

Informal Romanization across Languages and Scripts

Maria Ryskina



Carnegie Mellon University
Language
Technologies
Institute

Carnegie Mellon University
Language Technologies Institute

Non-standard language & NLP

From standardized language...

```
(S (NP-SBJ (NN Compound)
     (NNS yields))
  (VP (VBP assume)
    (UCP (NP (NP (NN reinvestment))
              (PP (IN of)
                  (NP (NNS dividends))))))
  (CC and)
  (SBAR (IN that)
    (S (NP-SBJ (DT the)
          (JJ current)
          (NN yield))
     (VP (VBZ continues)
       (PP-TMP (IN for)
         (NP (DT a)
             (NN year)))))))
```

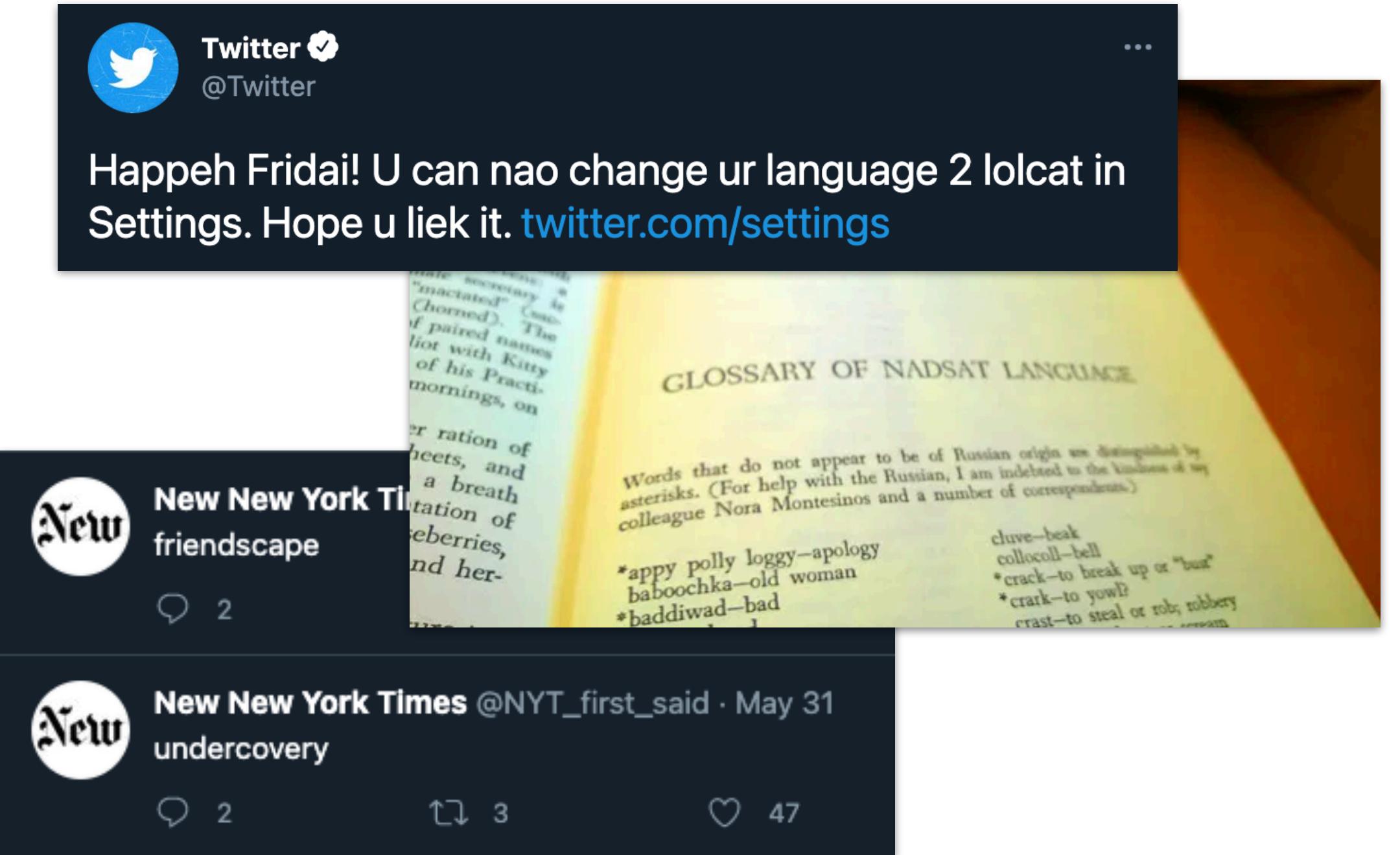
(. .)

)

e.g. PTB:

- Newswire
- Finance-related
- Formal

...To creative language



- Variety of genres
- Variety of domains
- Spectrum of formality

Linguistic innovation: What

- Non-standard, novel linguistic items...
 - Lexical: new word forms (*brony*)
 - Morphological: new morphemes (-gate) or derivatives (*prolifeness*)
 - Orthographic: non-standard spellings (*2nite*)
- ... before they become attested (*tweet*)
- People can infer their meaning, but NLP systems largely treat them as noise

Linguistic innovation: How

- **Q1: How do people process non-standard items?**
 - Shared knowledge or perception: $2 = \text{'two'} = /tu/$
 - Compositionality: $2nite = \text{'two'} + \text{'nite'} = /tu/ + /naɪt/ \approx /tənaɪt/$
 $\text{antivehicleness} = \text{'anti'} + \text{'vehicle'} + \text{'ness'}$
- **Q2: How can we get our NLP systems to that level?**
 - Text normalization: $2nite \rightarrow tonight$ (Baldwin et al., 2015: W-NUT shared task)
 - Improving robustness to noise & ‘noise’ (Li et al., 2019: WMT shared task)
 - Maybe we can encode creative reasoning into them?

Linguistic innovation: Why

- Creativity: expressing extralinguistic information
- Necessity: bypassing constraints or lack of tools

БАСК
ИИ
THE
Ц22Я



Привет,
как дела?
Privet,
kak dela?

Informal romanization

- *Romanization*: rendering non-Latin-script languages in Latin alphabet
- *Informal*: used online, arises out of Unicode/keyboard issues

Russian	человек	<i>chelovek, 4elovek, ceJloBek, ...</i>
Arabic	صباح	<i>saba7, sba7, sabah, ...</i>
Greek	ξένος	<i>xenos, ksenos, 3enos, ...</i>

Informal romanization

- Idiosyncratic representation: character substitutions up to the user
 - Substitution choices can convey social meaning (Nguen et al., 2021)

Russian	человек	<i>chelovek, 4elovek, ceJloBek, ...</i>
Arabic	صباح	<i>saba7, sba7, sabah, ...</i>
Greek	ξένος	<i>xenos, ksenos, 3enos, ...</i>

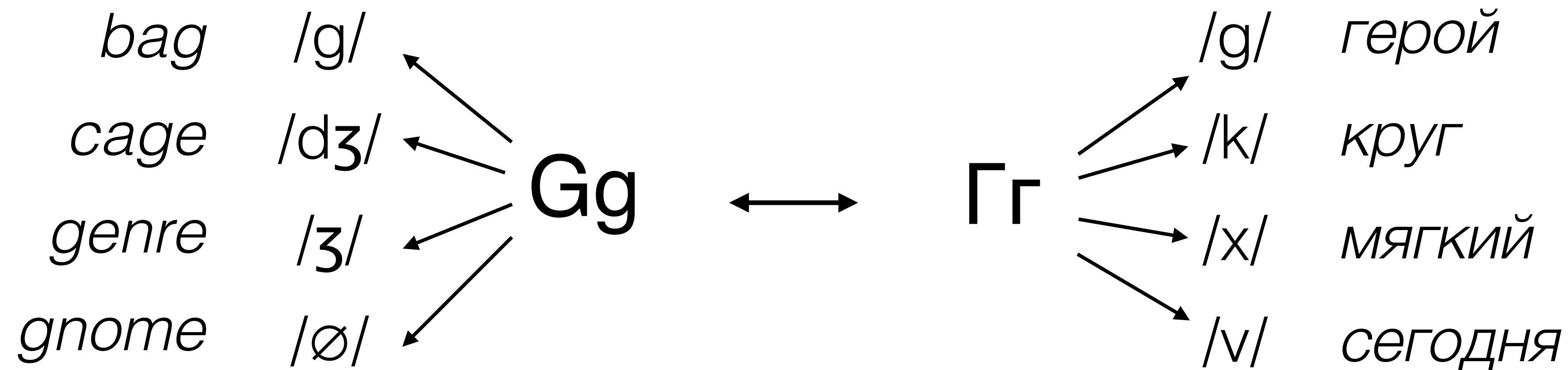
Informal romanization

- Idiosyncratic representation: character substitutions up to the user
 - Substitution choices can convey social meaning (Nguen et al., 2021)
- Most substitutions are based on **phonetic** or **visual** similarity
[but not all: ъ→', θ→u (Chalamandaris et al., 2006)]

Russian	человек	<i>chelovek, 4elovek, ceJloBek, ...</i>
Arabic	صباح	<i>saba7, sba7, sabah, ...</i>
Greek	ξένος	<i>xenos, ksenos, 3enos, ...</i>

Phonetic romanization

- What does it mean for two characters to be phonetically similar?



- This is just in one language each!

Phonetic romanization

- What does it mean for two characters to be phonetically similar?
- Out-of-context grapheme-phoneme association: $\Gamma \sim /g/ \rightarrow g$



Every letter makes a sound:
'A' says /eɪ/!*

*and /a/

Phonetic romanization

- What does it mean for two characters to be phonetically similar?
- Out-of-context grapheme-phoneme association: ر~/g/→g
- Phoneme produced in context: انتي /enti/→enty, صباح /sabaħ/→saba7

Visual romanization

- Broad similarity between glyph shapes $a\sim/a/\rightarrow a, \Gamma\sim/g/\rightarrow r$
- Single characters can map to bi-/trigraphs $\acute{y}\rightarrow bl, \dot{x}\rightarrow }\|{$
- Can be conditioned on a transformation $\mathcal{E}\rightarrow 3, \mathcal{L}\rightarrow v$
- Can be applied to a part of a glyph $\acute{i}\rightarrow 2$

Character alignment

- Monotonic alignment that depends on the writing system of the language

Alphabet

хорошо

|||||

xorosho

~ one-to-one

Abjad
(consonantal)

كريم

krym

/|\\|

kareem

~ one-to-one + null

Abugida
(alphasyllabary)

బెలగితు

/\\|\\|

belagitu

~ one-to-many

Task framing

- Convert romanized text to the conventional orthography of the language

Russian

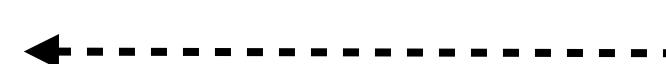
конгресс не одобрил бюджет



kongress ne odobril biudjet

Egyptian
Arabic

انا حأعدى عليك بكرة على 8 كده



ana h3dyy 3lek bokra 3la 8 kda

Kannada

ಮನ ಬೆಳಗಿತು



mana belagitu

Outline

- **Part 2:** Unsupervised WFST (Ryskina et al., ACL 2020)
 - Proposed model and parameterization
 - Similarity-based inductive bias
- **Part 3:** Unsupervised seq2seq, error analysis (Ryskina et al., SIGMORPHON 2021)
 - Seq2seq model and combining it with the WFST
 - Data and experimental results
 - Comparative error analysis

Task framing

- Convert romanized text to the conventional orthography of the language

Russian

конгресс не одобрил бюджет



kongress ne odobril biudjet

Egyptian
Arabic

انا حأعدى عليك بكرة على 8 كده



ana h3dyy 3lek bokra 3la 8 kda

Kannada

ಮನ ಬೆಳಗಿತು



mana belagitu

latent
(what they meant)

observed
(what they typed)

Task framing

- Parallel data does not occur naturally ⇒ **unsupervised** learning
- Perceptions of similarity are shared across users and even languages!
 - But does it mean we can train **language-independent** models? 

**Linguistically Naïve != Language Independent:
Why NLP Needs Linguistic Typology**

Emily M. Bender
University of Washington
Seattle, WA, USA
`ebender@u.washington.edu`

M Ryskina, MR Gormley, T Berg-Kirkpatrick. Phonetic and Visual Priors for Decipherment of Informal Romanization. ACL 2020.

Task framing

- Parallel data does not occur naturally ⇒ **unsupervised** learning
- Perceptions of similarity are shared across users and even languages!
 - But does it mean we can train **language-independent** models? 
- **Hypothesis:** **inductive bias** encoding these similarity notions provides signal that can somewhat **approximate human supervision**
 - We rely on **manually-curated resources** to operationalize it

M Ryskina, MR Gormley, T Berg-Kirkpatrick. Phonetic and Visual Priors for Decipherment of Informal Romanization. ACL 2020.

Decipherment

- Can be viewed as a decipherment task (Knight et al., 2006)

...When I look at an article in Russian, I say: “This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.”

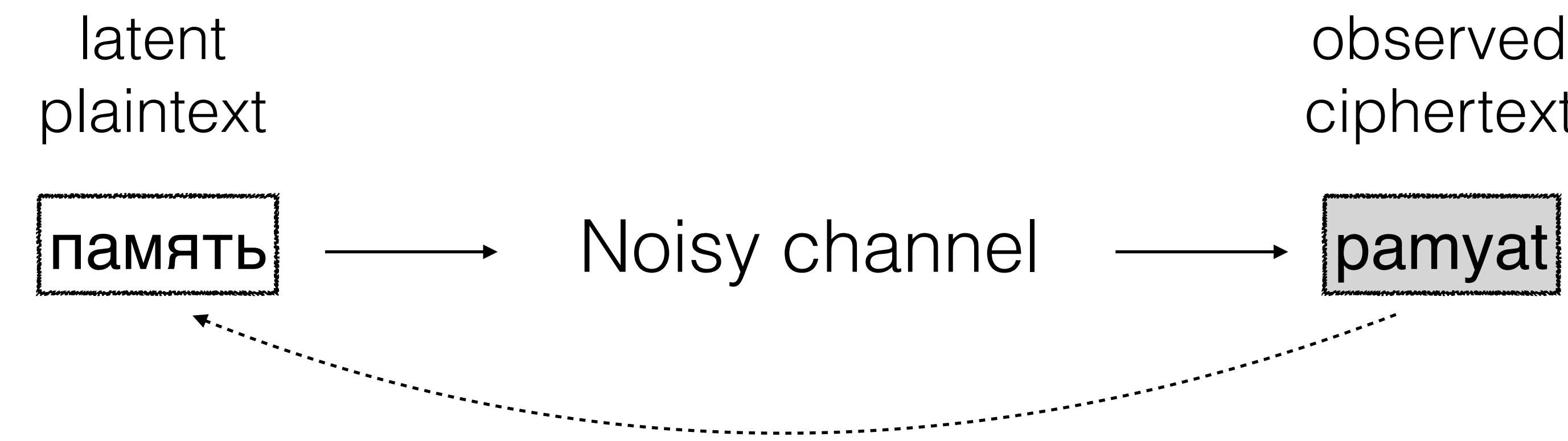
— Warren Weaver, 1947

Decipherment

- Can be viewed as a decipherment task (Knight et al., 2006)

...When I look at an article in **romanized Russian**, I say: “This is really written in **Cyrillic**, but it has been coded in some strange symbols. I will now proceed to decode.”

— Warren Weaver, 1947



Noisy-channel model

latent $n = \text{п а м я т ъ}$

observed $r = \text{p а m y a t}$

$$p(r) = \sum p(n; \gamma) \cdot p(r|n; \theta) \cdot p_{\text{prior}}(\theta; \alpha)$$

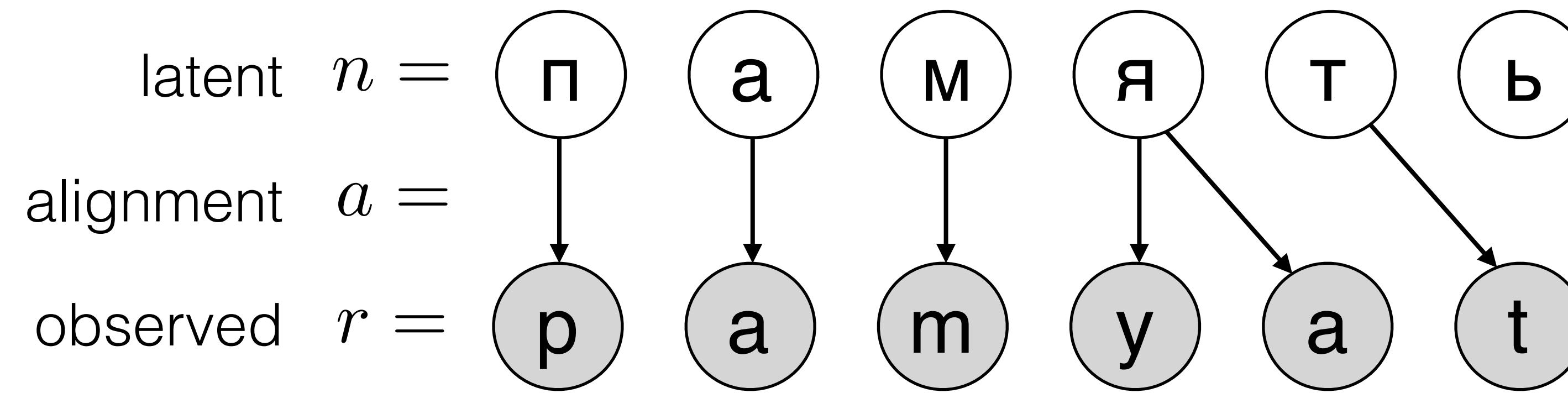
all possible
native script
sequences

n
transition probabilities

emission probabilities

θ
prior on parameters

Noisy-channel model



$$p(r) = \sum_{n,a} p(n; \gamma) \cdot p(r|n, a; \theta) \cdot p_{\text{prior}}(\theta; \alpha)$$

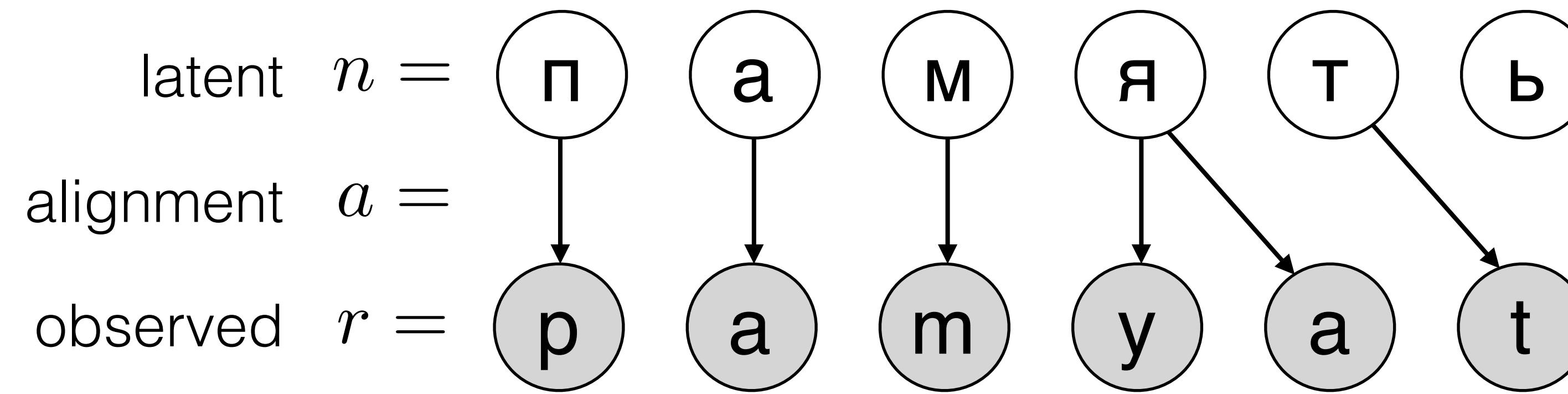
all possible
native script
sequences and
alignments

n, a
transition probabilities

/
emission probabilities

\
prior on parameters

Noisy-channel model



$$p(r) = \sum_{n,a} p(n; \gamma) \cdot p(r|n, a; \theta) \cdot p_{\text{prior}}(\theta; \alpha)$$

/ |
transition probabilities emission probabilities
prior on parameters

Phonetic bias

- Phonetic priors: mappings off **phonetic keyboard layouts**



https://en.wikipedia.org/wiki/Phonetic_keyboard_layout,

24 <https://arabic.omaralzabir.com/>, <http://kaulonline.com/uninagari/kannada/>

Phonetic bias

- Phonetic priors: mappings off **phonetic keyboard layouts**
 - One-to-one mapping constraints lead to spurious mappings



https://en.wikipedia.org/wiki/Phonetic_keyboard_layout,

25 <https://arabic.omaralzabir.com/>, <http://kaulonline.com/uninagari/kannada/>

Visual bias

- Visual priors: mappings off the **Unicode confusables list**
- Designed to combat spoofing attacks

y	ȸ	Y	Ȳ	y	ȳ	Ȳ	ȶ
0079 LATIN SMALL LETTER Y	0263 LATIN SMALL LETTER GAMMA	028F LATIN LETTER SMALL CAPITAL Y	03B3 GREEK SMALL LETTER GAMMA	0443 CYRILLIC SMALL LETTER U	04AF CYRILLIC SMALL LETTER STRAIGHT U	10E7 GEORGIAN LETTER QAR	
p	ρ	϶	پ	ϙ	ϙ	ϙ	پ
0070 LATIN SMALL LETTER P	03C1 GREEK SMALL LETTER RHO	03F1 GREEK RHO SYMBOL	0440 CYRILLIC SMALL LETTER ER	2374 APL FUNCTIONAL SYMBOL RHO	2CA3 COPTIC SMALL LETTER RO	1D429 MATHEMATICAL BOLD SMALL P	

sigtyp.io

sigtyp.io

Visual bias

- Visual priors: mappings off the **Unicode confusables list**
- Designed to combat spoofing attacks



y	γ	Y	Ƴ	y	Y	y
0079 LATIN SMALL LETTER Y	0263 LATIN SMALL LETTER GAMMA	028F LATIN LETTER SMALL CAPITAL Y	03B3 GREEK SMALL LETTER GAMMA	0443 CYRILLIC SMALL LETTER U	04AF CYRILLIC SMALL LETTER STRAIGHT U	10E7 GEORGIAN LETTER QAR
p	ρ	ε	پ	ρ	P	p
0070 LATIN SMALL LETTER P	03C1 GREEK SMALL LETTER RHO	03F1 GREEK RHO SYMBOL	0440 CYRILLIC SMALL LETTER ER	2374 APL FUNCTIONAL SYMBOL RHO	2CA3 COPTIC SMALL LETTER RO	1D429 MATHEMATICAL BOLD SMALL P

sigtyp.io

The site you just tried to visit looks fake. Attackers sometimes mimic sites by making small, hard-to-see changes to the URL.

Visual bias

- Visual priors: mappings off the **Unicode confusables list**
 - Designed to combat spoofing attacks
 - Hardly any mappings for Arabic and Kannada!

y	ȝ	Y	ȝ	y	ȝ	Y	y
0079 LATIN SMALL LETTER Y	0263 LATIN SMALL LETTER GAMMA	028F LATIN LETTER SMALL CAPITAL Y	03B3 GREEK SMALL LETTER GAMMA	0443 CYRILLIC SMALL LETTER U	04AF CYRILLIC SMALL LETTER STRAIGHT U	10E7 GEORGIAN LETTER QAR	
p	ƿ	ƿ	p	ƿ	P	ƿ	p
0070 LATIN SMALL LETTER P	03C1 GREEK SMALL LETTER RHO	03F1 GREEK RHO SYMBOL	0440 CYRILLIC SMALL LETTER ER	2374 APL FUNCTIONAL SYMBOL RHO	2CA3 COPTIC SMALL LETTER RO	1D429 MATHEMATICAL BOLD SMALL P	

Fully visual



Fully phonetic



Russian

Arabic

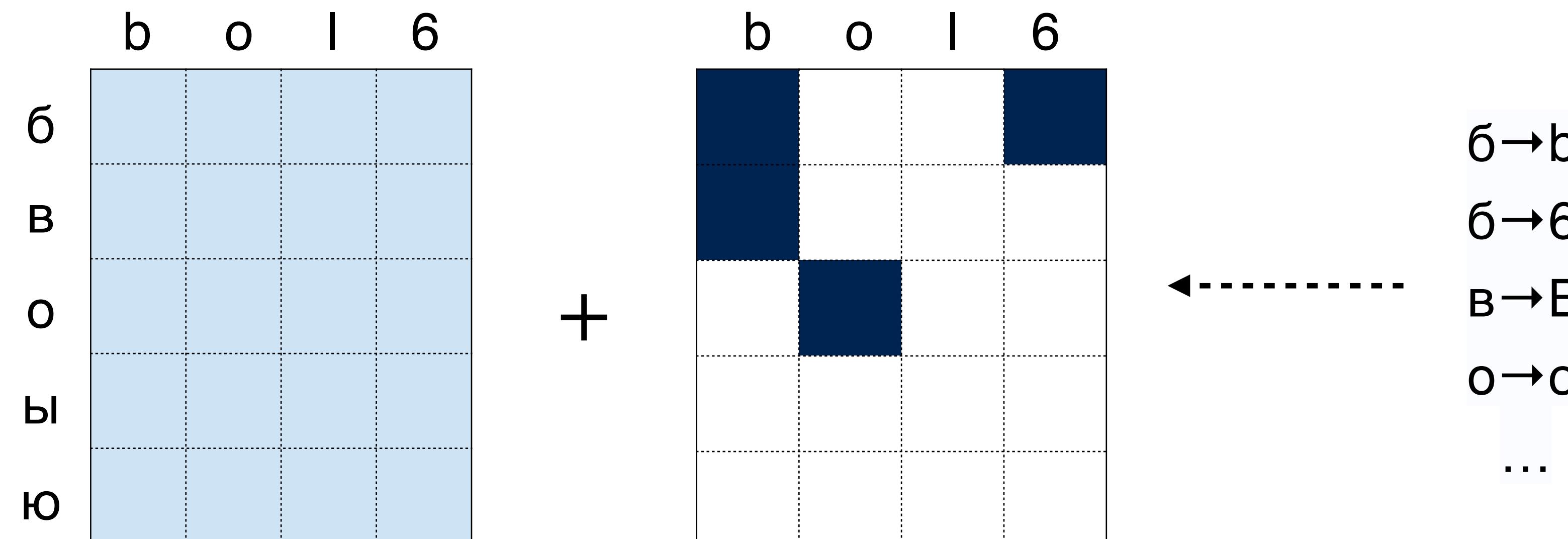
Kannada

Informative priors

- Use mappings of similar characters as **priors on emission parameters**

$$c_r | c_n \sim \text{Mult}(\theta_{c_n})$$

$$\theta \sim \text{Dir}(\alpha)$$

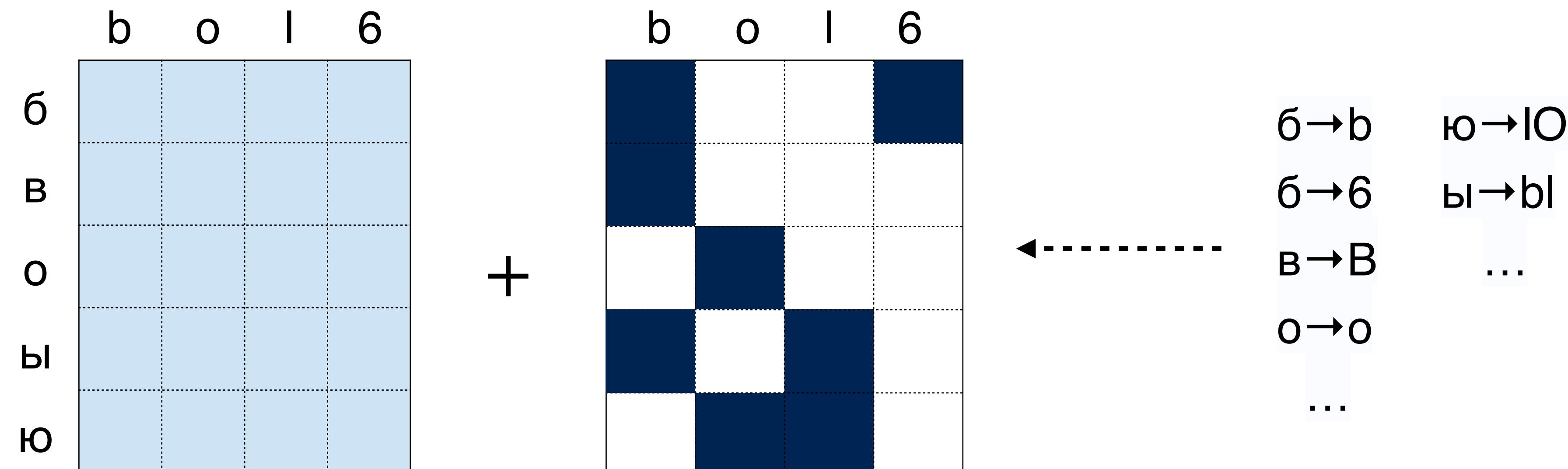


Informative priors

- Use mappings of similar characters as **priors on emission parameters**

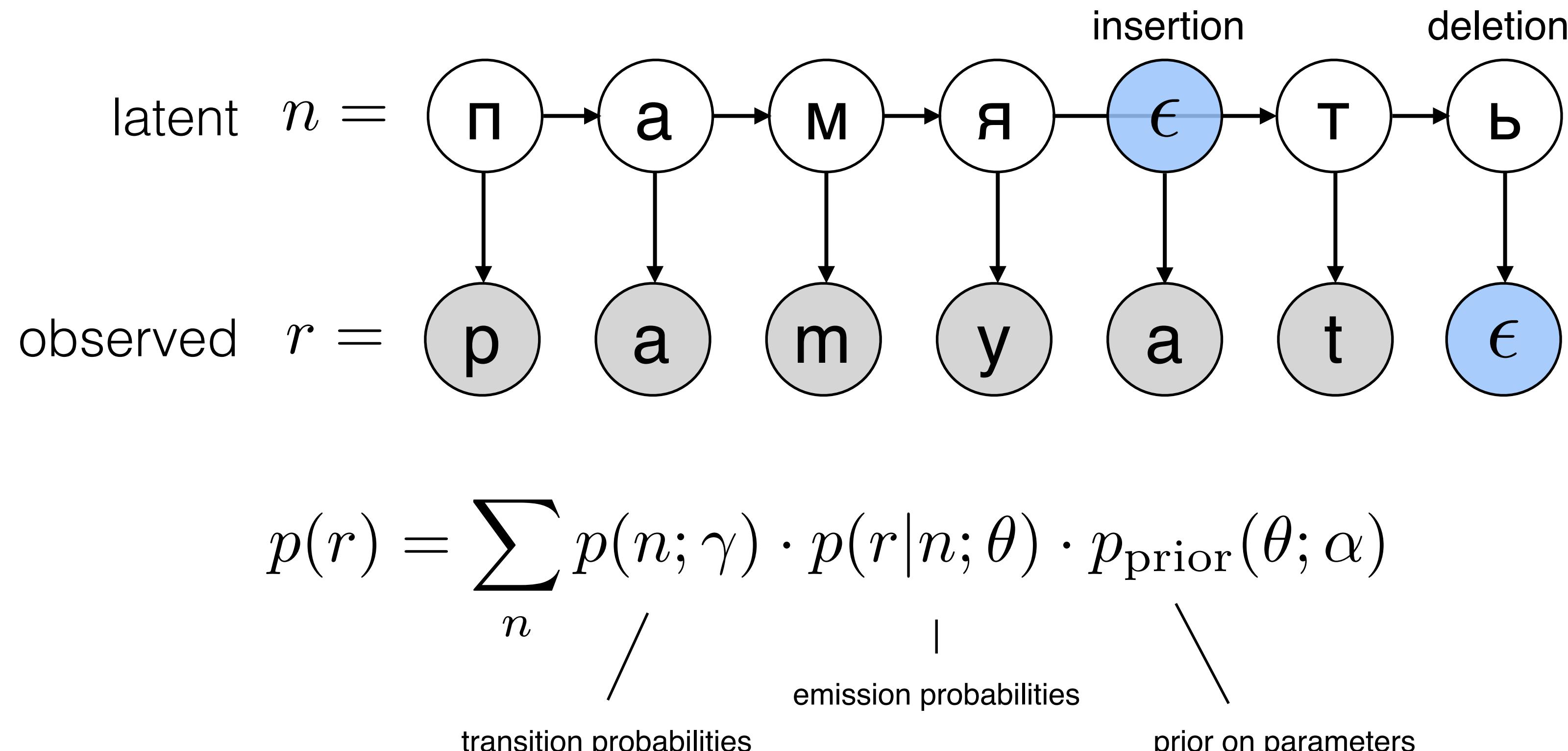
$$c_r | c_n \sim \text{Mult}(\theta_{c_n})$$

$$\theta \sim \text{Dir}(\alpha)$$



Noisy-channel model

- Representing latent alignments via **insertions and deletions**
 - Natural match for null alignments, but will be used for character n-grams too!



Proposed WFST

- Transition WFSA T
 - 6-gram LM built with OpenGrm (Roark et al., 2012)

Proposed WFST

- Transition WFSA T
 - 6-gram LM built with OpenGrm (Roark et al., 2012)
- Emission WFST E
 - Needs to support insertions and deletions

Proposed WFST

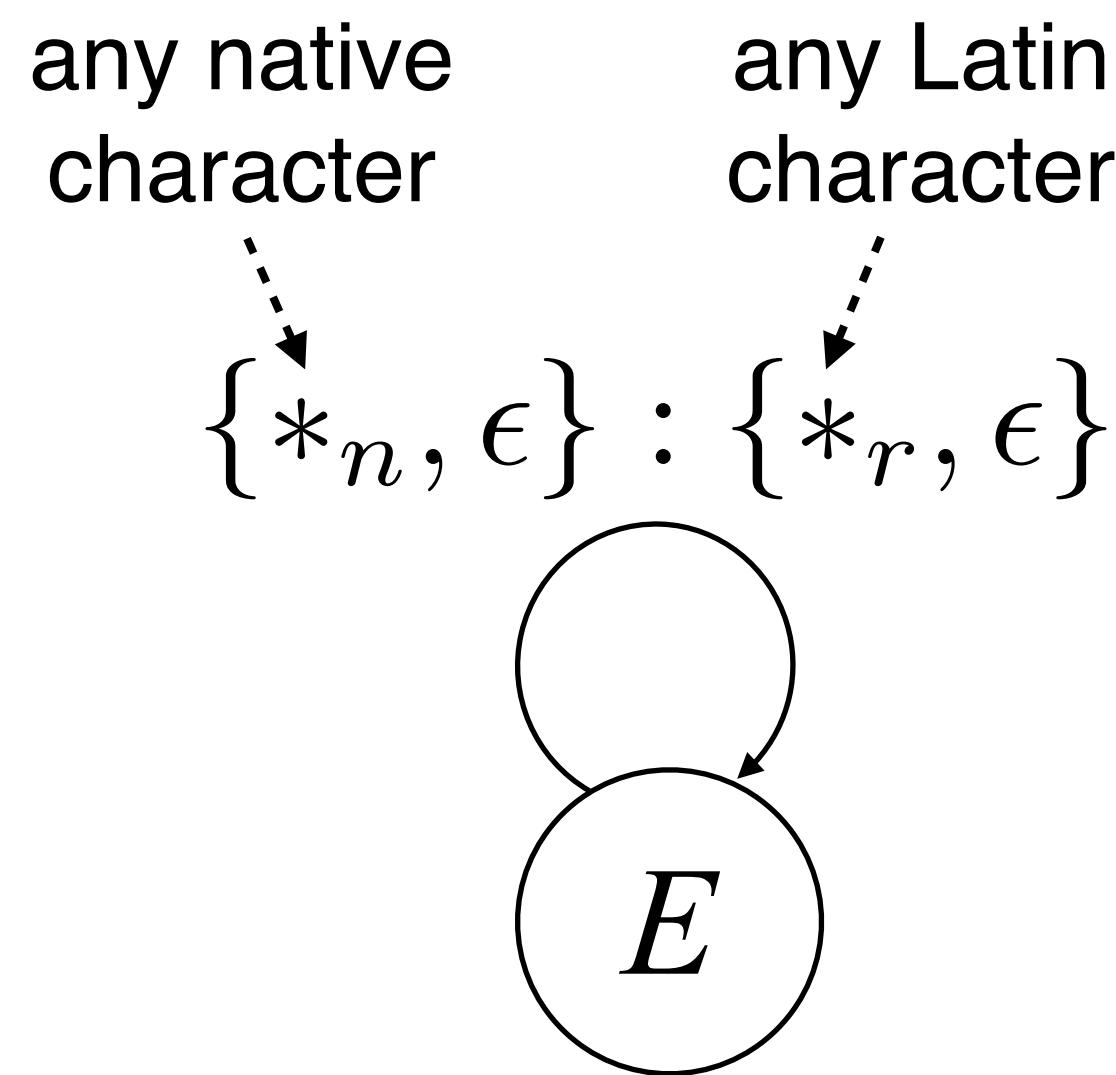
- Transition WFSA T
 - 6-gram LM built with OpenGrm (Roark et al., 2012)
- Emission WFST E
 - Needs to support insertions and deletions
- Input/output acceptors A

Proposed WFST

- Transition WFSA T
 - 6-gram LM built with OpenGrm (Roark et al., 2012)
- Emission WFST E
 - Needs to support insertions and deletions
- Input/output acceptors A
- Training and inference via finite-state methods in $T \circ E \circ A(r)$
 - OpenFst (Allauzen et al., 2007)

Emission model

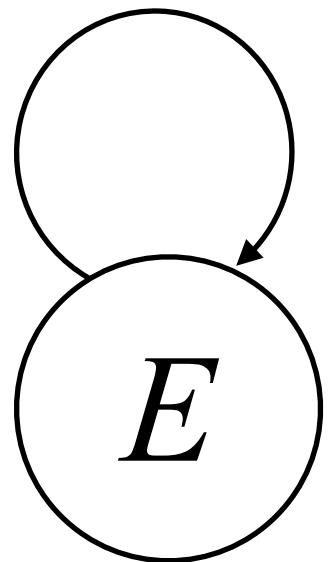
- Needs to support substitutions, insertions and deletions



Emission model

- Needs to support substitutions, insertions and deletions

$$\{ \text{I} , \epsilon \} : \{ \text{O} , \epsilon \}$$



$$E \circ A('O')$$

$$\text{I} : \text{O}$$

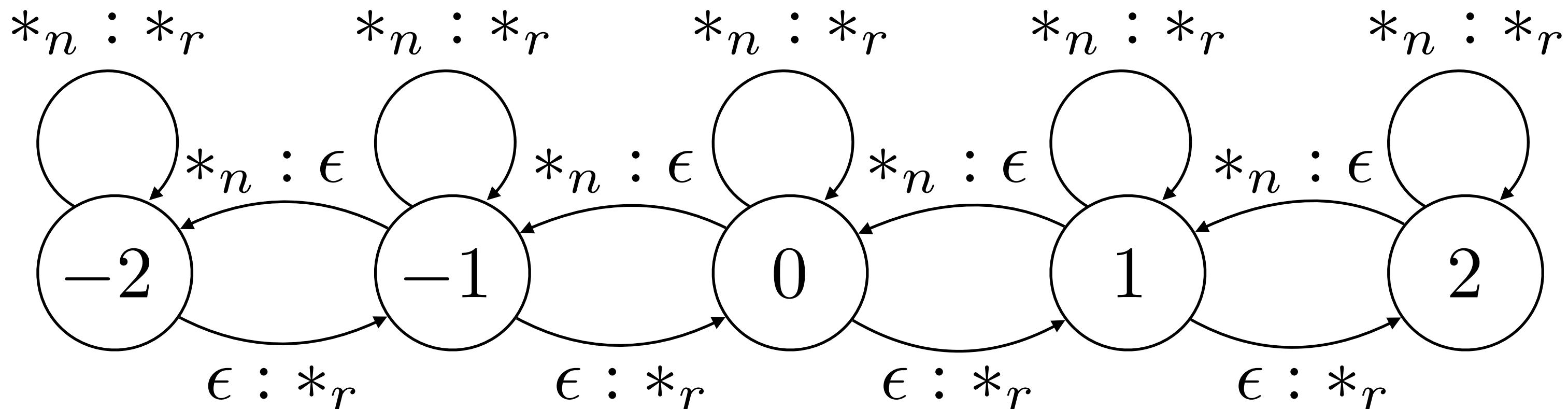
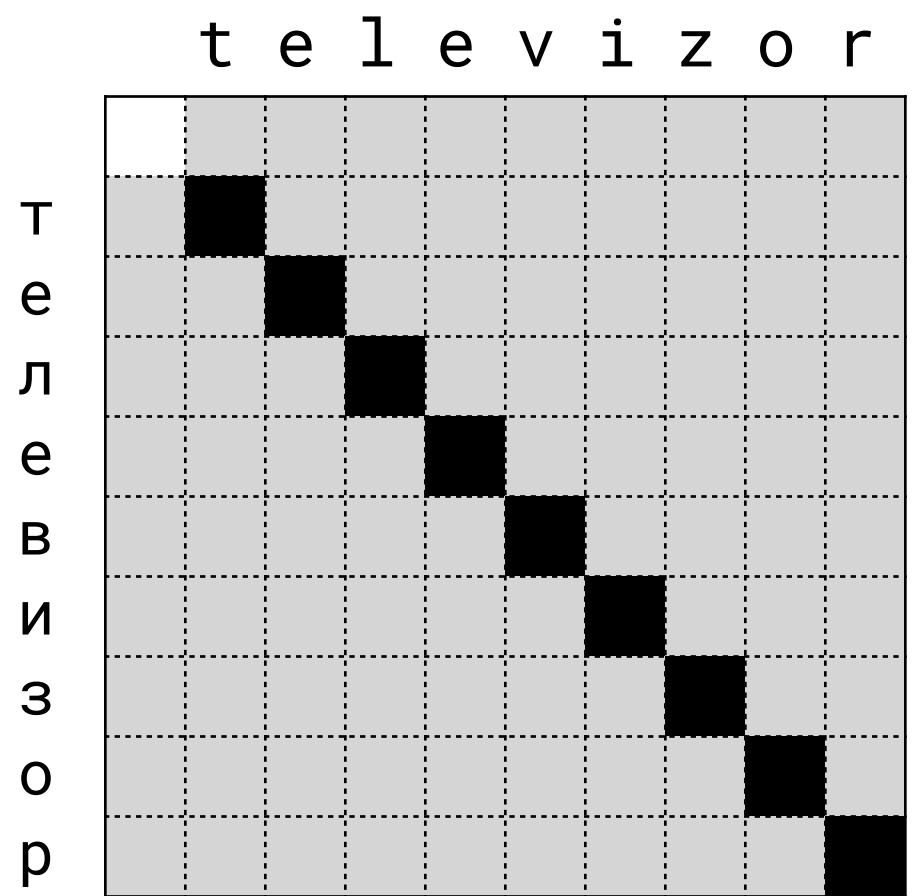
$$\epsilon : \text{O}$$

$$\text{I} : \epsilon \quad \epsilon : \text{O}$$

$$\text{I} : \epsilon \quad \text{I} : \epsilon \quad \dots \quad \epsilon : \text{O}$$

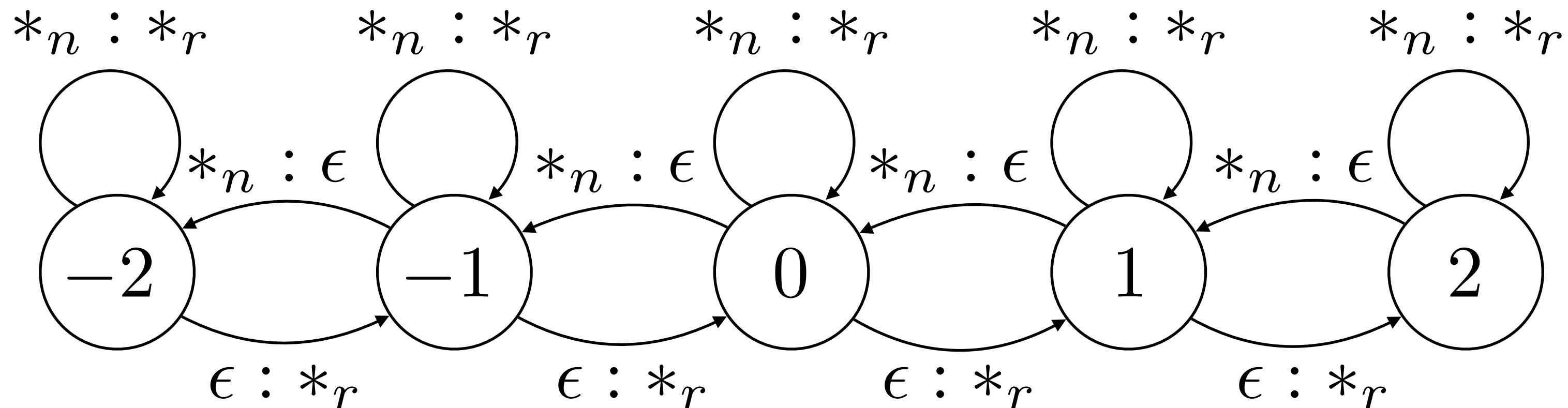
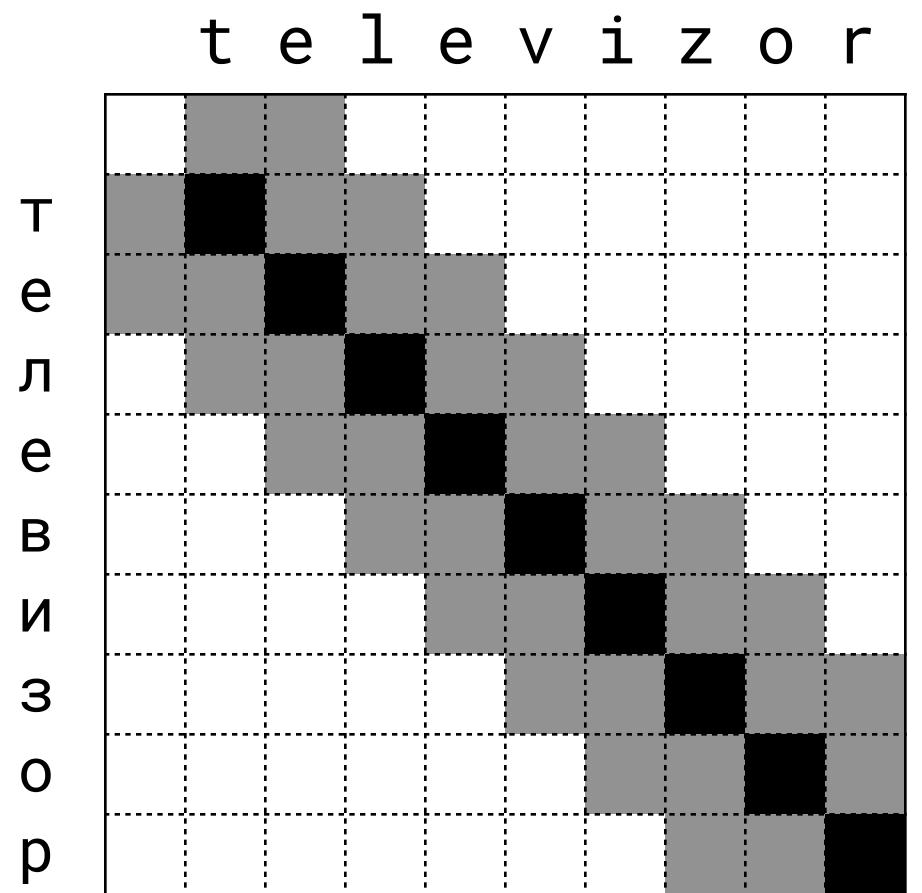
Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $| \# \text{ of insertions} - \# \text{ of deletions} |$



Emission model

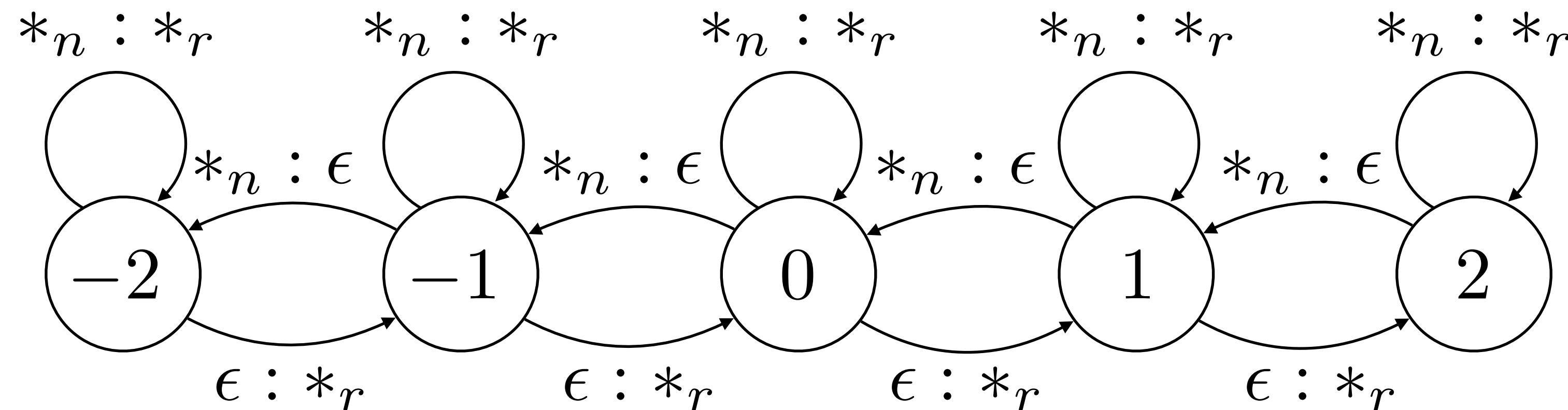
- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $| \# \text{ of insertions} - \# \text{ of deletions} |$



Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $|\# \text{ of insertions} - \# \text{ of deletions}|$

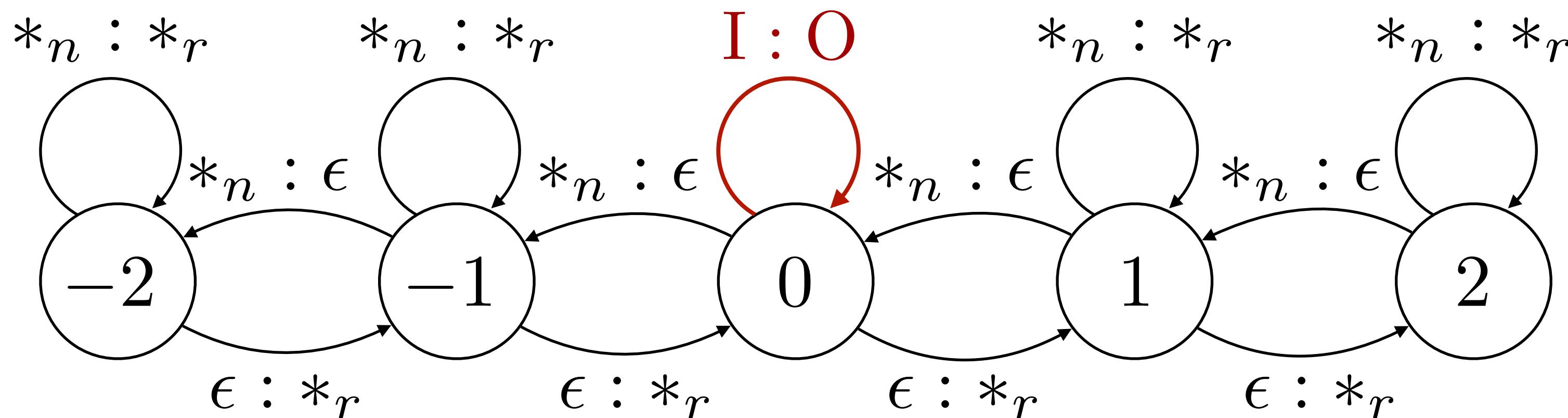
$E \circ A('O')$



Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $| \# \text{ of insertions} - \# \text{ of deletions} |$

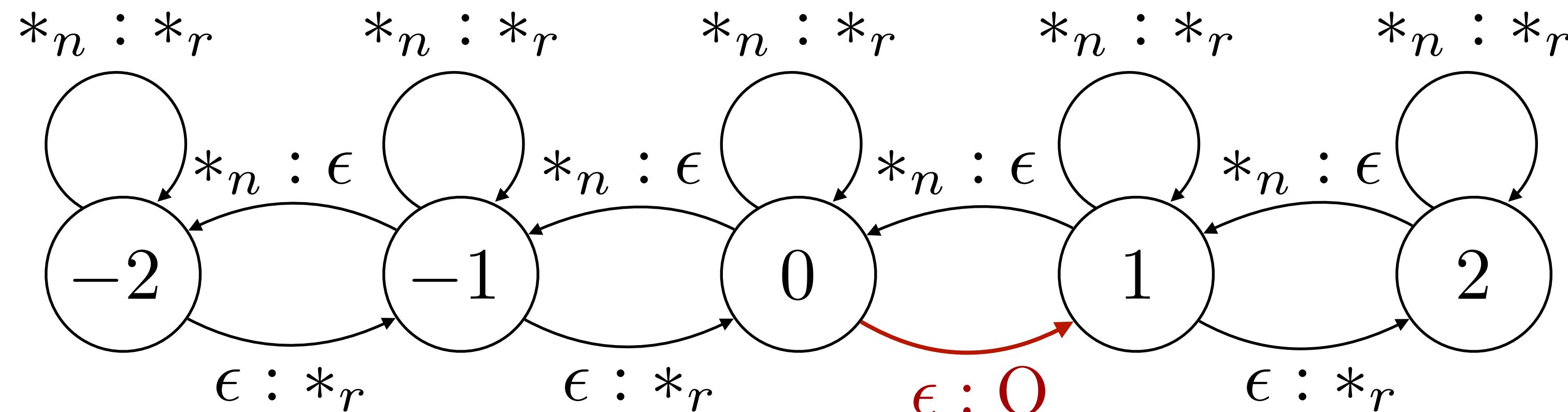
$E \circ A('O')$
I : O



Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $|\# \text{ of insertions} - \# \text{ of deletions}|$

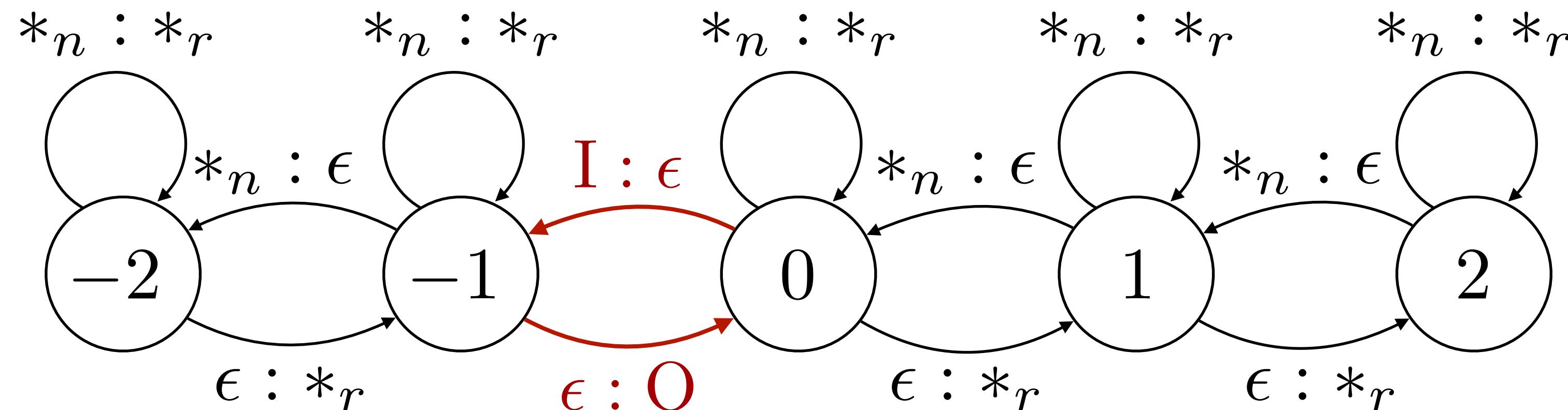
$$\begin{aligned} & E \circ A('O') \\ & \epsilon : O \end{aligned}$$



Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $| \# \text{ of insertions} - \# \text{ of deletions} |$

$$\begin{array}{l} E \circ A('O') \\ I : \epsilon \quad \epsilon : O \end{array}$$

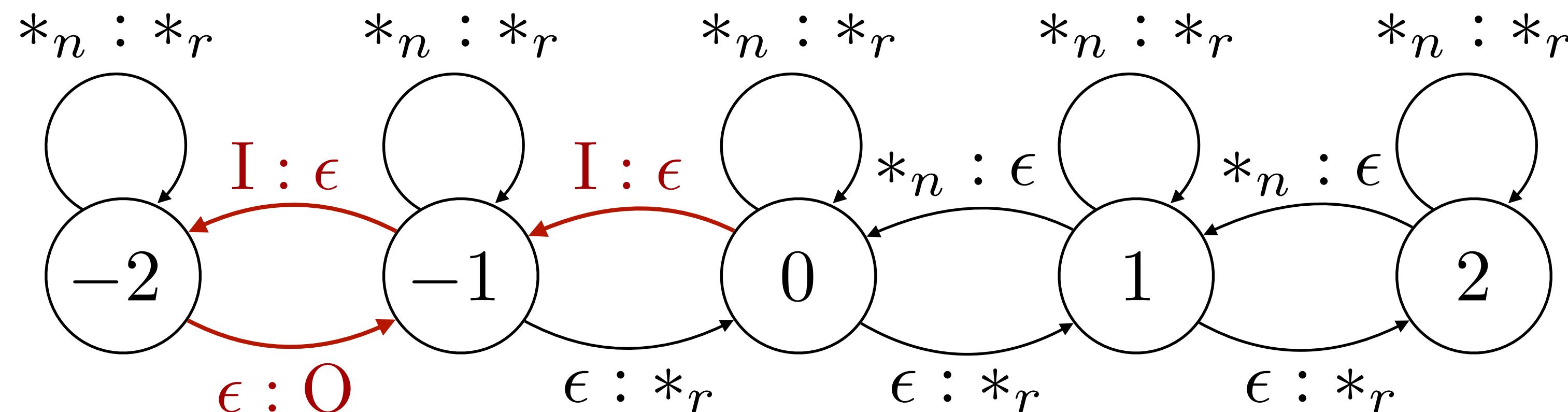


Emission model

- Needs to support substitutions, insertions and deletions
- Fixed limit on delay: $| \# \text{ of insertions} - \# \text{ of deletions} |$

$E \circ A('O')$

I : ϵ I : ϵ ... ϵ : O



Training and inference

- Training with EM algorithm
 - E-step: shortest distance in expectation semiring (Eisner, 2002)
 - M-step: parameter reestimation
- Many tricks to speed up training!
 - Stepwise batched EM (Liang and Klein, 2009)
 - Curriculum learning: shortest sequences first
 - Increasing LM order as training progresses
 - Pruning emission arcs during training
- Inference: shortest path in $(\max, +)$ semiring

Datasets

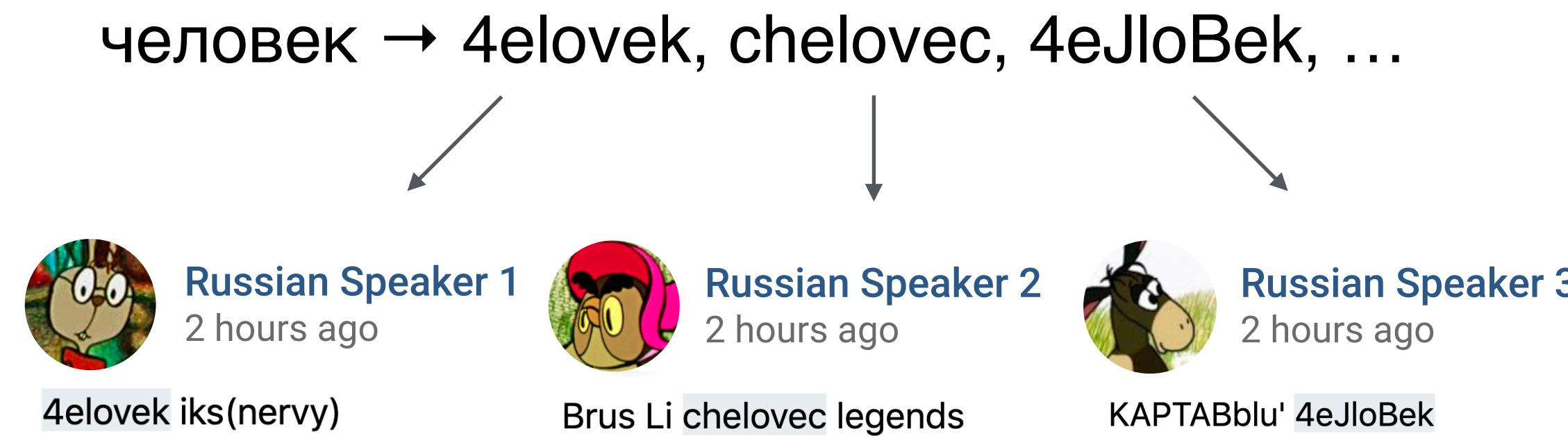
- Arabic: LDC BOLT dataset (Bies et al., 2014)
 - Arabizi SMS/chat dialogs, converted to CODA (Habash et al., 2012)
- Kannada: Dakshina dataset (Roark et al., 2020)
 - Kannada Wikipedia, romanizations elicited from native speakers
- Russian:
 - Romanized: collected and partly annotated data from social media
 - Native: Taiga corpus (Shavrina & Shapovalova, 2017), comments in political forums

Saba7 el 5eir!
Ezayeeky?



Russian data

- Romanizations of common words used as queries (Darwish, 2014)



- Manually removed sentences in other languages (e.g. Polish)
- Annotated validation and test with minor error correction

Source: proishodit s prirodoy **4to to very very bad**

Filtered: proishodit s prirodoy **4to to <...>**

Target: происходит с природой **ЧТО-ТО <...>**

Defining vocabulary

- Characters ~ Unicode codepoints?
 - Diacritics treated as separate characters
 - Non-printing symbols (e.g. ZWNJ) treated as separate characters
- ક = ka
ક + ઊ = ku
ક + ઊંગું = k

રાજુકુમારુ

rāj|kumār

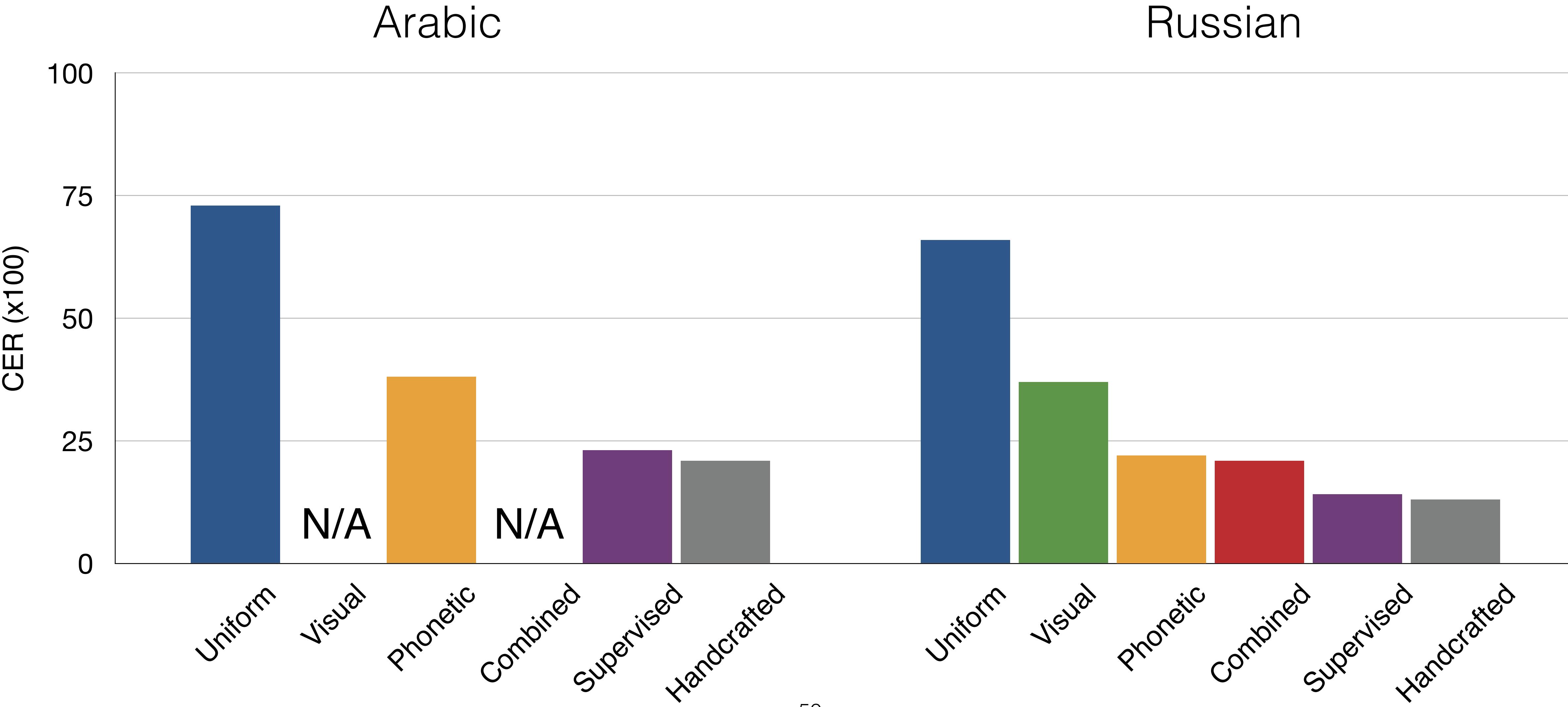
ર ા જ ઊ ક ઊ મ ા ર ઊ
Ra Aa Ja Ø [ZW_NJ] Ka U Ma Aa Ra Ø

Defining vocabulary

- Characters ~ Unicode codepoints?
 - Diacritics treated as separate characters
 - Non-printing symbols (e.g. ZWNJ) treated as separate characters
 - Combined characters may have multiple Unicode representations

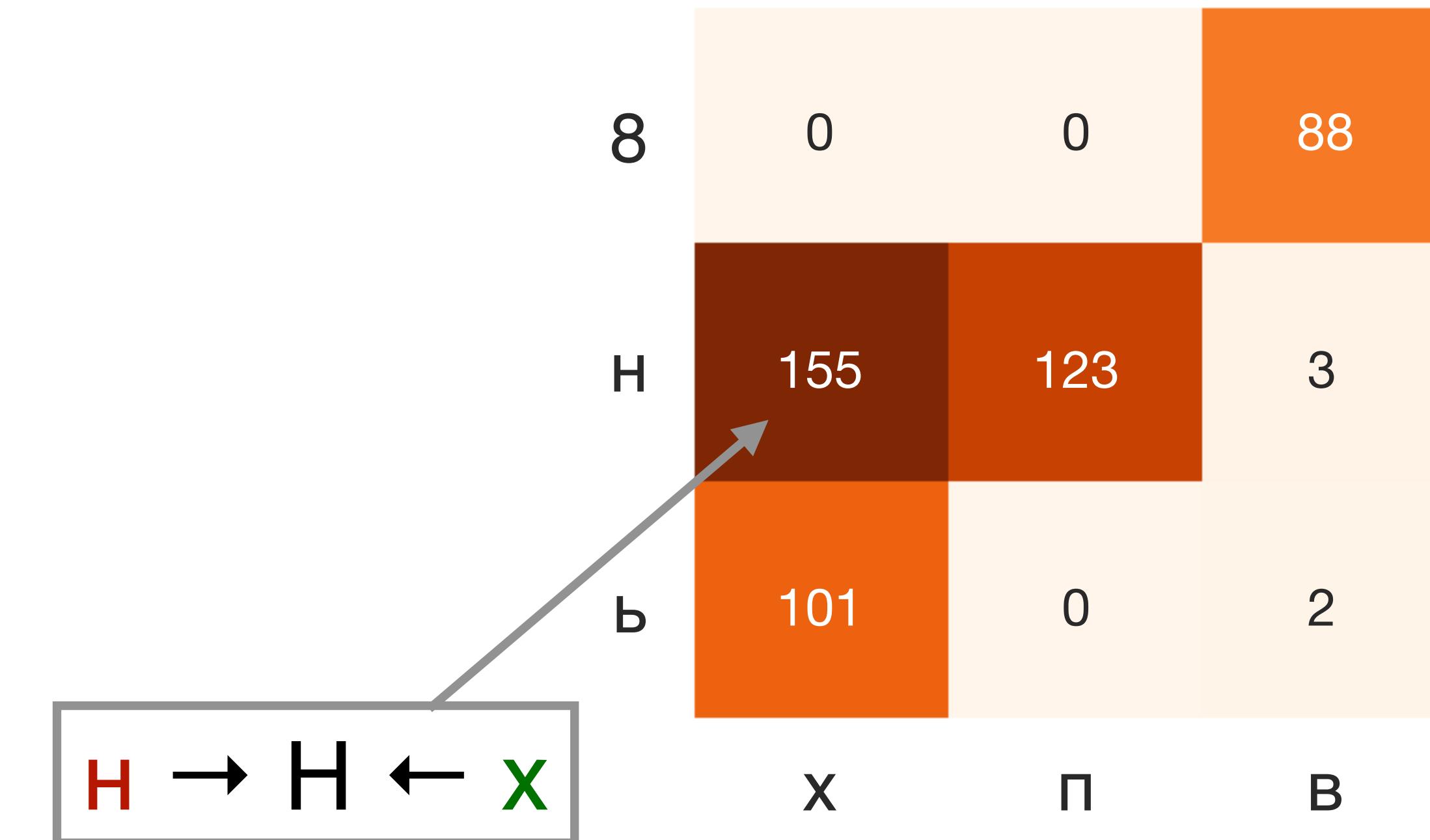
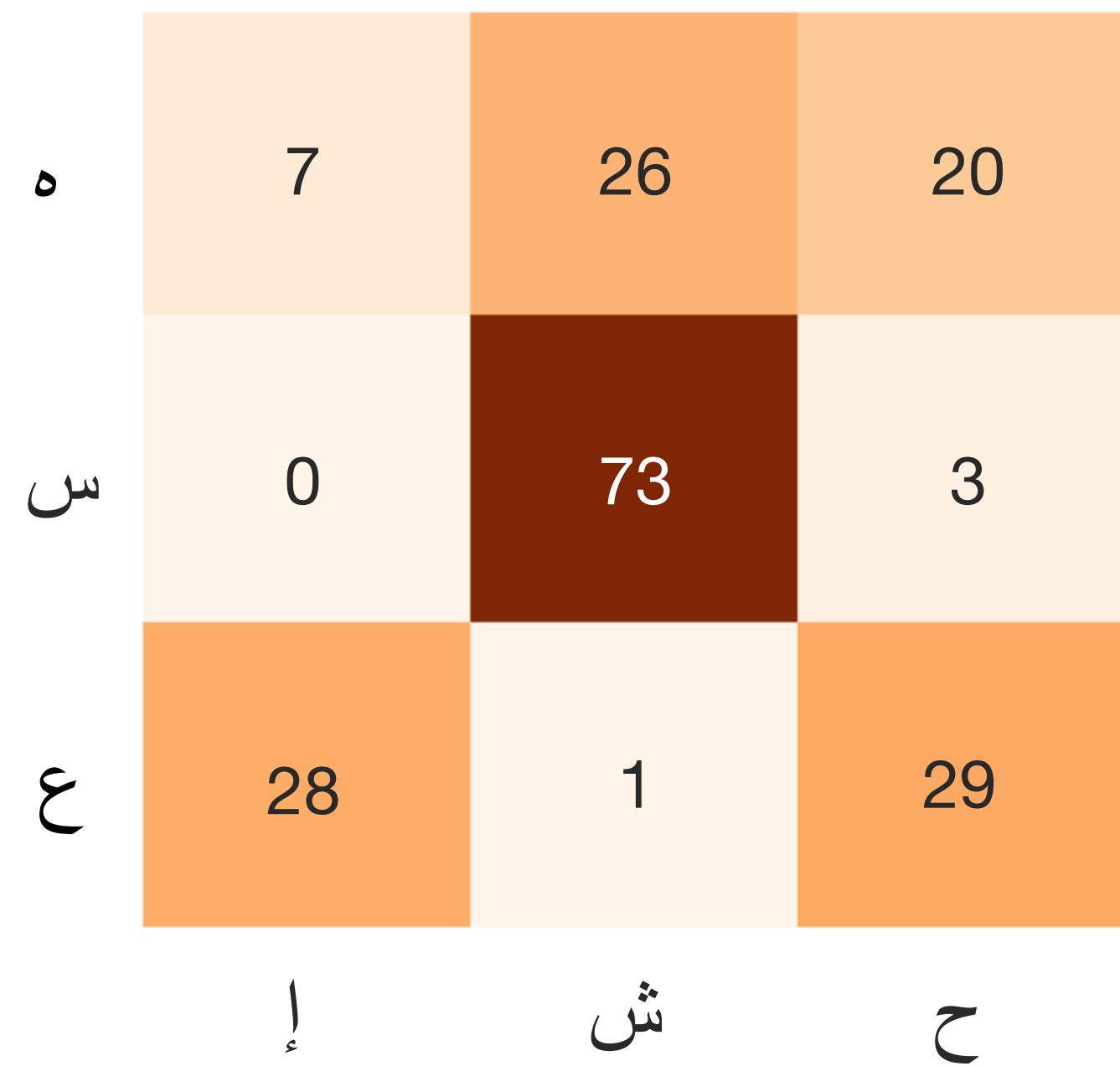
For	Use	Do Not Use
ନ୍ଦ୍ର	0C8A	<0C89, 0CBE>
ଚ୍ଛ	0C94	<0C92, 0CCC>
ବ୍ରା	0CE0	<0C8B, 0CBE>

WFST results



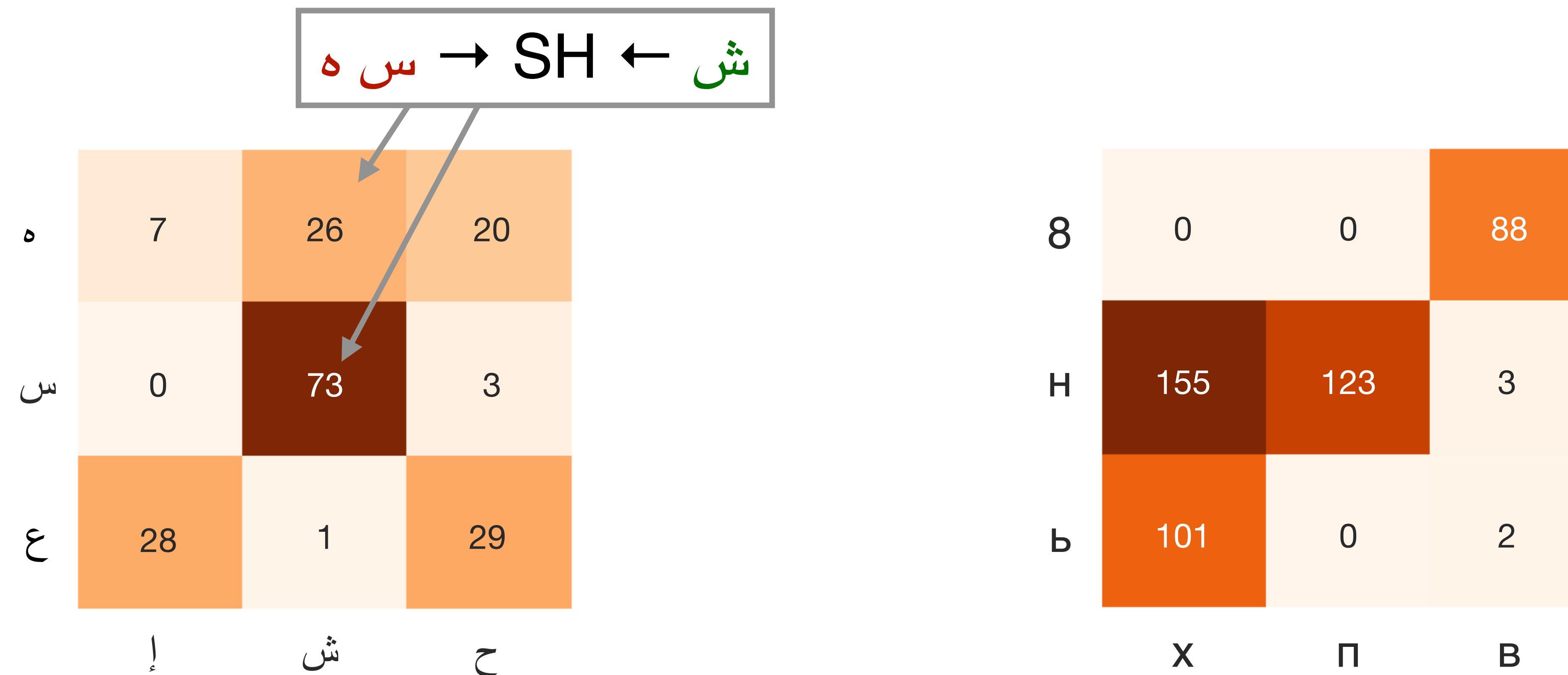
WFST error analysis

- Incorrect choice of plausible de-romanization (e.g. visual instead of phonetic)



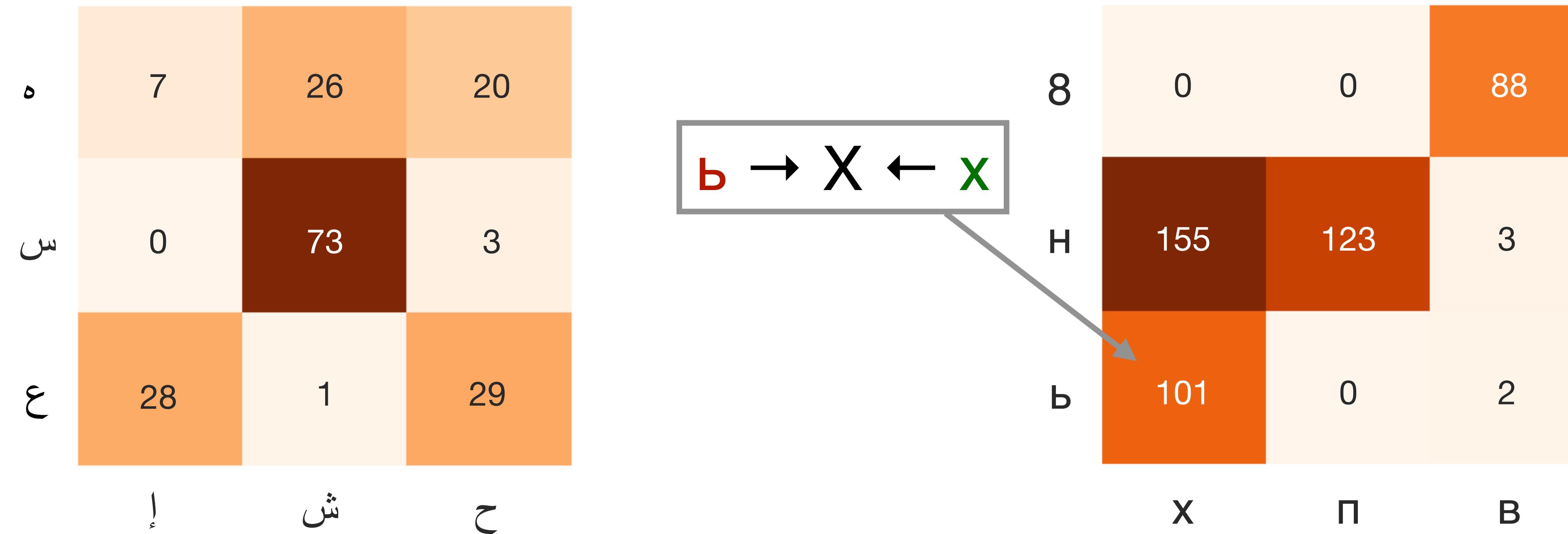
WFST error analysis

- Inability to handle digraphs like SH



WFST error analysis

- Distracted by spurious mappings in priors



Model classes

WFSTs are **structured**

- ✓ Easy to encode constraints
- ✓ Can learn from small data
- ✗ Slow exact maximization
- ✗ Weak n-gram language model

Seq2seqs are **powerful**

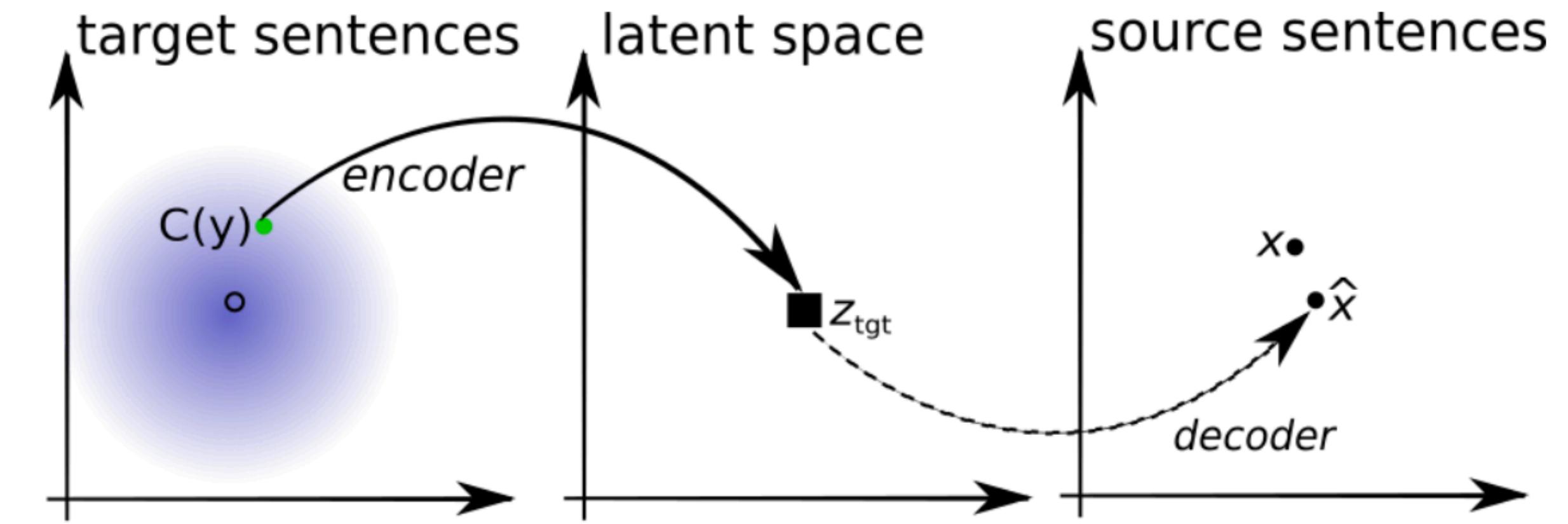
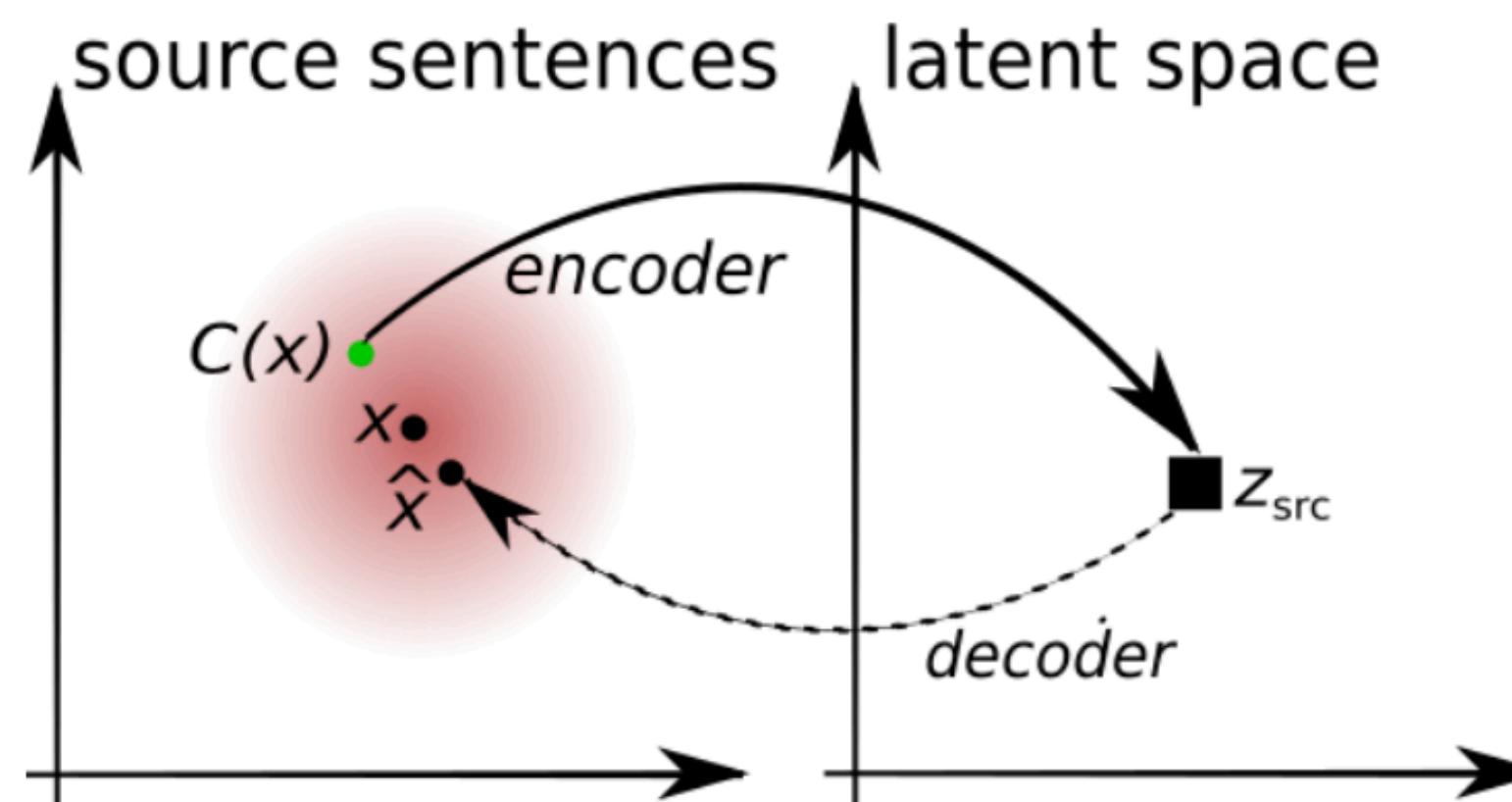
- ✓ Strong language model
- ✓ Faster batch processing
- ✗ Need large training data
- ✗ Hallucinations and search errors

In our case, both are trained **unsupervised!**

M Ryskina, E Hovy, T Berg-Kirkpatrick, MR Gormley. Comparative Error Analysis in Neural and Finite-state Models for Unsupervised Character-level Transduction. SIGMORPHON 2021.

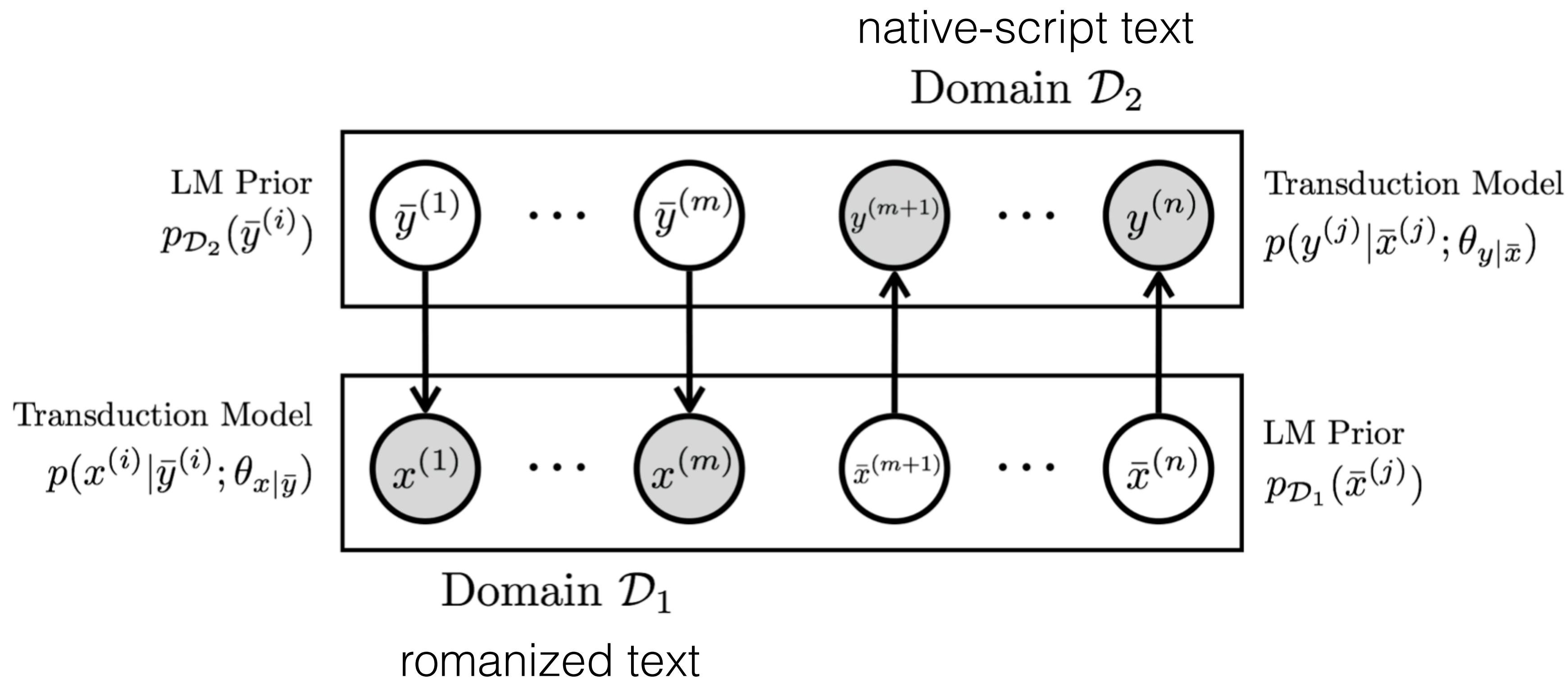
Unsupervised seq2seq

- Unsupervised neural machine translation (UNMT; Lample et al., 2018)
 - Auto-encoding: reconstructing a sentence from its noisy version
 - Back-translation: round trip through the latent space
 - Adversarial: discriminating between sentences in two domains



Unsupervised seq2seq

- Probabilistic formulation of UNMT: deep latent sequence model (He et al., 2020)

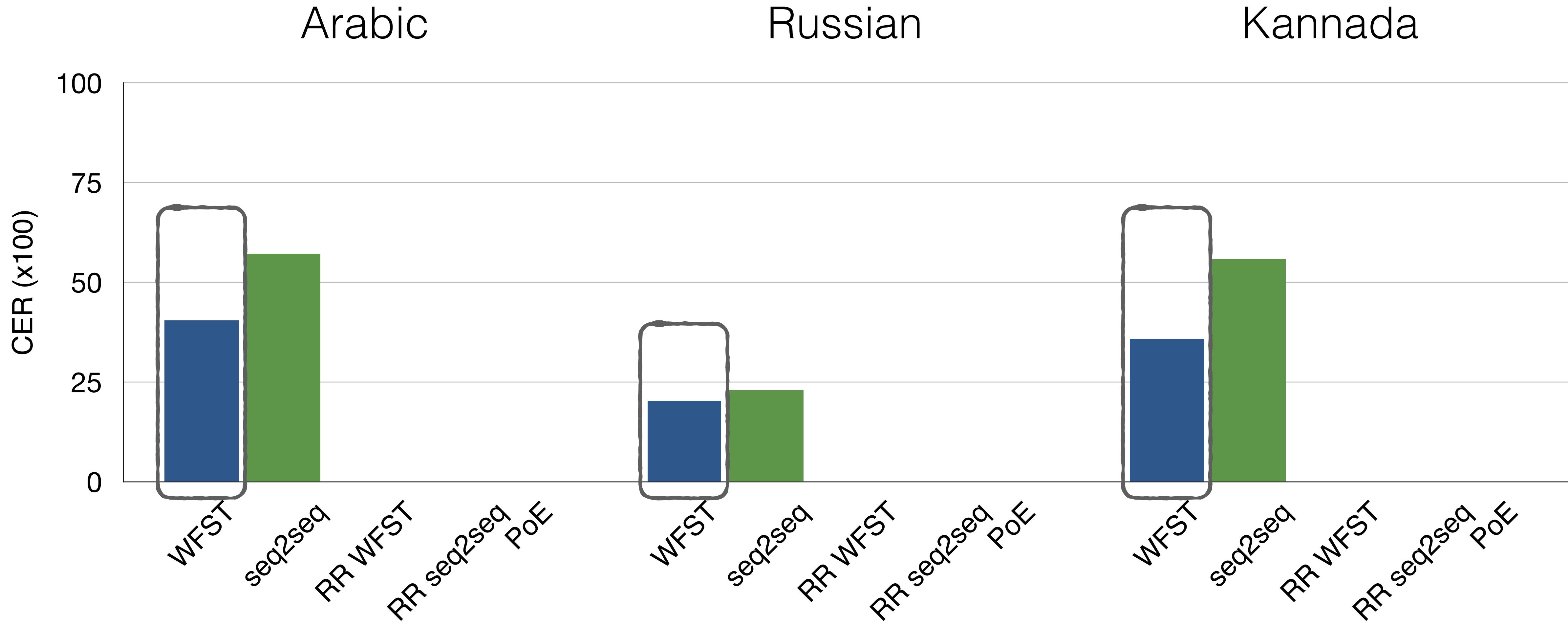


Model combinations

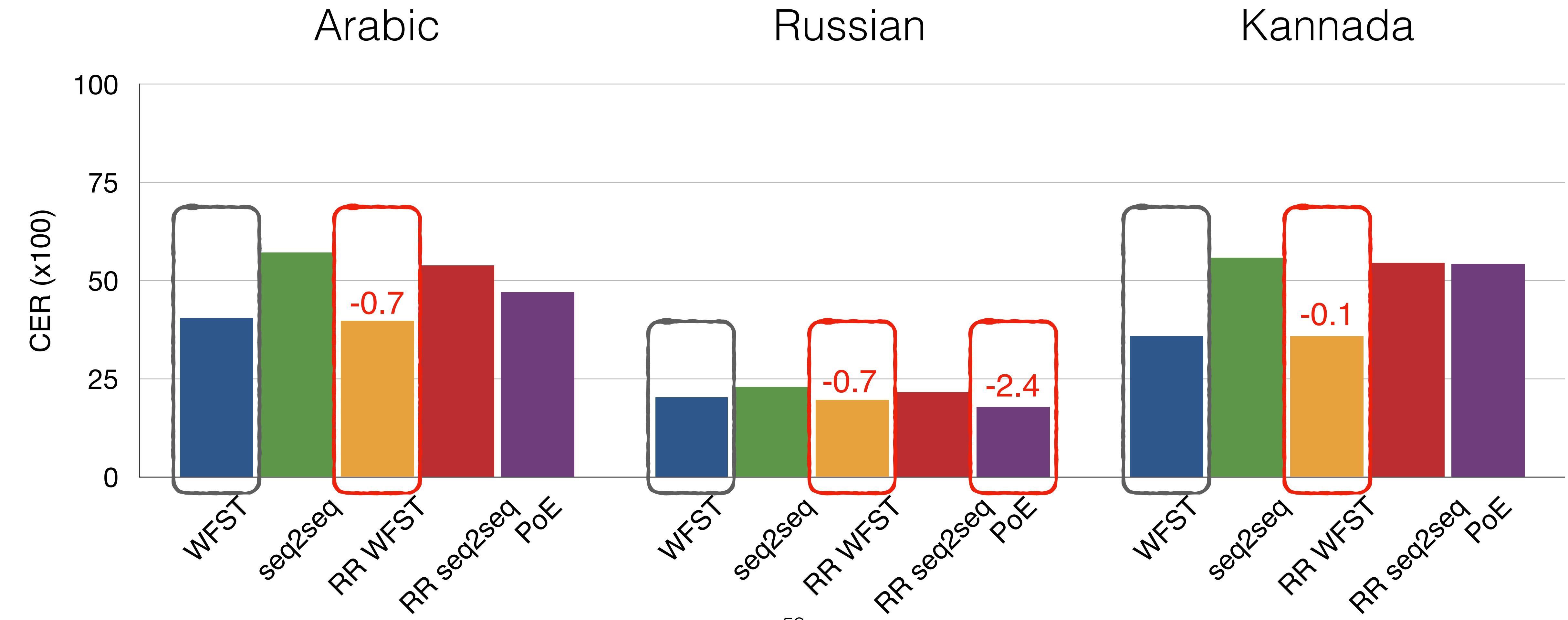
- Reranking
 - M1 generates top k candidate outputs
 - M2 selects the highest-scoring candidate
- Product of experts
 - Beam search on the WFST lattice
 - WFST arcs reweighted with Seq2seq softmax at the corresponding timestep
 - Deletions of input characters are not reweighted
 - Candidates are grouped by consumed input length
- We train the models separately and combine at test time

Results

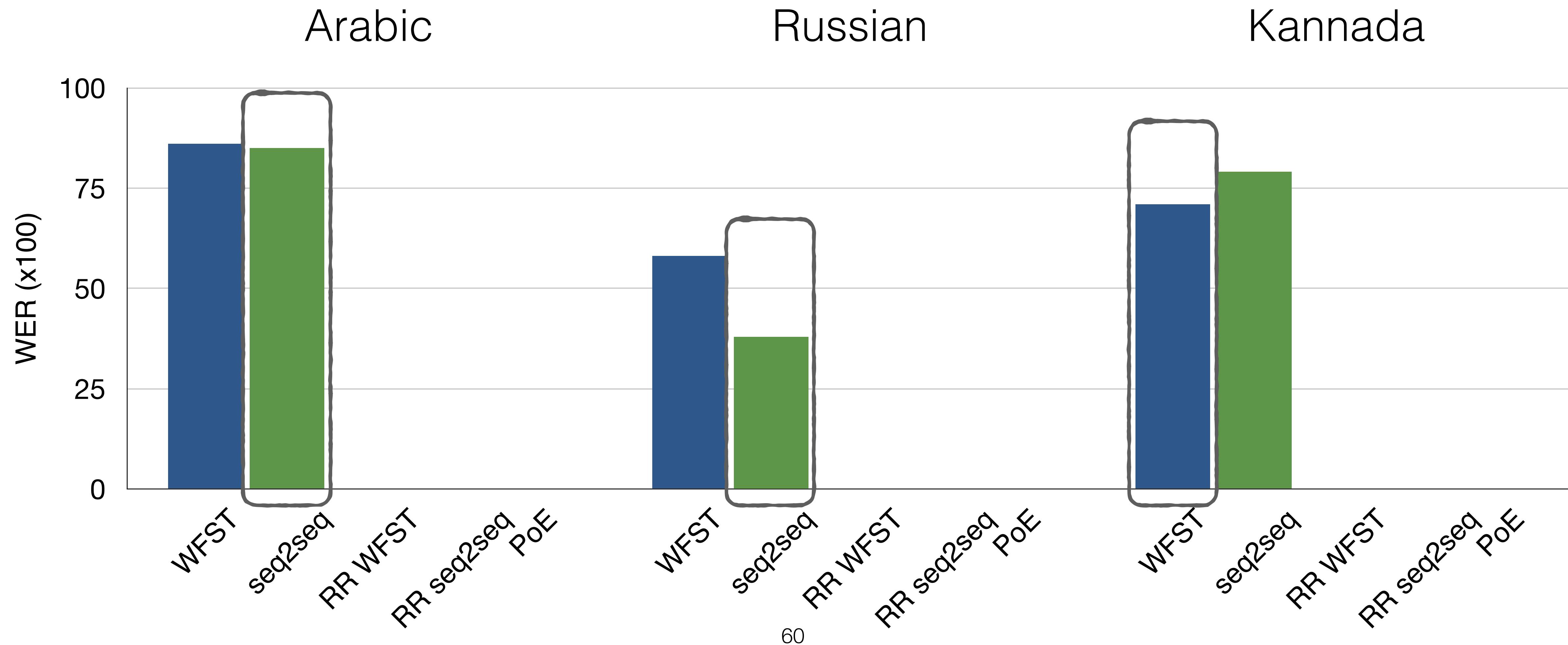
Base models are trained on different amounts of data!



Results

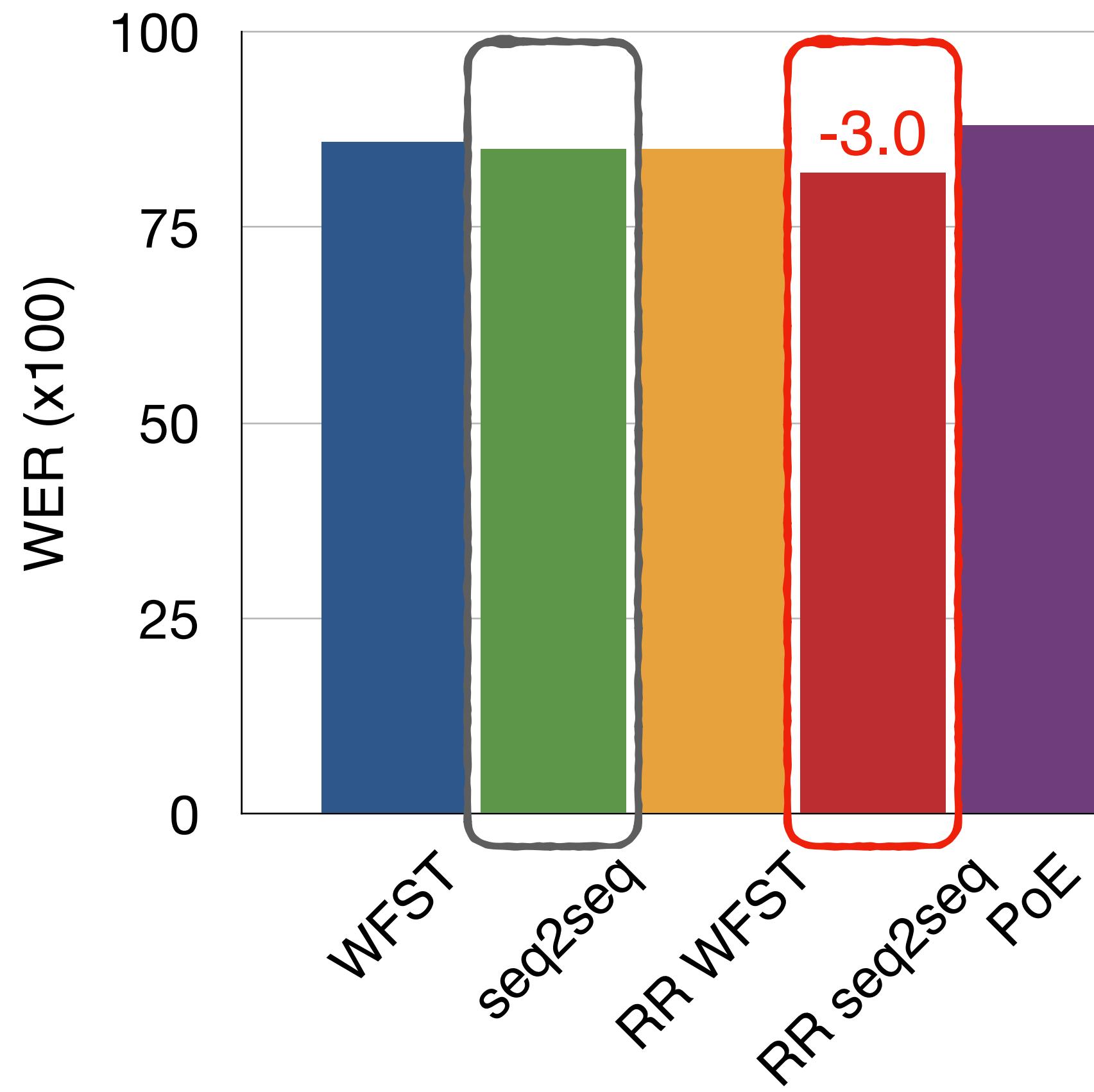


Results



Results

Arabic

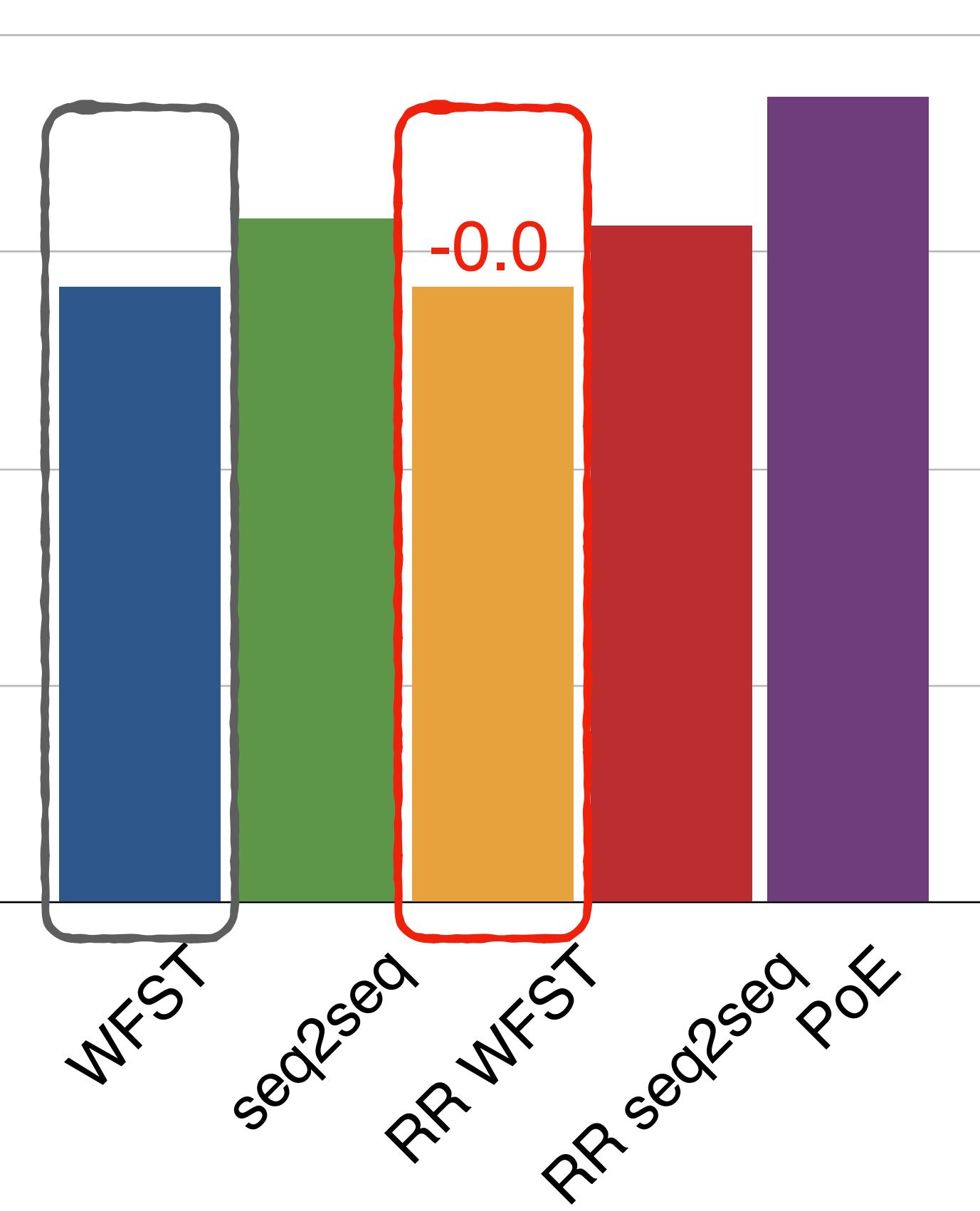


Russian

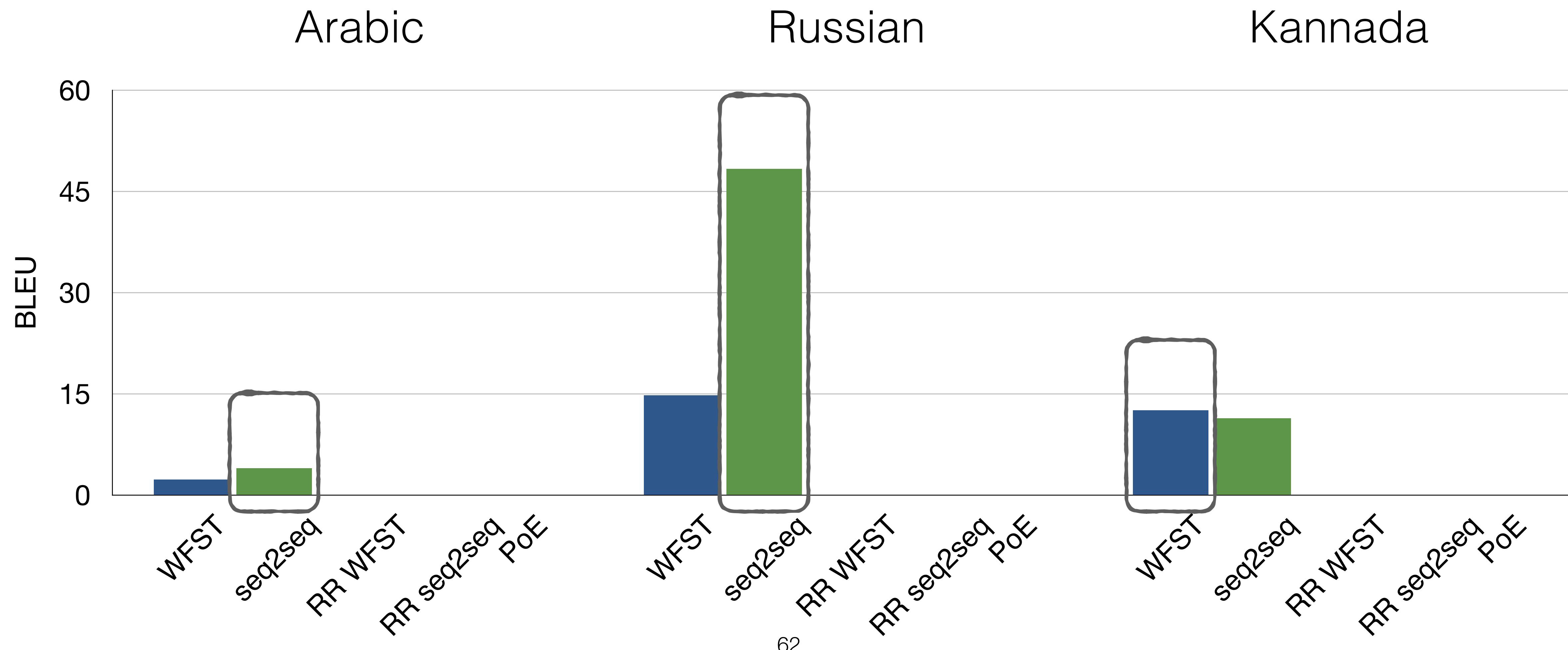
Russian

Russian

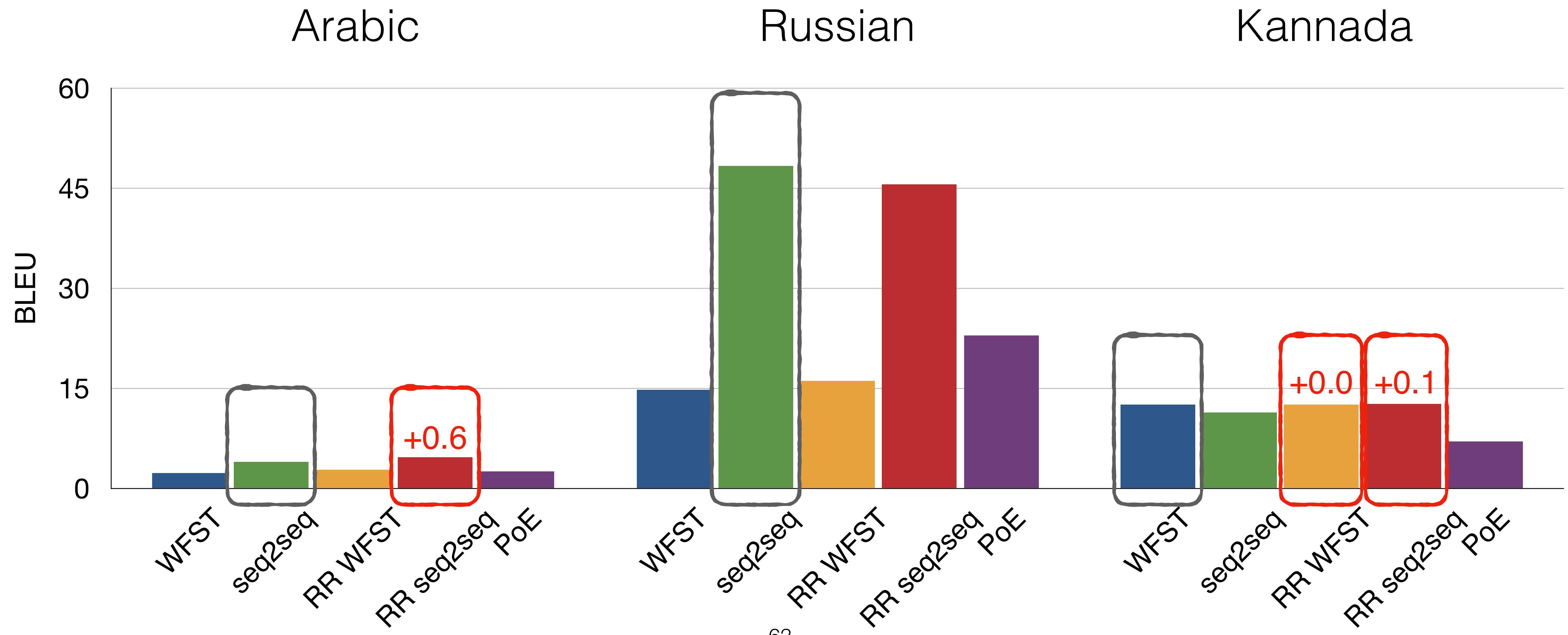
Kannada



Results



Results



Error analysis

Input	kongress ne odobril biudjet dlya osuchestvleniye "bor'bi s kommunizmom" v yuzhniy amerike.	
Ground truth	конгресс не одобрил бюджет для осуществления "борьбы с коммунизмом" в южной америке.	kongress ne odobril bjudžet dlja osuščestvlenija "bor'by s kommunizmom" v južnoj amerike.
WFST	конгресс не одобрил виудет для осу[с]чествениы[e] "бор#би с коммунизмом" в уузнани америке.	kongress ne odobril viudet dla osuščestvleniye "bor#bi s kommunizmom" v uuznani amerike.
Reranked WFST	конгресс не одобрил видет дела осу[с]чествениы[e] "бор#би с коммунизмом" в уузнани америке.	kongress ne odobril videt dela osuščestvleniye "bor#bi s kommunizmom" v uuznani amerike. #=UNK
Seq2Seq	конгресс не одобрил бы удивительно с коммунизмом" в южный америке.	kongress ne odobril by udivitel'no s kommunizmom" v južnyj amerike.
Reranked Seq2Seq	конгресс не одобрил бюджет для осуществление "борьбы с коммунизмом" в южный америке.	kongress ne odobril bjudžet dlja osuščestvlenie "bor'by s kommunizmom" v južnyj amerike.
Product of experts	конгресс не одобрил бидет для а осуществениы[e] "борьбы с коммунизмом" в уузник амери	kongress ne odobril bidet dlja a osuščestvleniye "bor'by s kommunizmom" v uuznnik ameri

Error analysis

Input	kongress ne odobril biudjet dlya osuchestvleniye "bor'bi s kommunizmom" v yuzhniy amerike.		
Ground truth	конгресс не одобрил бюджет для осуществления "борьбы с коммунизмом" в южной америке.	kongress ne odobril bjudžet dlja osuščestvlenija "bor'by s kommunizmom" v južnoj amerike.	
WFST	конгресс не одобрил виудет для осу[с]чествениы[e] "бор#би с коммунизмом" в уузнани америке.	kongress ne odobril viudet dla osuščestvleniye "bor#bi s kommunizmom" v uuznani amerike.	
Reranked WFST	конгресс не одобрил видет дела осу[с]чествениы[e] "бор#би с коммунизмом" в уузнани америке.	kongress ne odobril videt dela osuščestvleniye "bor#bi s kommunizmom" v uuznani amerike.	
Seq2Seq	конгресс не одобрил бы удивительно с коммунизмом" в южный америке.	kongress ne odobril by udivitel'no s kommunizmom" v južnyj amerike.	Hallucination
Reranked Seq2Seq	конгресс не одобрил бюджет для осуществление "борьбы с коммунизмом" в южный америке.	kongress ne odobril bjudžet dlja osuščestvlenie "bor'by s kommunizmom" v južnyj amerike.	Incorrect but faithful
Product of experts	конгресс не одобрил бидет для а осуществениы[e] "борьбы с коммунизмом" в уузник амери	kongress ne odobril b1det dlja a osuščestvleniye "bor'by s kommunizmom" v uuznnik ameri	

High-level takeaways

- Model combinations **still suffer from search issues**

Source: `eto uzhe (strashno skazat') stariy rolik.`

Target: `это уже (страшно сказать) старый ролик`

Gloss: ‘By now this is (I’m almost afraid to say it) an old video’

Final beam hypotheses and reranker scores:

456.7, `единая россия уже #страшно сказать) старый`

502.0, `единоросы уже #страшно сказать) старый рол`

482.0, `единороссы уже #страшно сказать) старый ро`

456.8, `единую россию уже #страшно сказать) старый`

449.8, `единой россии уже #страшно сказать) старый`

High-level takeaways

- Model combinations **still suffer from search issues**

Source: **eto uzhe (strashno skazat') stariy rolik.**

Target: **это уже (страшно сказать) старый ролик**

Gloss: ‘This’ **ow this is (I’m almost afraid to say it) an old video’**

Final beam hypotheses and reranker scores:

456.7, **единая россия** уже #страшно сказать) старый

502.0, **единоросы** уже #страшно сказать) старый рол

482.0, **единороссы** уже #страшно сказать) старый ро

456.8, **единую россию** уже #страшно сказать) старый

449.8, **единой россии** уже #страшно сказать) старый

‘United Russia’

High-level takeaways

- Model combinations **still suffer from search issues**

Source: **eto uzhe (strashno skazat') stariy rolik.**

Target: **это уже (страшно сказать) старый ролик**

Gloss: ‘By now this is (I’m almost afraid to say it) an old video’

Final beam hypotheses and reranker scores:

456.7, **единая россия уже #страшно сказать) старый**

502.0, **единоросы уже #страшно сказать) старый рол**

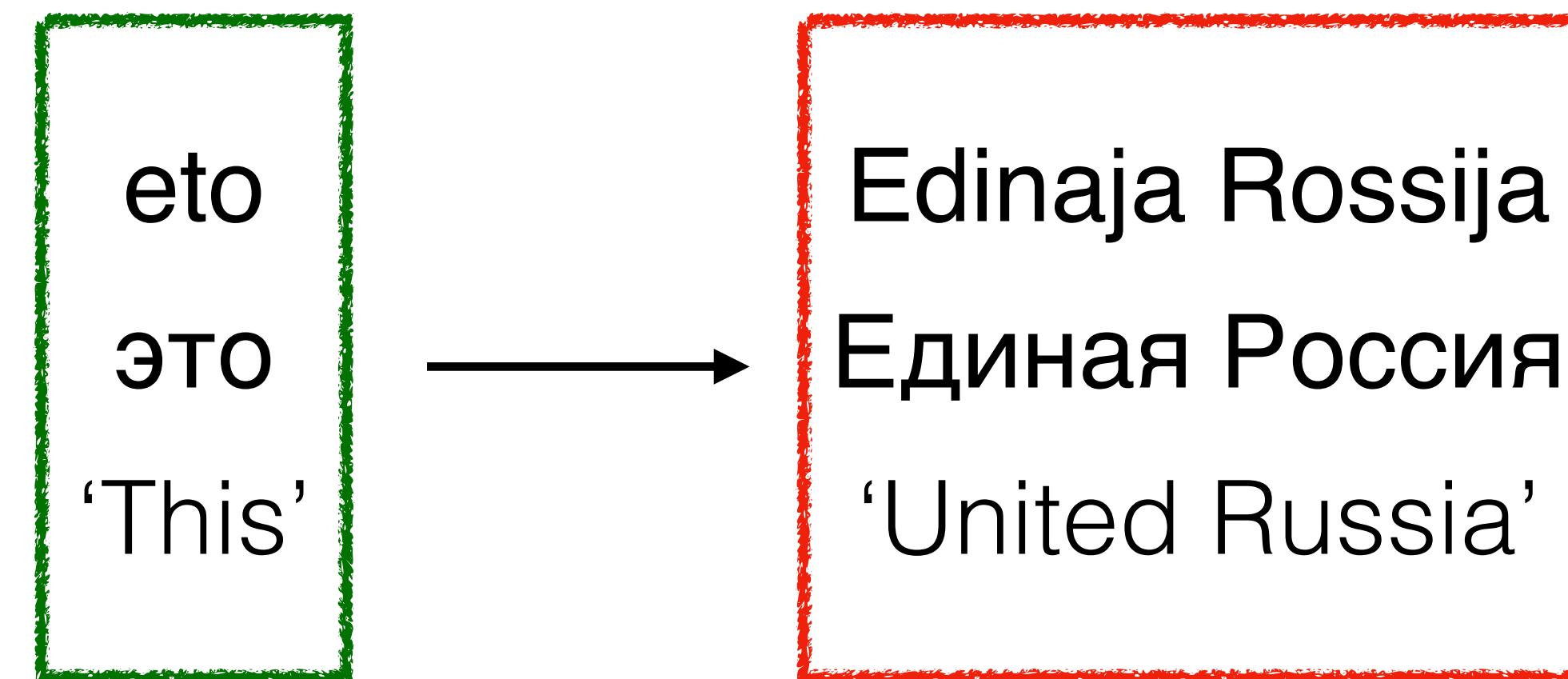
482.0, **единороссы уже #страшно сказать) старый ро**

456.8, **единую россию уже #страшно сказать) старый**

449.8, **единой россии уже #страшно сказать) старый**

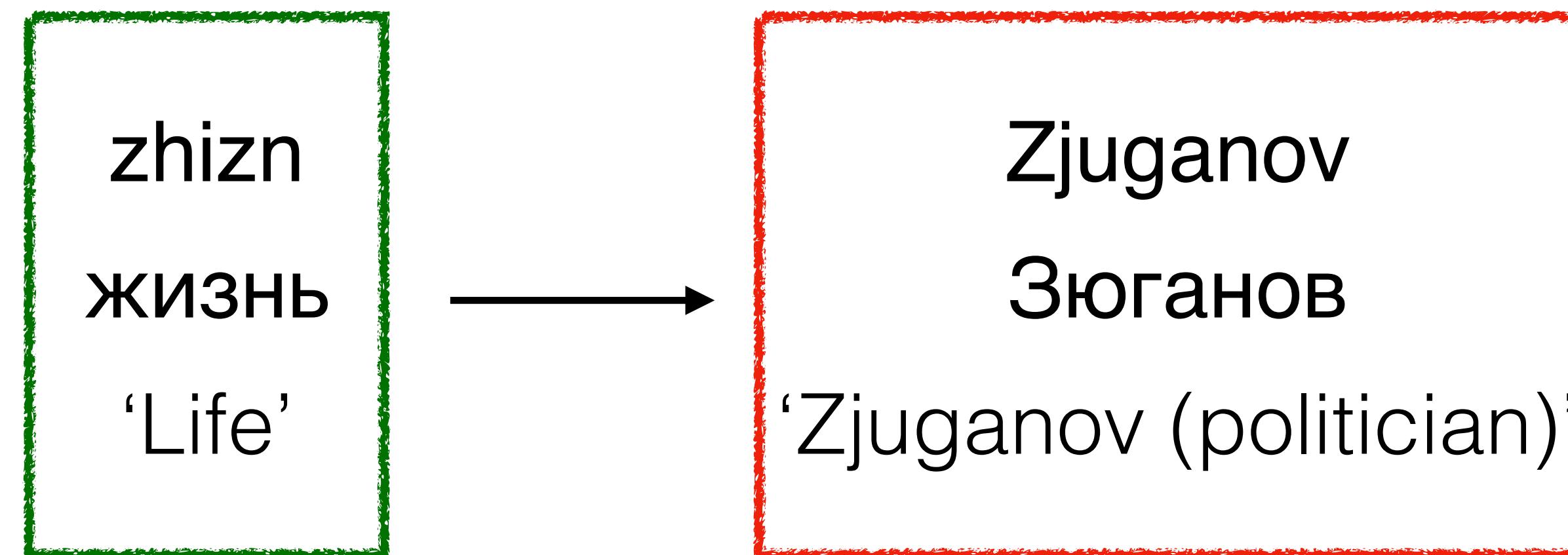
High-level takeaways

- Seq2seq is more sensitive to **distributional shifts**
 - Remember that our Cyrillic data comes from political discussion groups



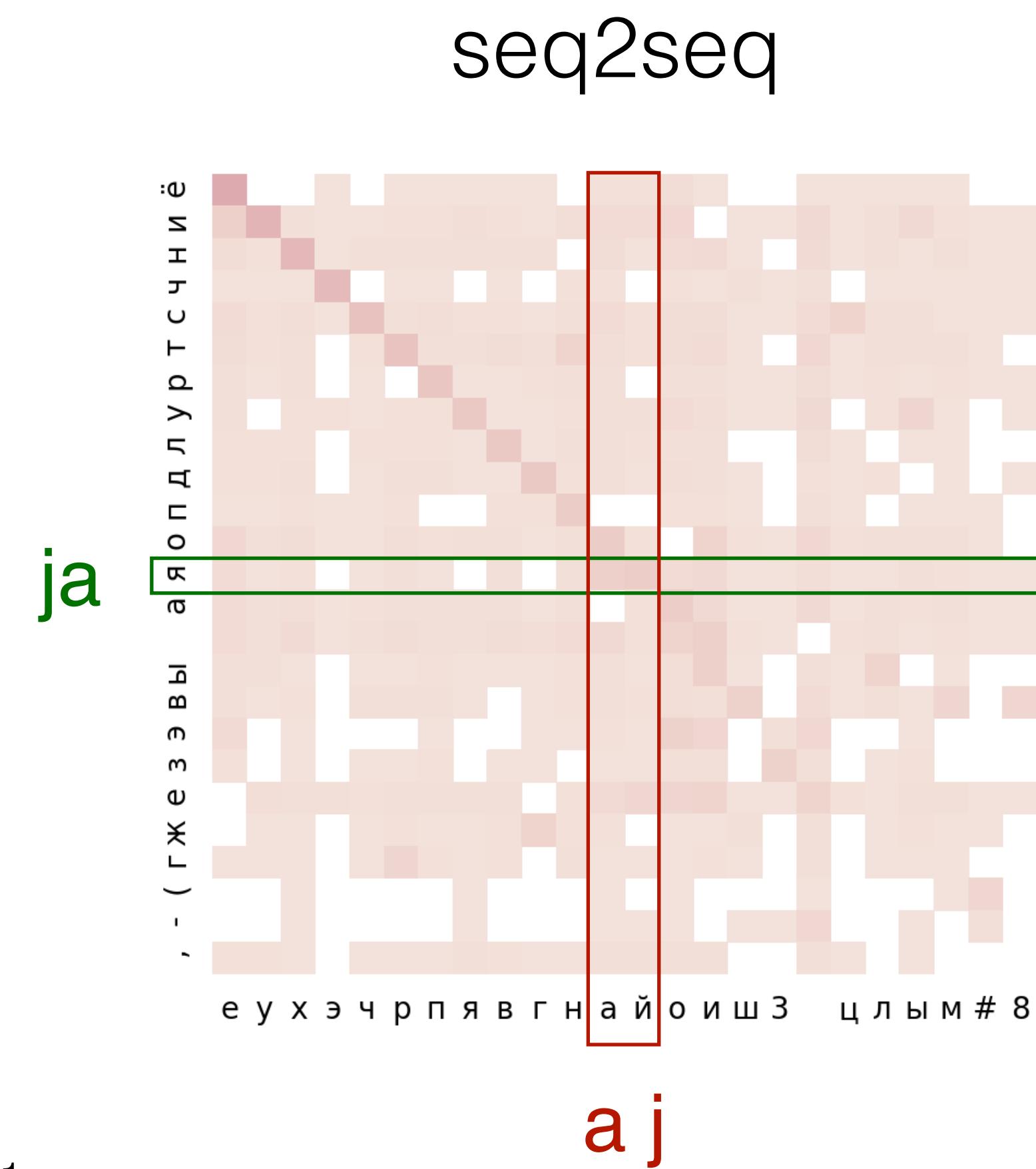
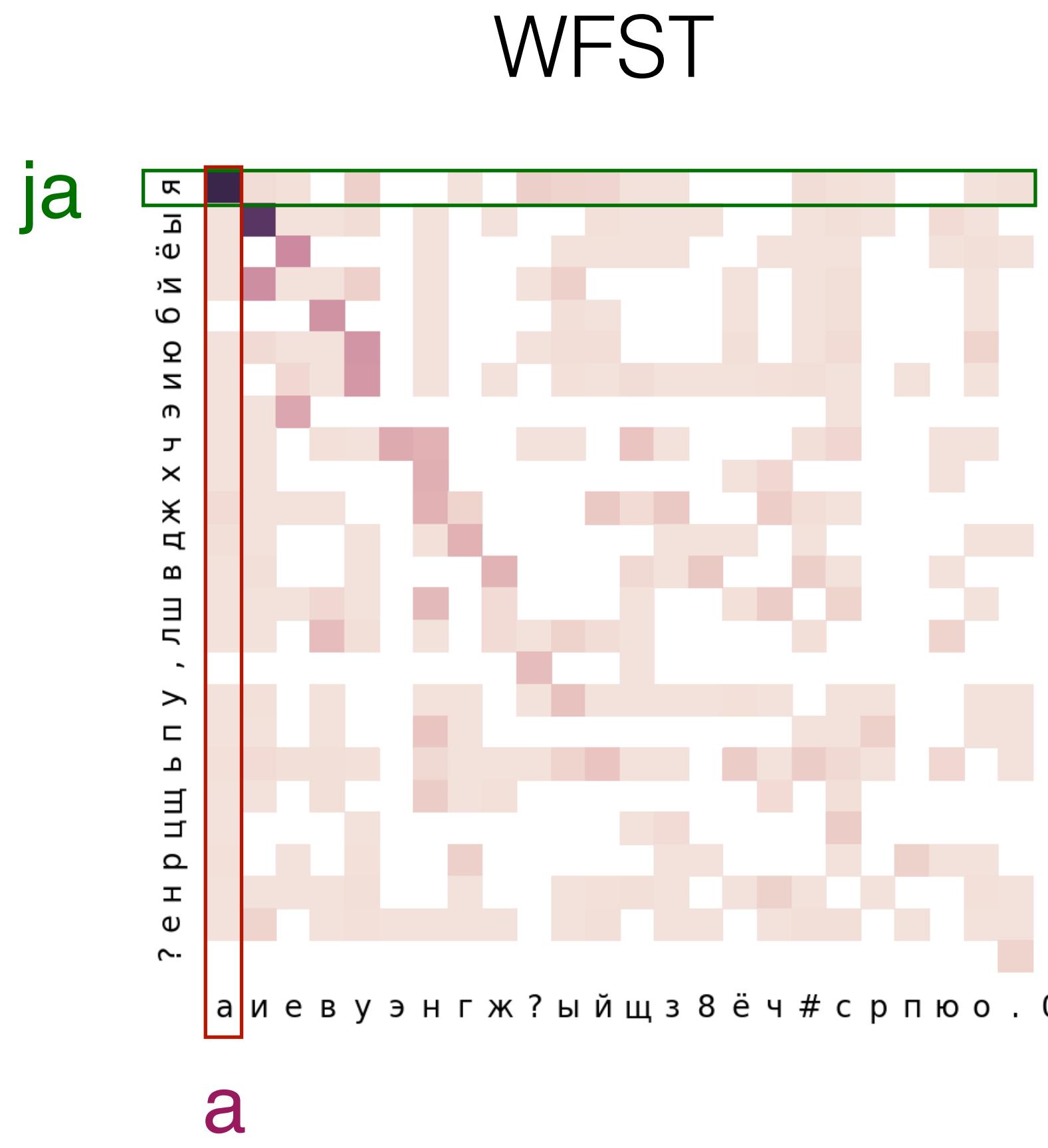
High-level takeaways

- Seq2seq is more sensitive to **distributional shifts**
 - Remember that our Cyrillic data comes from political discussion groups
 - 25% of most frequent substitutions under the seq2seq are caused by domain mismatch, compared to 3% for WFST!



High-level takeaways

- WFST makes **more repetitive errors**
 - Suggests that WFST outputs might be easier to correct with rule-based postprocessing



Future work

Combining unsupervised finite-state and neural models

- Our decoding-time combinations didn't help much, but I still believe in combining the two model classes!
- It could be joint training...
- ...Or holistic structural combinations...
 - Encoding WFST-like constraints into the attention mechanism (Aharoni & Goldberg, 2017; Makarov et al., 2017; Wu et al., 2018; Wu & Cotterell, 2019)
 - Neural reweighting of WFSTs (Rastogi et al., 2016; Lin et al., 2019)
- ...Or 'softer' biasing of one model towards another model's behavior

Future work

Combining unsupervised finite-state and neural models

- Can enough ‘power’ replace ‘structure’?
 - Transformer can learn character-level transduction without structural constraints (Wu et al., 2021)
 - But less likely to be sufficient in unsupervised or low-data settings!

Future work

Making sense of user preferences

- Users tend to have consistent substitution preferences
- These preferences can be correlated with the author's background
 - Where does the author live?



English: $B \rightarrow v$



German: $B \rightarrow w$

Future work

Making sense of user preferences

- Users tend to have consistent substitution preferences
- These preferences can be correlated with the author's background
 - Where does the author live?
 - What are other dominant languages in that area?

Former British colonies
(English)

ش → sh

Former French colonies
(French)

ش → ch

Future work

Making sense of user preferences

- Users tend to have consistent substitution preferences
- These preferences can be correlated with the author's background
 - Where does the author live?
 - What are other dominant languages in that area?
 - How old might the author be?

Older standard
(~German)

Ц→с, ю→ju

Newer standard
(~English)

Ц→ts, ю→yu

Future work

Making sense of user preferences

- Users tend to have consistent substitution preferences
- These preferences can be correlated with the author's background
 - Where does the author live?
 - What are other dominant languages in that area?
 - How old might the author be?
- Could be used to study creative language variation and change
- But also might be a privacy risk!

Final thoughts

- Typology is important!
 - One needs to carefully account for differences between languages
 - These differences extend to seemingly ‘technical’ side as well
- Rules can be better than data
 - Handcrafted converters, even simple ones, outperformed the supervised WFST
- Normalization of creative phenomena is lossy
 - Ideally, annotation should be performed with feedback from the author
 - Developing non-destructive normalization methods to preserve social meaning

References

- R. Aharoni, Y. Goldberg. Morphological inflection generation with hard monotonic attention. ACL 2017.
- C. Allauzen, M. Riley, J. Schalkwyk, W. Skut, M. Mohri. OpenFst: A general and efficient weighted finite-state transducer library. CIAA 2007.
- A. Bies, Z. Song, M. Maamouri, S. Grimes, H. Lee, J. Wright, S. Strassel, N. Habash, R. Eskander, O. Rambow. Transliteration of Arabizi into Arabic orthography: Developing a parallel annotated Arabizi-Arabic script SMS/chat corpus. ANLP 2014.
- T. Baldwin, M. C. de Marneffe, B. Han, Y.-B. Kim, A. Ritter, W. Xu. Shared tasks of the 2015 Workshop on Noisy User-generated Text: Twitter lexical normalization and named entity recognition. W-NUT 2015.
- E. M. Bender. Linguistically naïve != Language independent: Why NLP needs linguistic typology. ILCL 2009
- A. Chalamandaris, A. Protopapas, P. Tsiakoulis, S. Raptis. All Greek to me! An automatic Greeklish to Greek transliteration system. LREC 2006
- K. Darwish. Arabizi detection and conversion to Arabic. ANLP 2014.
- J. Eisner. Parameter estimation for probabilistic finite-state transducers. ACL 2002.
- N. Habash, M. Diab, O. Rambow. Conventional orthography for dialectal Arabic. LREC 2012.
- J. He, X. Wang, G. Neubig, T. Berg-Kirkpatrick. A probabilistic formulation of unsupervised text style transfer. ICLR 2020.
- K. Knight, A. Nair, N. Rathod, K. Yamada. Unsupervised analysis for decipherment problems. COLING/ACL 2006.
- G. Lample, A. Conneau, L. Denoyer, M. Ranzato. Unsupervised machine translation using monolingual corpora only. ICLR 2018.
- X. Li, P. Michel, A. Anastasopoulos, Y. Belinkov, N. Durrani, O. Firat, P. Koehn, G. Neubig, J. Pino, H. Sajjad. Findings of the first shared task on machine translation robustness. WMT 2019

References

- P. Liang, D. Klein. Online EM for unsupervised models. NAACL 2009.
- C.-C. Lin, H. Zhu, M. R. Gormley, J. Eisner. Neural finite-state transducers: Beyond rational relations. NAACL 2019.
- P. Makarov, T. Ruzsics, S. Clematide. Align and copy: UZH at SIGMORPHON 2017 shared task for morphological reinflection. SIGMORPHON 2017.
- D. Nguyen, L. Rosseel, J. Grieve. On learning and representing social meaning in NLP: A sociolinguistic perspective. NAACL 2021.
- P. Rastogi, R. Cotterell, J. Eisner. Weighting finite-state transductions with neural context. NAACL 2016.
- B. Roark, R. Sproat, C. Allauzen, M. Riley, J. Sorensen, T. Tai. The OpenGrm open-source finite-state grammar software libraries. ACL 2012.
- B. Roark, L. Wolf-Sonkin, C. Kirov, S. J. Mielke, C. Johny, I. Demirsahin, K. Hall. Processing South Asian languages written in the Latin script: The Dakshina dataset. LREC 2020
- M. Ryskina, M. R. Gormley, T. Berg-Kirkpatrick. Phonetic and visual priors for decipherment of informal romanization. ACL 2020.
- M. Ryskina, E. Hovy, T. Berg-Kirkpatrick, M. R. Gormley. Comparative error analysis in neural and finite-state models for unsupervised character-level transduction. To appear at SIGMORPHON 2021.
- T. Shavrina, O. Shapovalova. To the methodology of corpus construction for machine learning: Taiga syntax tree corpus and parser. CORPORA 2017.
- S. Wu, P. Shapiro, R. Cotterell. Hard non-monotonic attention for character-level transduction. EMNLP 2018.
- S. Wu, R. Cotterell. Exact hard monotonic attention for character-level transduction. ACL 2019.
- S. Wu, R. Cotterell, M. Hulden. Applying the Transformer to character-level transduction. EACL 2021.