

# Unsupervised Decipherment of Informal Romanization

Maria Ryskina



Carnegie Mellon University  
Language  
Technologies  
Institute

Carnegie Mellon University  
Language Technologies Institute

# 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

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

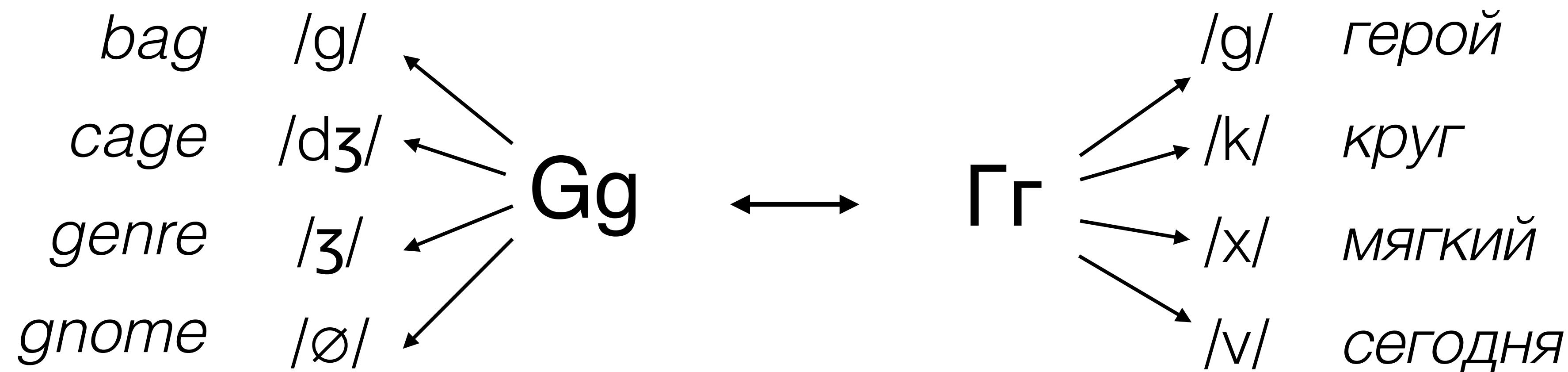
# Informal romanization

- Idiosyncratic representation: character substitutions up to the user
- Most substitutions are based on **phonetic** or **visual** similarity

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

# Character alignment

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

Alphabet

хорошо

|||||

xorosho

~ one-to-one

Abjad  
(consonantal)

کریم

krym

/|\\|

kareem

~ one-to-one + null

Abugida  
(alphasyllabary)

బెలగితు

Unicode: బ ల గ త ట ఱ

\|\|/\|/\|/

belagitu

~ one-to-one + 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

latent  
(what they meant)

observed  
(what they typed)

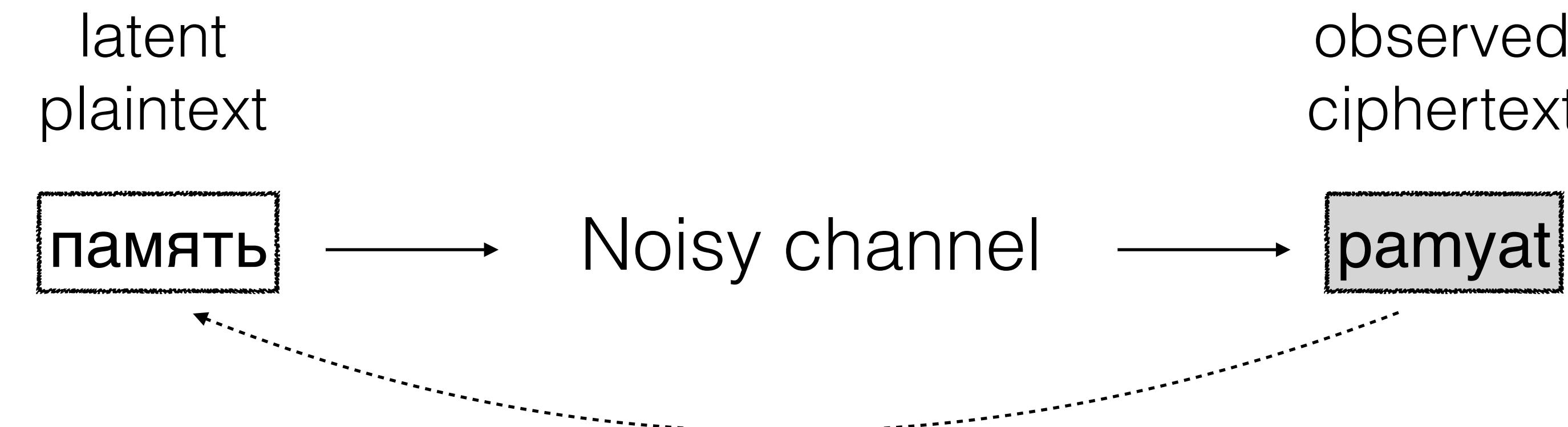
# Task framing

- Parallel data does not occur naturally ⇒ **unsupervised** learning
- Perceptions of similarity are shared across users and even languages!
- **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)



# Noisy-channel model

latent  $n = \text{п а м я т ъ}$

observed  $r = \text{p а m y a t}$

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

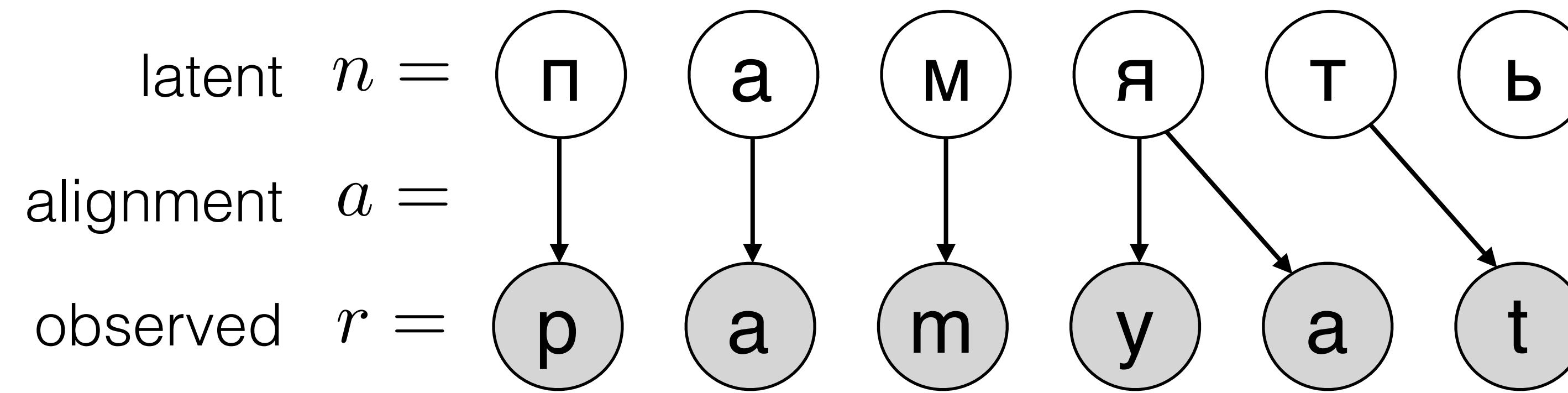
all possible  
native script  
sequences

$n$   
transition probabilities

emission probabilities

$\theta$   
prior on parameters

# Noisy-channel model



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

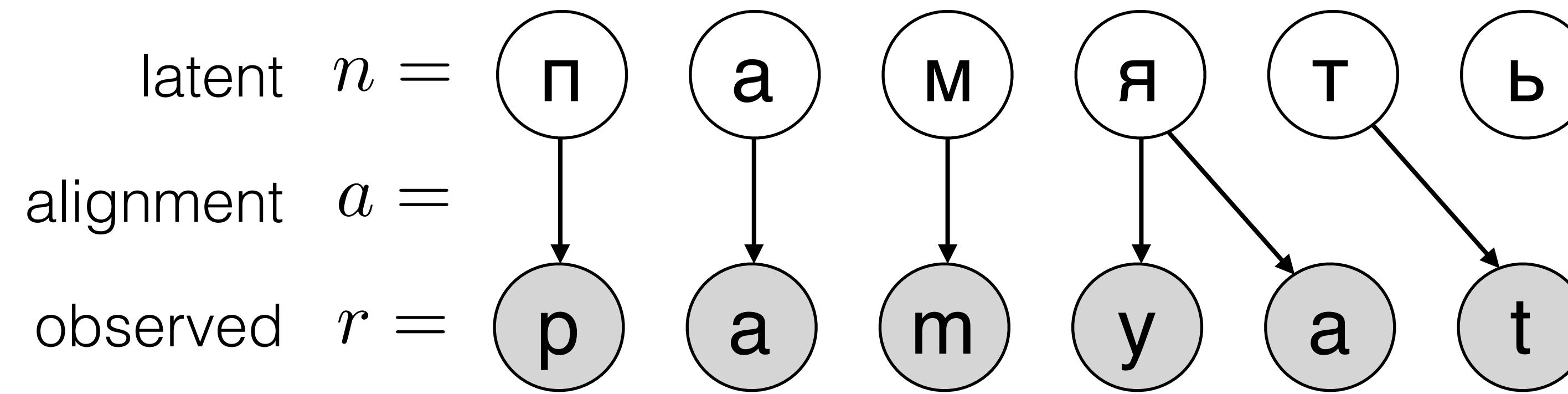
all possible  
native script  
sequences and  
alignments

$n, a$   
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**

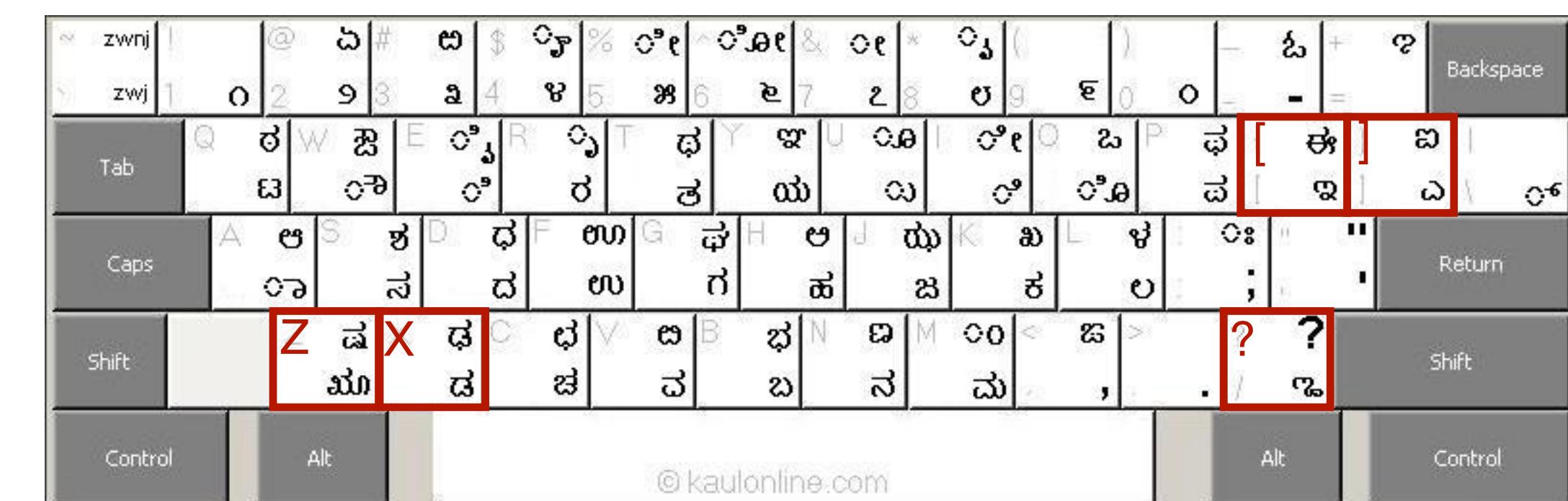


[https://en.wikipedia.org/wiki/Phonetic\\_keyboard\\_layout](https://en.wikipedia.org/wiki/Phonetic_keyboard_layout),

17 <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



# Visual bias

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

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	ρ	϶	پ	ϙ	ϙ	P	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!

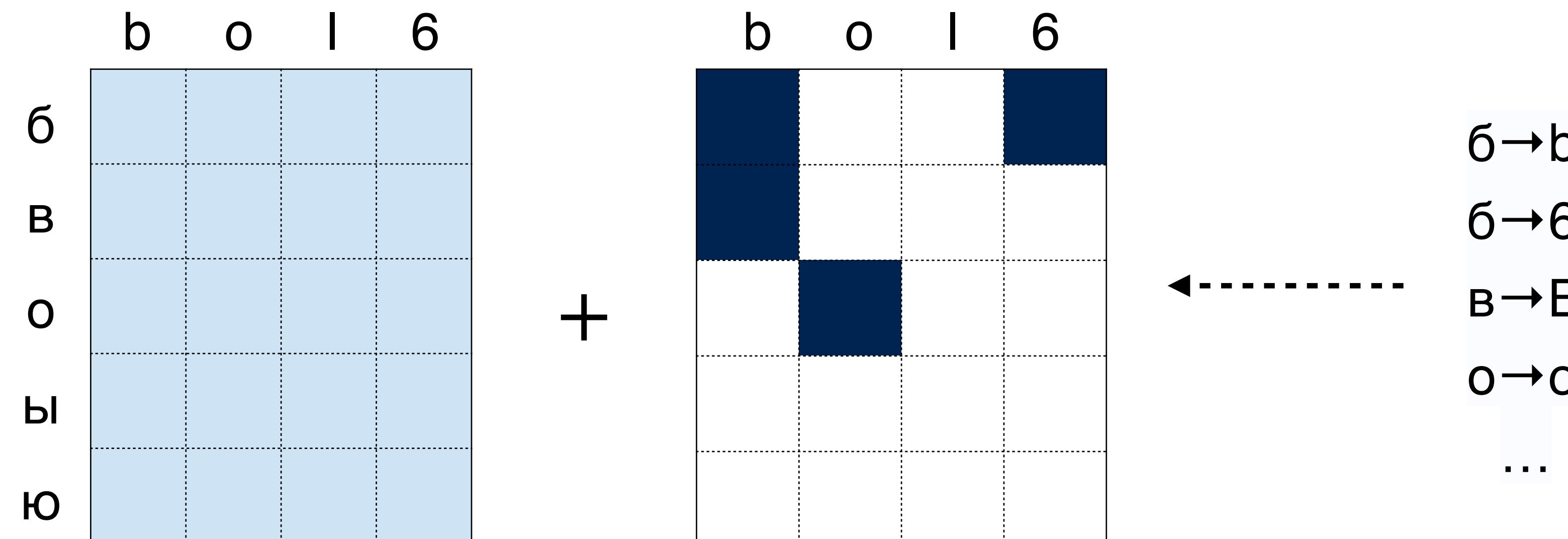
y	ȝ	Y	ȝ	y	ȝ	Y	y
0079 LATIN SMALL LETTER Y	0263 LATIN SMALL LETTER GAMMA	028F LATIN LETTER SMALL CAPITAL Y	03B3 GREEK SMALL LETTER GAMMA	0443 CYRILLIC SMALL LETTER U	04AF CYRILLIC SMALL LETTER STRAIGHT U	10E7 GEORGIAN LETTER QAR	
p	ƿ	ƿ	ƿ	ƿ	ƿ	P	p
0070 LATIN SMALL LETTER P	03C1 GREEK SMALL LETTER RHO	03F1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

# Informative priors

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

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

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

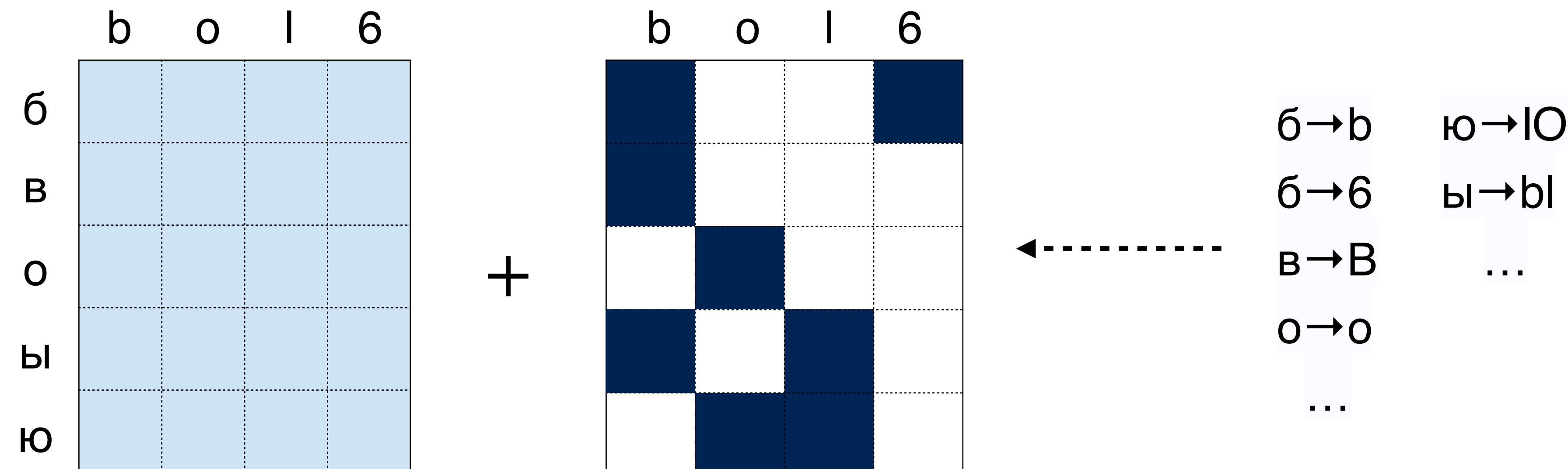


# Informative priors

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

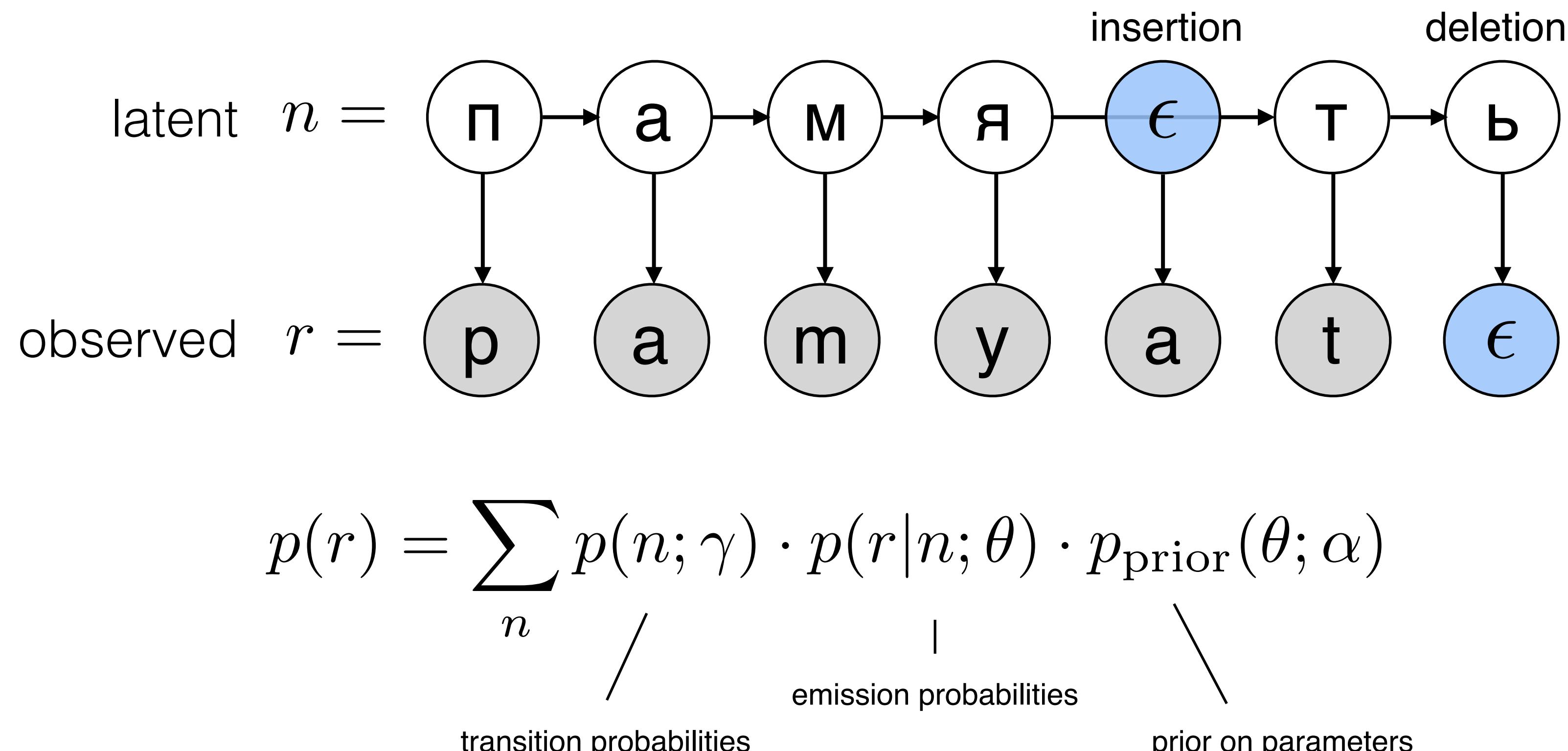
$$c_r | c_n \sim \text{Mult}(\theta_{c_n})$$

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



# Noisy-channel model

- Representing latent alignments via **insertions and deletions**

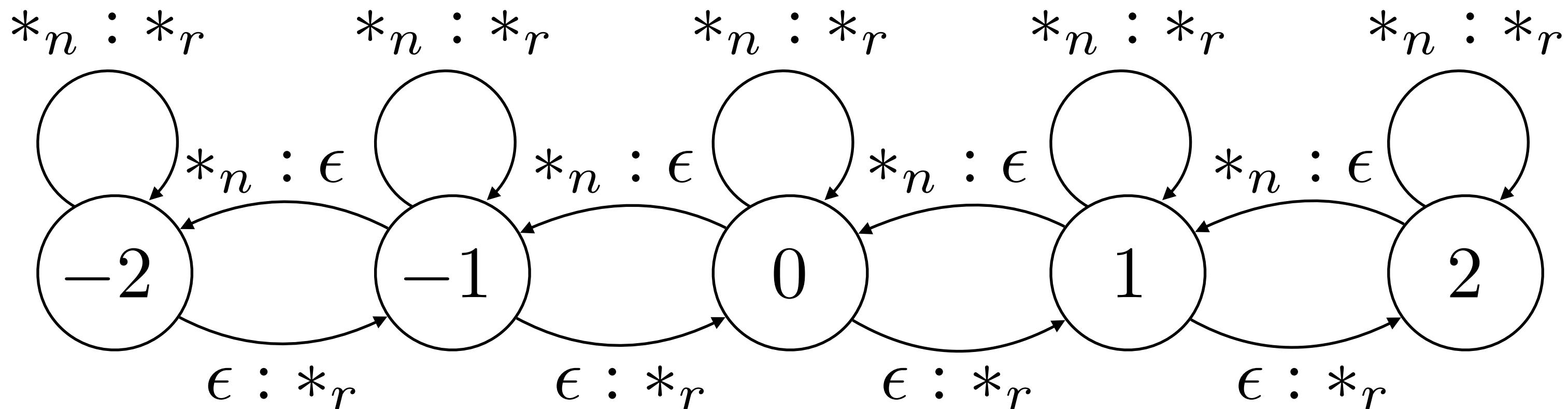
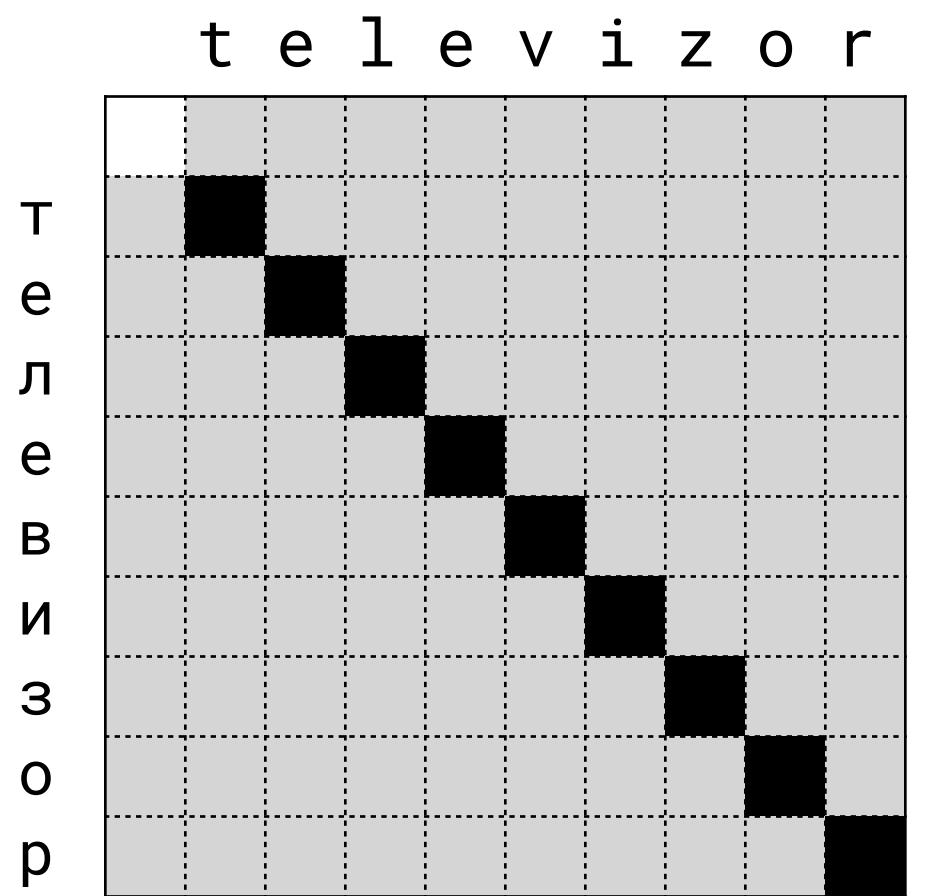


# WFST cascade

- Transition WFSA
  - 6-gram LM built with OpenGrm (Roark et al., 2012)
- Emission WFST
  - Supports all substitutions, insertions and deletions

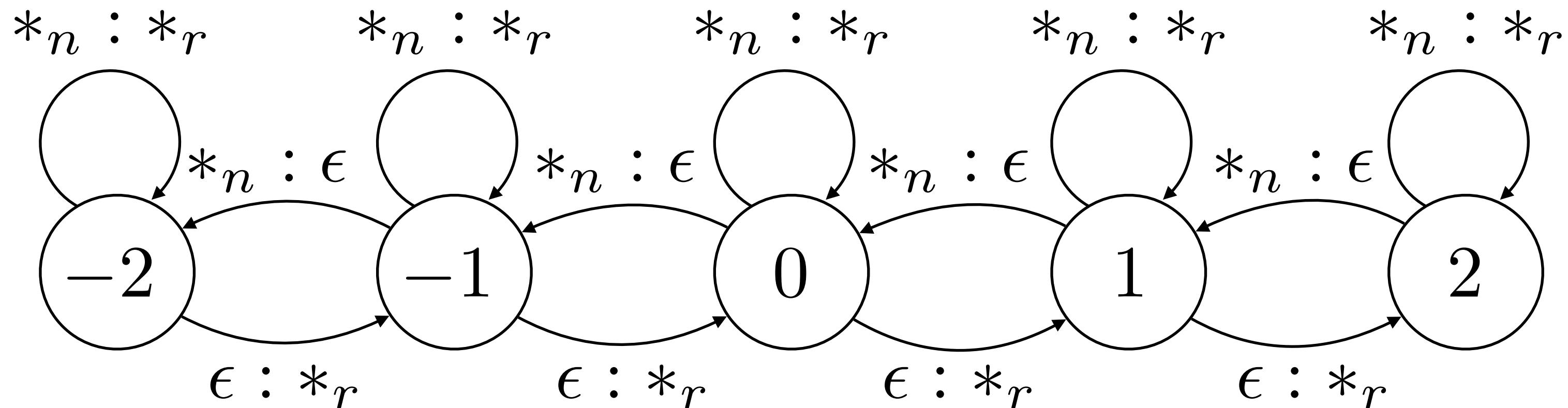
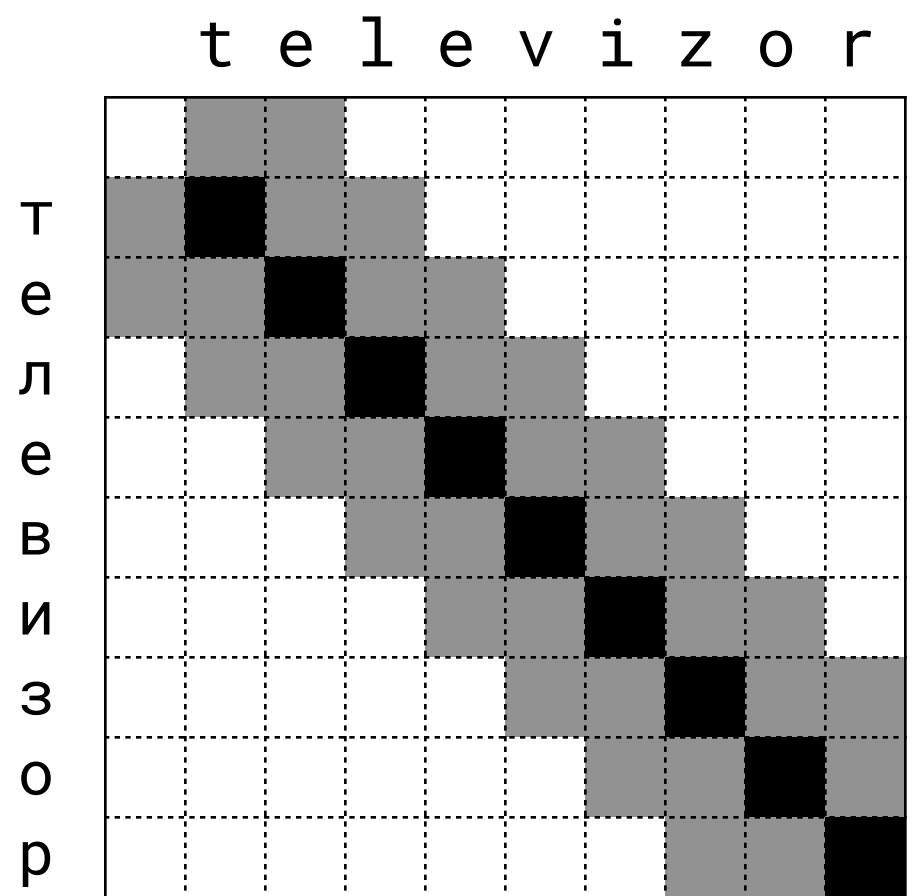
# Emission model

- 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:  $| \# \text{ of insertions} - \# \text{ of deletions} |$



# WFST cascade

- Transition WFSA
  - 6-gram LM built with OpenGrm (Roark et al., 2012)
- Emission WFST
  - Supports all substitutions, insertions and deletions
- Trained with EM algorithm
  - OpenFst (Allauzen et al., 2007)
  - Speedup tricks: stepwise training, curriculum learning, pruning...

# Datasets

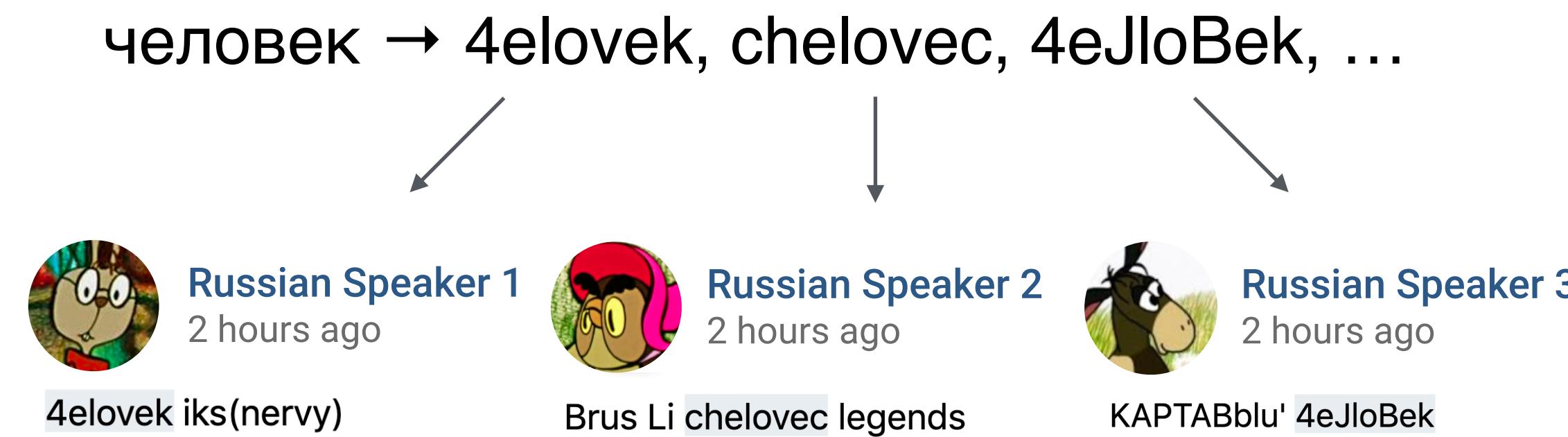
- Arabic: LDC BOLT dataset (Bies et al., 2014)
  - Arabizi SMS/chat dialogs, converted to CODA (Habash et al., 2012)
- Kannada: Dakshina dataset (Roark et al., 2020)
  - Kannada Wikipedia, romanizations elicited from native speakers
- Russian:
  - Romanized: collected and partly annotated data from social 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)



- Annotated validation and test with minor error correction

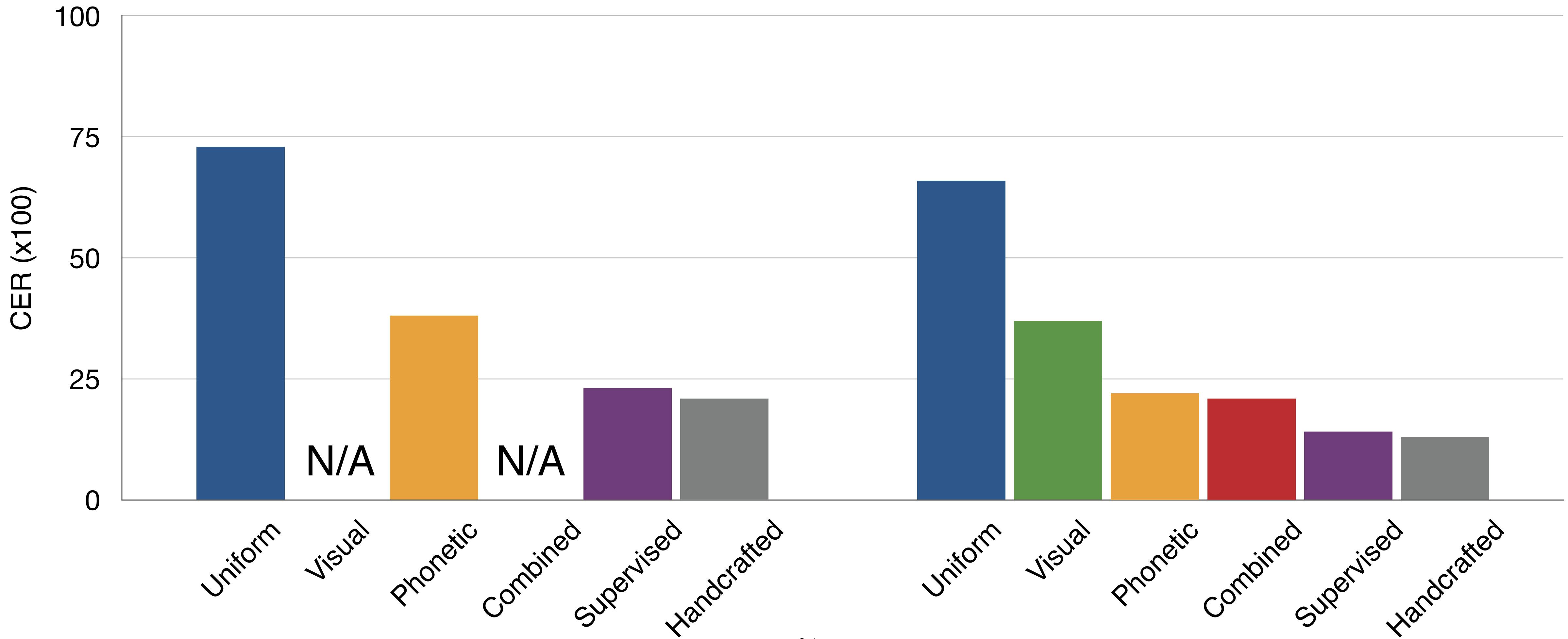
Source: *proishodit s prirodoy 4to to very very bad*

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

Target: *происходит с природой что-то <...>*

# WFST results

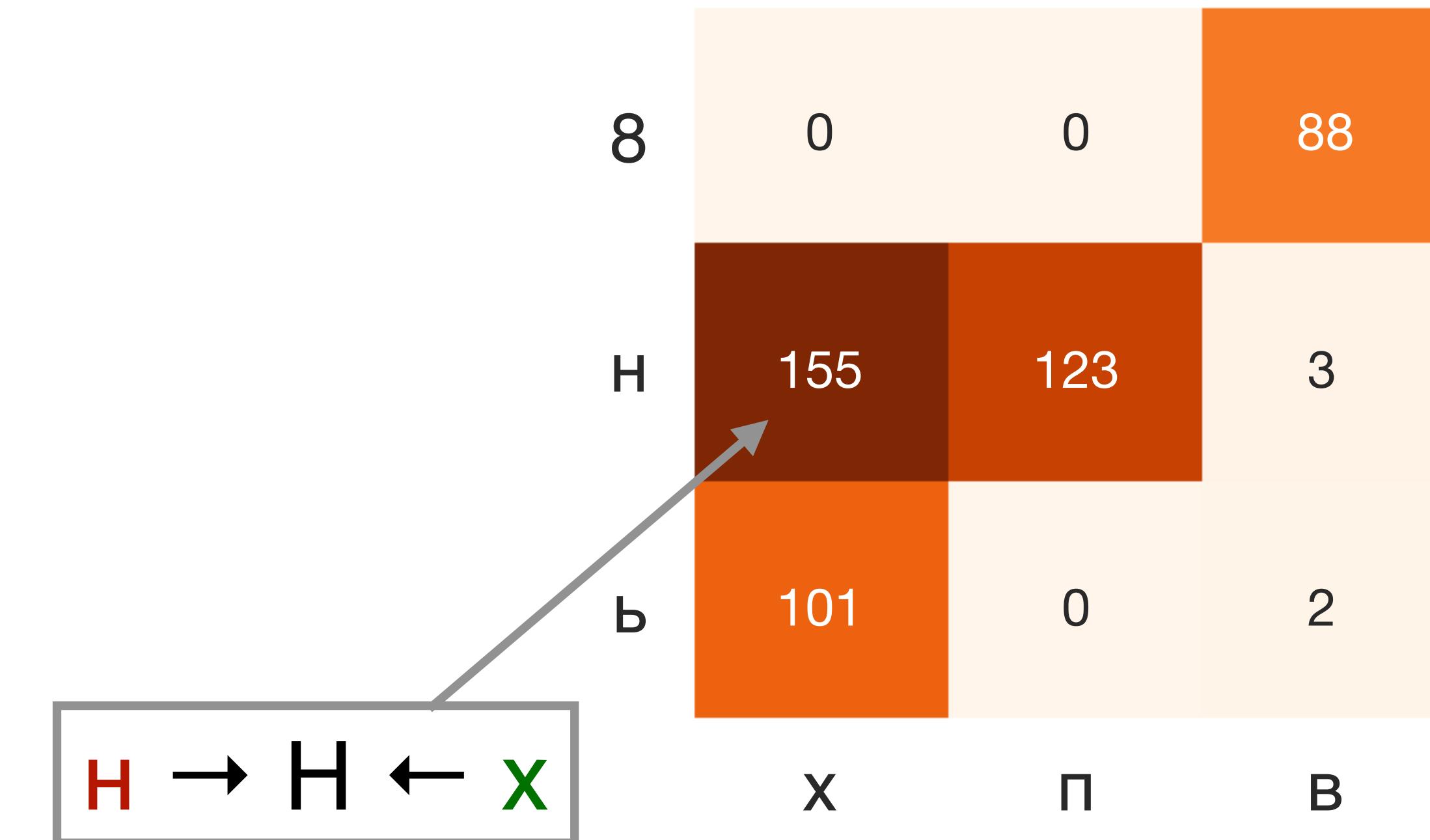
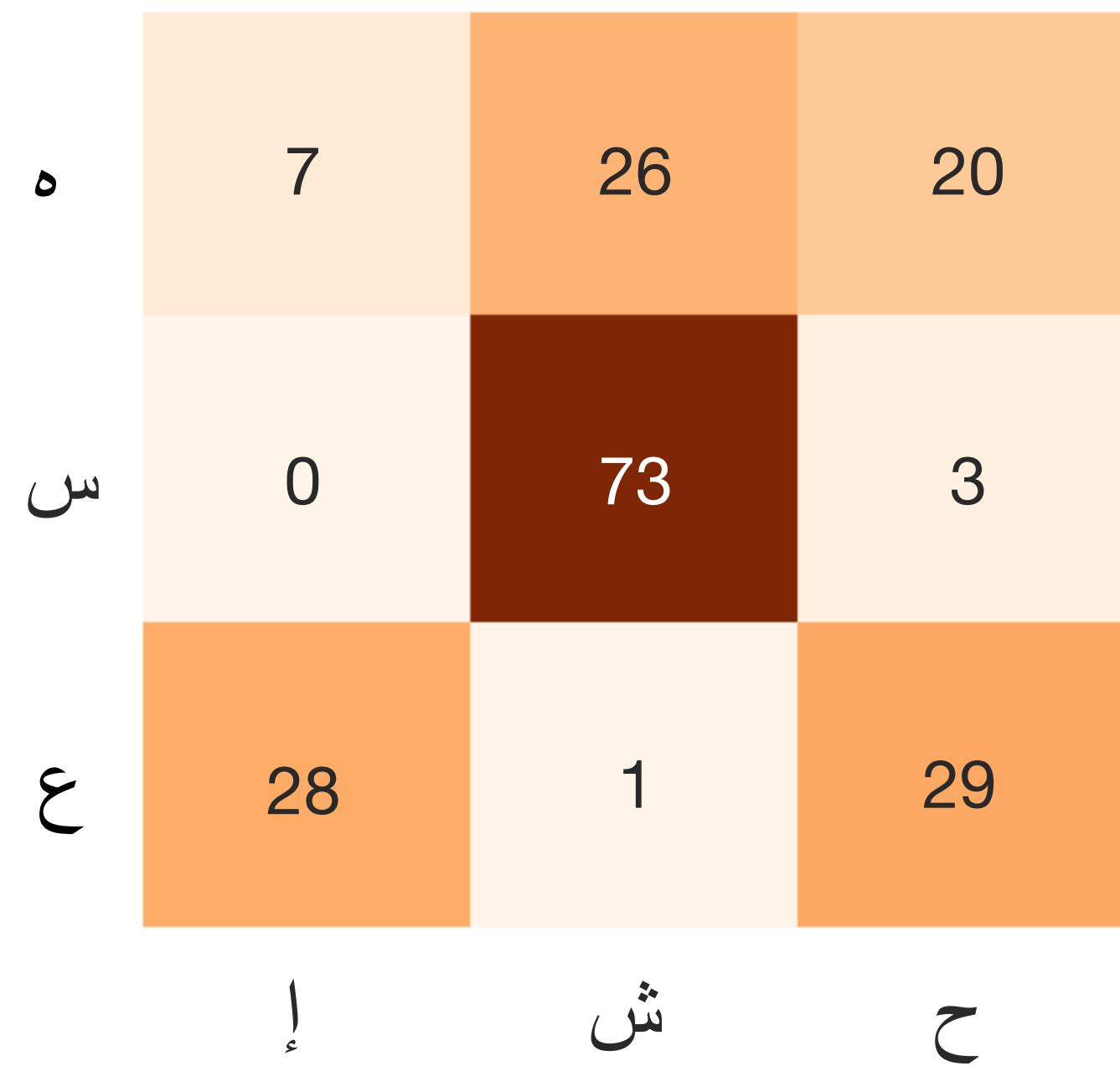
Arabic



Russian

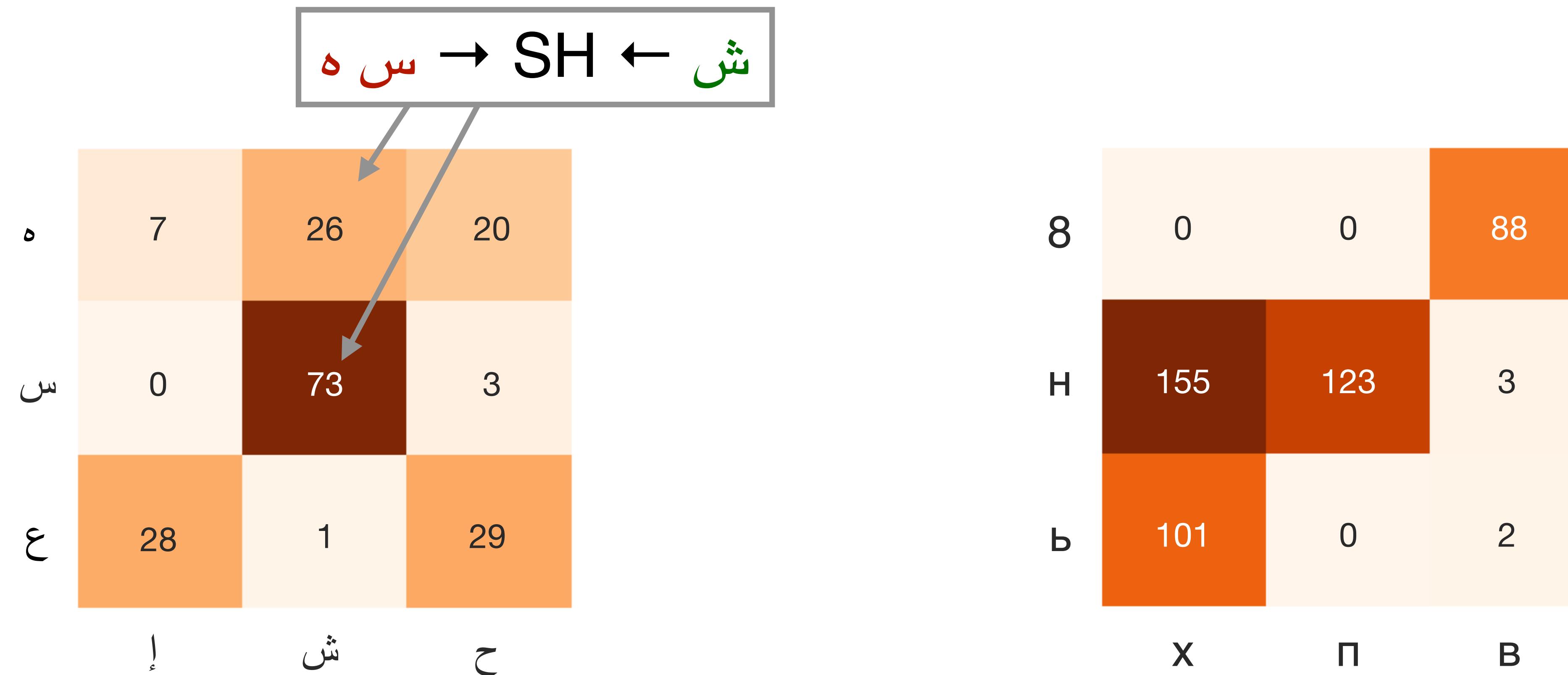
# WFST error analysis

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



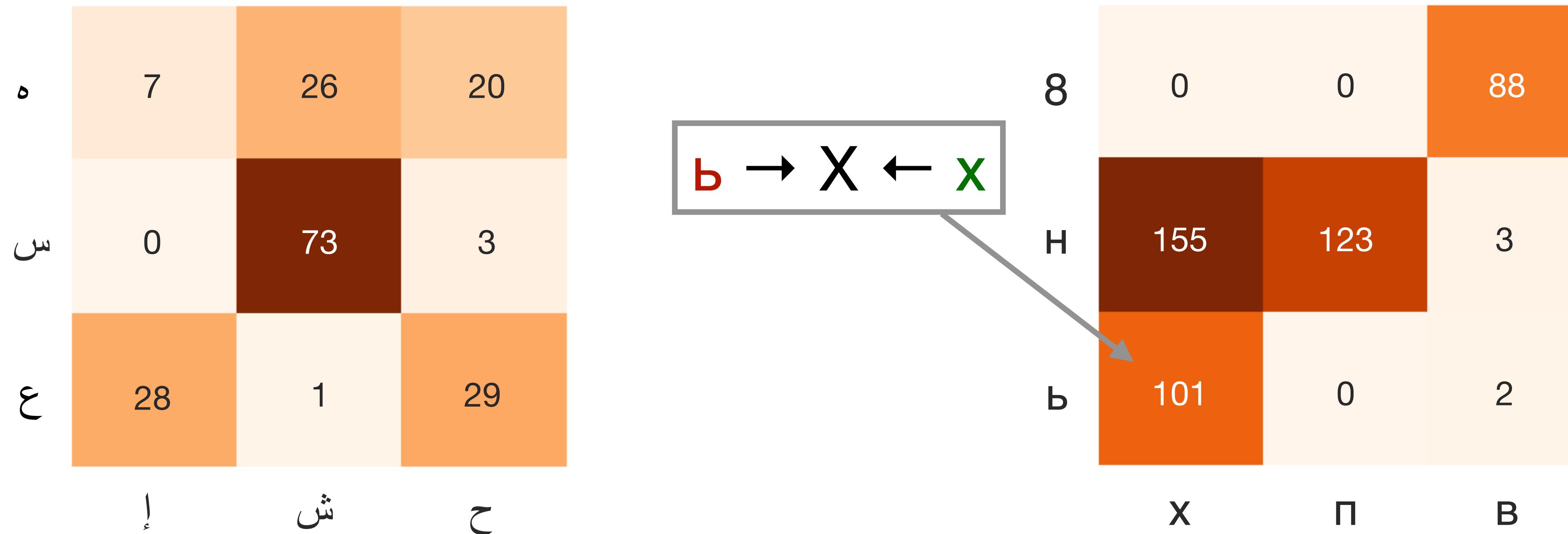
# WFST error analysis

- 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**

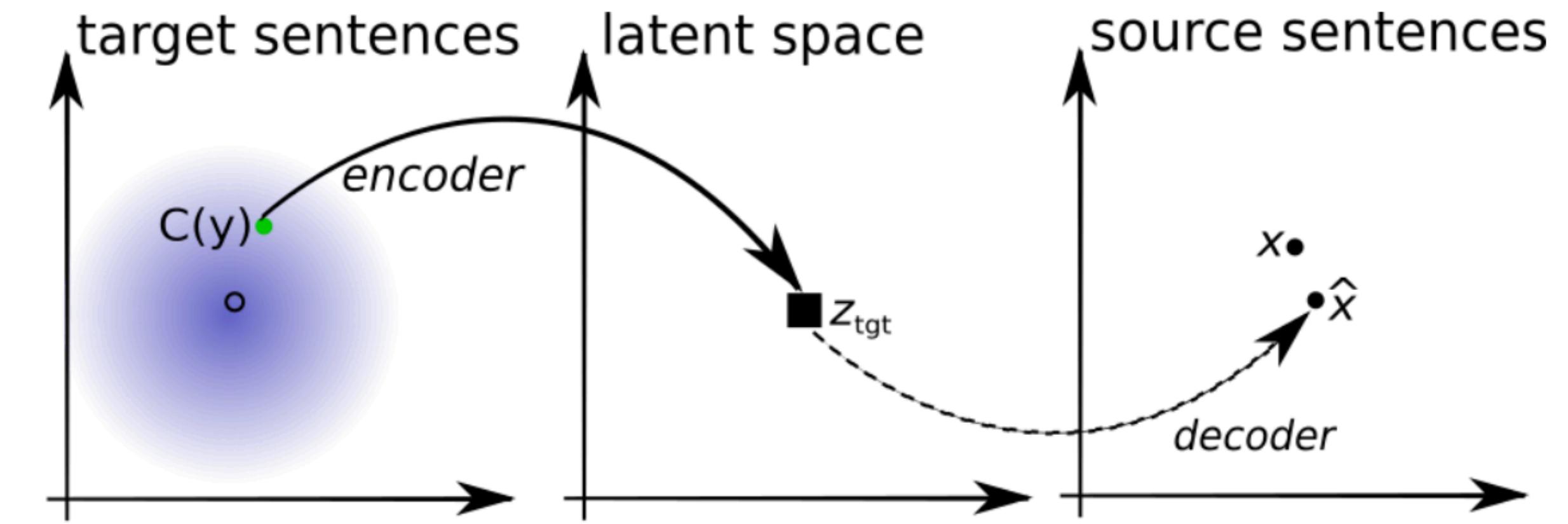
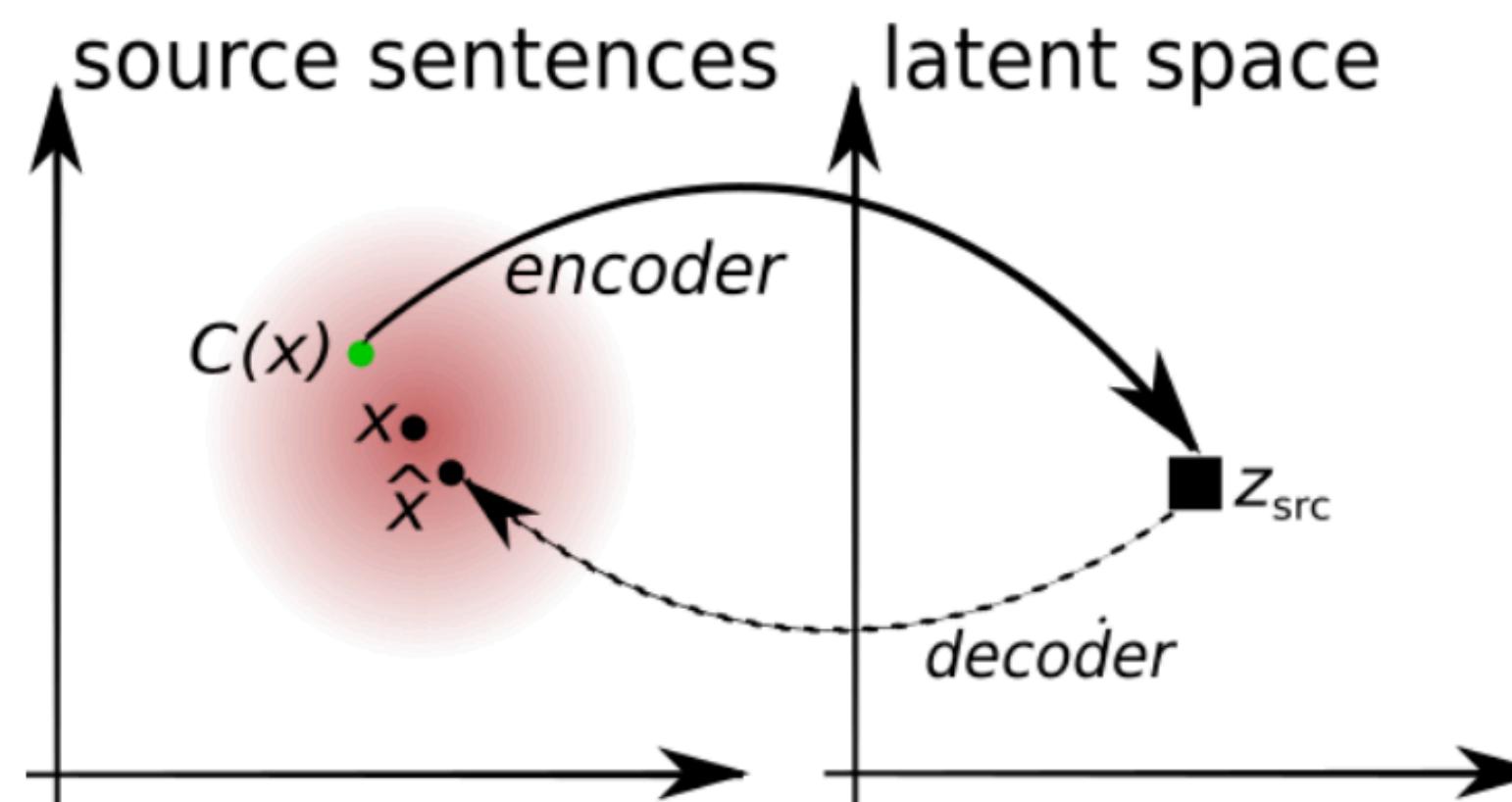
- ✓ Strong language model
- ✓ Faster batch processing
- ✗ Need large training data
- ✗ Hallucinations and search errors

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.

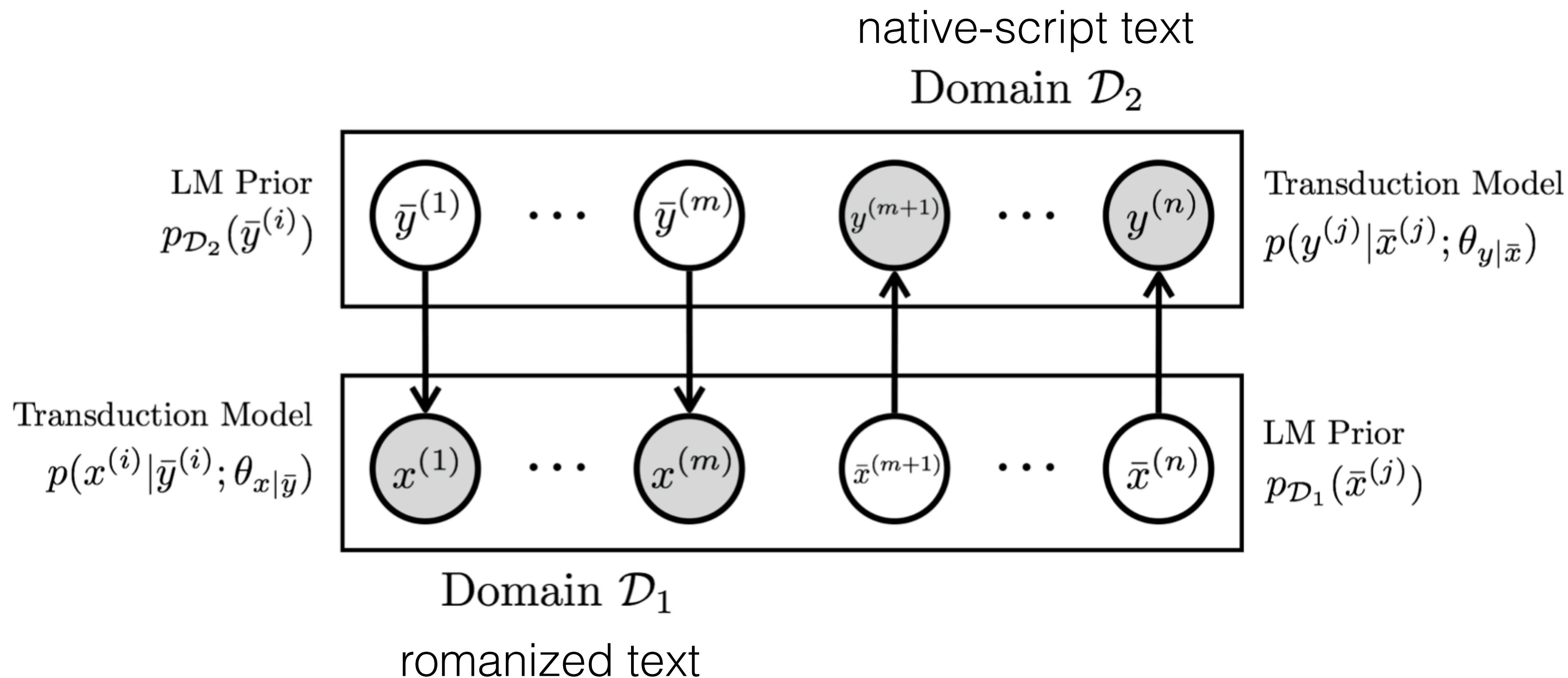
# Unsupervised seq2seq

- 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)

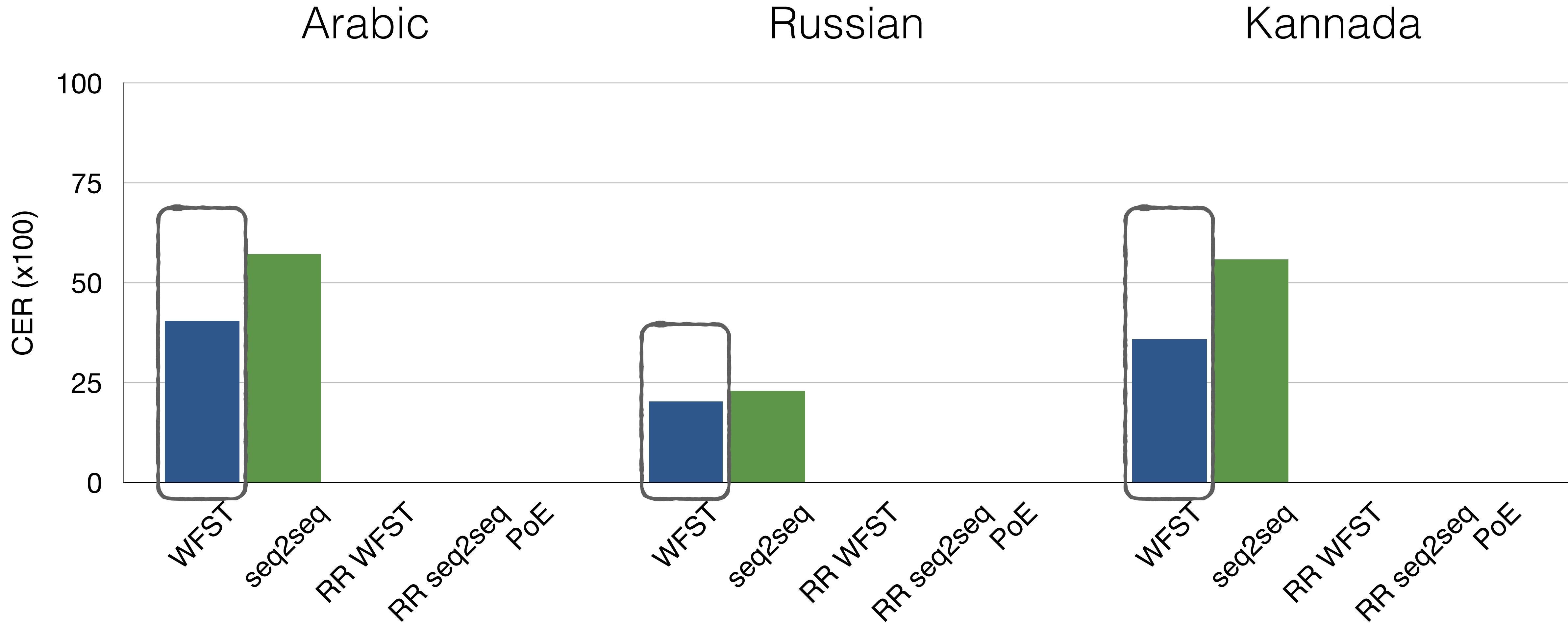


# Model combinations

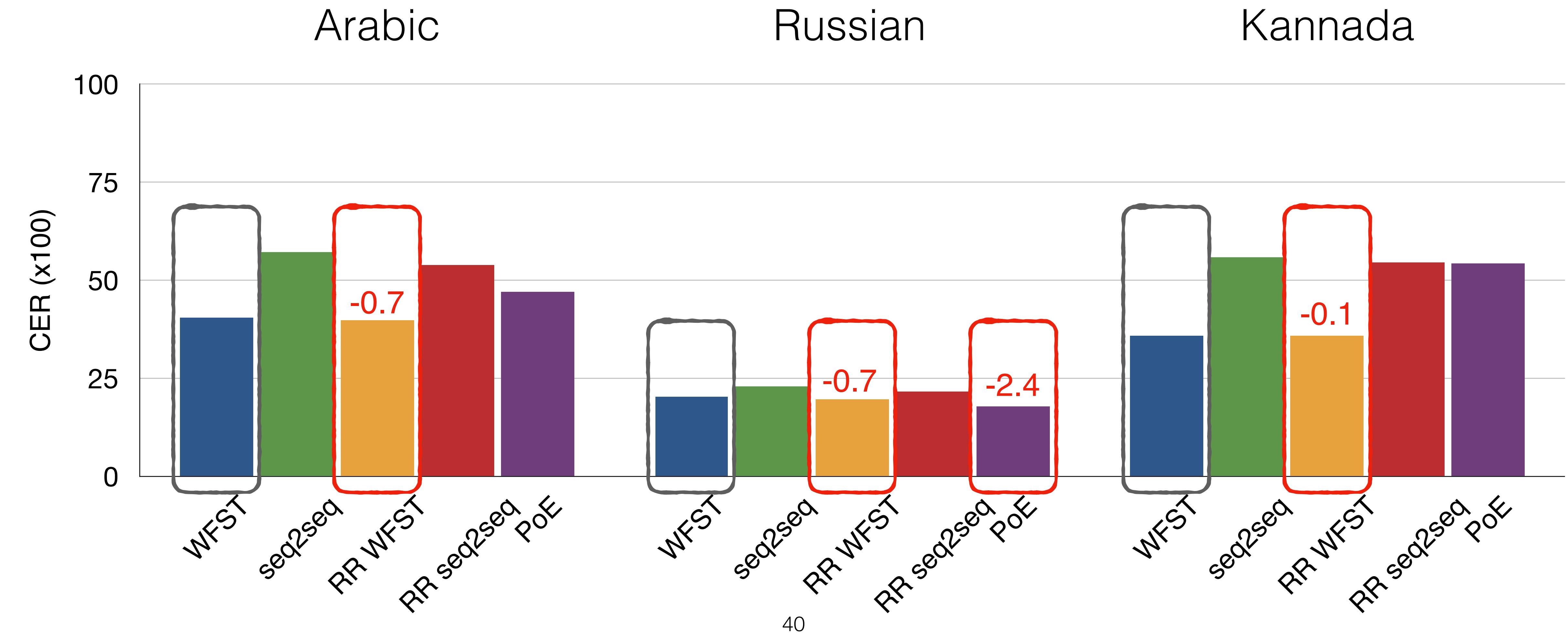
- Reranking
  - M1 generates top k candidate outputs
  - M2 selects the highest-scoring candidate
- Product of experts
  - Beam search on the WFST lattice
  - WFST arcs reweighted with Seq2seq softmax at the corresponding timestep
  - Deletions of input characters are not reweighted
  - Candidates are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

# Results

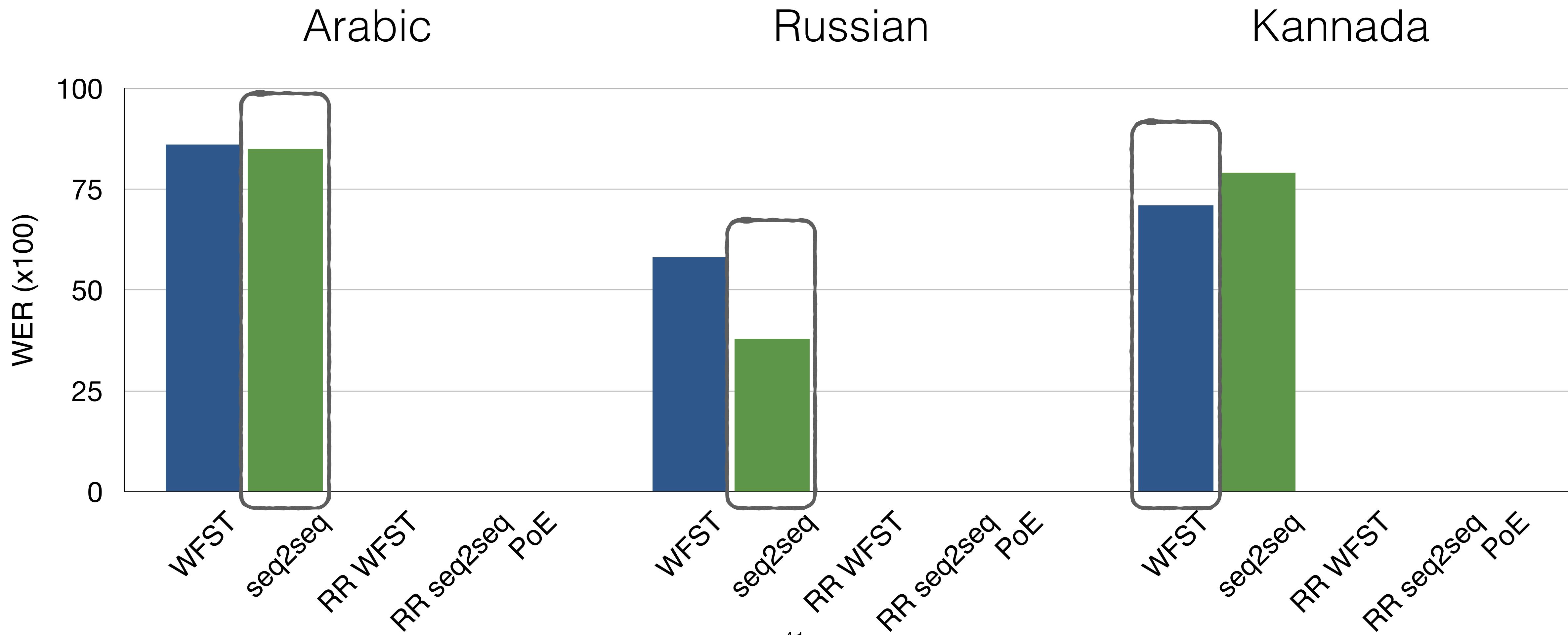
Base models are trained on different amounts of data!



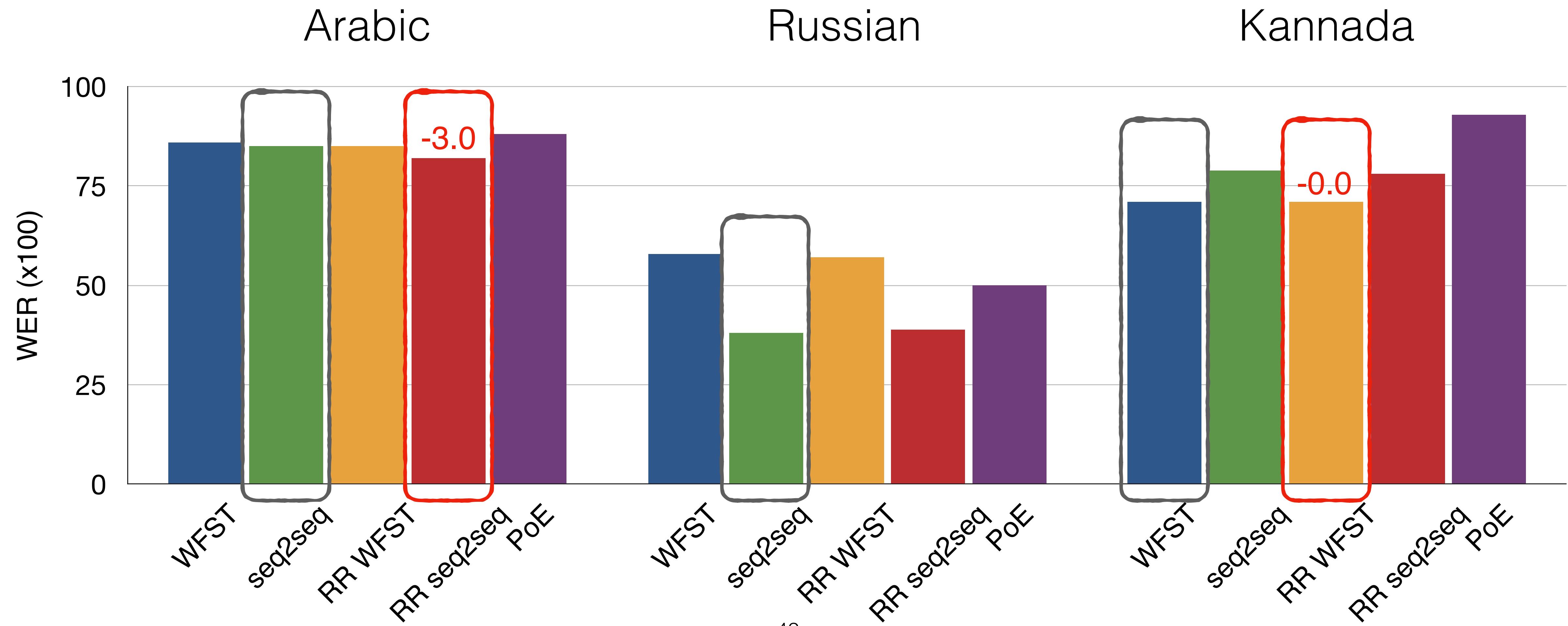
# Results



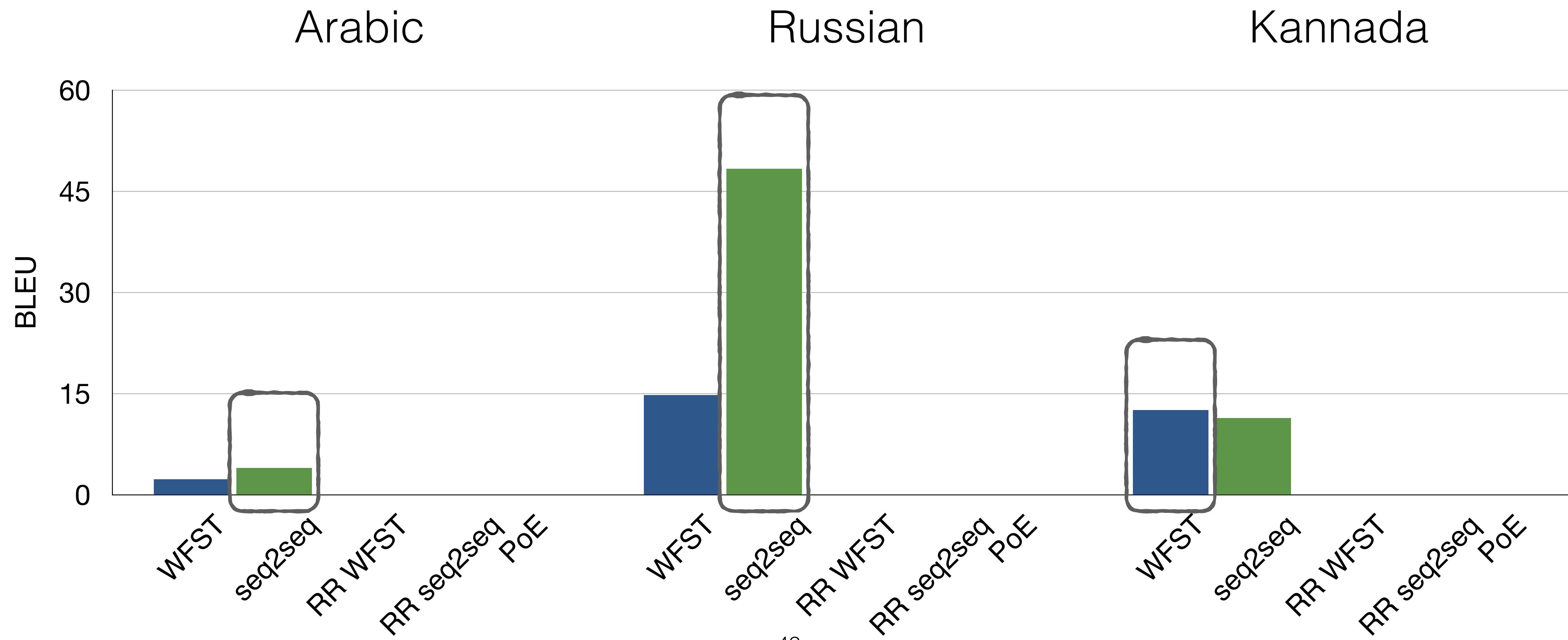
# Results



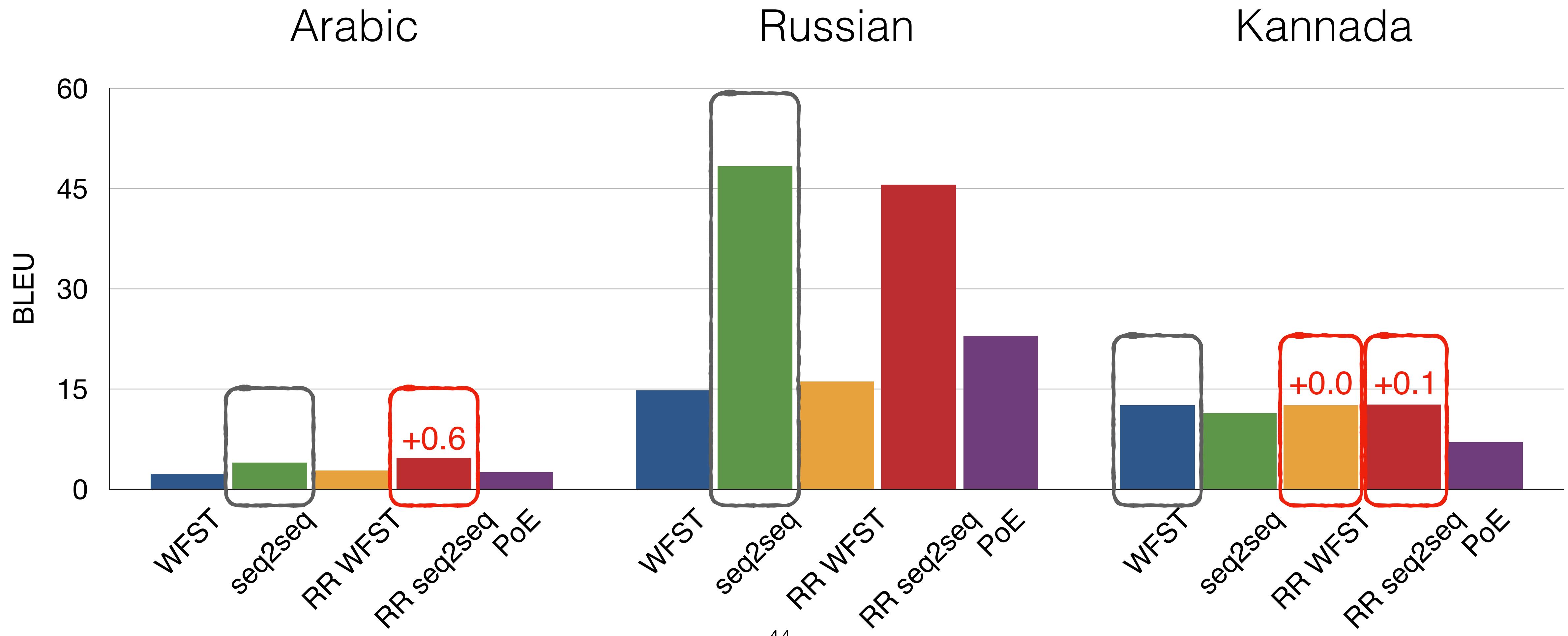
# Results



# 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**

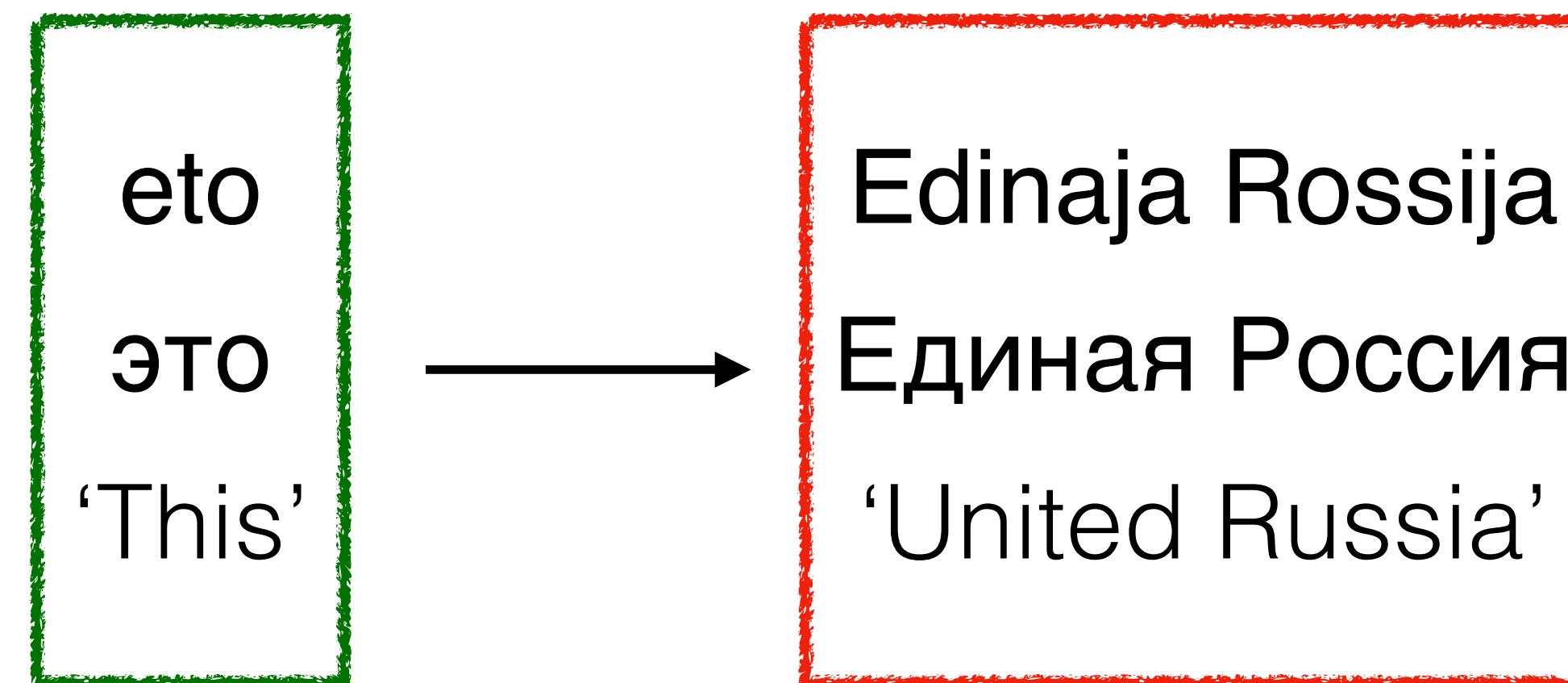
482.0, **единороссы уже #страшно сказать) старый ро**

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

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

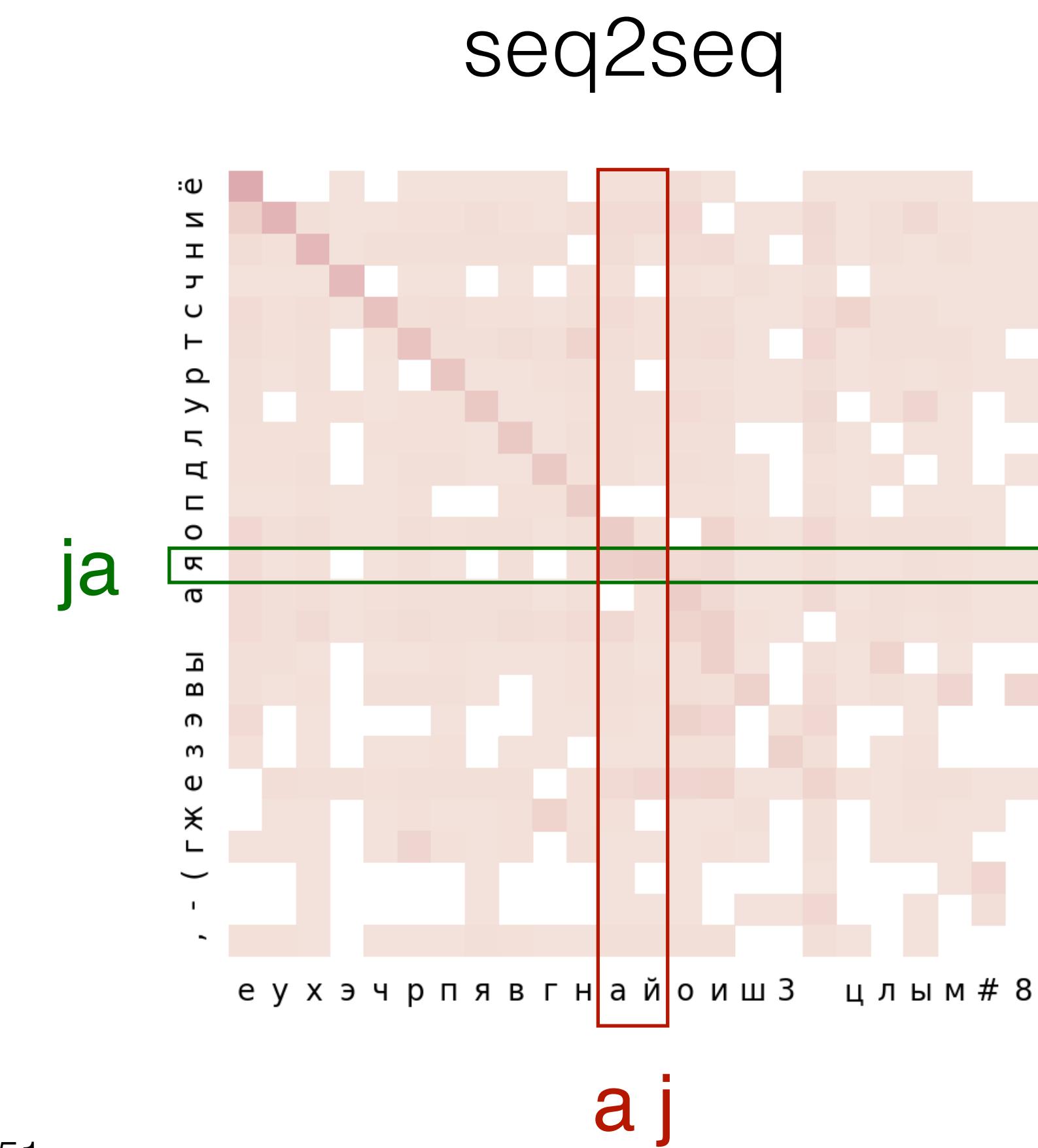
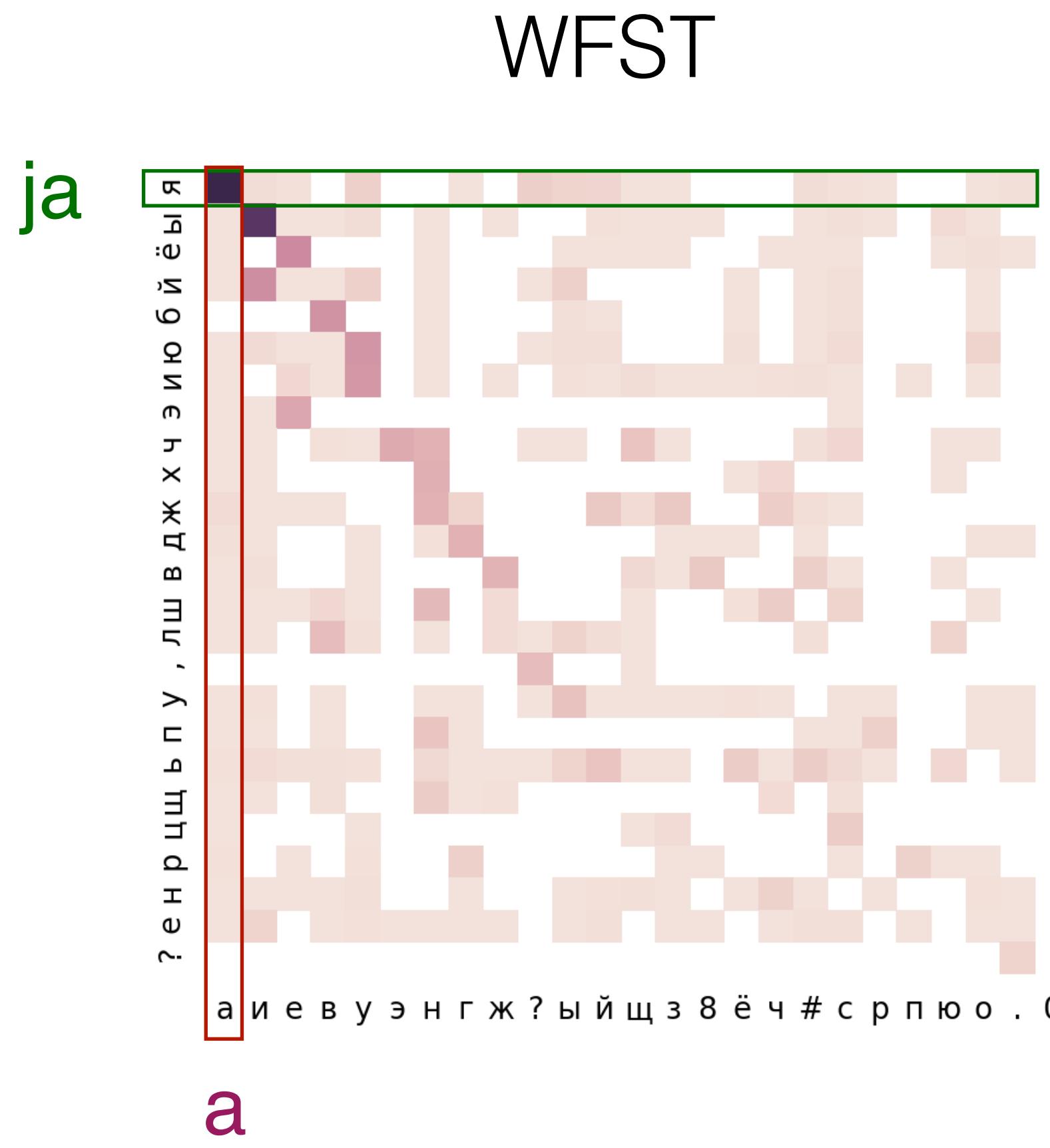
# High-level takeaways

- Seq2seq is more sensitive to **distributional shifts**
  - Remember that our Cyrillic data comes from political discussion groups
  - 25% of common word-level errors in seq2seq are of this type!



# 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 at training time
  - Joint training of two separate models
  - Holistic structural combinations
- Analysis of user preferences
  - Users tend to be consistent in their preferences
  - Substitution choices can reveal user background
  - Might want to obscure these attributes for privacy reasons
- Readability analysis instead of automated metrics

# References

- 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.
- K. Darwish. [Arabizi detection and conversion to Arabic](#). ANLP 2014.
- 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.
- 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.