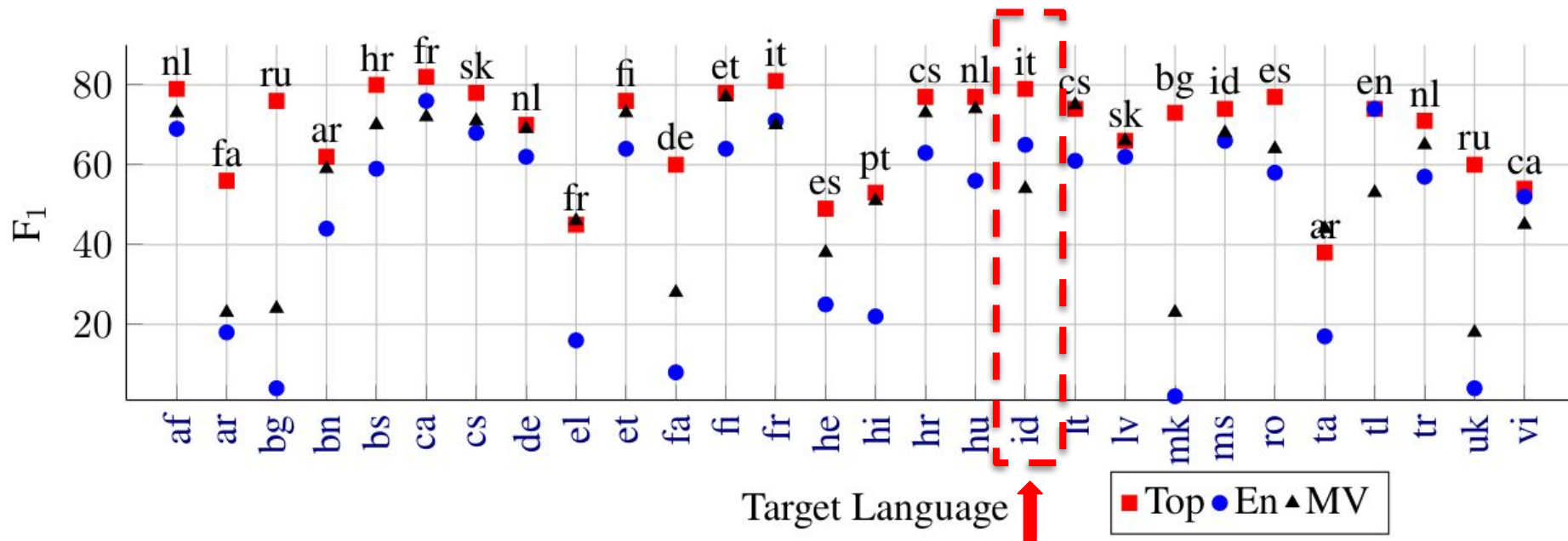
*Pre-trained NER source models*

Input: Unlabelled sentences in the target language encoded with cross-lingual embeddings

6

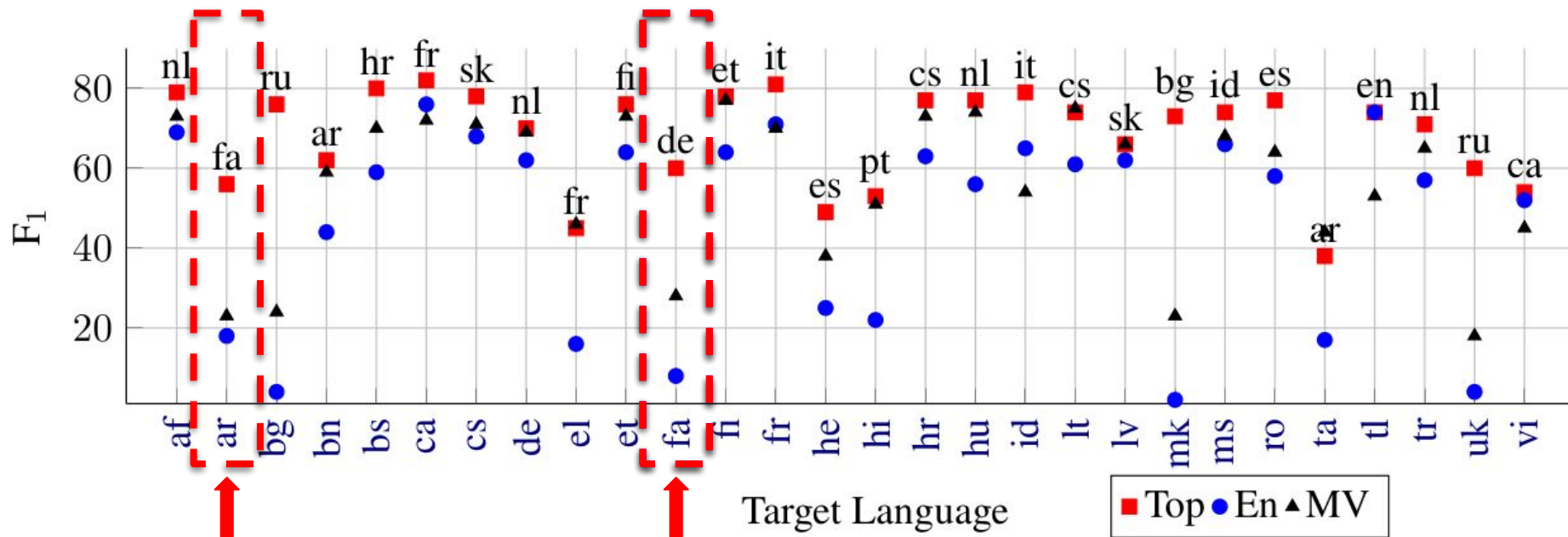# Direct Transfer Results (NER F1 score, WikiANN)



unsuprising

# Direct Transfer Results (NER F1 score, WikiANN)
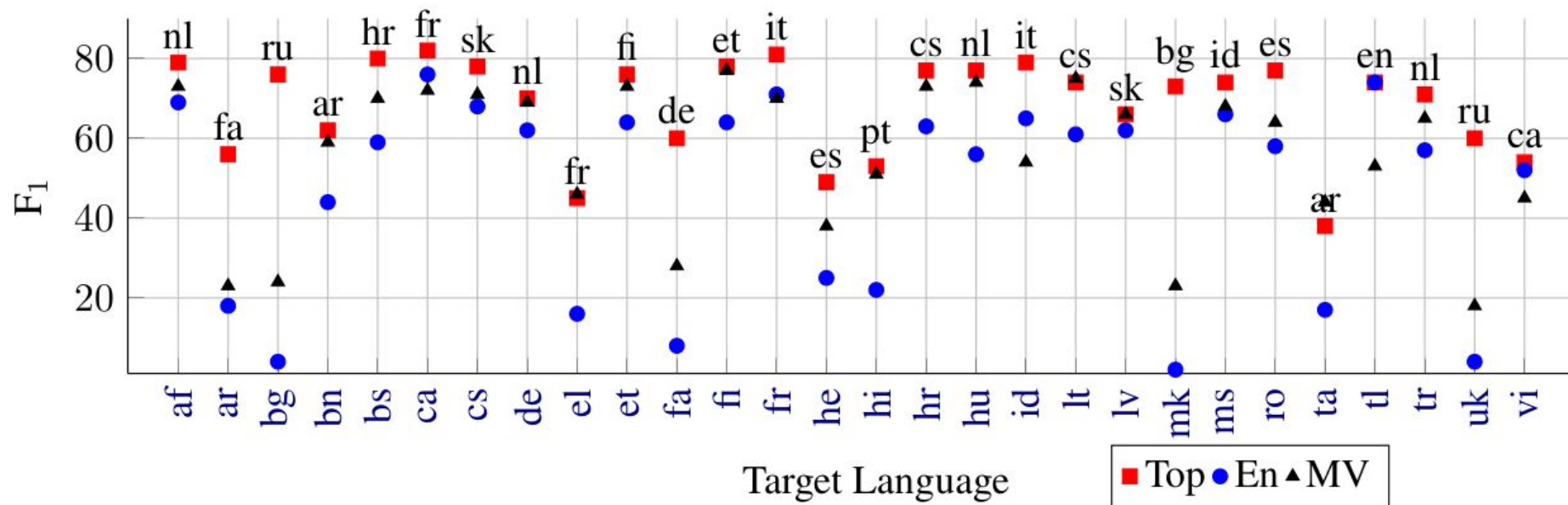


unrelated

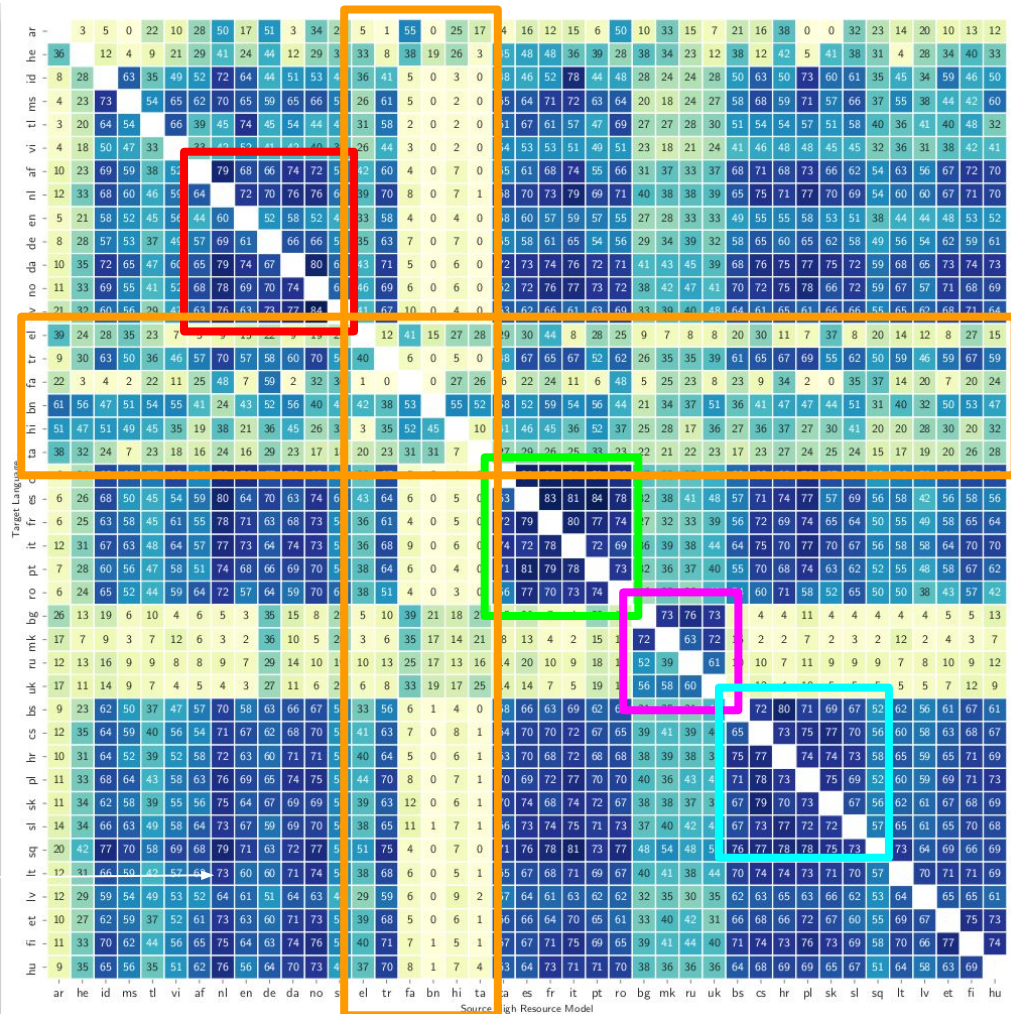# Direct Transfer Results (NER F1 score, WikiANN)



asymmetry

# Voting & English are often poor!

# General findings

- Transfer strongest within language family (**Germanic**, **Roman**, **Slavic-Cyr**, **Slavic-Latin**)
- Asymmetry between use as source vs target language (**Slavic-Cyr, Greek/Turkish/...**)
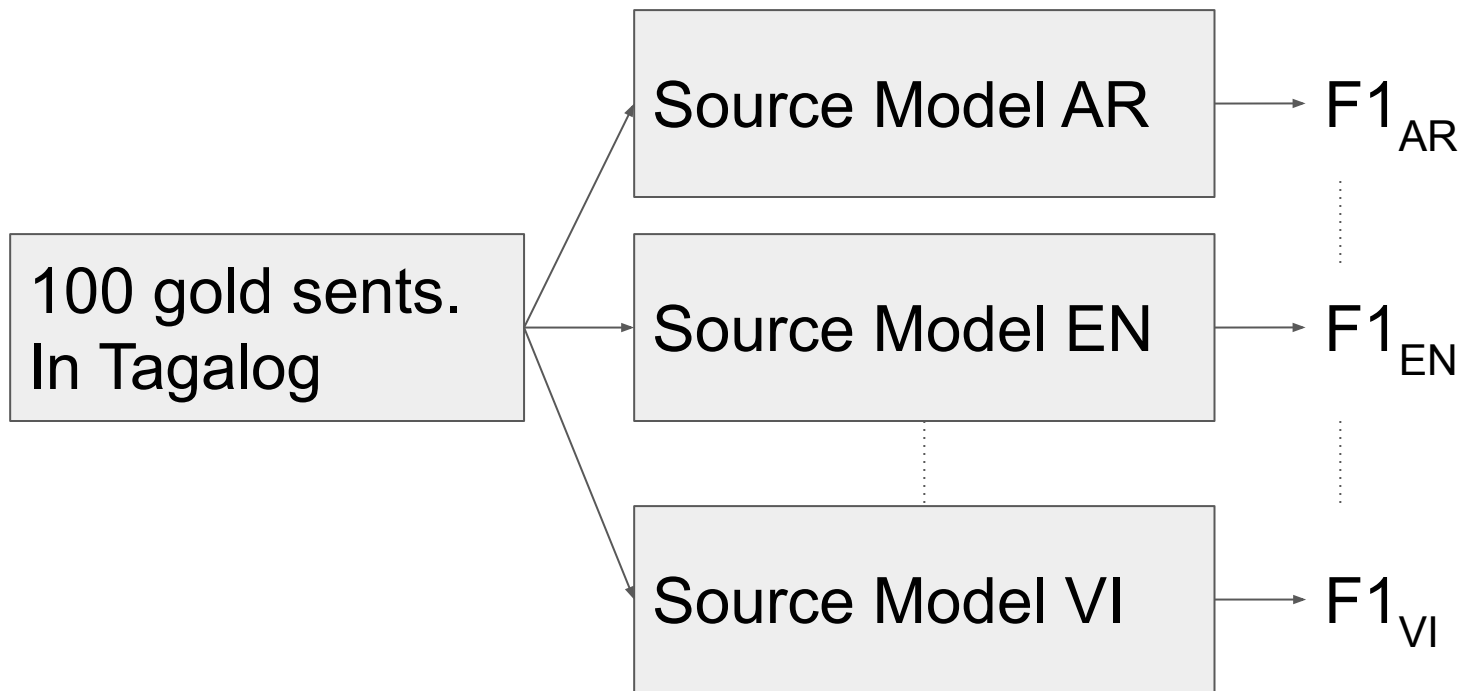- But lots of odd results & overall highly noisy

# Problem Statement

**Input:**

- N black-box source models
- Unlabelled data in target language
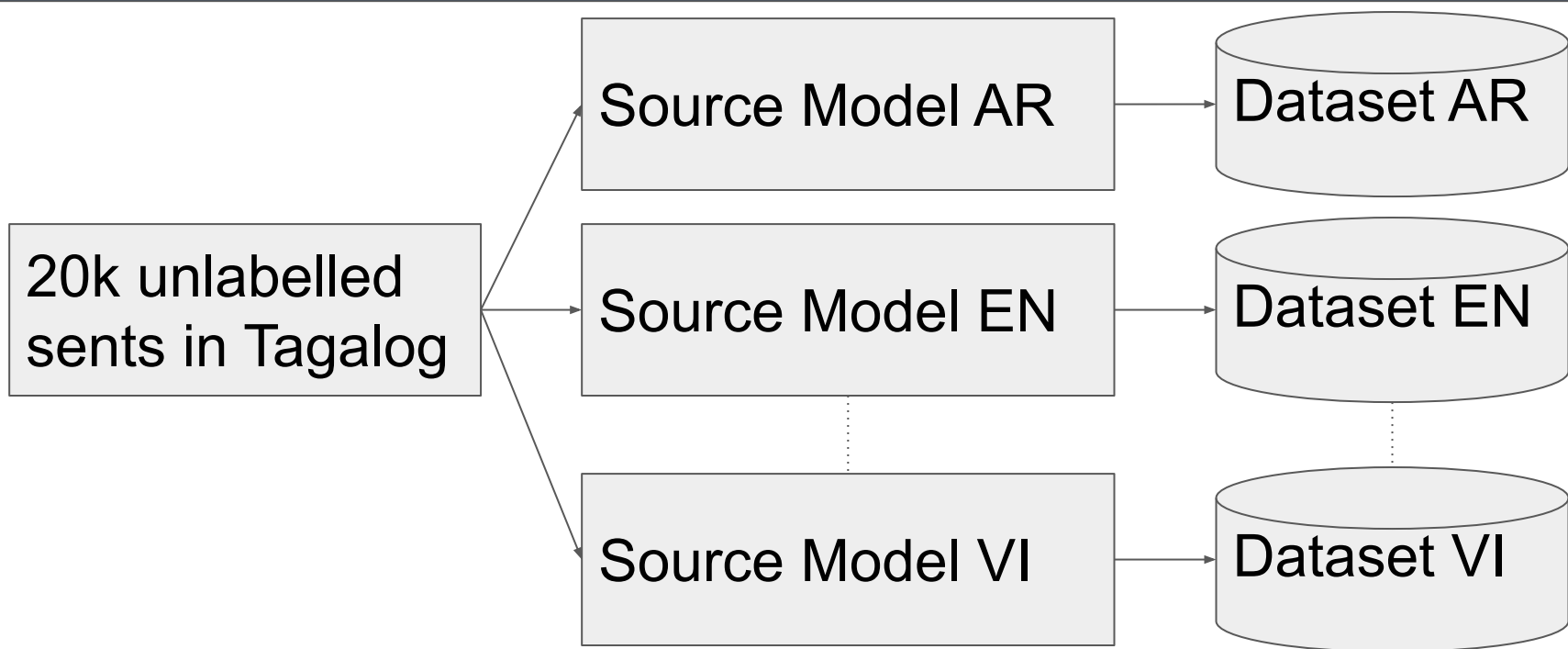- Little or no labelled data (few shot and zero shot)

**Output:**

- Good predictions in the target language

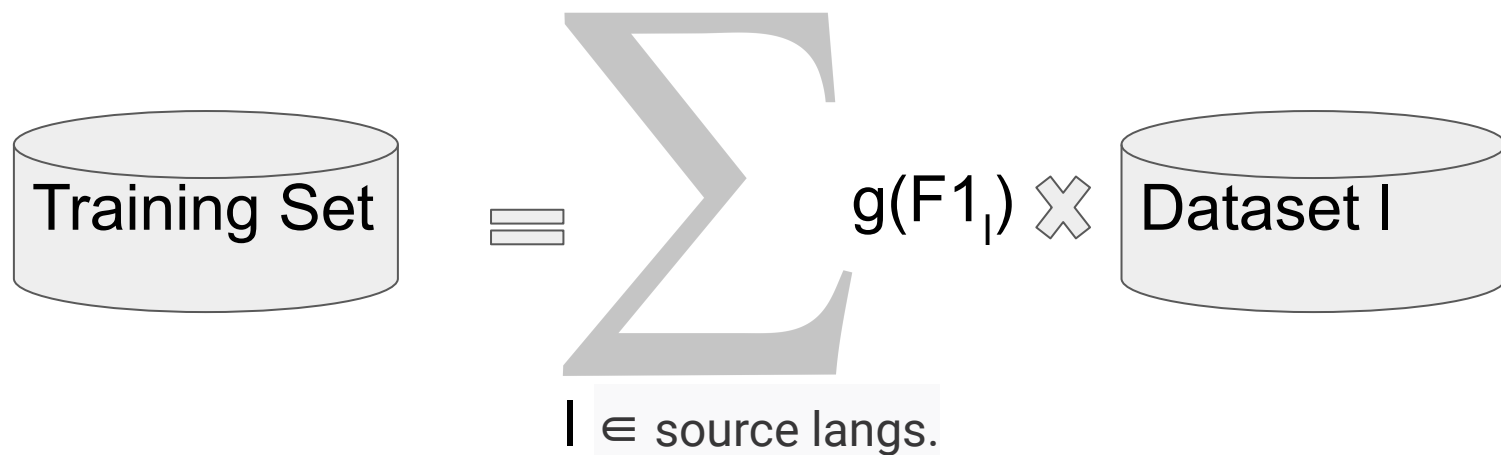# Model 1: Few Shot Ranking and Retraining (RaRe)



100 gold sents. In Tagalog → Source Model AR → $F1_{AR}$

100 gold sents. In Tagalog → Source Model EN → $F1_{EN}$

100 gold sents. In Tagalog → Source Model VI → $F1_{VI}$

**Source model qualities**

# Model 1: Few Shot Ranking and Retraining (RaRe)

20k unlabelled sents in Tagalog

Source Model AR → Dataset AR

Source Model EN → Dataset EN

Source Model VI → Dataset VI

**N training sets in Tagalog**

# Model 1: Few Shot Ranking and Retraining (RaRe)

$$\text{Training Set} = \sum_{l \, \in \, \text{source langs.}} g(F1_l) \times \text{Dataset } l$$

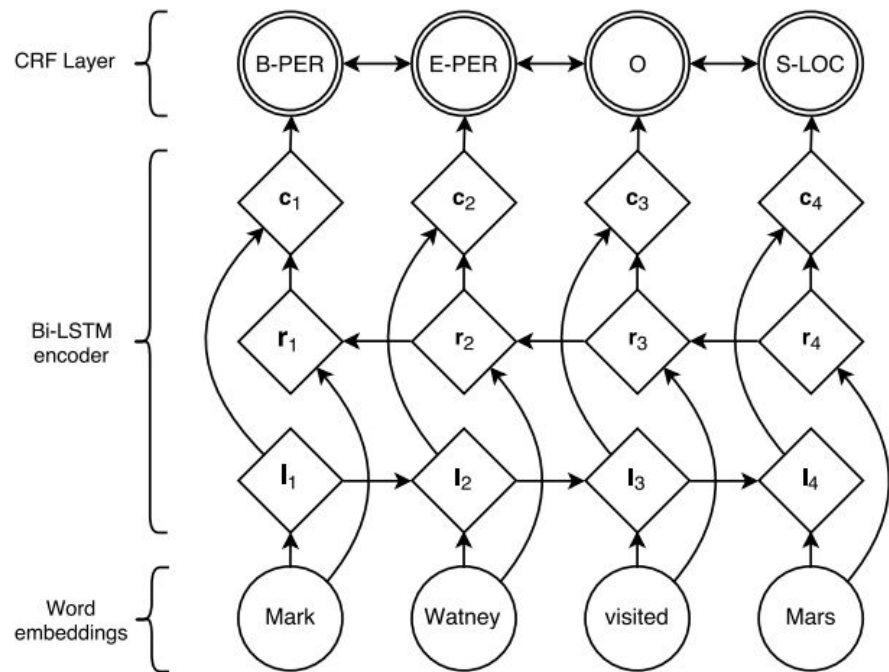**Final training set, a mixture of distilled knowledge**

# Model 1: Few Shot Ranking and Retraining (RaRe)

1. Train an NER model on the mixture datasets.
2. Fine-tune on 100 gold samples.

Zero-shot variant: uniform sampling without fine-tuning (**RaRe**$_{\textbf{uns}}$)

# Hierarchical BiLSTM-CRF as model



Lample et al., (2016)

Our method is **independent** of model choice.
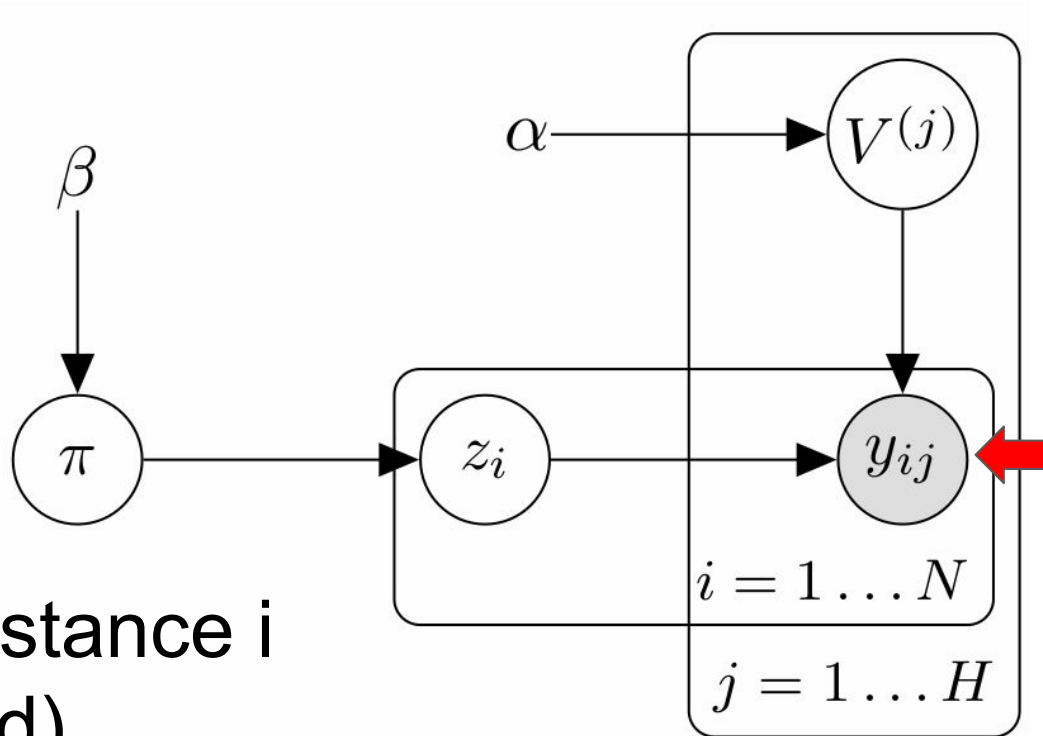
# Model 2: Zero Shot Transfer (BEA)

What if no gold labels are available?

1. Treat gold labels Z as hidden variables
2. Estimate Z that best explains all the observed predictions
3. Re-estimate the quality of source models

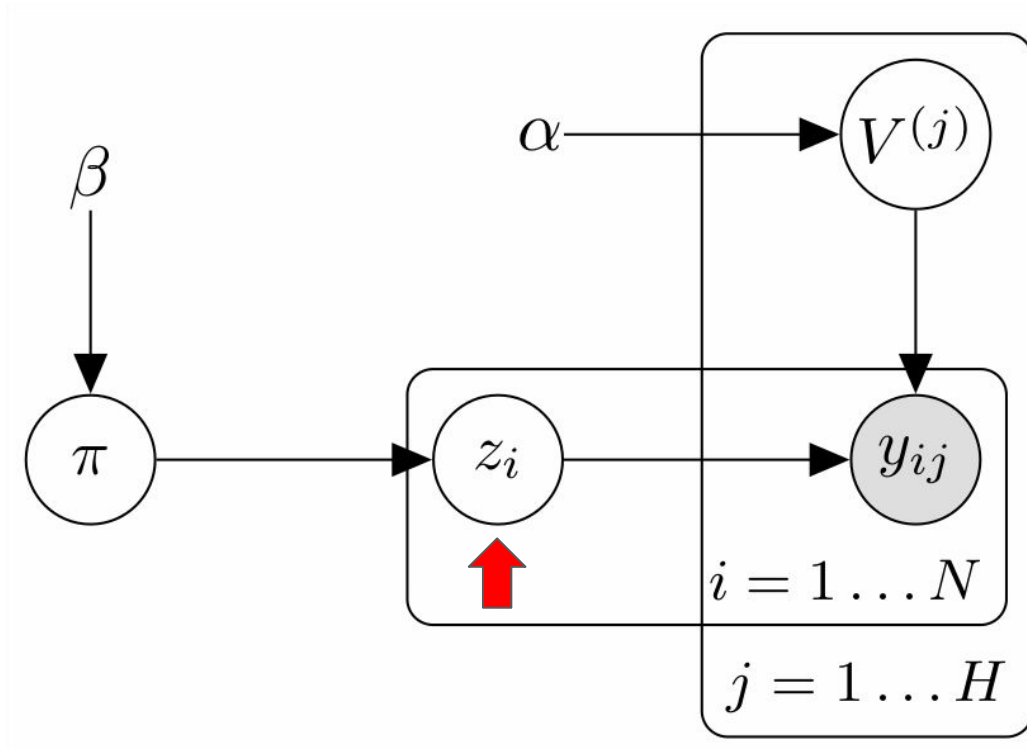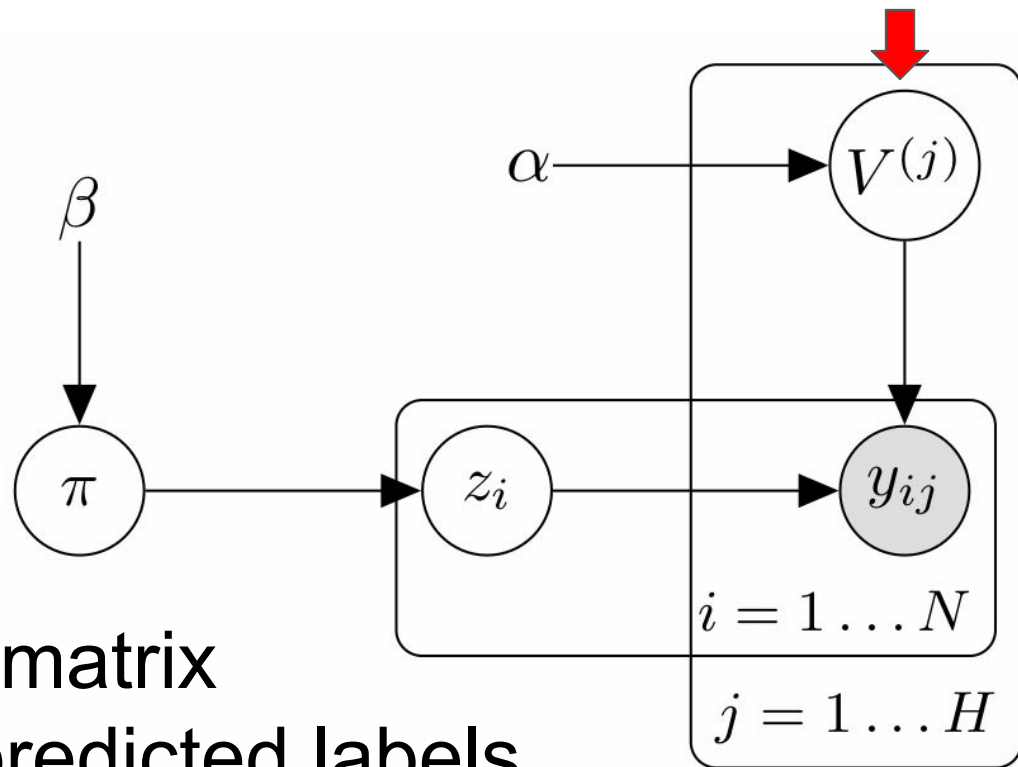Inspired by Kim and Ghahramani (2012)

# Model 2: Zero Shot Transfer (BEA)



Predicted label of instance i
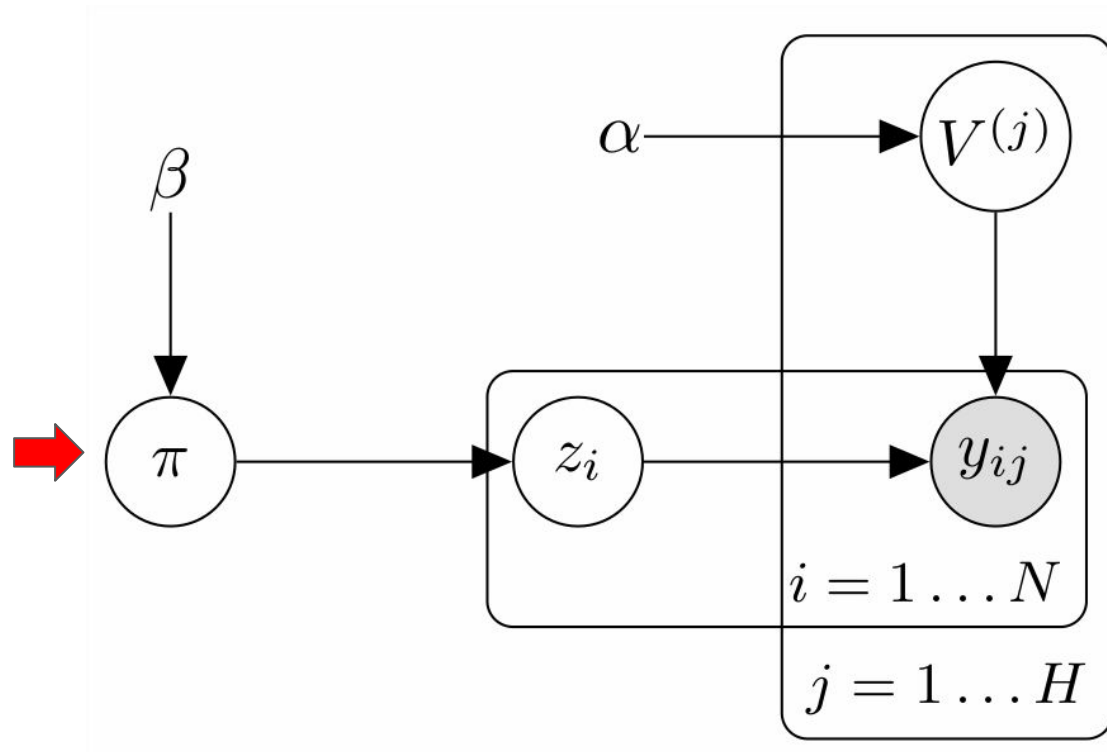by model j (observed)

# Model 2: Zero Shot Transfer (BEA)



True label of instance i

Model j's confusion matrix between True and predicted labels.
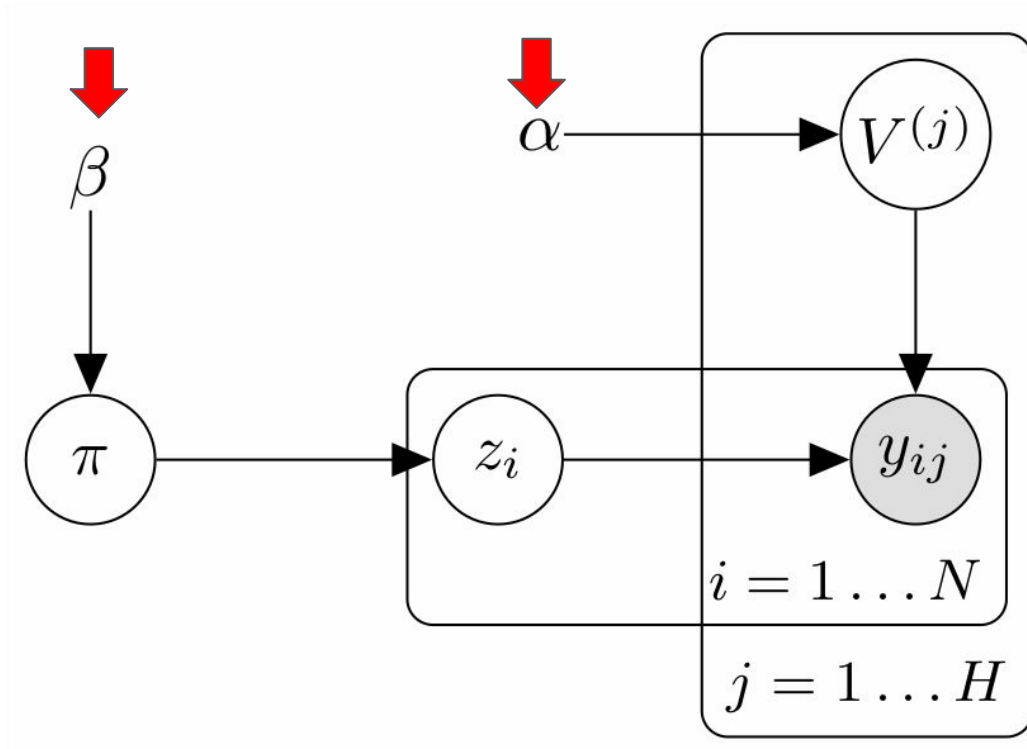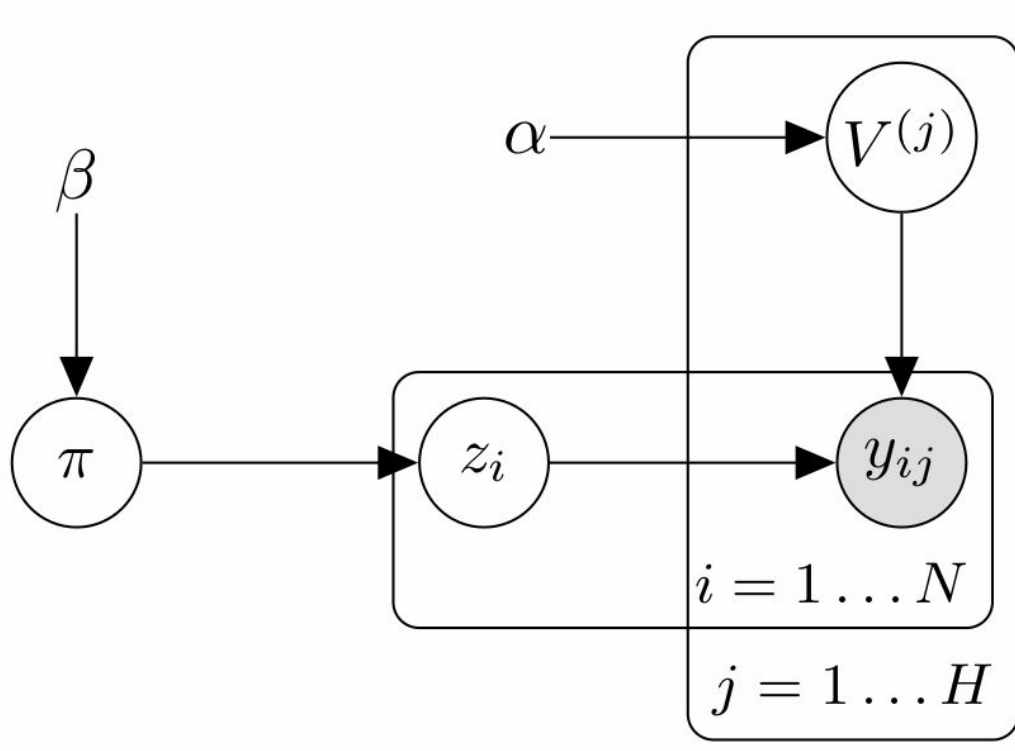
Categorical
Distribution

Uninformative
Dirichlet Priors

# Model 2: Zero Shot Transfer (BEA)

Find Z to maximises P(Z|Y,$\alpha$,$\beta$), using variational mean-field approx.

Warm-start with MV.

# Extensions to BEA

1. **Spammer removal:**
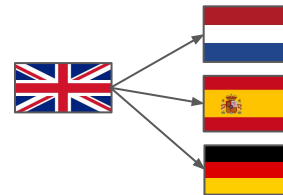After running BEA, estimate source model qualities and remove bottom k, run BEA again (**BEA$_{unsx2}$**)

2. **Few shot scenario:**
Given 100 gold sentences, estimate source model confusion matrices, then run BEA (**BEA$_{sup}$**)

3. **Token vs Entity application**

# Benchmark: BWET (Xie et al., 2018)

Single source annotation projection with bilingual dictionaries from cross-lingual word embeddings

- Transfer english training data to German, Dutch, and Spanish.

- Train a transformer NER on the projected training data.

State-of-the-art on zero-shot NER transfer (orthogonal to this)

# CoNLL Results (avg F1 over de, nl, es)



Täckström et al. (2012)
Nothman et al. (2013)
Tsai et al. (2016)
Ni et al. (2017)
Mayhew et al. (2017)

Use parallel data, dictionary or wikipedia

0          25          50          75          100

AVG F1 over de, nl and es

# CoNLL Results (avg F1 over de, nl, es)



Täckström et al. (2012)
Nothman et al. (2013)
Tsai et al. (2016)
Ni et al. (2017)
Mayhew et al. (2017)
Xie et al. (2018) → Zero shot

0        25        50        75        100

AVG F1 over de, nl and es

# CoNLL Results (avg F1 over de, nl, es)
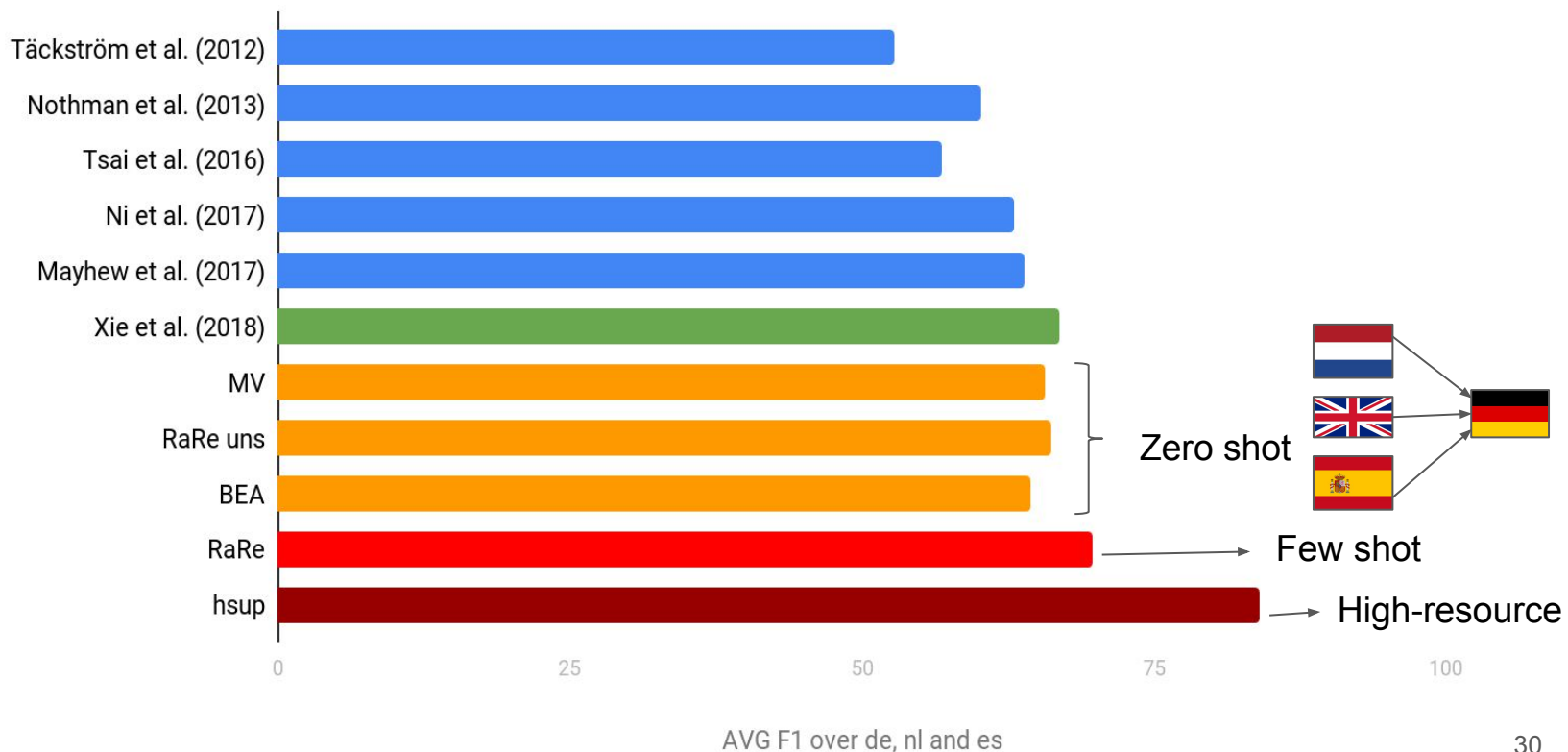

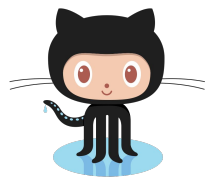
AVG F1 over de, nl and es

# CoNLL Results (avg F1 over de, nl, es)



AVG F1 over de, nl and es

# WIKIANN NER Datasets (Pan et al., 2017)

- Silver annotations from Wikipedia for **282** languages.
- We picked **41** languages based on availability of bilingual dictionaries.
- Created balanced training/dev/test partitions
  (varying size of training according to data availability)

github.com/afshinrahimi/mmner

# L.O.O. over 41 languages

**L.O.O. over 41 languages**

**Transfer from 40 source languages**

Tagalog

# L.O.O. over 41 languages

**L.O.O. over 41 languages**

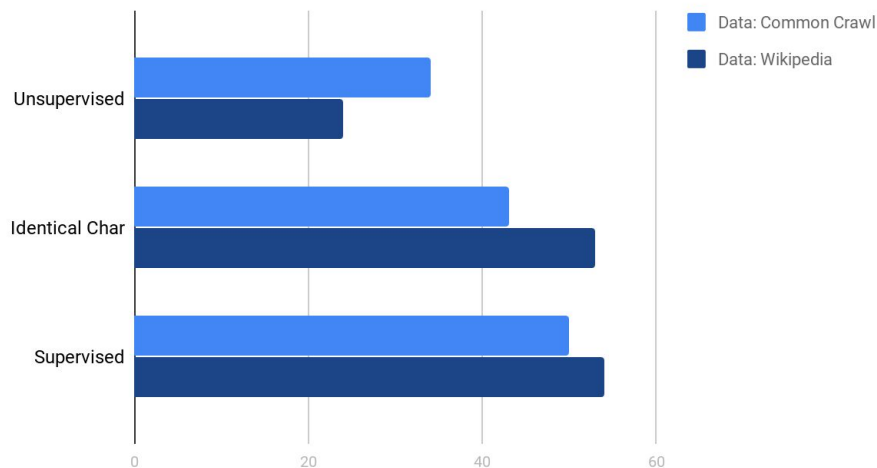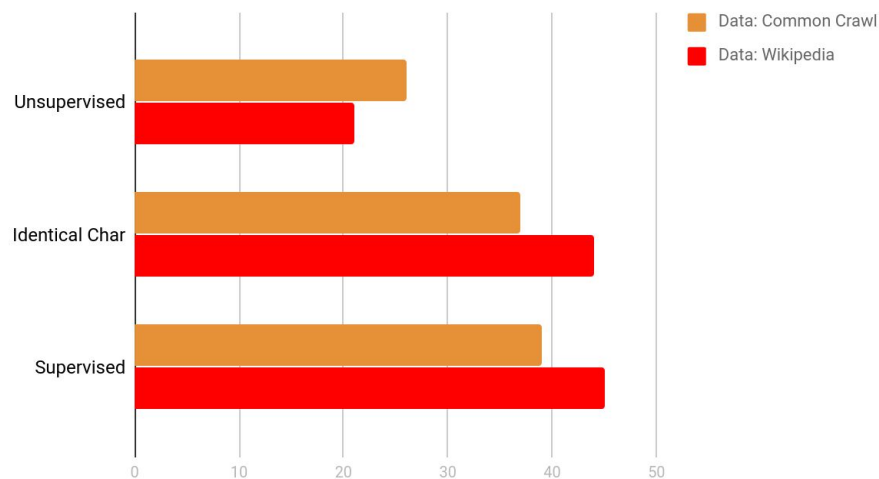**Transfer from 40 source languages**

Tamil

# Word representation: FastText/MUSE

Use **fasttext** monolingual **wiki** embeddings mapped to English space using **Identical Character Strings**.

Bilingual Induction Accuracy

■ Data: Common Crawl
■ Data: Wikipedia

Direct Transfer F1 averaged over 41 languages

■ Data: Common Crawl
■ Data: Wikipedia

Conneau et al. (2017)

# Results: WikiANN

Supervised: no transfer



LSup        Low-resource

HSup        High-resource

0    25    50    75    100

AVG F1 over 41 languages

# Results: WikiANN

Many low quality source models



MV — Zero shot

LSup — Low-resource

HSup — High-resource

0    25    50    75    100

AVG F1 over 41 languages

# Results: WikiANN

Single source (en)



AVG F1 over 41 languages

# Results: WikiANN

Bayesian ensembling



AVG F1 over 41 languages

# Results: WikiANN



+spammer removal

MV
BWET
BEA
BEA (spam.)

Zero shot

LSup — Low-resource
HSup — High-resource

0    25    50    75    100

AVG F1 over 41 languages

# Results: WikiANN



MV between top 3 sources

AVG F1 over 41 languages

# Results: WikiANN



Estimate BEA confusion & prior from annotations

AVG F1 over 41 languages

43

# Results: WikiANN



Ranking Retraining Method (using character info)

Chart axis labels (top to bottom): MV, BWET, BEA, BEA (spam.), MV (sup), BEA (sup), RaRe, LSup, HSup

Right side groupings: Zero shot, Few shot, Low-resource, High-resource

X-axis: 0, 25, 50, 75, 100

AVG F1 over 41 languages

44

# Effect of increasing #source languages

Methods robust to many varying quality source languages.

Even better with few-shot supervision.

**Transfer from <span style="color:red">multiple source languages</span> helps because for many languages we don't know the best source language.**

**takeaway** / noun [uk/aus/nz]: *a meal cooked and bought at a shop or restaurant but taken somewhere else...*
Cambridge English Dictionary

46

With multiple source languages, you need to <span style="color:red">estimate their qualities</span> because uniform voting doesn't perform well.

**takeaway** / noun [uk/aus/nz]: *a meal cooked and bought at a shop or restaurant but taken somewhere else...*
Cambridge English Dictionary

47

A **small training set** in target language helps, and can be done cheaply and quickly (Garrette and Baldridge, 2013).

# Thank you!

Datasets & code
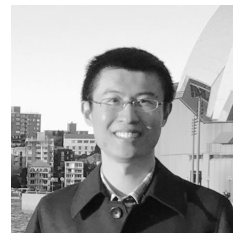
github.com/afshinrahimi/mmner

# Future Work

- Map all scripts to IPA or Roman alphabet
  (good for shared embeddings and character-level transfer)
    - uroman: Hermjakob et al. (2018)
    - epitran:   Mortensen et al. (2018)
- Can we estimate the quality of source models/languages
  for a specific target language based on language
  characteristics (Littell et al., 2017)?
- Technique should apply beyond NER to other tasks.