

Romanization with Friends: Deciphering Informally Romanized Text

Maria Ryskina

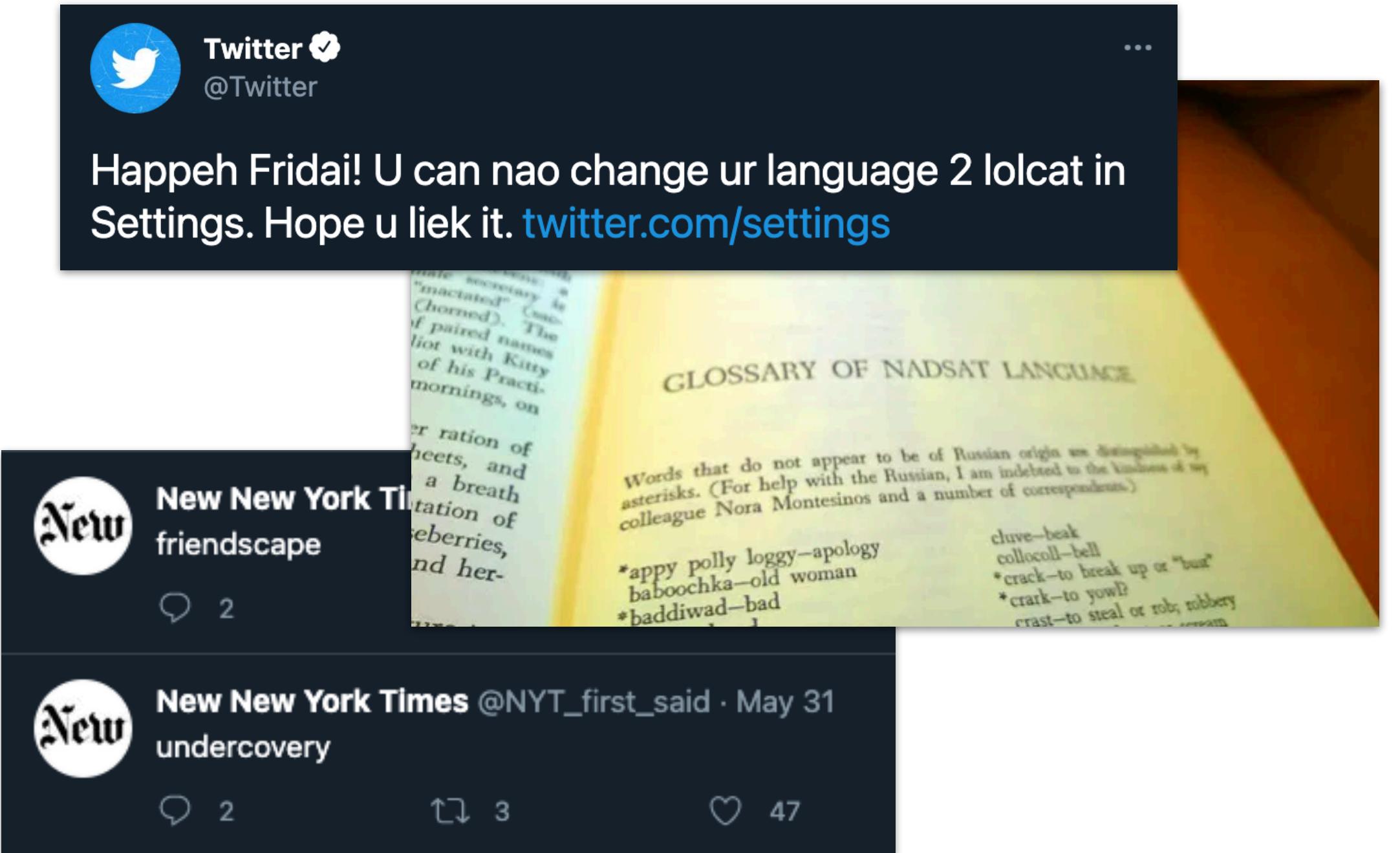


Carnegie Mellon University
Language
Technologies
Institute

Carnegie Mellon University
Language Technologies Institute

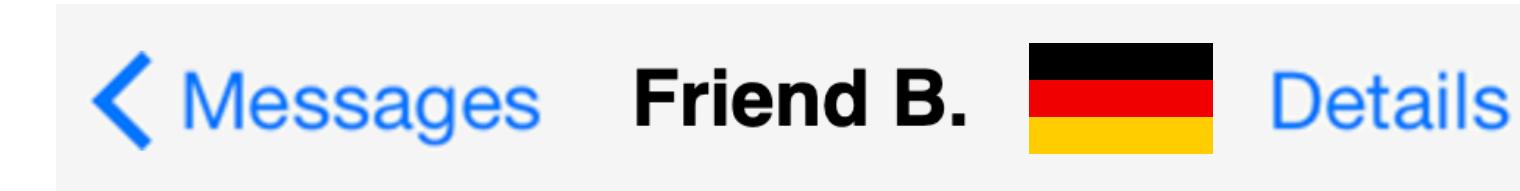
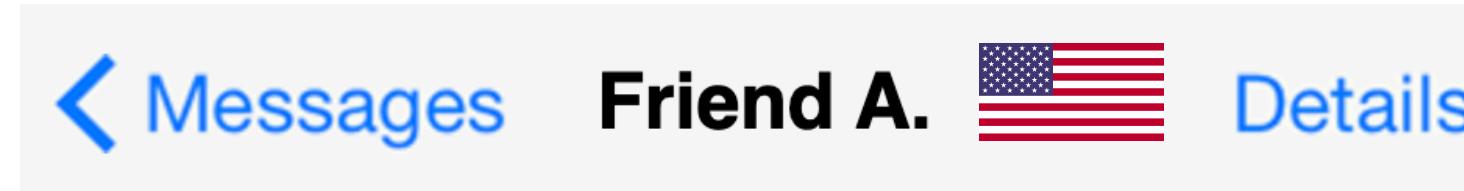
Non-standard language & NLP

- Some domains of text contain non-standard language:
 - New word forms (*brony*)
 - New morphemes (-gate) or derivatives (*prolifeness*)
 - Non-standard spellings (*2nite*)

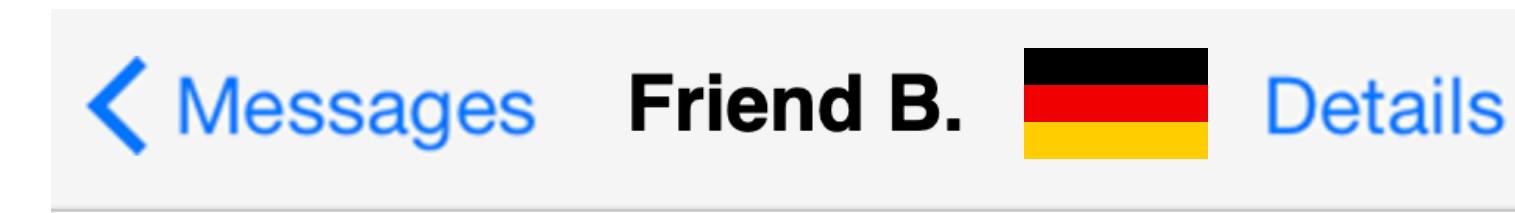
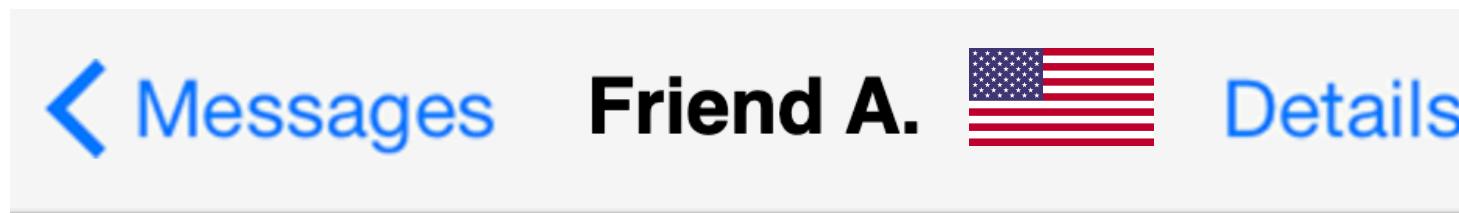


- People can infer their meaning even when seeing them for the first time
- Can we teach our NLP systems to do the same?

Romanization with friends

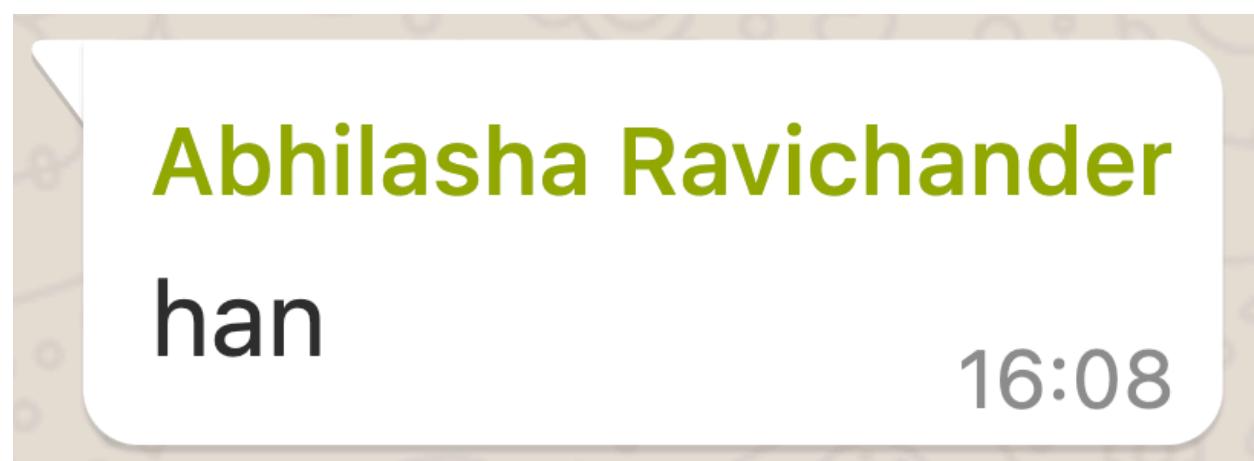
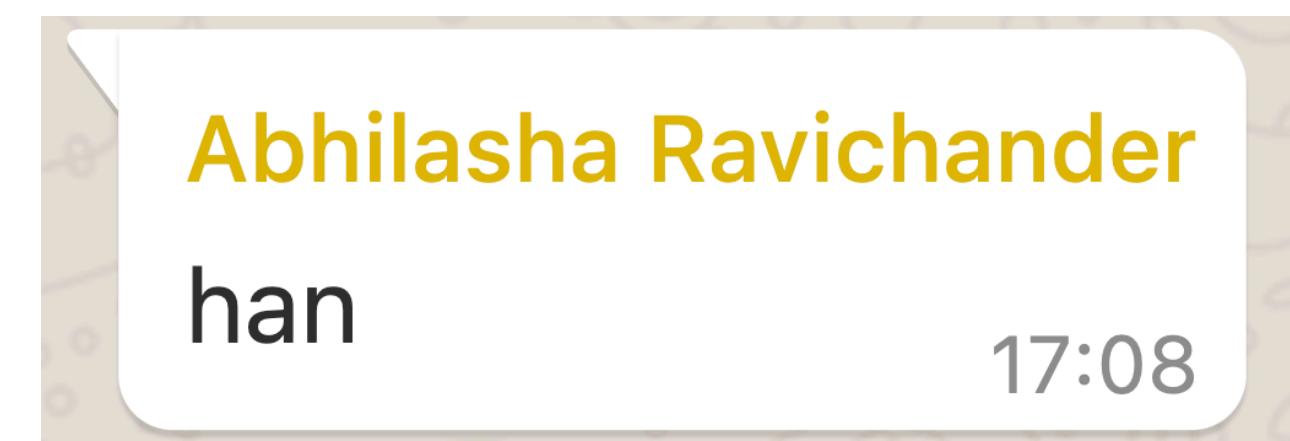
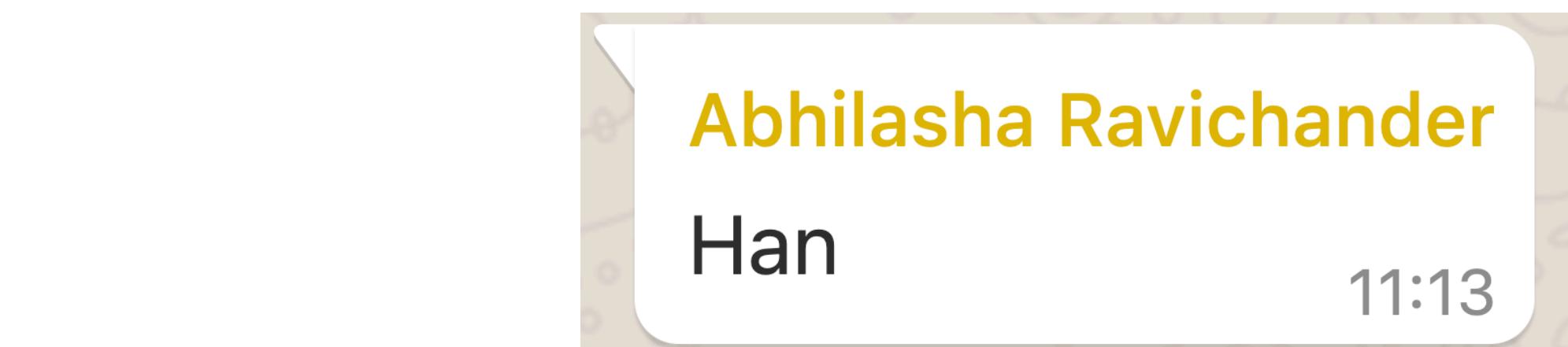


Romanization with friends

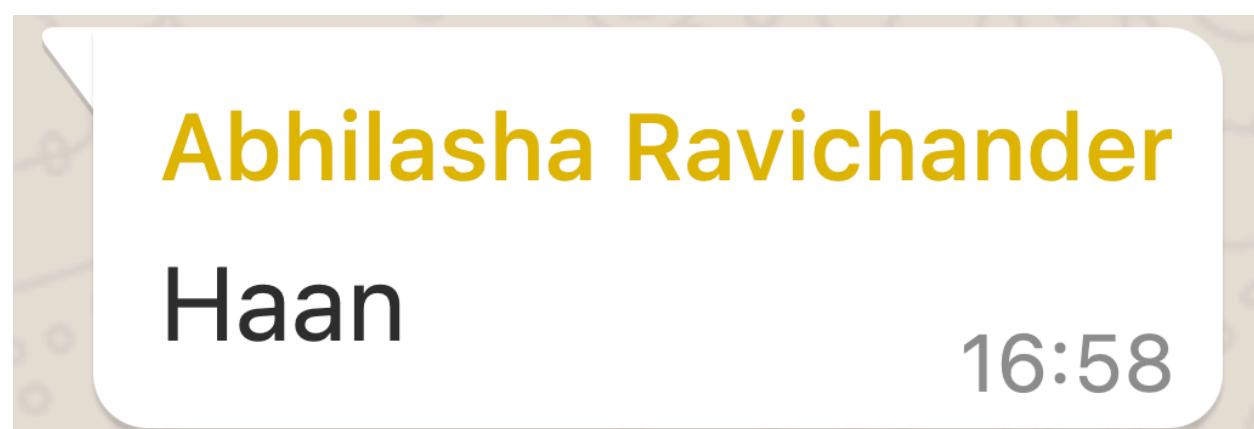
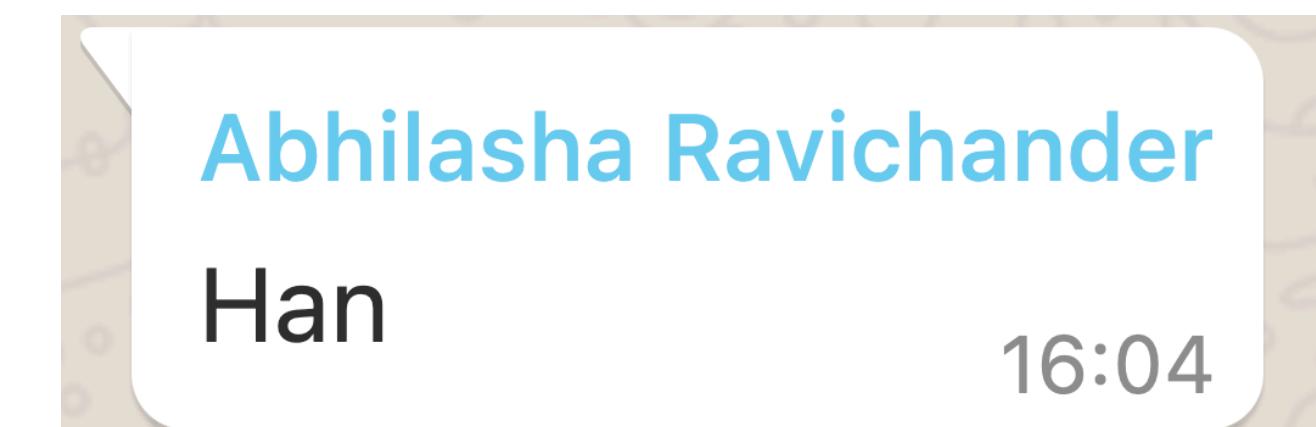


хорошо
[horošo]

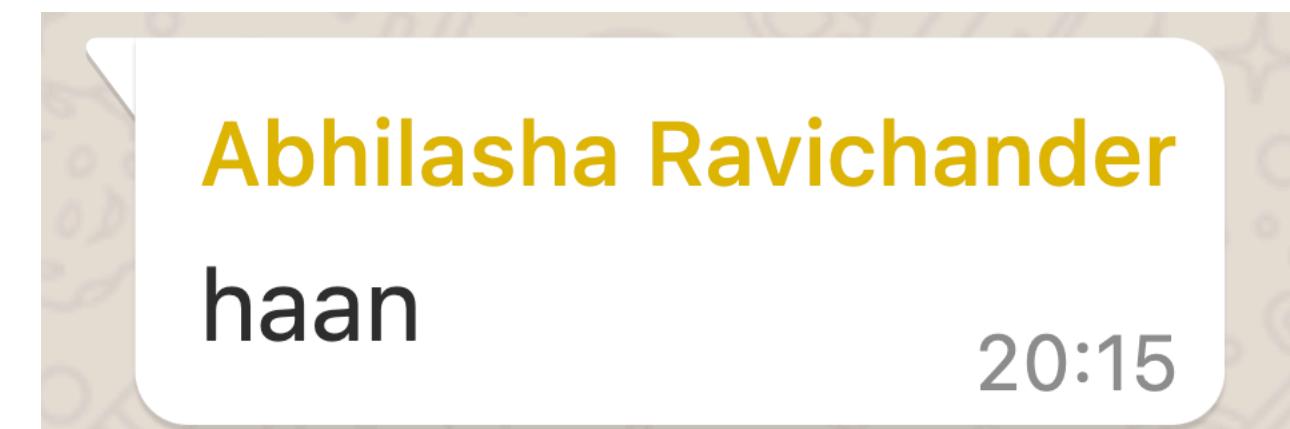
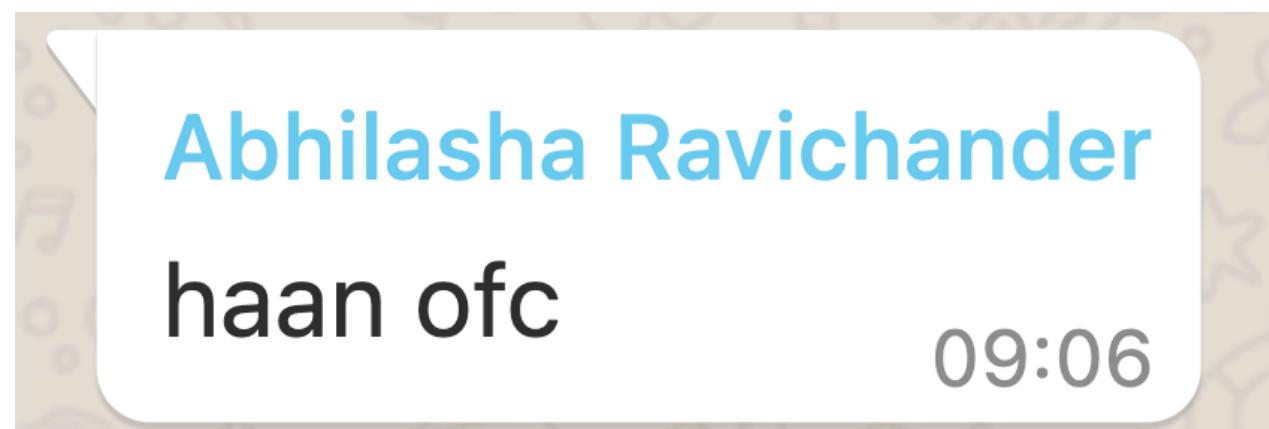
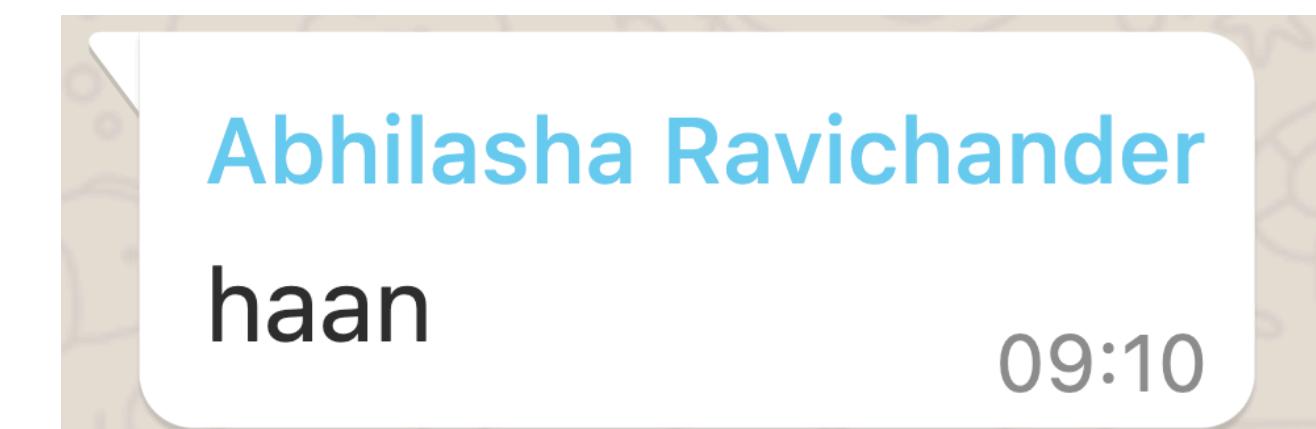
Romanization with friends



हाँ



[hān]



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> |

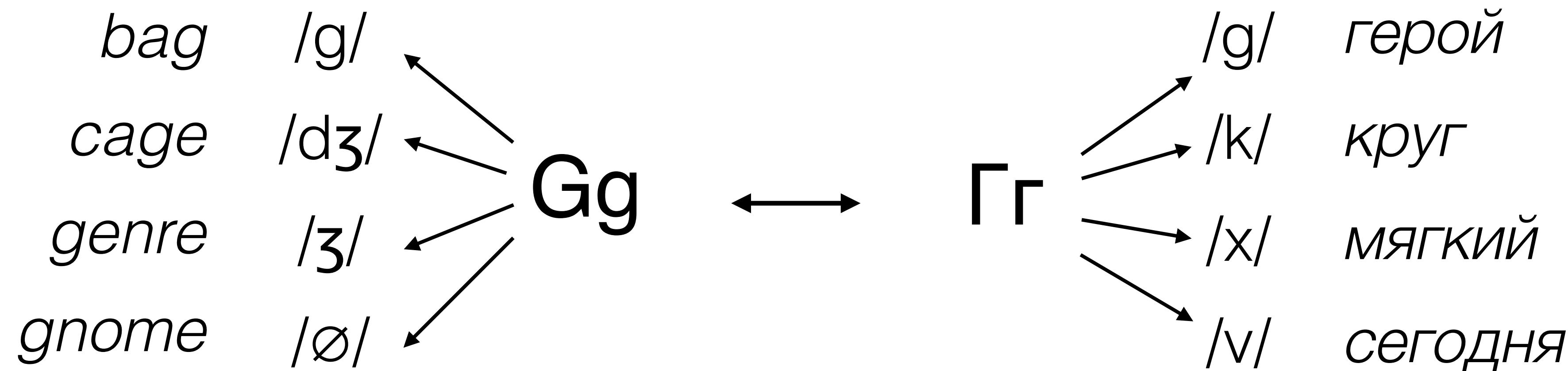
Informal romanization

- 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)

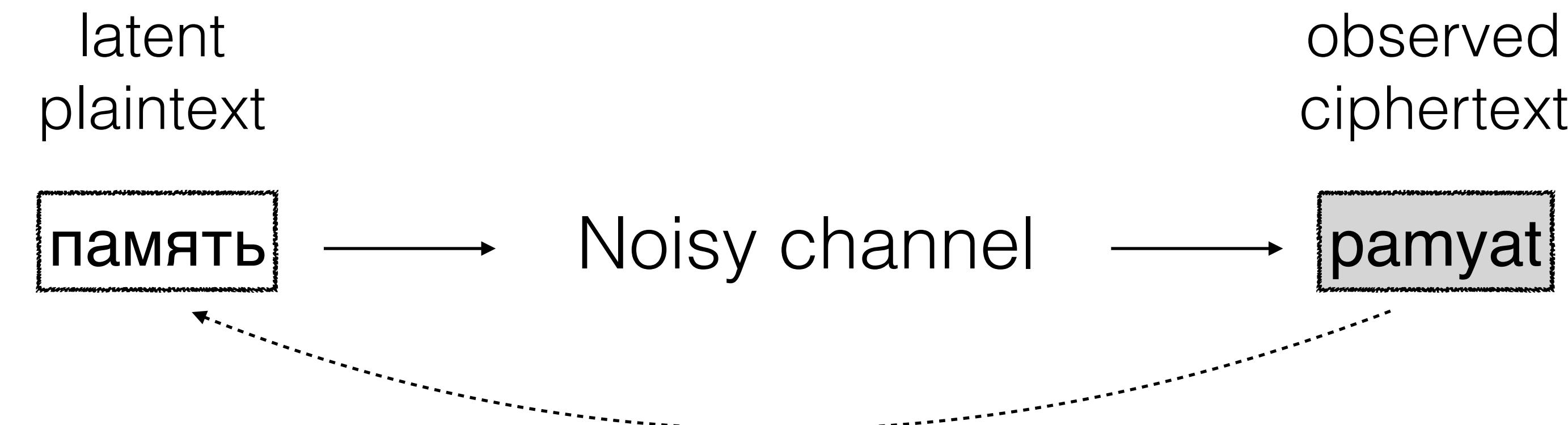
Task framing

- 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 $r = \text{p а m y a t}$

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

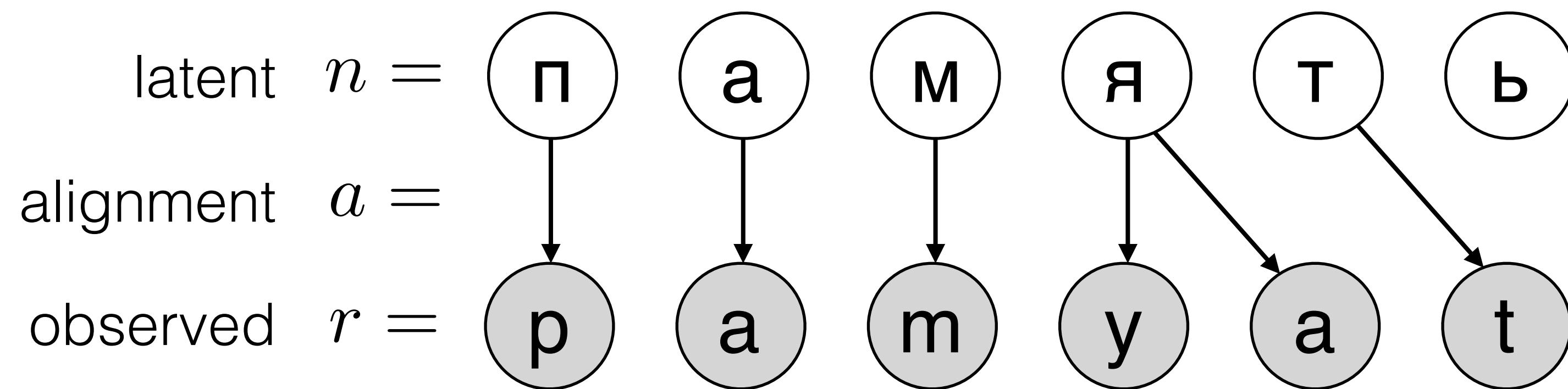
all possible
native script
sequences

n
transition probabilities

emission probabilities

θ
prior on parameters

Noisy-channel model



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

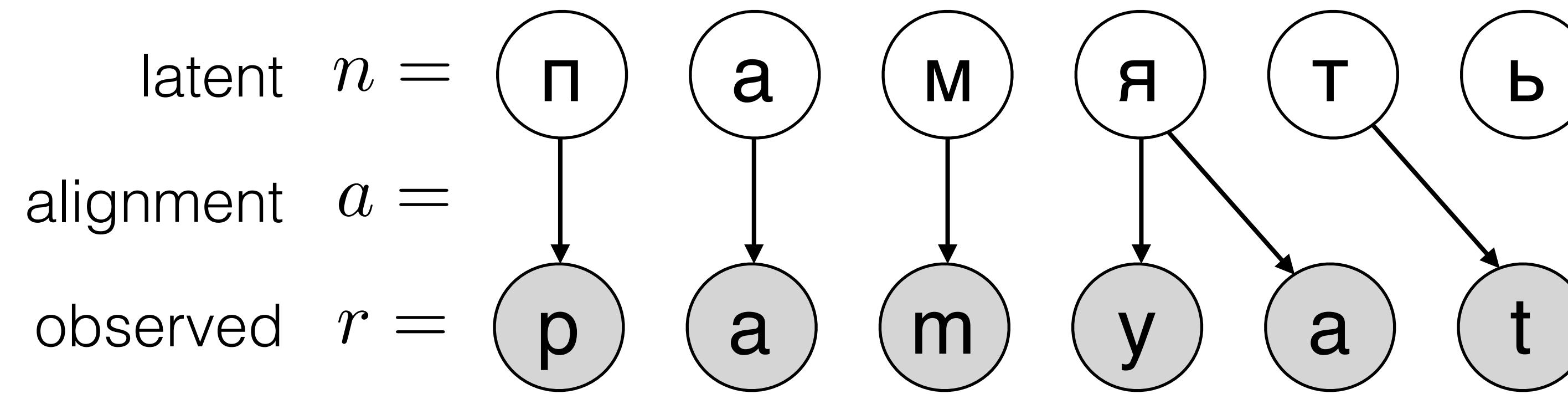
all possible native script sequences and alignments

n, a

emission probabilities

prior on parameters

Noisy-channel model

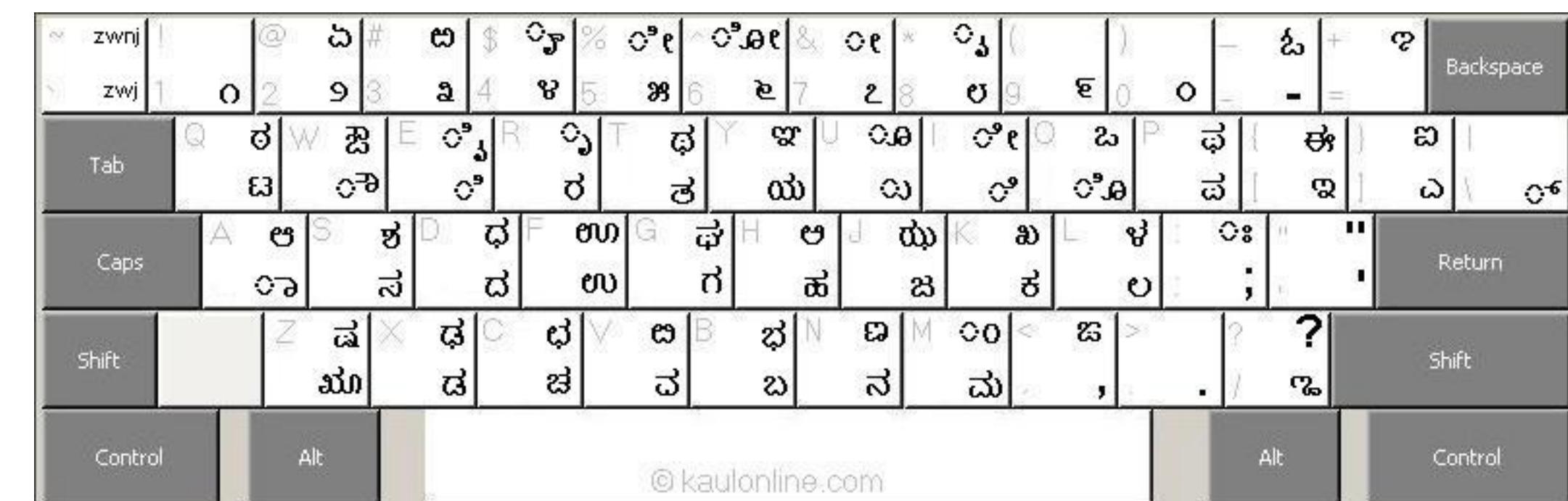


$$p(r) = \sum_{n,a} p(n; \gamma) \cdot p(r|n, a; \theta) \cdot p_{\text{prior}}(\theta; \alpha)$$

/ |
transition probabilities emission probabilities
prior on parameters

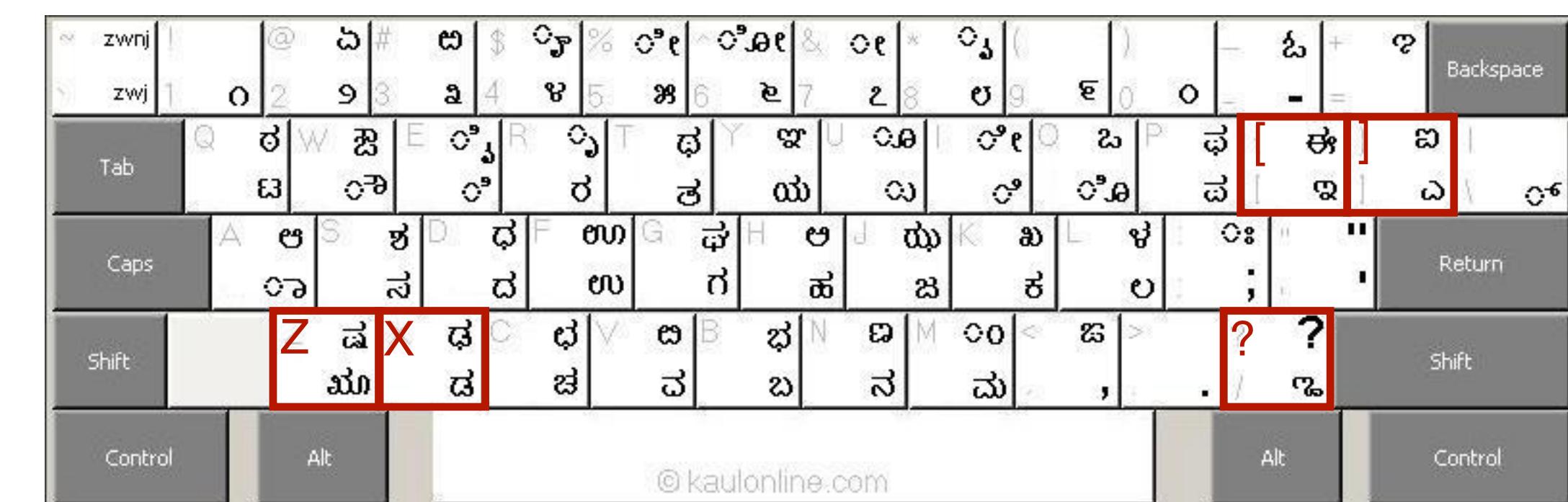
Phonetic bias

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



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

| | | | | | | |
|---------------------------------------|---|---|--|--|--|---|
| p | ρ | ε | p | ρ | P | ප |
| 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 |
| e | ε | ε | e | ε | ϶ | □ |
| 0065 LATIN SMALL LETTER E | 0435 CYRILLIC SMALL LETTER IE | 04BD CYRILLIC SMALL LETTER ABKHASIAN CHE | 212E ESTIMATED SYMBOL | 212F SCRIPT SMALL E | 2147 DOUBLE- STRUCK ITALIC SMALL E | AB32 LATIN SMALL LETTER BLACKLETTER E |

nlpwithfriends.com

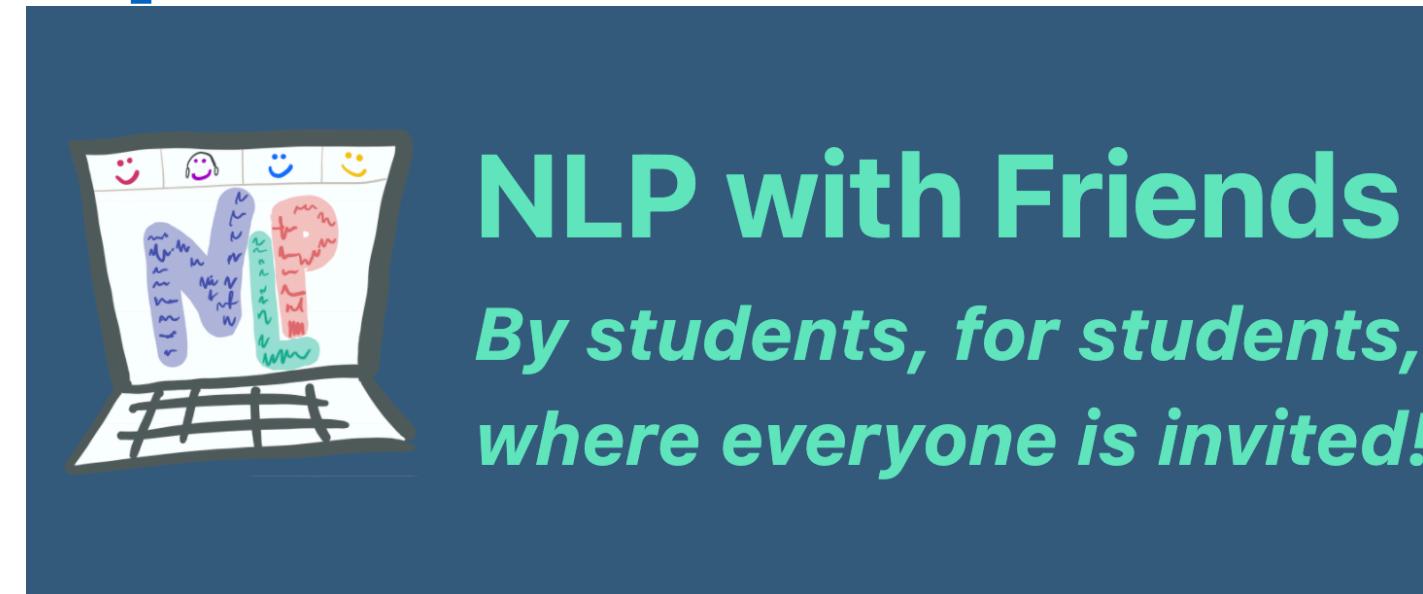
nlpwithfriends.com

Visual bias

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

| | | | | | | |
|--|--|--|---|---|---|--|
| p 0070 LATIN SMALL LETTER P | p 03C1 GREEK SMALL LETTER RHO | e 03F1 GREEK RHO SYMBOL | p 0440 CYRILLIC SMALL LETTER ER | p 2374 APL FUNCTIONAL SYMBOL RHO | P 2CA3 COPTIC SMALL LETTER RO | p 1D429 MATHEMATICAL BOLD SMALL P |
| e 0065 LATIN SMALL LETTER E | e 0435 CYRILLIC SMALL LETTER IE | e 04BD CYRILLIC SMALL LETTER ABKHAZIAN CHE | e 212E ESTIMATED SYMBOL | e 212F SCRIPT SMALL E | € 2147 DOUBLE- STRUCK ITALIC SMALL E | □ AB32 LATIN SMALL LETTER BLACKLETTER E |

nlpwithfriends.com



NLP with Friends
*By students, for students,
where everyone is invited!*

nlpwithfriends.com

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!

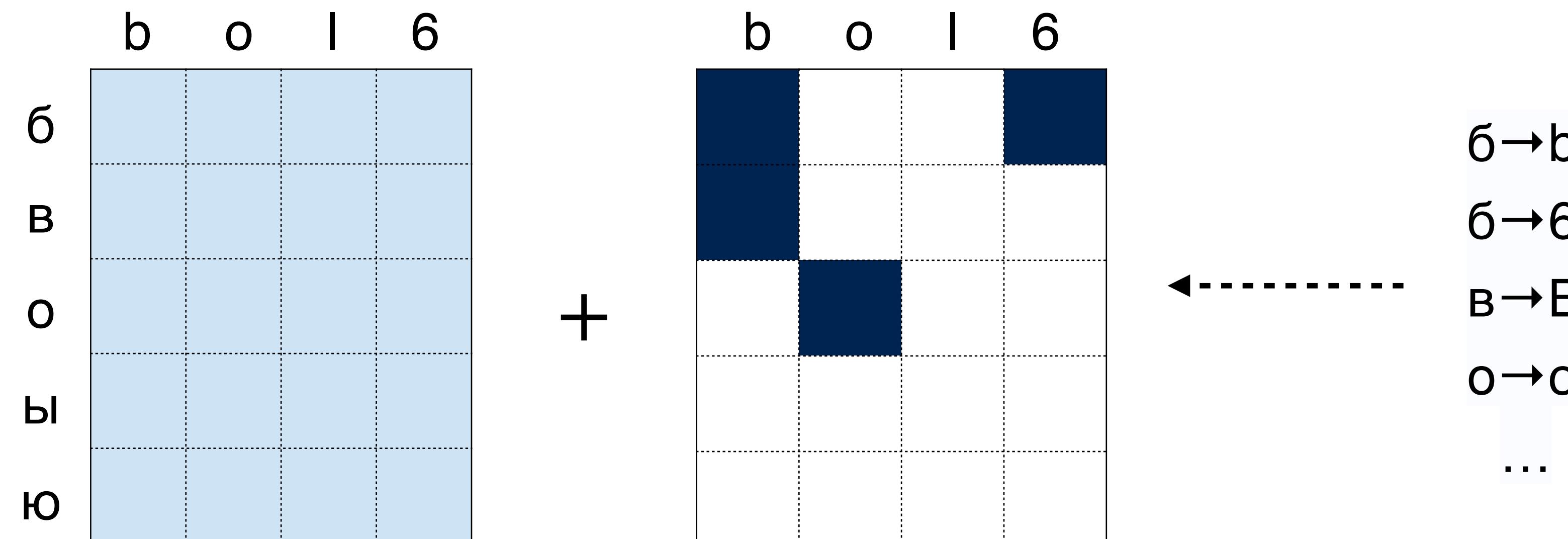
| | | | | | | |
|---------------------------------------|---|---|--|--|--|---|
| p | ρ | ε | p | ρ | P | ප |
| 0070 LATIN SMALL LETTER P | 03C1 GREEK SMALL LETTER RHO | 03F1 GREEK RHO SYMBOL | 0440 CYRILLIC SMALL LETTER ER | 2374 APL FUNCTIONAL SYMBOL RHO | 2CA3 COPTIC SMALL LETTER RO | 1D429 MATHEMATICAL BOLD SMALL P |
| e | ε | € | e | € | € | □ |
| 0065 LATIN SMALL LETTER E | 0435 CYRILLIC SMALL LETTER IE | 04BD CYRILLIC SMALL LETTER ABKHAZIAN CHE | 212E ESTIMATED SYMBOL | 212F SCRIPT SMALL E | 2147 DOUBLE- STRUCK ITALIC SMALL E | AB32 LATIN SMALL LETTER BLACKLETTER E |

Informative priors

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

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

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

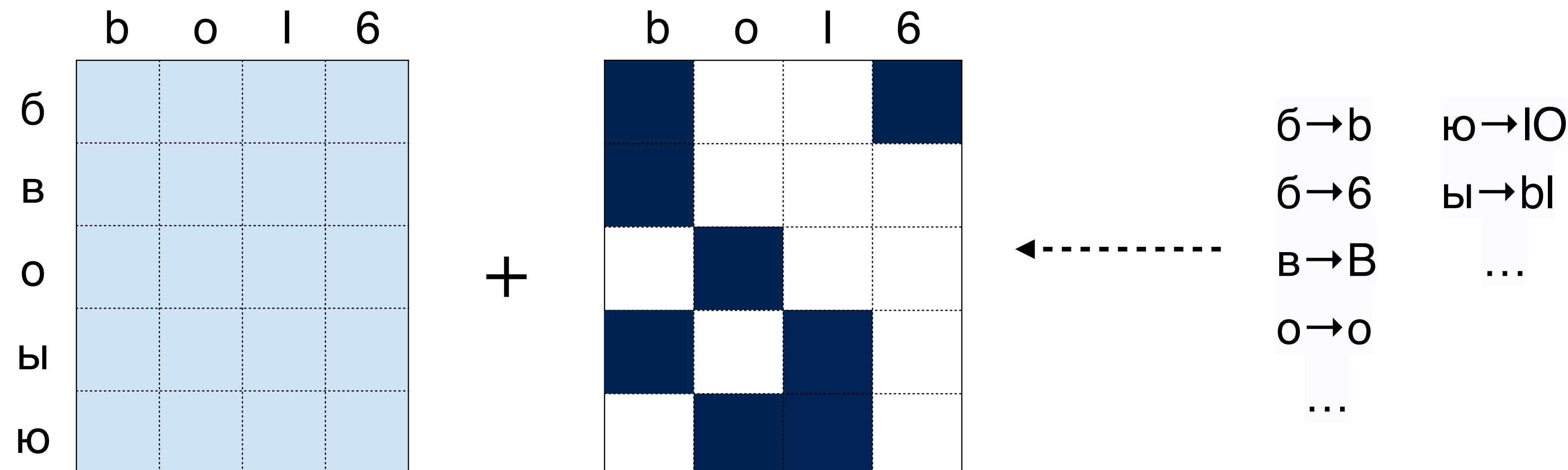


Informative priors

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

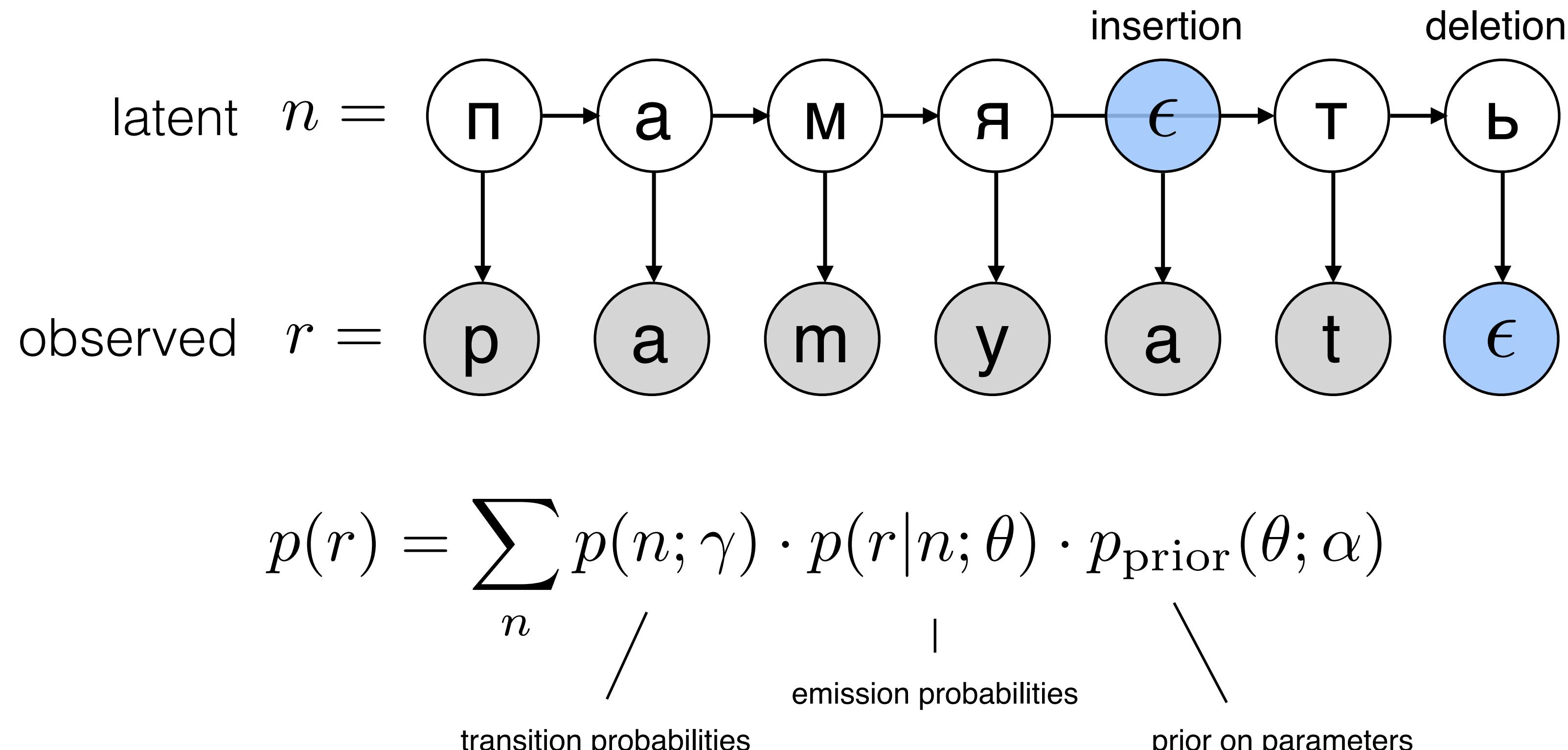
$$c_r | c_n \sim \text{Mult}(\theta_{c_n})$$

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



Noisy-channel model

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

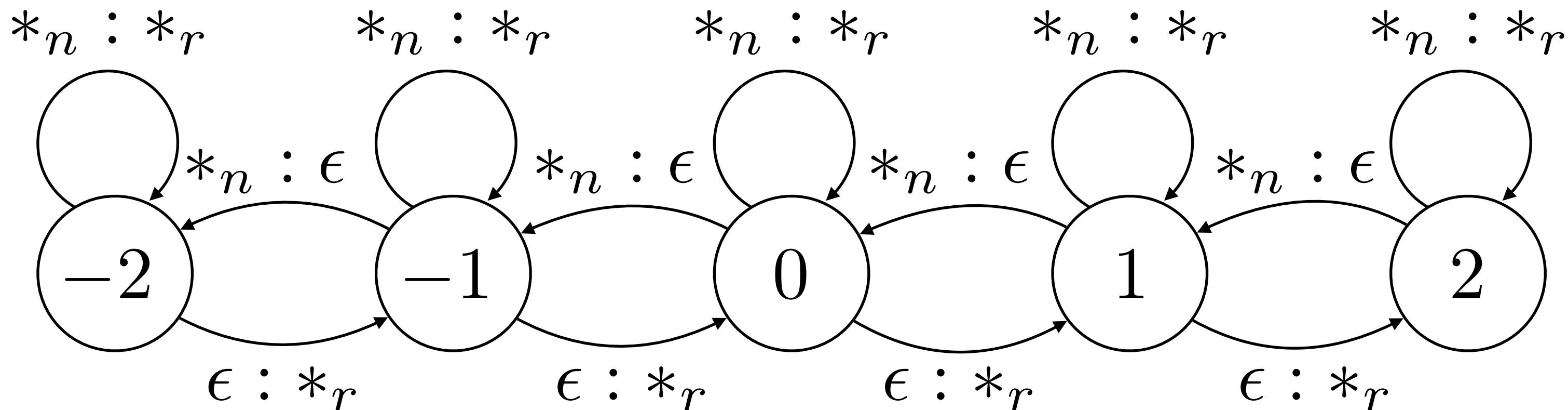
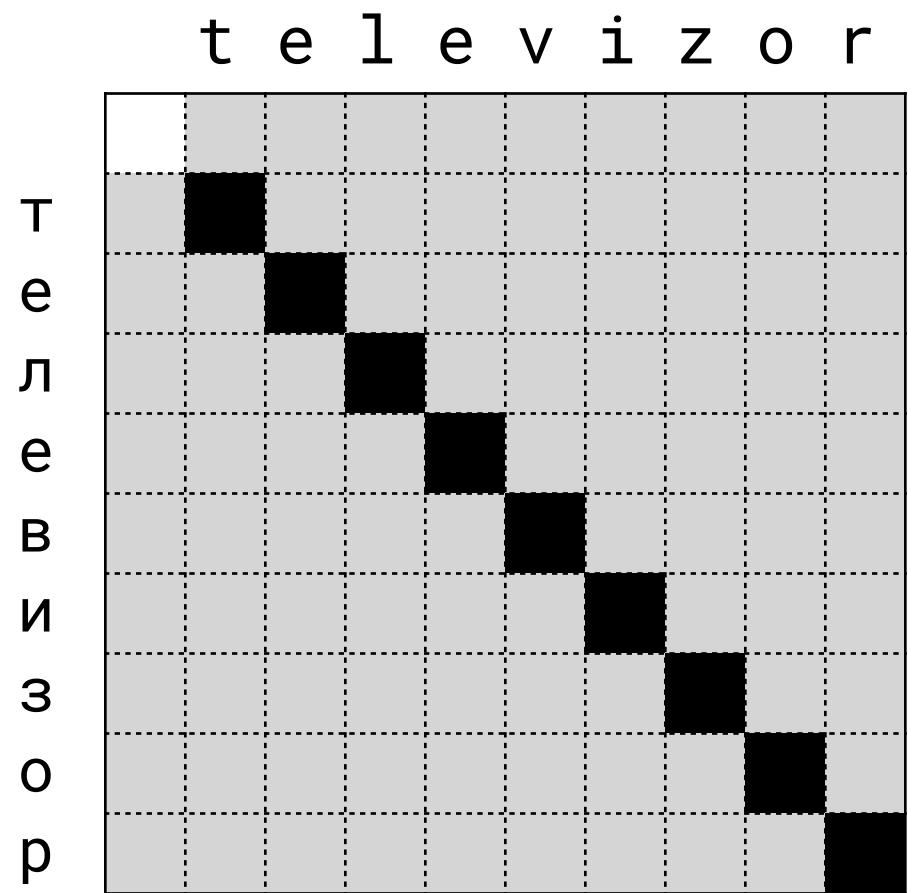


WFST cascade

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

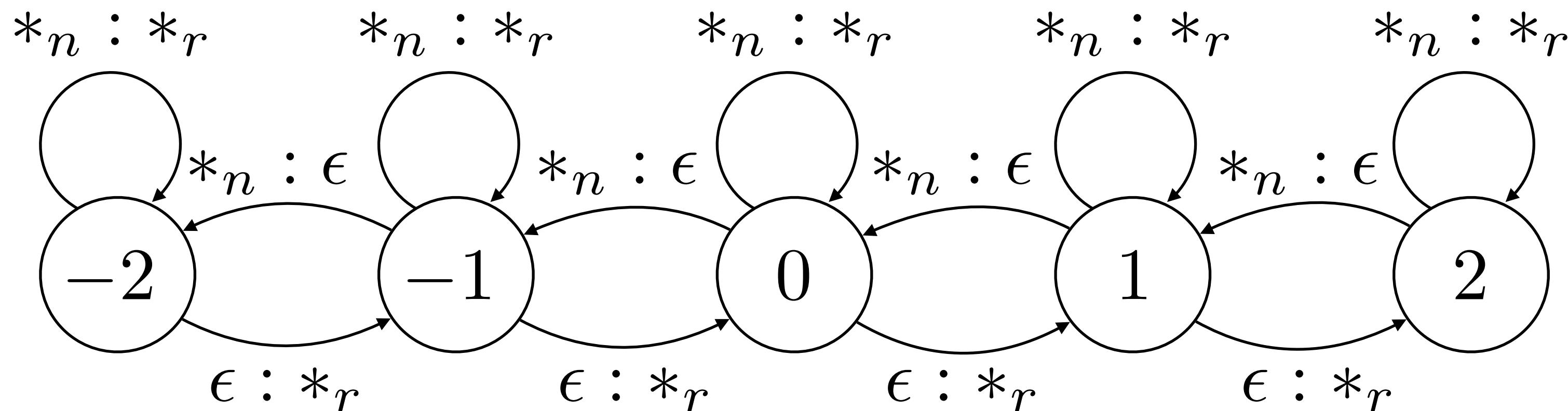
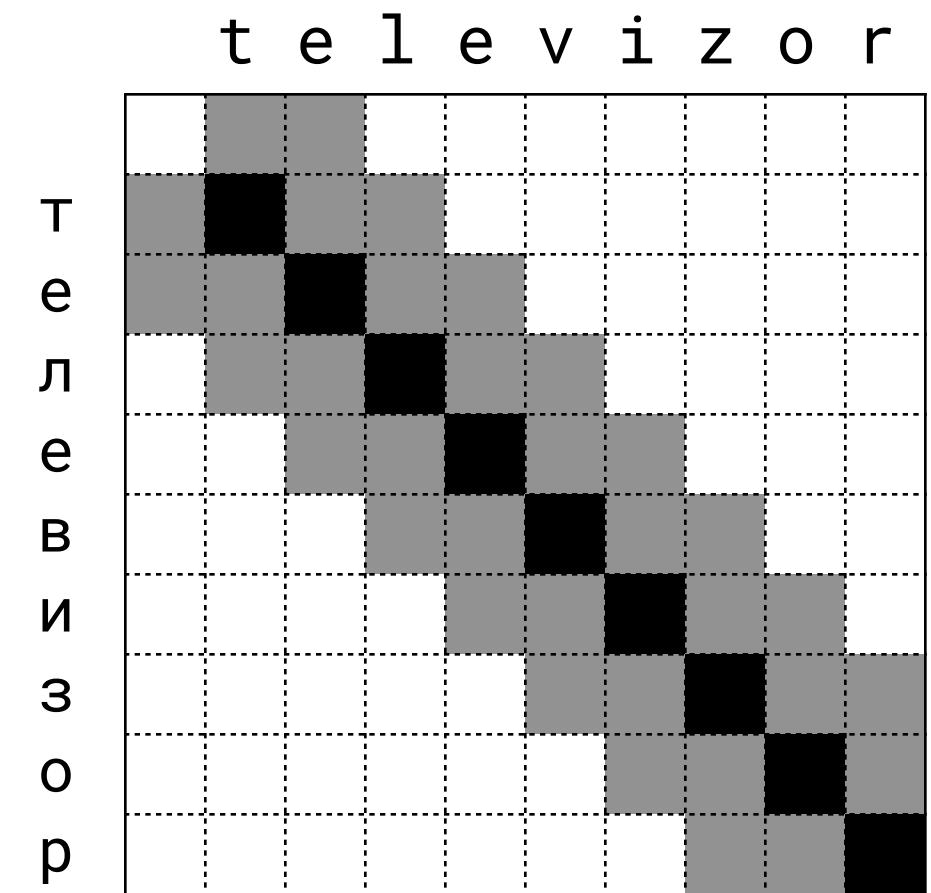
Emission model

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



Training and inference

- Training with EM algorithm
 - E-step: shortest distance in expectation semiring (Eisner, 2002)
 - M-step: parameter reestimation

Training and inference

- Training with EM algorithm
 - E-step: shortest distance in expectation semiring (Eisner, 2002)
  Easy bookkeeping
 - M-step: parameter reestimation
  Slow training
- 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

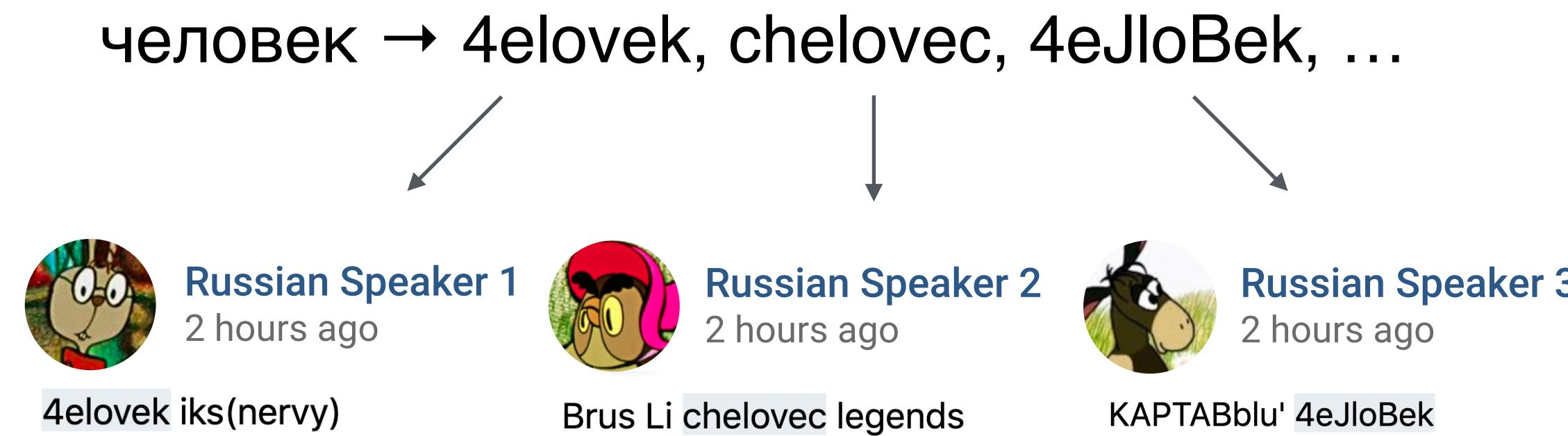
- Arabic: LDC BOLT dataset (Bies et al., 2014)
 - Arabizi SMS/chat dialogs, converted to CODA (Habash et al., 2012)
- Kannada: Dakshina dataset (Roark et al., 2020)
 - Kannada Wikipedia, romanizations elicited from native speakers
- Russian:
 - Romanized: collected and partly annotated data from social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



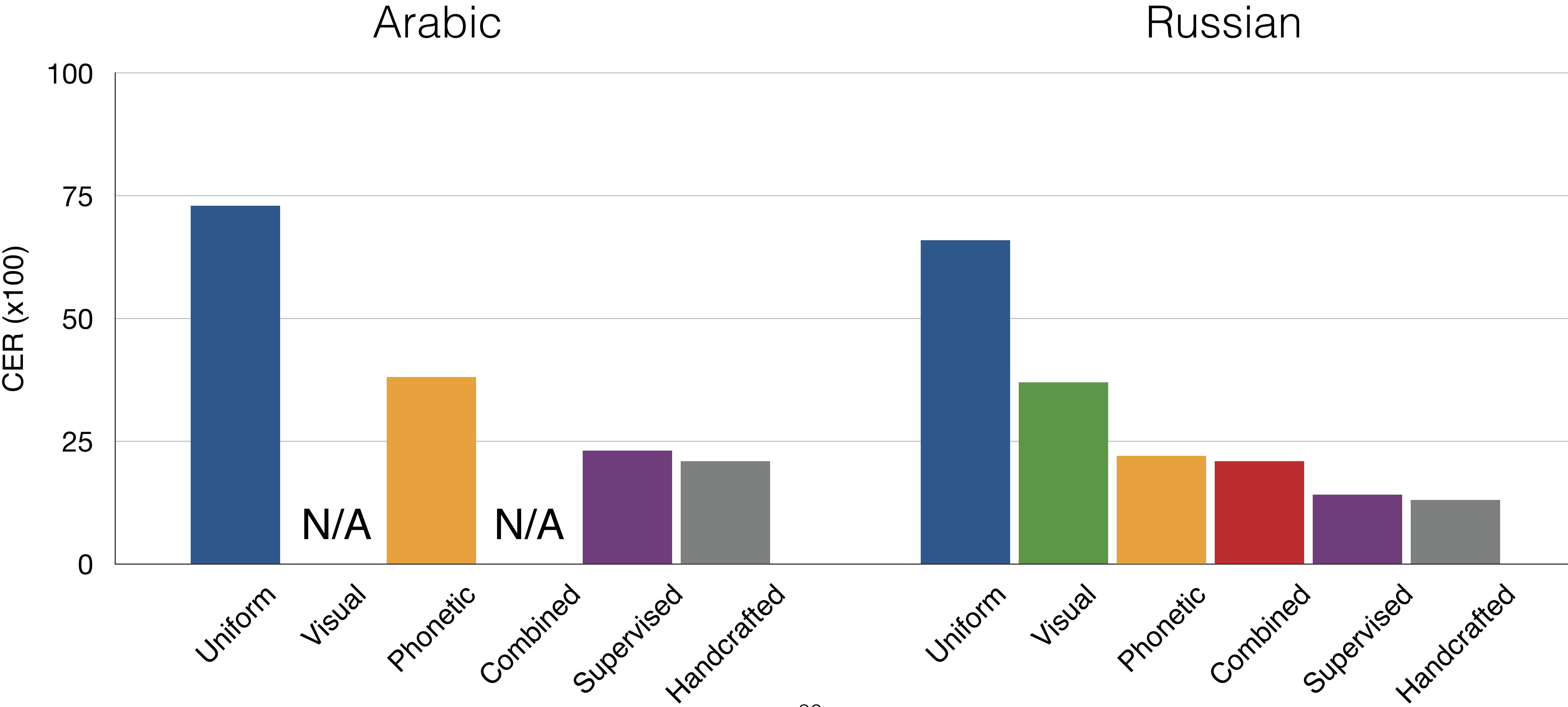
- Manually removed sentences in other languages (e.g. Polish)
- Annotated validation and test with minor error correction

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

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

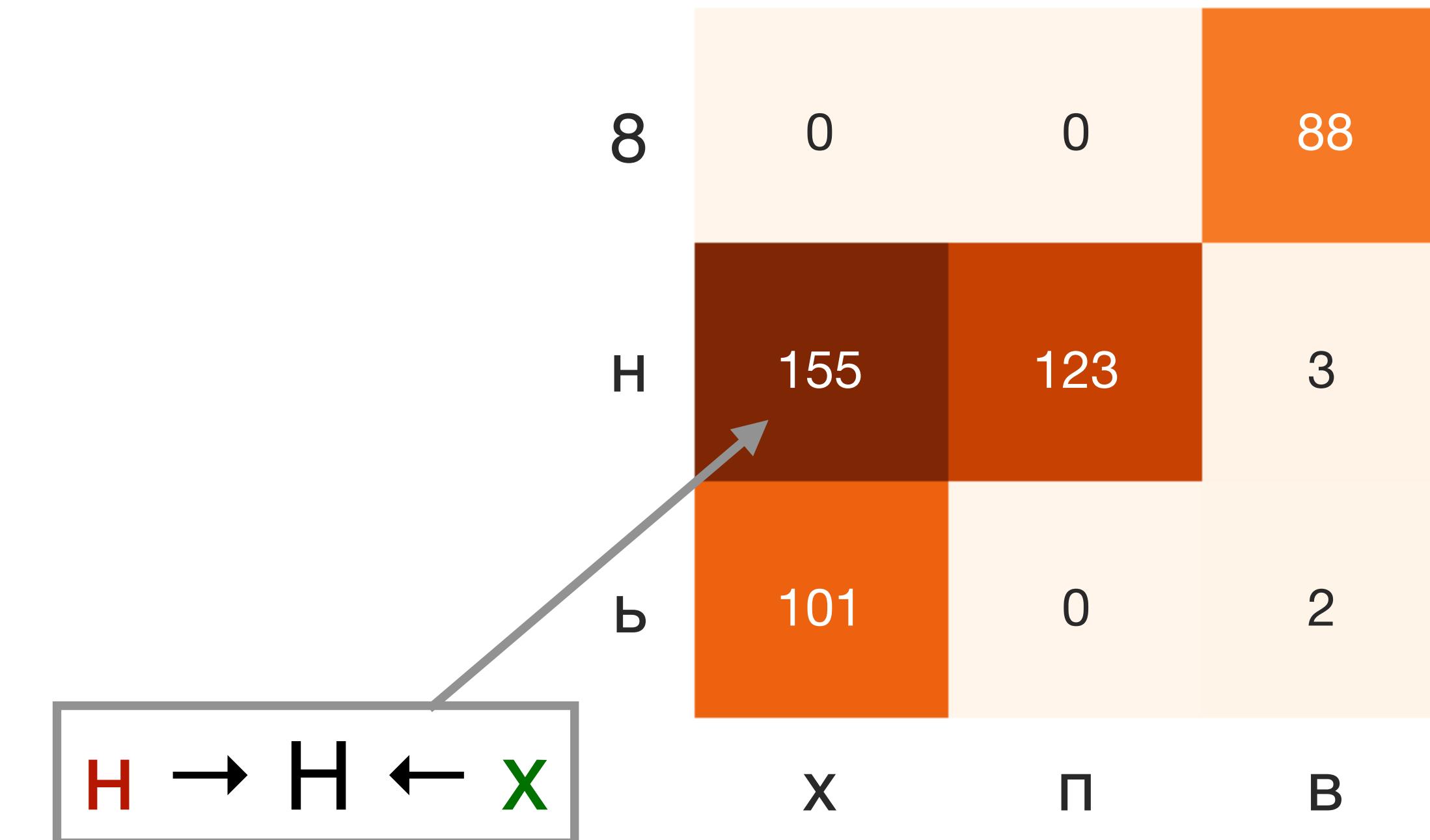
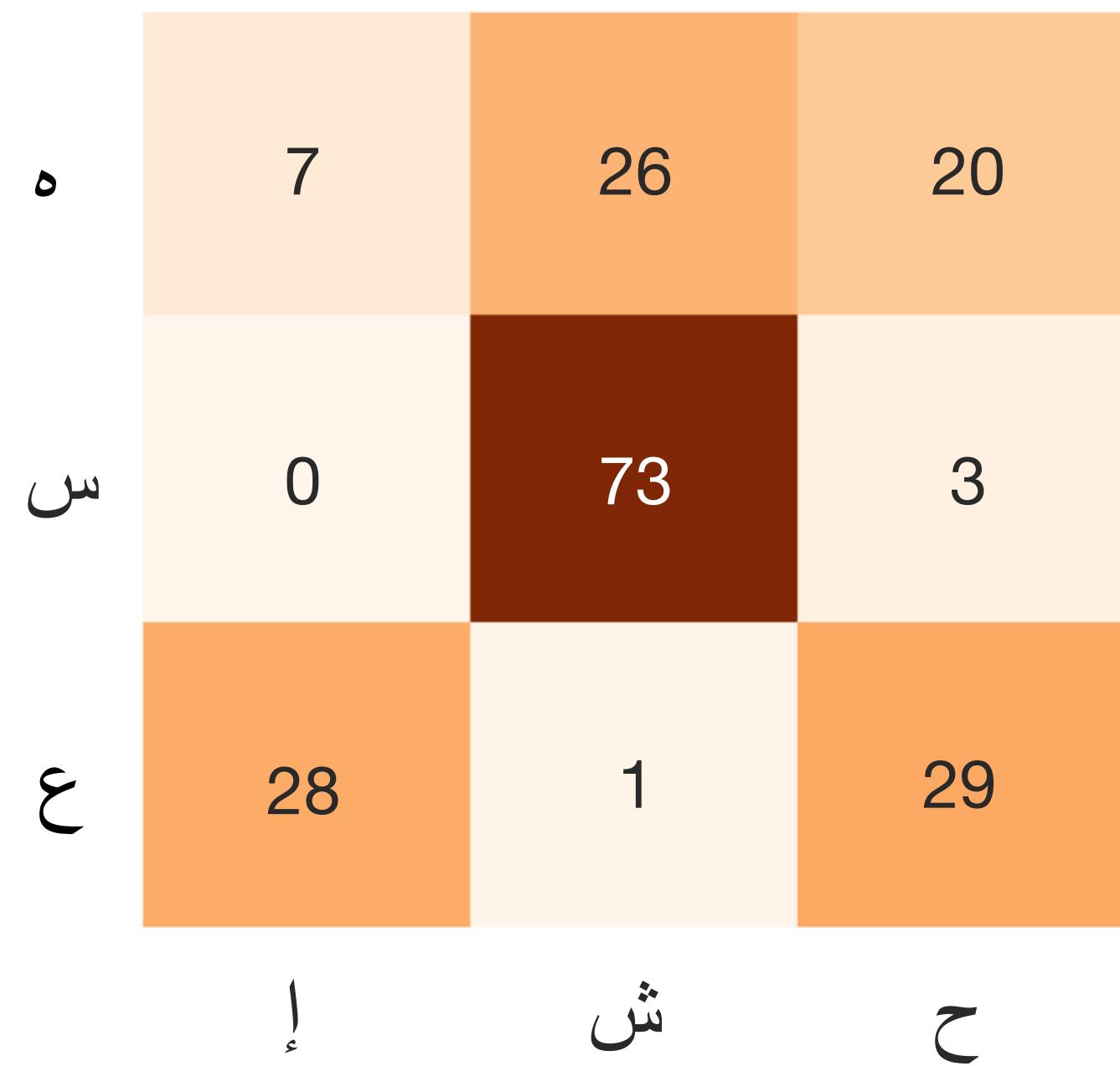
Target: происходит с природой ЧТО-ТО <...>

WFST results



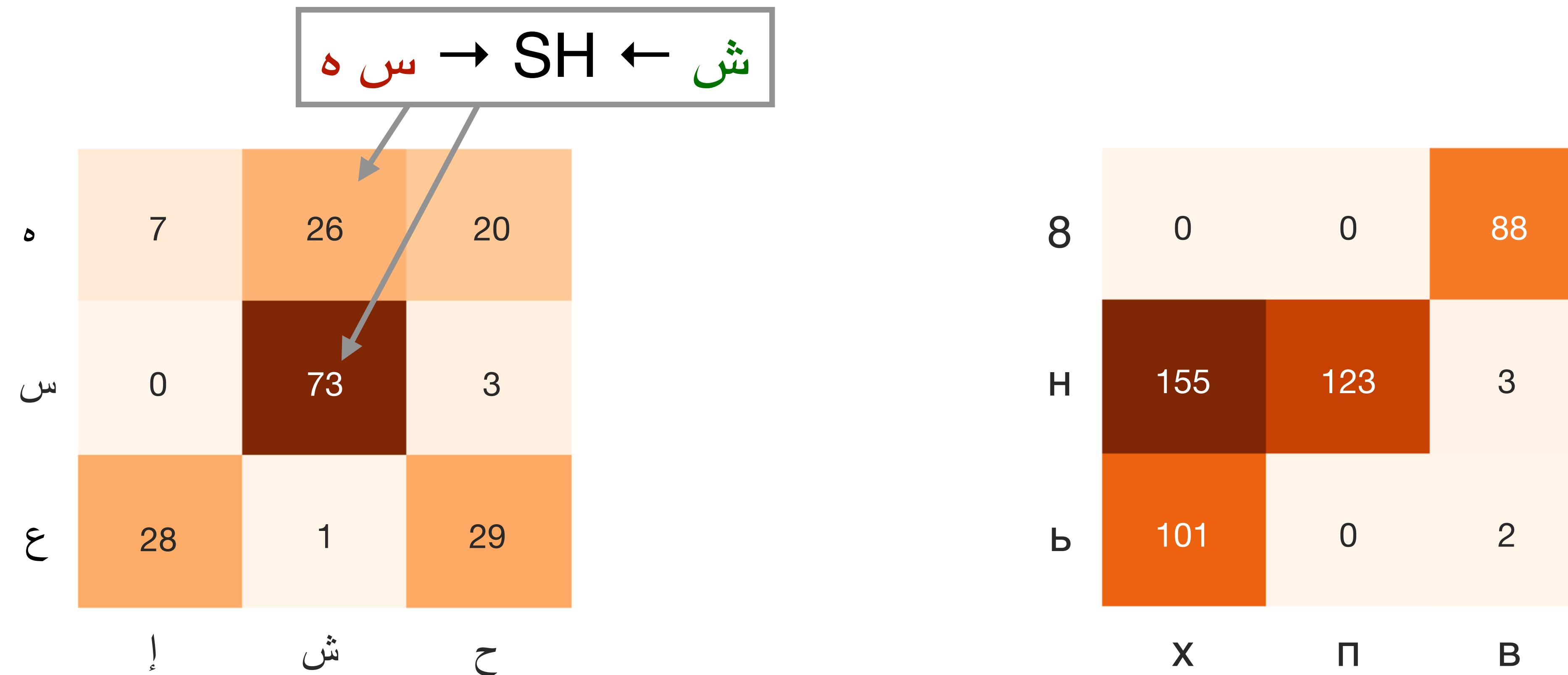
WFST error analysis

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



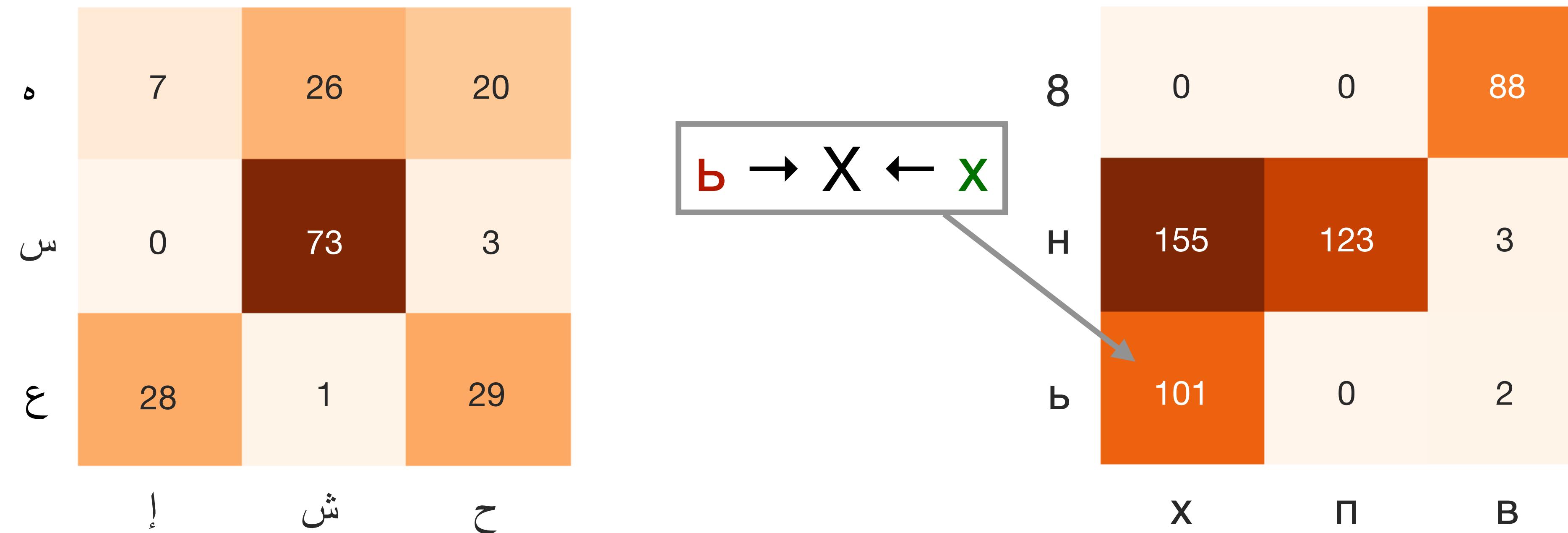
WFST error analysis

- 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**

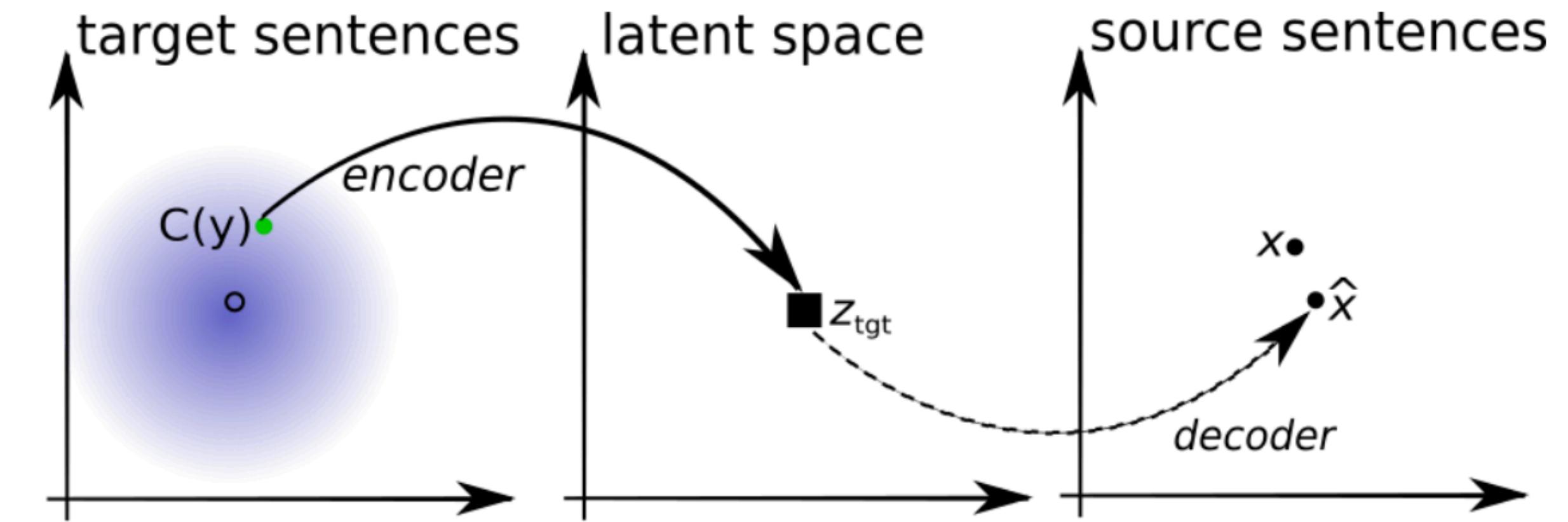
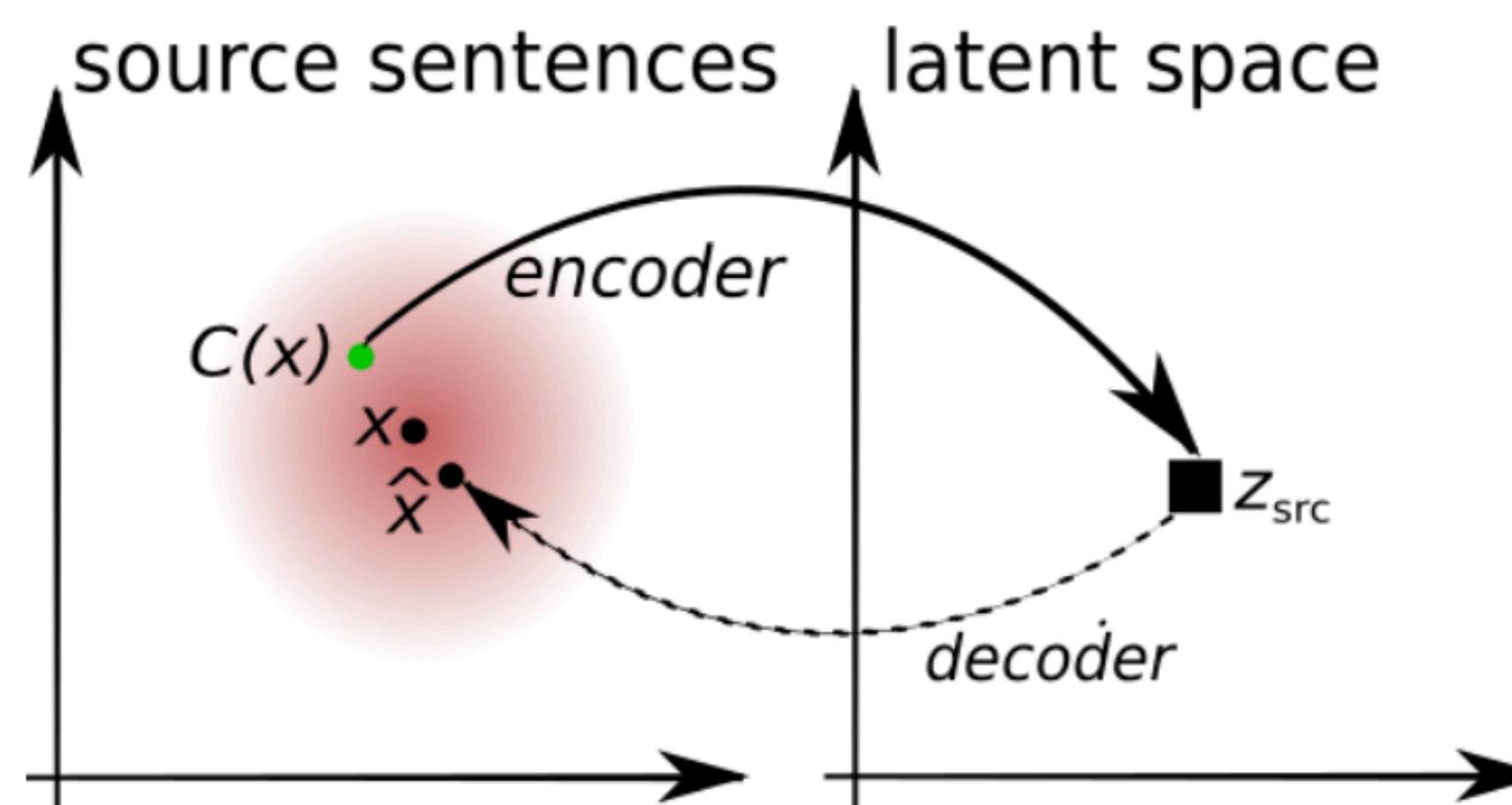
- ✓ Strong language model
- ✓ Faster batch processing
- ✗ Need large training data
- ✗ Hallucinations and search errors

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.

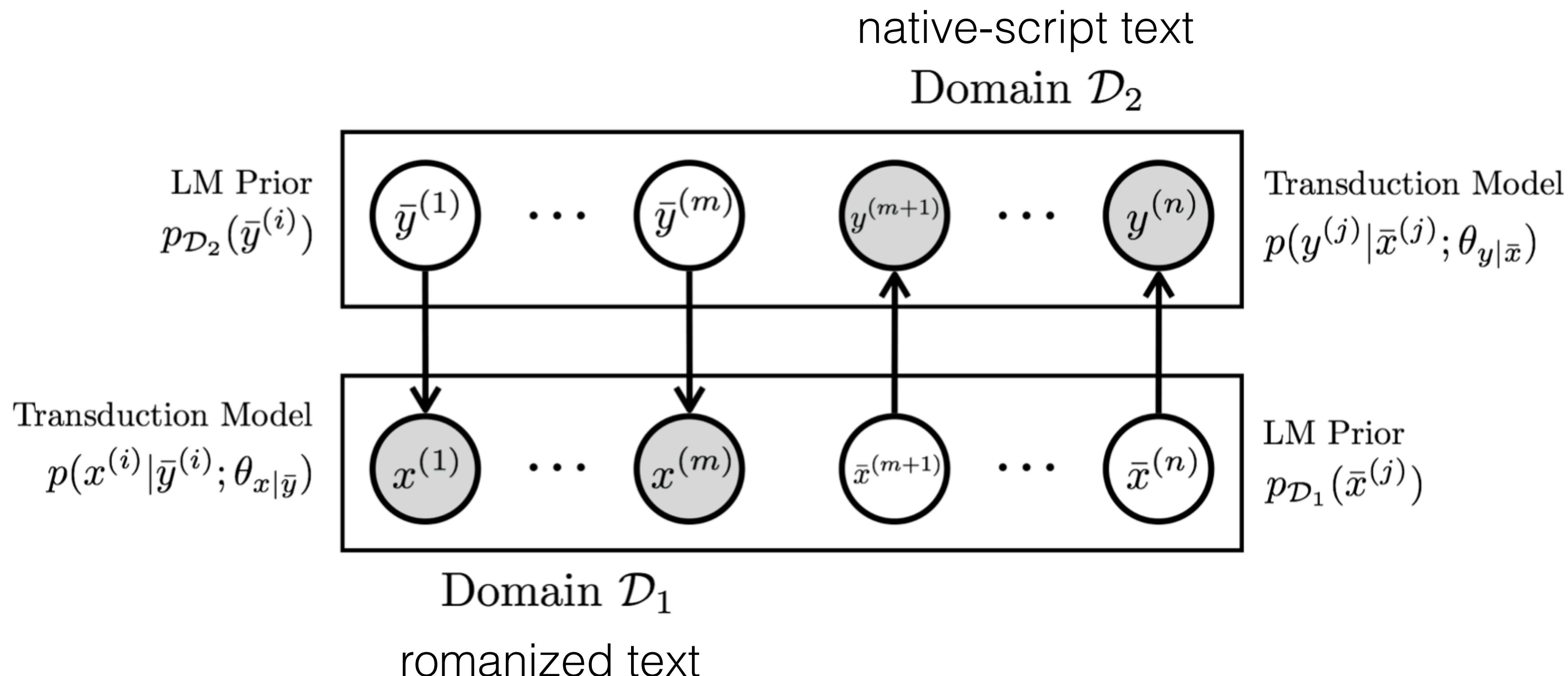
Unsupervised seq2seq

- 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)

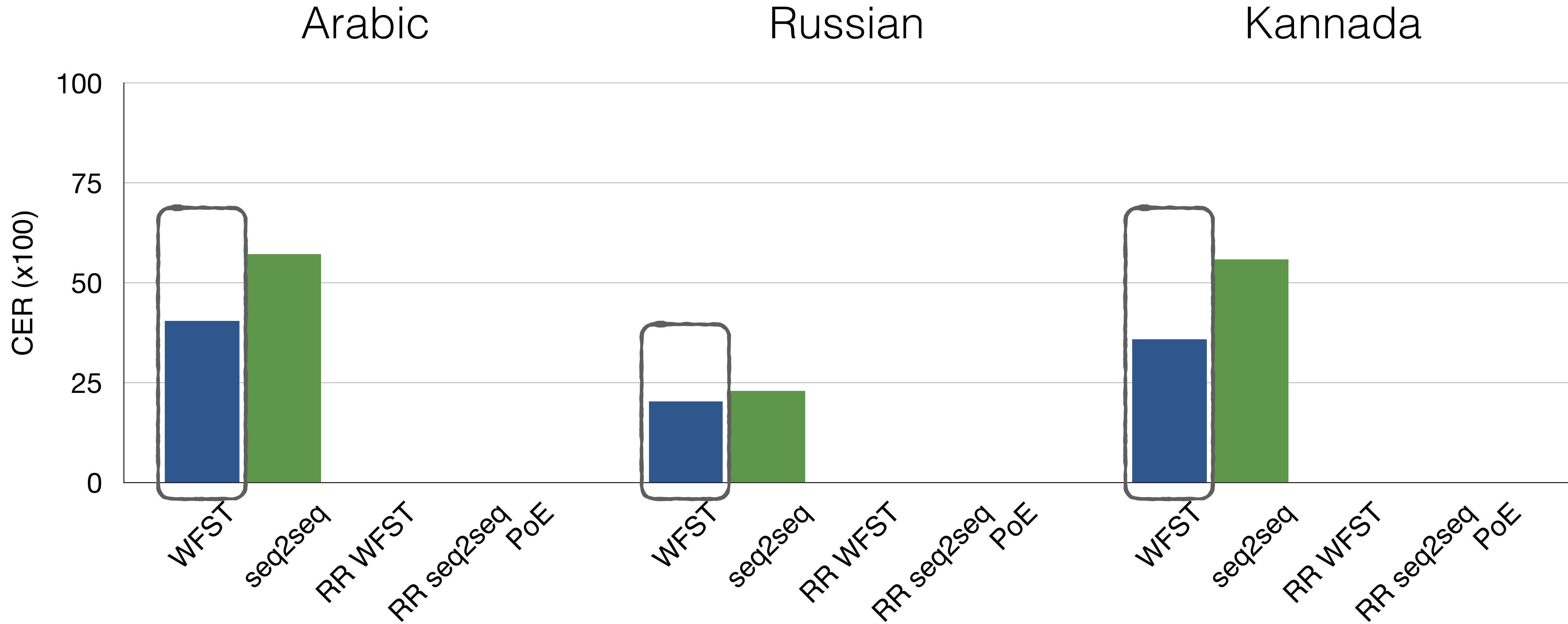


Model combinations

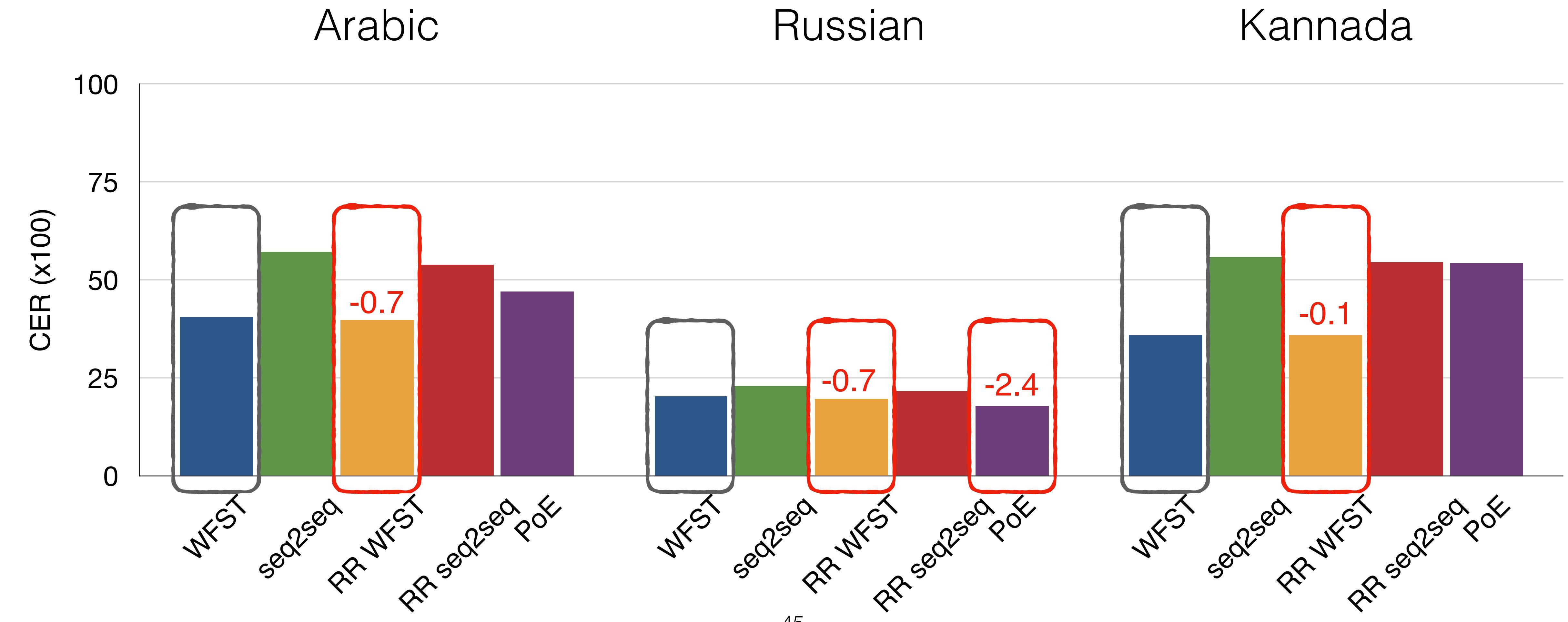
- Reranking
 - M1 generates top k candidate outputs
 - M2 selects the highest-scoring candidate
- Product of experts
 - Beam search on the WFST lattice
 - WFST arcs reweighted with Seq2seq softmax at the corresponding timestep
 - Deletions of input characters are not reweighted
 - Candidates are grouped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

Results

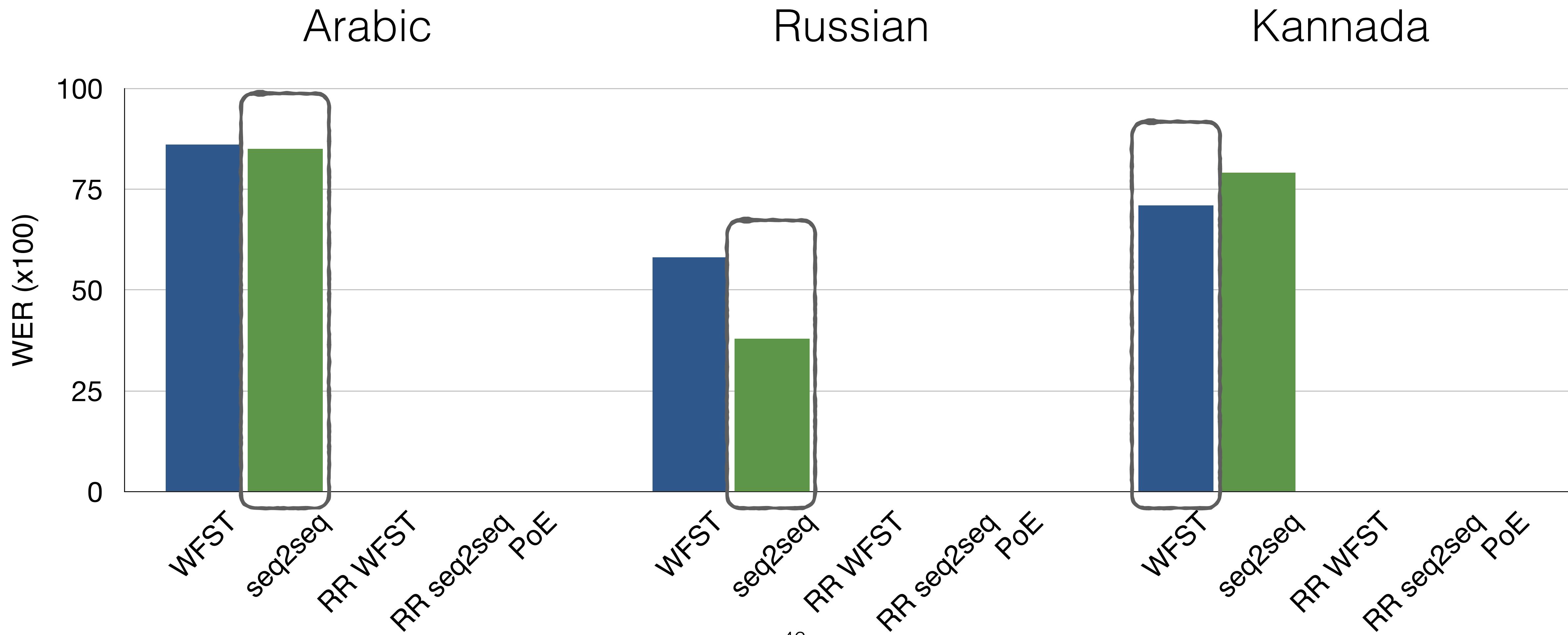
Base models are trained on different amounts of data!



Results

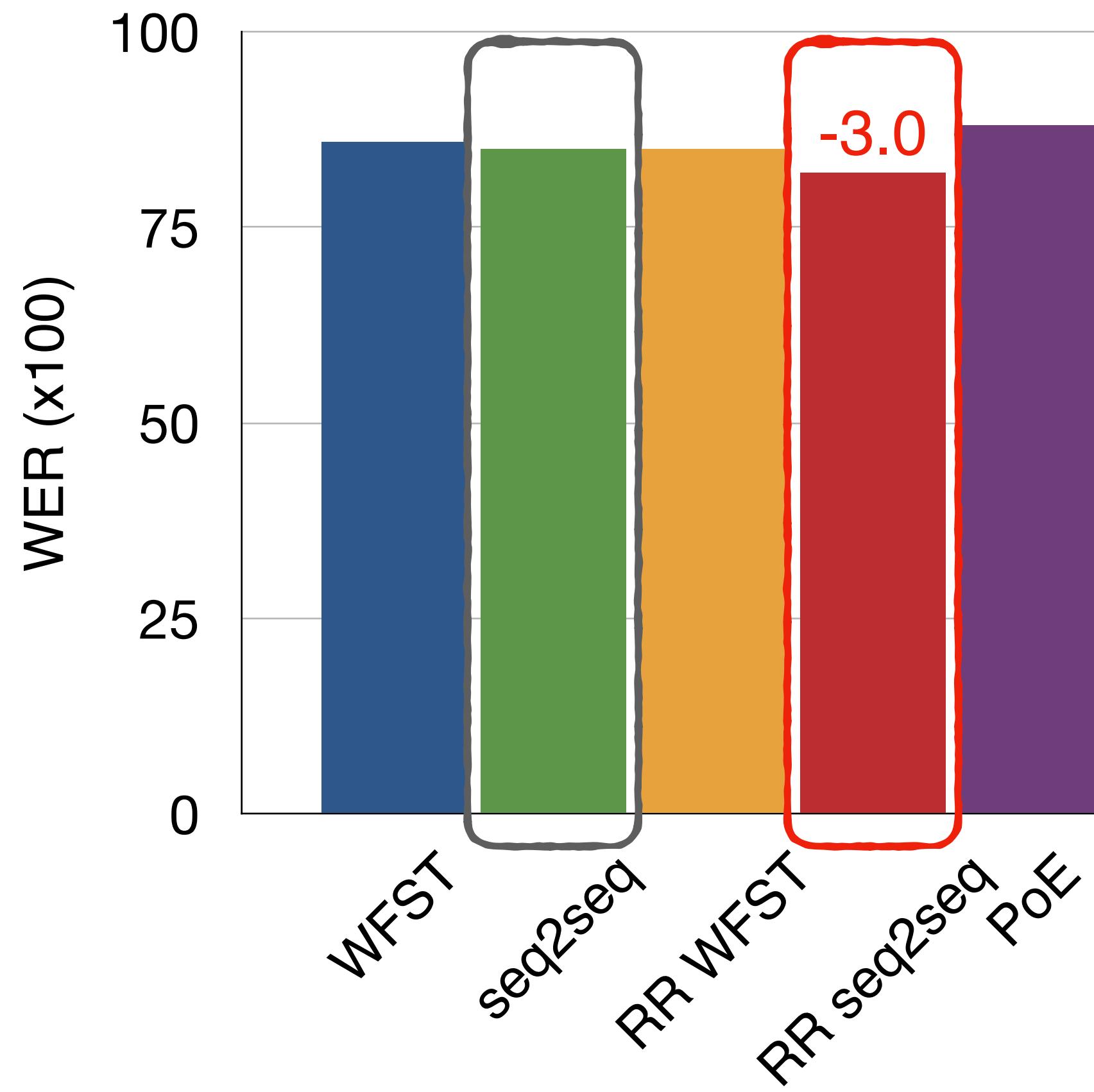


Results



Results

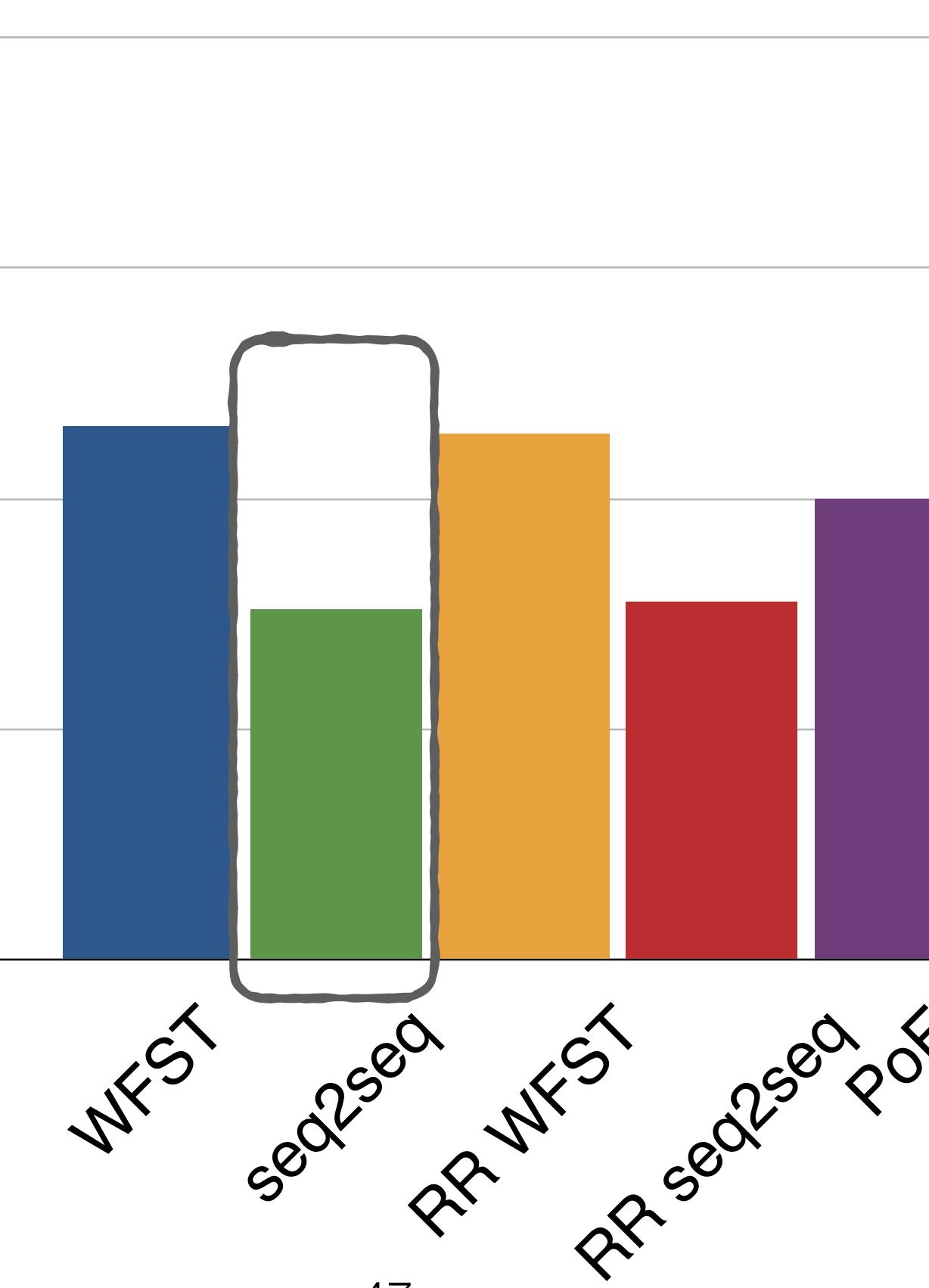
Arabic



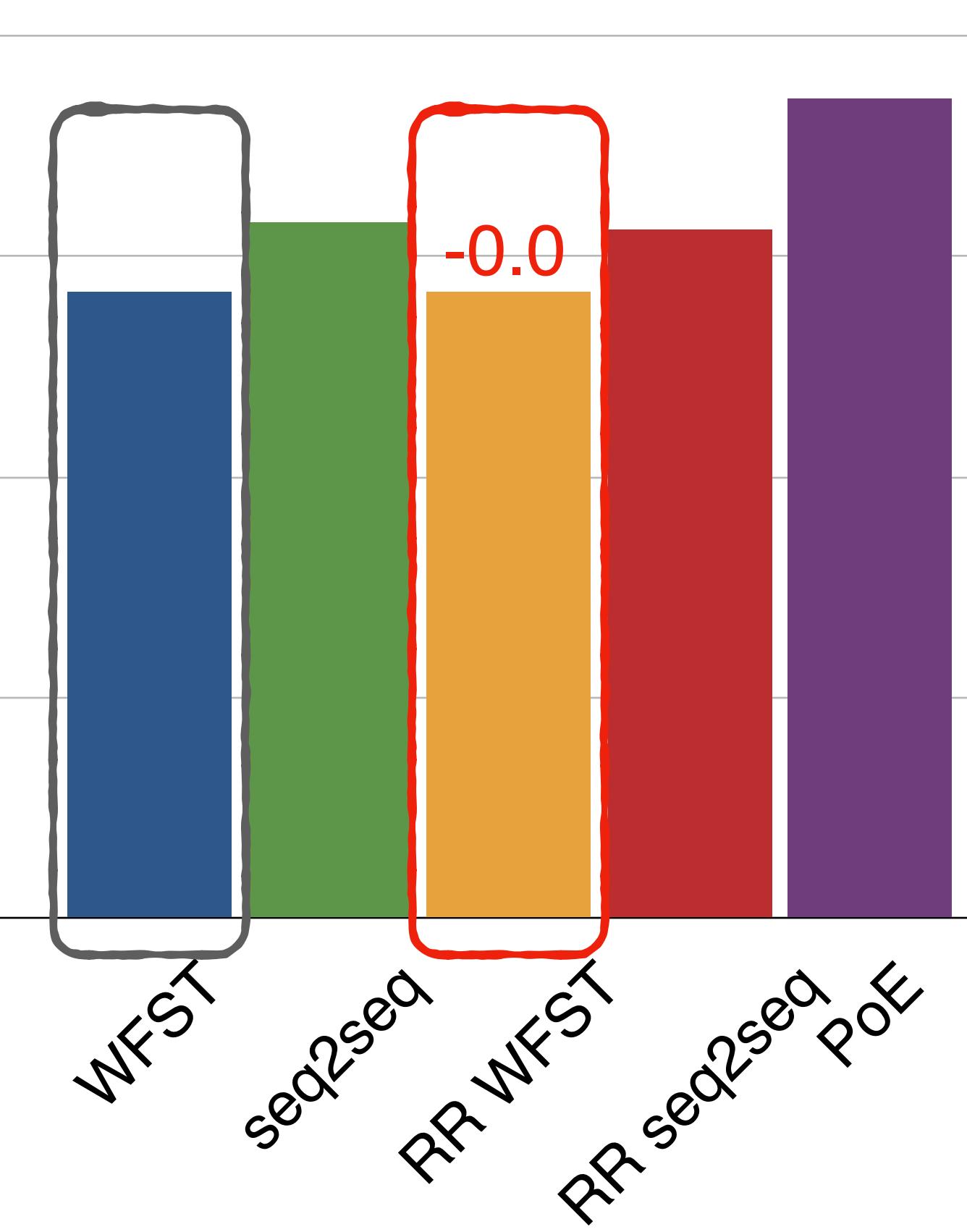
Russian

Russian

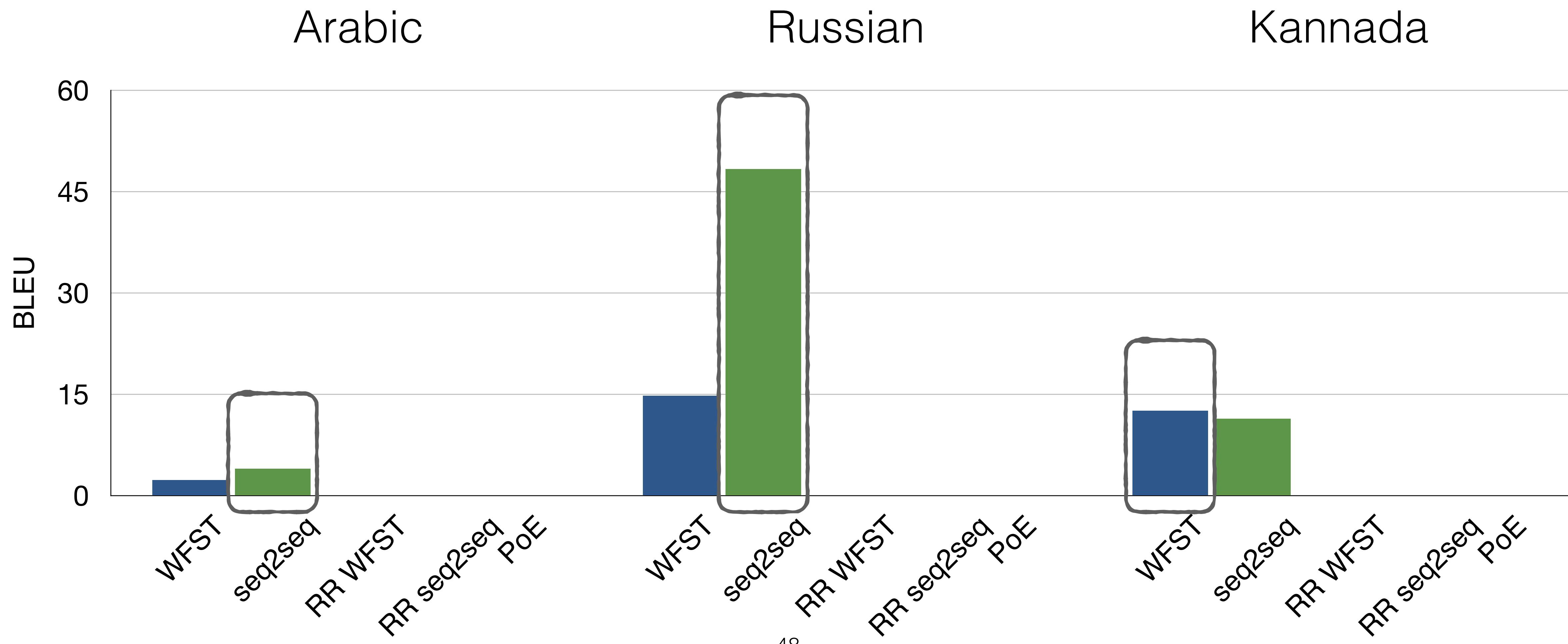
Russian



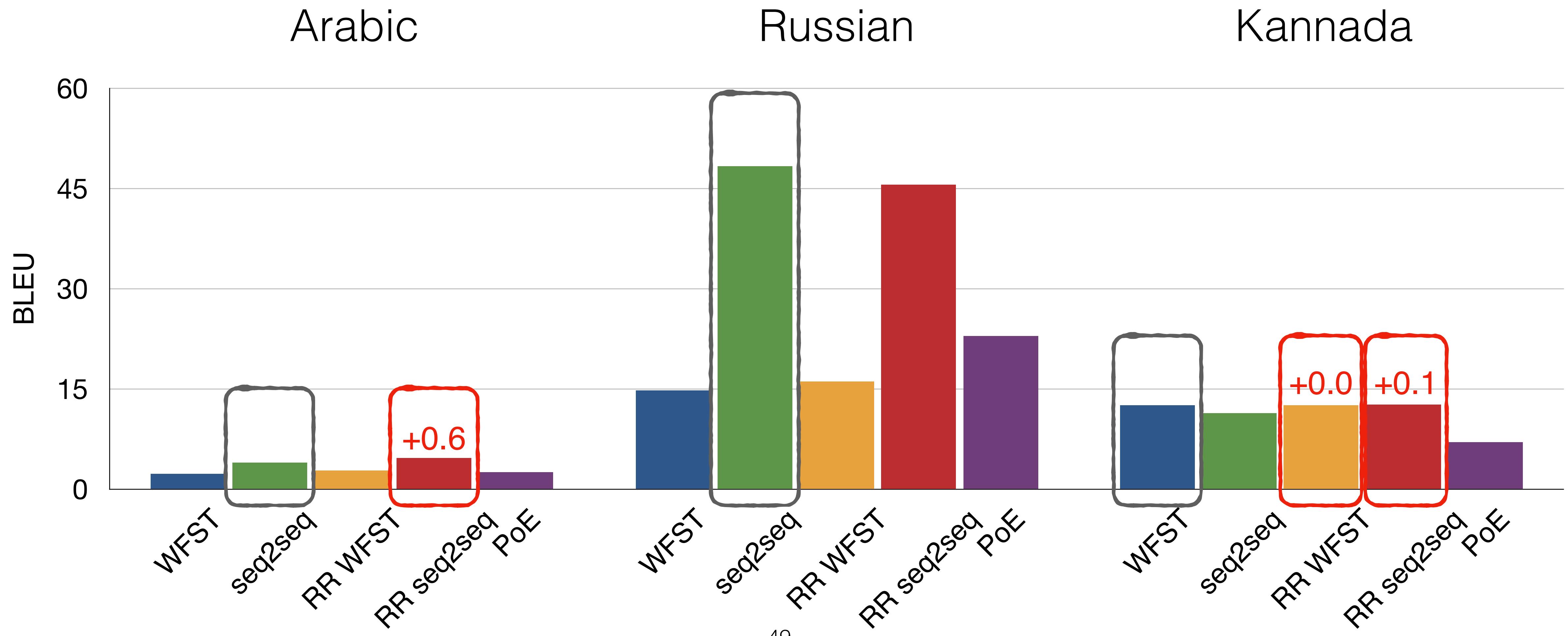
Kannada



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 осуϲчествлиниые "бор#би с коммунизмом" в уузнани америке. | kongress ne odobril viudet dla osuſčestvleniye "bor#bi s kommunizmom" v uuznani amerike. | |
| Reranked WFST | конгресс не одобрил видет дела осуϲчествлиниые "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 осуществлиниые "борьбы с коммунизмом" в уузник амери | kongress ne odobril b1d et 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**

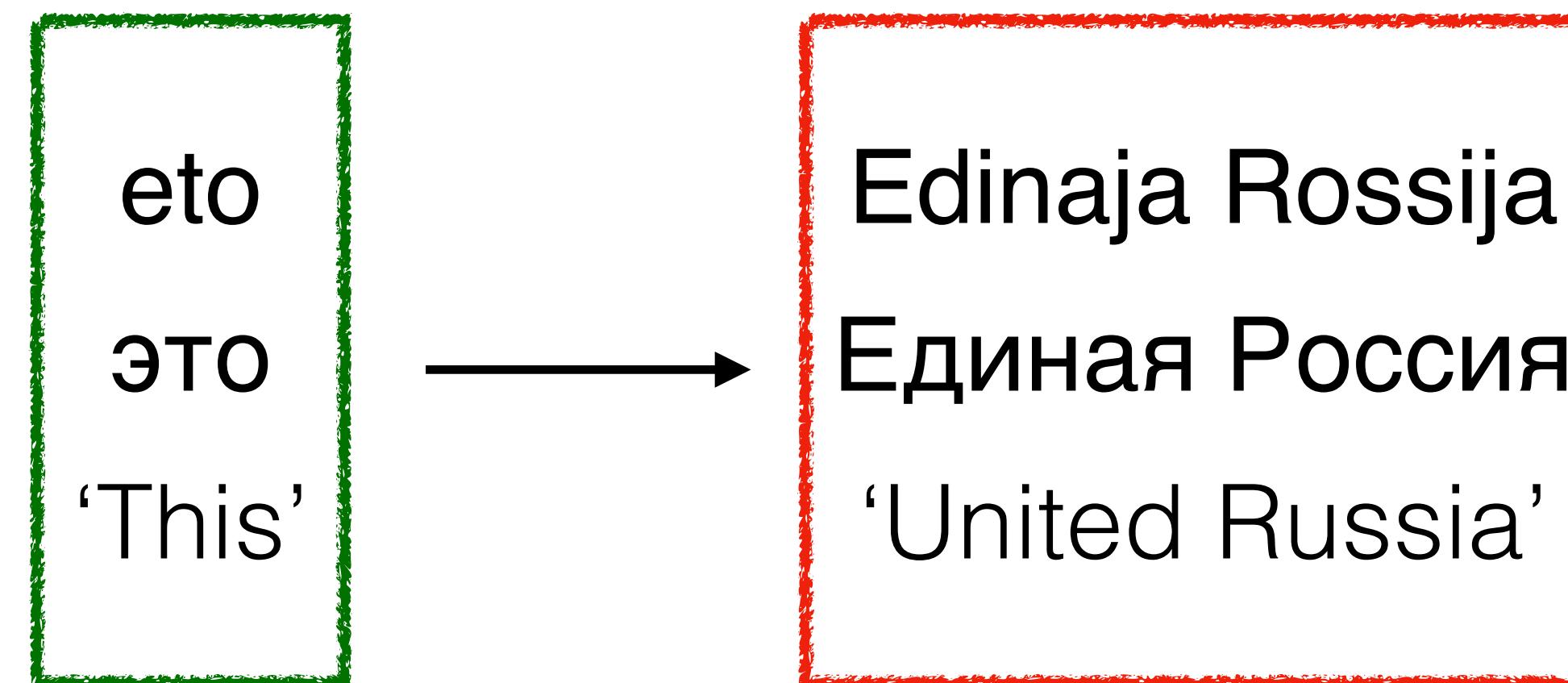
482.0, **единороссы уже #страшно сказать) старый ро**

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

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

High-level takeaways

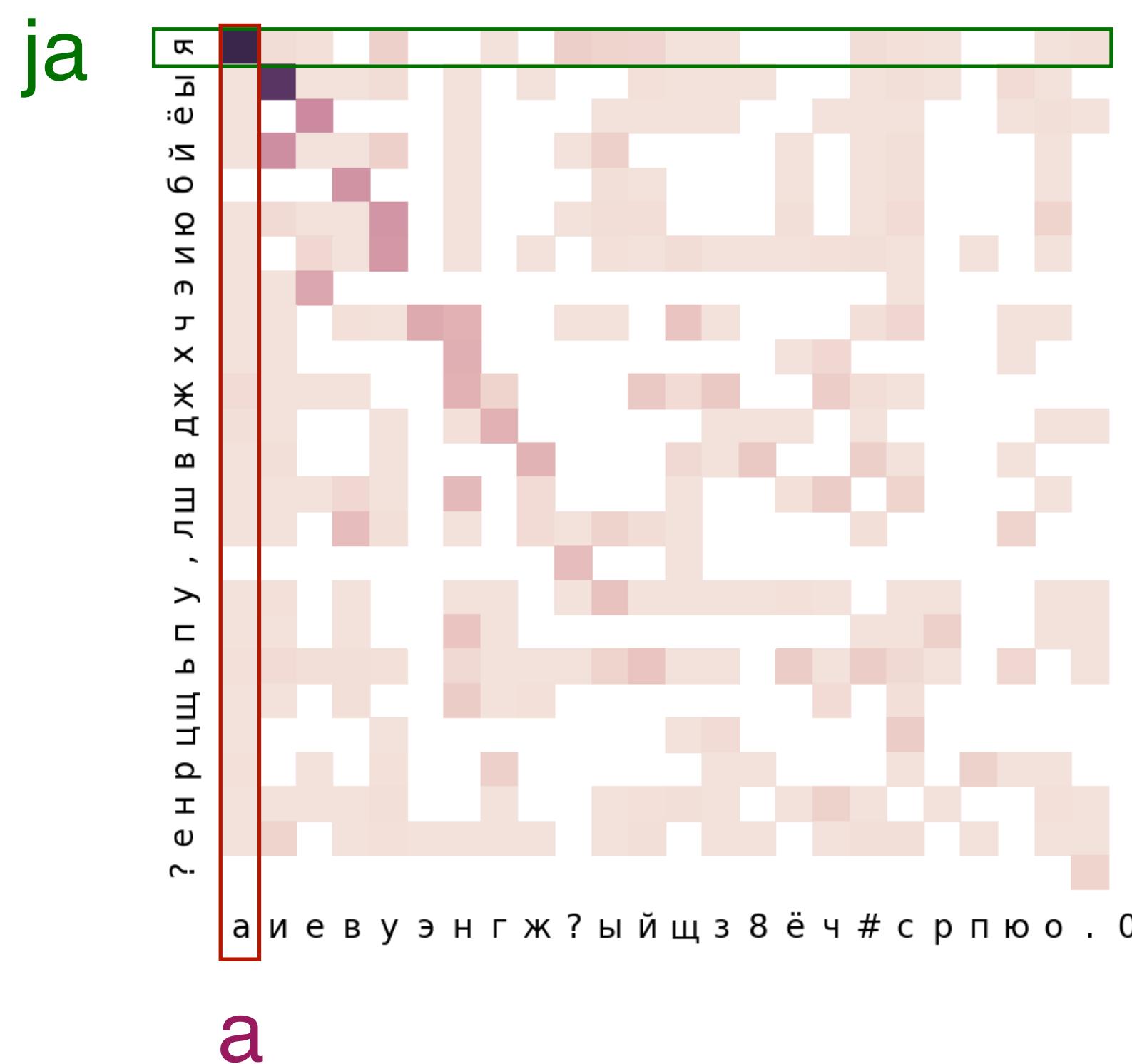
- Seq2seq is more sensitive to **distributional shifts**
 - Remember that our Cyrillic data comes from political discussion groups
 - 25% of common word-level errors in seq2seq are of this type!



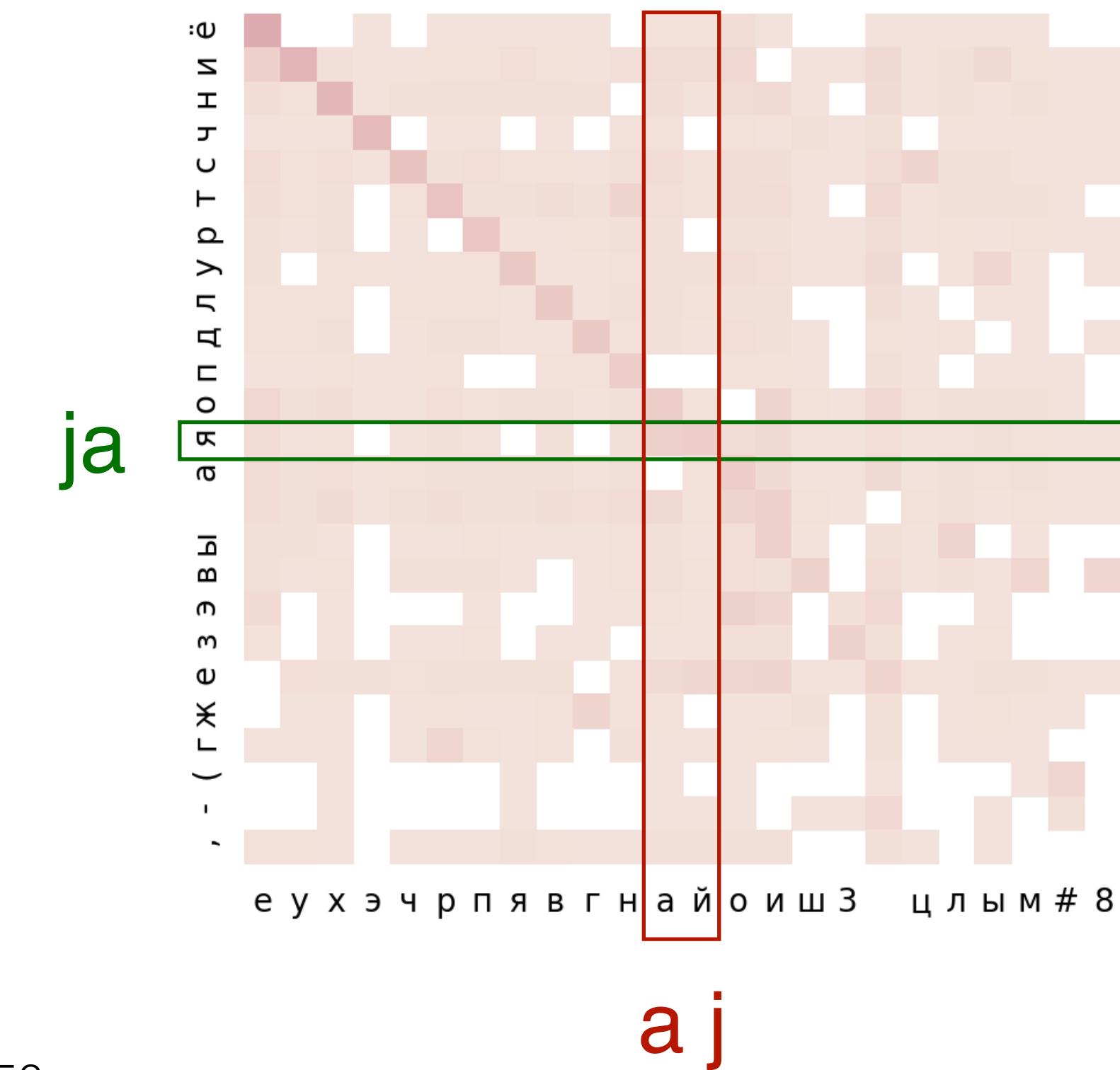
High-level takeaways

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

WFST



seq2seq



Final points

- Linguistic typology is important!
 - Intuitions from one language don't necessarily extend to others
- Normalization of creative phenomena is lossy
 - Ideally, annotation should be performed with feedback from the author
- FSTs and seq2seq models have complementary strengths
 - Our decoding-time combinations didn't help much, but there's hope!
 - It could be joint training...
 - ...Or holistic structural combinations...
 - ...Or 'softer' biasing of one model towards another model's behavior

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.
- J. Eisner. [Parameter estimation for probabilistic finite-state transducers](#). ACL 2002.
- 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](#). SIGMORPHON 2021.
- T. Shavrina, O. Shapovalova. [To the methodology of corpus construction for machine learning: Taiga syntax tree corpus and parser](#). CORPORA 2017.