

Carnegie Mellon University

Unlocking Resources for Under-resourced Languages

Graham Neubig

Shruti Rijhwani, Xinyi Wang



Antonios Anastasopoulos, Daisy Rosenblum,
Sebastian Ruder

Language technologies in a multilingual world

Language technologies in a multilingual world

Considerable recent progress in expanding NLP to many languages!

Language technologies in a multilingual world

Considerable recent progress in expanding NLP to many languages!

Multilingual benchmark
datasets for NLP tasks

WMT MasakhaNER
TyDiQA XNLI FLORES-101
UD Treebank XTREME

Language technologies in a multilingual world

Considerable recent progress in expanding NLP to many languages!

Multilingual benchmark datasets for NLP tasks

WMT MasakhaNER
TyDiQA XNLI FLORES-101
UD Treebank XTREME

Pretrained multilingual language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Language technologies in a multilingual world

Considerable recent progress in expanding NLP to many languages!

Multilingual benchmark datasets for NLP tasks

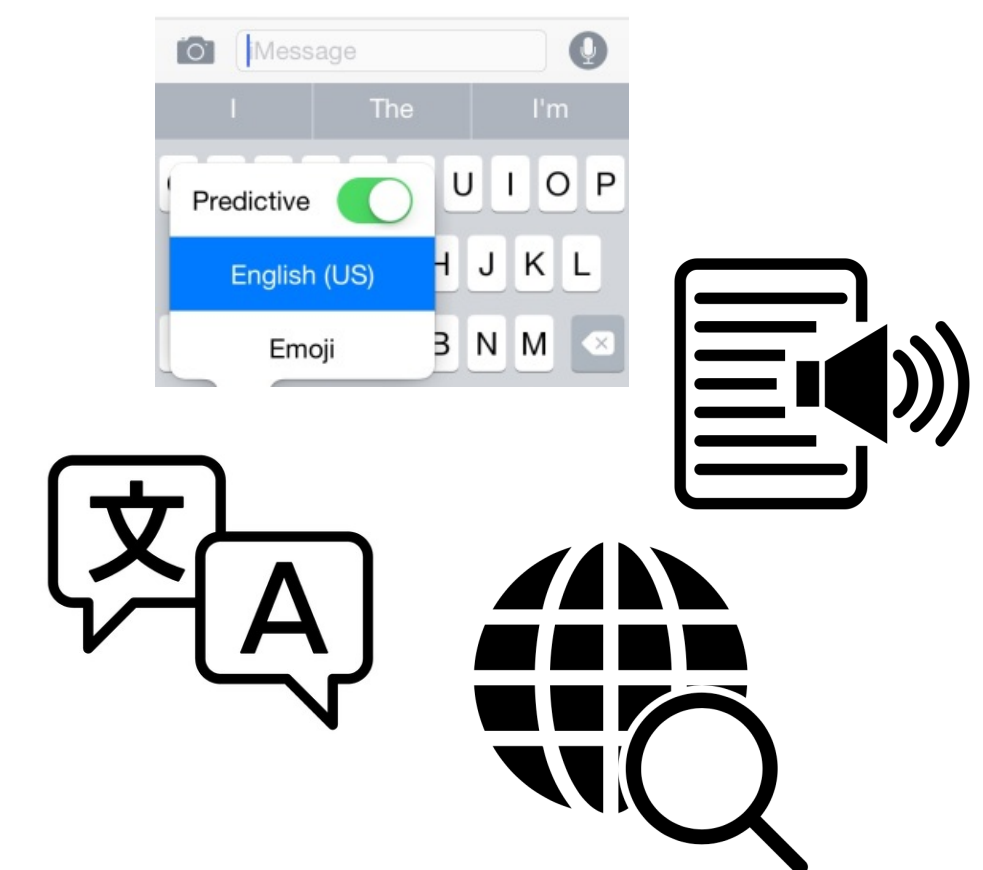
WMT MasakhaNER
TyDiQA XNLI FLORES-101
UD Treebank XTREME

Pretrained multilingual language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Commercial models that support many languages

Voice assistants
Predictive keyboards Translation
Web search Text analytics



Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Trained on unlabeled text corpora
(e.g., Wikipedia and Common Crawl)

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages



Enables NLP applications
through cross-lingual transfer

- Named entity recognition
- Entity linking
- Web search
- Machine translation
- Question answering

...

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

Multiple societal benefits of NLP
that includes many languages!

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

Wikipedia

From Wikipedia, the free encyclopedia

*This article is about the online encyclopedias in
other languages, see [List of Wikipedias](#).*

Wikipedia (/ˌwɪkɪˈpiːdiə/ [ⓘ] [ⓘ] listen) *wik-ih-PEE-encyclopedia* written and maintained by a co

Automatic translation
of information

Multiple societal benefits of NLP
that includes many languages!

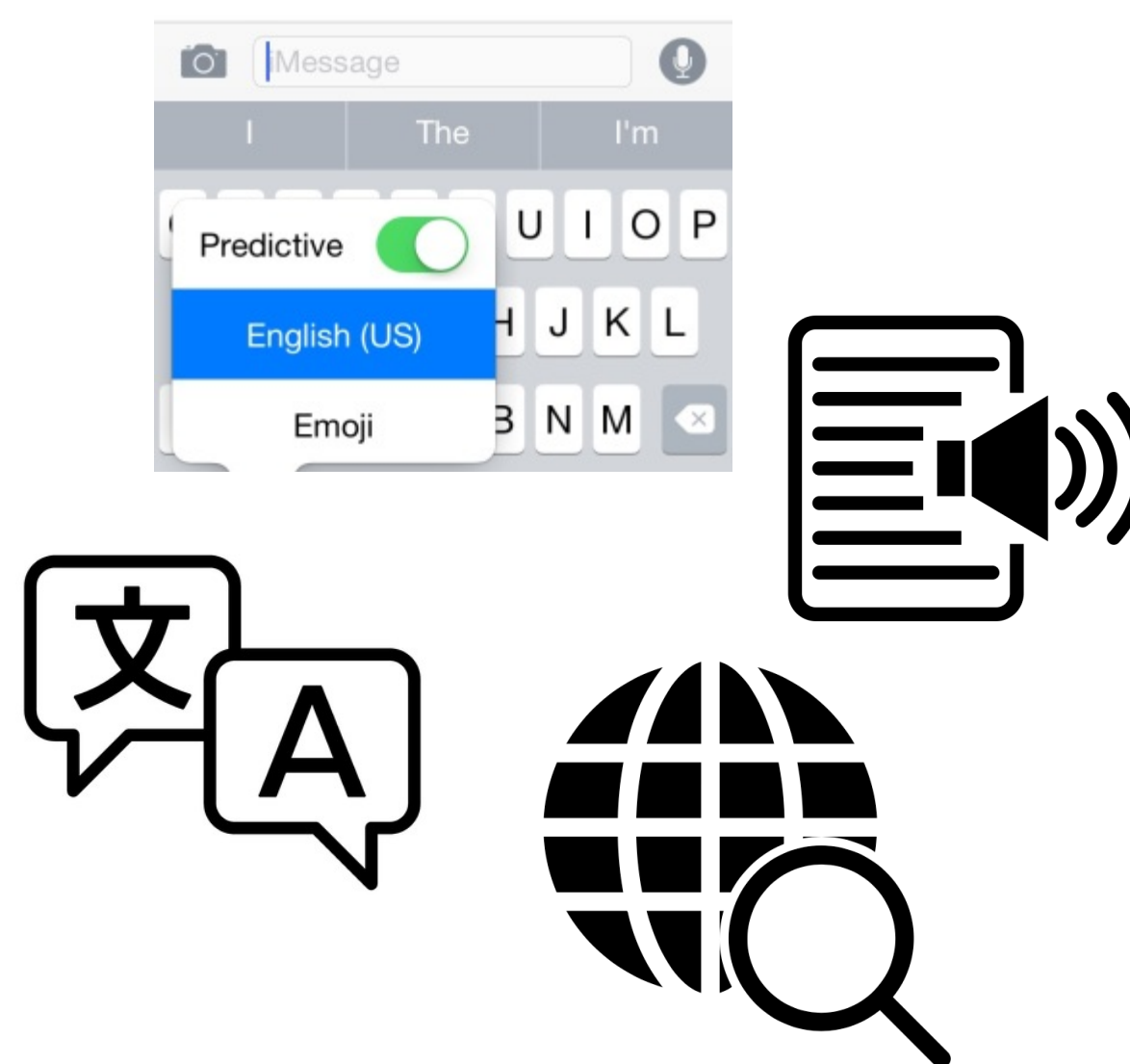
Access to information and
education from other languages

Language technologies in a multilingual world

Pretrained multilingual
language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages



Multiple societal benefits of NLP
that includes many languages!

Access to information and
education from other languages

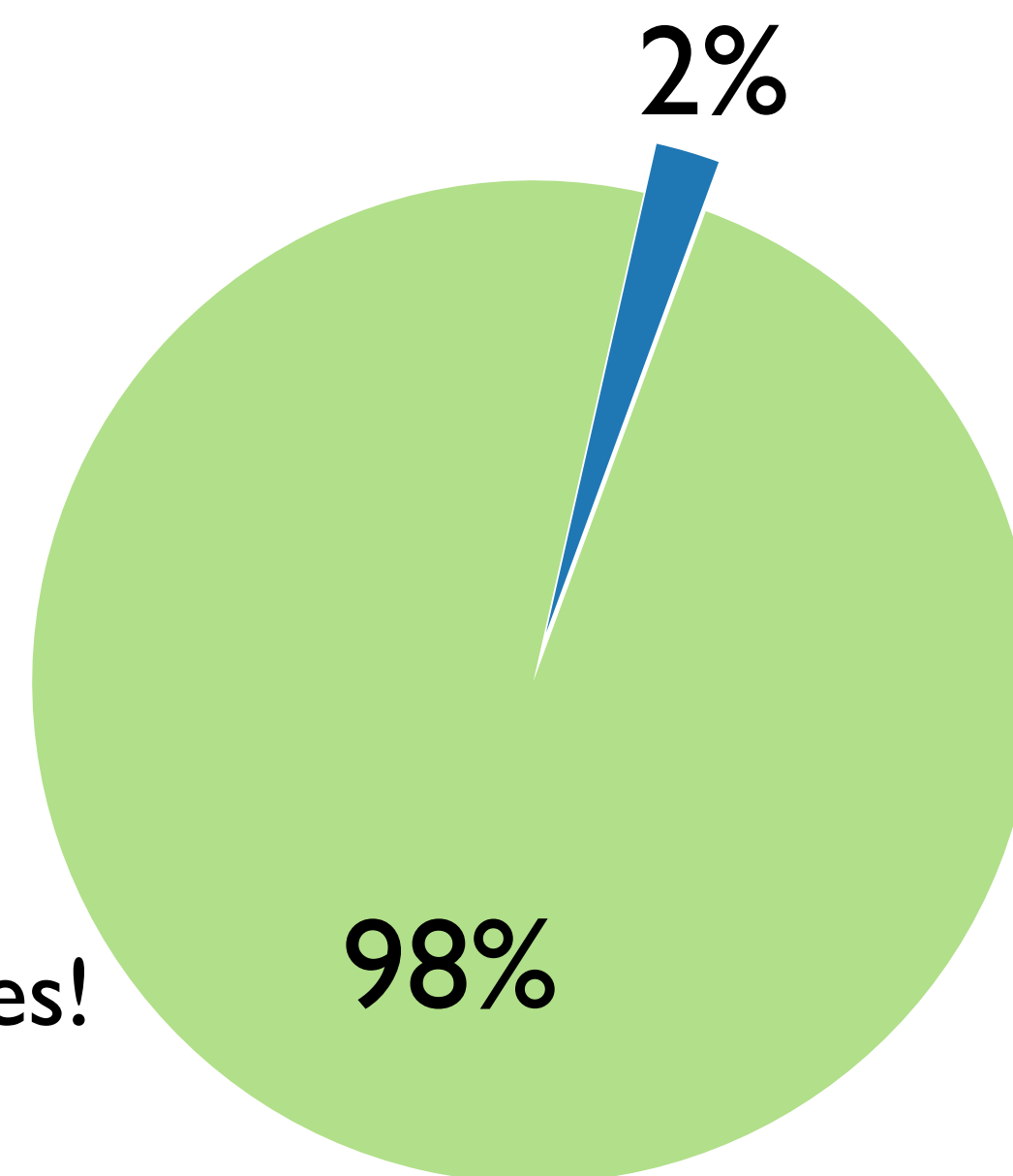
Language technologies that serve
many more people!

Language technologies in a multilingual world

Pretrained multilingual language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages



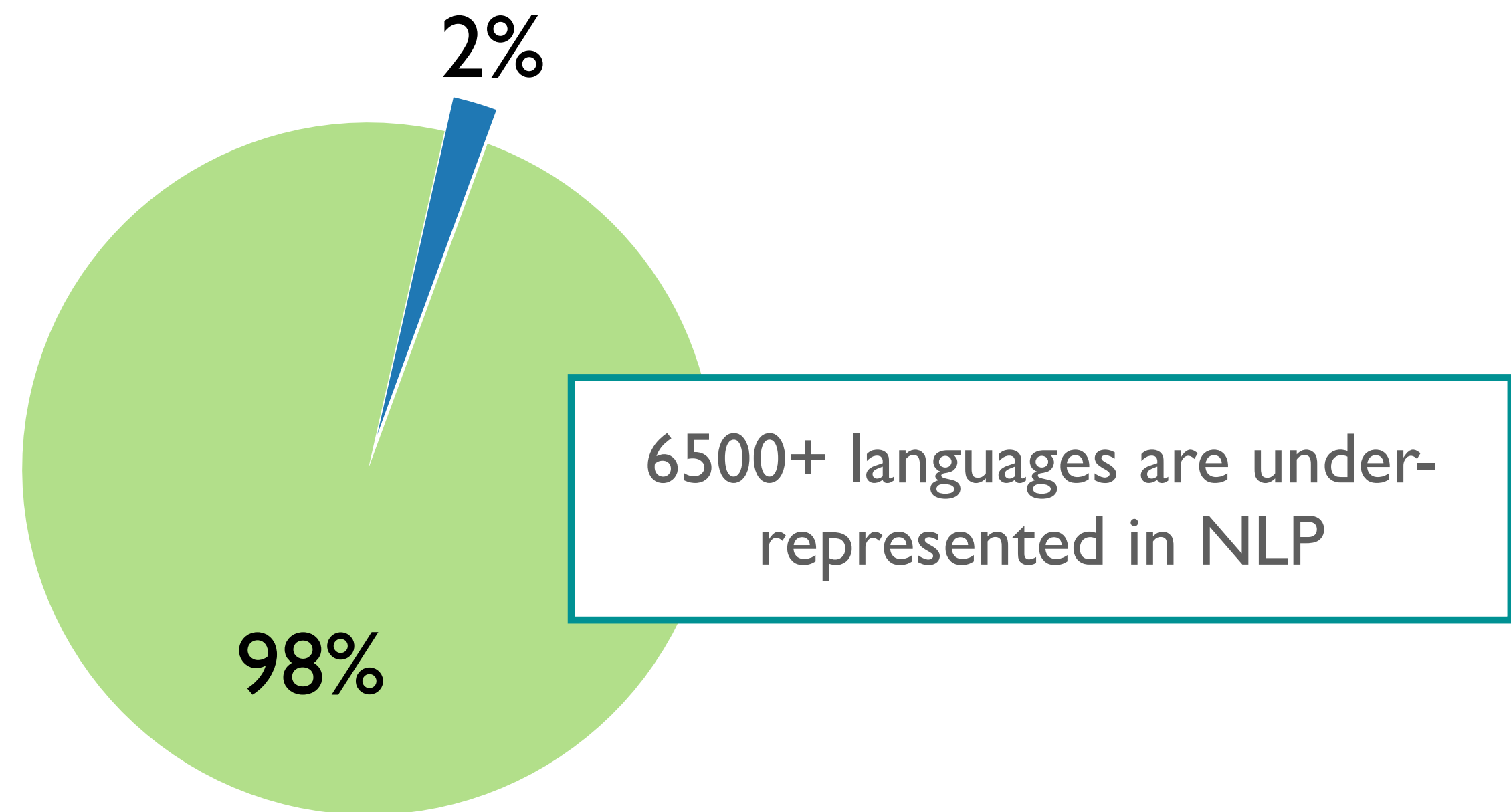
There are over 7000 living languages!

Language technologies in a multilingual world

Pretrained multilingual language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

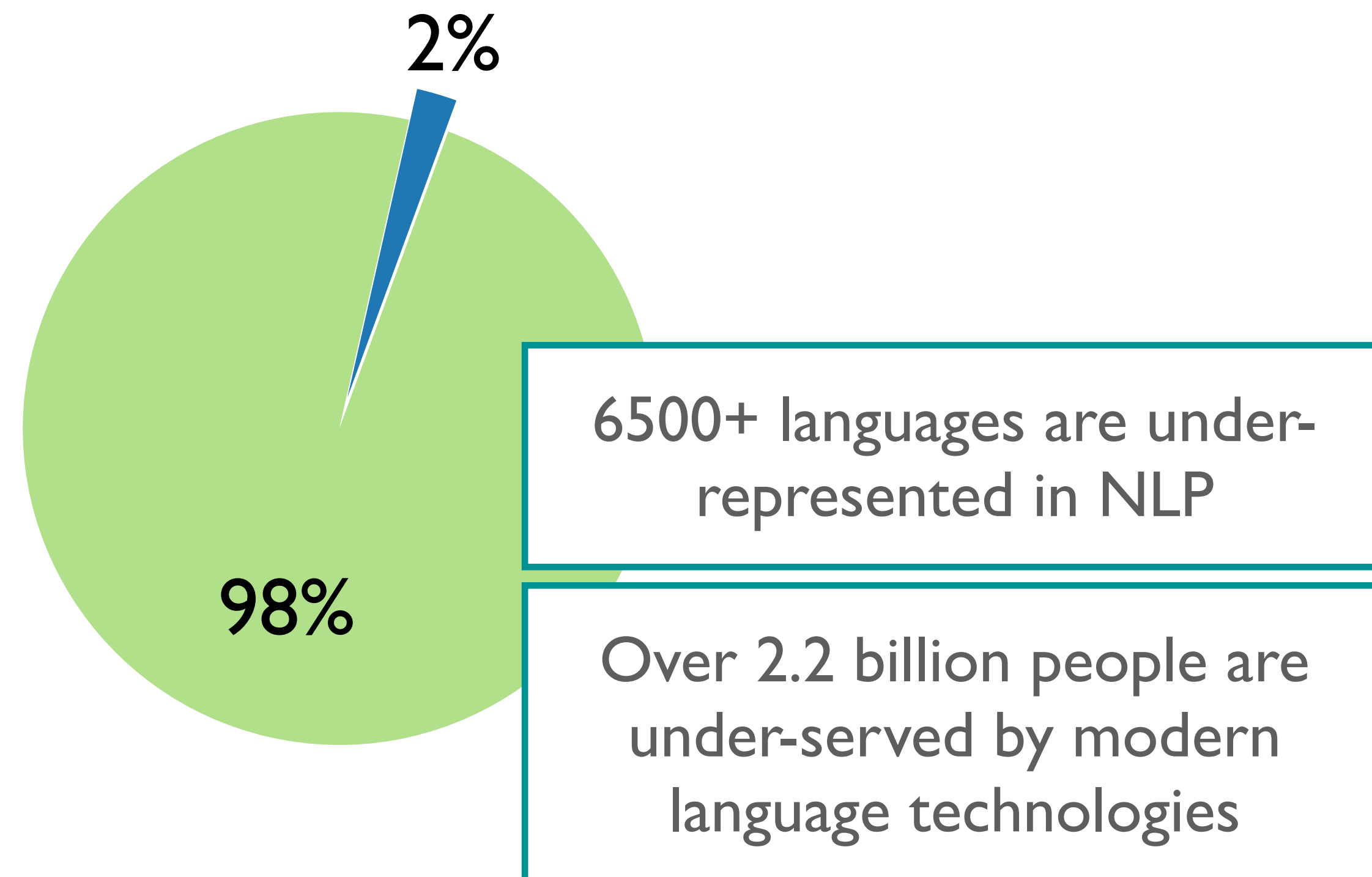


Language technologies in a multilingual world

Pretrained multilingual language models

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

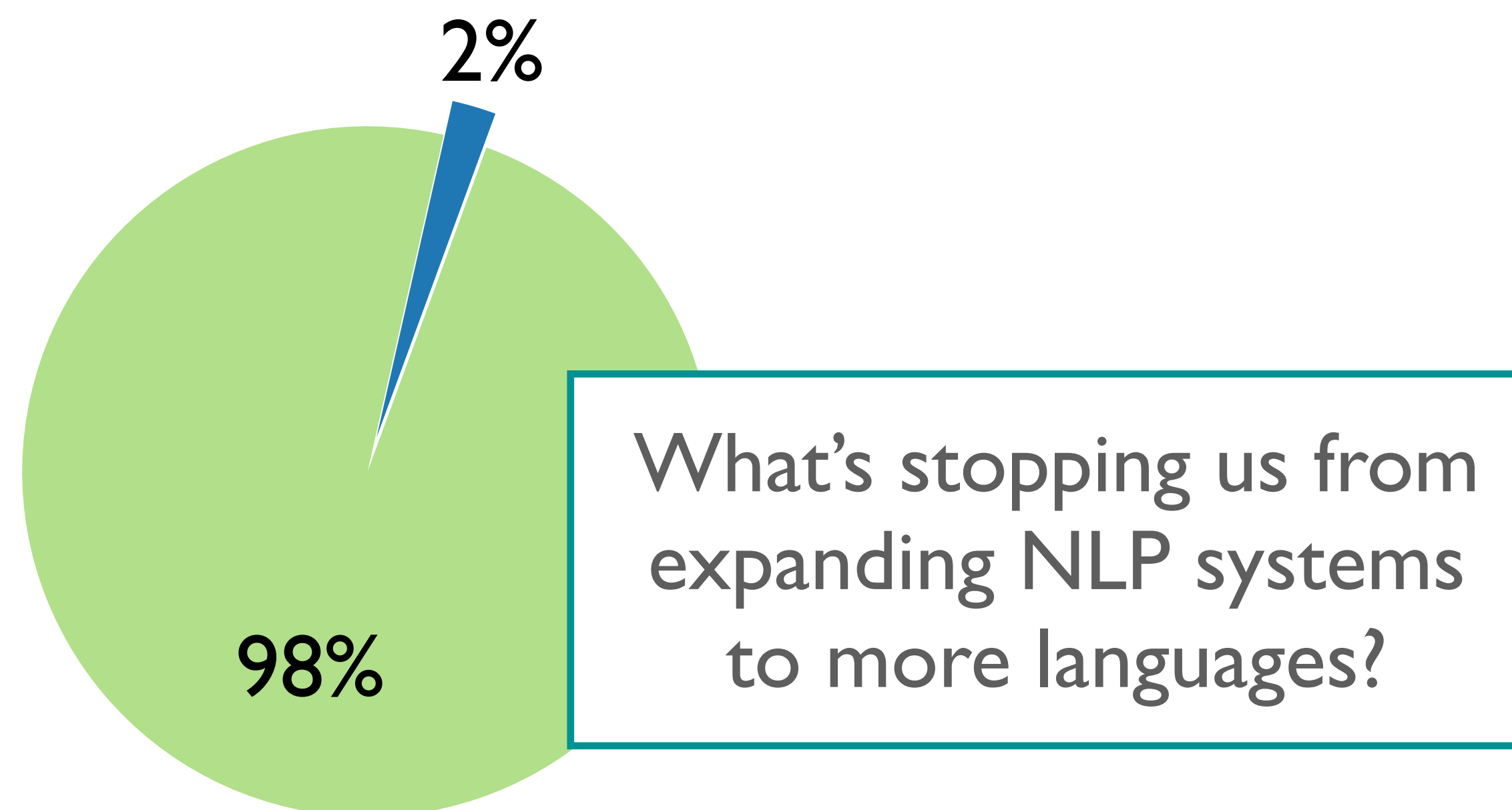


Language technologies in a multilingual world

Pretrained multilingual
language models

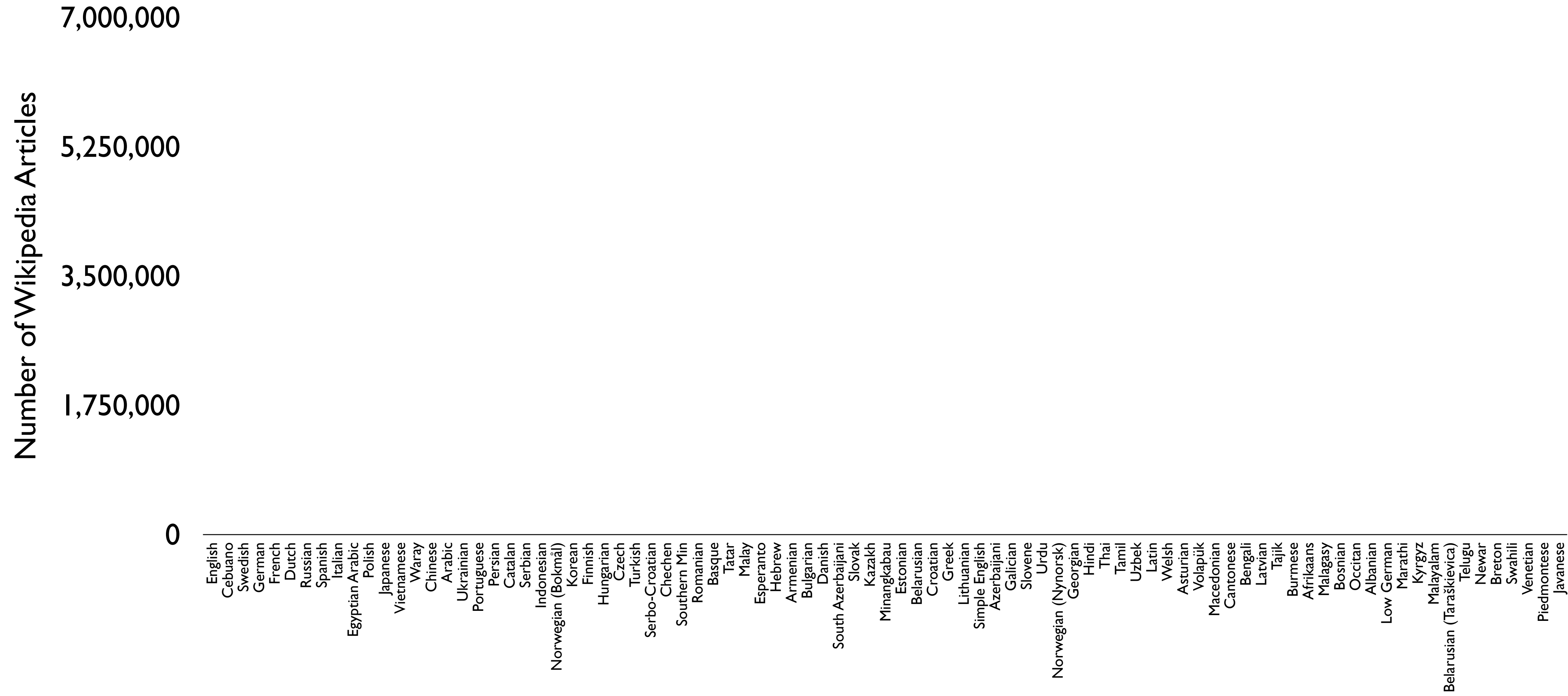
XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Support 100 – 200 languages

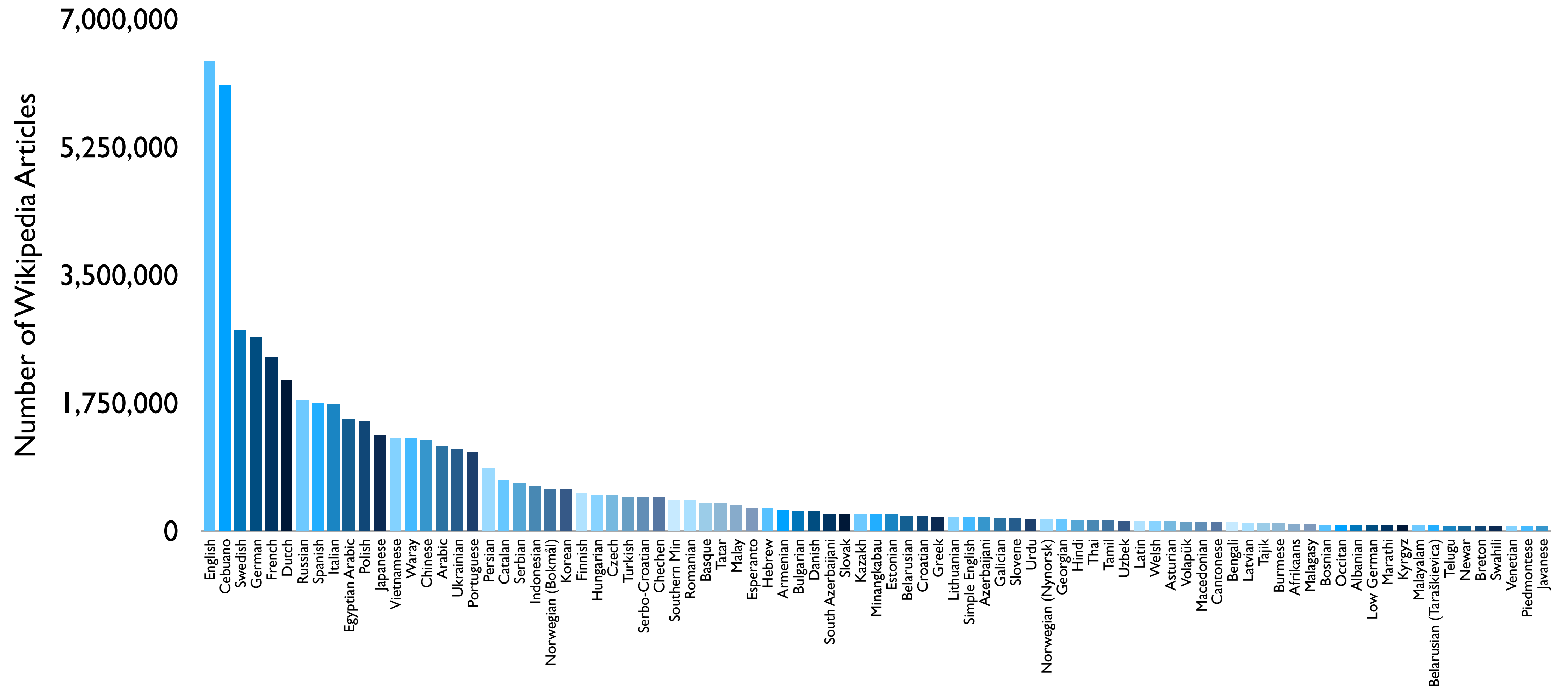


The unlabeled text bottleneck

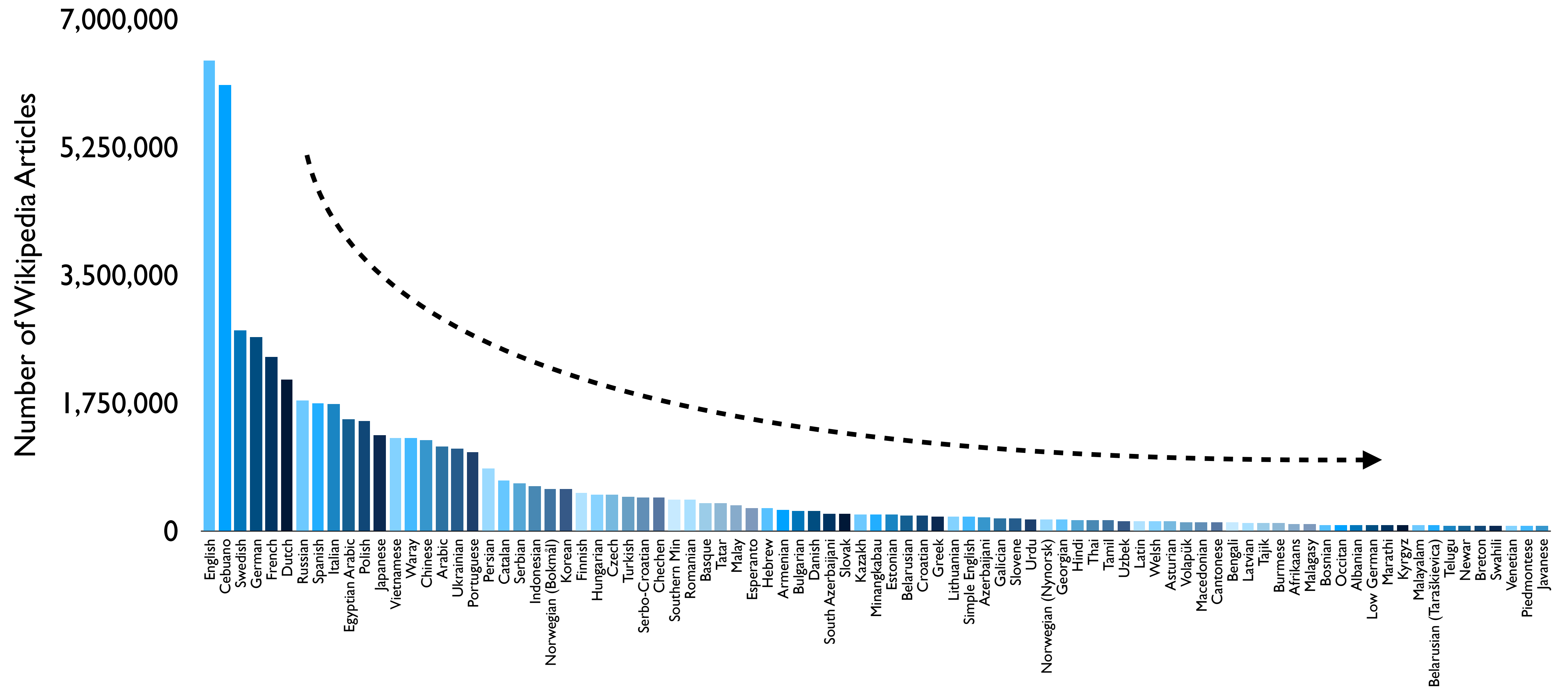
The unlabeled text bottleneck



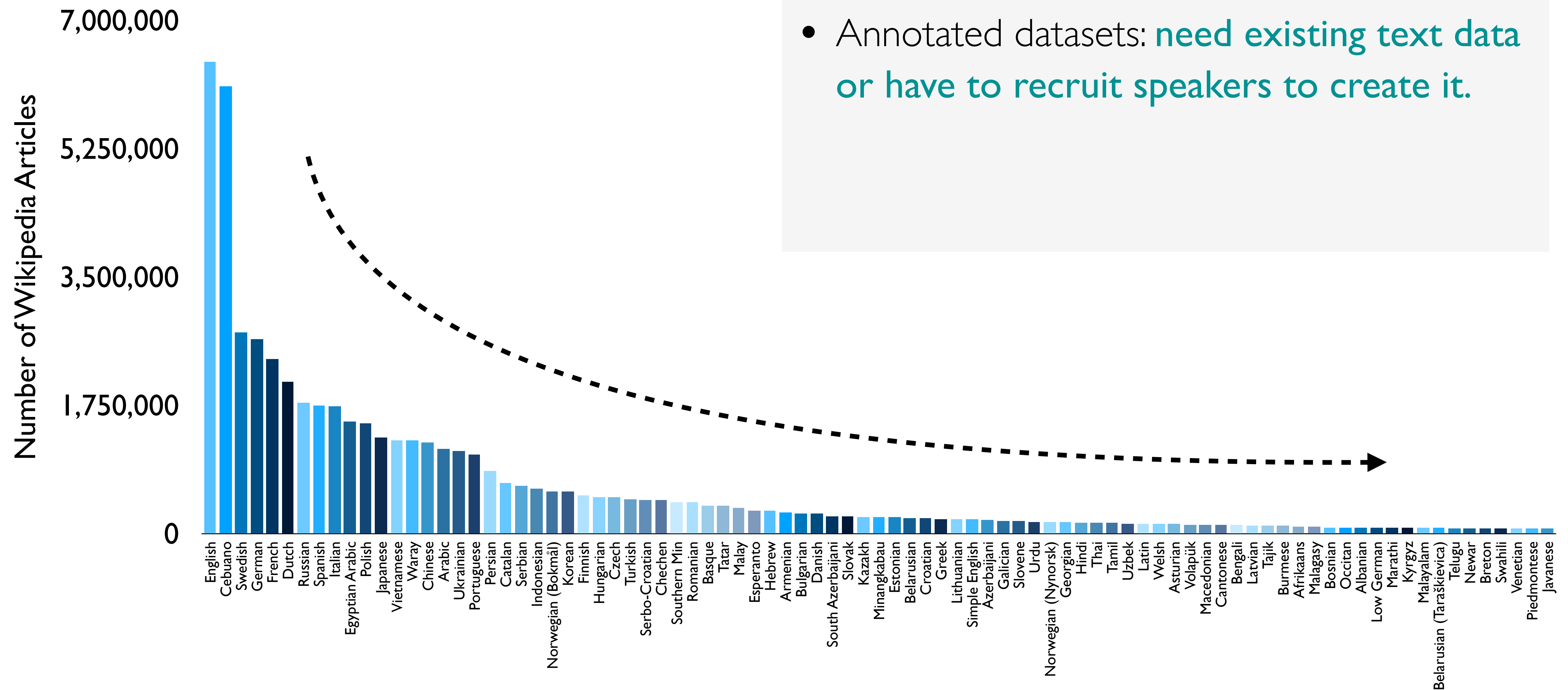
The unlabeled text bottleneck



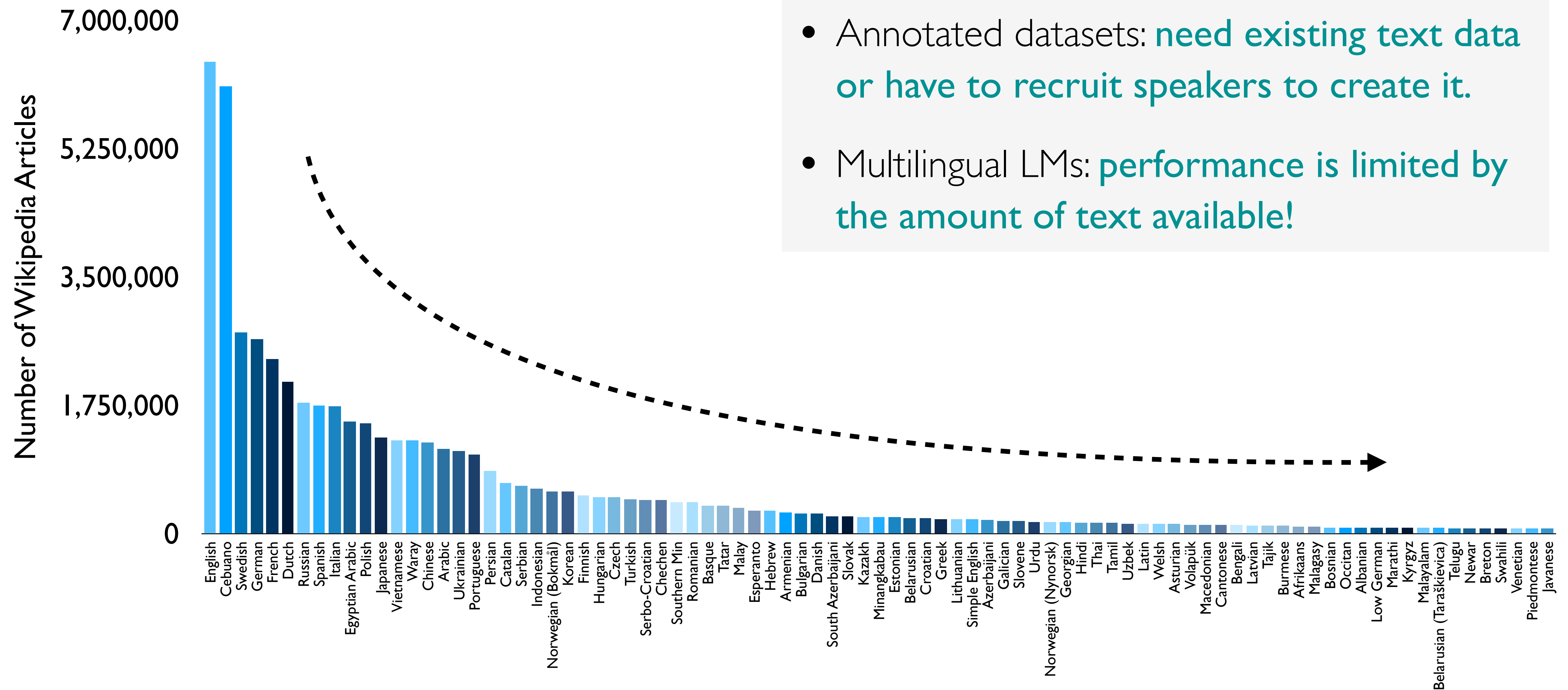
The unlabeled text bottleneck



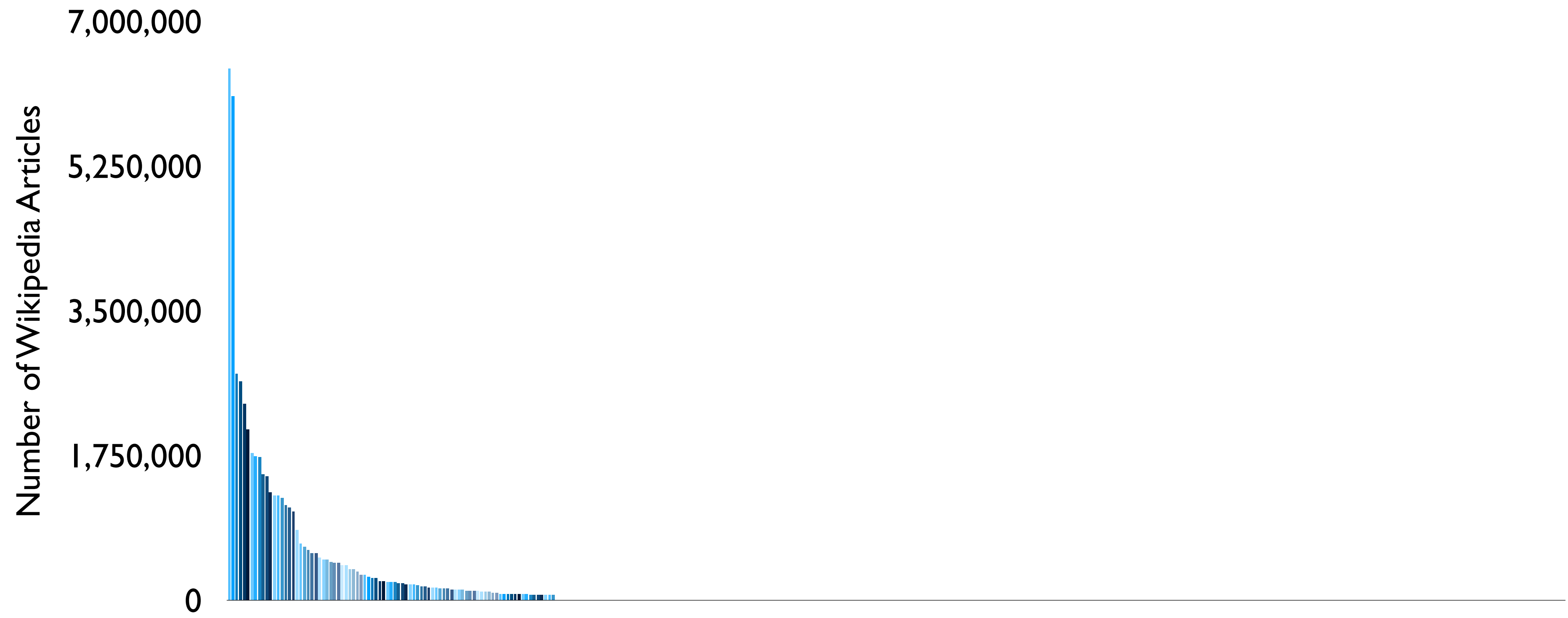
The unlabeled text bottleneck



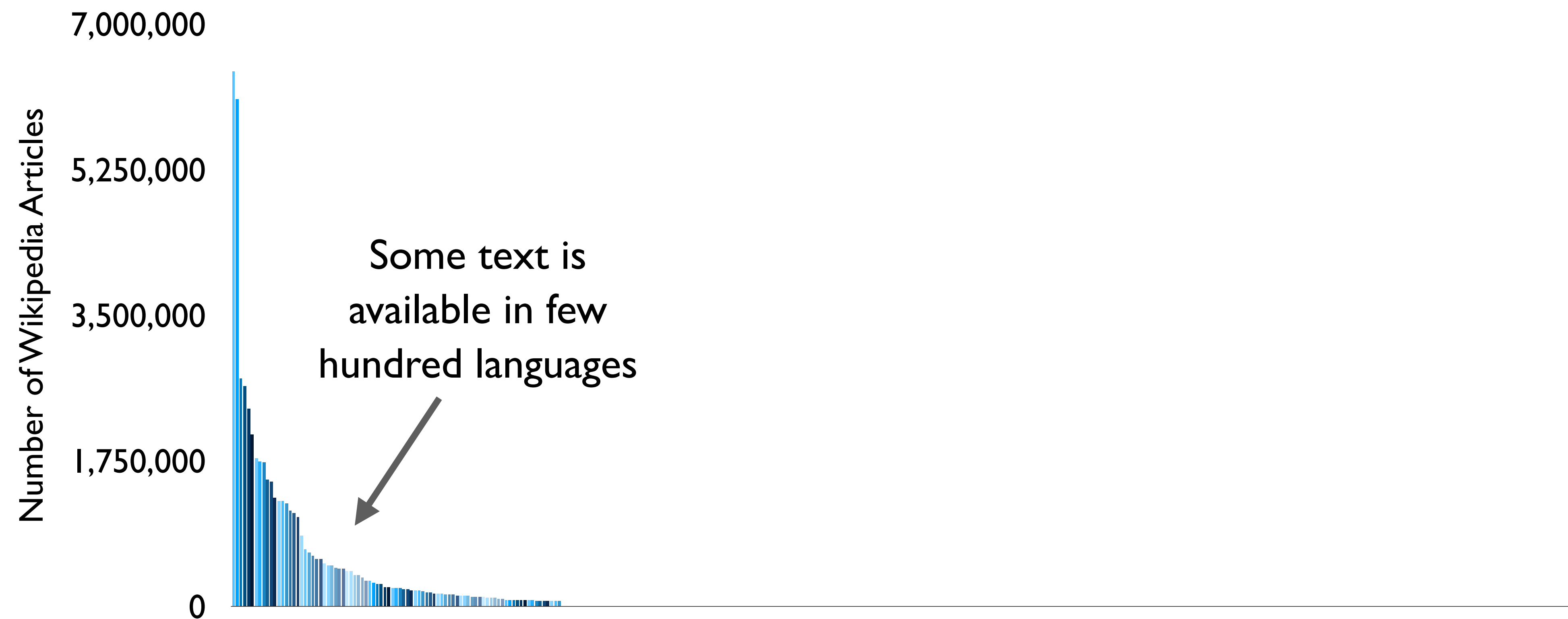
The unlabeled text bottleneck



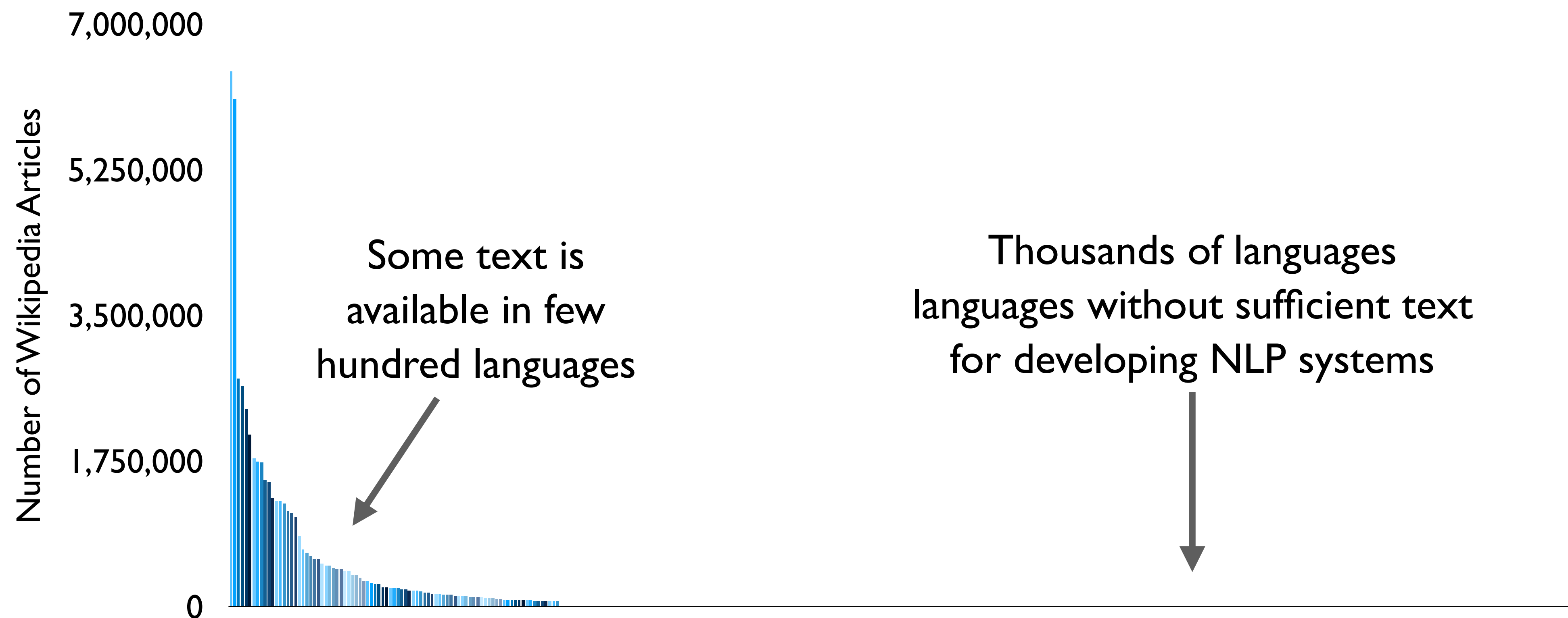
The unlabeled text bottleneck



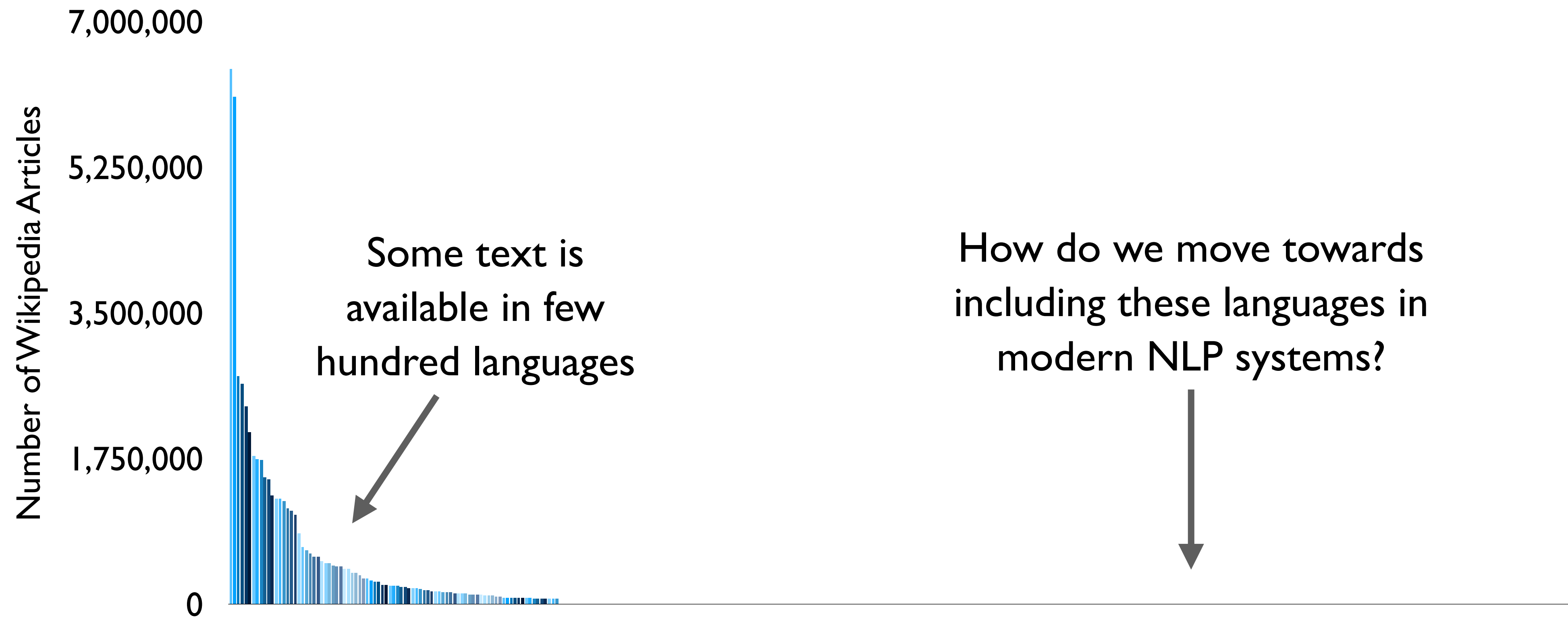
The unlabeled text bottleneck



The unlabeled text bottleneck



The unlabeled text bottleneck



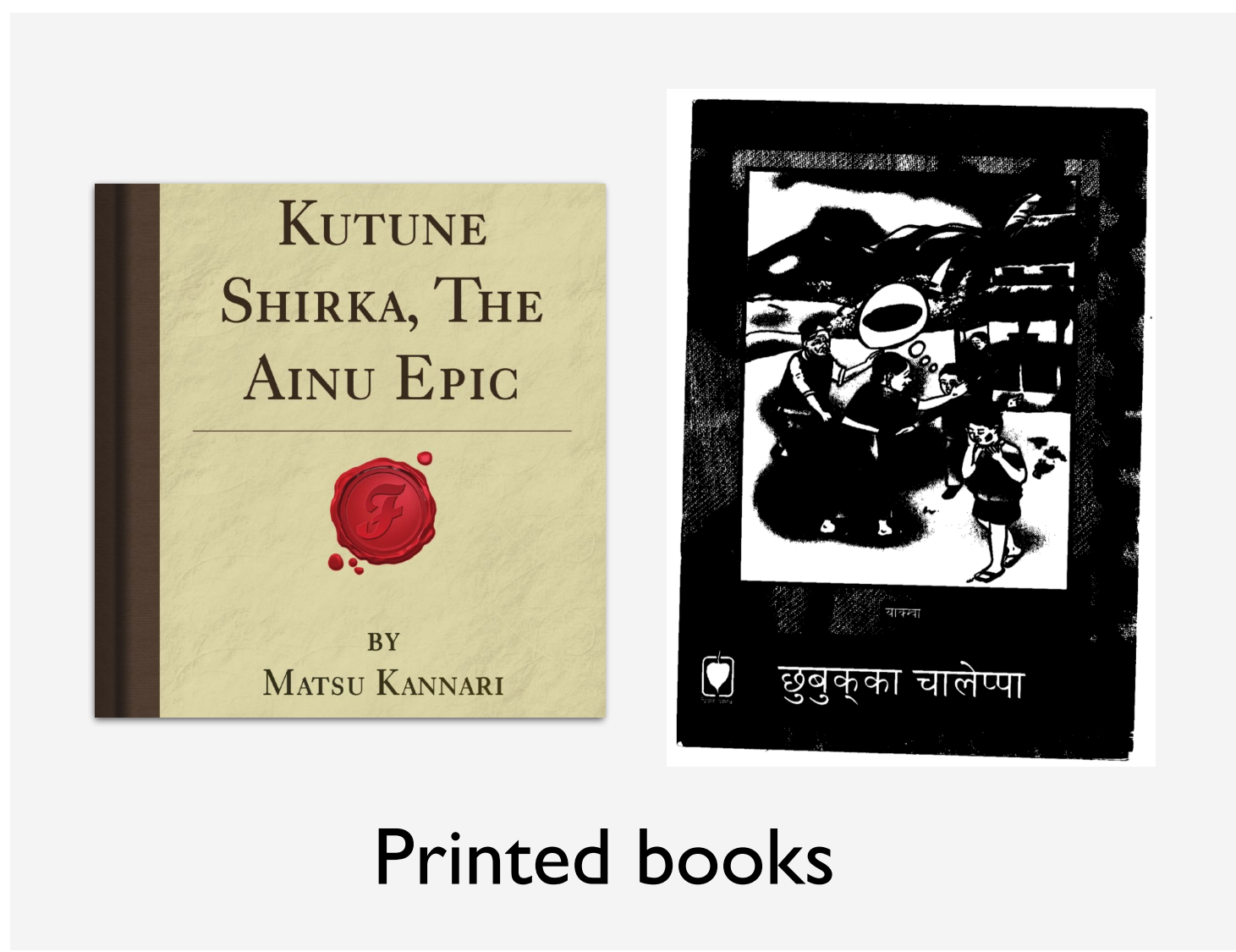
Text resources do exist in many more languages!

Text resources do exist in many more languages!

But locked away in formats that are not machine-readable

Text resources do exist in many more languages!

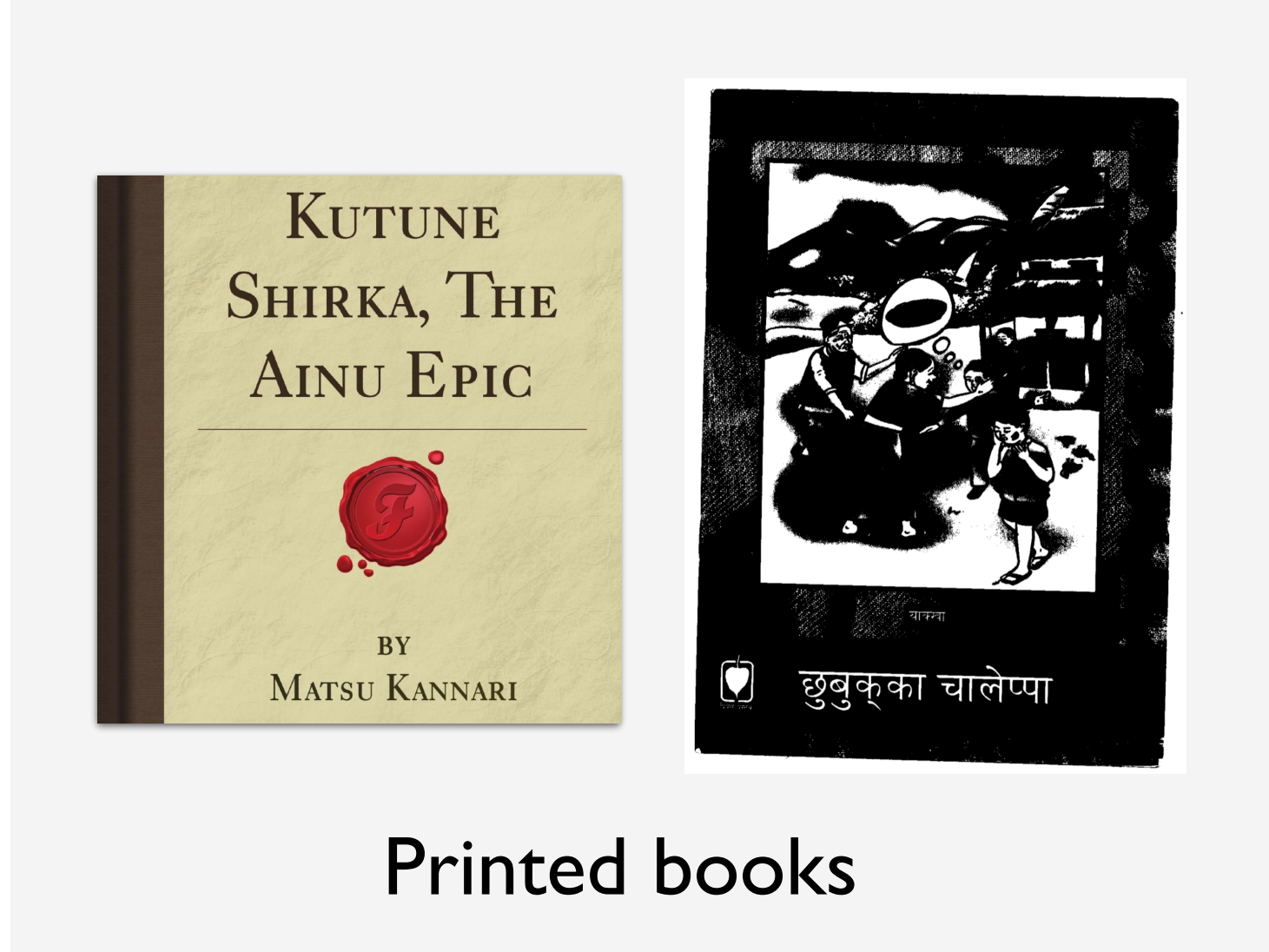
But locked away in formats that are not machine-readable



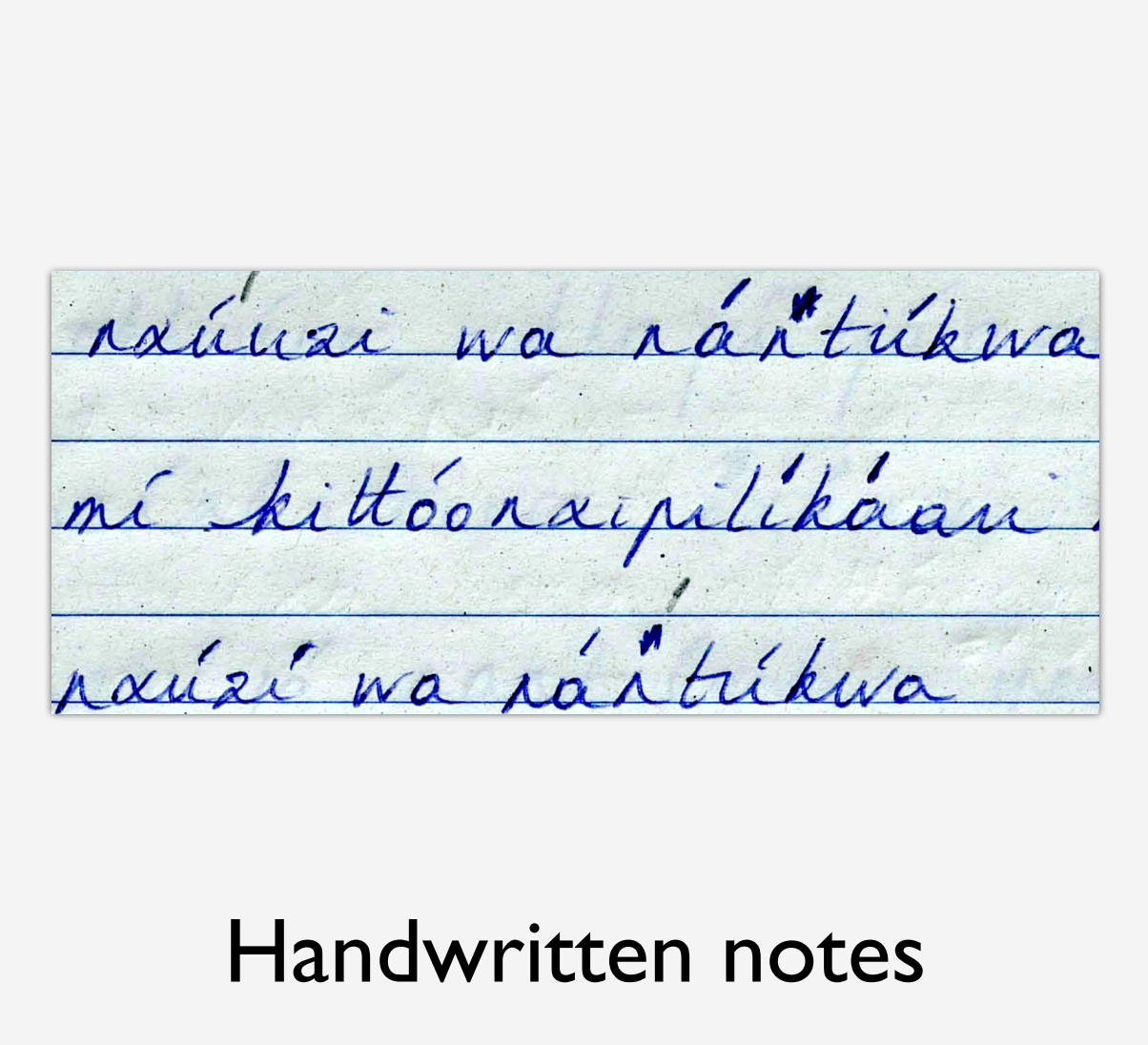
Printed books

Text resources do exist in many more languages!

But locked away in formats that are not machine-readable



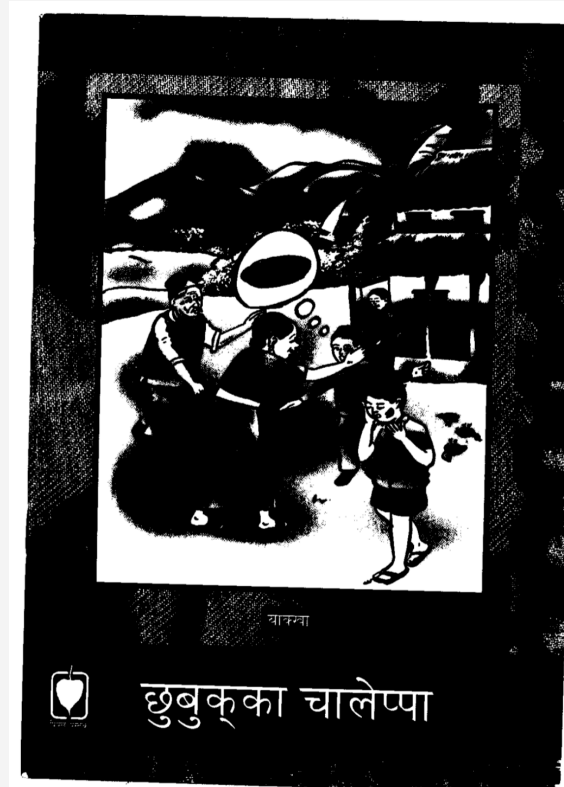
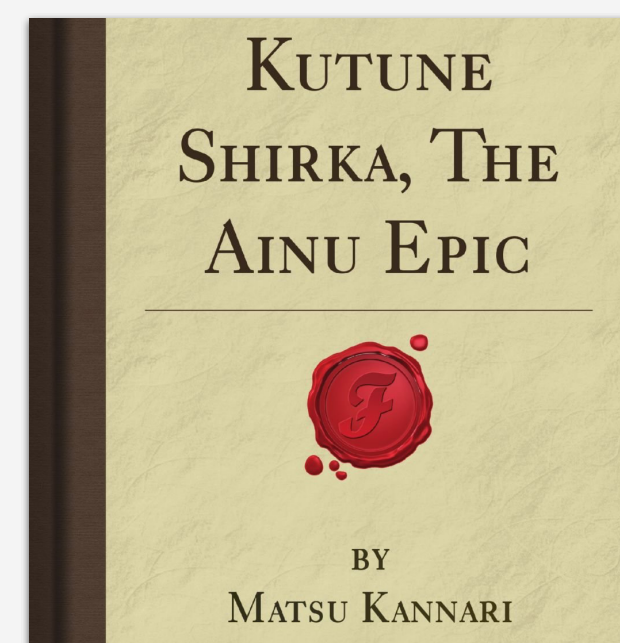
Printed books



Handwritten notes

Text resources do exist in many more languages!

But locked away in formats that are not machine-readable



Printed books

naúuzi wa nántúkwá
mí kittóonaxipilikáan
naúuzi wa nántúkwá

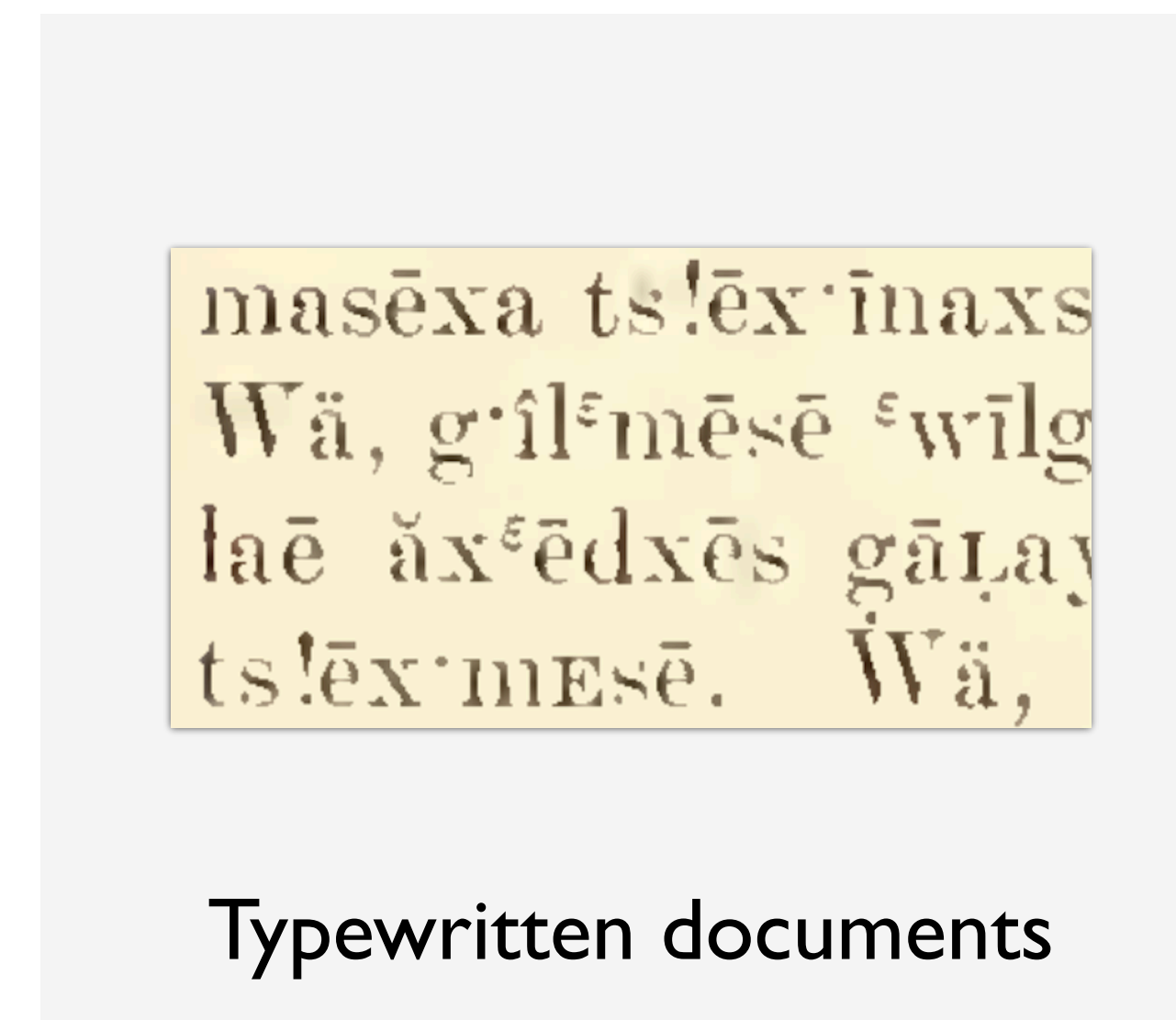
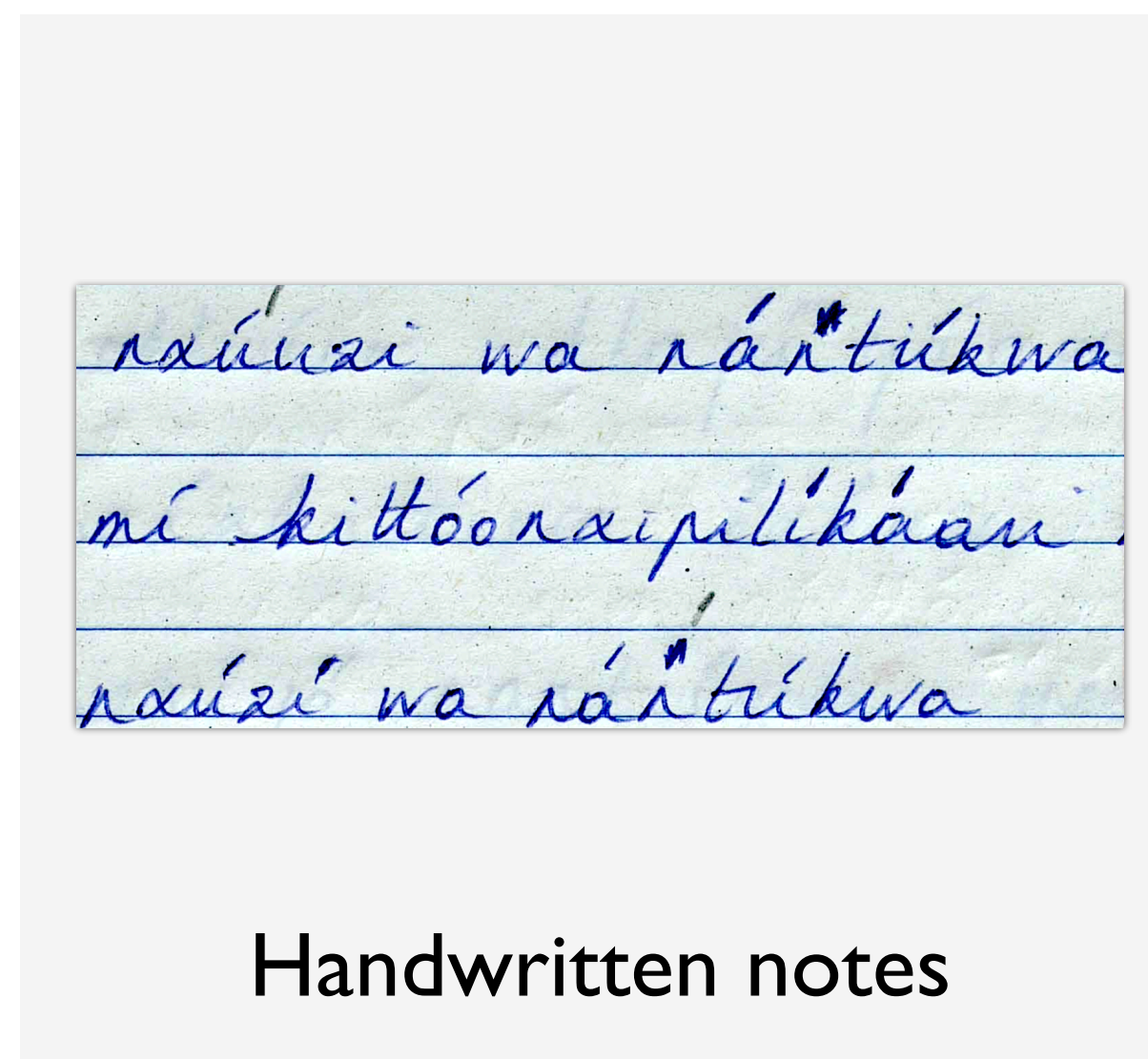
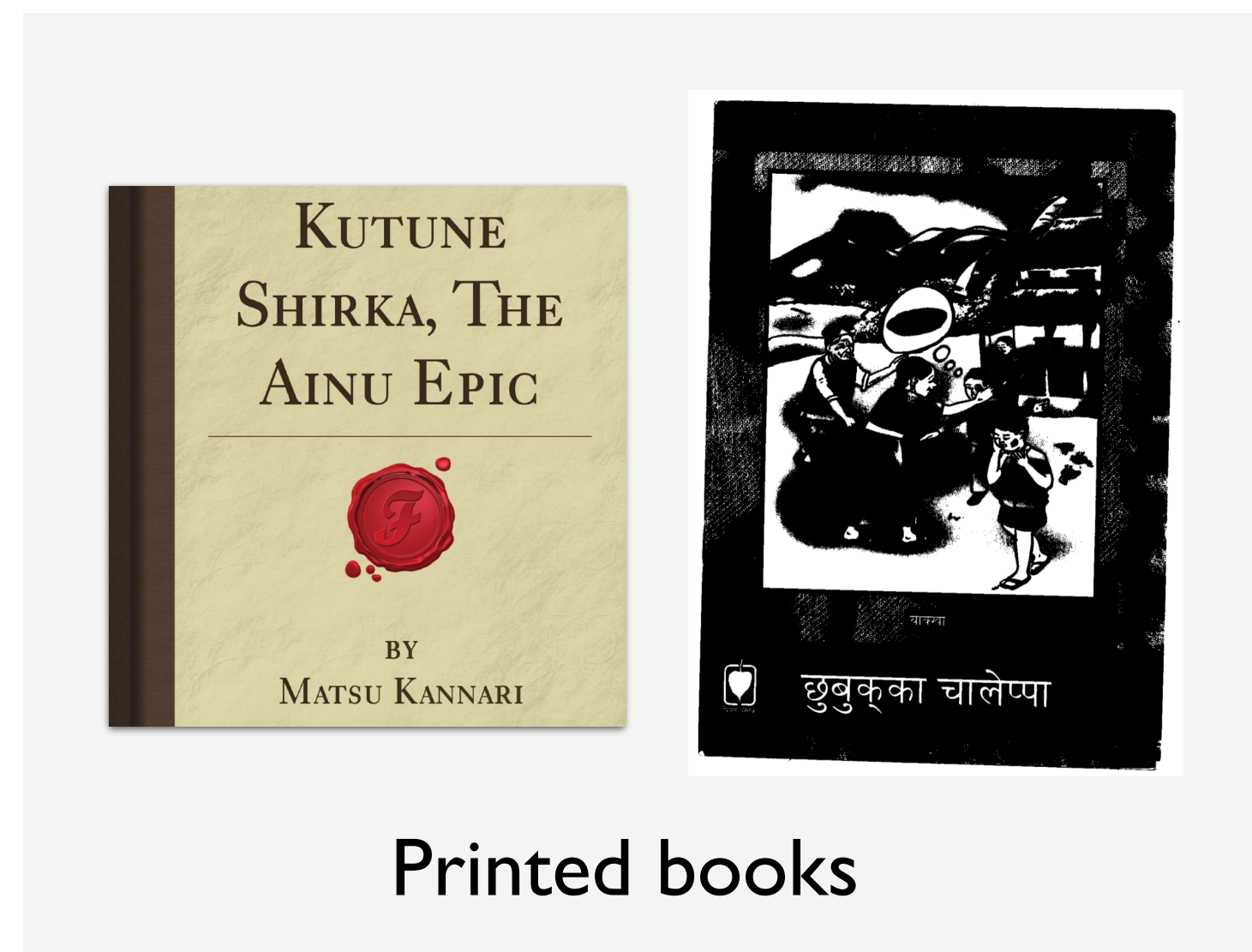
Handwritten notes

masēxa ts!ēx·inaxs
 Wä, g·il^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mēsē. Wä,

Typewritten documents

Text resources do exist in many more languages!

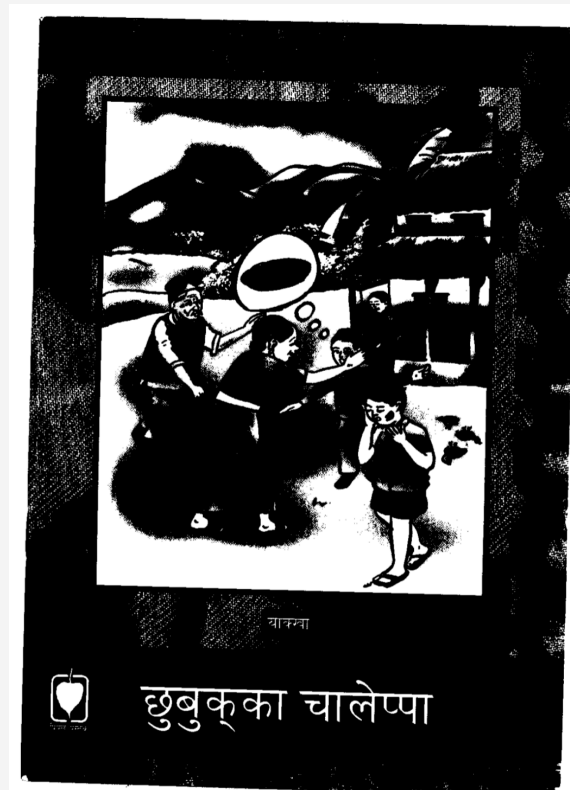
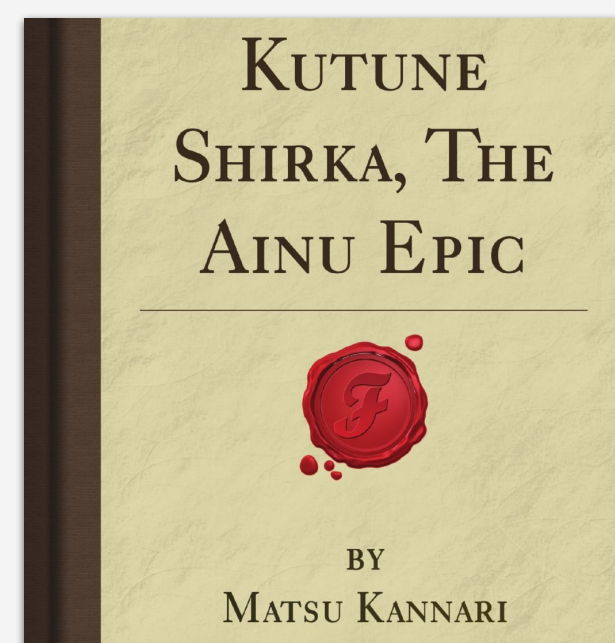
But locked away in formats that are not machine-readable



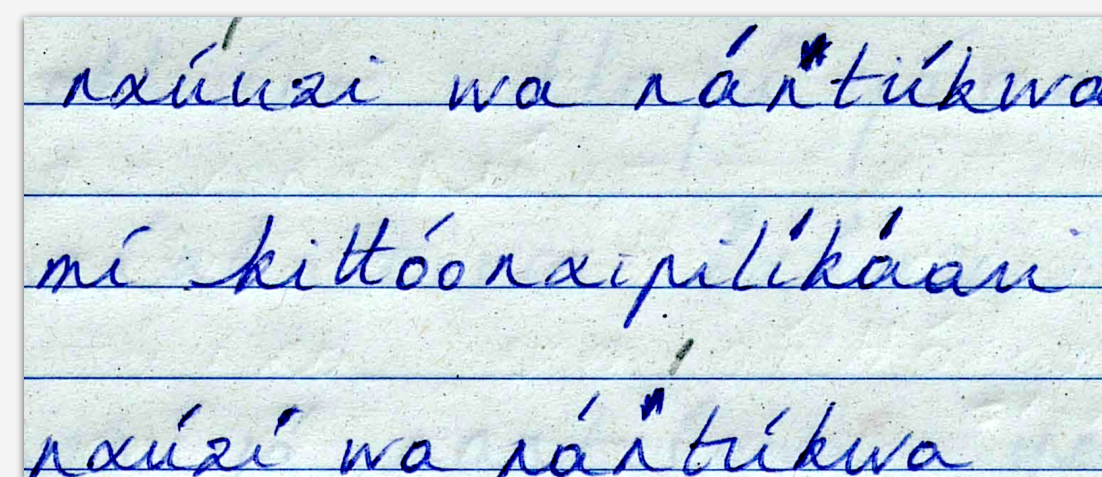
or in other formats such as bilingual lexicons

Text resources do exist in many more languages!

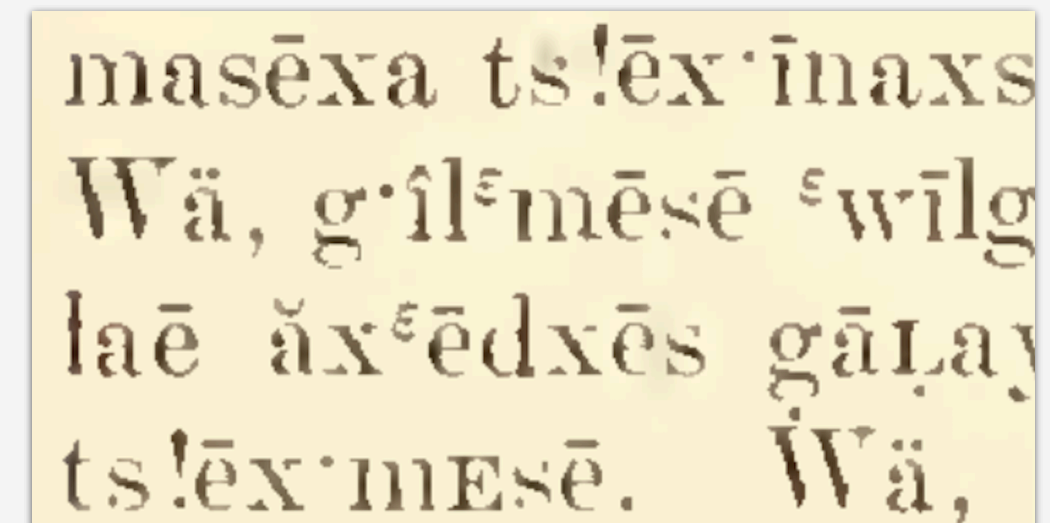
But locked away in formats that are not machine-readable



Printed books

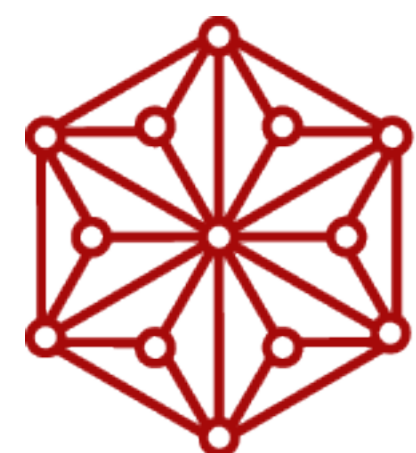


Handwritten notes



Typewritten documents

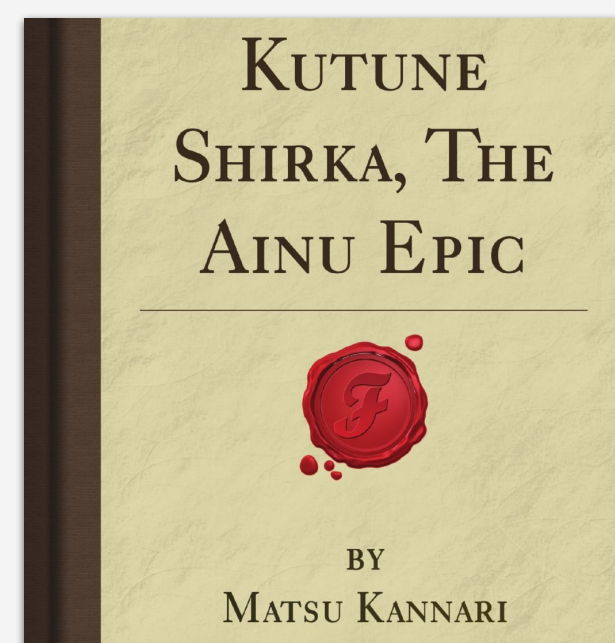
or in other formats such as bilingual lexicons



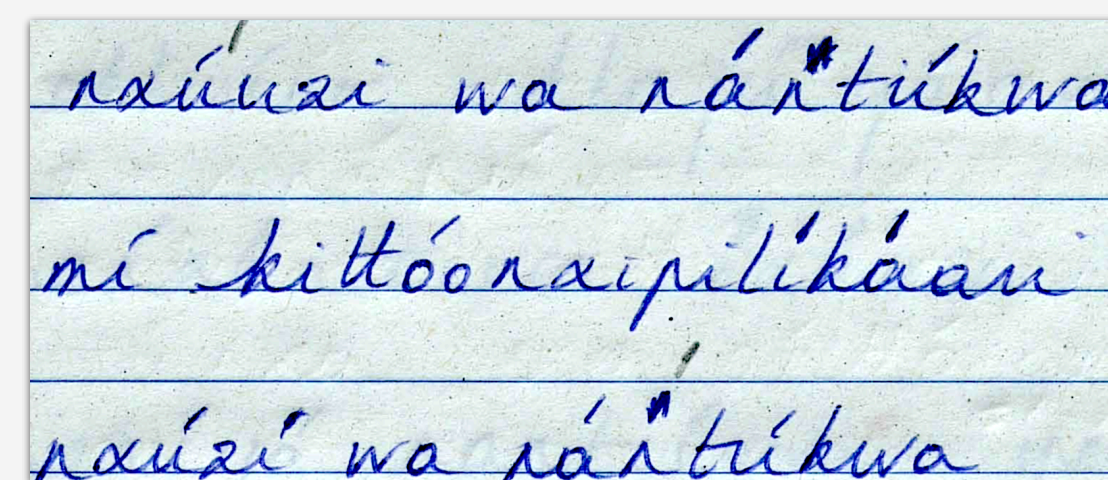
PANLEX

Text resources do exist in many more languages!

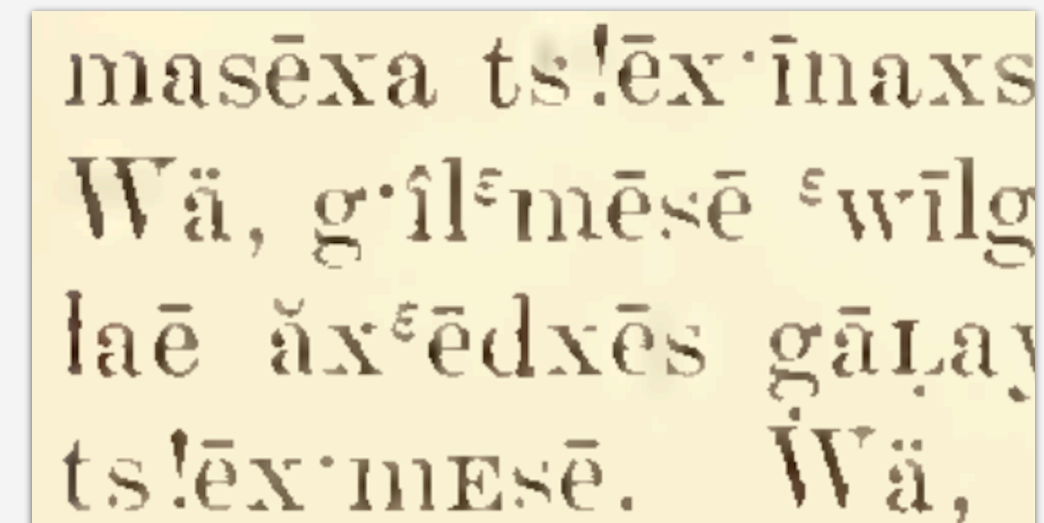
But locked away in formats that are not machine-readable



Printed books

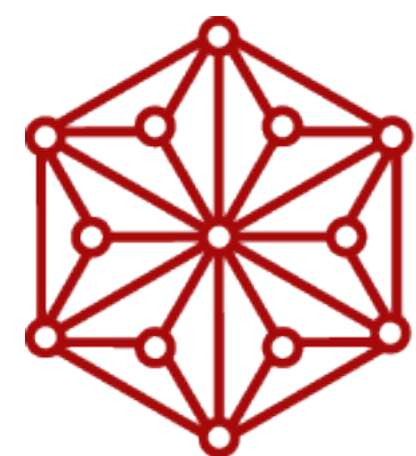


Handwritten notes



Typewritten documents

or in other formats such as bilingual lexicons

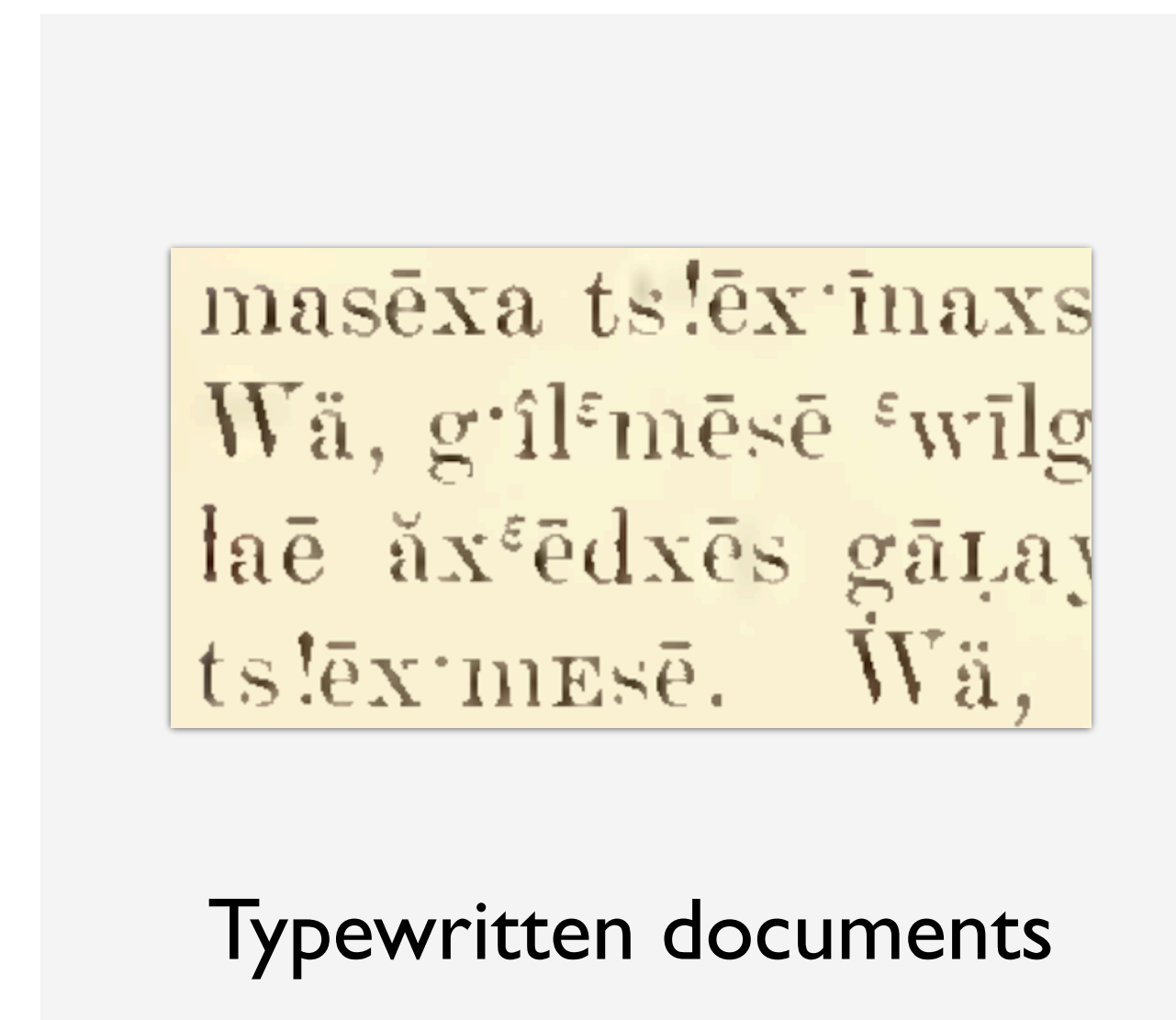
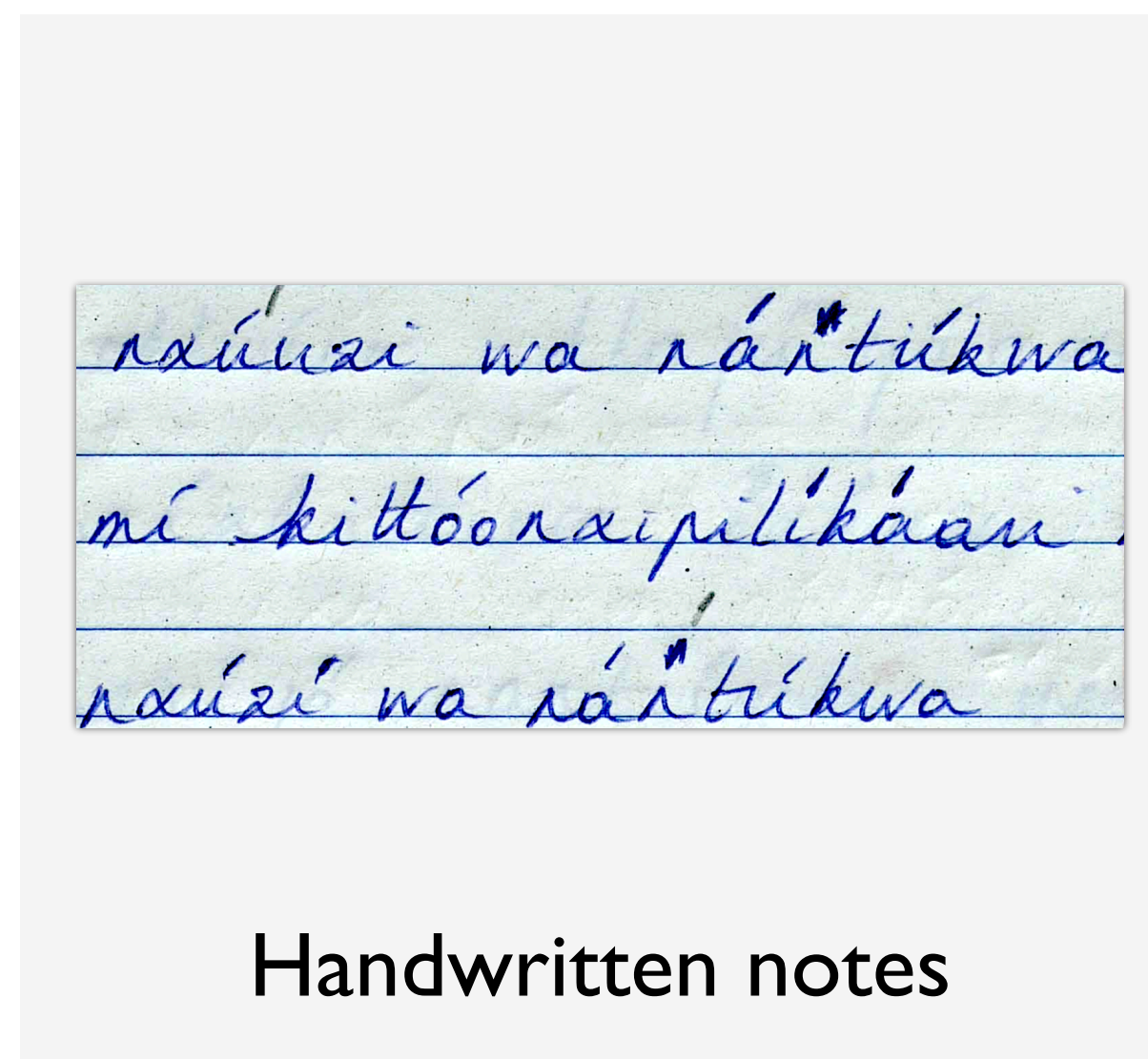
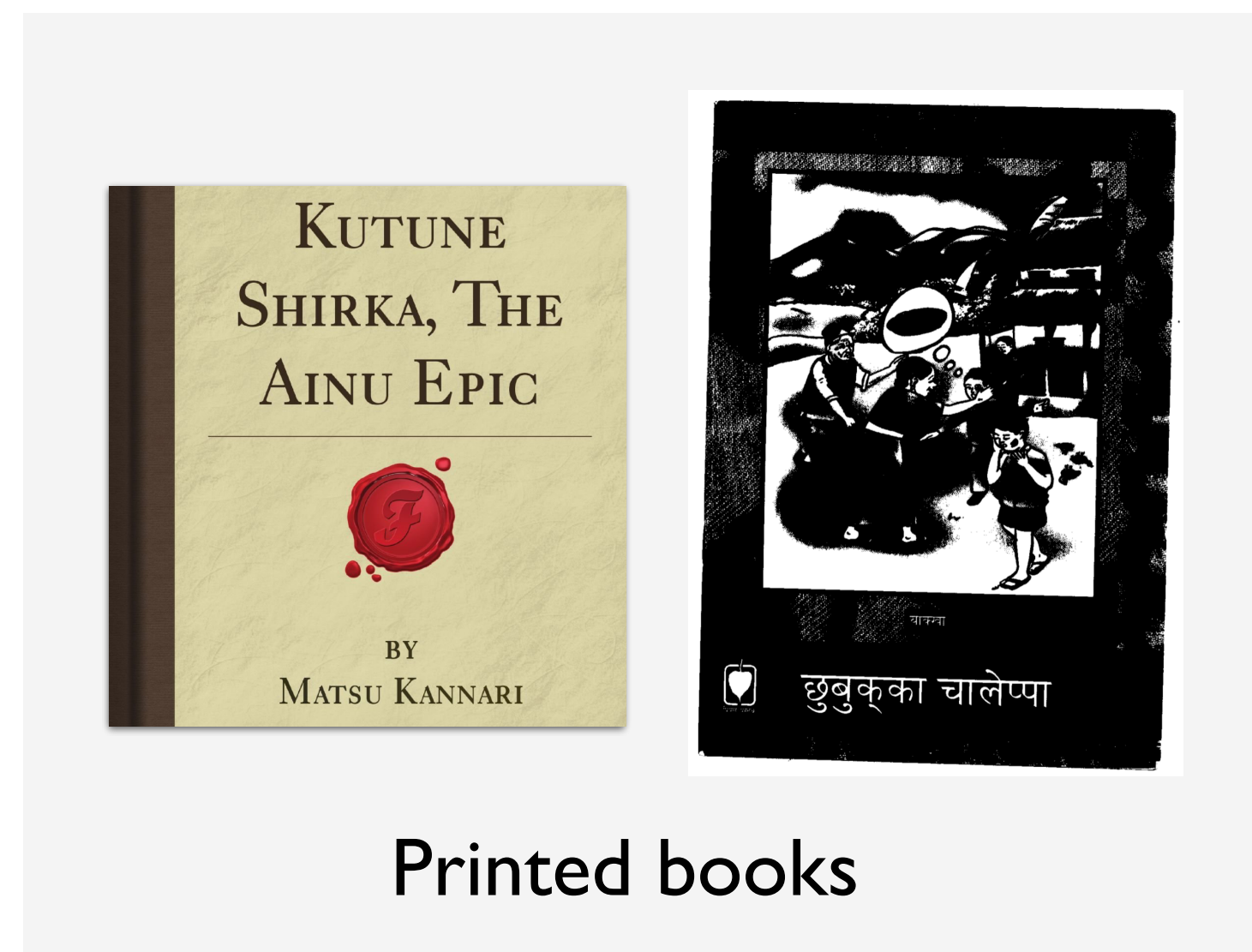


PANLEX

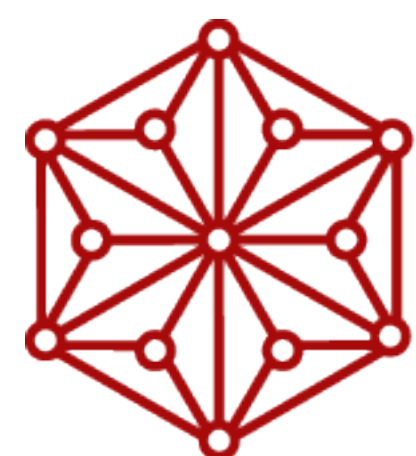


Text resources do exist in many more languages!

But locked away in formats that are not machine-readable



or in other formats such as bilingual lexicons



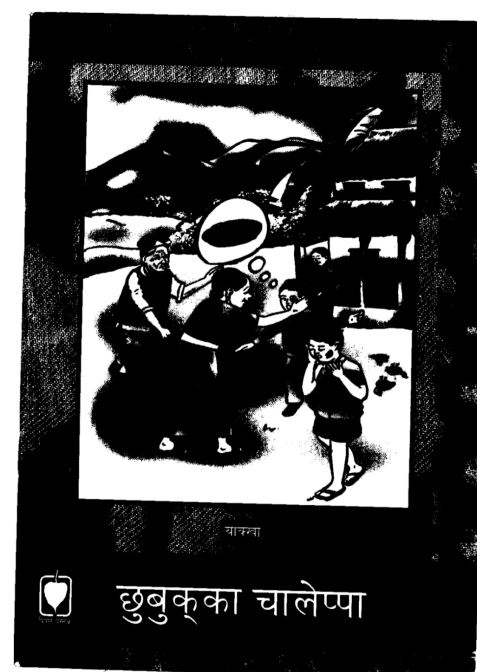
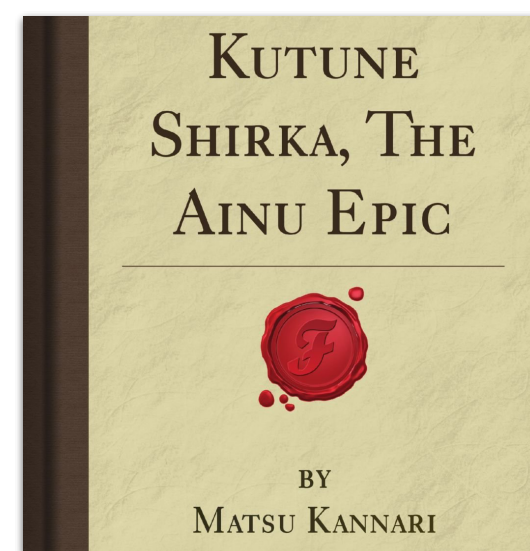
PANLEX



what can we do?!

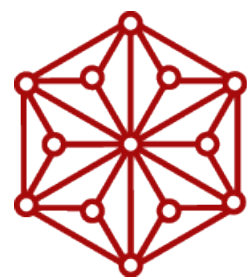
Text resources do exist in many more languages!

Text resources do exist in many more languages!



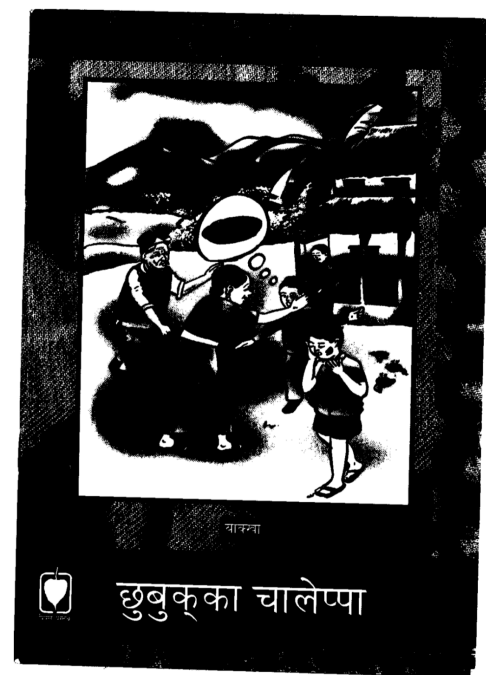
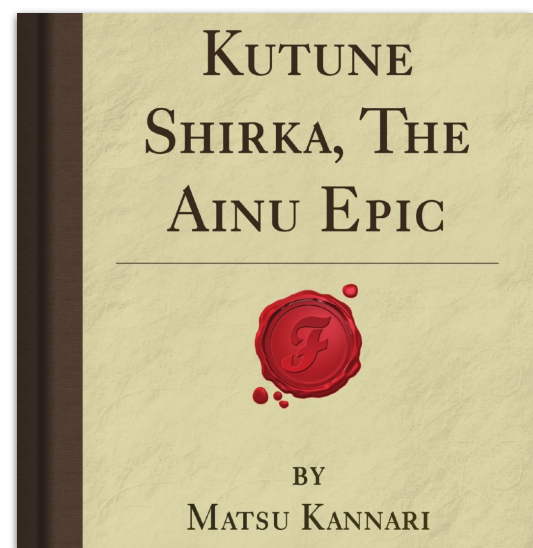
naúzi wa nántúkwa
mí kittóonáipilikáan
naúzi wa nántúkwa

masēxa ts!ēx·inaxs
 Wä, g·íl^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mēsē. Wä,



PANLEX

Text resources do exist in many more languages!



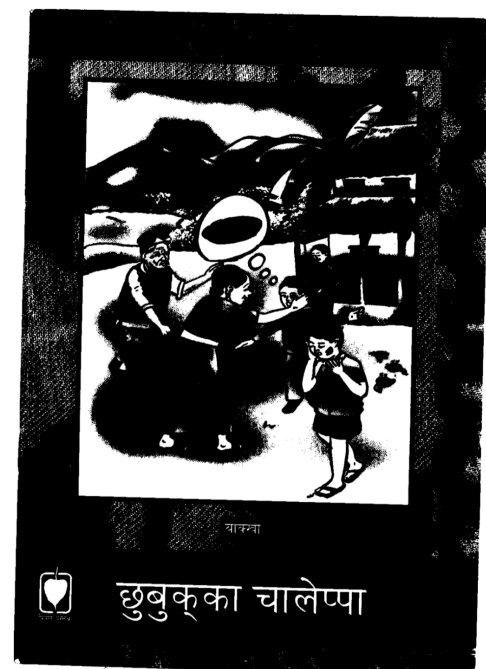
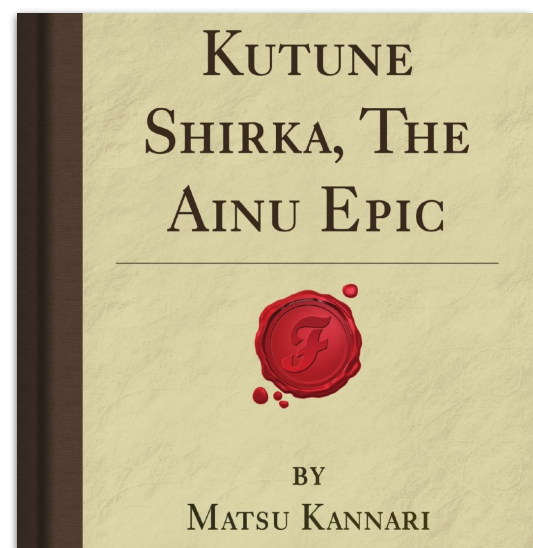
*naúzi wa nántúkwa
mí kittóonaipilikáani
naúzi wa nántúkwa*

masēxa ts!ēx·inaxs
Wä, g·íl^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,

Unlocking non-traditional resources



Text resources do exist in many more languages!



*naúzi wa nántúkwa
mí kittóonaipilikáani
naúzi wa nántúkwa*

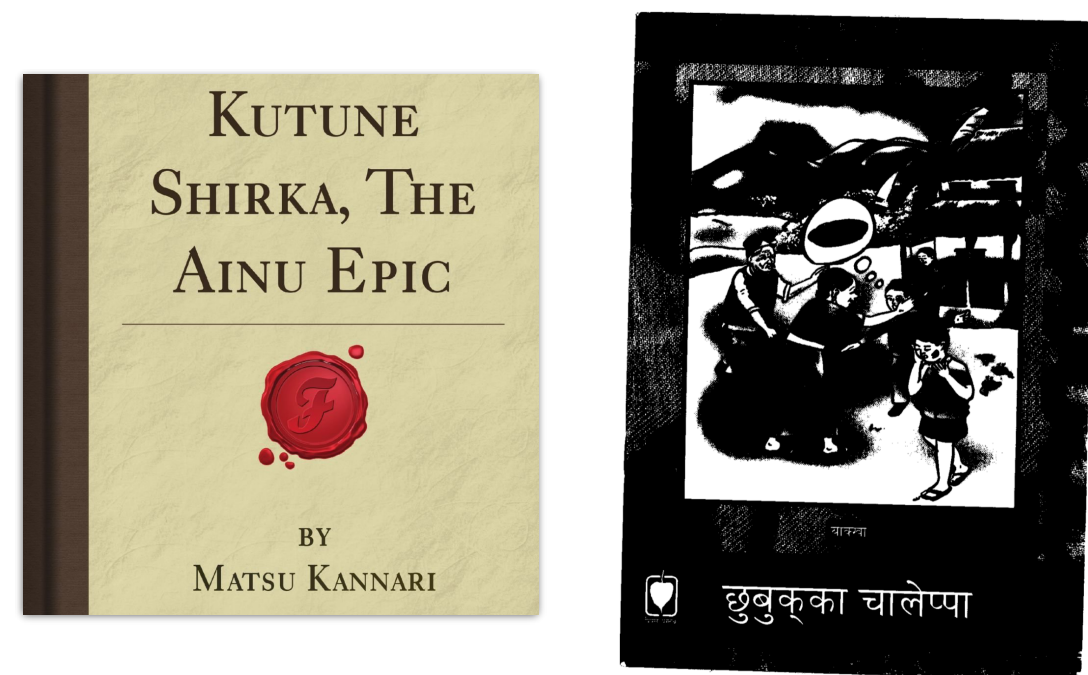
masēxa ts!ēx·inaxs
Wä, g·íl^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mesē. Wä,

Unlocking non-traditional resources

Enable NLP for under-resourced languages



Text resources do exist in many more languages!



*naúzi wa nántúkwa
 mí kittóonaipilikáani
 naúzi wa nántúkwa*

masēxa ts!ēx·inaxs
 Wä, g·íl^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mēsē. Wä,



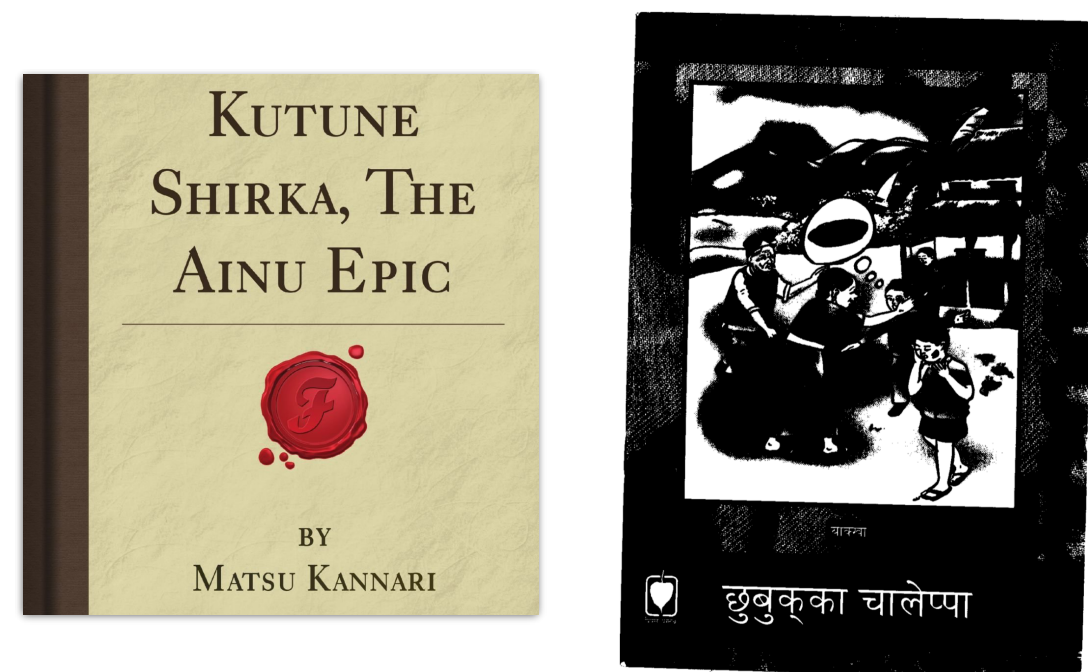
Unlocking non-traditional resources

Enable NLP for under-resourced languages

Expand multilingual LMs to more languages

XLM-R mBERT
 mT5 mBART ERNIE-M
 Turing ULR

Text resources do exist in many more languages!



*naúzi wa nántúkwa
 mí kittóonaipilikáani
 naúzi wa nántúkwa*

masēxa ts!ēx·inaxs
 Wä, g·il^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mesē. Wä,



Unlocking non-traditional resources

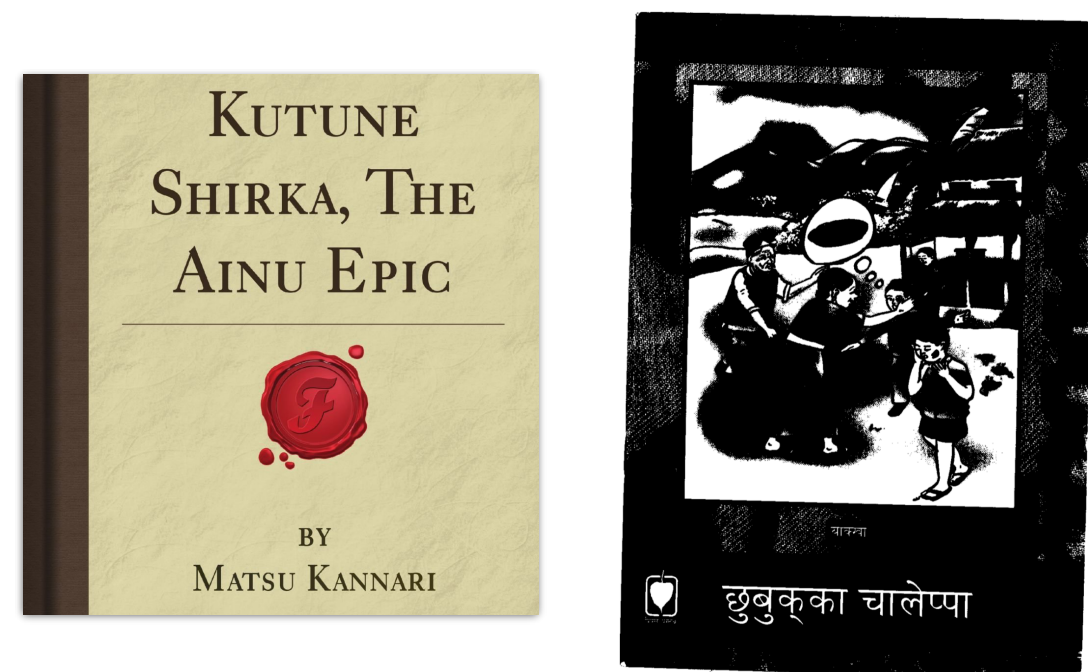
Enable NLP for under-resourced languages

Expand multilingual LMs to more languages

XLM-R mBERT
 mT5 mBART ERNIE-M
 Turing ULR

Annotate datasets for downstream NLP tasks

Text resources do exist in many more languages!



*naúzi wa nántúkwa
 mí kittóonaipilikáani
 naúzi wa nántúkwa*

masēxa ts!ēx·inaxs
 Wä, g·il^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mesē. Wä,



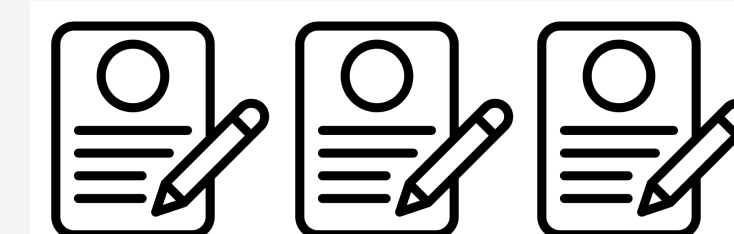
Unlocking non-traditional resources

Enable NLP for under-resourced languages

Expand multilingual LMs to more languages

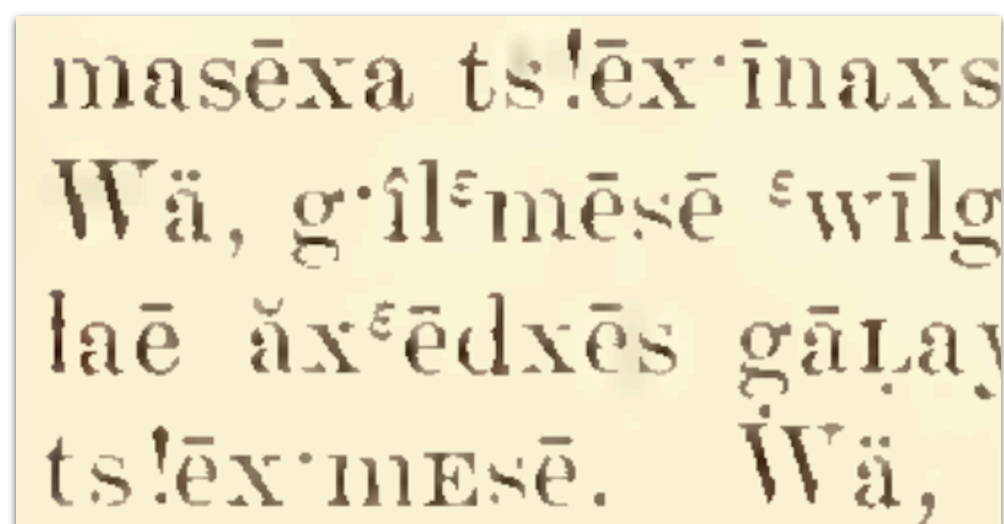
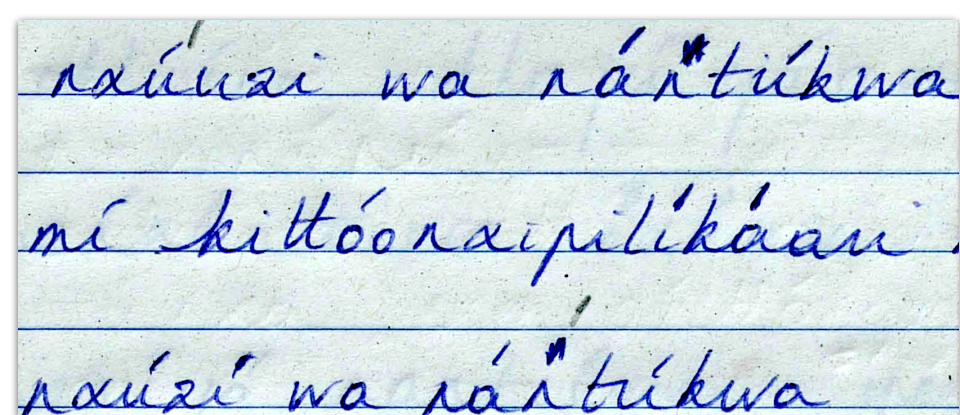
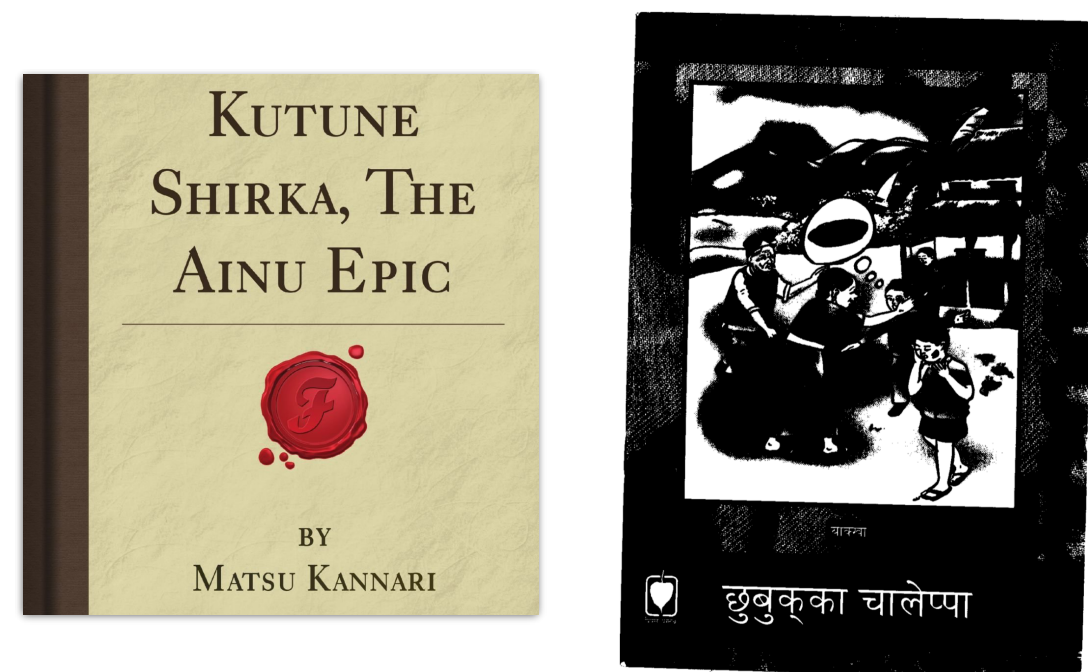
XLM-R mBERT
 mT5 mBART ERNIE-M
 Turing ULR

Annotate datasets for downstream NLP tasks



Support communities that speak these languages

Text resources do exist in many more languages!



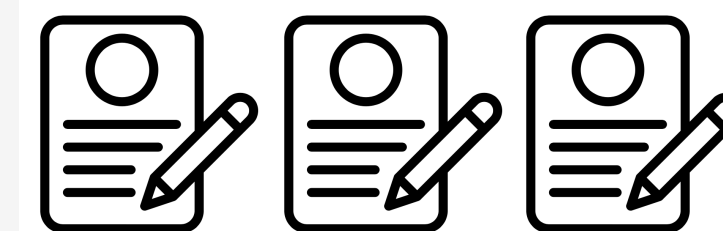
Unlocking non-traditional resources

Enable NLP for under-resourced languages

Expand multilingual LMs to more languages

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Annotate datasets for downstream NLP tasks

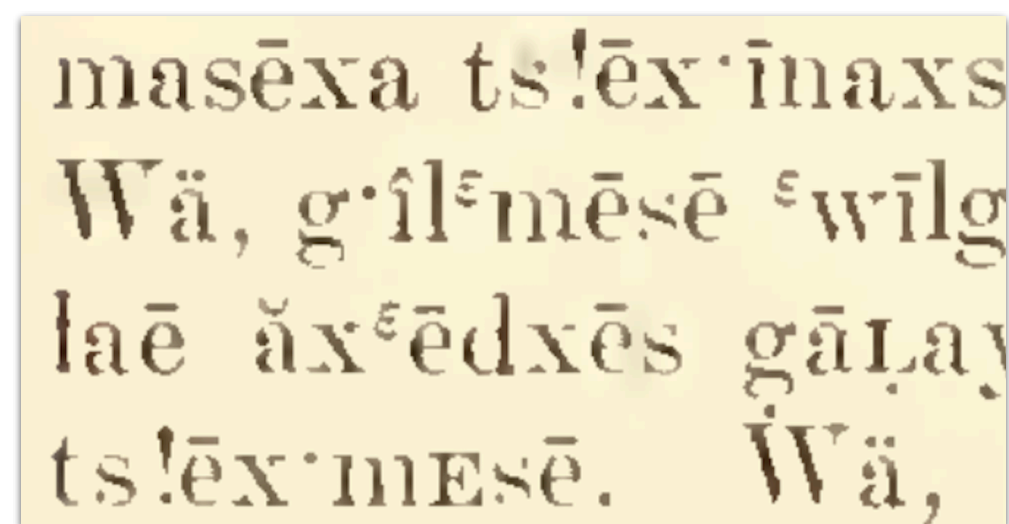
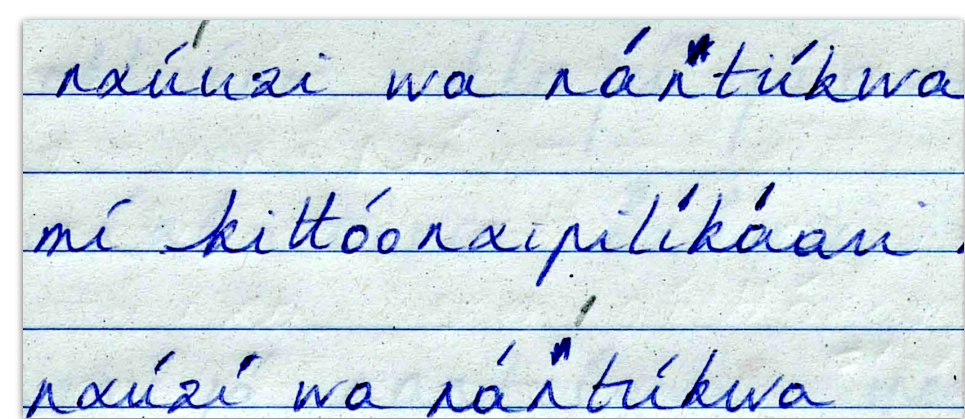
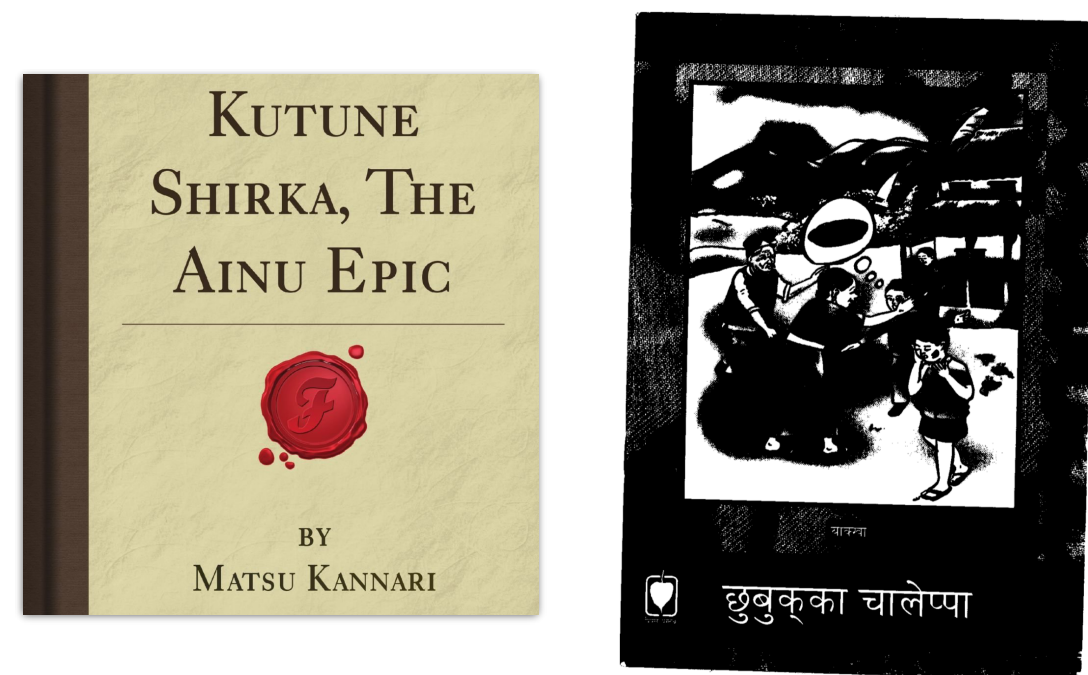


Support communities that speak these languages

Make native texts digitally accessible and searchable



Text resources do exist in many more languages!



Unlocking non-traditional resources

Enable NLP for under-resourced languages

Expand multilingual LMs to more languages

XLM-R mBERT
mT5 mBART ERNIE-M
Turing ULR

Annotate datasets for downstream NLP tasks



Support communities that speak these languages

Make native texts digitally accessible and searchable



Aid language researchers, educators, libraries...

Unlocking Un-digitized Text



Shruti Rijhwani, Antonios Anastasopoulos, Graham Neubig.
OCR Post-Correction for Endangered Language Texts.
EMNLP 2020.

Shruti Rijhwani, Daisy Rosenblum, Antonios Anastasopoulos, Graham Neubig.
Lexically-Aware Semi-Supervised Learning for OCR Post-Correction.
TACL 2021.

Extracting text from scanned documents

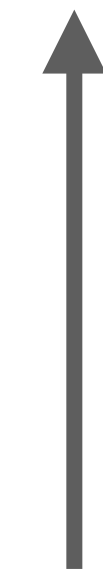
”Ma ti exi’ pu klei’?”
”Iklèo *ka* itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Scanned document

Extracting text from scanned documents

”Ma ti exi’ pu klei’?”
”Iklèo *ka* itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Scanned document



Scan from a book of folk tales in Griko

Extracting text from scanned documents

”Ma ti exi’ pu klei’?”
”Iklèo *ka* itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Scanned document

Extracting text from scanned documents

”Ma ti exi’ pu klei’?”
”Iklèo *ka* itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



Extracting text from scanned documents

"Ma ti exi' pu klei'?"
"Iklèo *ka* itela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti exi' pu klei'?"
"Iklèo ka itela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

Extracting text from scanned documents

"Ma ti exi' pu klei'?"
"Iklèo *ka* itela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti exi' pu klei'?"
"Iklèo ka itela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèò *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèò ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- High accuracy on languages that have easily available resources!

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèo *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- High accuracy on languages that have easily available resources!
- Off-the-shelf tools support many scripts and languages

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Machine readable text

- High accuracy on languages that have easily available resources!
- Off-the-shelf tools support many scripts and languages

Google Vision
Tesseract
EasyOCR
...

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- High accuracy on languages that have easily available resources!
- Off-the-shelf tools support many scripts and languages

Support 80-100 languages

Google Vision
Tesseract
EasyOCR
...

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèò *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèò ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- High accuracy on languages that have easily available resources!
- Off-the-shelf tools support many scripts and languages
- Little to no prior work on very low-resourced settings

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèò *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèò ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- Little to no prior work on very low-resourced settings

Extracting text from scanned documents

"Ma ti eḡi' pu klei'?"
"Iklèò *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
Recognition (OCR)



"Ma ti eḡi' pu klei'?"
"Iklèò ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Machine readable text

- Little to no prior work on very low-resourced settings

Evaluation dataset

Promises and pitfalls of existing methods

Neural models for improving OCR performance in low-resource settings

Rijhwani, Anastasopoulos, Neubig. EMNLP 2020.

Extracting text from scanned documents

”Ma ti eḡi’ pu klei’?”
 ”Iklèò *ka* ìtela n’armastò.”
 I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
 Recognition (OCR)



”Ma ti eḡi’ pu klei’?”
 ”Iklèò ka ìtela n’armastò.”
 I vèkkia àggale tria dattilìtia:

Machine readable text

- Little to no prior work on very low-resourced settings

Evaluation dataset

Promises and pitfalls of existing methods

Neural models for improving OCR performance in low-resource settings

Rijhwani, Anastasopoulos, Neubig. EMNLP 2020.

Extracting text from scanned documents

”Ma ti eḡi’ pu klei’?”
 ”Iklèò *ka* ìtela n’armastò.”
 I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
 Recognition (OCR)



”Ma ti eḡi’ pu klei’?”
 ”Iklèò ka ìtela n’armastò.”
 I vèkkia àggale tria dattilìtia:

Machine readable text

- Little to no prior work on very low-resourced settings

Evaluation dataset

Promises and pitfalls of existing methods

Neural models for improving OCR performance in low-resource settings

Rijhwani, Anastasopoulos, Neubig. EMNLP 2020.

Extracting text from scanned documents

”Ma ti eḡi’ pu klei’?”
 ”Iklèò *ka* ìtela n’armastò.”
 I *vèkkia* àggale tria dattilìtia:

Scanned document

Optical Character
 Recognition (OCR)



”Ma ti eḡi’ pu klei’?”
 ”Iklèò ka ìtela n’armastò.”
 I vèkkia àggale tria dattilìtia:

Machine readable text

- Little to no prior work on very low-resourced settings

Evaluation dataset

Promises and pitfalls of existing methods

Neural models for improving OCR performance in low-resource settings

Semi-supervised learning to improve performance with unlabeled images

Rijhwani, Anastasopoulos, Neubig. EMNLP 2020.

Rijhwani, Rosenblum, Anastasopoulos, Neubig. TACL 2021.

Evaluation dataset for low-resource OCR

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèo *ka* itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèo ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèò ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोकडागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèo ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexayē
lexelāsa nekwāxa nek!ūlē. Wā, hēⁿ
wā, lā hēlēda ⁿemsgemē; wā, hē^{mi}
lexelās. Wā hēem lēgemsa ^wālēga^s

- Orthographically, typologically, geographically diverse

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Latin

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèò ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

- Orthographically, typologically, geographically diverse

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Latin

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèò ka ìtela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Latin+Greek

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

- Orthographically, typologically, geographically diverse

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Latin

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèò ka ìtela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Latin+Greek

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Devanagari

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

- Orthographically, typologically, geographically diverse

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Latin

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèo ka itela n’armastò.”
I vèkkia àggale tria dattilìtia:

Latin+Greek

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोकडागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Devanagari

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

Boas

- Orthographically, typologically, geographically diverse

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèo ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोकडागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

- Orthographically, typologically, geographically diverse
- The languages currently have:
 - No Wikipedia/Common Crawl text
 - Not supported by multilingual LMs
 - No easily accessible bilingual lexica

Evaluation dataset for low-resource OCR

Ainu
(Japan)

kira-an patek
aeyairamshitne⁽¹⁾
5760 hushkotoi wano⁽²⁾
iki-an aine

Griko
(Italy)

”Ma ti exi’ pu klei’?”
”Iklèò ka itela n’armastò.”
I *vèkkia* àggale tria dattilìtia:

Yakkha
(Nepal)

मा, ना चिगा निड्वामाड् ओम,
हाखोक्डागो लेम्साड् खा?ला लुया,
“पिछानाछा लेड्माहोड प्याक छो छो
लाप्लाप मेन्जोकमाहा।”

Kwak’wala
(Canada)

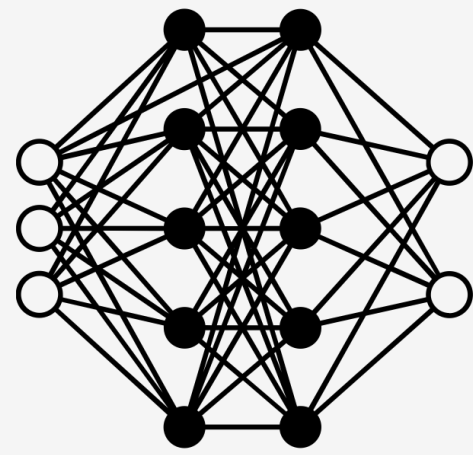
q!âLElax gwēg’ilasasa lexēlāxa lexā^εyē
lexelāsa nekwāxa nek!ūlē. Wā, hē^εn
wā, lā hēlēda ^εnemsgemē; wā, hē^εmi
lexelās. Wā hēem lēgemsa ^εwālēga^ε

- Orthographically, typologically, geographically diverse
- The languages currently have:
 - No Wikipedia/Common Crawl text
 - Not supported by multilingual LMs
 - No easily accessible bilingual lexica
- <1000 transcribed lines per language

Existing OCR methods

Existing OCR methods

Supervised

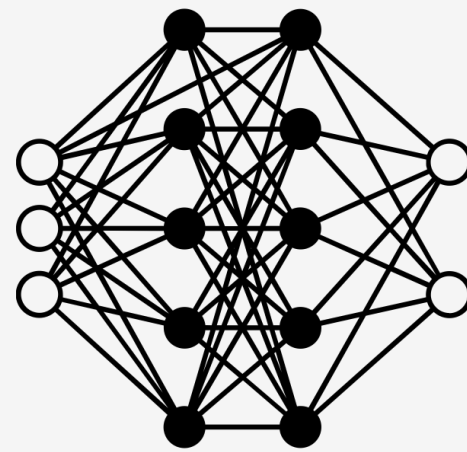


Large neural networks

Requires: **10000s of transcribed images**

Existing OCR methods

Supervised



Large neural networks

Requires: 10000s of
transcribed images

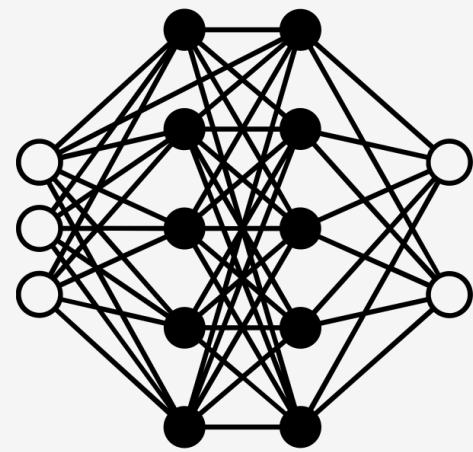
Unsupervised

- Unlabeled images
- Language model

Requires: text corpus
or lexicon in the
target language

Existing OCR methods

Supervised



Large neural networks

Requires: 10000s of transcribed images

Unsupervised

- Unlabeled images
- Language model

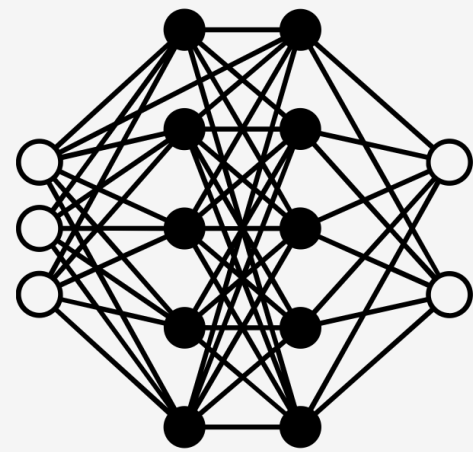
Requires: text corpus or lexicon in the target language

Off-the-shelf

- Support ~100 languages
- Not trained on our target languages
- Can act as a **general character recognizer** for many scripts

Existing OCR methods

Supervised



Large neural networks

Requires: 10000s of transcribed images

Unsupervised

- Unlabeled images
- Language model

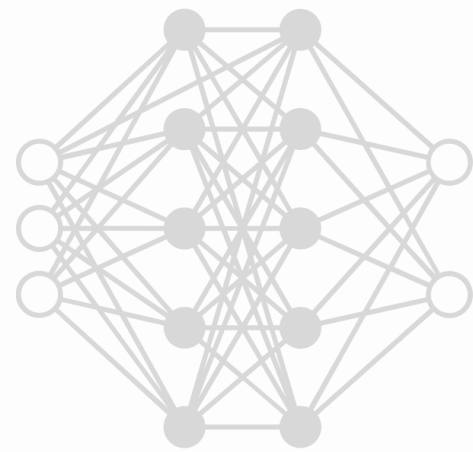
Requires: text corpus or lexicon in the target language

Off-the-shelf

- Support ~100 languages
- Not trained on our target languages
- Can act as a general character recognizer for many scripts

Existing OCR methods

Supervised



Large neural networks

Requires: 10000s of transcribed images

Unsupervised

- Unlabeled images
- Language model

Requires: **text corpus** or **lexicon** in the target language

Off-the-shelf

- Support ~100 languages
- Not trained on our target languages
- Can act as a **general character recognizer** for many scripts

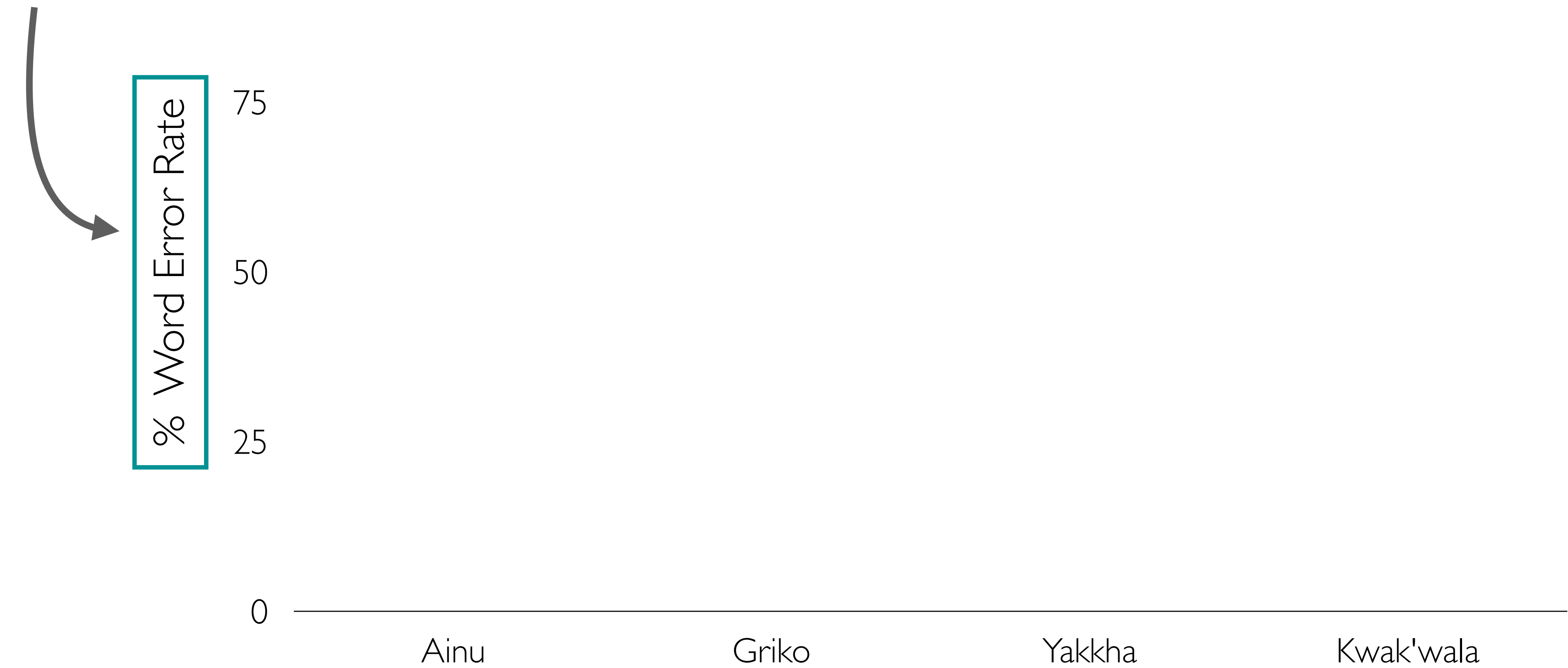
Existing OCR methods: promises and pitfalls

Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls

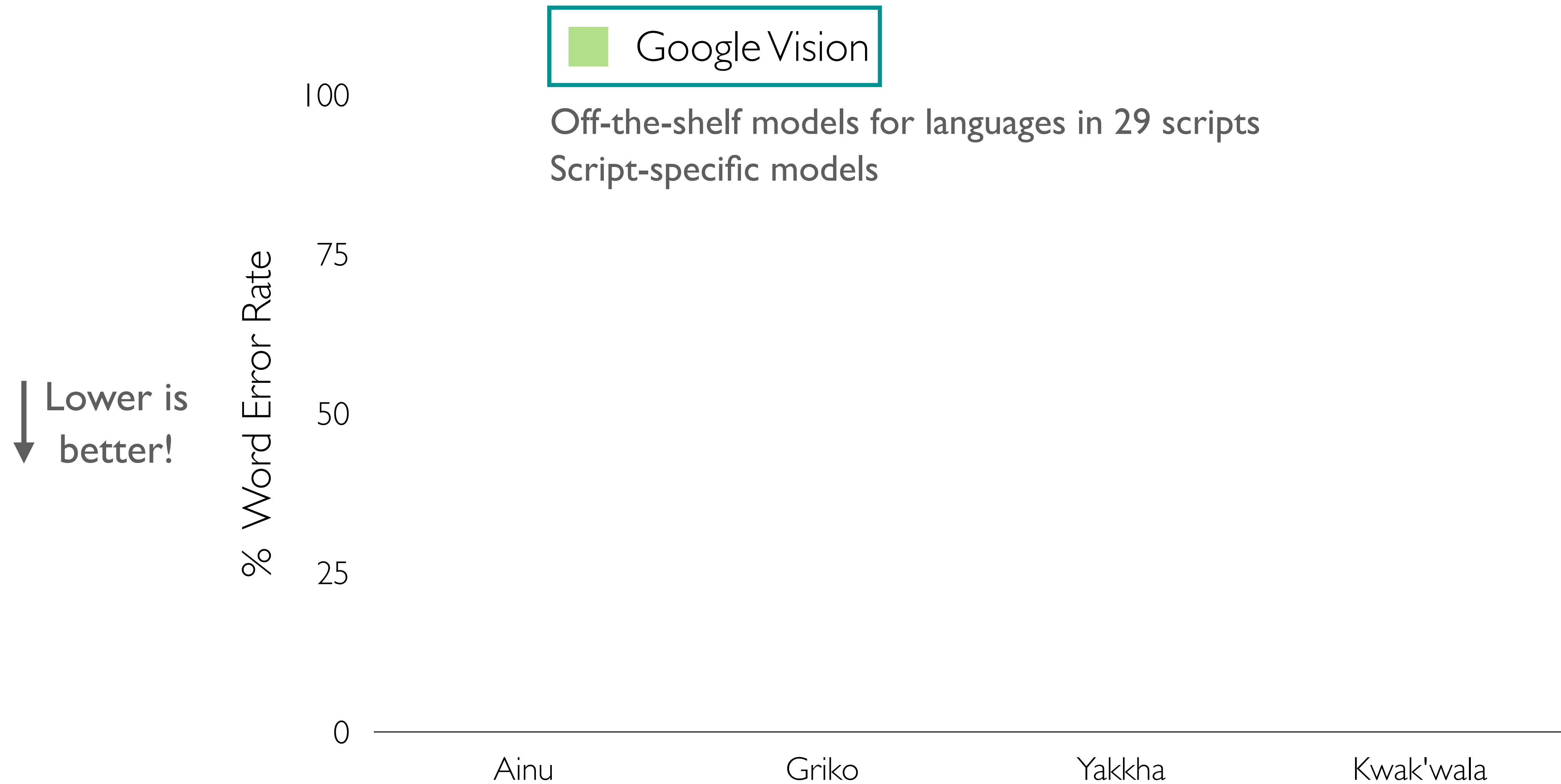
word edit distance between prediction and reference
number of words in reference



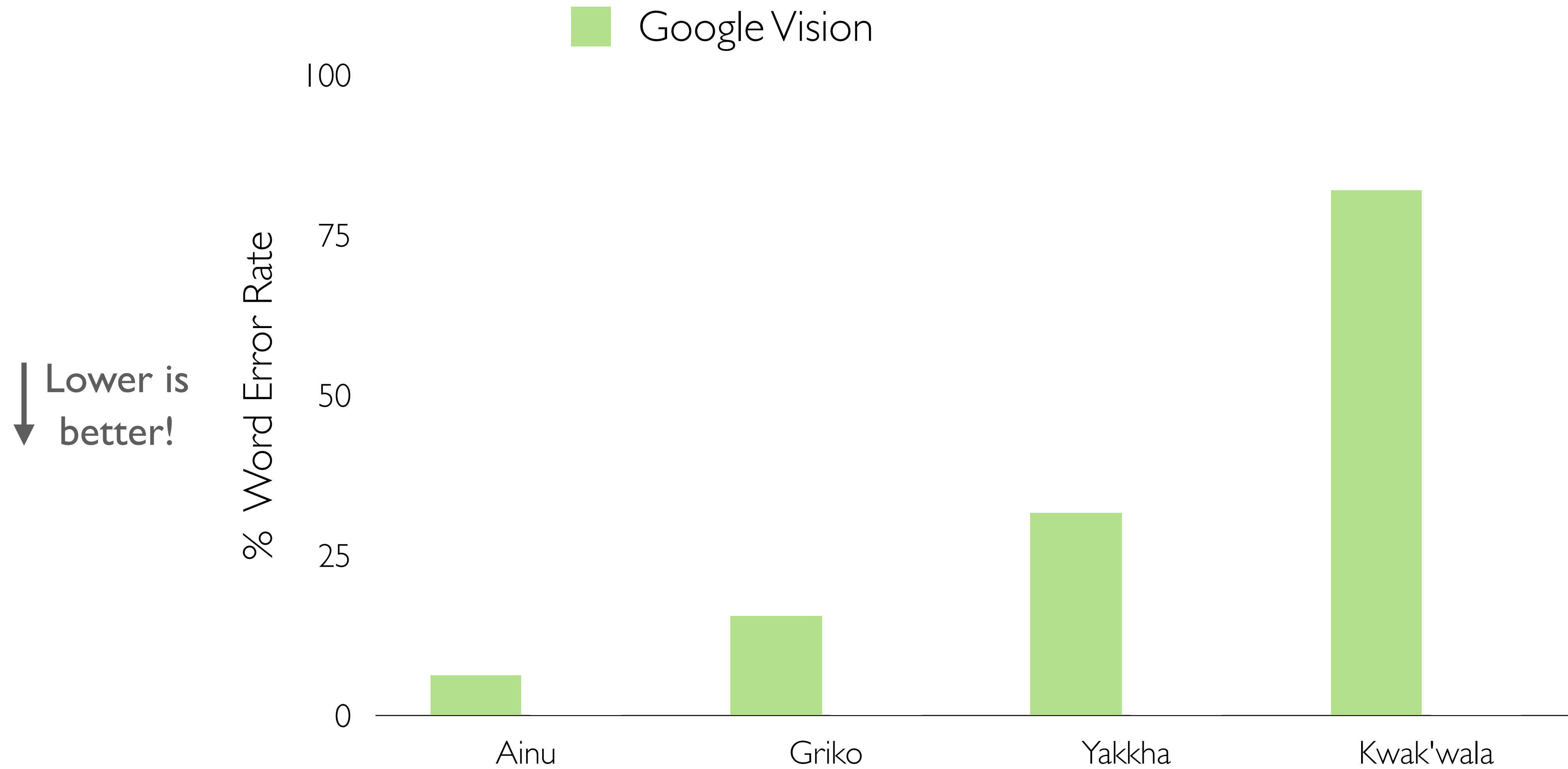
Existing OCR methods: promises and pitfalls



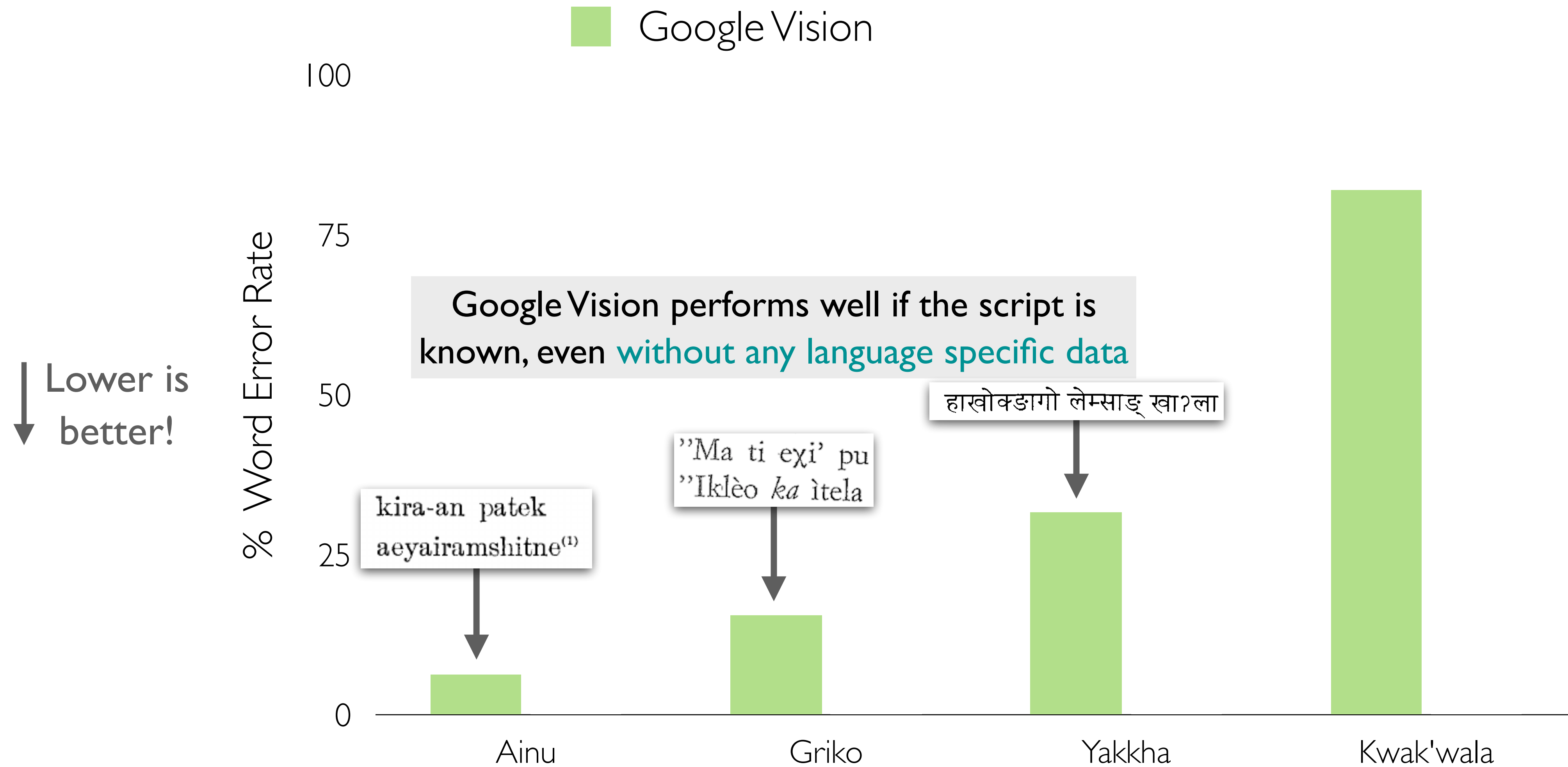
Existing OCR methods: promises and pitfalls



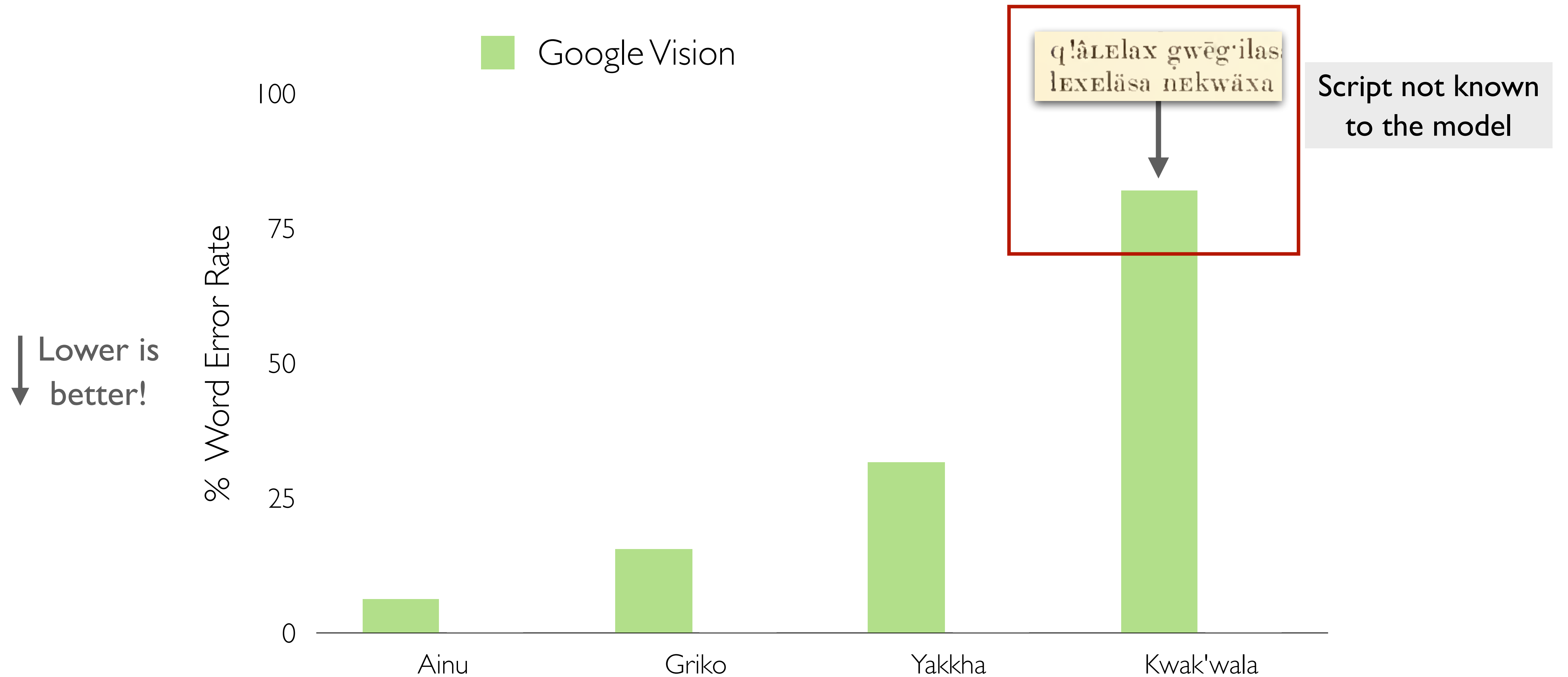
Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls

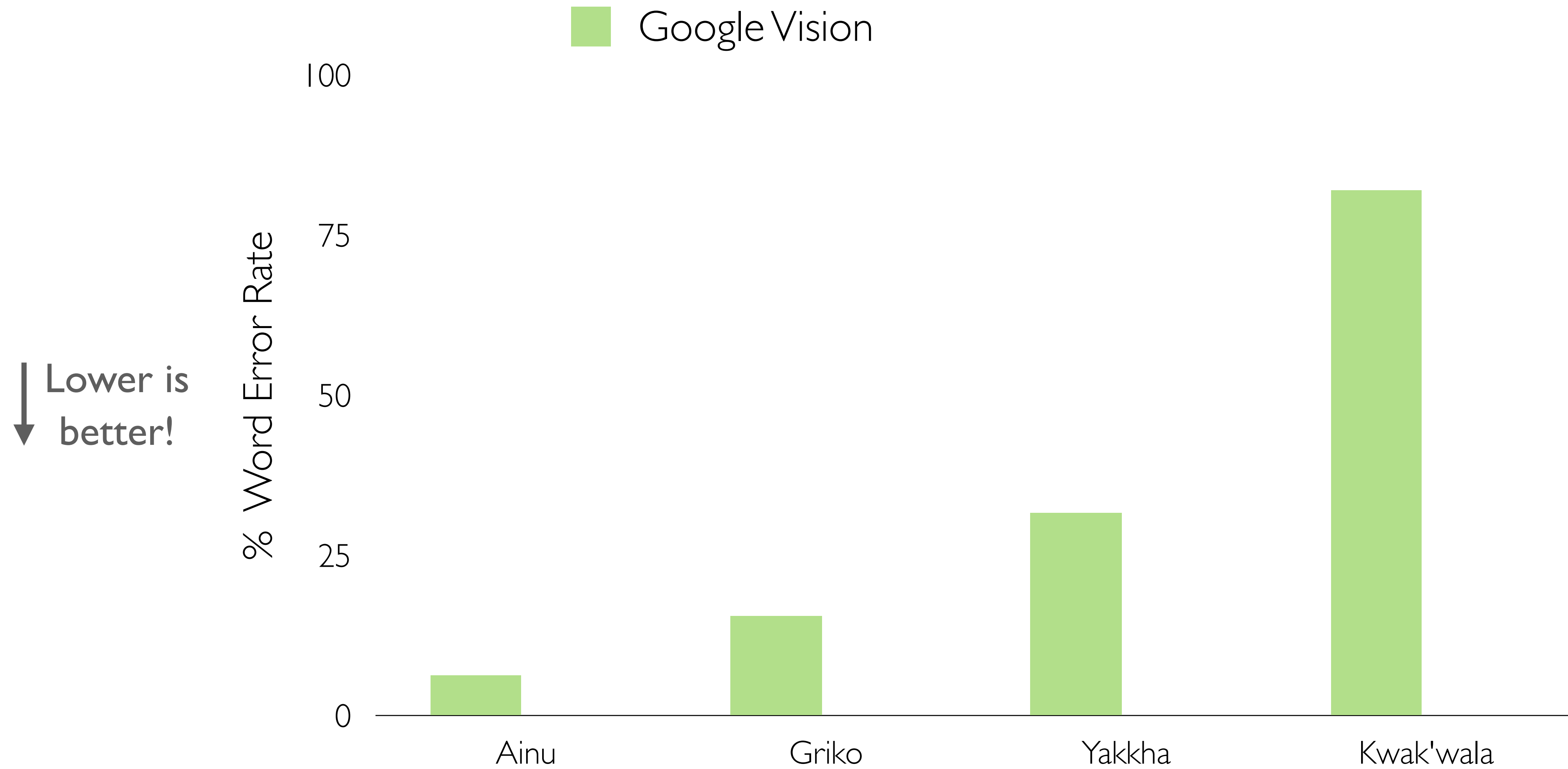


Existing OCR methods: promises and pitfalls

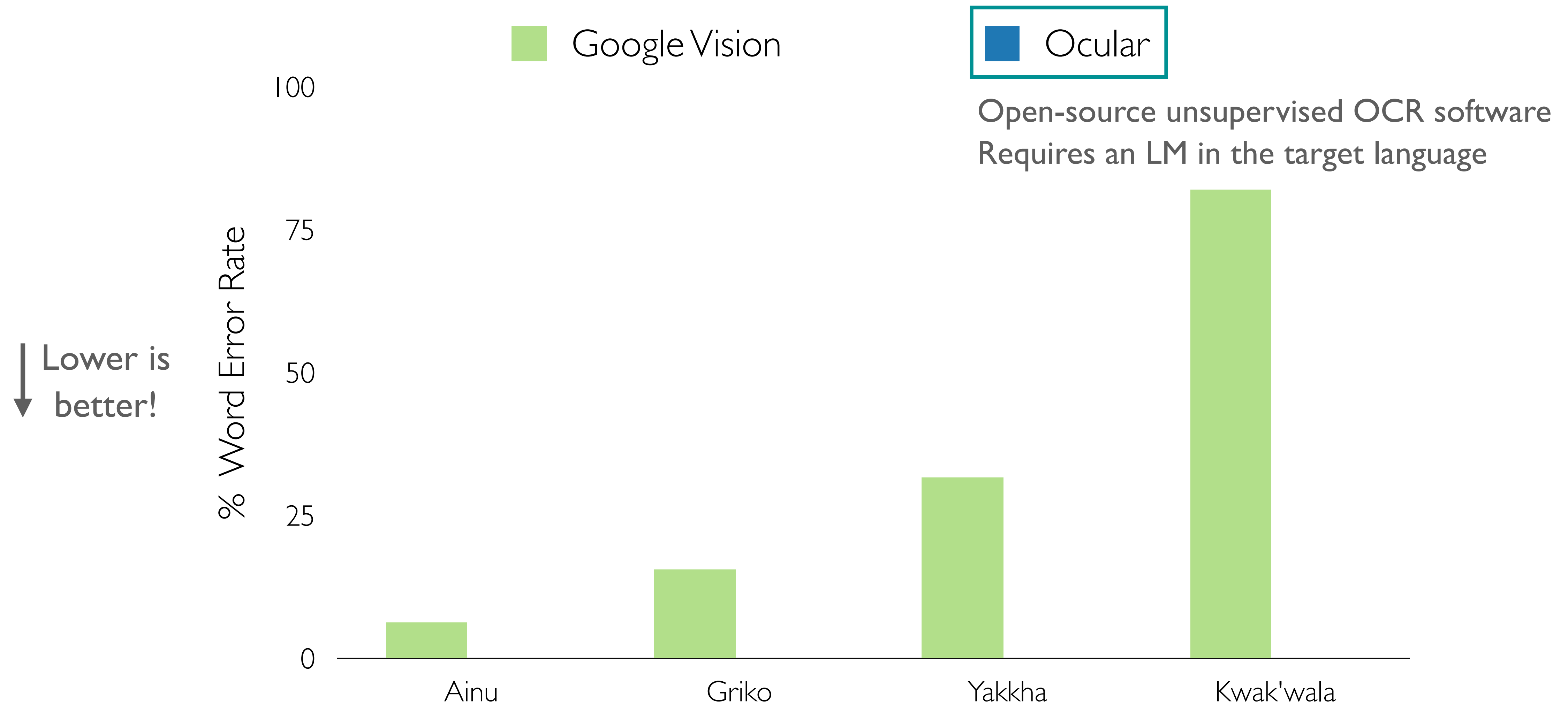


Existing OCR methods: promises and pitfalls

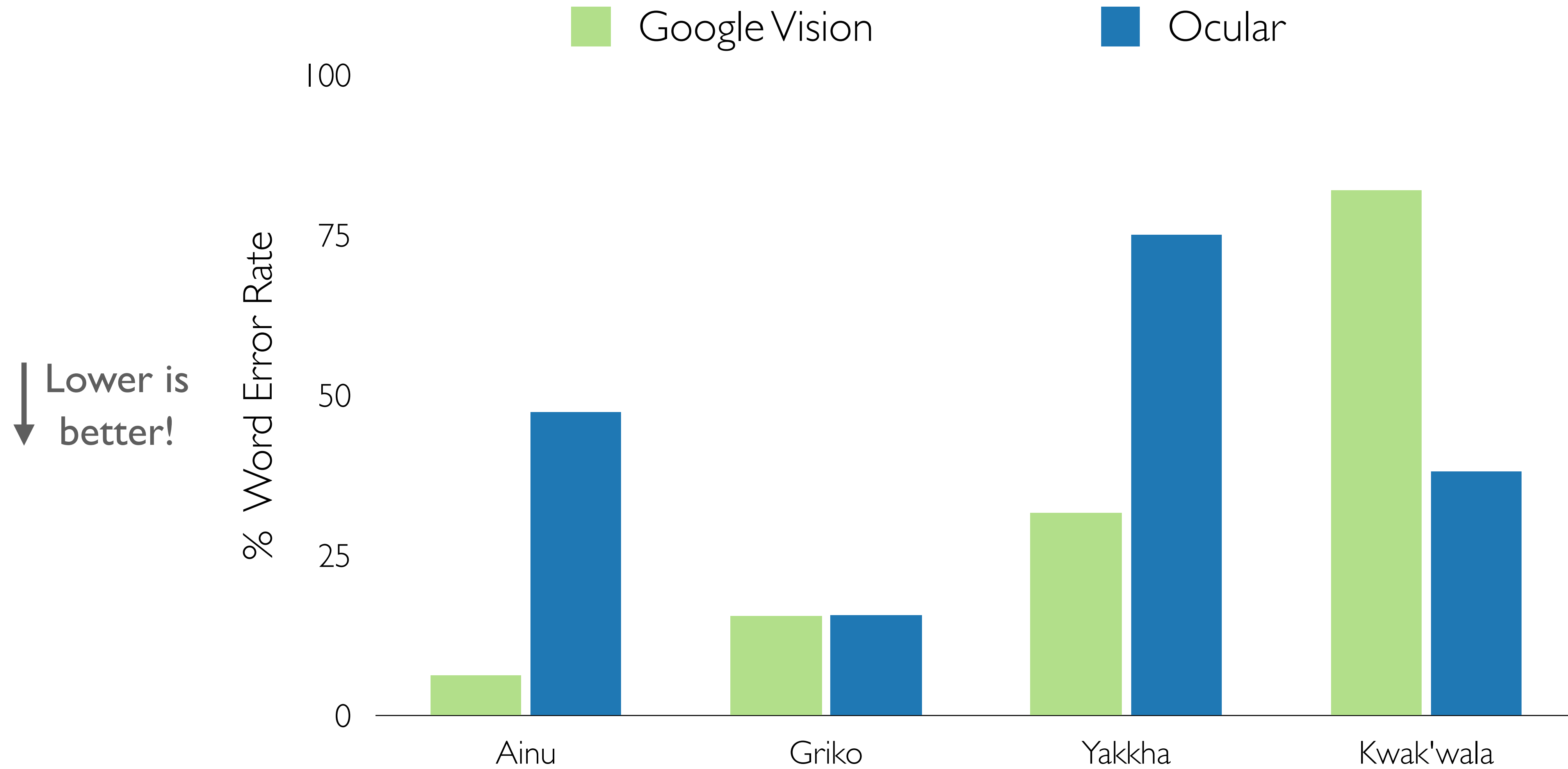
Existing OCR methods: promises and pitfalls



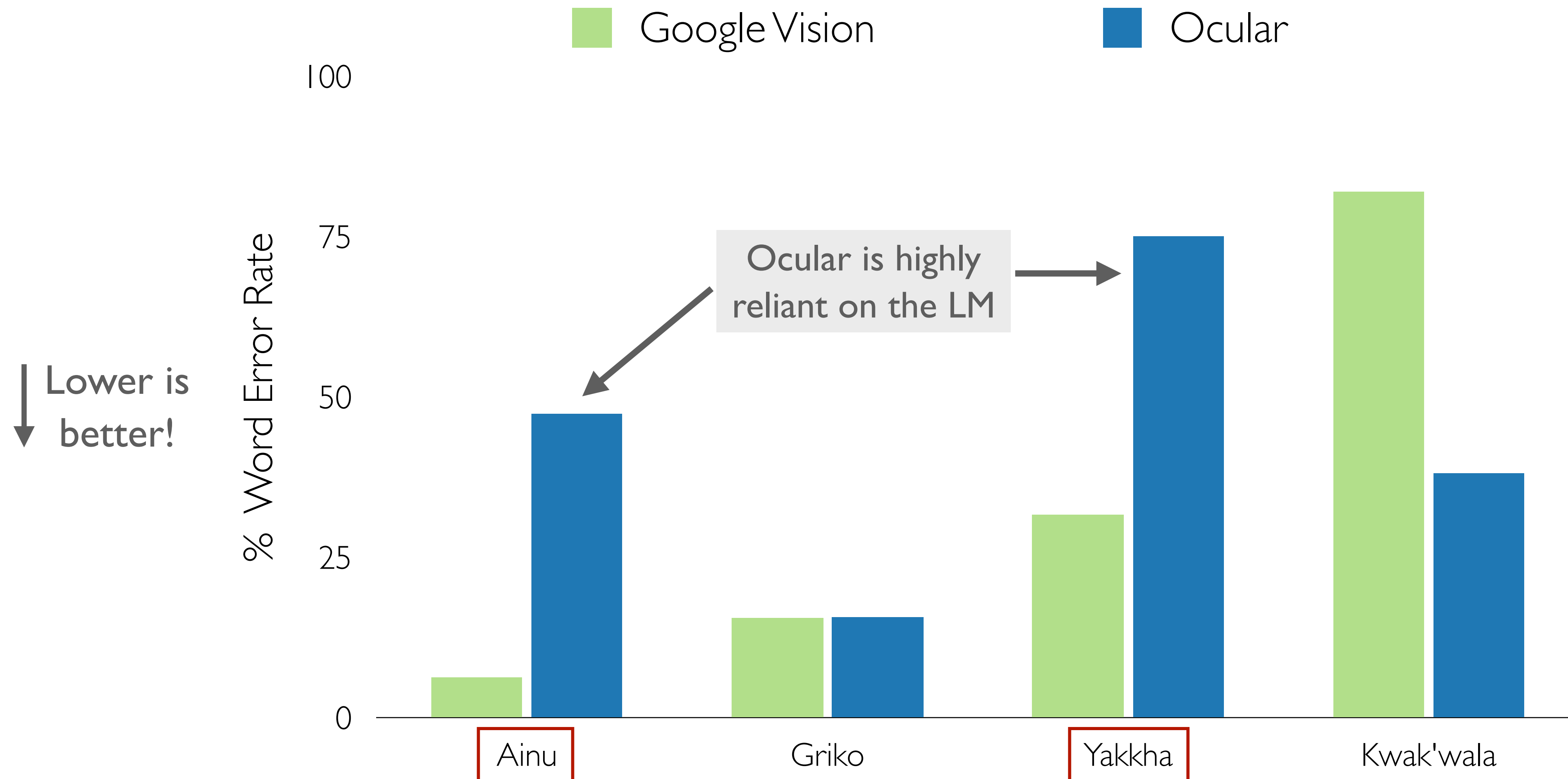
Existing OCR methods: promises and pitfalls



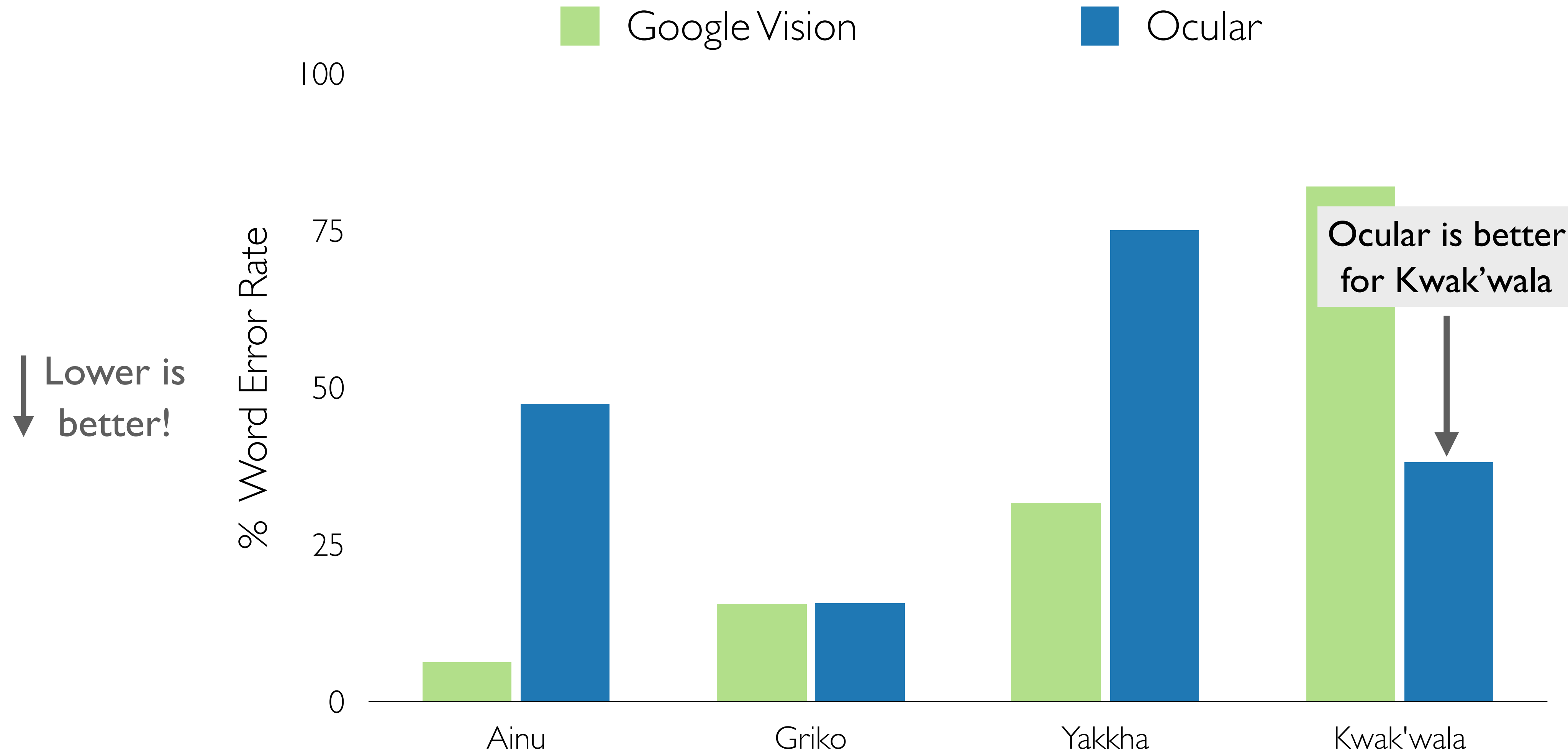
Existing OCR methods: promises and pitfalls



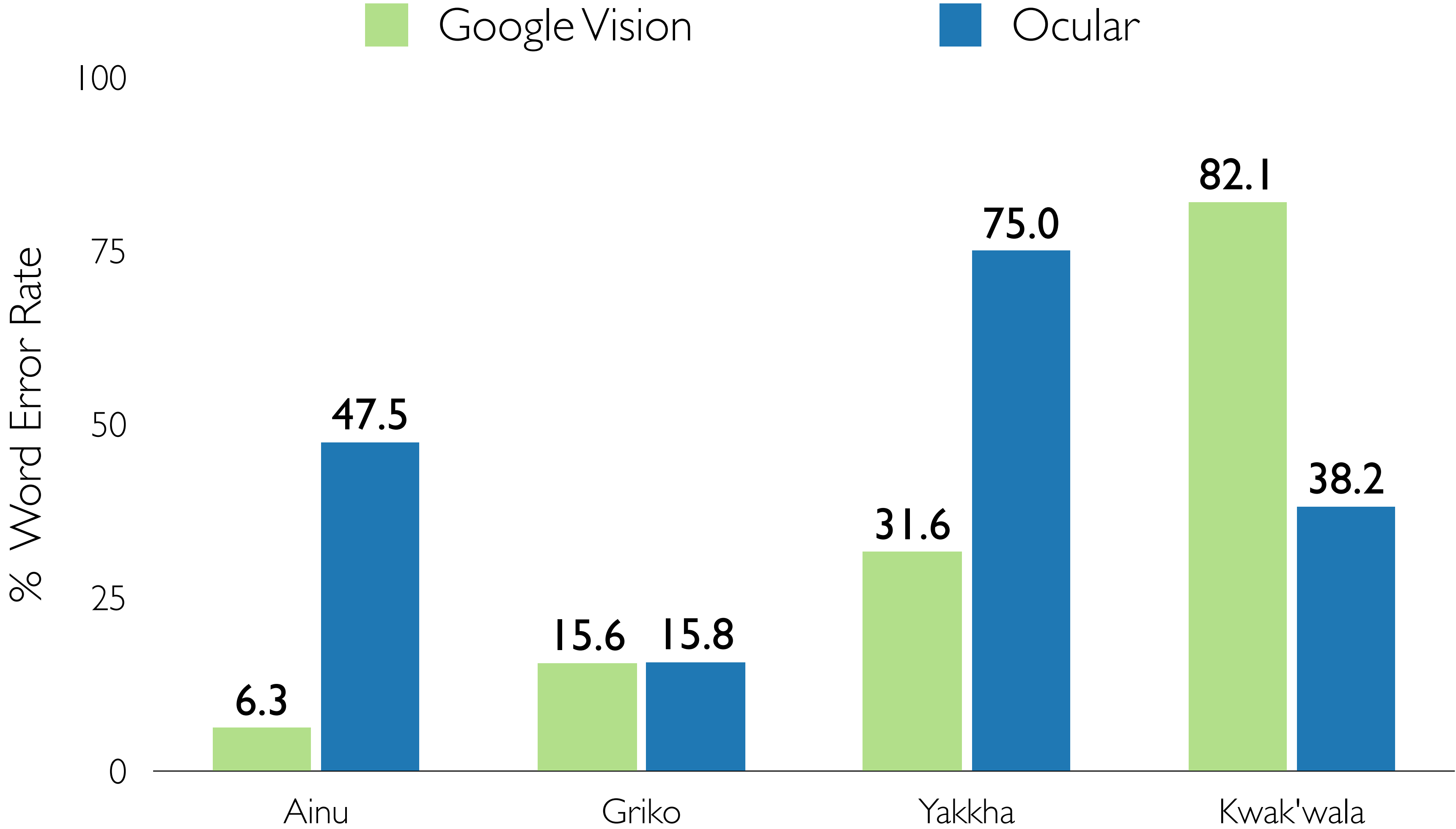
Existing OCR methods: promises and pitfalls



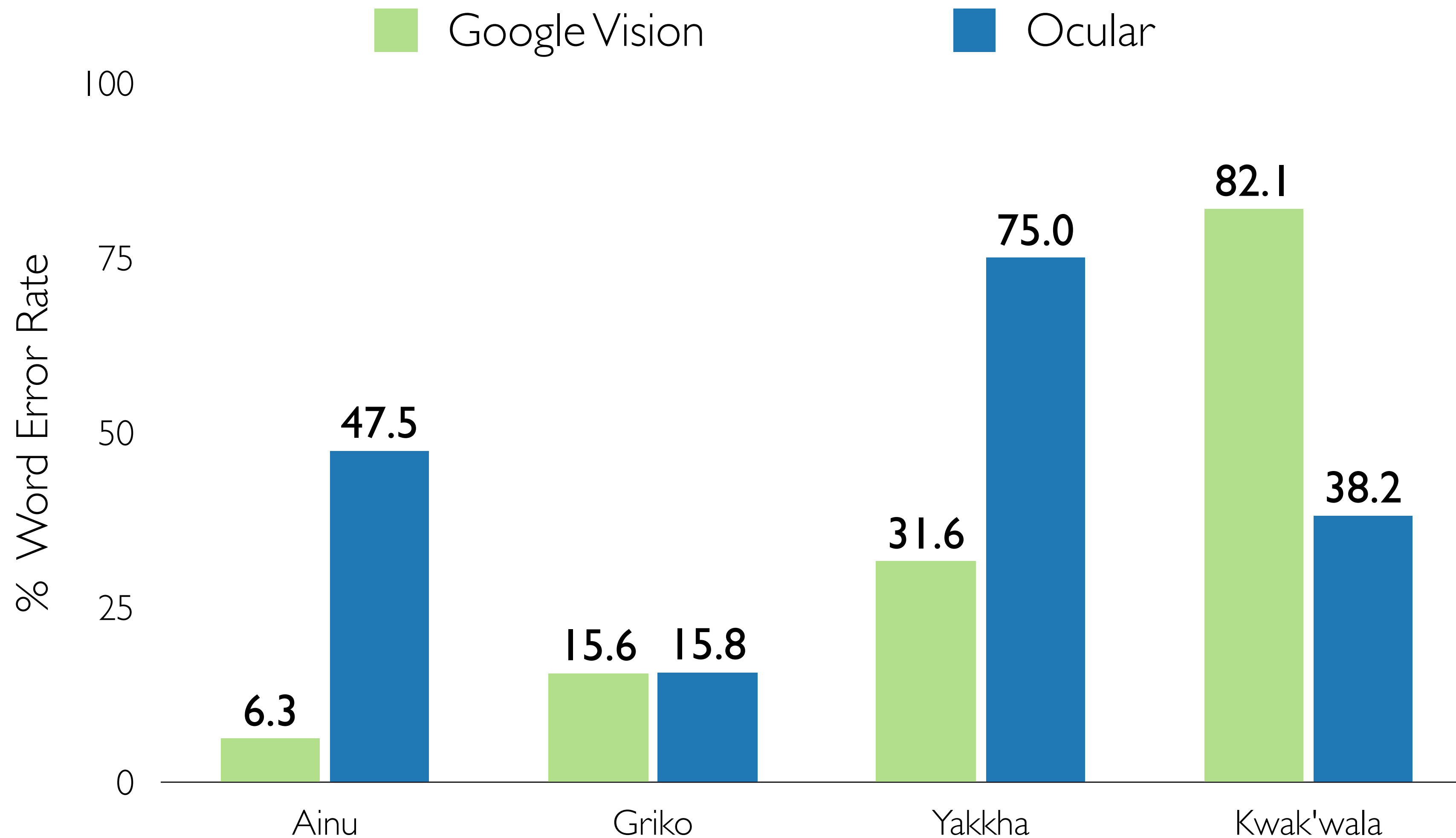
Existing OCR methods: promises and pitfalls



Existing OCR methods: promises and pitfalls

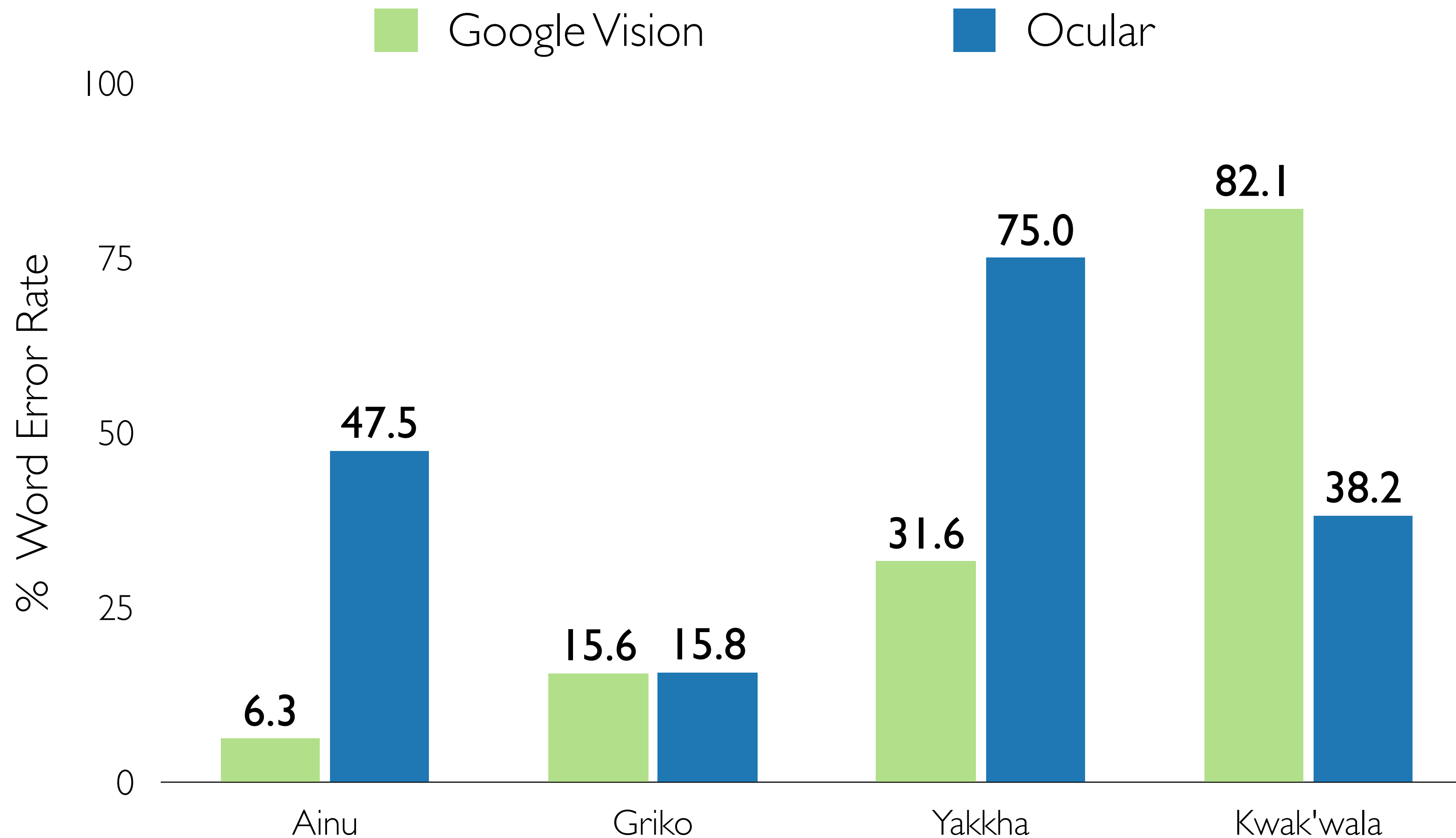


Existing OCR methods: promises and pitfalls



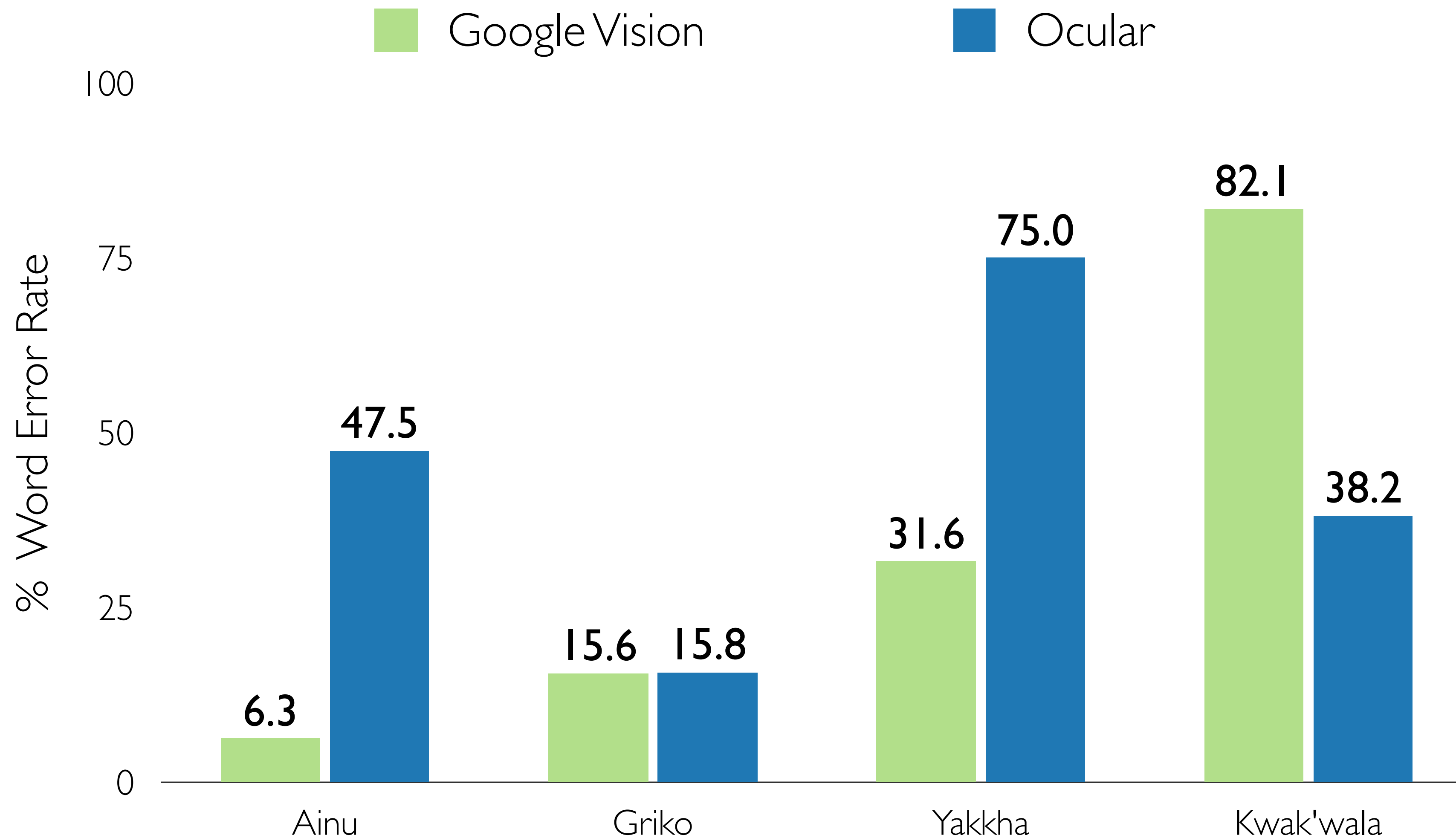
- **Considerable room for improvement** compared to high-resource languages

Existing OCR methods: promises and pitfalls



- **Considerable room for improvement** compared to high-resource languages
- Recognizes the majority of words correctly

Existing OCR methods: promises and pitfalls



- **Considerable room for improvement** compared to high-resource languages
- Recognizes the majority of words correctly
- **Reliable starting point** for further improvements

Improving the results of existing OCR systems

Improving the results of existing OCR systems

”Ma ti exi’ pu klei’?”

”Iklèo *ka* itela n’armastò.”

I *vèkkia* àggale tria dattilìtia:

Improving the results of existing OCR systems

”Ma ti exi’ pu klei’?”
”Iklèo ka ìtela n’armastò.”
I vèkkia àggale tria dattilìtia:

”Ma ti exi’ pu klei’?”
”Ikleo ka ìtela _armastò.”
I vekkia aggale tria dattilitia:

OCR output (“first pass”)

Improving the results of existing OCR systems

"Ma ti e*x*i' pu klei'?"
"Iklèo *ka* itela n'armastò."
I *vèkkia* àggale tria dattilìtia:

"Ma ti e*x*i' pu klei'?"
"Ikleo ka itela armastò."
I vekkia aaggale tria dattilitia:

OCR output ("first pass")



OCR output has
some errors

Improving the results of existing OCR systems

"Ma ti exi' pu klei'?"
"Iklèo ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

"Ma ti exi' pu klei'?"
"Ikleo ka ìtela _armastò."
I vekkia aggale tria dattilitia:

OCR output ("first pass")

Automatic OCR
Post-Correction



Improving the results of existing OCR systems

"Ma ti e*x*i' pu klei'?"
"Iklèo *ka* ìtela n'armastò."
I *vèkkia* àggale tria dattilìtia:

"Ma ti e*x*i' pu klei'?"
"Ikleo ka ìtela _armastò."
I vekkia aaggale tria dattilitia:

OCR output ("first pass")

Automatic OCR
Post-Correction



"Ma ti e*x*i' pu klei'?"
"Iklèo ka ìtela n'armastò."
I vèkkia àggale tria dattilìtia:

Corrected transcription

Improving the results of existing OCR systems

”Ma ti eḡi’ pu klei’?”
 ”Iklèo ka ìtela n’armastò.”
 I vèkkia àggale tria dattilìtia:

”Ma ti e*x*i’ pu klei’?”
 ”Ikl*e*o ka ìtela armastò.”
 I v*e*kkia *a*ggale tria dattil*i*tia:

OCR output (“first pass”)

Automatic OCR
 Post-Correction

”Ma ti e*x*i’ pu klei’?”
 ”Iklè*o* ka ìtela *n*’armastò.”
 I vèkkia àggale tria dattilìtia:

Corrected transcription

Previous work: improve results for
 unseen fonts, layouts, domains.
 This talk: low-resourced languages.

Improving the results of existing OCR systems

”Ma ti exi’ pu klei?”
”Ikleo ka itela _armastò.”
I vekkia aggale tria dattilitia:

OCR output (“first pass”)

Automatic OCR
Post-Correction

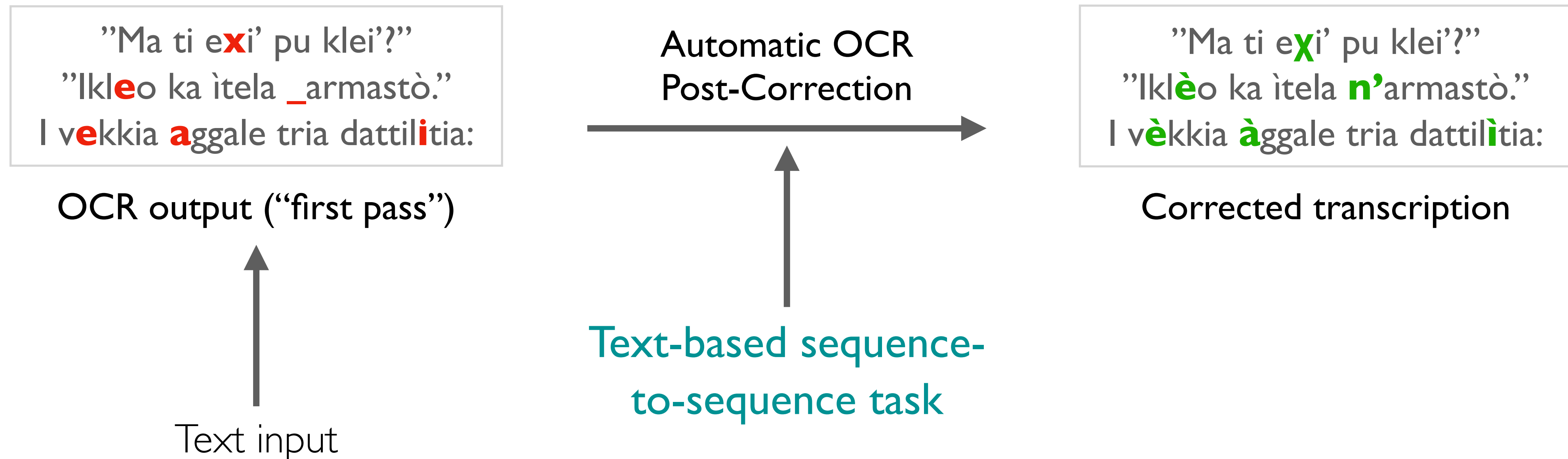


”Ma ti exi’ pu klei?”
”Iklèo ka itela n’armastò.”
I vèkkia àggale tria dattilitia:

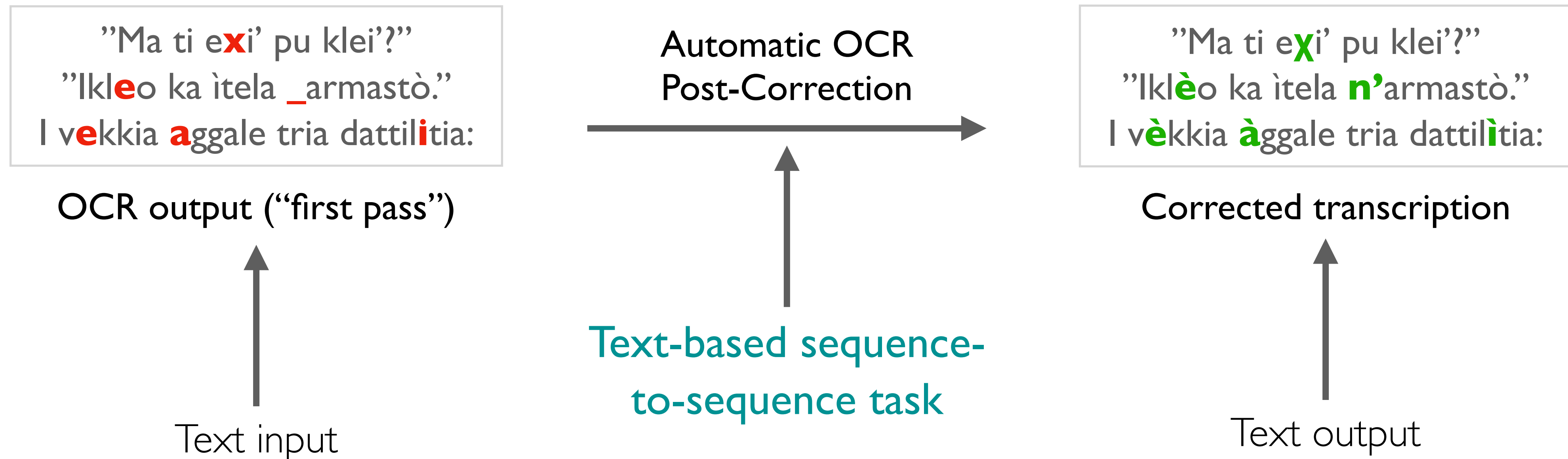
Corrected transcription

Text-based sequence-
to-sequence task

Improving the results of existing OCR systems



Improving the results of existing OCR systems



Adapting to low-resource settings

Prior work: character-level encoder-decoder with attention

- Add structural biases to the model
 - Diagonal attention loss, copy mechanism, coverage mechanism

Adapting to low-resource settings

Prior work: character-level encoder-decoder with attention

- Add structural biases to the model
 - Diagonal attention loss, copy mechanism, coverage mechanism
- Leverage additional information from the source document

What additional information is available?

Matiaxh	Khunik	jos.om	marimpa
Mathias	John	work wood (tv).agent	marimba
Matiaxh	Khunik	wood-worker	(of) marimbass

ruwe-ne noine
 poro ape are wa
 hekota rok wa
 uweneusar⁽¹⁾
 kor okai.
 Inkar ne wa
 akip ne korka
 ine-ap-kusu
 arushka wa

ものの如く
 澤山に火を焚きて
 そこに向ひて坐して
 昔噺やお伽などをし
 つゝ暮らし居たりき。
 只それを見るのみ
 にはあれど
 いかばかり
 わが腹立たしくて

ཐག་ཉེ་ (thag-ngea) adv. near,
 close or at a shorter distance.

ཐག་ཀོ་ (thag-ko) n. rope.

ཐག་ཚོད་ (thag-choth) v. to be
 dedicated/settled/resolved.

ཐག་མཚོང་ (thag-chong) n. Rope
 skipping, jumping. v. to skip, to
 jump.

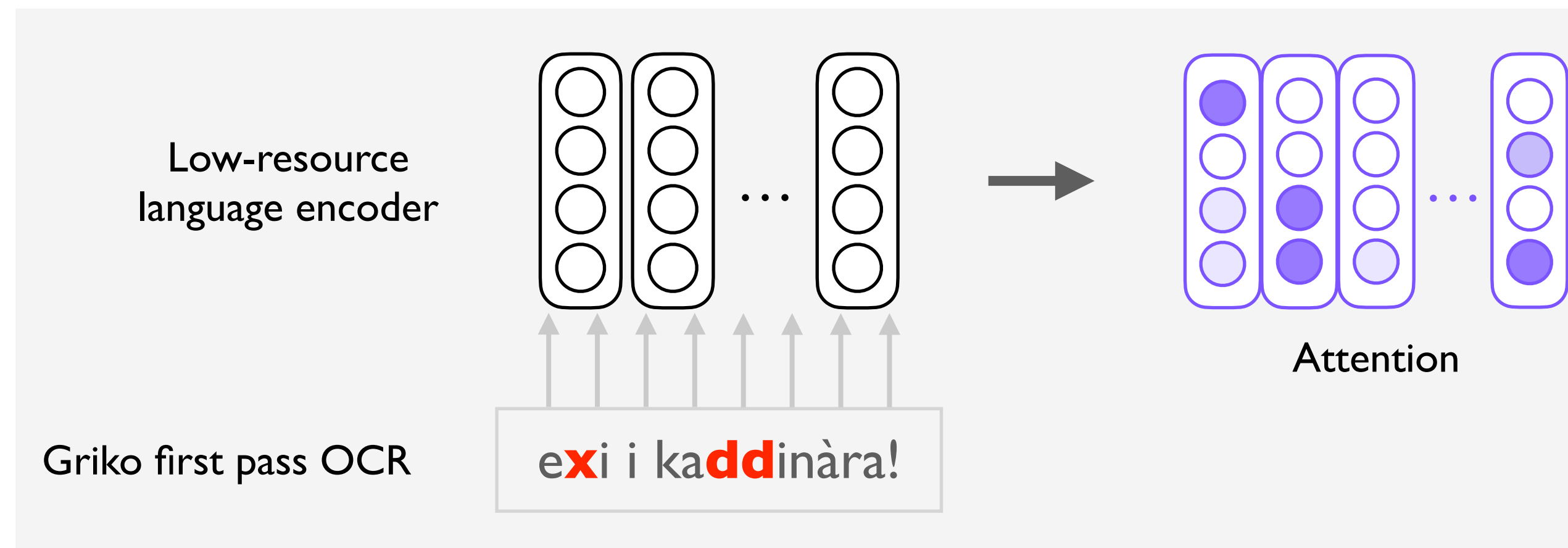
- Many documents containing text in low-resource languages also **contains a translation of the text**
- Interlinear glosses, dictionaries, linguistic documentation, language learning material...

Seall thall thar an aiseig am fasnadh nan craobh,
 Am bothan beag glan ud, 's e gealaicht' le aol ;
 Sud agaibh mo dhachaidh : 's i dachaidh mo ghaoil,
 Gun chaisteal 'san t-saoghal a 's feàrr leam.

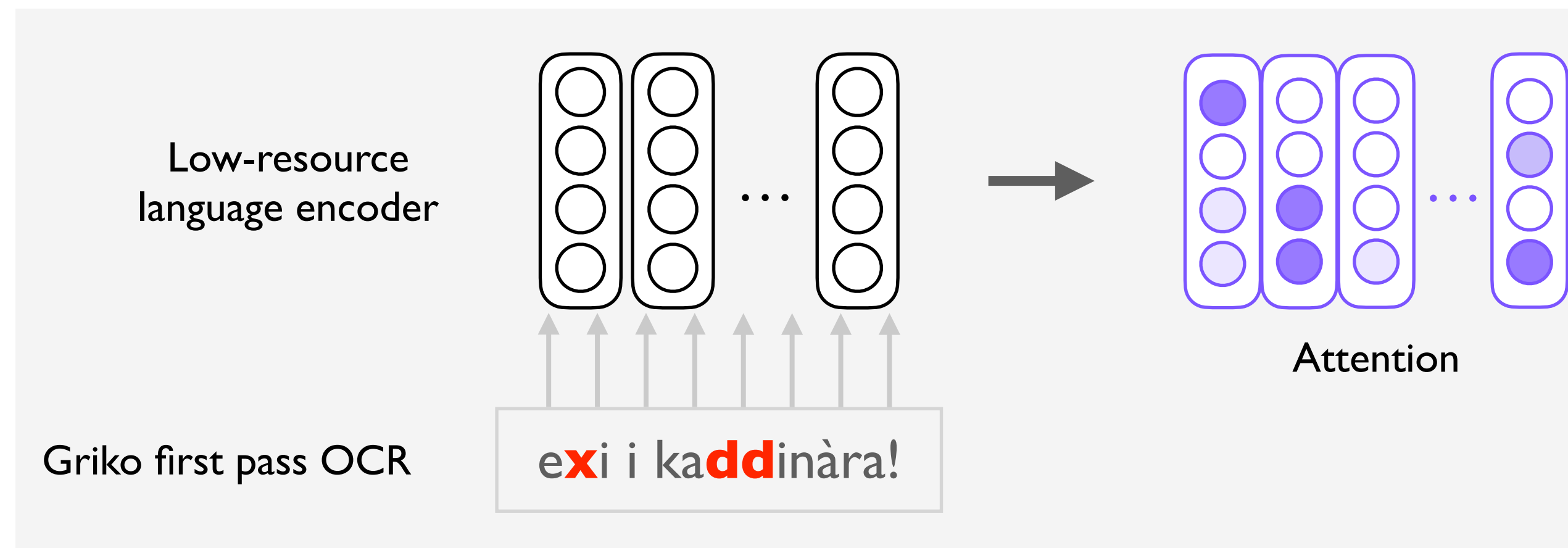
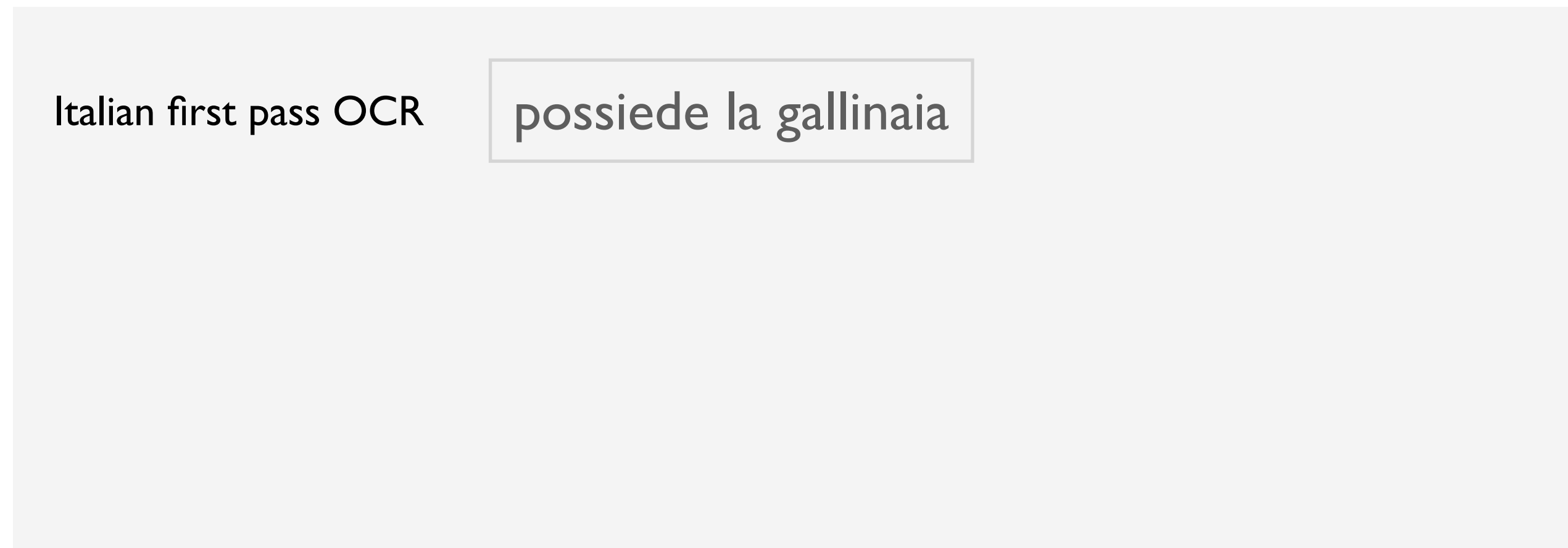
Ayont by the ferry, whaur woodlands are green,
 My cantie cot housie stan's tidy an' clean ;
 I envy nae laird in his castle, I ween,
 I'm happy an' bien in my ain house.

Multi-source model for post-correction

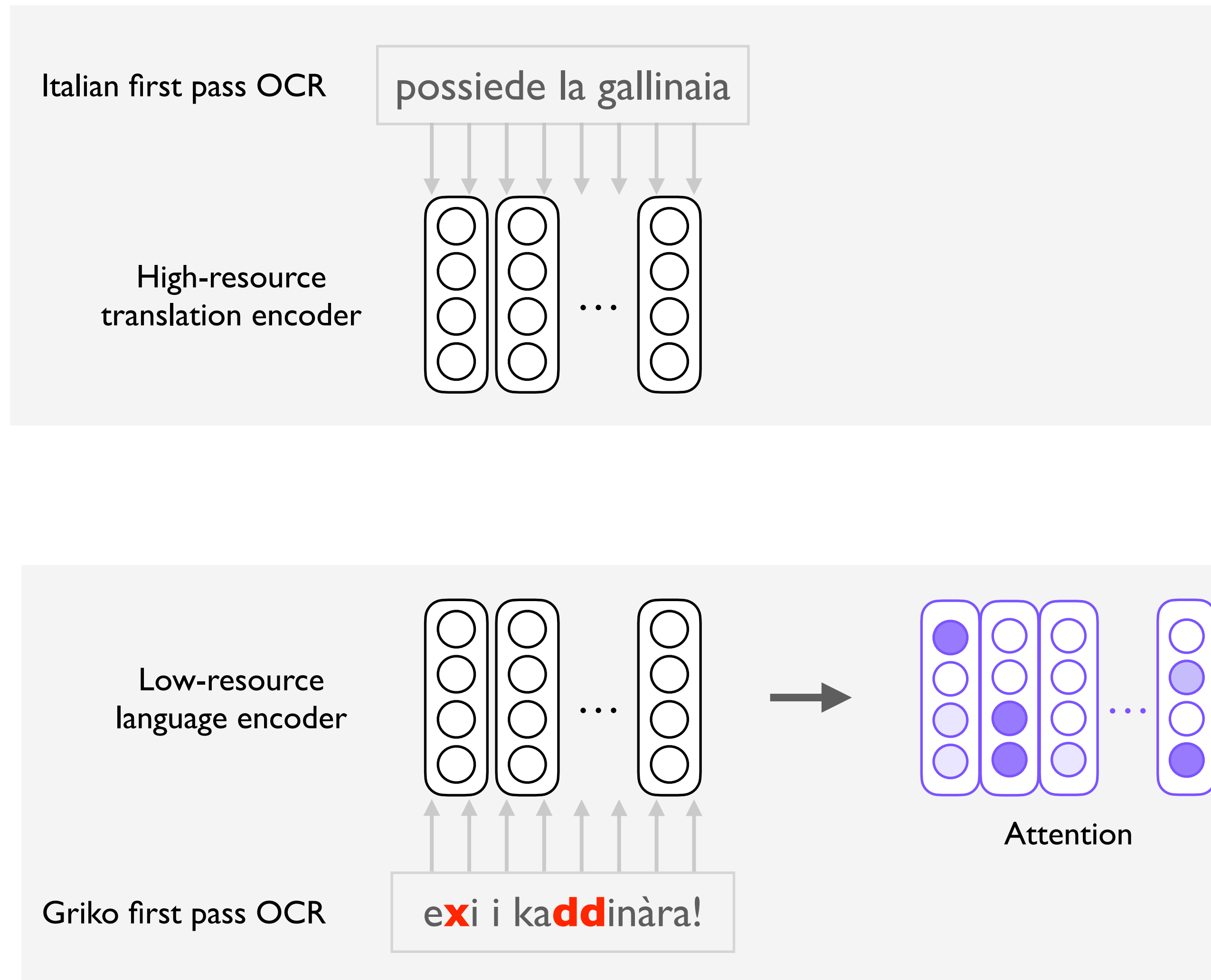
Multi-source model for post-correction



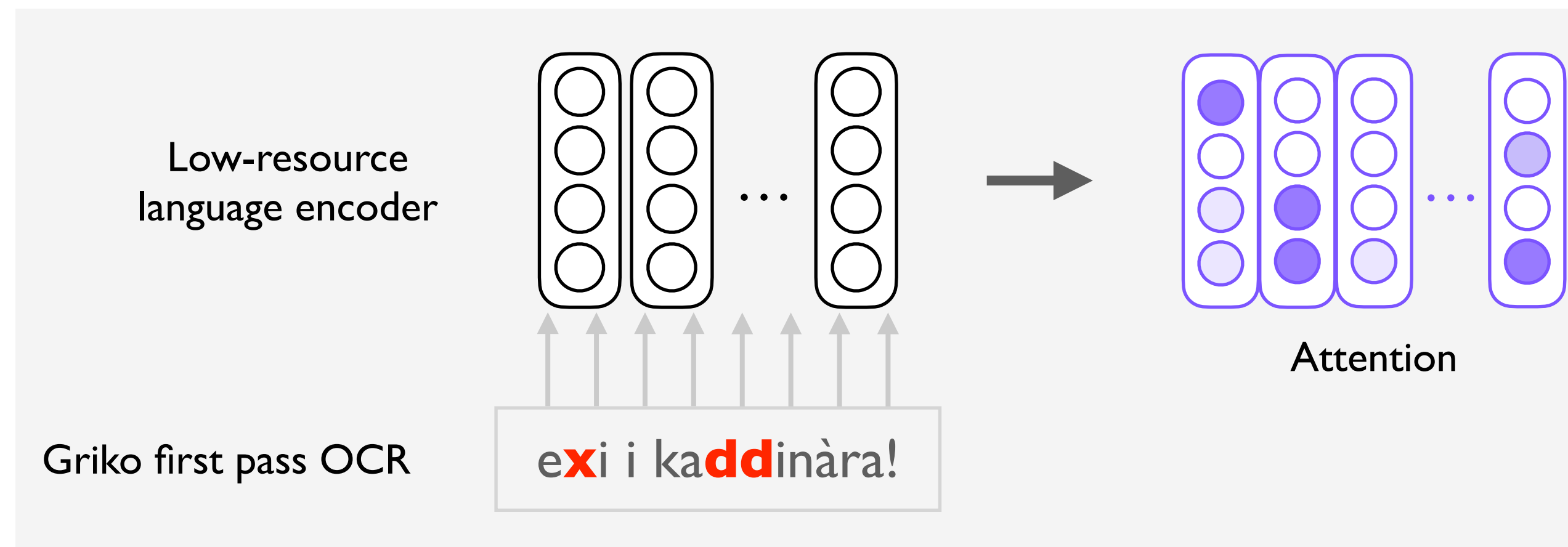
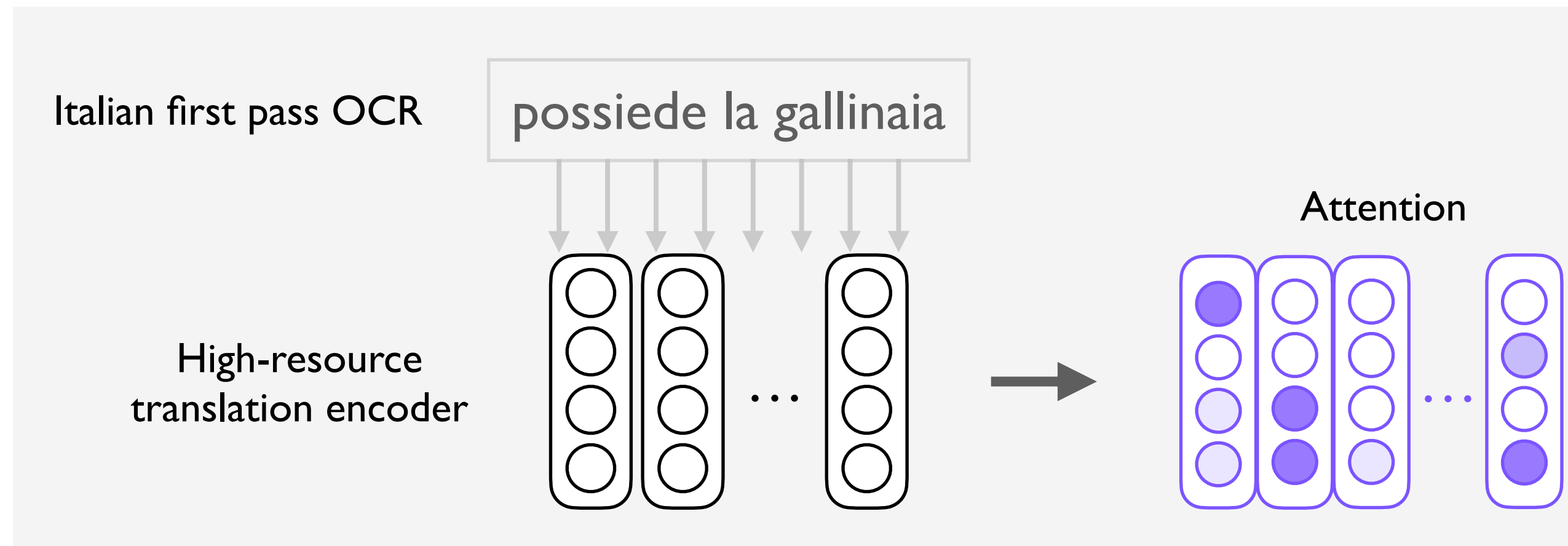
Multi-source model for post-correction



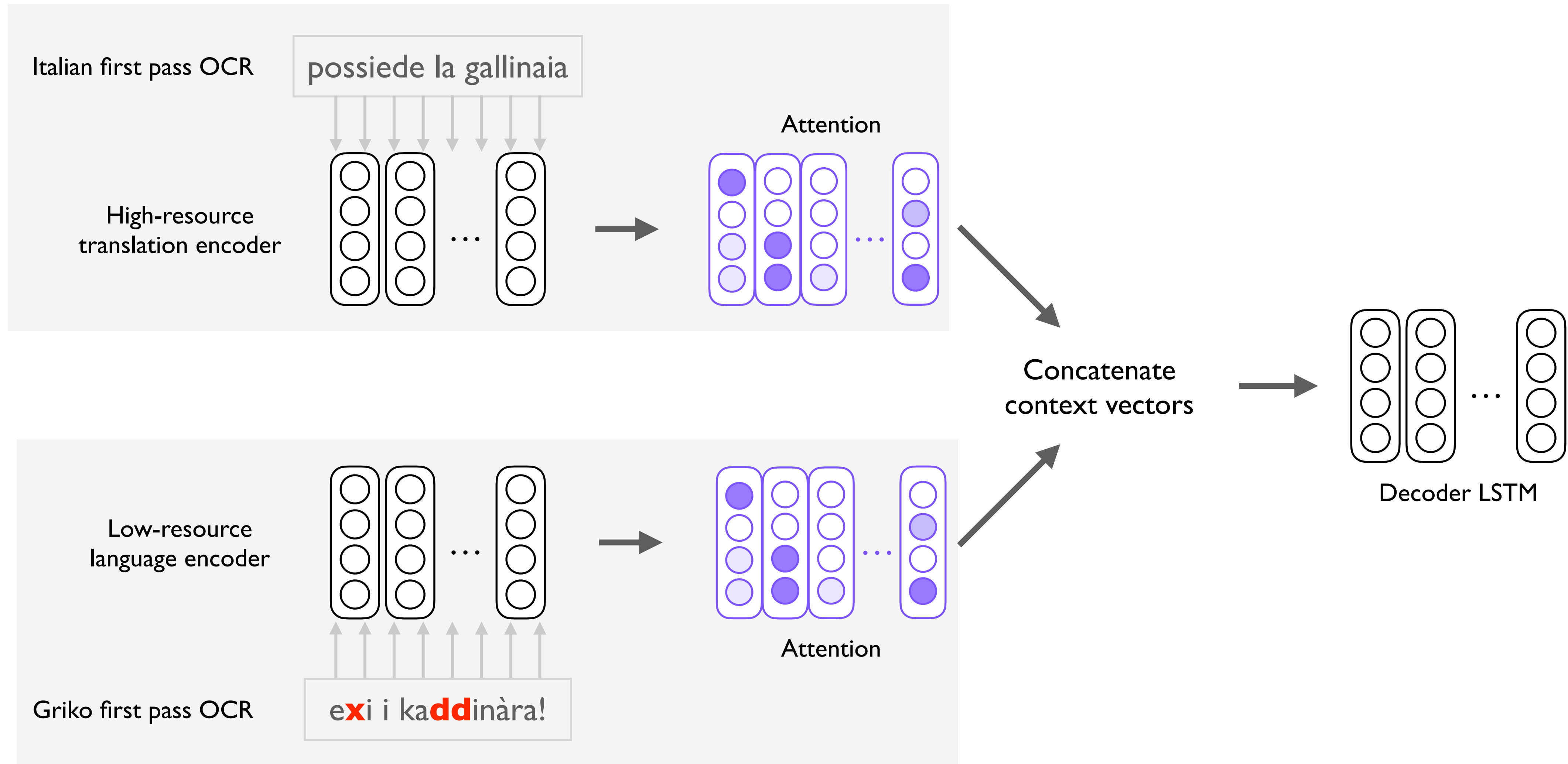
Multi-source model for post-correction



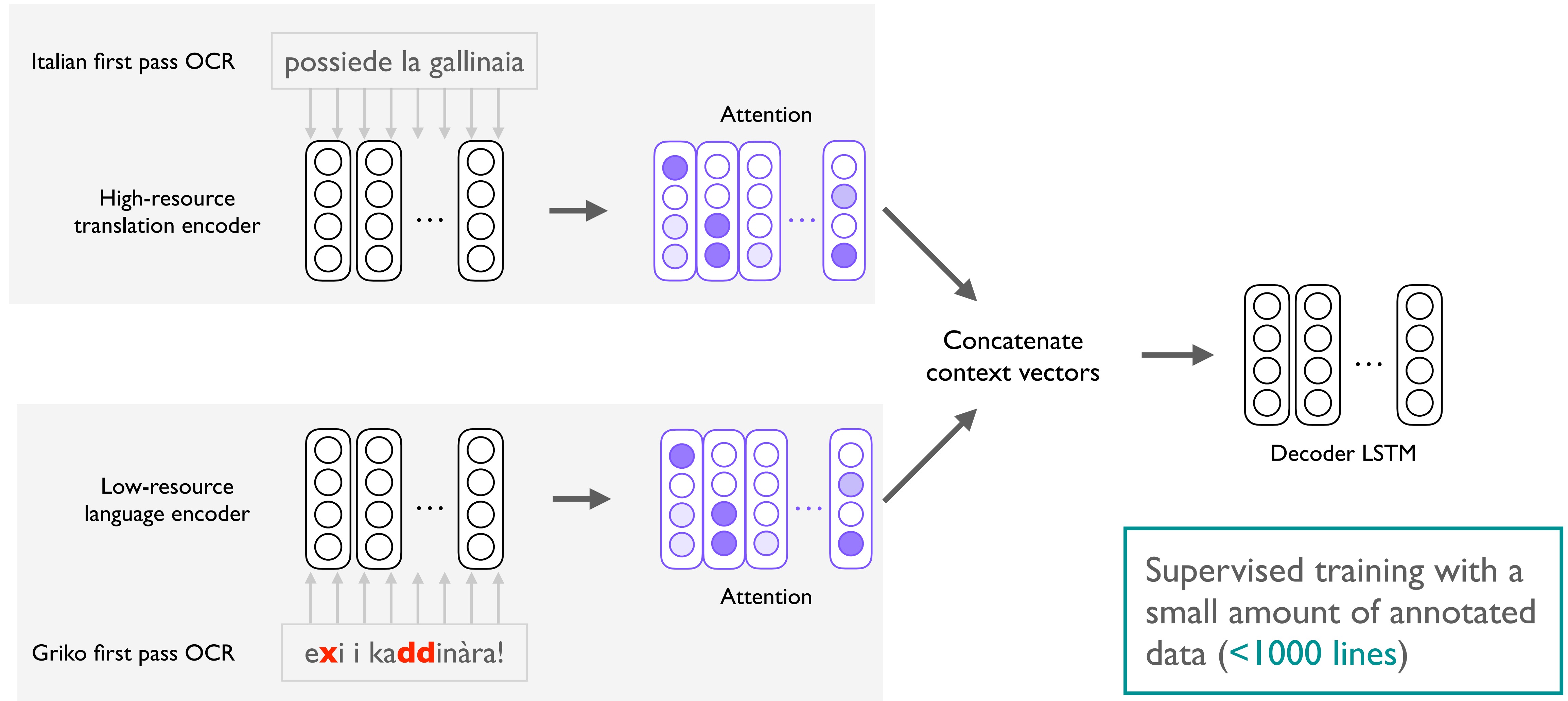
Multi-source model for post-correction



Multi-source model for post-correction

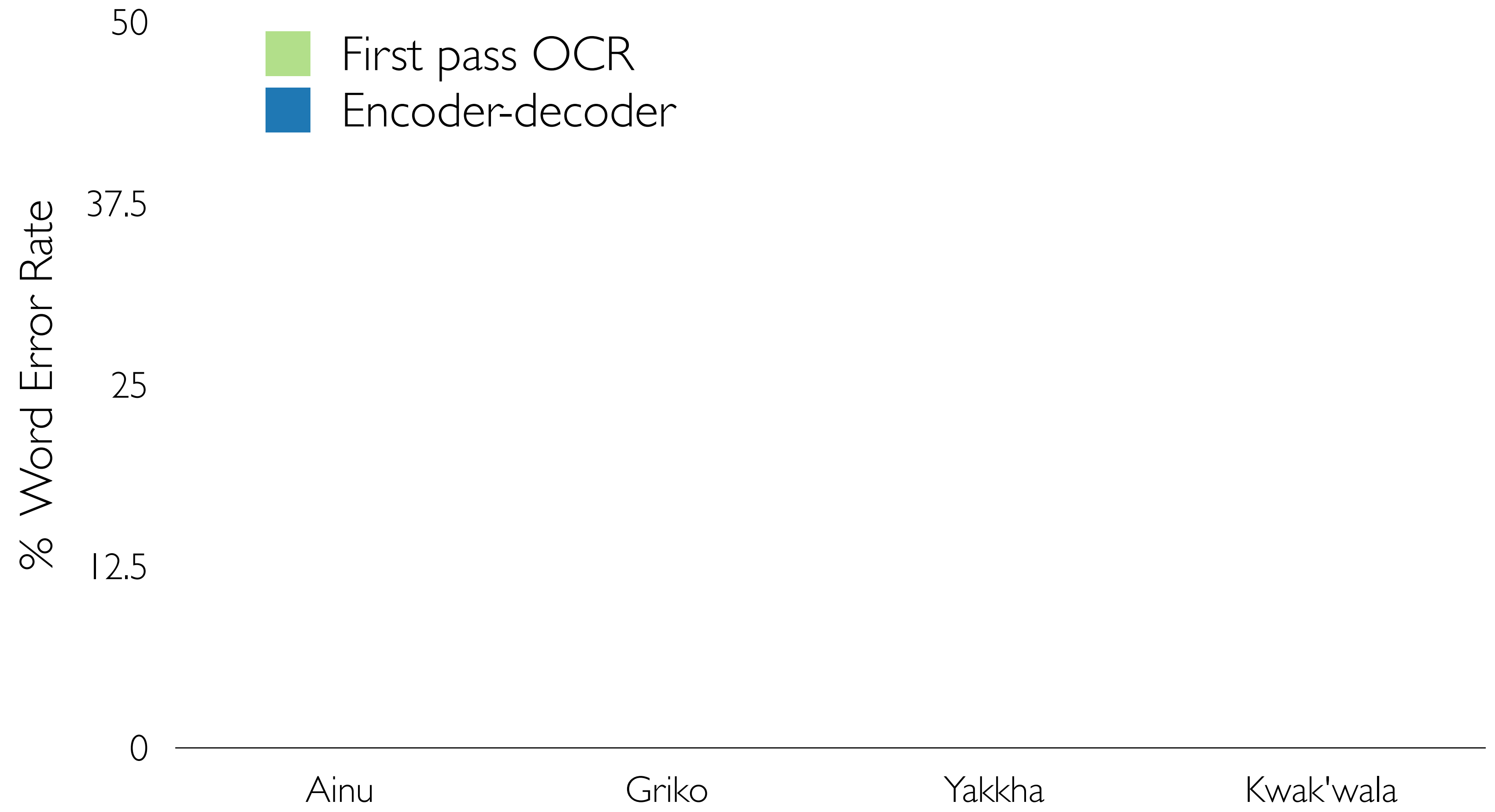


Multi-source model for post-correction

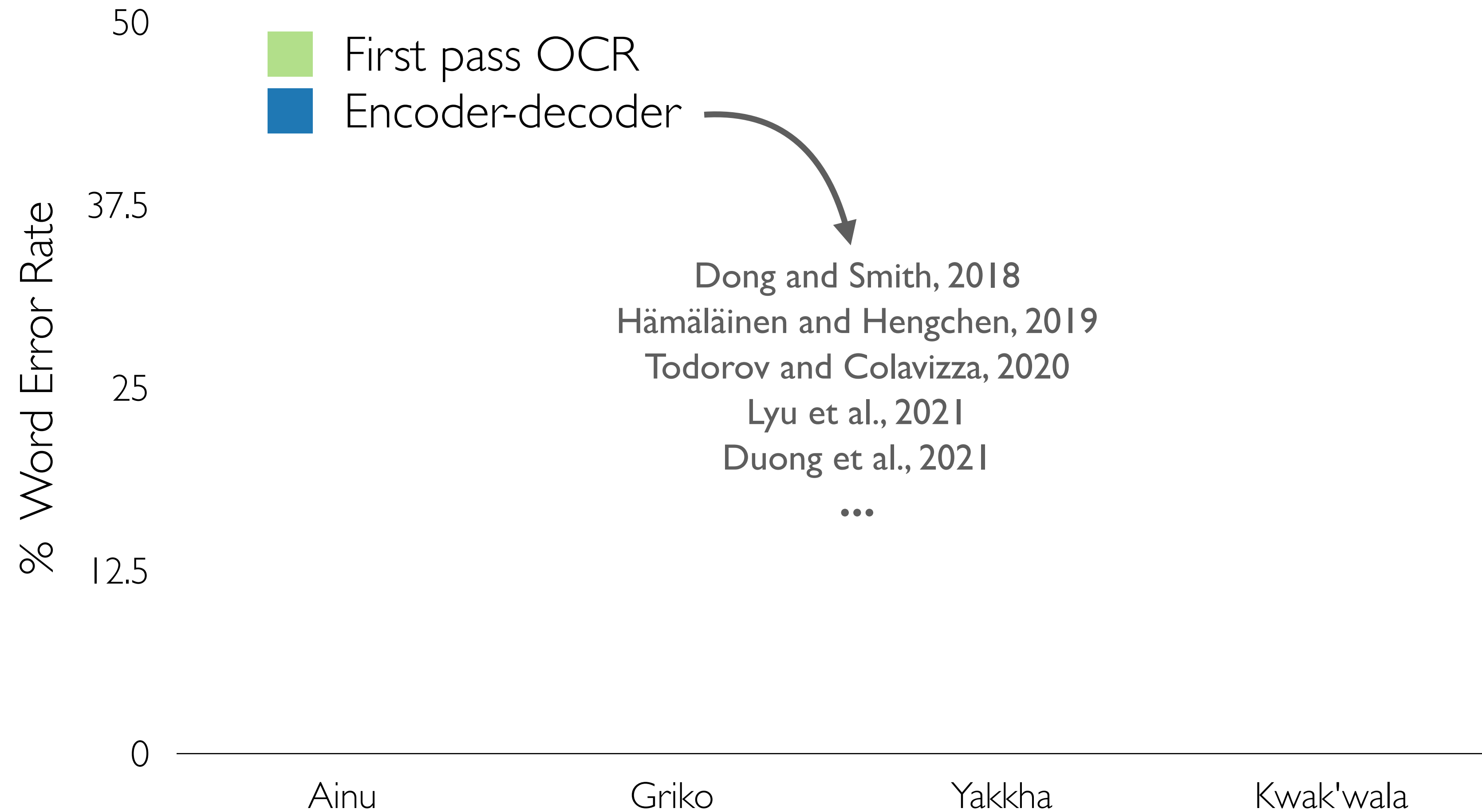


Experiments: how do existing post-correction methods perform?

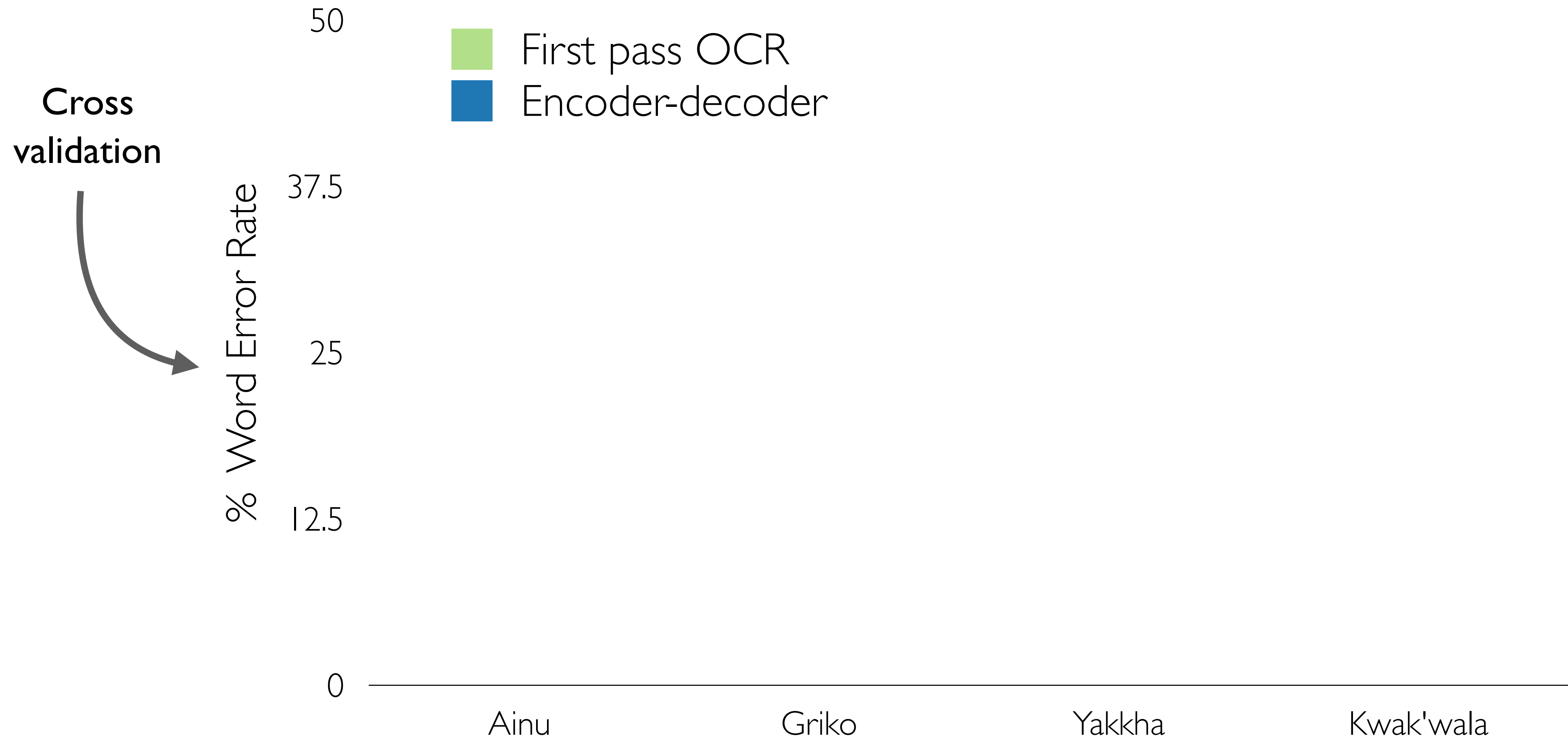
Experiments: how do existing post-correction methods perform?



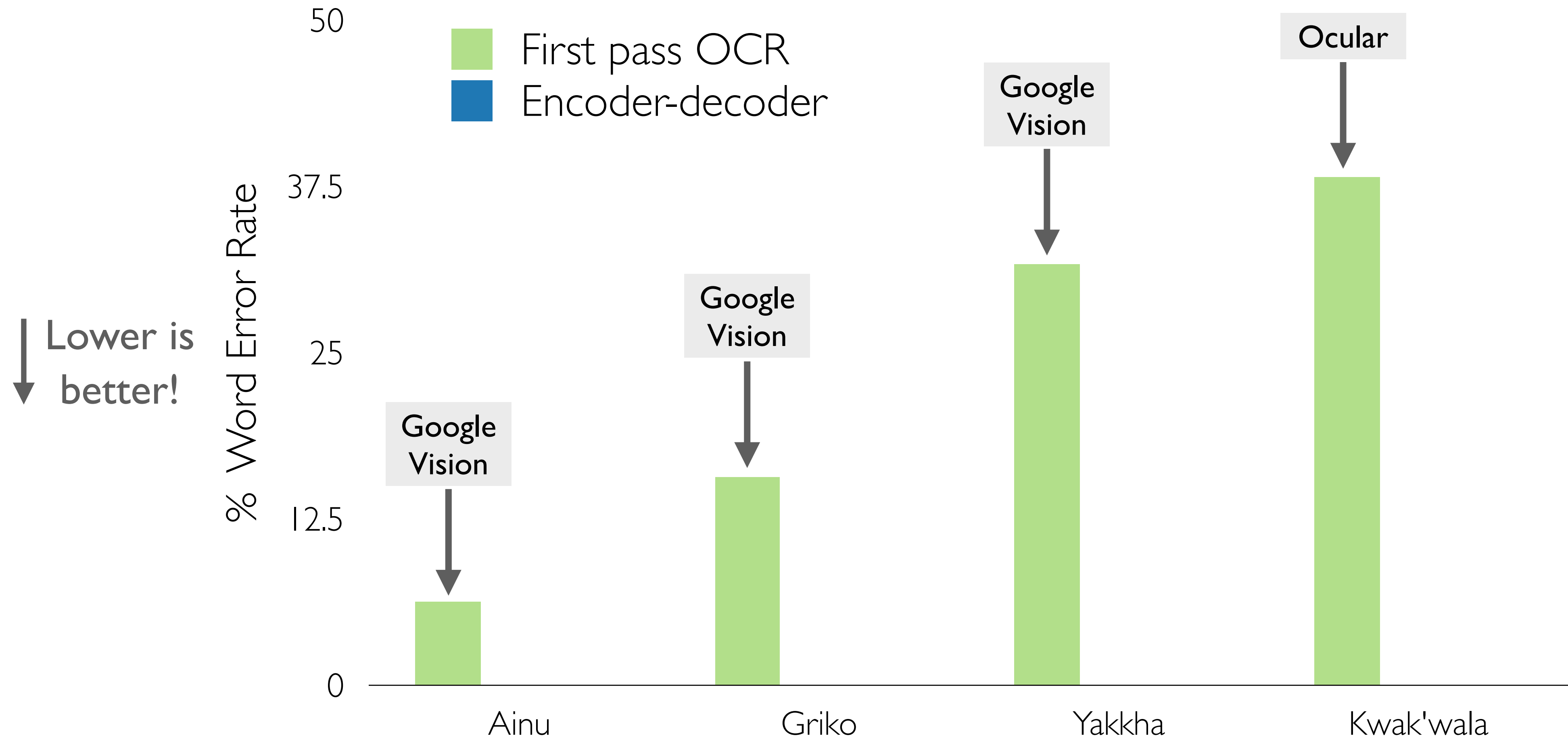
Experiments: how do existing post-correction methods perform?



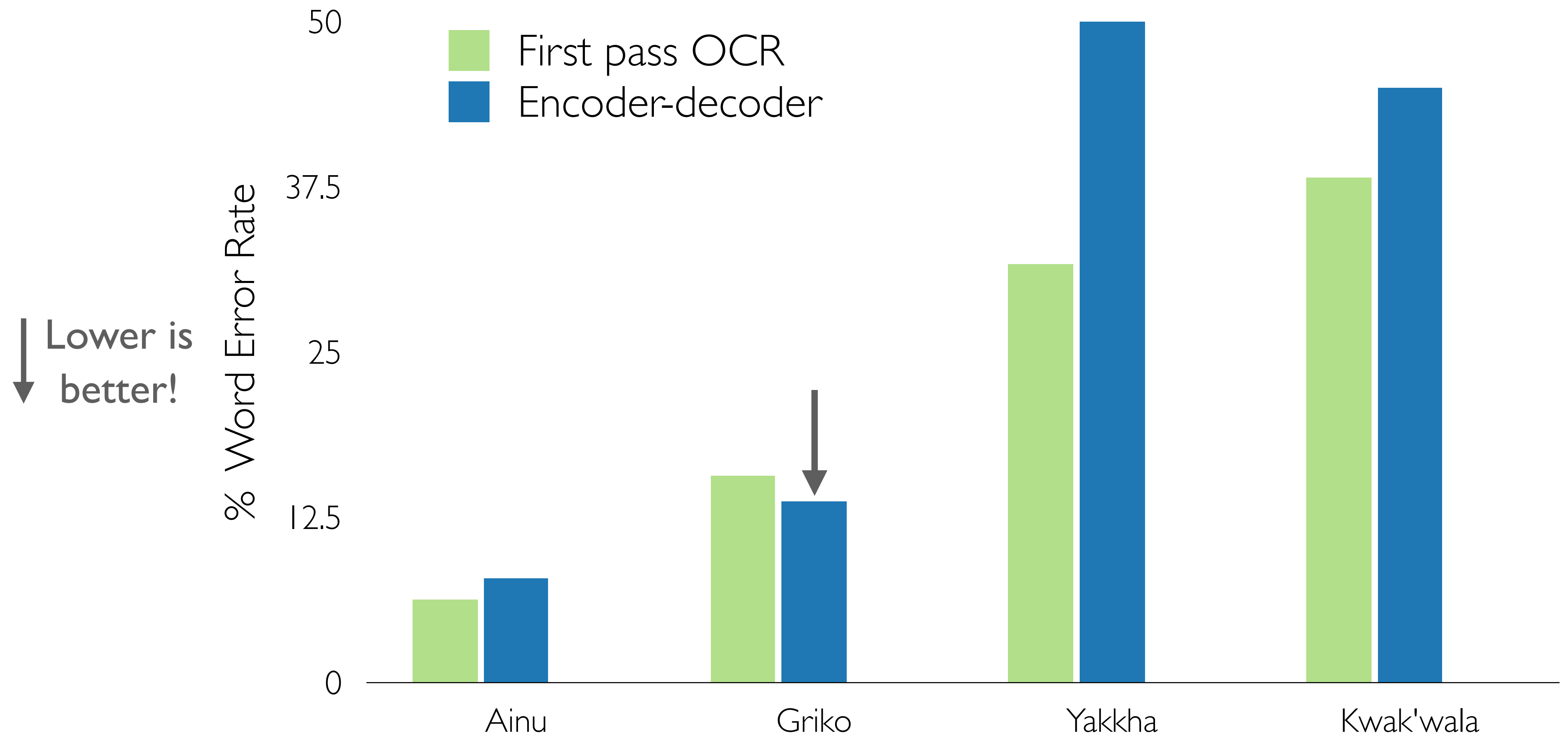
Experiments: how do existing post-correction methods perform?



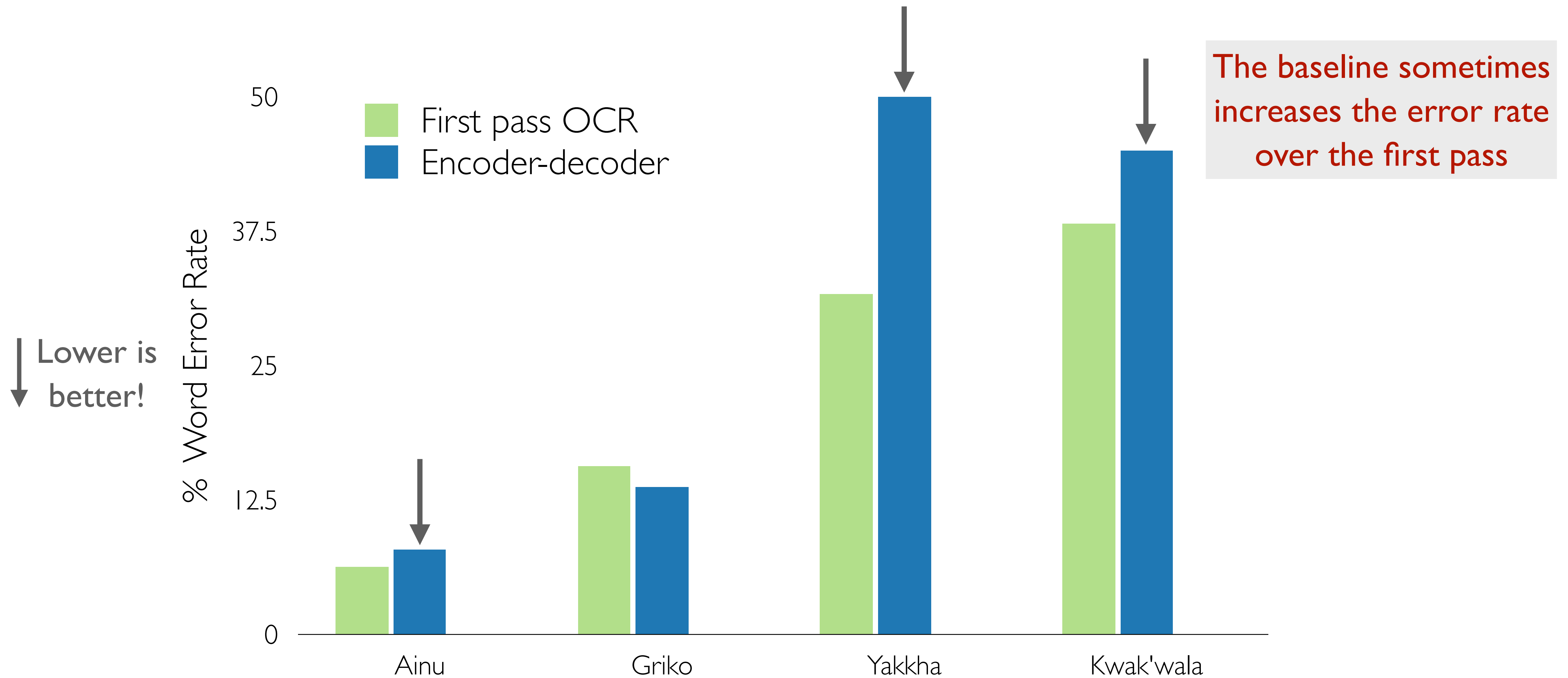
Experiments: how do existing post-correction methods perform?



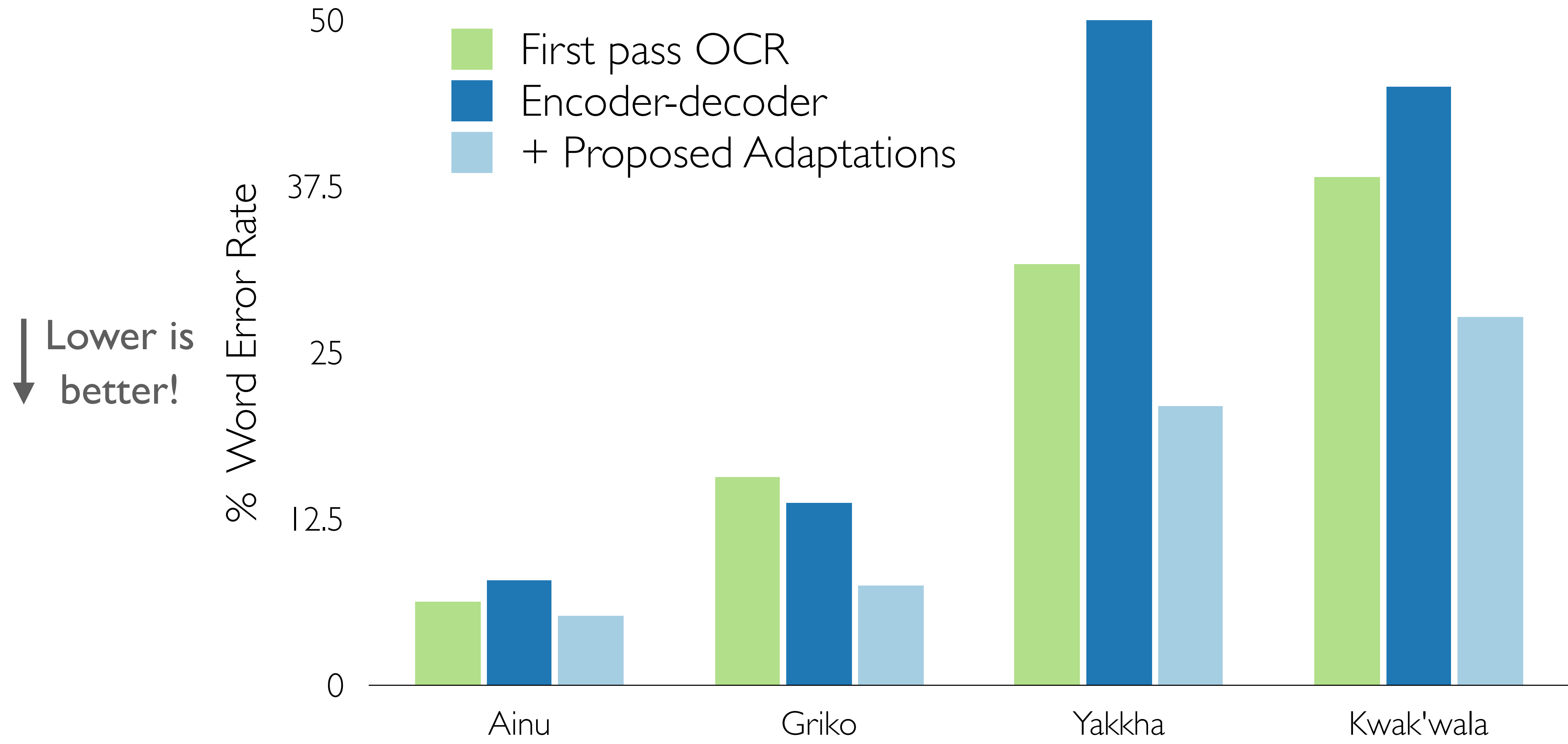
Experiments: how do existing post-correction methods perform?



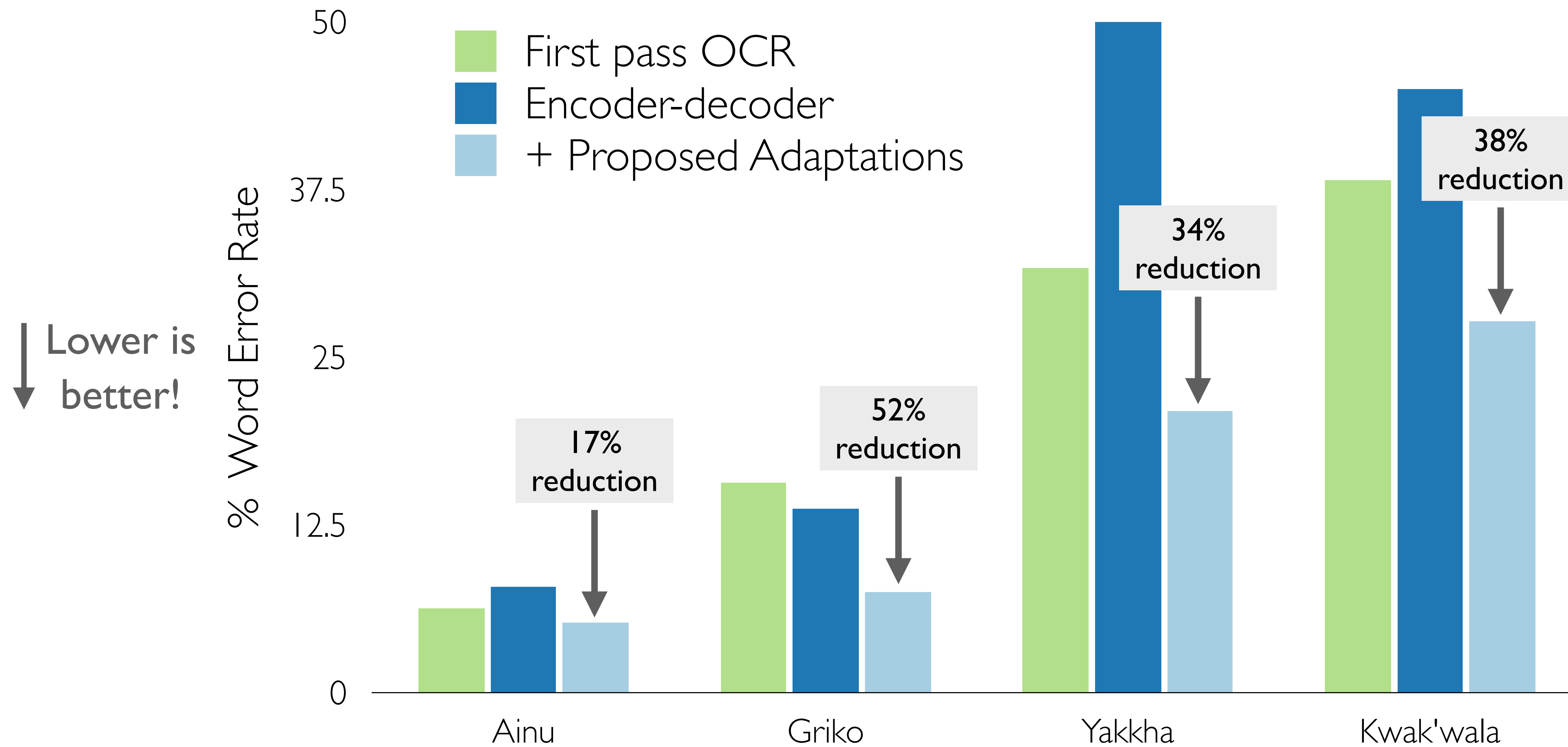
Experiments: how do existing post-correction methods perform?



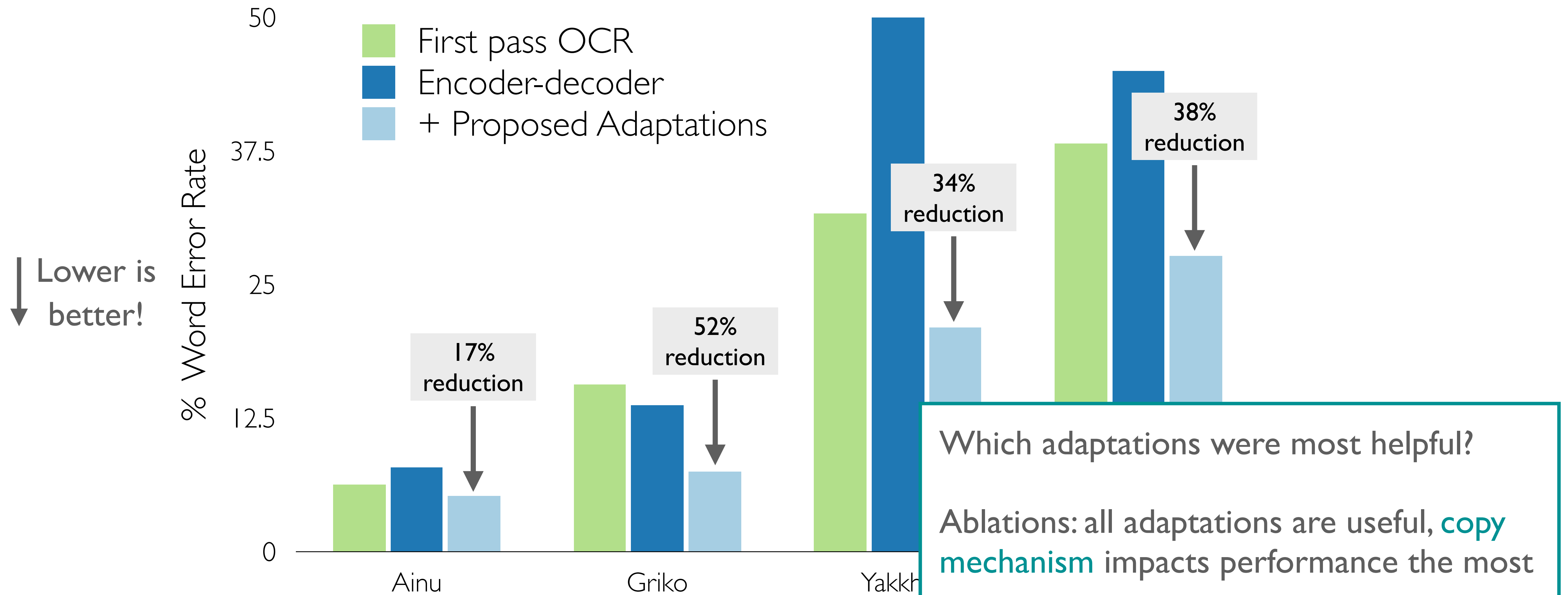
Experiments: do the adaptations help low-resource learning?



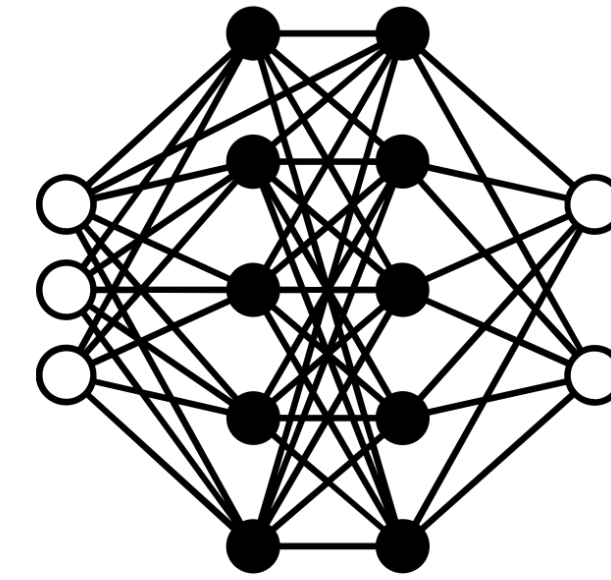
Experiments: do the adaptations help low-resource learning?



Experiments: do the adaptations help low-resource learning?

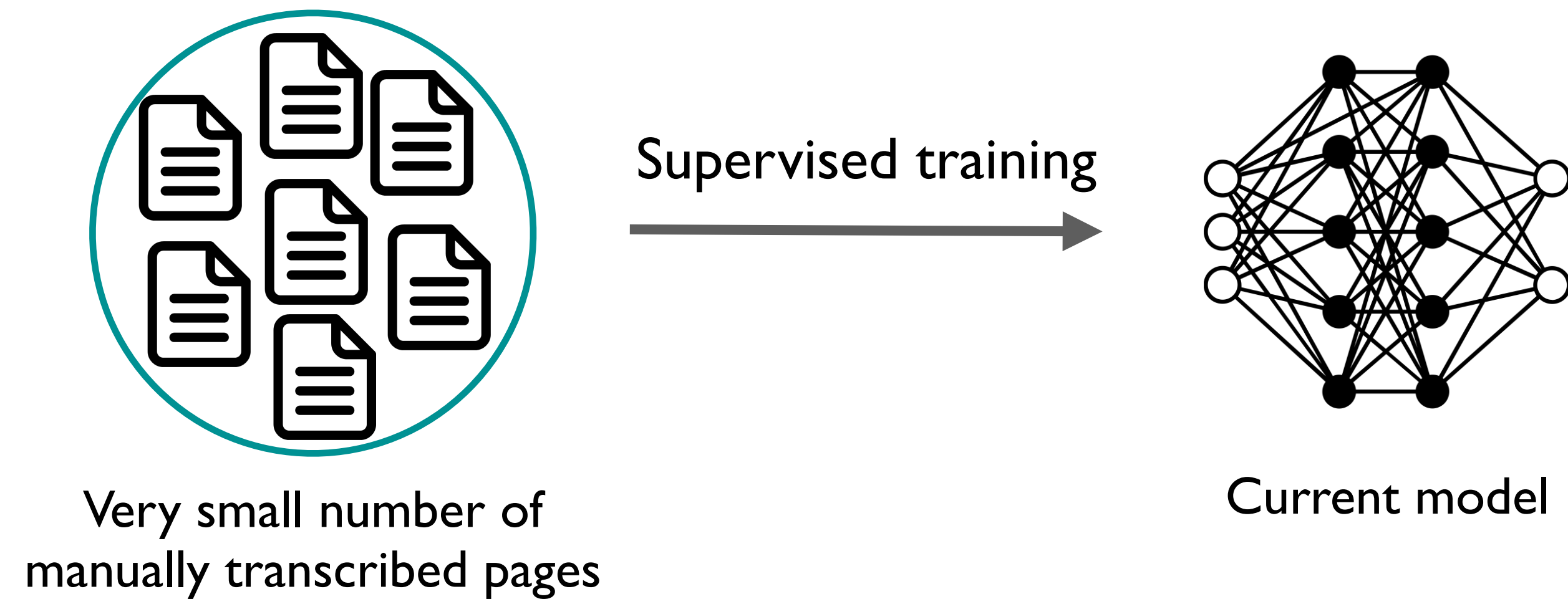


Improving performance without additional annotation



Current model

Improving performance without additional annotation



Improving performance without additional annotation

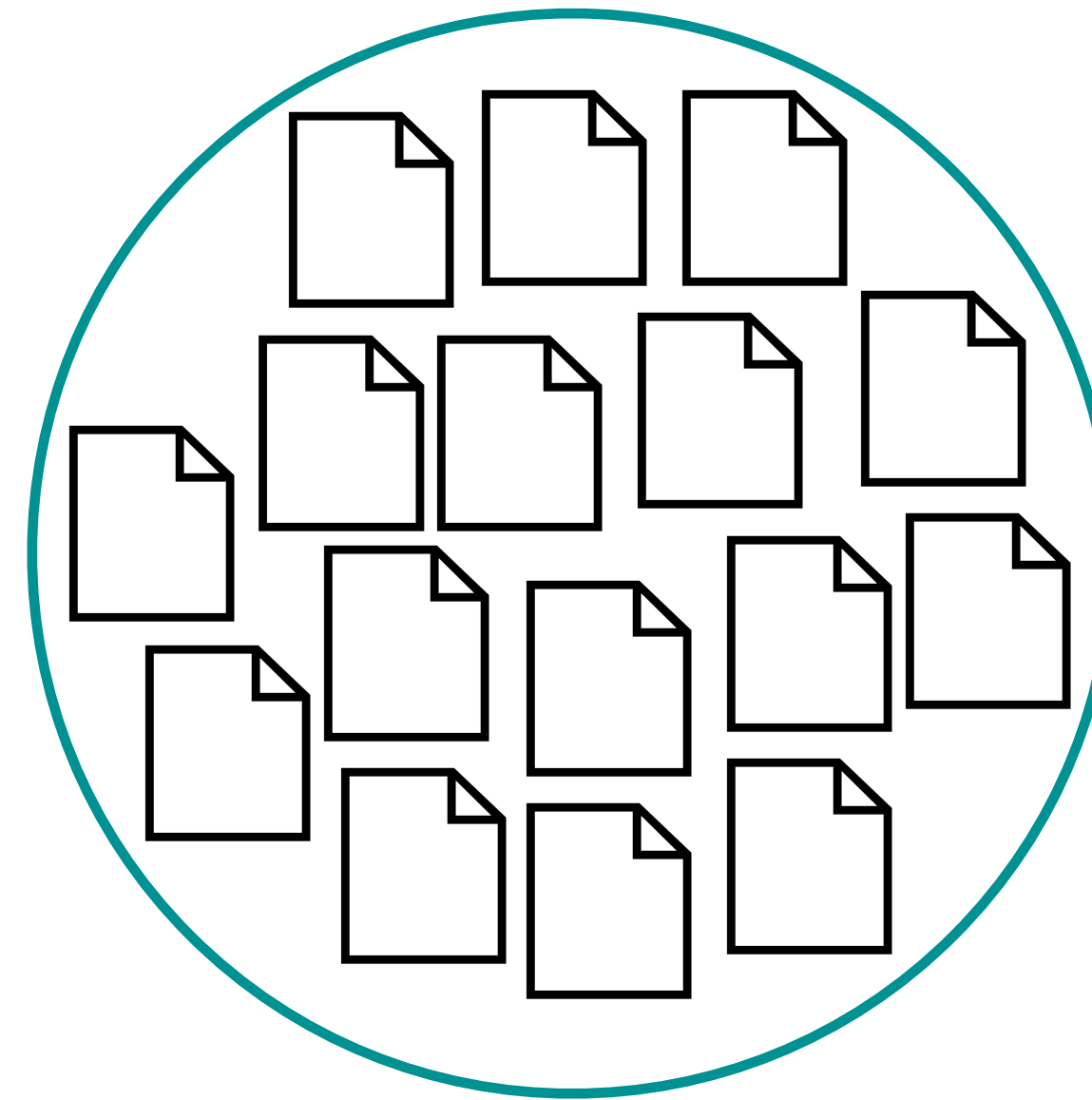


Very small number of
manually transcribed pages

Improving performance without additional annotation



Very small number of
manually transcribed pages

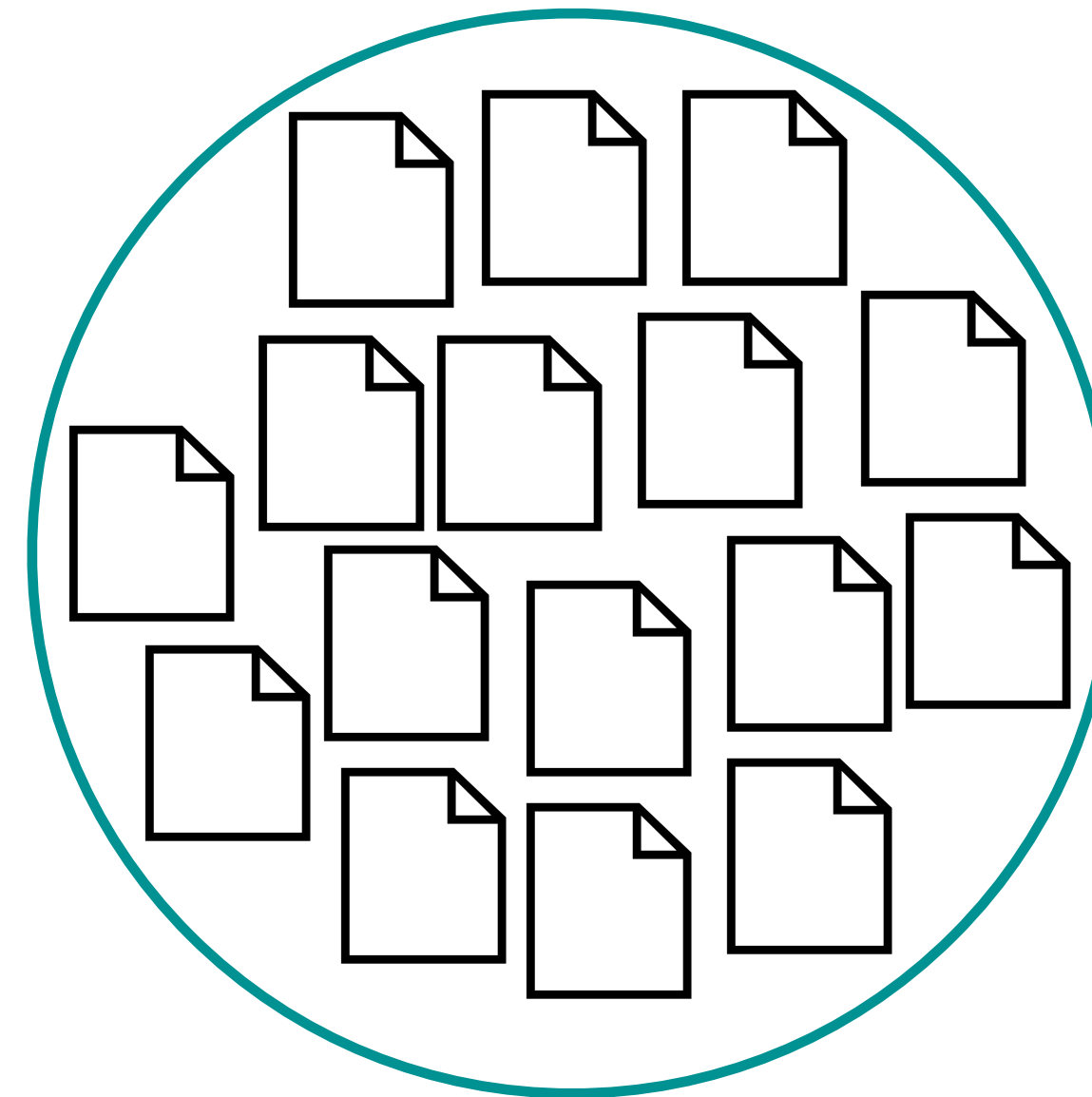


Relatively larger number
of raw images that need
to be digitized

Improving performance without additional annotation



Very small number of manually transcribed pages



Relatively larger number of raw images that need to be digitized

Our dataset:

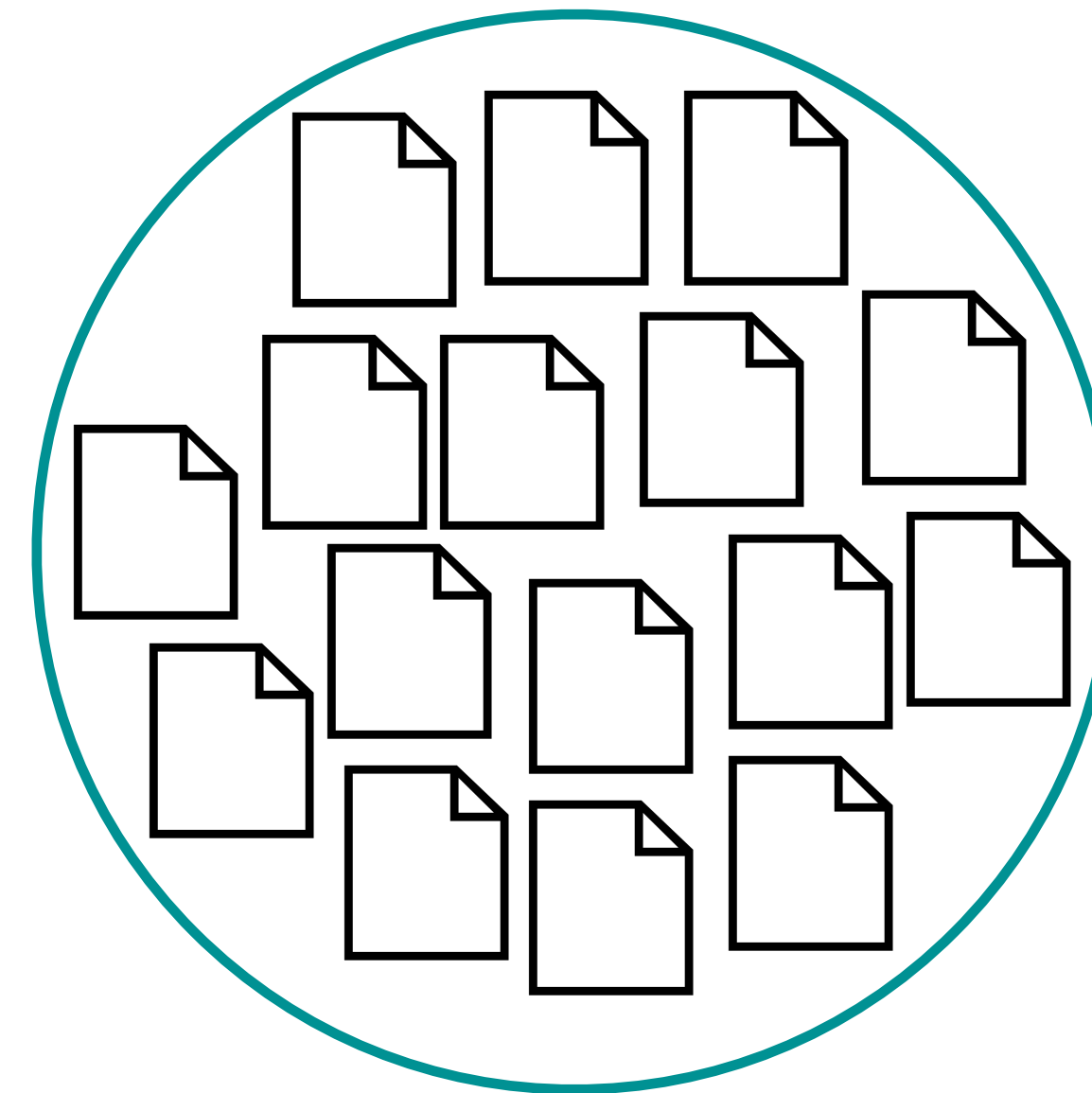
Documents contain 300 – 800 pages

Only ~30 are manually transcribed

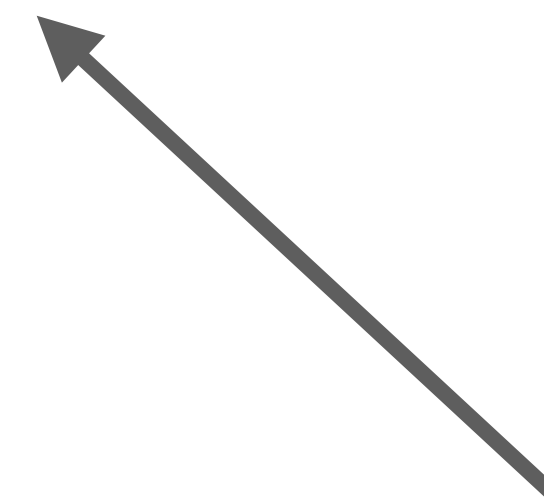
Improving performance without additional annotation



Very small number of manually transcribed pages



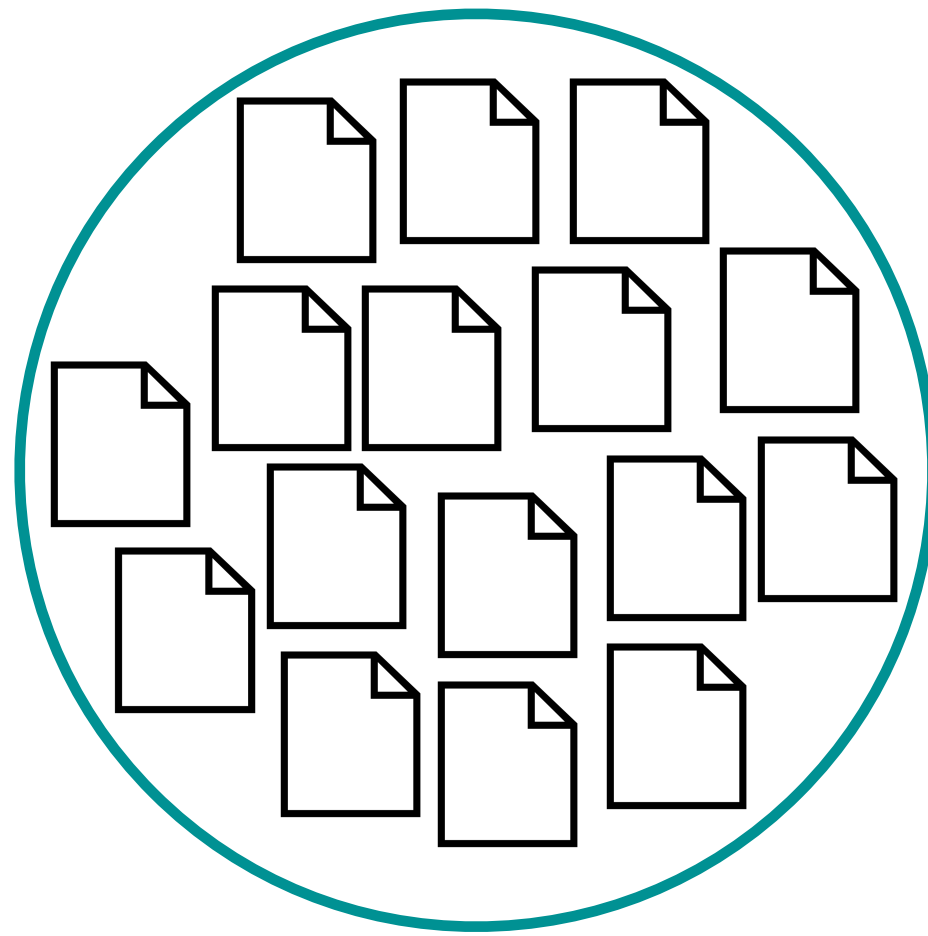
Relatively larger number of raw images that need to be digitized



Semi-supervised learning for efficient use of the unlabeled images

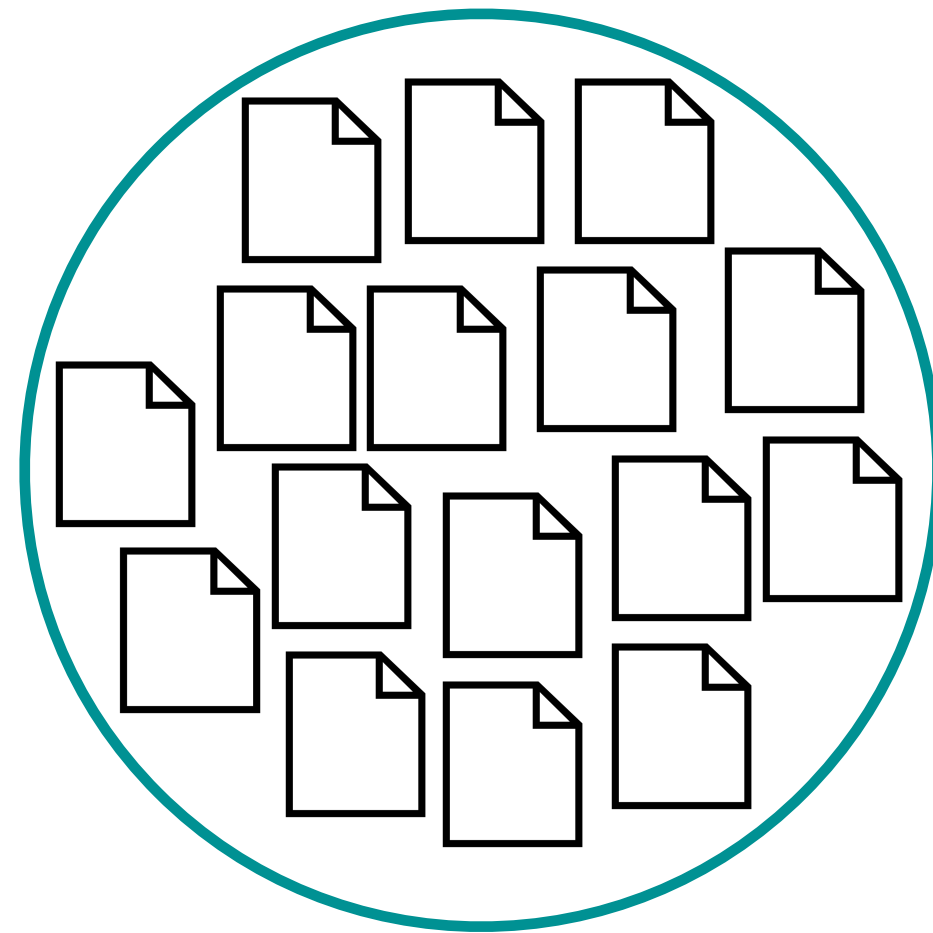
Self-training for OCR post-correction

Self-training for OCR post-correction

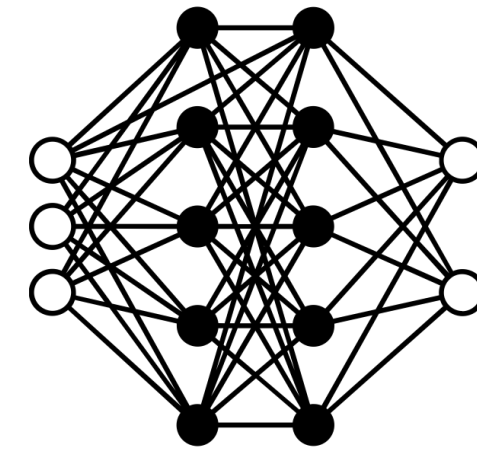


First pass OCR on
unlabeled images

Self-training for OCR post-correction

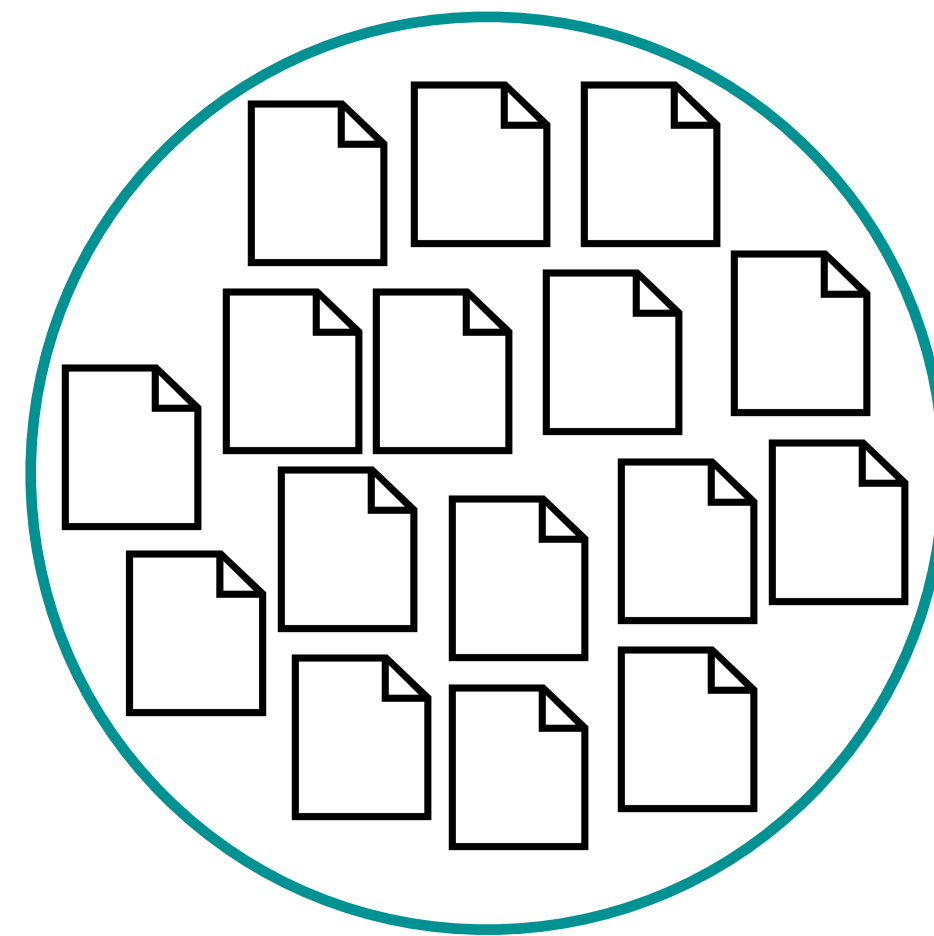


First pass OCR on
unlabeled images

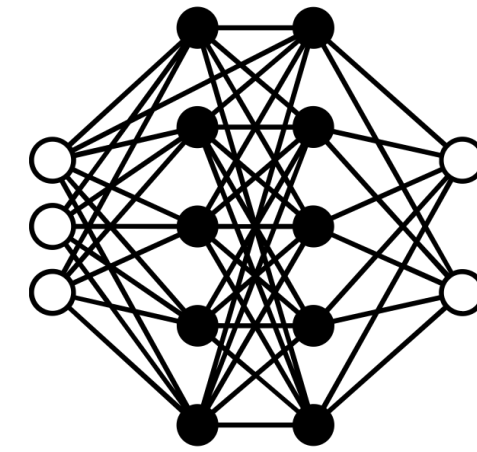


Trained post-
correction model

Self-training for OCR post-correction



First pass OCR on
unlabeled images

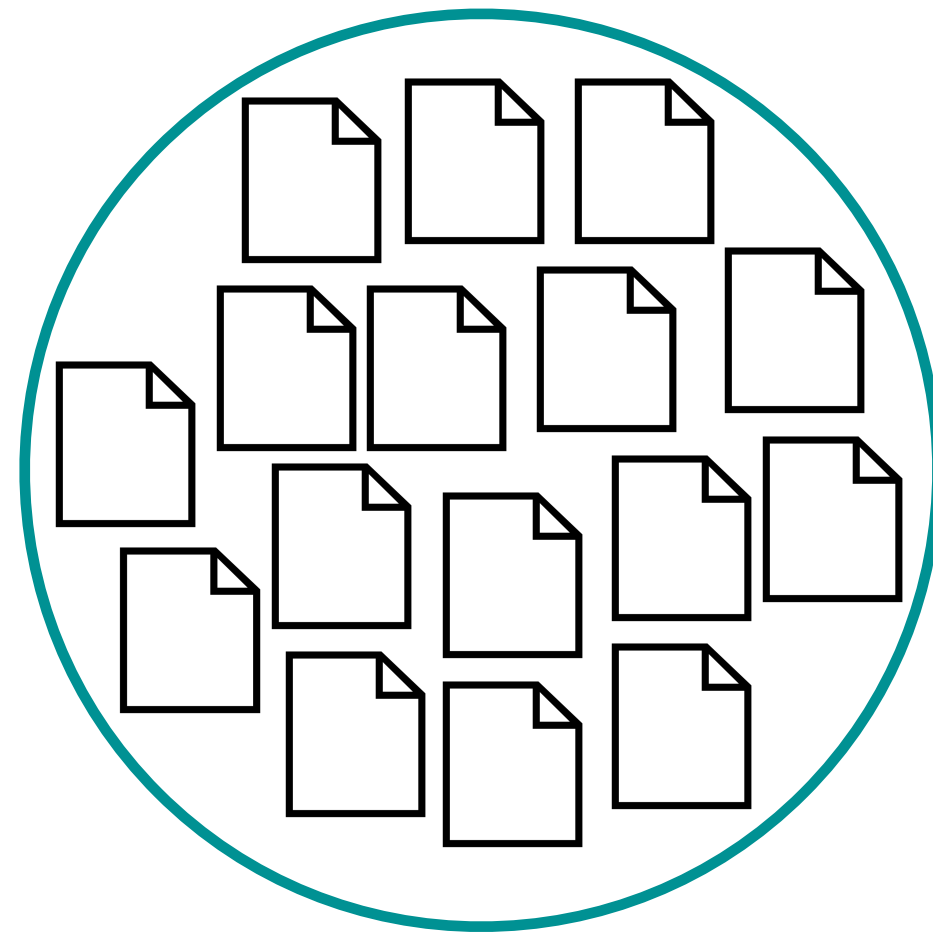


Trained post-
correction model

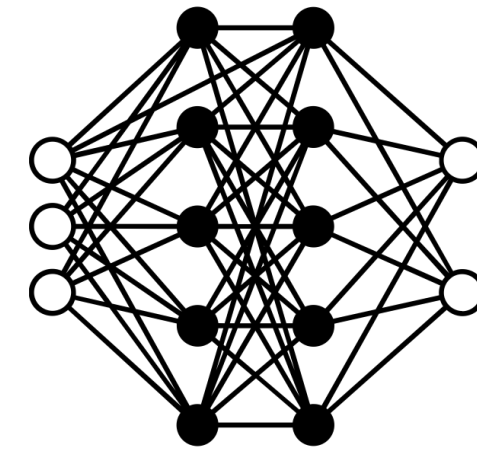


Previous best
supervised model

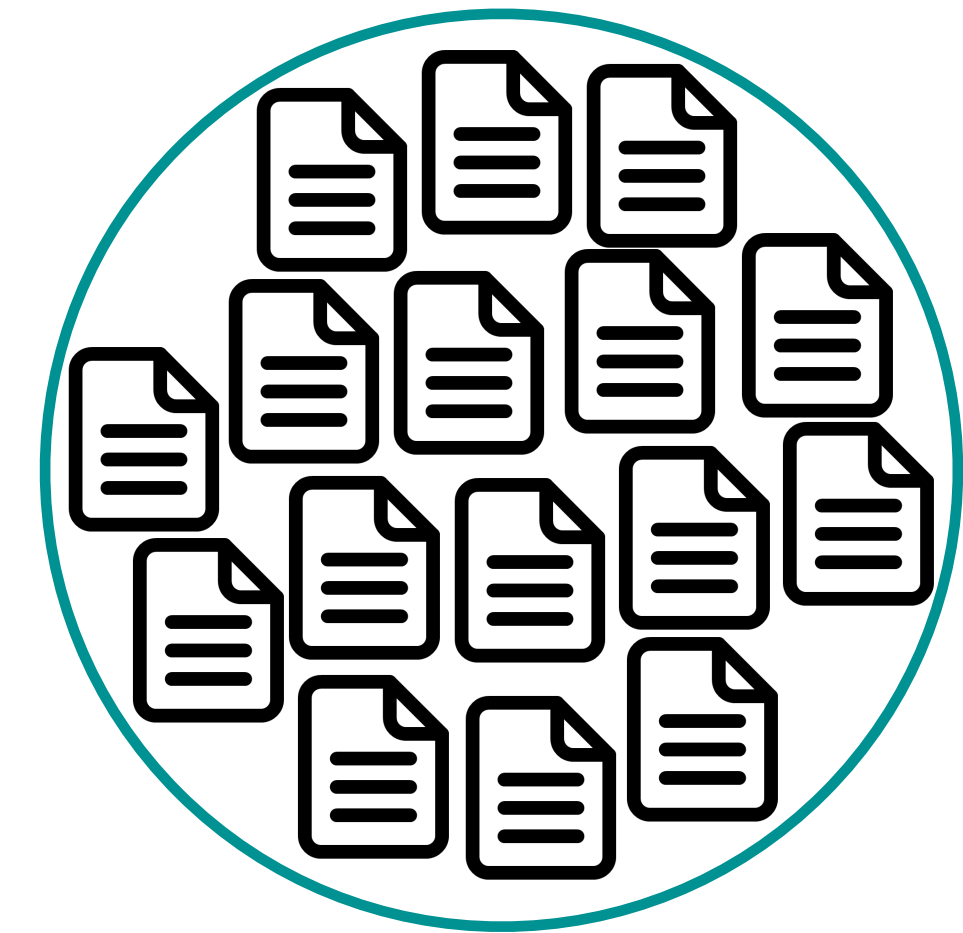
Self-training for OCR post-correction



First pass OCR on
unlabeled images

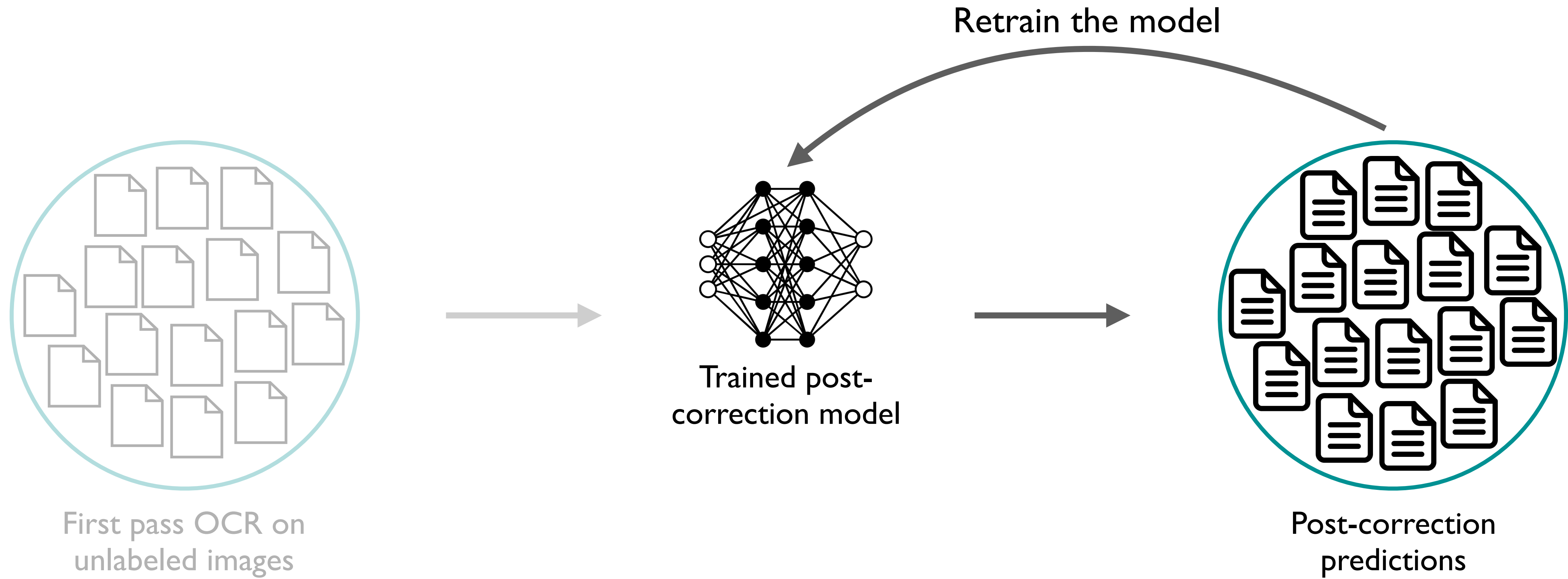


Trained post-
correction model



Post-correction
predictions

Self-training for OCR post-correction



Self-training may introduce noise

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē



g·îl^εmēsē (7)

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē



g·îl^εmēsē (7)



g·il^εmēsē (5)

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē



g·îl^εmēsē (7)



g·il^εmēsē (5)



g·îl^εm^εsē (2)

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

✓ g·îl^εmēsē (7)

✗ g·il^εmēsē (5)

✗ g·îl^εm^εsē (2)

✗ g·îl^εmi^εsē (2)

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

✓ g·îl^εmēsē (7)

✗ g·il^εmēsē (5)

✗ g·îl^εm^εsē (2)

✗ g·îl^εmi^εsē (2)

} Different subsets
of characters are
incorrect

↑ Empirical observations

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāg·alis lāx

...

g·îl^εmēsē lāg·aa lāqēxs laē

✓ g·îl^εmēsē (7)

✗ g·il^εmēsē (5)

✗ g·îl^εm^εsē (2)

✗ g·îl^εmi^εsē (2)

Noise from self-training
is typically inconsistent
at the word-level

Empirical observations

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē ḡwālamasqēxs laē

. . . Wä, g·îl^εmēsē ^εnāx

Wä, g·îl^εmēsē lāḡ·alis lāx

...

g·îl^εmēsē lāḡ·aa lāqēxs laē

✓ g·îl^εmēsē (7)

✗ g·il^εmēsē (5)

✗ g·îl^εm^εsē (2)

✗ g·îl^εmi^εsē (2)

Correct form of the word ends up being more frequent

Noise from self-training is typically inconsistent at the word-level

Empirical observations

Self-training may introduce noise

Can we bias post-correction towards generating correct words?

denema lāq. Wä, g·îl^εmēsē

g·îl^εmēsē gwālamasqēxs laē

Can we use the word frequency information to bias the model towards correct forms?

...

g·îl^εmēsē lāg·aa lāqēxs laē

✓ g·îl^εmēsē (7)

✗ g·il^εmēsē (5)

✗ g·îl^εm^εsē (2)

✗ g·îl^εmi^εsē (2)

← Correct form of the word ends up being more frequent

Incorporating word frequency information

Incorporating word frequency information

$$P(y) = p_{\text{lstm}}(y)$$

Incorporating word frequency information

$$P(y) = p_{\text{lstm}}(y)$$



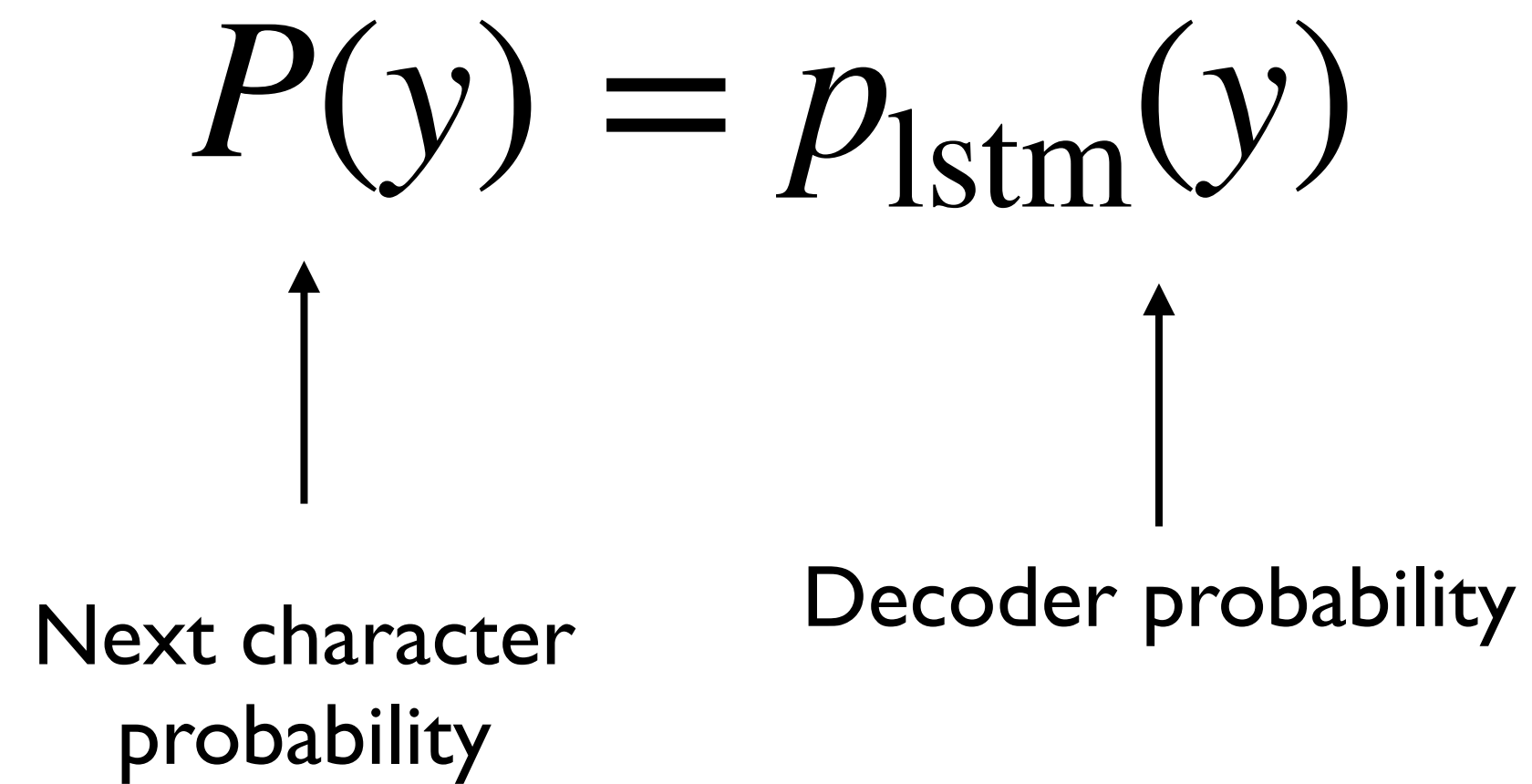
Next character
probability

Incorporating word frequency information

$$P(y) = p_{\text{lstm}}(y)$$

↑ ↑

Next character Decoder probability
probability

The diagram shows the equation $P(y) = p_{\text{lstm}}(y)$ centered at the top. Below the left side of the equation, an upward-pointing arrow connects the text 'Next character probability' to the $P(y)$ term. Below the right side of the equation, an upward-pointing arrow connects the text 'Decoder probability' to the $p_{\text{lstm}}(y)$ term.

Incorporating word frequency information

$$P(y) = p_{\text{lstm}}(y)$$

$$P_{\text{freq}}$$


Frequency-based
probability to explicitly
bias the model

Incorporating word frequency information

$$P(y) = p_{\text{lstm}}(y)$$

p_{freq}



How do we get
probabilities based on
word frequency?

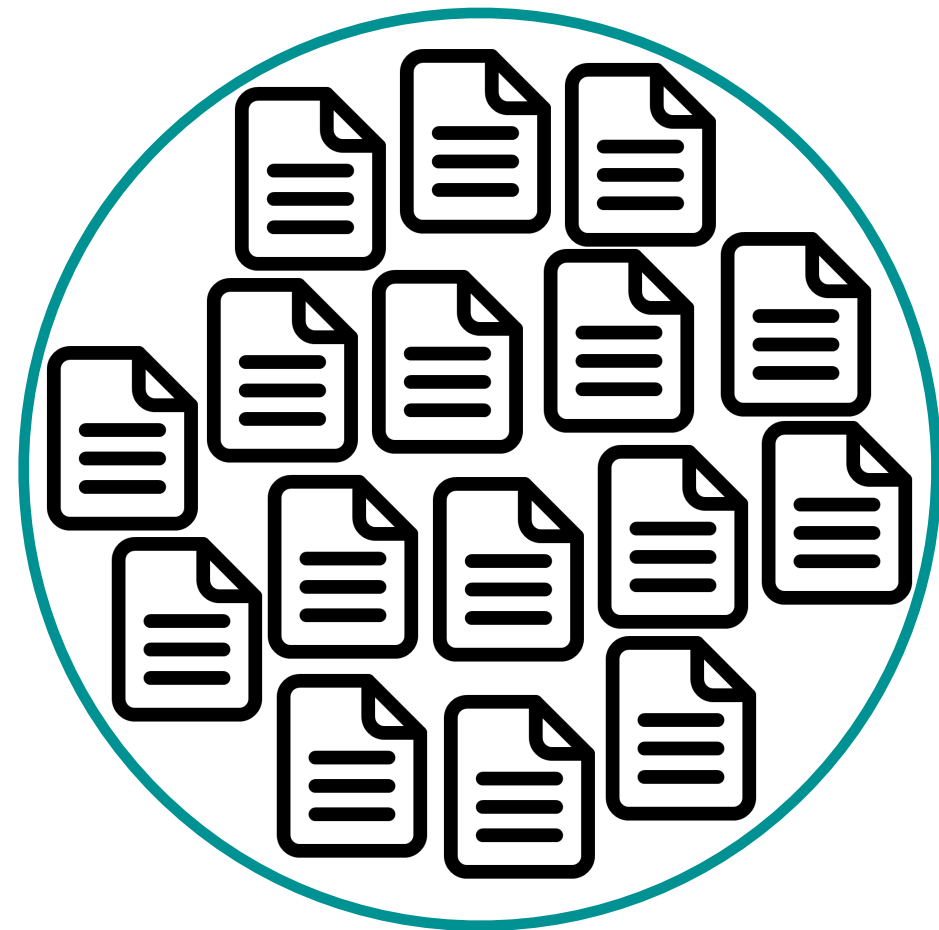
Modeling word frequency

Modeling word frequency

Simple model for word frequency: **count-based language model**

Modeling word frequency

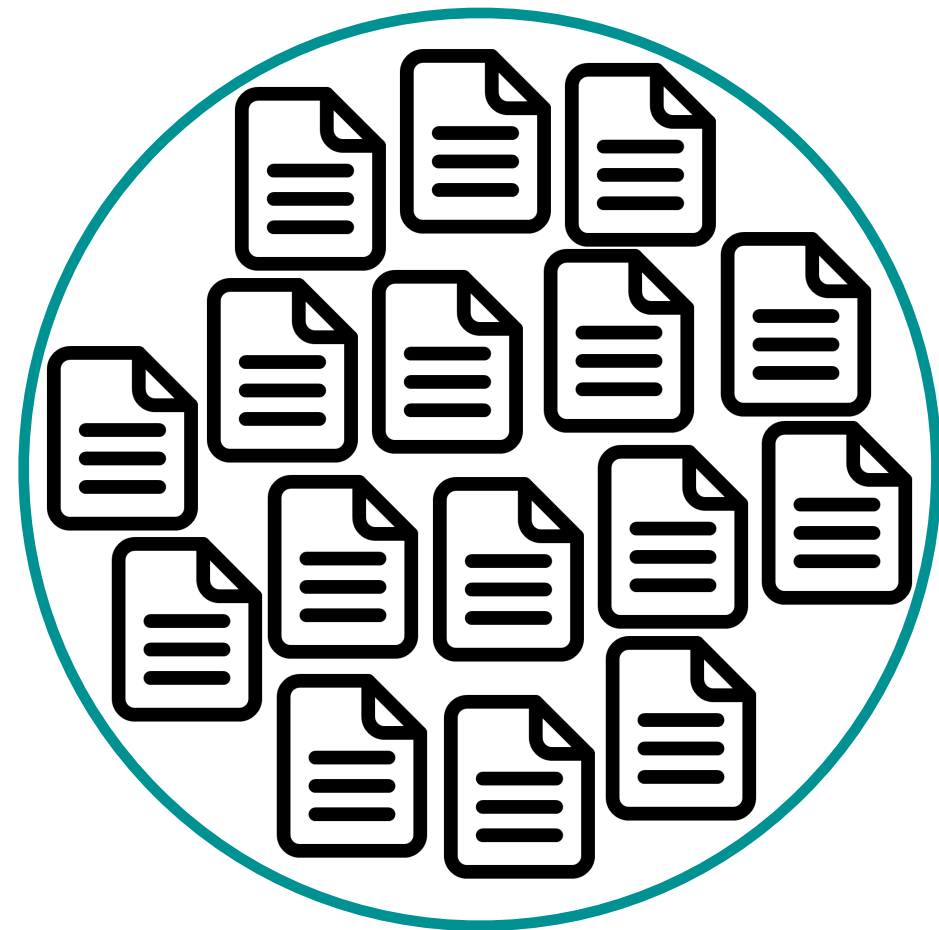
Simple model for word frequency: **count-based language model**



Predictions from self-training

Modeling word frequency

Simple model for word frequency: **count-based language model**



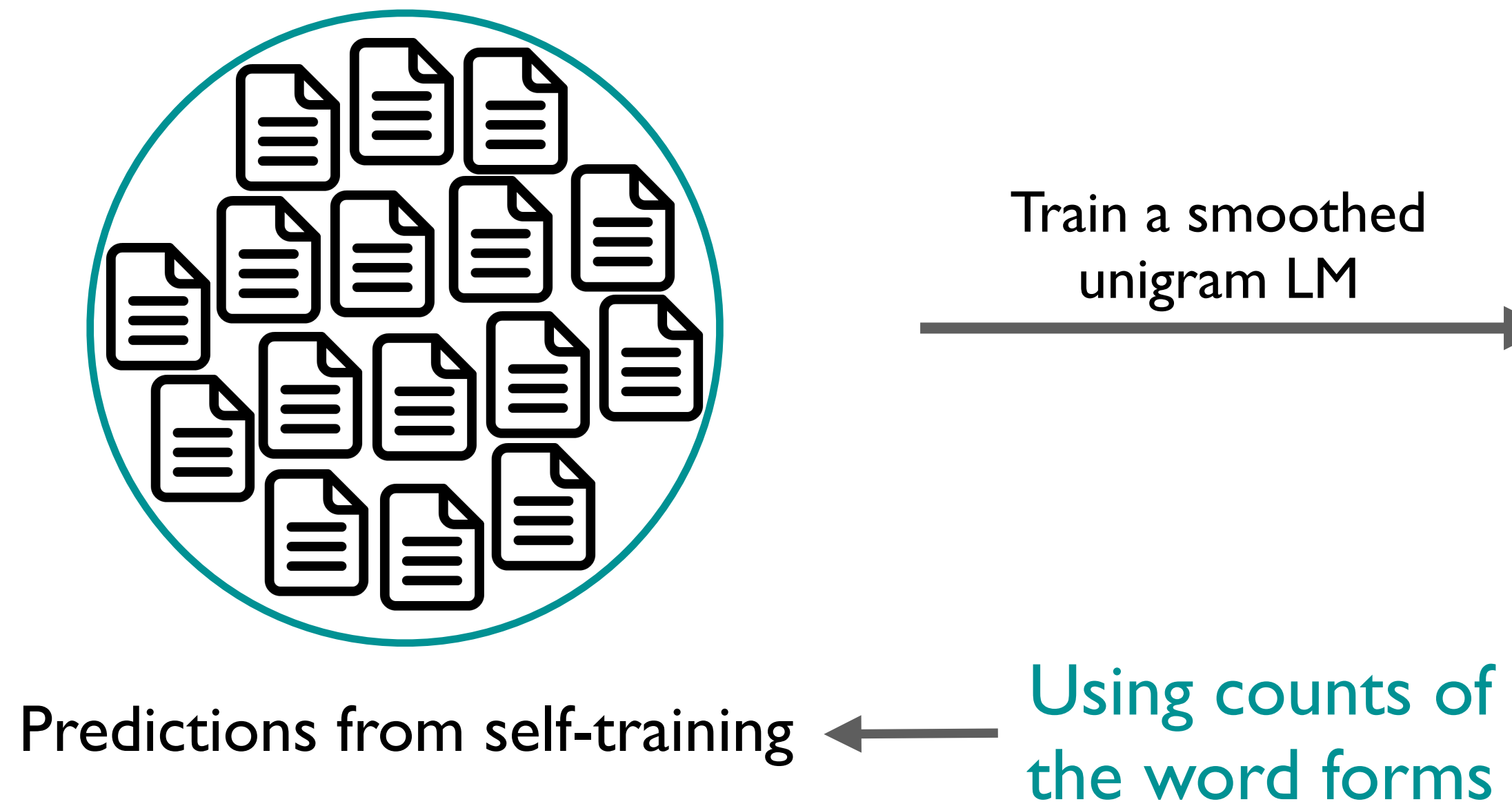
Predictions from self-training

Train a smoothed
unigram LM



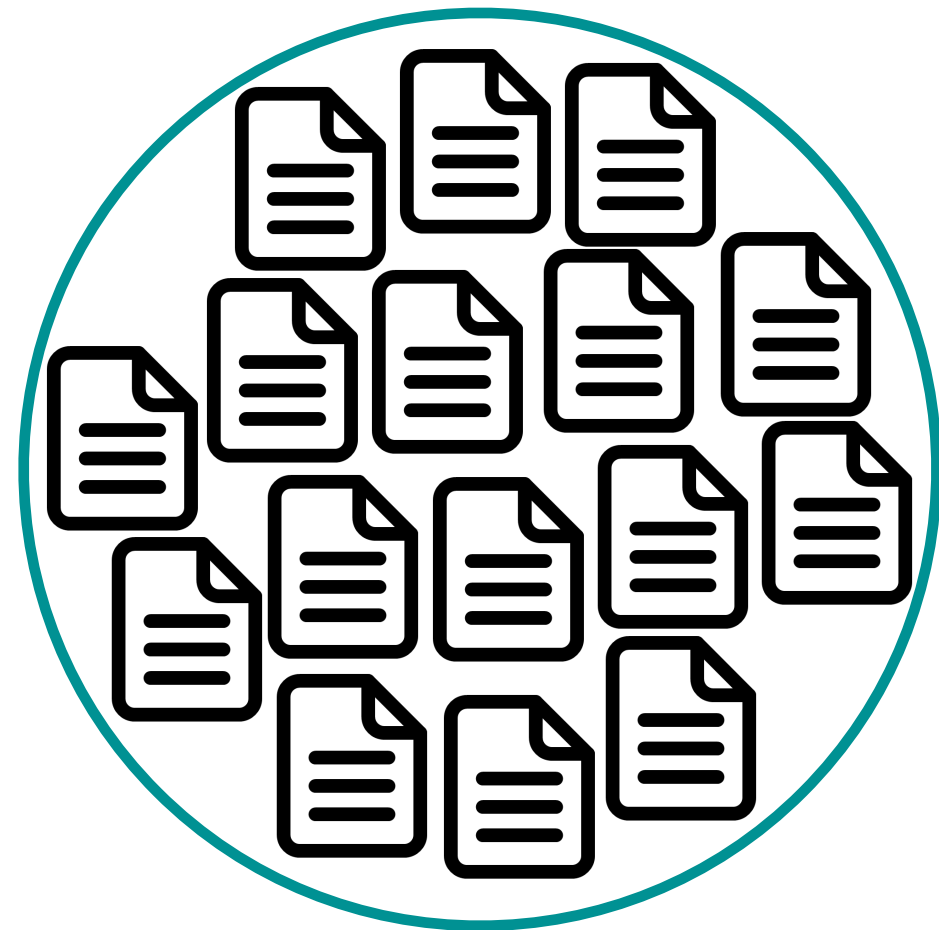
Modeling word frequency

Simple model for word frequency: **count-based language model**



Modeling word frequency

Simple model for word frequency: **count-based language model**



Predictions from self-training

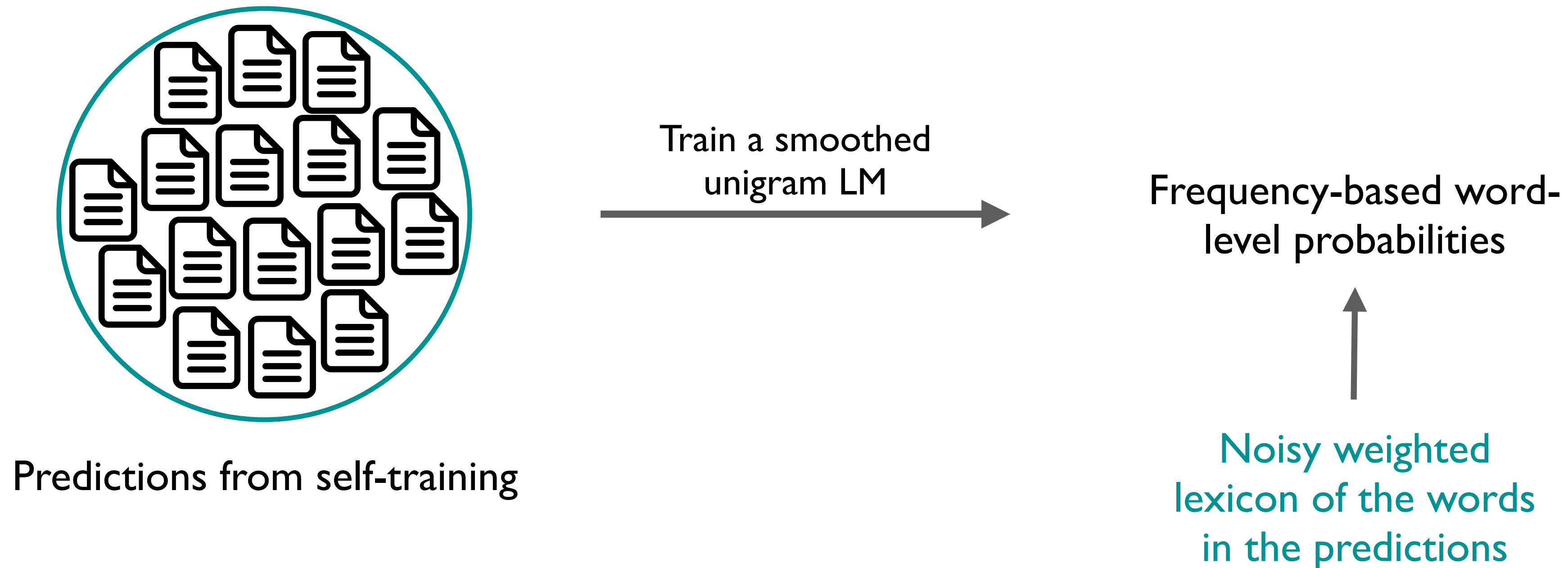
Train a smoothed
unigram LM



Frequency-based word-
level probabilities

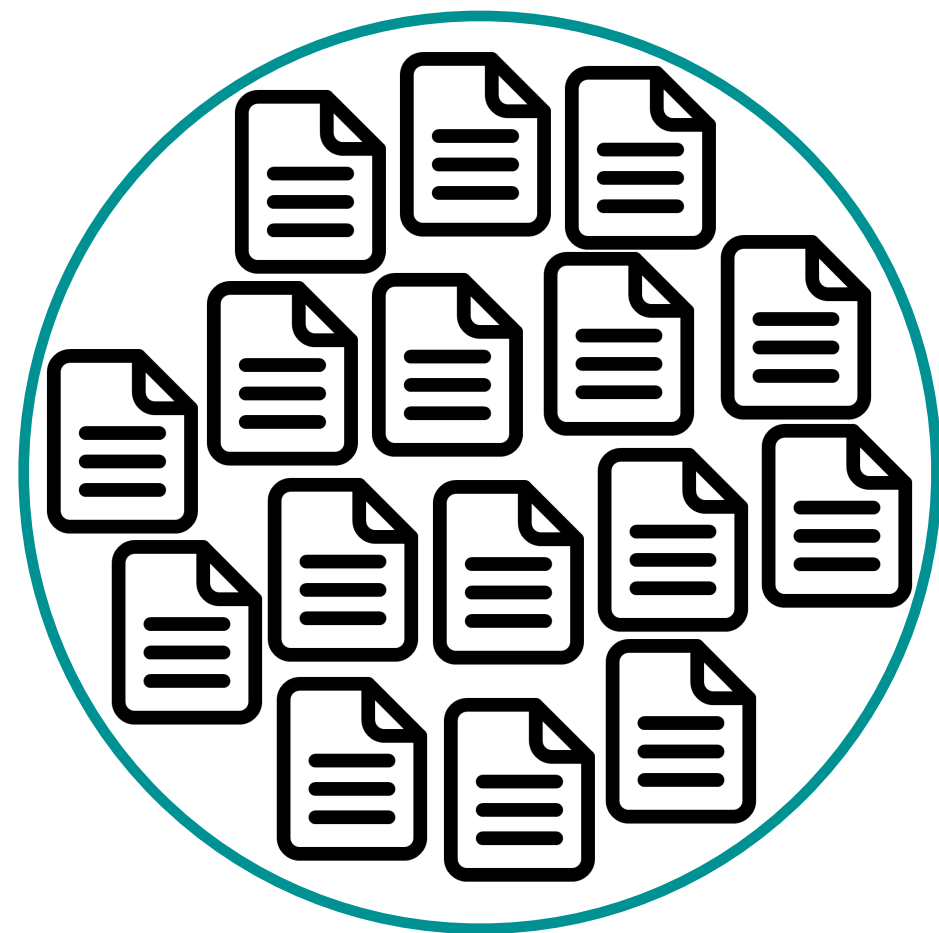
Modeling word frequency

Simple model for word frequency: **count-based language model**



Modeling word frequency

Simple model for word frequency: **count-based language model**



Predictions from self-training

Train a **smoothed**
unigram LM



Frequency-based word-
level probabilities

Lexically-aware decoding for post-correction

Lexically-aware decoding for post-correction

$$P(y) = p_{\text{lstm}}(y) \quad p_{\text{freq}}$$

Lexically-aware decoding for post-correction

$$P(y) = p_{\text{lstm}}(y)$$

p_{freq}



We have frequency-
based probabilities
from the unigram LM!

Lexically-aware decoding for post-correction

$$P(y) = p_{\text{lstm}}(y)$$

p_{freq}



But these are at the word-level: how we get character-level scores?

Scoring at the character-level

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM



Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Set of states with
weighted transitions

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

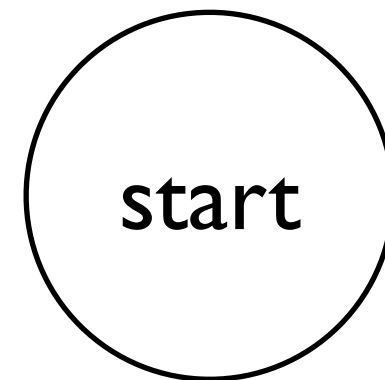
Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM



Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$



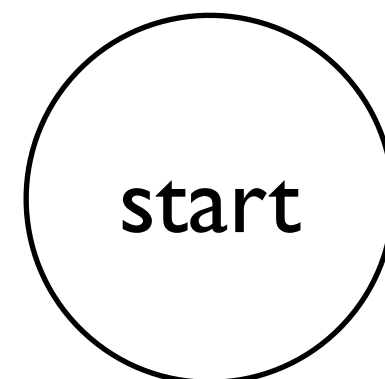
Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM



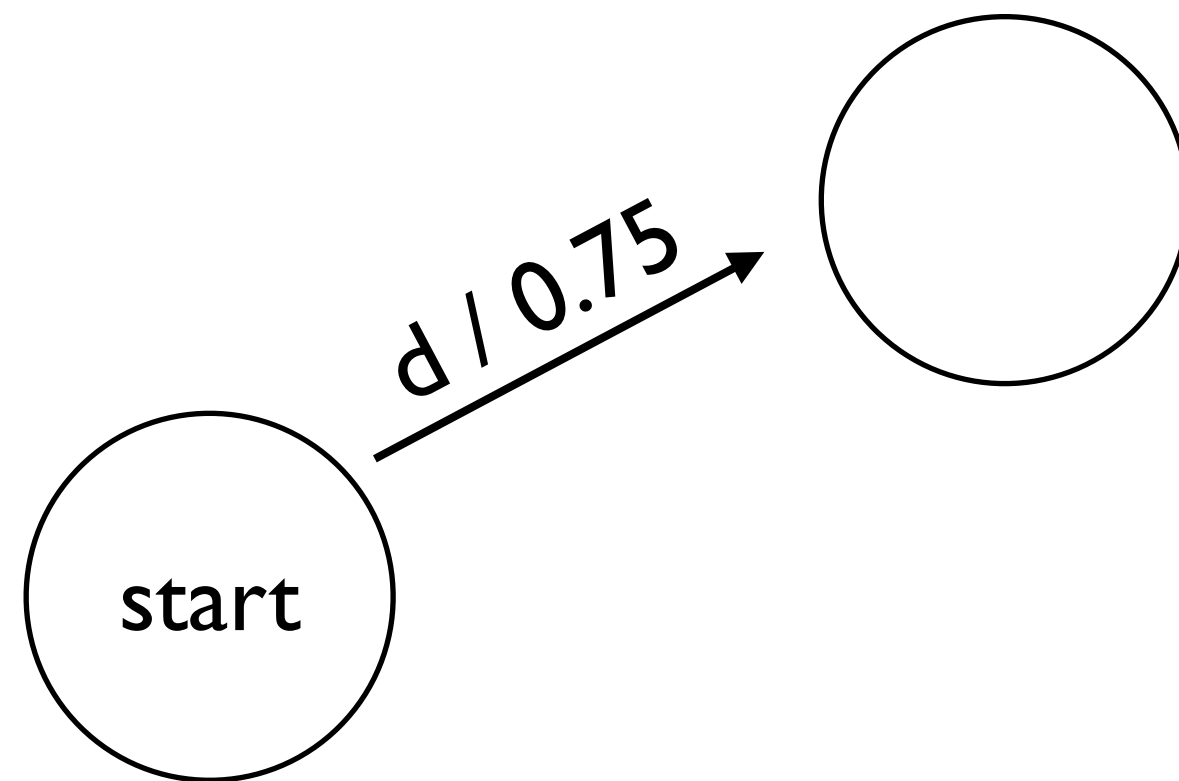
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$



Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM



Example LM with two words:

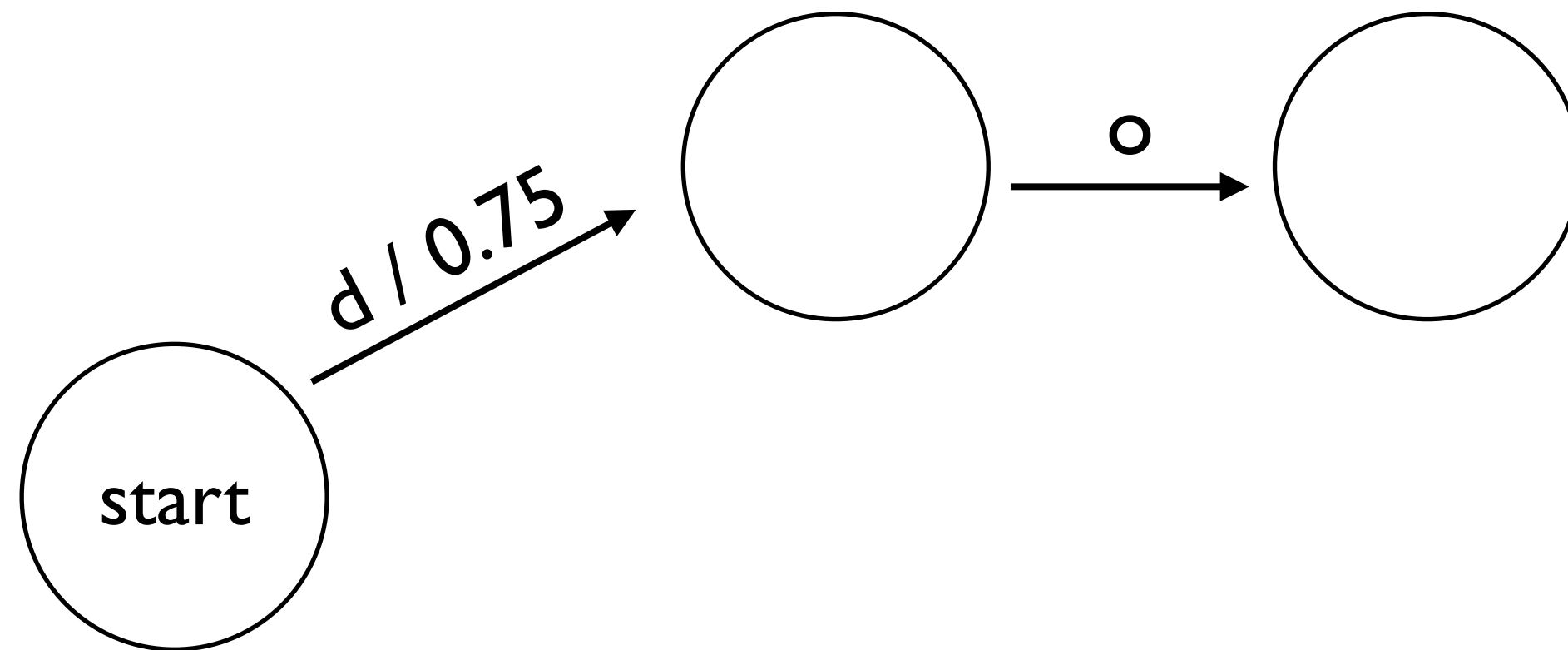
- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

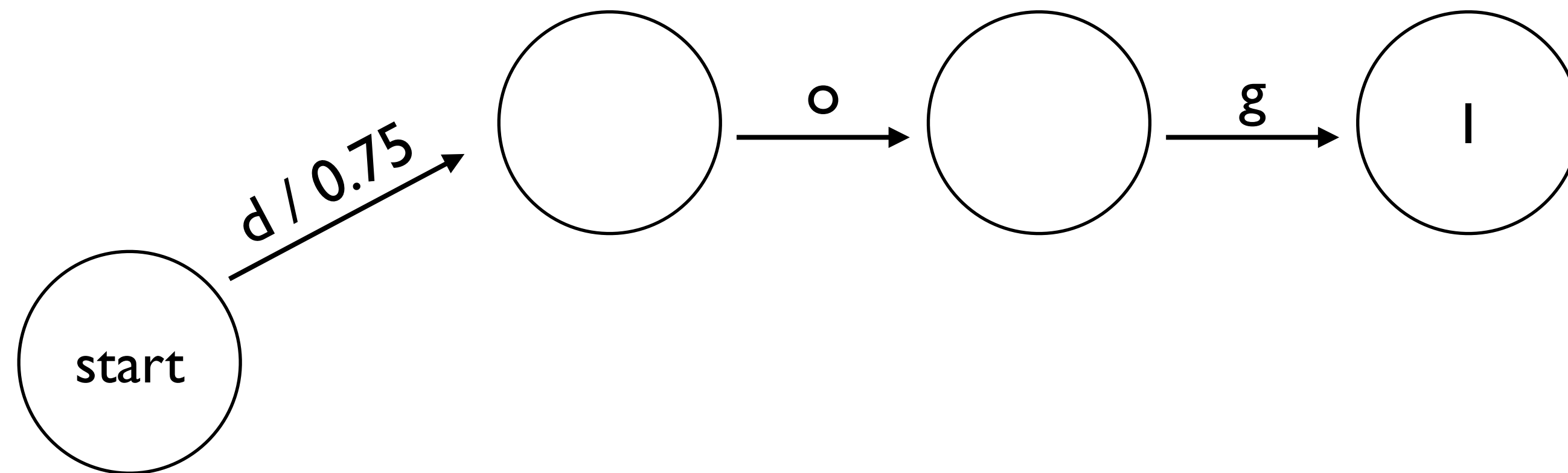


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

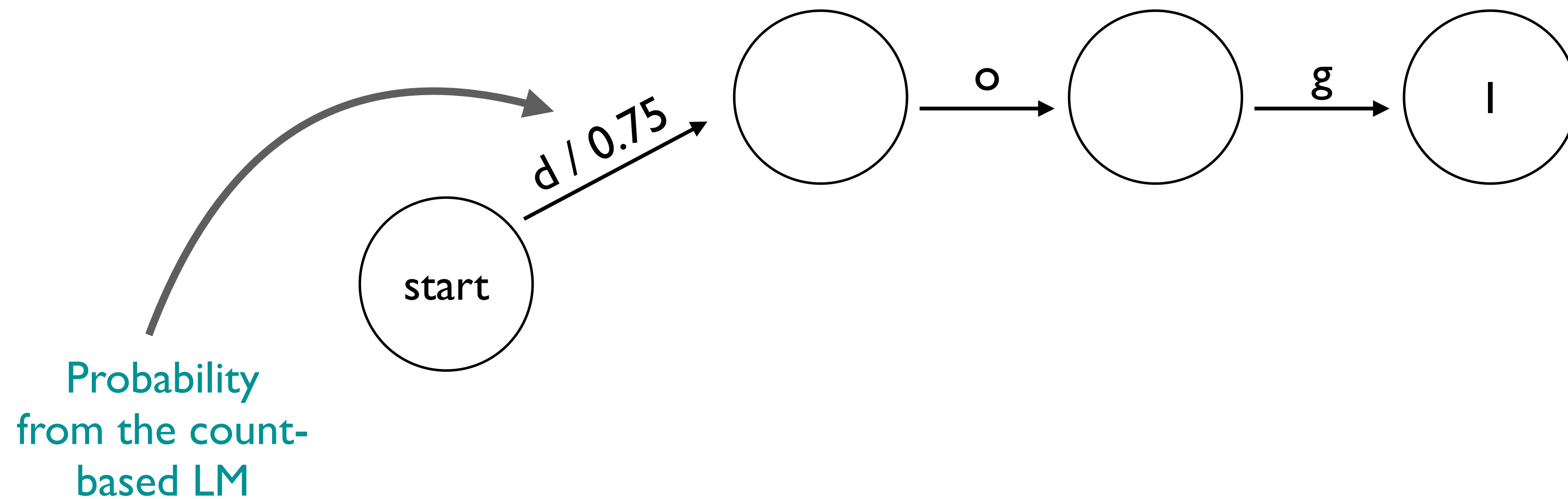


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

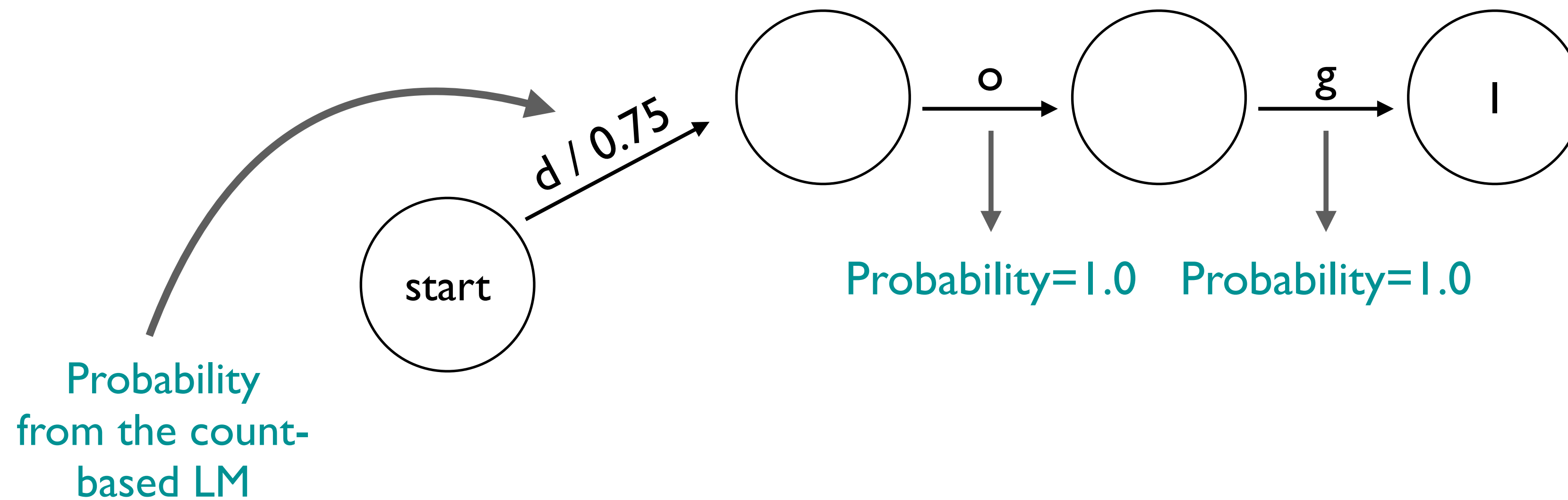


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

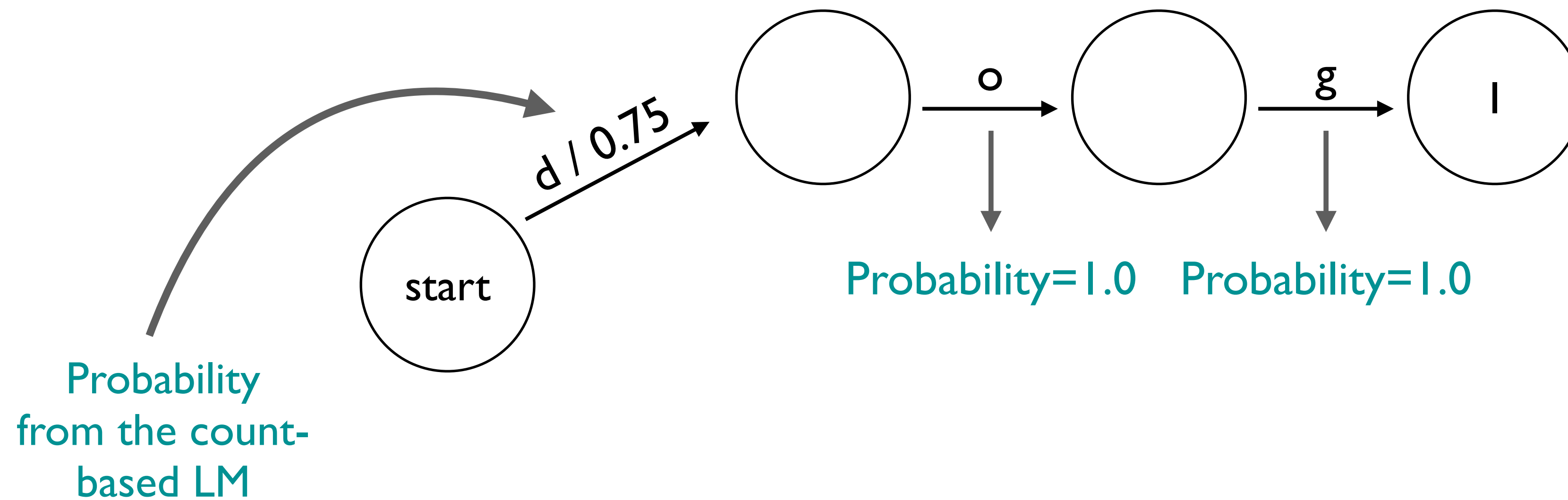


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$



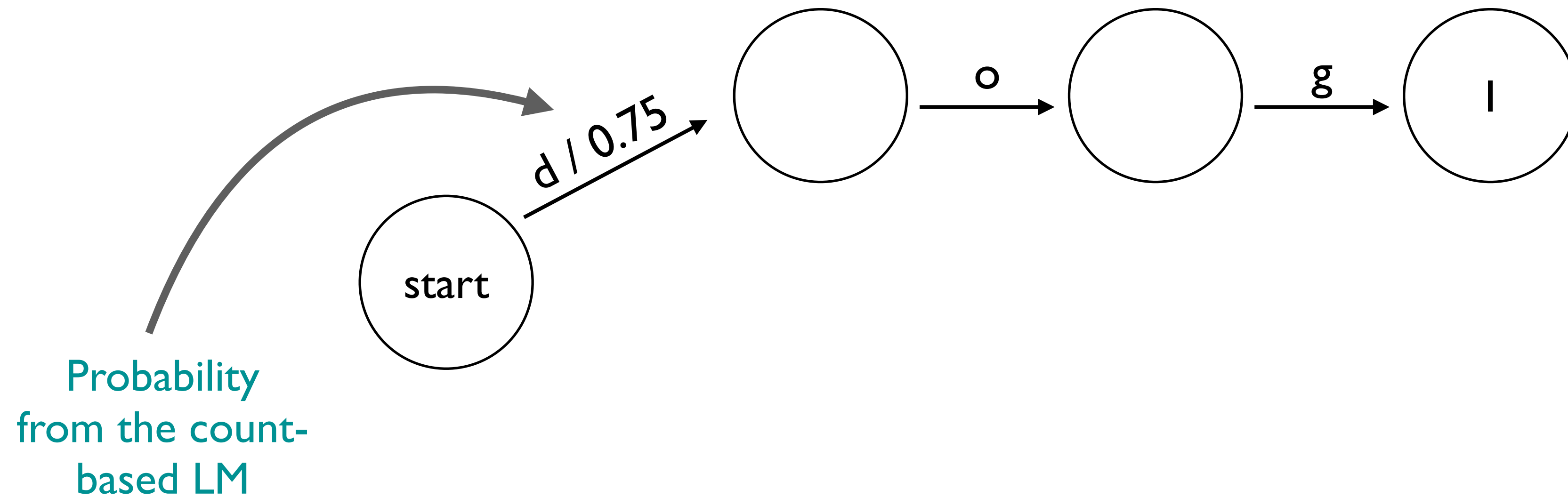
Weight of path for "dog" = 0.75
Same as the word-level LM!

Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

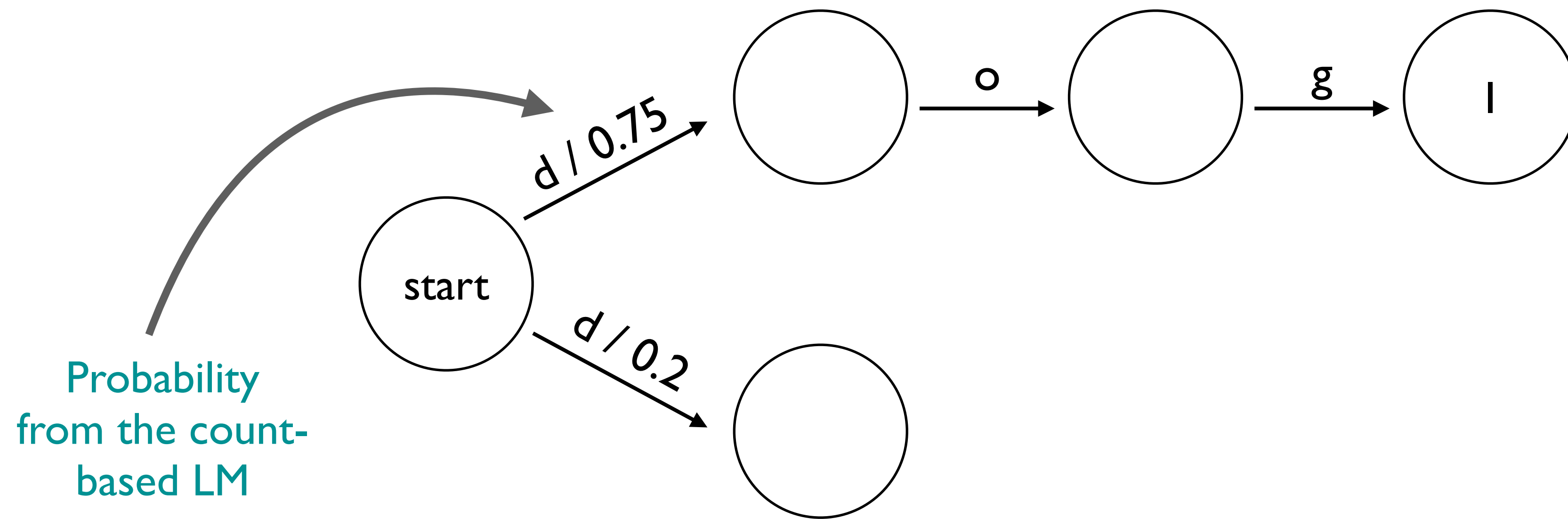


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

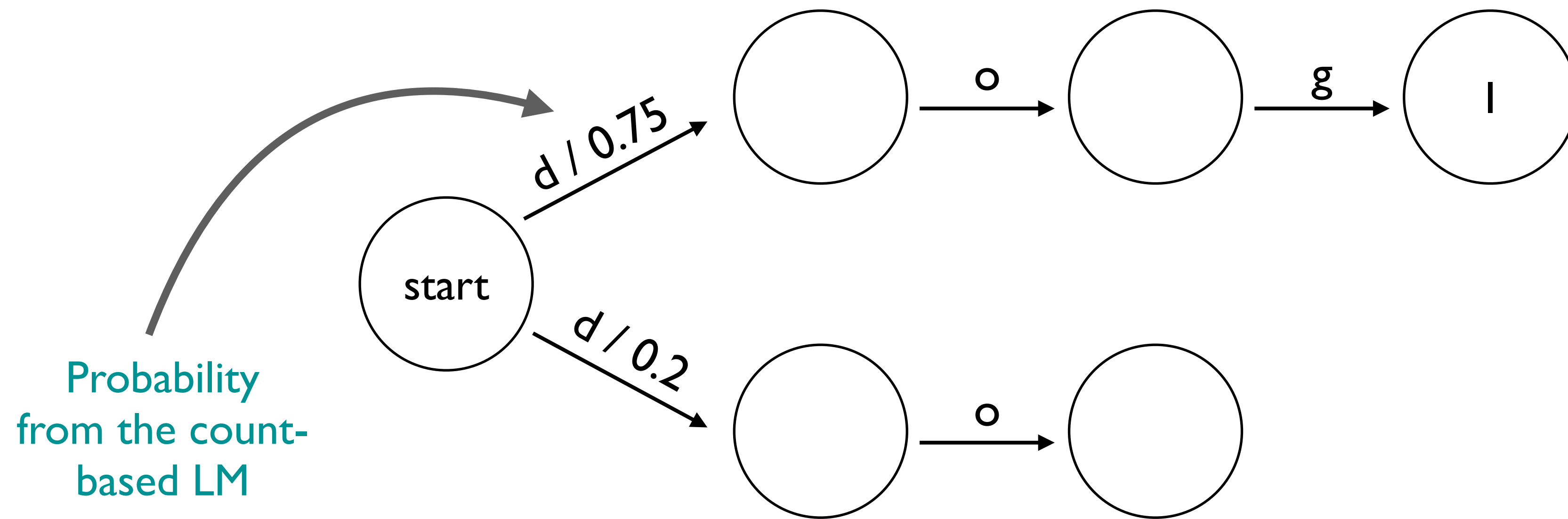


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

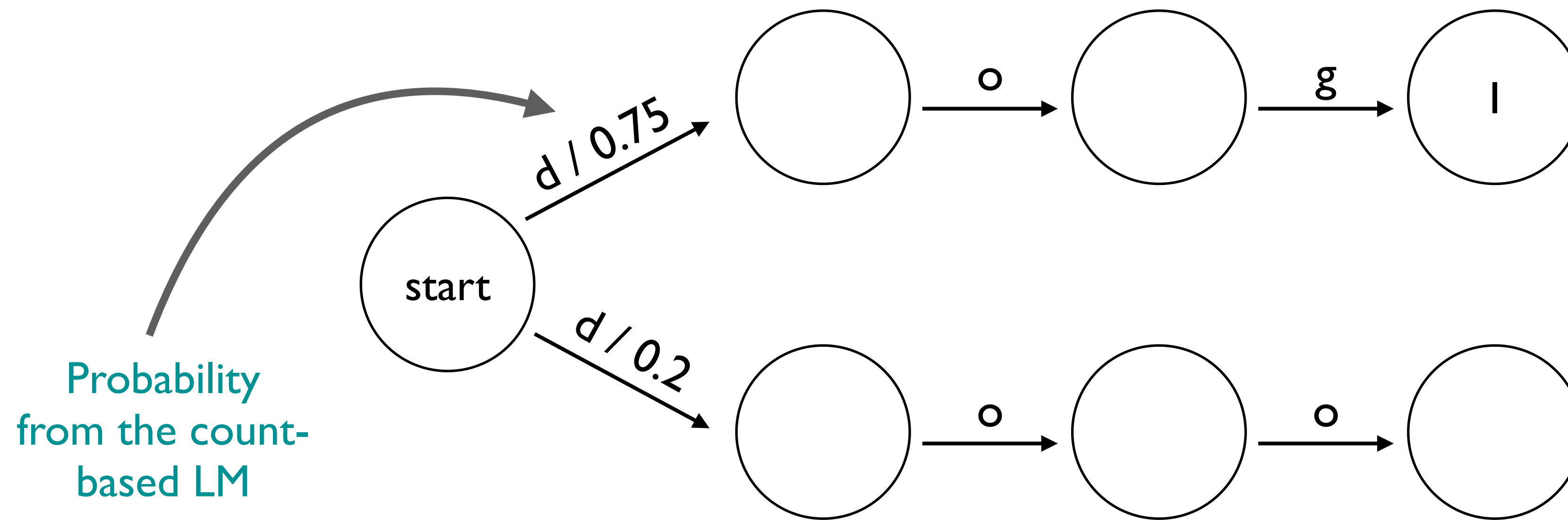


Scoring at the character-level

Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

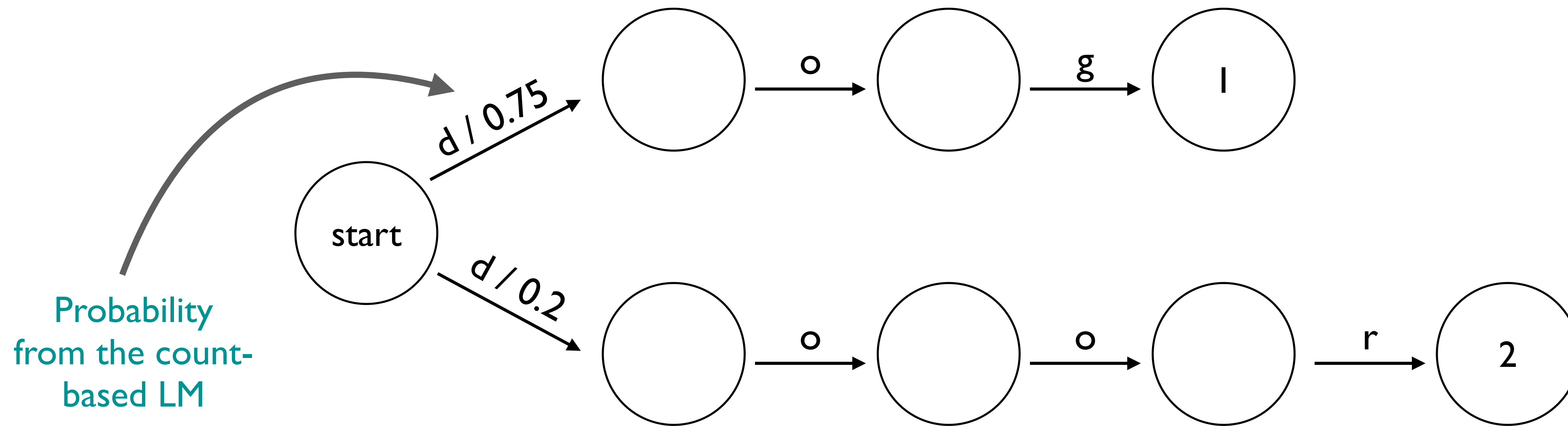


Scoring at the character-level

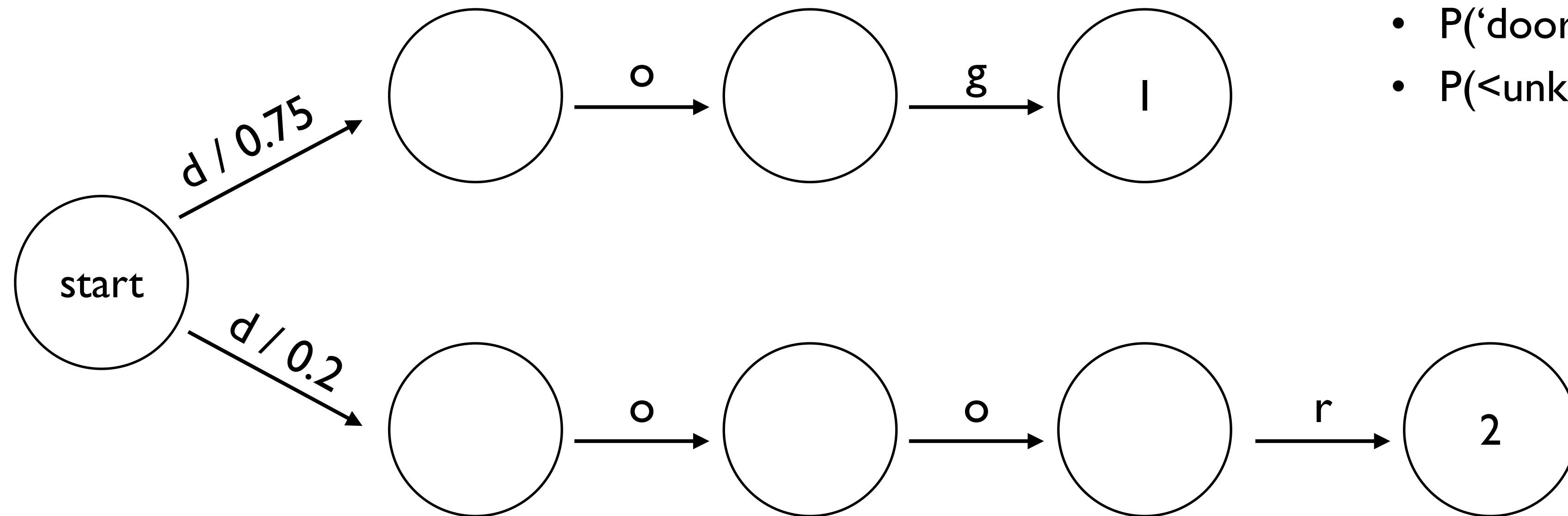
Weighted Finite State Automaton (WFSA)
representation of the LM

Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$



Scoring at the character-level



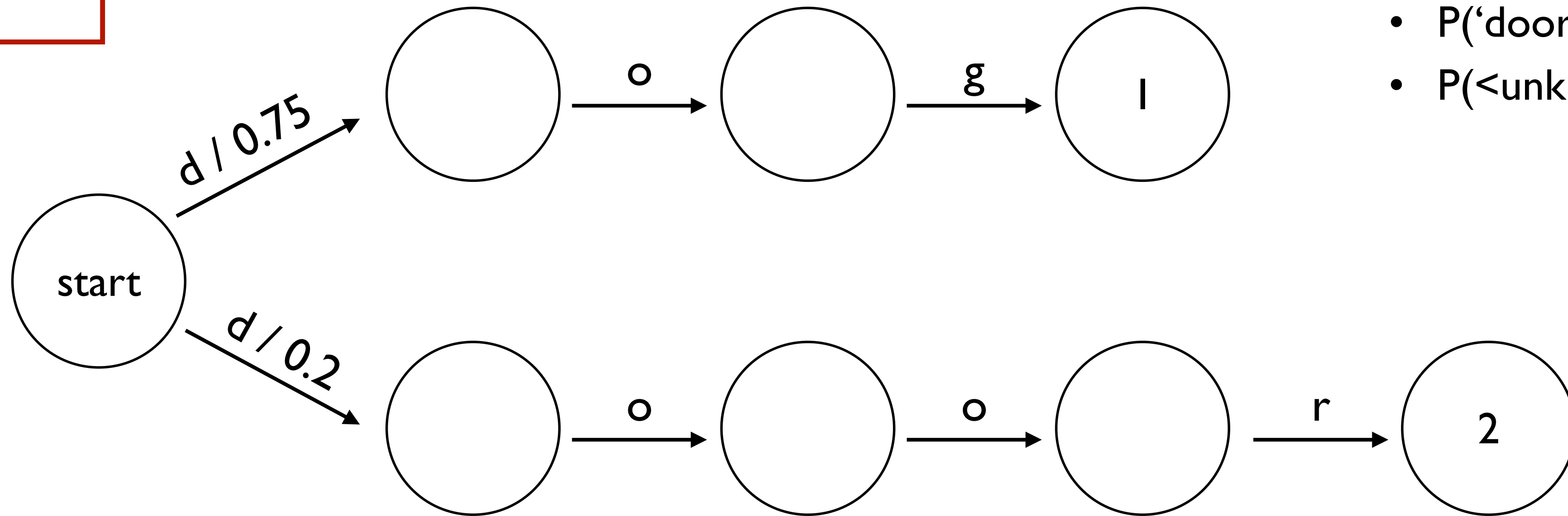
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

d o o r



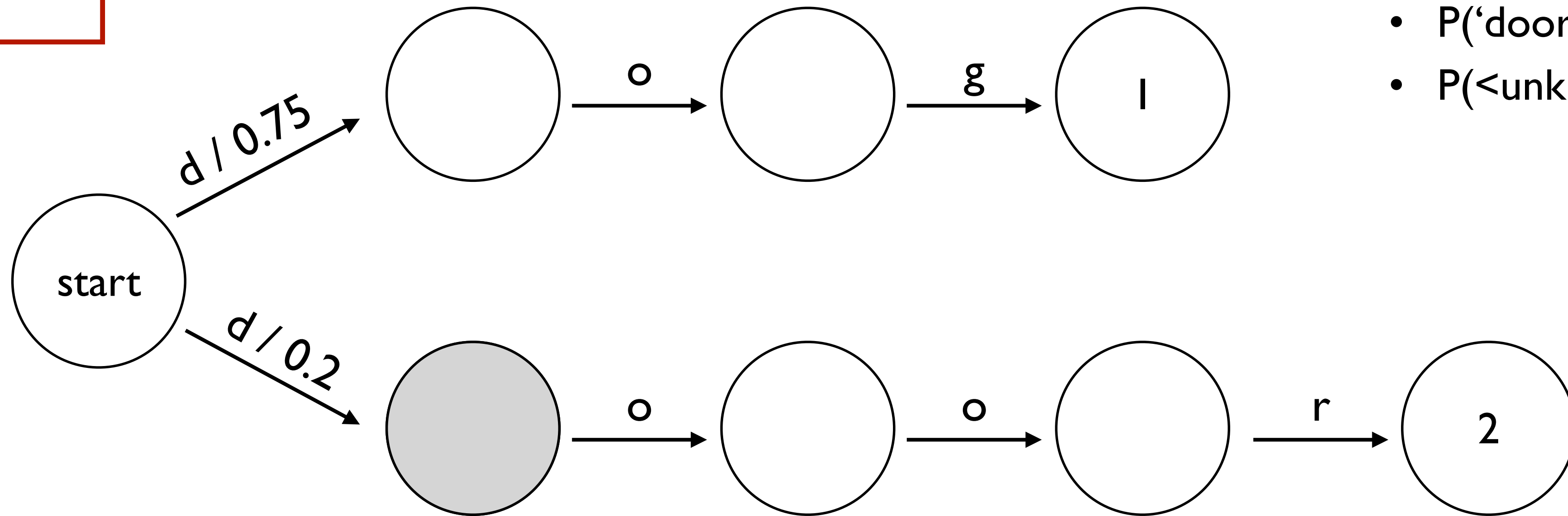
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

d o o r



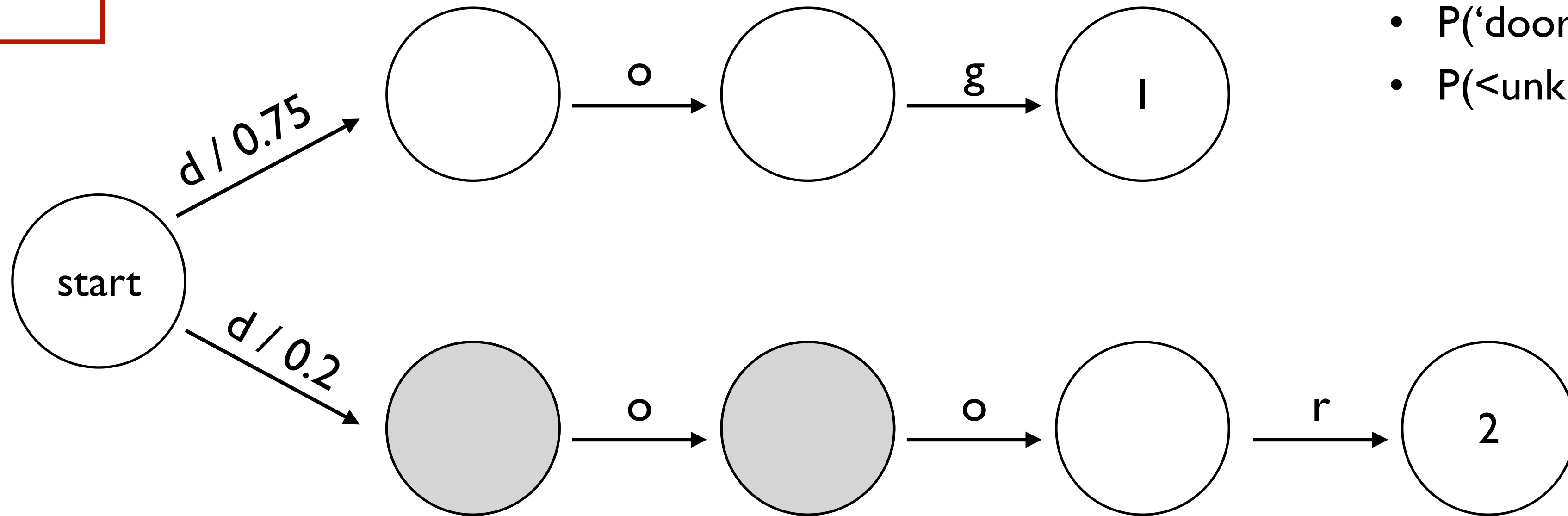
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

d o o r



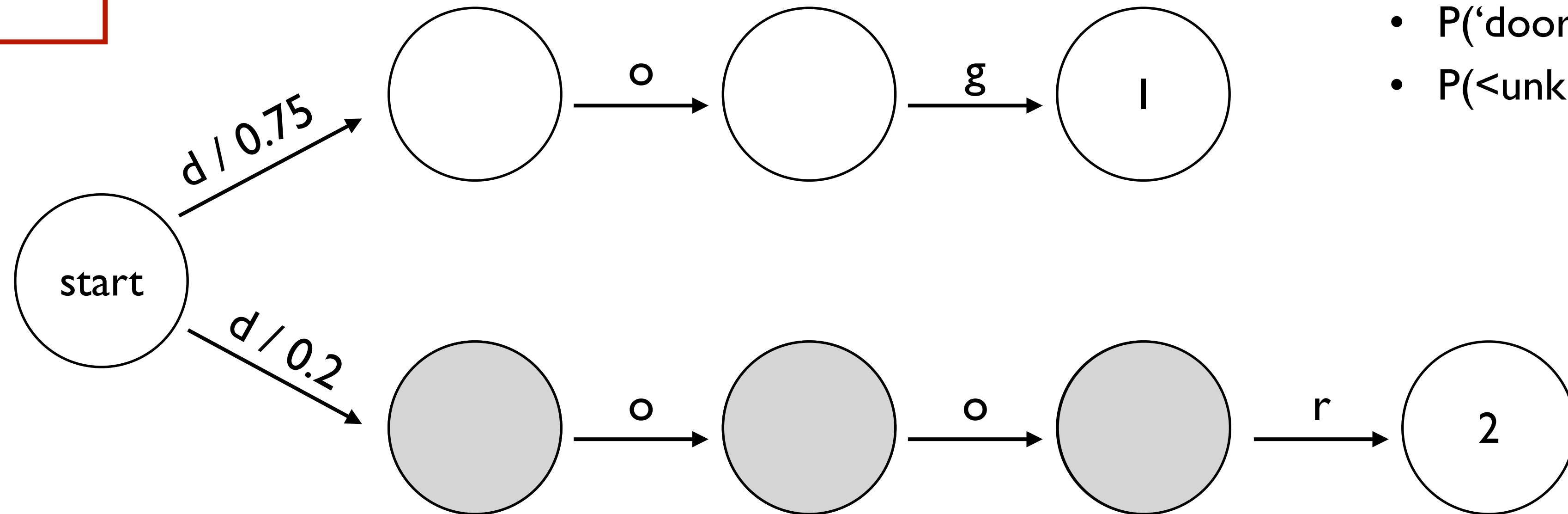
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

d o o r



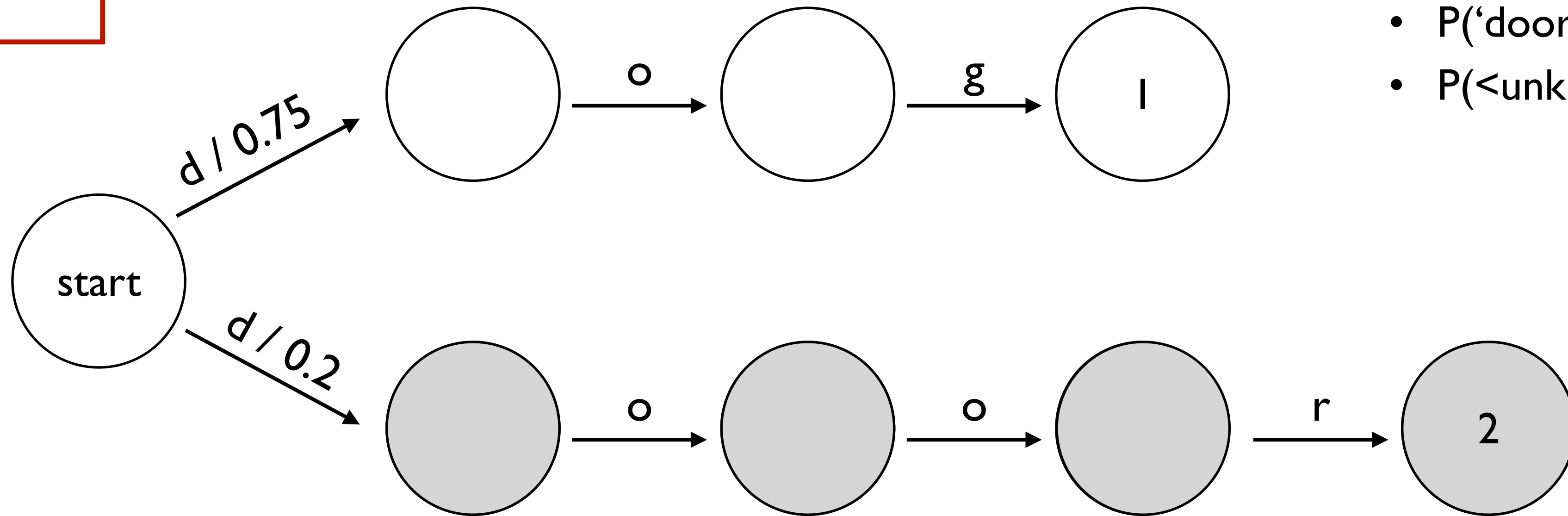
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

d o o r



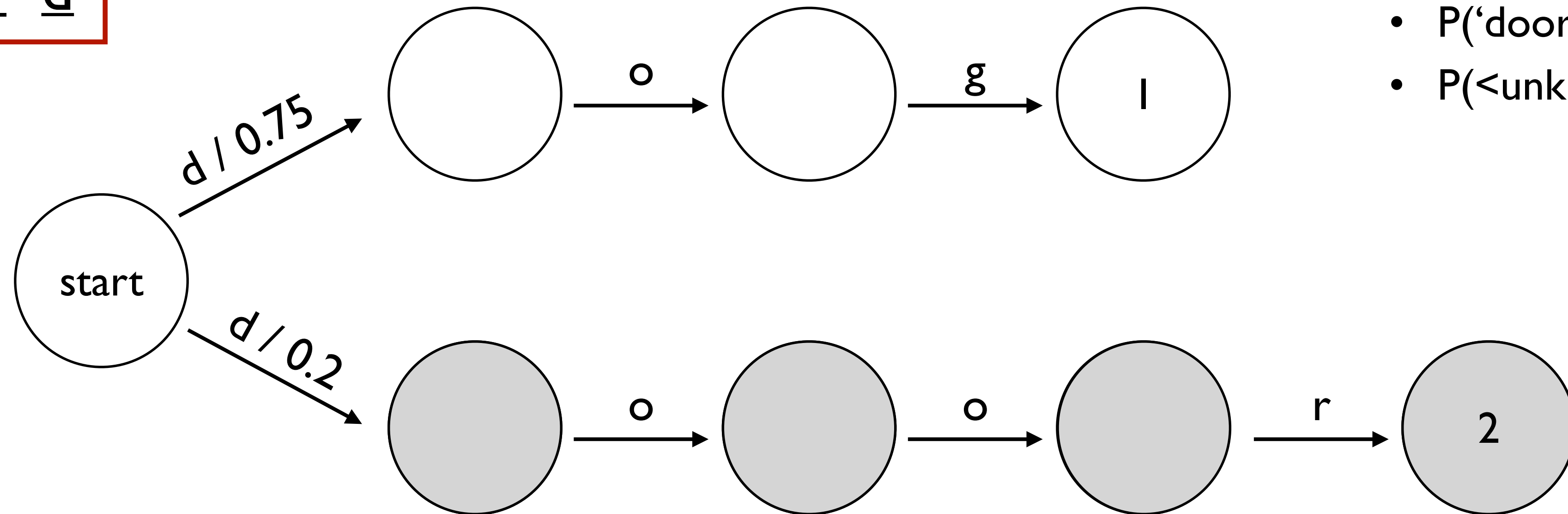
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Output:

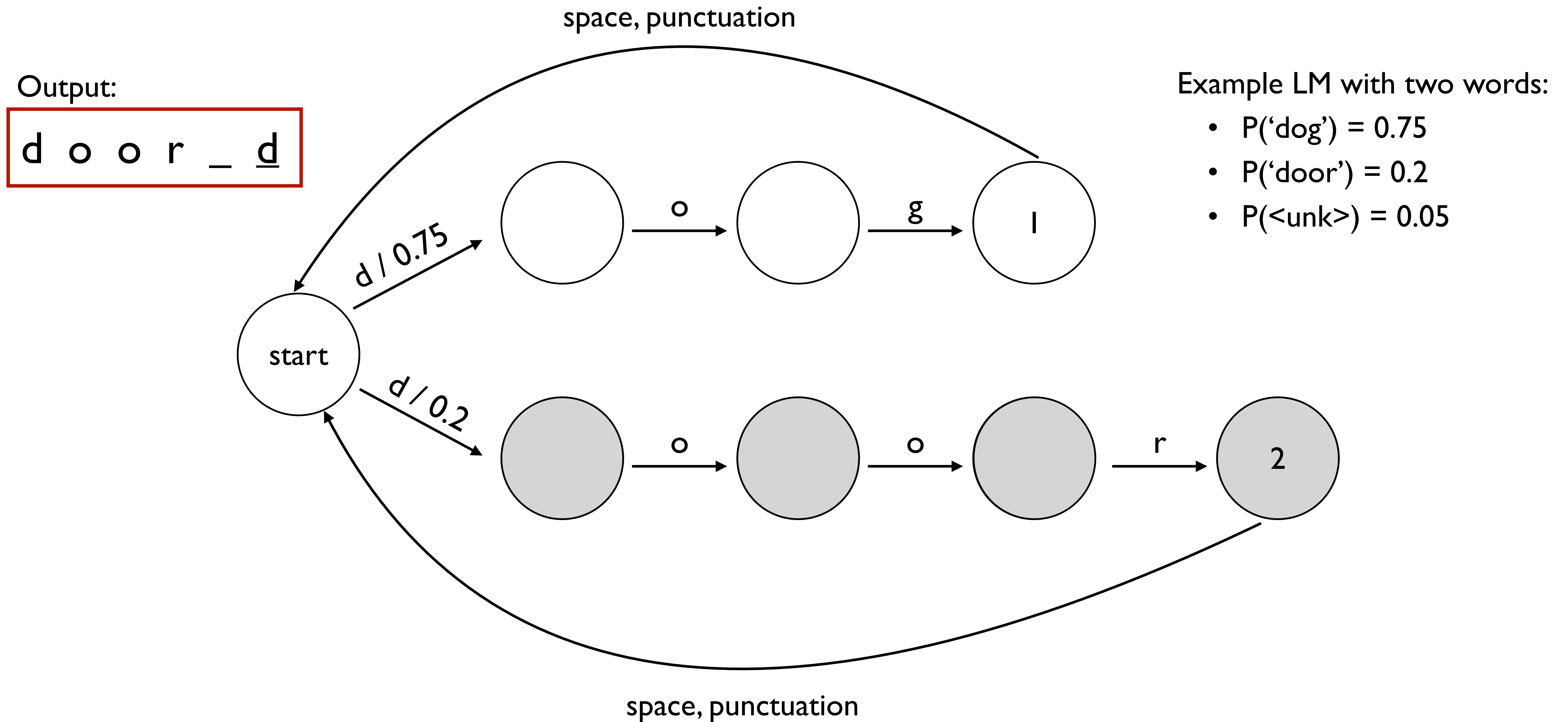
d o o r _ d



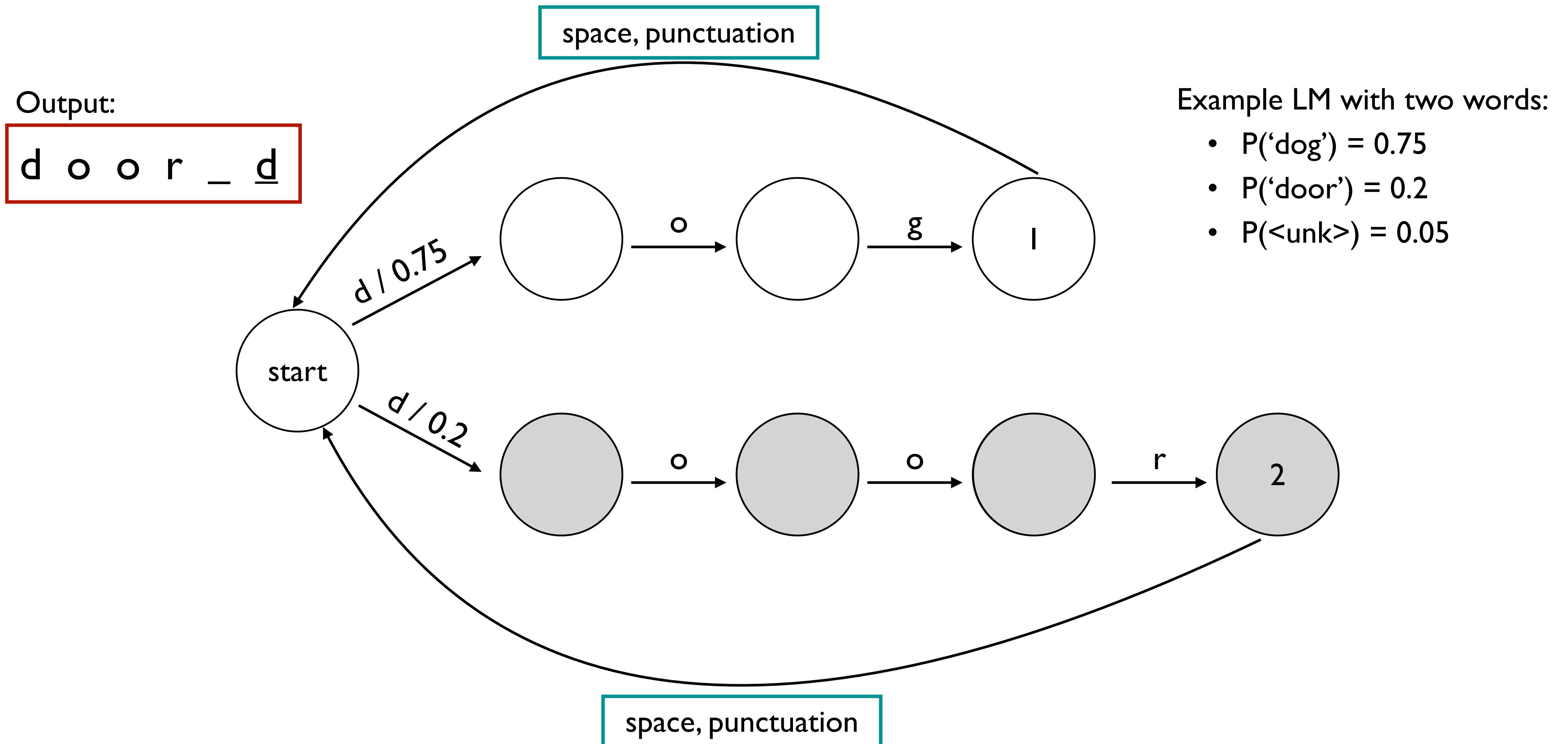
Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

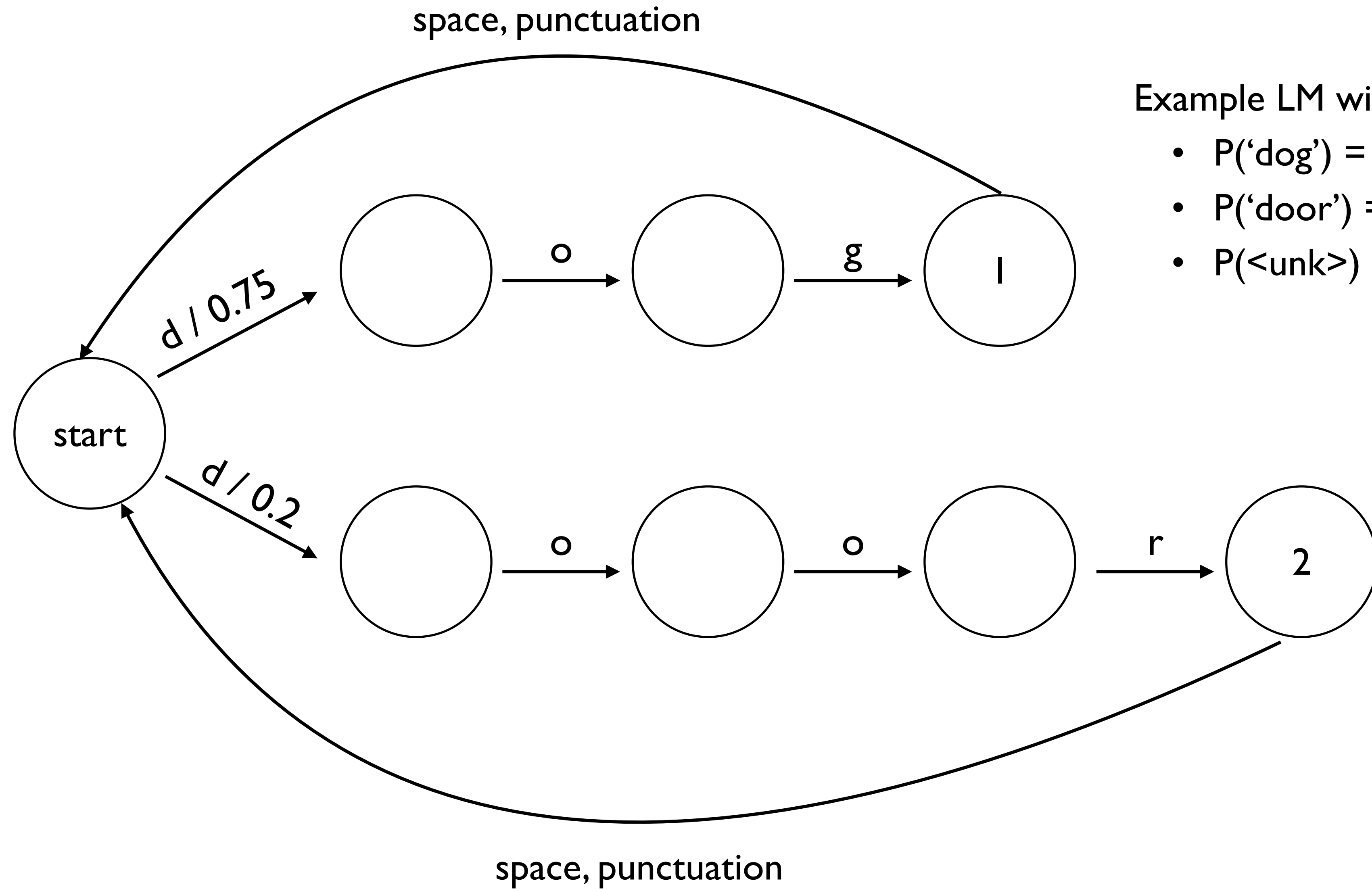
Scoring at the character-level



Scoring at the character-level



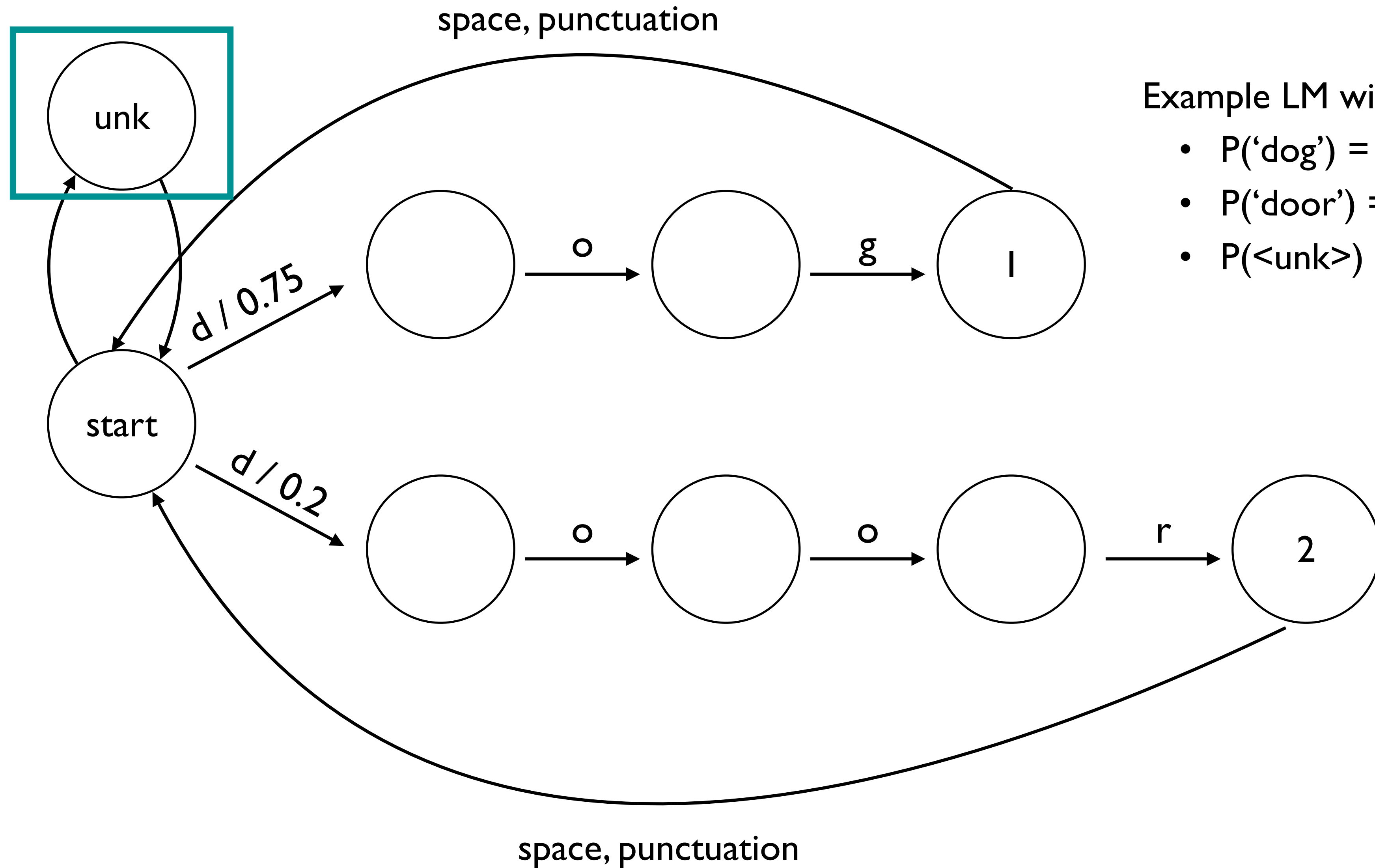
Scoring at the character-level



Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

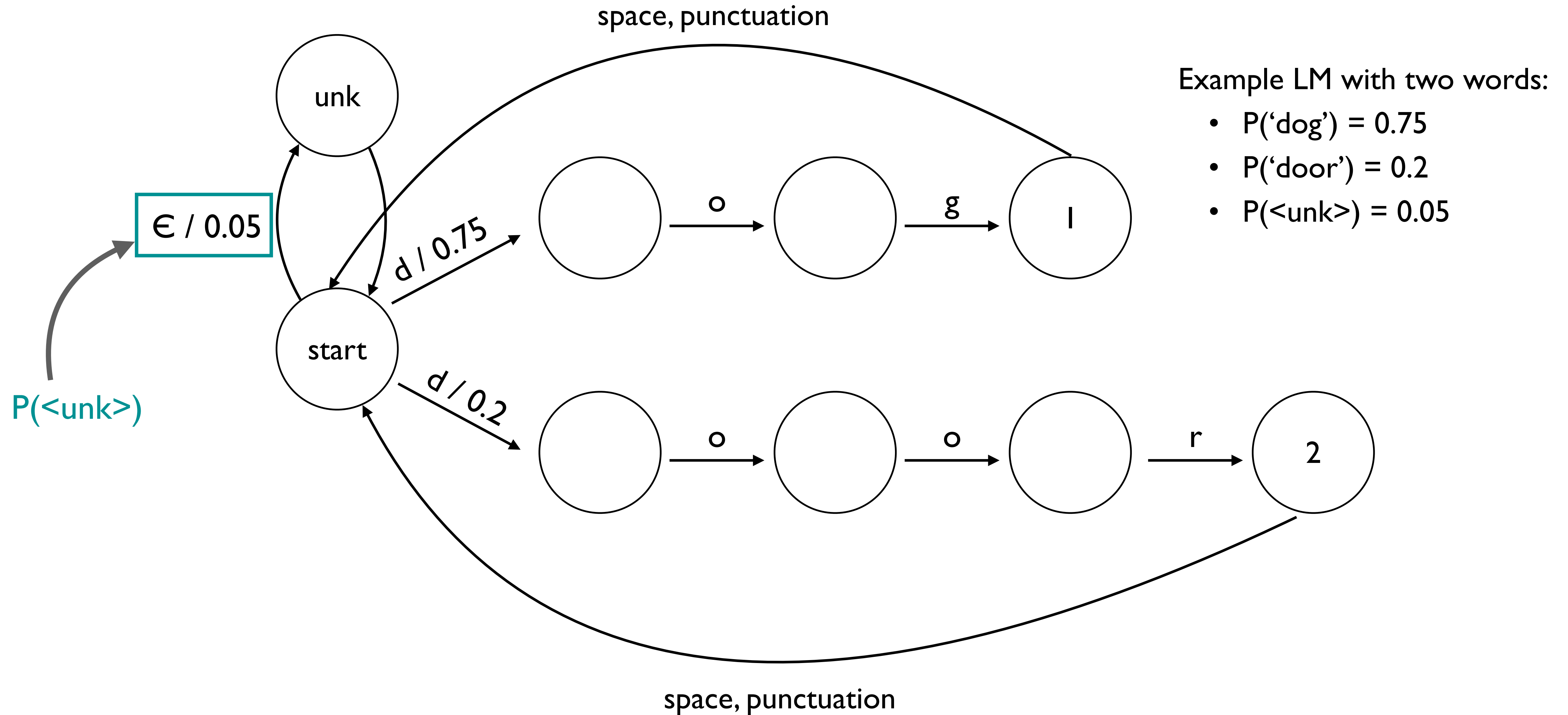
Scoring at the character-level



Example LM with two words:

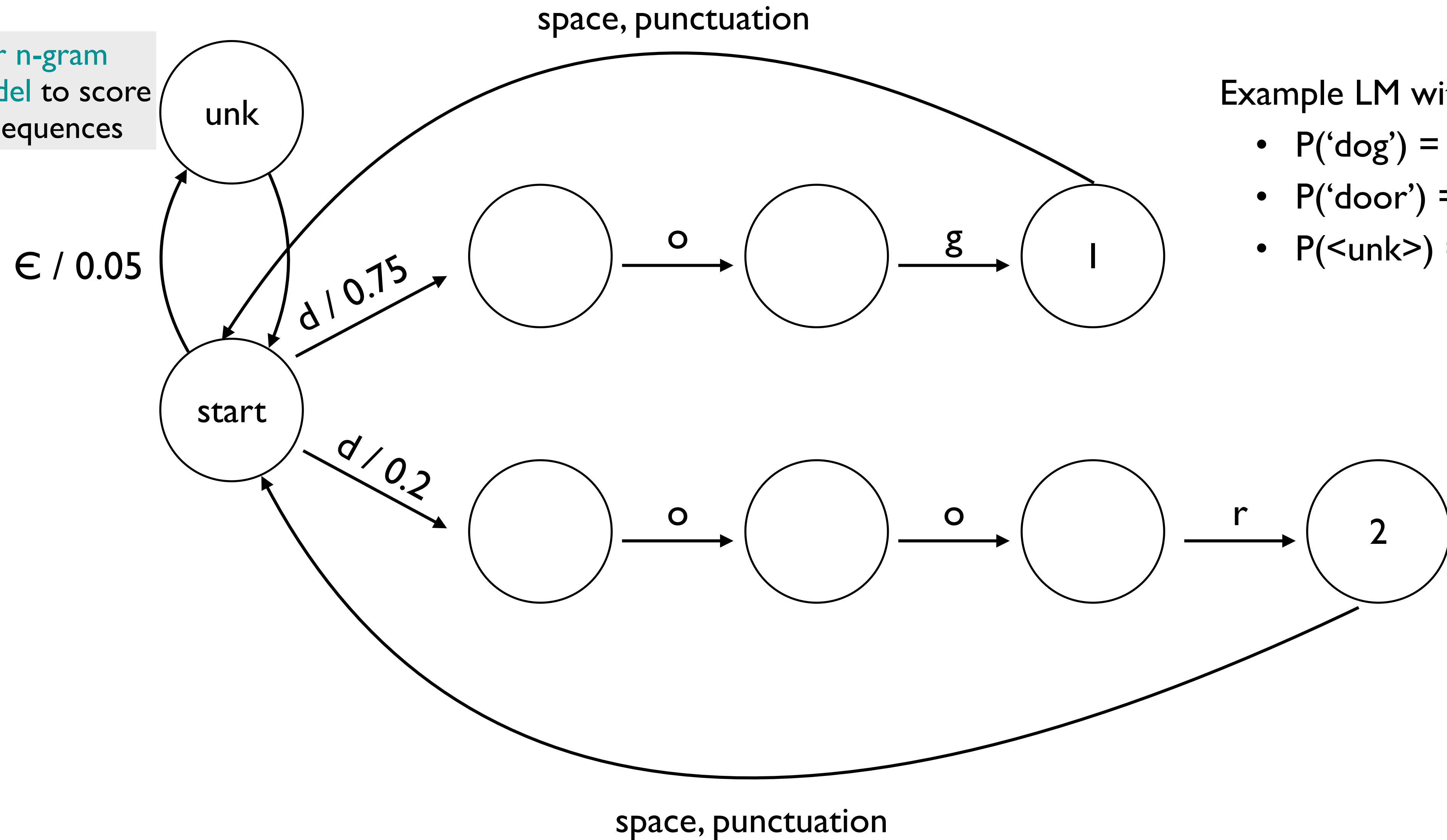
- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level



Scoring at the character-level

Character n-gram
language model to score
unknown sequences

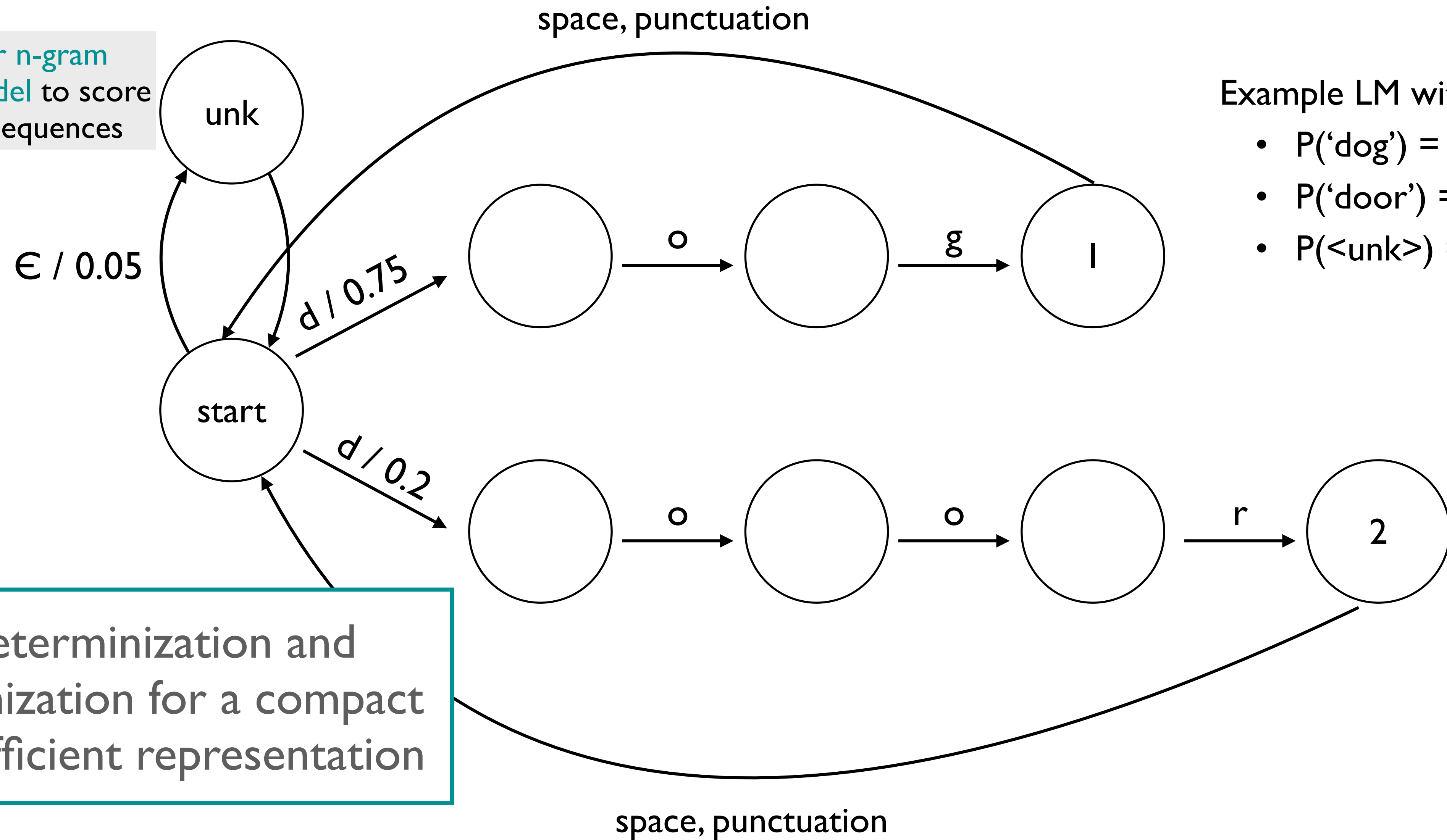


Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\text{'<unk>'}) = 0.05$

Scoring at the character-level

Character n-gram
language model to score
unknown sequences



Example LM with two words:

- $P(\text{'dog'}) = 0.75$
- $P(\text{'door'}) = 0.2$
- $P(\langle \text{unk} \rangle) = 0.05$

Determinization and
minimization for a compact
and efficient representation

space, punctuation

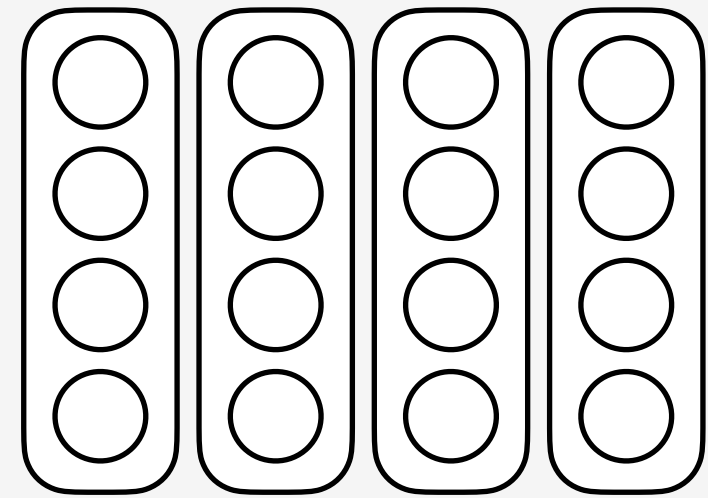
Lexically-aware decoding for post-correction

Lexically-aware decoding for post-correction

$$P(y) = p_{\text{lstm}}(y) \quad p_{\text{freq}}(y)$$

Lexically-aware decoding for post-correction

$$P(y) = p_{\text{lstm}}(y) \quad p_{\text{freq}}(y)$$



Lexically-aware decoding for post-correction

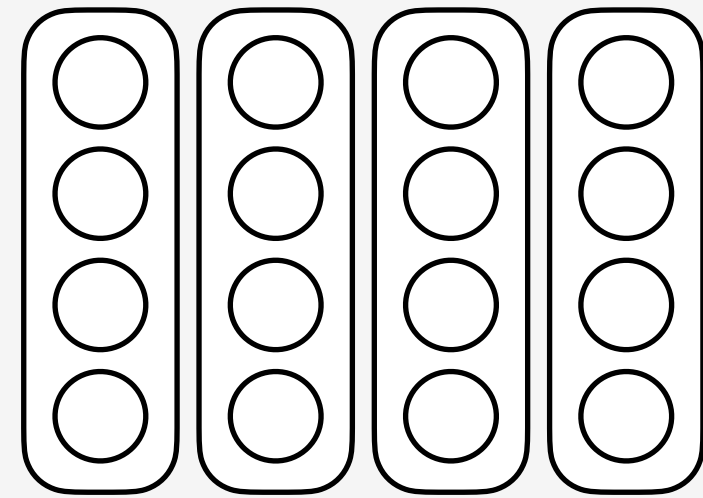
$$P(y) =$$


$$p_{\text{lstm}}(y)$$

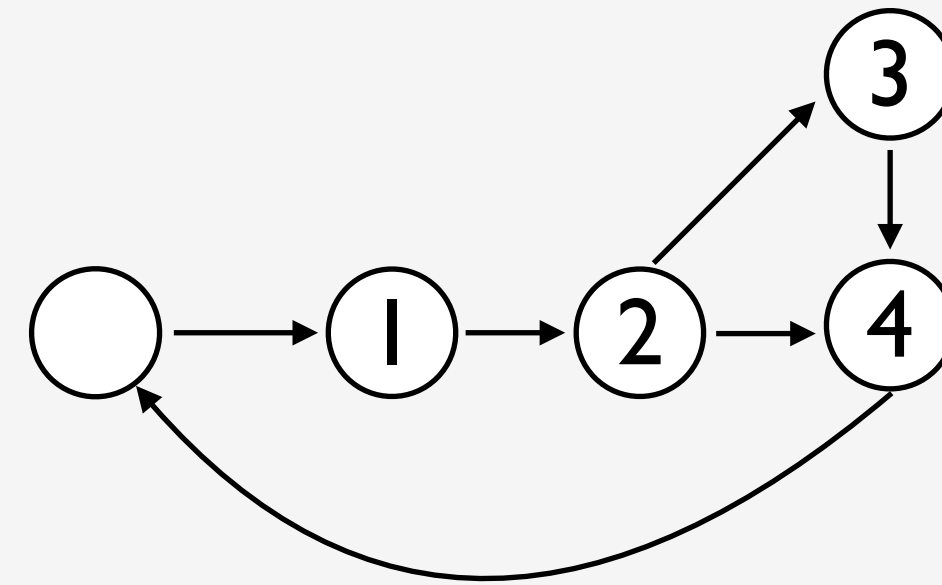
$$p_{\text{freq}}(y)$$

Lexically-aware decoding for post-correction

$$P(y) =$$



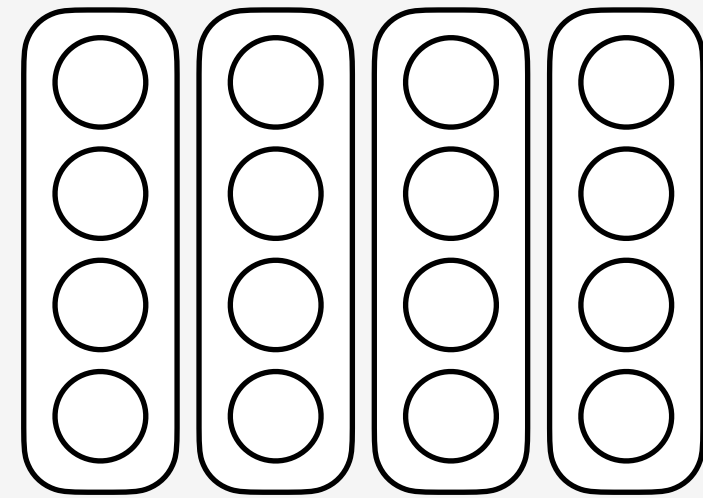
$$p_{\text{lstm}}(y)$$



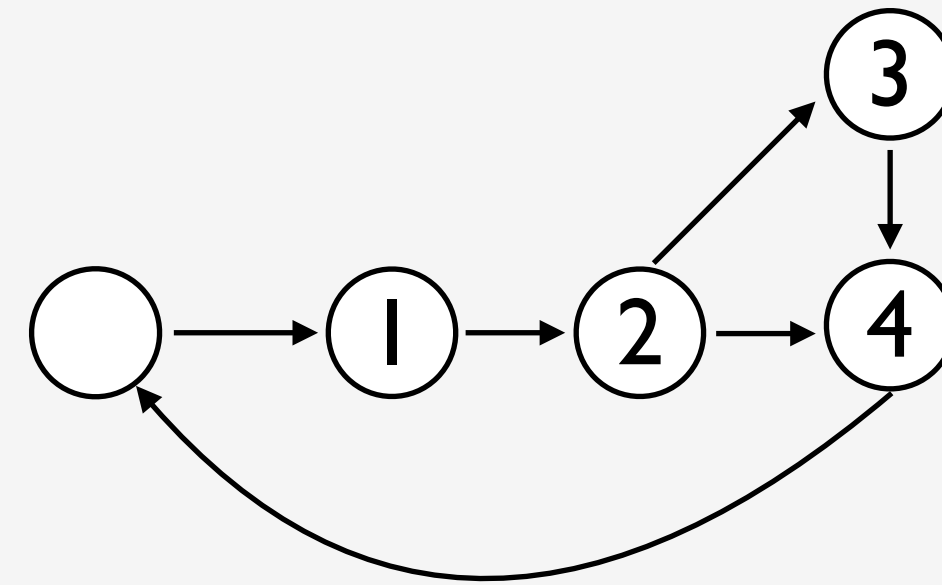
$$p_{\text{wfsa}}(y)$$

Lexically-aware decoding for post-correction

$$P(y) =$$



$$p_{\text{lstm}}(y)$$

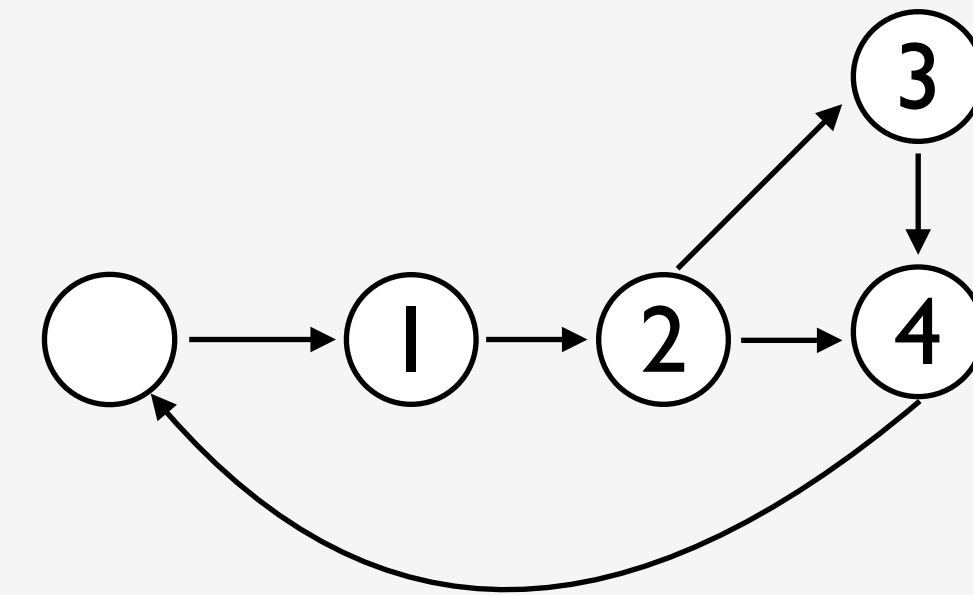
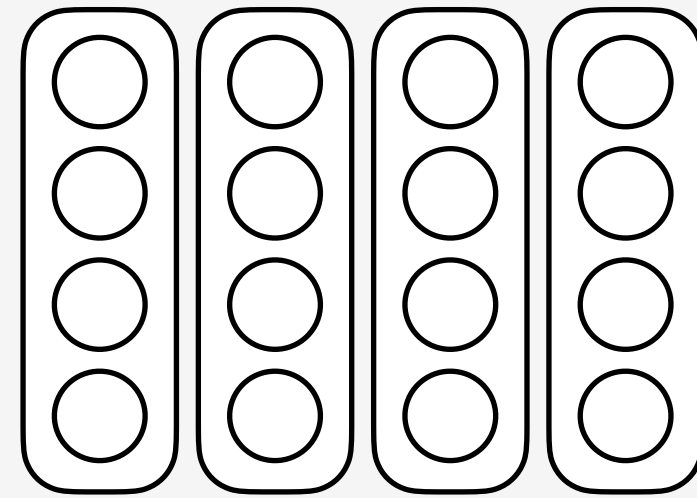


$$p_{\text{wfsa}}(y)$$

WFSA representation gives
character-level scores

Lexically-aware decoding for post-correction

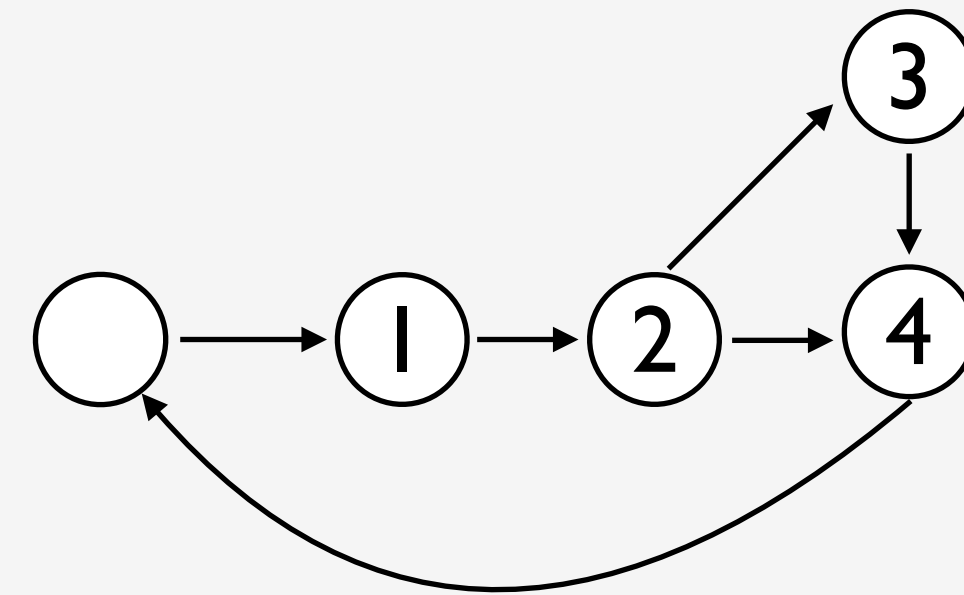
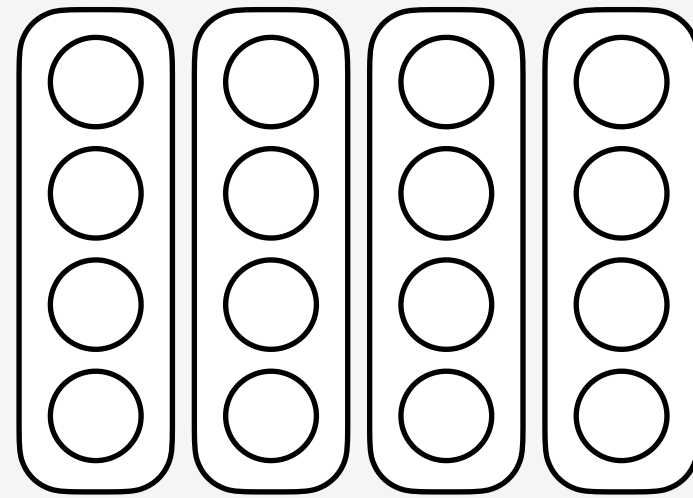
Linear interpolation to combine the probabilities for joint inference



$$P(y) = (1 - \lambda) * p_{\text{lstm}}(y) + \lambda * p_{\text{wfsa}}(y)$$

Lexically-aware decoding for post-correction

Linear interpolation to combine the probabilities for joint inference

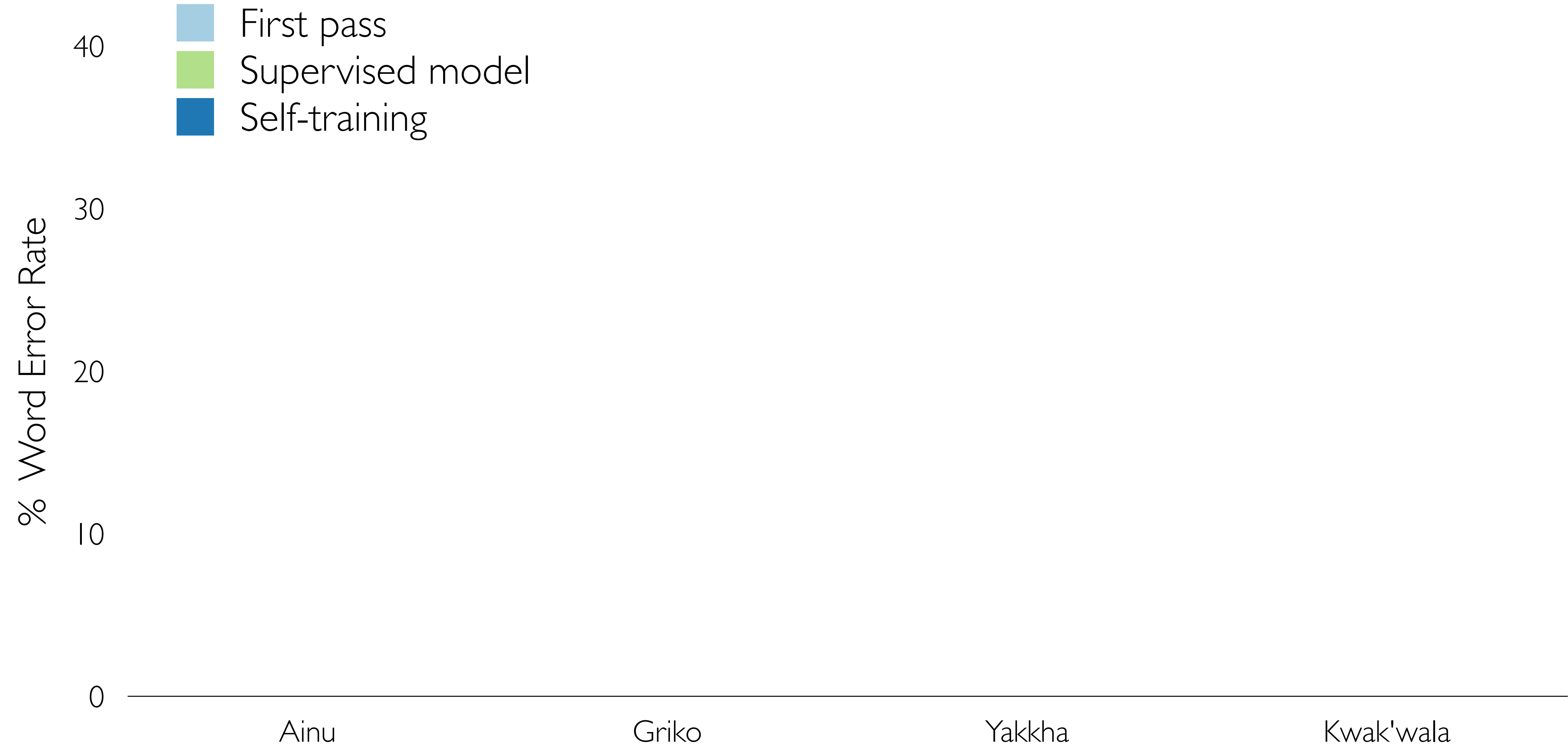


$$P(y) = (1 - \lambda) * p_{\text{lstm}}(y) + \lambda * p_{\text{wfsa}}(y)$$

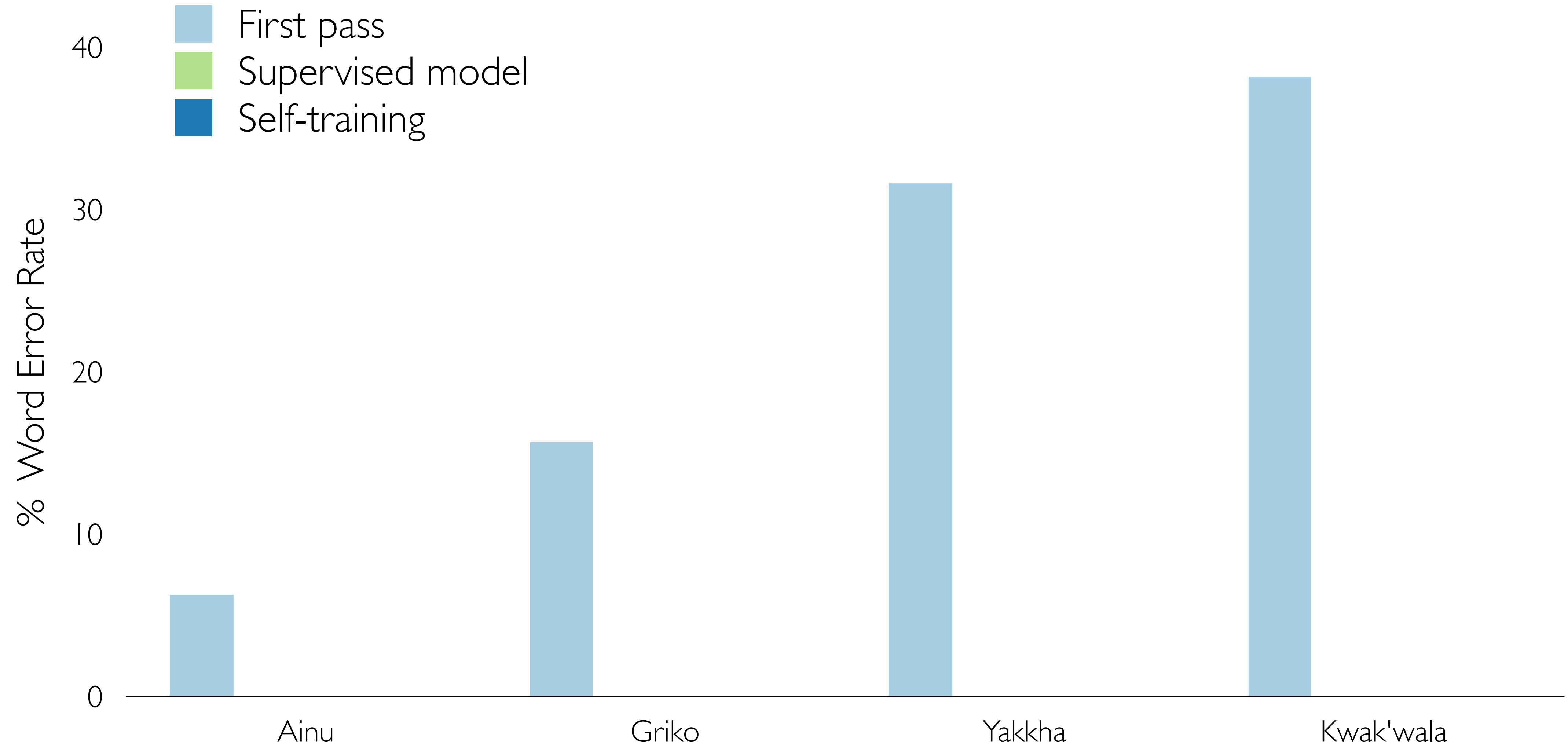
Hyperparameter for
linear interpolation

Experiments: does self-training improve performance?

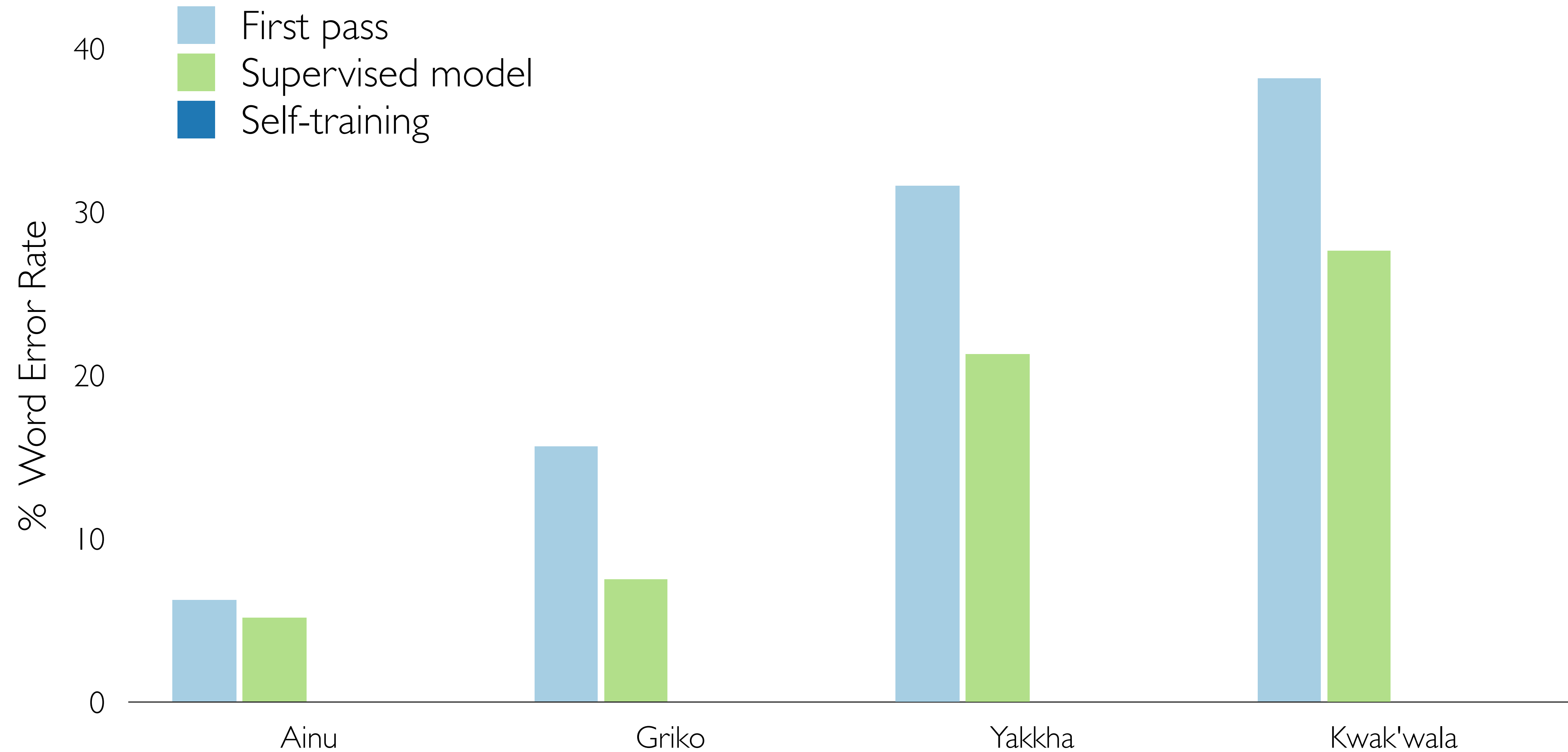
Experiments: does self-training improve performance?



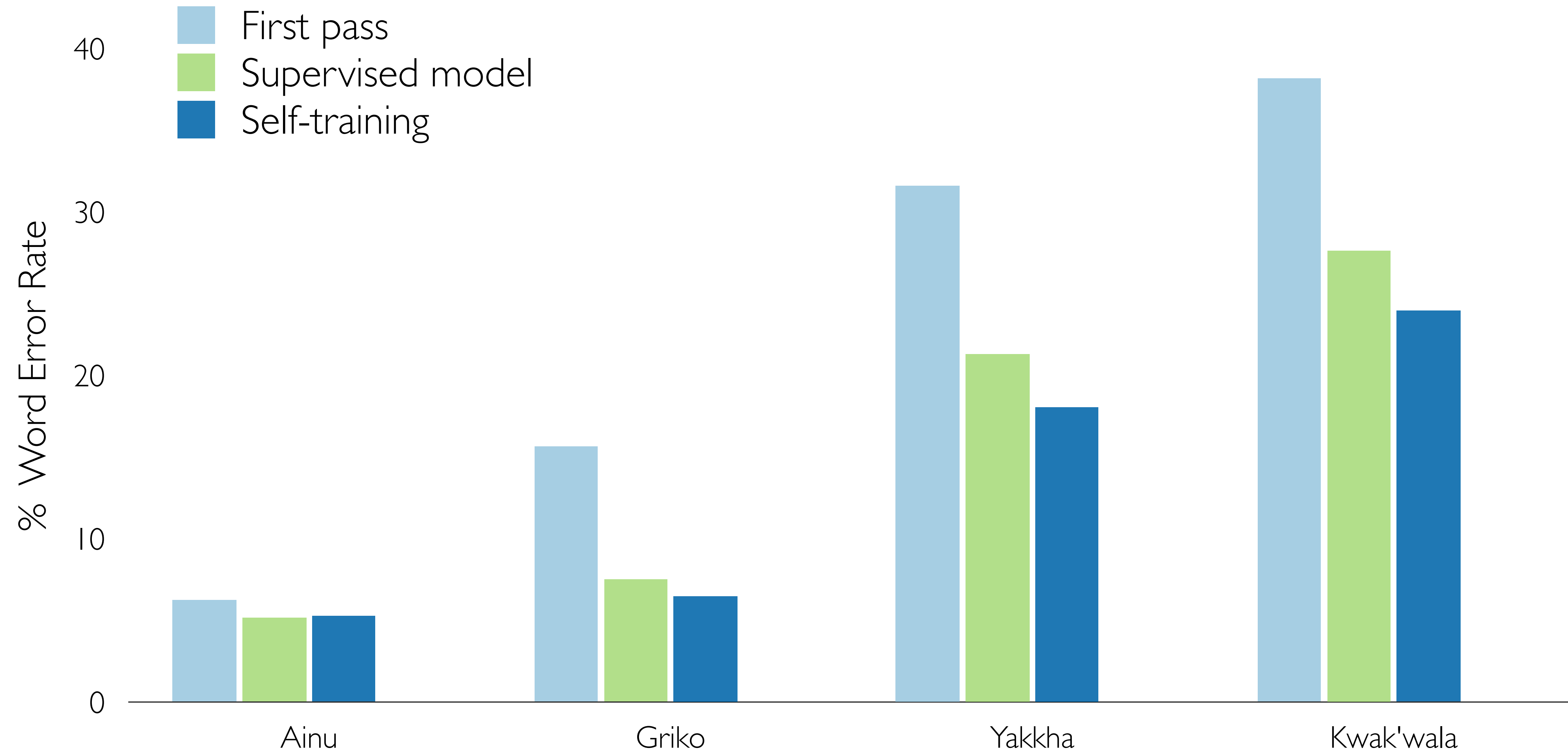
Experiments: does self-training improve performance?



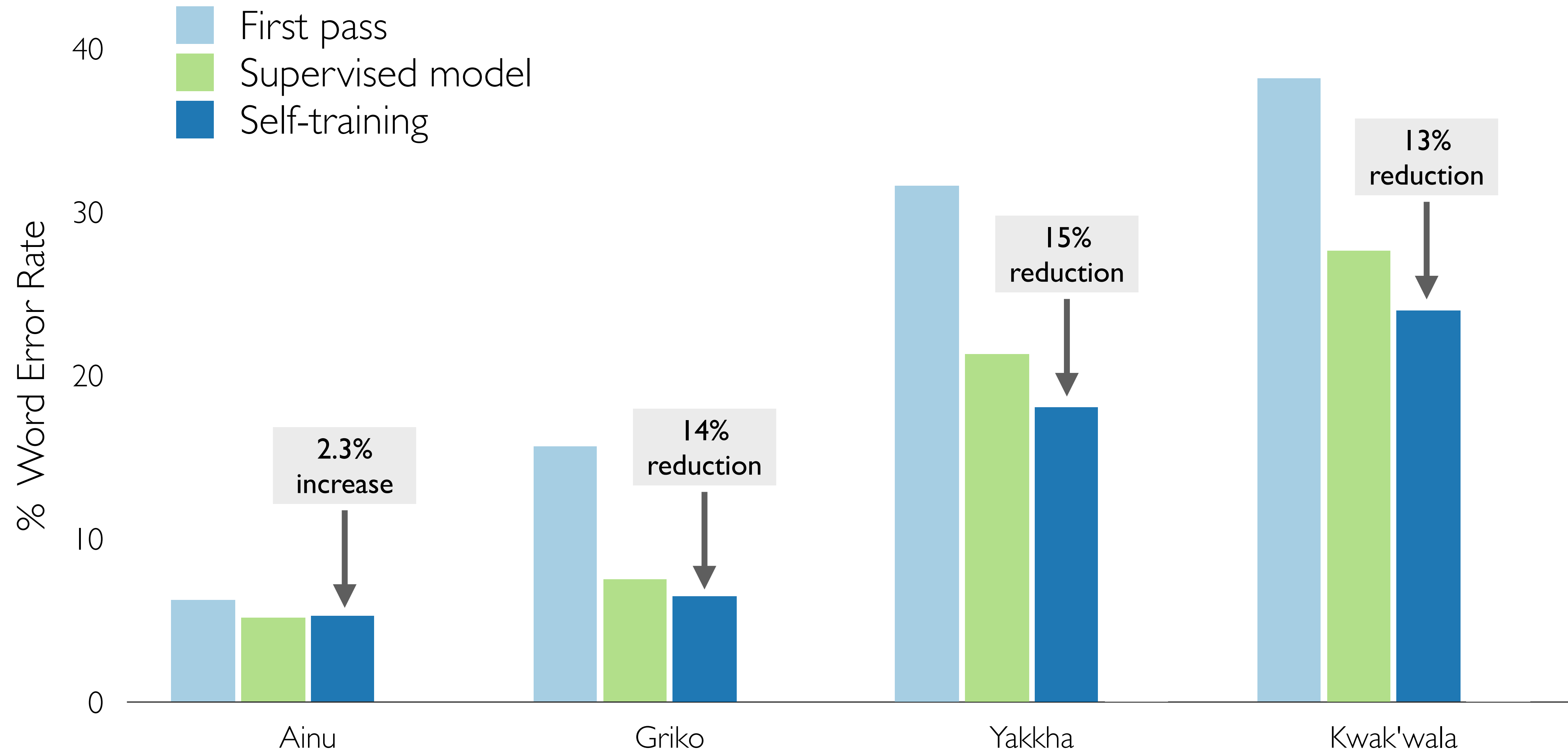
Experiments: does self-training improve performance?



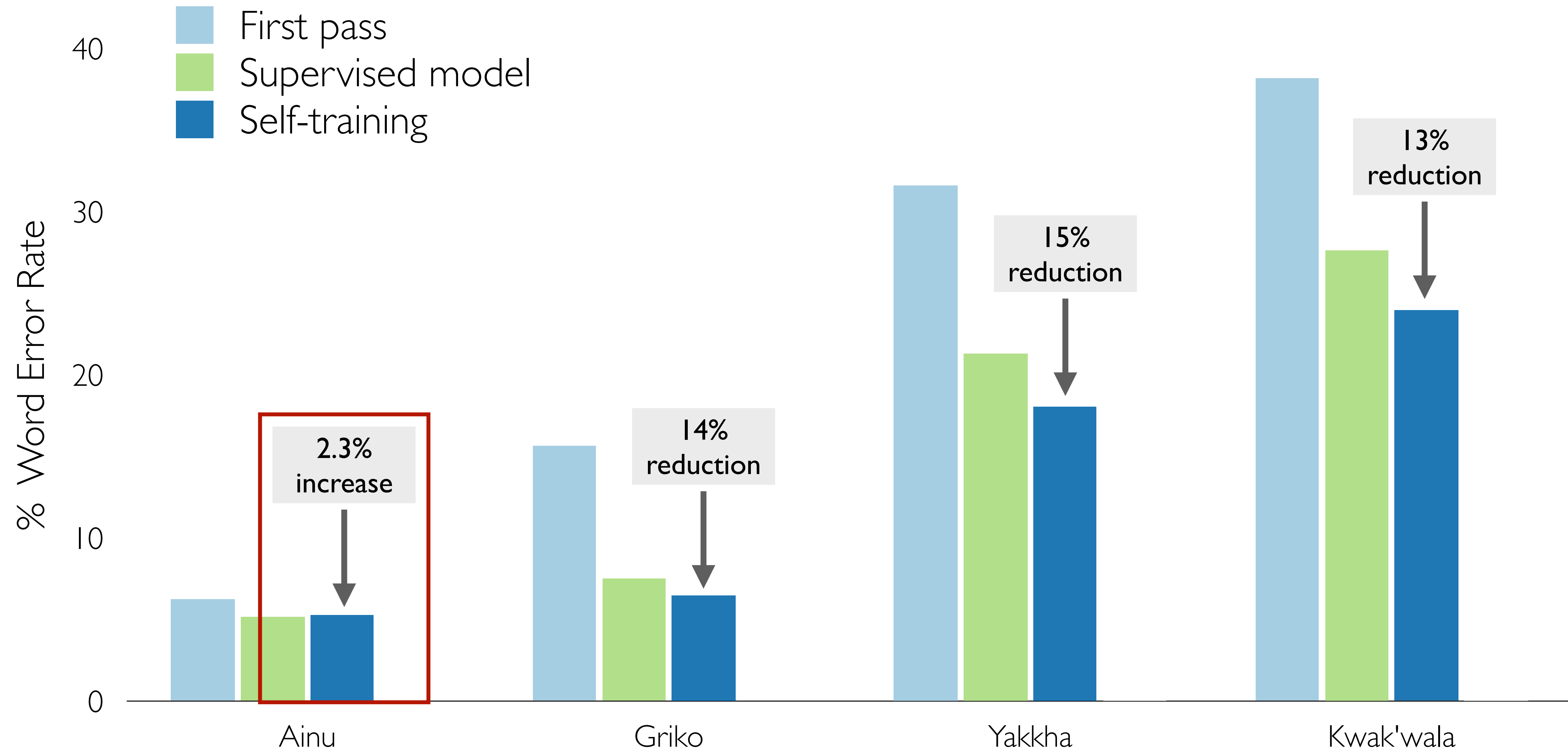
Experiments: does self-training improve performance?



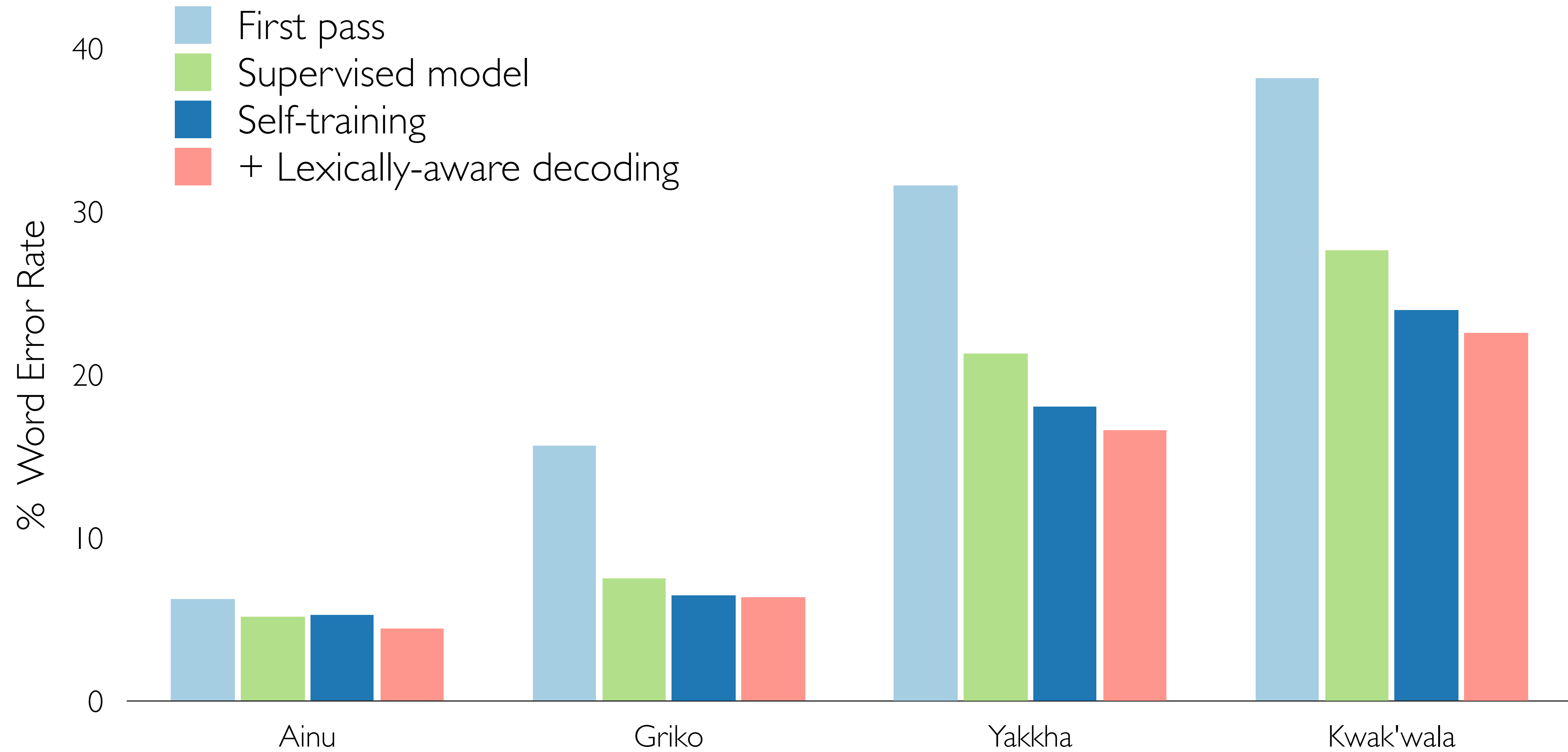
Experiments: does self-training improve performance?



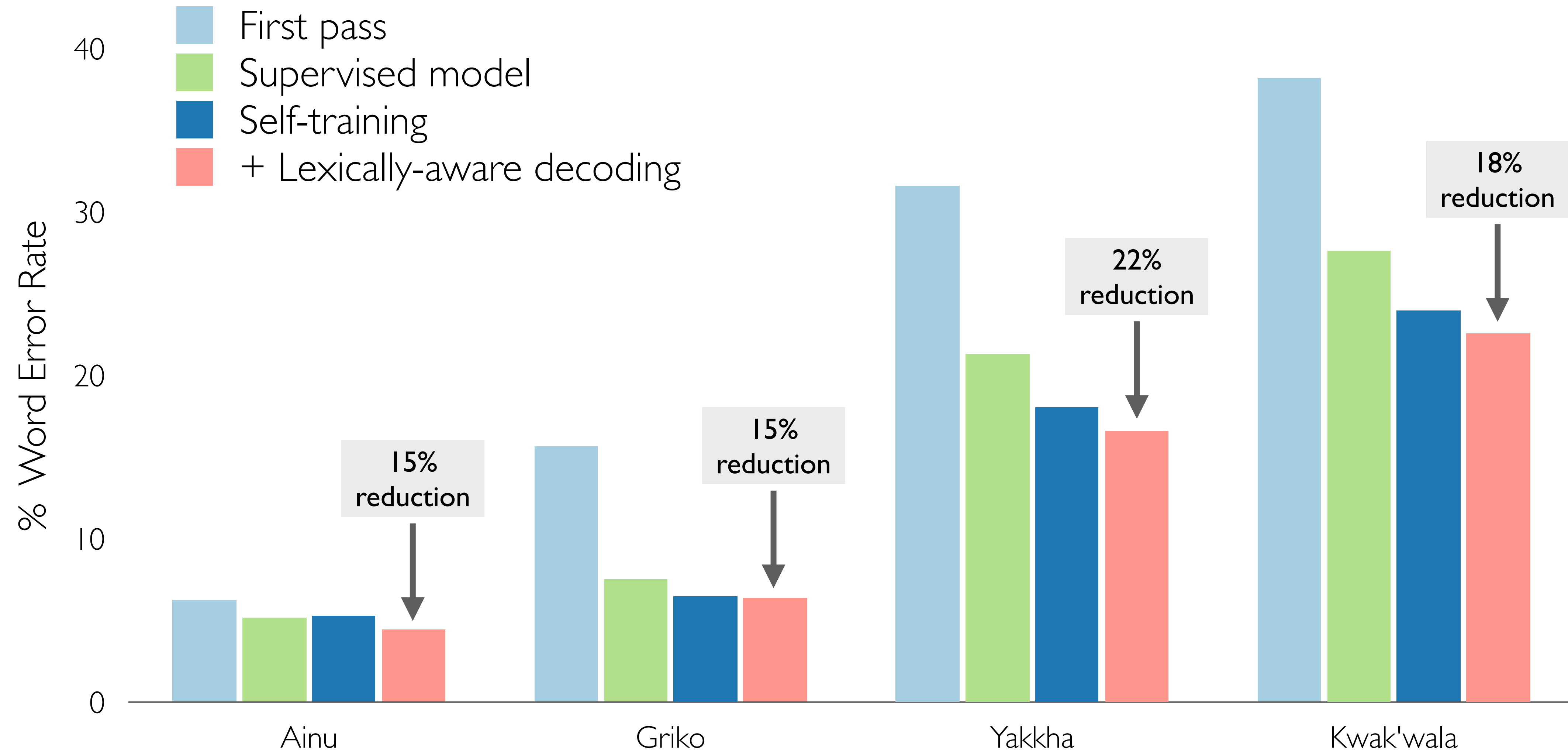
Experiments: does self-training improve performance?



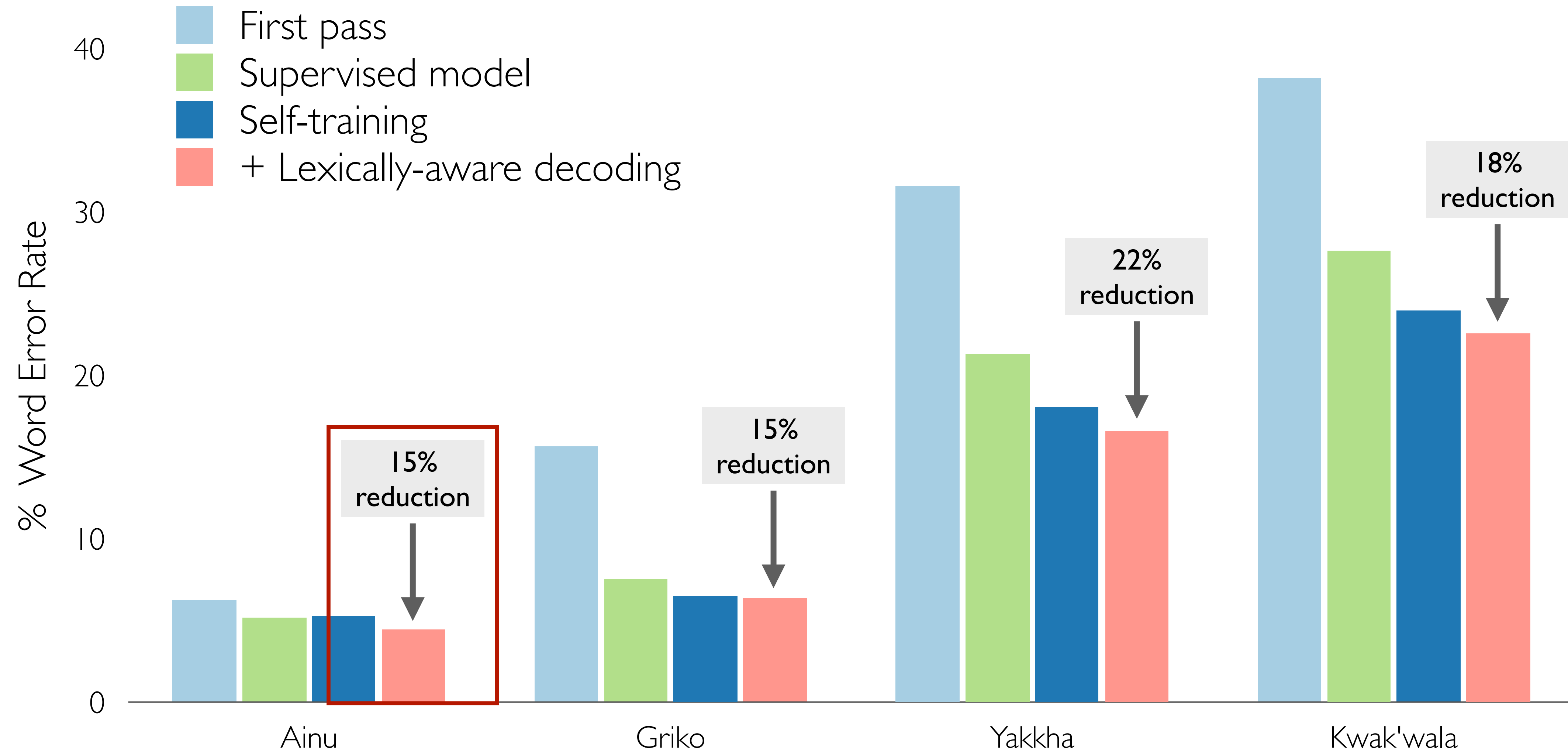
Experiments: does lexically-aware decoding counteract noise?



Experiments: does lexically-aware decoding counteract noise?



Experiments: does lexically-aware decoding counteract noise?



Summary

Summary

- Thousands of languages do not have easily accessible text to build NLP models

Summary

- Thousands of languages do not have easily accessible text to build NLP models
 - Text data does exist in many of these languages!

Summary

- Thousands of languages do not have easily accessible text to build NLP models
 - Text data does exist in many of these languages!
 - Locked away in non-machine-readable formats like printed books

Summary

- Thousands of languages do not have easily accessible text to build NLP models
 - Text data does exist in many of these languages!
 - Locked away in non-machine-readable formats like printed books
- OCR post-correction improves text extraction in very low-resourced settings

Summary

- Thousands of languages do not have easily accessible text to build NLP models
 - Text data does exist in many of these languages!
 - Locked away in non-machine-readable formats like printed books
- OCR post-correction improves text extraction in very low-resourced settings

Multi-source model: ↓ WER 17% – 52%

Summary

- Thousands of languages do not have easily accessible text to build NLP models
 - Text data does exist in many of these languages!
 - Locked away in non-machine-readable formats like printed books
- OCR post-correction improves text extraction in very low-resourced settings

Multi-source model: ↓ WER 17% – 52%

Semi-supervised with lexically-aware decoding: ↓ WER 29% – 59%

Impact case study: Kwak'wala

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

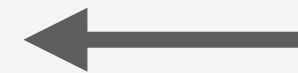
Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!ex^éid xült!ēdex^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qex^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaem ānēx^éēdxa qlēxa^élō qa^{és} t!ēqwa-
 w gw k· pela. Wä, g'íl^émēsē lāxa qlēxa^élaxs laē mōgwalīas lāx māg'īn
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexā^éyē qa^{és} lä lents!ēs lāx
 m äx^éxō L!ema^éisasēs g'ōkwē. Wä, lä xē^éx"ts!ālasa hā^éyāl^éa t!ēsem lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éūsdešelaq qa^{és} lä k'!ō-
 wū yo gwīlelaq lāxēs wülē^élasē g'ōkwaxēs wūlasē^éwē g'ōkwa qa^{és} lä gūgē-
 lāx^ém nōlīsas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!ema^éisē k'!ōx-
 ēy k'!ōtelaxēs t!ägats!ē lexā^éya. Wä, laxaē ēt!ēd t!äx^éts!ālasa t!ēsemē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xē^éx"ts!ālasa t!ēsemē lāxēs
 xegwats!ē t!ēsema. Wä, g'íl^éemxaāwisē gwanāla lōk"sēxs laē

Produced by Franz
Boas in 1921



Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qEX^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxa qlēxa^élō qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lāxa qlēxa^élaxs laē mōgwalīlas lāx māg'īn
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lEXa^éyē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!ālasa hā^éyāl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwīlelaq lāxēs wülē^élasē g'ōkwaxēs wūlasē^éwē g'ōkwa qa^{és} lä gūgE-
 lāx^ém nōlīsas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
 ēy k'!ōtelaxēs t!ägats!ē lEXa^éya. Wä, laxaē ēt!ēd t!äxts!ālasa t!ēSEMē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xE^éx"ts!ālasa t!ēSEMē lāxēs
 xEGwats!ē t!ēSEMa. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- Tremendous cultural and linguistic value!

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qEX^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxa qlēxa!lē qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lāxa qlēxa!laxs laē mōgwalīlas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexayē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!ālasa hā^éyāl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwīlelaq lāxēs wülē^élasē g'ōkwaxēs wūlasē^éwē g'ōkwa qa^{és} lä gūgE-
 lāx^ém nōlīsas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
 ēy k'!ōtelaxēs t!ägats!ē lexayā. Wä, laxaē ēt!ēd t!äxts!ālasa t!ēSEMē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xE^éx"ts!ālasa t!ēSEMē lāxēs
 xEgwats!ē t!ēSEMA. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- Tremendous cultural and linguistic value!
- Minimally accessible to researchers

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qEX^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxa qlēxa!lē qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lälxa qlēxa!laxs laē mōgwalilas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexayē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!älasa hä^éyäl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwILElaq lāxēs wülē^élasē g'ōkwaxēs wülase^éwē g'ōkwa qa^{és} lä gūgE-
 lāx^ém nōlissas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
 éy k'!ōtelaxēs t!ägats!ē lexayā. Wä, laxaē ēt!ēd t!äxts!älasa t!ēSEMē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xE^éx"ts!älasa t!ēSEMē lāxēs
 xEgwats!ē t!ēSEMA. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- Tremendous cultural and linguistic value!
- Minimally accessible to researchers
- Manual search in scanned images

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'il^εmēsē gwālexs laē äx^εälēlōdxēs menyayowē qa^εs k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^εid xült!ēDEX^εwālag'i-
 ol las Wä, g'il^εmēsē gwāla laē äx^εēdxa selbekwē dewēxa qa^εs qEX^εälē-
 m g'il^εlōc g'il^εmēsē gwālexs laē hōx^εidaem äñēx^εēdxa qlēxa!lē qa^εs t!ēqwa-
 w gw k' pela. Wä, g'il^εmēsē lälxa qlēxa!laxs laē mōgwalilas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^εēdxa lexayē qa^εs lä lents!ēs lāx
 m äx^εxō L!EMa^εisasēs g'ōkwē. Wä, lä xE^εx"ts!älasa hä^εyäl^εa t!ēSEM lāq.
 la gw W Wä, g'il^εmēsē gwanāla lōk"sēxs laē k'!ōx^εüsdēslaq qa^εs lä k'!ō-
 wü yo gwILElaq lāxēs wülē^εlasē g'ōkwaxēs wülase^εwē g'ōkwa qa^εs lä gügE-
 lāx^εm nōlissas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^εisē k'!ōx-
 εy k'!ōtelaxēs t!ägats!ē lexayā. Wä, laxaē ēt!ēd t!äxts!älasa t!ēSEMē
 lāq. Wä, la^εnēk'ēda waōkwē bāk'lumas xE^εx"ts!älasa t!ēSEMē lāxēs
 xEgwats!ē t!ēSEMA. Wä, g'il^εEMxaāwisē gwanāla lōk"sēxs laē

- Tremendous cultural and linguistic value!
- Minimally accessible to researchers
- Manual search in scanned images
- Legacy orthography that is hard to read

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!ex^éid xült!ēdex^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qex^éale-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaem ānēx^éēdxa qlēxa^élō qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lāxa qlēxa^élaxs laē mōgwalīlas lāx māg'īn
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexā^éyē qa^{és} lä lents!ēs lāx
 m äx^éxō L!ema^éisasēs g'ōkwē. Wä, lä xē^éx"ts!ālasa hā^éyāl^éa t!ēsem lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wū yo gwīlelaq lāxēs wülē^élasē g'ōkwaxēs wūlasē^éwē g'ōkwa qa^{és} lä gūgē-
 lāx^ém nōlīsas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!ema^éisē k'!ōx-
 ēy k'!ōtelaxēs t!ägats!ē lexā^éya. Wä, laxaē ēt!ēd t!äxts!ālasa t!ēsemē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xē^éx"ts!ālasa t!ēsemē lāxēs
 xēgwats!ē t!ēsema. Wä, g'íl^éemxaāwisē gwanāla lōk"sēxs laē

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qEX^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxa qlēxa!lē qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lälxa qlēxa!laxs laē mōgwalilas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexayē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!älasa hä^éyäl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwILElaq lāxēs wülē^élasē g'ōkwaxēs wülase^éwē g'ōkwa qa^{és} lä gūgE-
 lāx^ém nōlissas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
 éy k'!ōtelaxēs t!ägats!ē lexayā. Wä, laxaē ēt!ēd t!äxts!älasa t!ēSEMē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xE^éx"ts!älasa t!ēSEMē lāxēs
 xEGwats!ē t!ēSEMA. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- **1500+ pages converted** to a machine-readable format

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qEX^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxa qlēxa^élō qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lälxa qlēxa^élaxs laē mōgwalilas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lEXa^éyē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!älasa hä^éyäl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwILElaq lāxēs wülē^élasē g'ōkwaxēs wülase^éwē g'ōkwa qa^{és} lä gügE-
 lāx^ém nōlissas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
 éy k'!ōtelaxēs t!ägats!ē lEXa^éya. Wä, laxaē ēt!ēd t!äxts!älasa t!ēSEMē
 lāq. Wä, lä^énēk'ēda waōkwē bāk'lumas xE^éx"ts!älasa t!ēSEMē lāxēs
 xEGwats!ē t!ēSEMA. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- **1500+ pages converted** to a machine-readable format
- **Searchable!**

Impact case study: Kwak'wala

- Collaborating with documentary linguists and Kwak'wala speakers
 - Identify documents that would be **most useful to extract text from**

Boas texts: 10 volumes of Kwak'wala language and community documentation

Wä, g'íl^émēsē gwālexs laē äx^éälēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwälēg'a (*fig.*). Wä, lä xäl!EX^éid xült!ēDEX^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxā selbekwē dewēxa qa^{és} qEX^éälē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaEM äñēx^éēdxā qlēxa^élō qa^{és} t!ēqwa-
 w gw k' pela. Wä, g'íl^émēsē lälxa qlēxa^élaxs laē mōgwalilas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxā lexā^éyē qa^{és} lä lents!ēs lāx
 m äx^éxō L!EMa^éisasēs g'ōkwē. Wä, lä xE^éx"ts!älasa hä^éyäl^éa t!ēSEM lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éüsdēslaq qa^{és} lä k'!ō-
 wü yo gwILElaq lāxēs wülē^élasē g'ōkwaxēs wülase^éwē g'ōkwa qa^{és} lä gügE-
 lāx^ém nōliskas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!EMa^éisē k'!ōx-
^éy k'!ōtelaxēs t!ägats!ē lexā^éya. Wä, laxaē ēt!ēd t!äxts!älasa t!ēSEMē
 lāq. Wä, lä ^énēk'ēda waōkwē bāk'lumas xE^éx"ts!älasa t!ēSEMē lāxēs
 xEgwats!ē t!ēSEMA. Wä, g'íl^éEMxaāwisē gwanāla lōk"sēxs laē

- **1500+ pages converted** to a machine-readable format
- **Searchable!**
- **Legacy orthography can be automatically transliterated** to modern writing systems

Impact case study: Kwak'wala

Wä, g'íl^émēsē gwālexs laē äx^éalēlōdxēs menyayowē qa^{és} k'at!al
 lē g'ōhalē, g'a gwālēg'a (*fig.*). Wä, lä xäl!ex^éid xült!ēdex^éwālag'i-
 ol las Wä, g'íl^émēsē gwāla laē äx^éēdxa selbekwē dewēxa qa^{és} qex^éalē-
 m g'íl^élōc g'íl^émēsē gwālexs laē hōx^éidaem ānēx^éēdxa qlēxa^élō qa^{és} t!ēqwa-
 w gw k· pela. Wä, g'íl^émēsē lāxa qlēxa^élaxs laē mōgwalīlas lāx māg'in
 n gw W walisasēs legwilē. Wä, lä äx^éēdxa lexayē qa^{és} lä lents!ēs lāx
 m äx^éxō L!ema^éisasēs g'ōkwē. Wä, lä xē^éx"ts!ālasa hā^éyāl^éa t!ēsem lāq.
 la gw W Wä, g'íl^émēsē gwanāla lōk"sēxs laē k'!ōx^éūsdēselaq qa^{és} lä k'!ō-
 wū yo gwīlelaq lāxēs wülē^élasē g'ōkwaxēs wūlasē^éwē g'ōkwa qa^{és} lä gūgē-
 lāx^ém nōlisas lāxēs legwilē. Wä, lä xwēlaqents!ēsa lāxa L!ema^éisē k'!ōx-
 ēy k'!ōtelaxēs t!ägats!ē lexayā. Wä, laxaē ēt!ēd t!äxts!ālasa t!ēsemē
 lāq. Wä, la^énēk'ēda waōkwē bāk'lumas xē^éx"ts!ālasa t!ēsemē lāxēs
 xēgwats!ē t!ēsema. Wä, g'íl^éemxaāwisē gwanāla lōk"sēxs laē

- **1500+ pages converted** to a machine-readable format
- **Searchable!**
- **Legacy orthography can be automatically transliterated** to modern writing systems

Impact and applications: beyond this talk

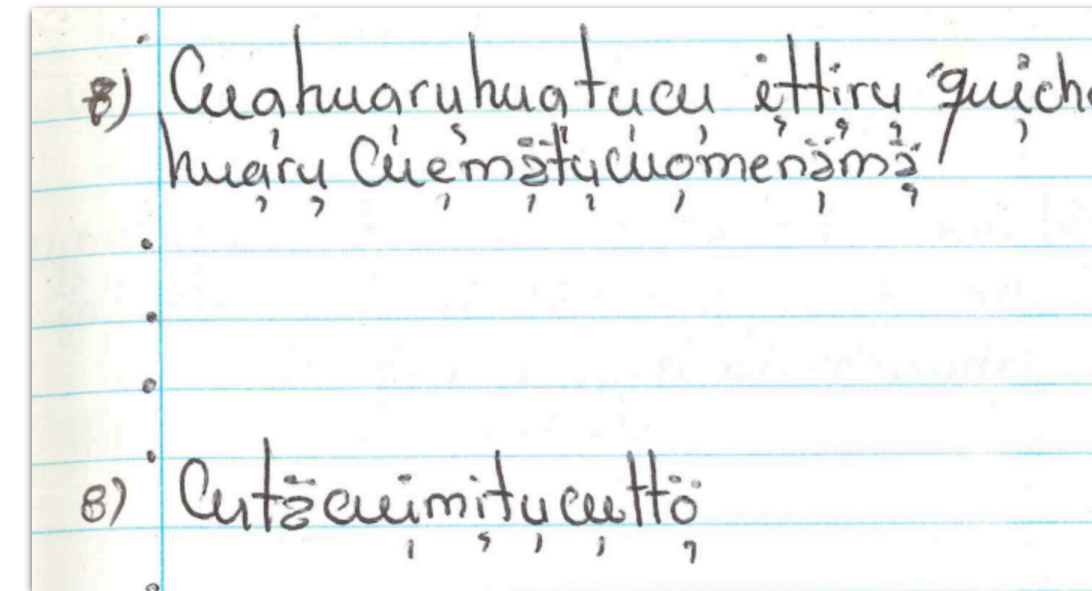
Impact and applications: beyond this talk

Our software is open-source and has been used on many other languages!

Impact and applications: beyond this talk

Our software is open-source and has been used on many other languages!

Bhutia
Sanskrit Quechua
Igbo Tibetan
Piaroa Secwepemctsin
Pintupi-Luritja



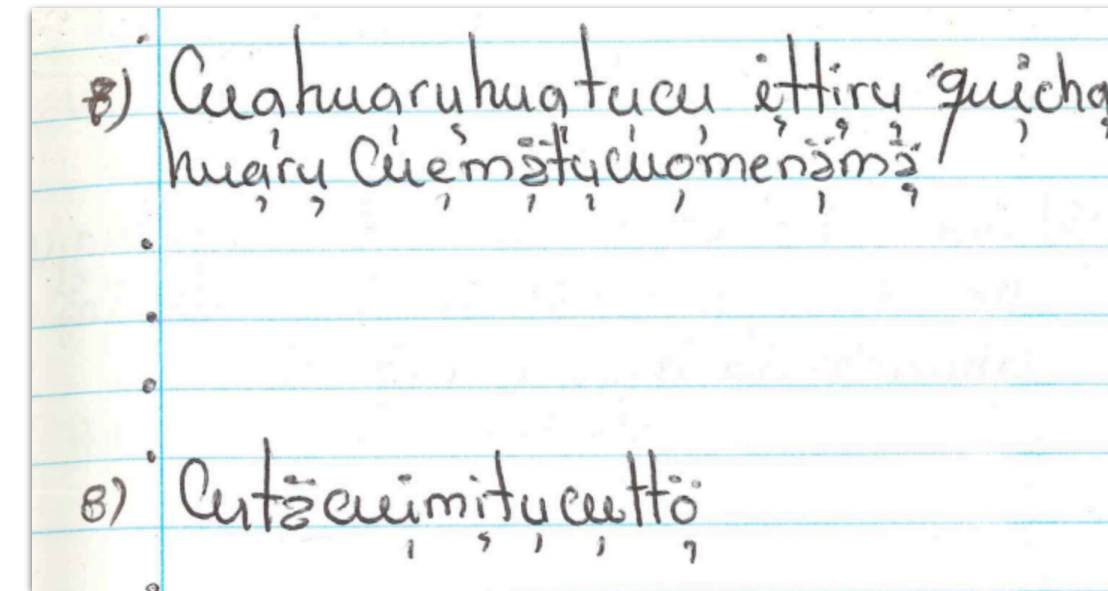
ana aka Igbo-
ana aka *n* [LL HH] twig; tree-branch. *var.* **aba aka**.
anaga *n* [HHH] surgical needle.
anagba *n* [HHH] anklet; bracelet.
anam *n* [LLL] cloth work loosely around the waist; loin cloth.
anambe *n* [LLL] (Mbieri) branching tuber of the cocoyam. *var.* **anünü**; **anünü-ede**.
anasī *n* [HHH] head-wife; first wife in a polygamous household; also called "nwanyī isi ci".

Tuyuta tjutangka kutjuya anu kutjupa tjuta tjutalingku nyinangu. Palunyatjanu kuunyi watjanu "Kala nyuntu ananyi ngurra nyuntup Utjula Ingkata kutjupalpi nyinaku. Tjana Palulanguruya watjalkulpi anangu tjutangka Kuunyi, Ingkata tjilpi paluru rawa nyinang ngalyanu Utjulakutu. Paluru rawa nyinangu nyinangu.
Palulangurulatju kala ngurrangkalpi tjarrp piitja nyangu tjilpi ulkumanu tjuta irriti ngurrara tjutanyatarra tjana papatayitja k Ngurra irrititjanuarra nyangu Ingkata irr

Impact and applications: beyond this talk

Our software is open-source and has been used on many other languages!

Bhutia
Sanskrit Quechua
Igbo Tibetan
Piaroa Secwepemctsin
Pintupi-Luritja



ana aka	Igbo-]
ana aka <i>n</i> [LL HH] twig; tree-branch. <i>var.</i> aba aka .	
anaga <i>n</i> [HHH] surgical needle.	
anagba <i>n</i> [HHH] anklet; bracelet.	
anam <i>n</i> [LLL] cloth work loosely around the waist; loin cloth.	
anambe <i>n</i> [LLL] (Mbieri) branching tuber of the cocoyam. <i>var.</i> anünü ; anünü-ede .	
anasī <i>n</i> [HHH] head-wife; first wife in a polygamous household; also called "nwanyī isi ci".	

Extracting text to **train machine translation** for Pintupi-Luritja



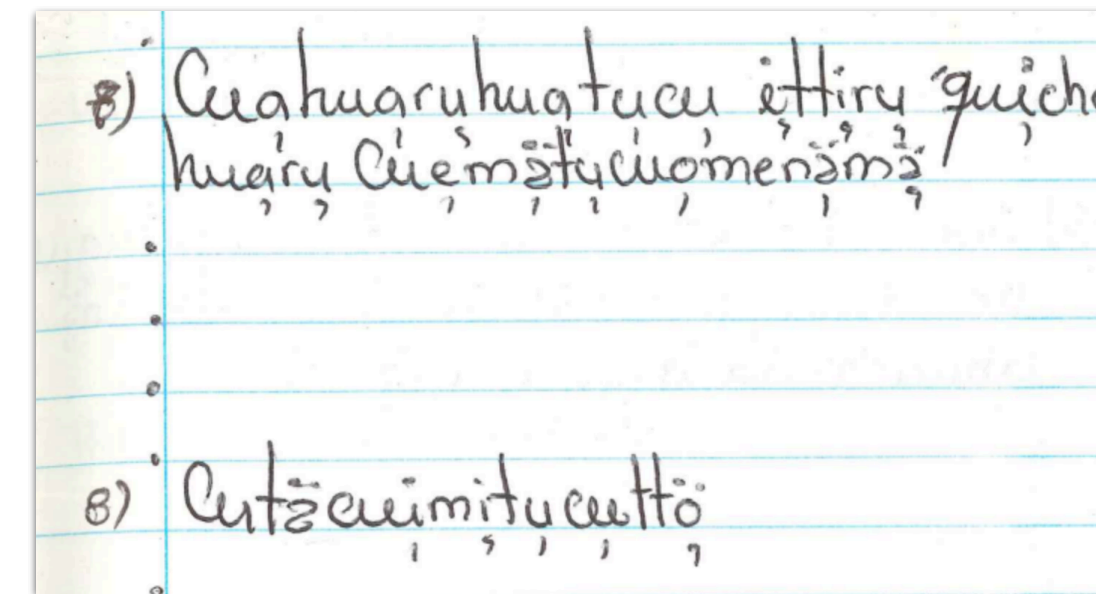
Tuyuta tjutangka kutjuya anu kutjupa tjuta tjutalingku nyinangu. Palunyatjanu kuunyi watjanu "Kala nyuntu ananyi ngurra nyuntup Utjula Ingkata kutjupalpi nyinaku. Tjana Palulanguruya watjalkulpi anangu tjutangka Kuunyi, Ingkata tjilpi paluru rawa nyinang ngalyanu Utjulakutu. Paluru rawa nyinangu nyinangu.
Palulangurulatju kala ngurrangkalpi tjarrp piitja nyangu tjilpi ulkumanu tjuta irriti ngurrara tjutanyatarra tjana papatayitja k Ngurra irrititjanutarra nyangu Ingkata irr

Impact and applications: beyond this talk

Our software is open-source and has been used on many other languages!

Bhutia
Sanskrit Quechua
Igbo Tibetan
Piaroa Secwepemctsin
Pintupi-Luritja

Automatic extraction of
handwritten speech
transcriptions in Piaroa



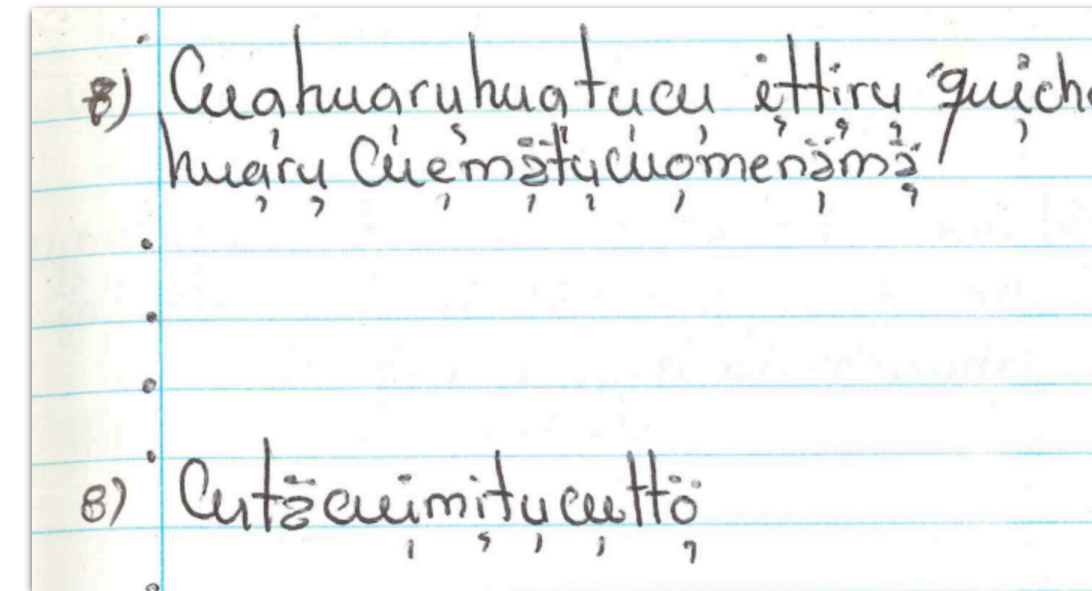
ana aka	Igbo-]
ana aka <i>n</i> [LL HH] twig; tree-branch. <i>var.</i> aba aka .	
anaga <i>n</i> [HHH] surgical needle.	
anagba <i>n</i> [HHH] anklet; bracelet.	
anam <i>n</i> [LLL] cloth work loosely around the waist; loin cloth.	
anambe <i>n</i> [LLL] (Mbieri) branching tuber of the cocoyam. <i>var.</i> anünü ; anünü-ede .	
anasī <i>n</i> [HHH] head-wife; first wife in a polygamous household; also called "nwanyī isi ci".	

Tuyuta tjutangka kutjuya anu kutjupa tjuta tjutalingku nyinangu. Palunyatjanu kuunyi watjanu "Kala nyuntu ananyi ngurra nyuntup Utjula Ingkata kutjupalpi nyinaku. Tjana Palulanguruya watjalkulpi anangu tjutangka Kuunyi, Ingkata tjilpi paluru rawa nyinang ngalyanu Utjulakutu. Paluru rawa nyinangu nyinangu.
Palulangurulatju kala ngurrangkalpi tjarrp piitja nyangu tjilpi ulkumanu tjuta irriti ngurrara tjutanyatarra tjana papatayitja k Ngurra irrititjanutarra nyangu Ingkata irr

Impact and applications: beyond this talk

Our software is open-source and has been used on many other languages!

Bhutia
Sanskrit Quechua
Igbo Tibetan
Piaroa Secwepemctsin
Pintupi-Luritja



ana aka	Igbo-]
ana aka <i>n</i> [LL HH] twig; tree-branch. <i>var.</i> aba aka .	
anaga <i>n</i> [HHH] surgical needle.	
anagba <i>n</i> [HHH] anklet; bracelet.	
anam <i>n</i> [LLL] cloth work loosely around the waist; loin cloth.	
anambe <i>n</i> [LLL] (Mbieri) branching tuber of the cocoyam. <i>var.</i> anünü ; anünü-ede .	
anasī <i>n</i> [HHH] head-wife; first wife in a polygamous household; also called "nwanyi isi ci".	

Tuyuta tjutangka kutjuya anu kutjupa tjuta tjutalingku nyinangu. Palunyatjanu kuunyi watjanu "Kala nyuntu ananyi ngurra nyuntup Utjula Ingkata kutjupalpi nyinaku. Tjana Palulanguruya watjalkulpi anangu tjutangka Kuunyi, Ingkata tjilpi paluru rawa nyinang ngalyanu Utjulakutu. Paluru rawa nyinangu nyinangu. Palulangurulatju kala ngurrangkalpi tjarrp piitja nyangu tjilpi ulkumanu tjuta irriti ngurrara tjutanyatarra tjana papatayitja k Ngurra irrititjanuarra nyangu Ingkata irr

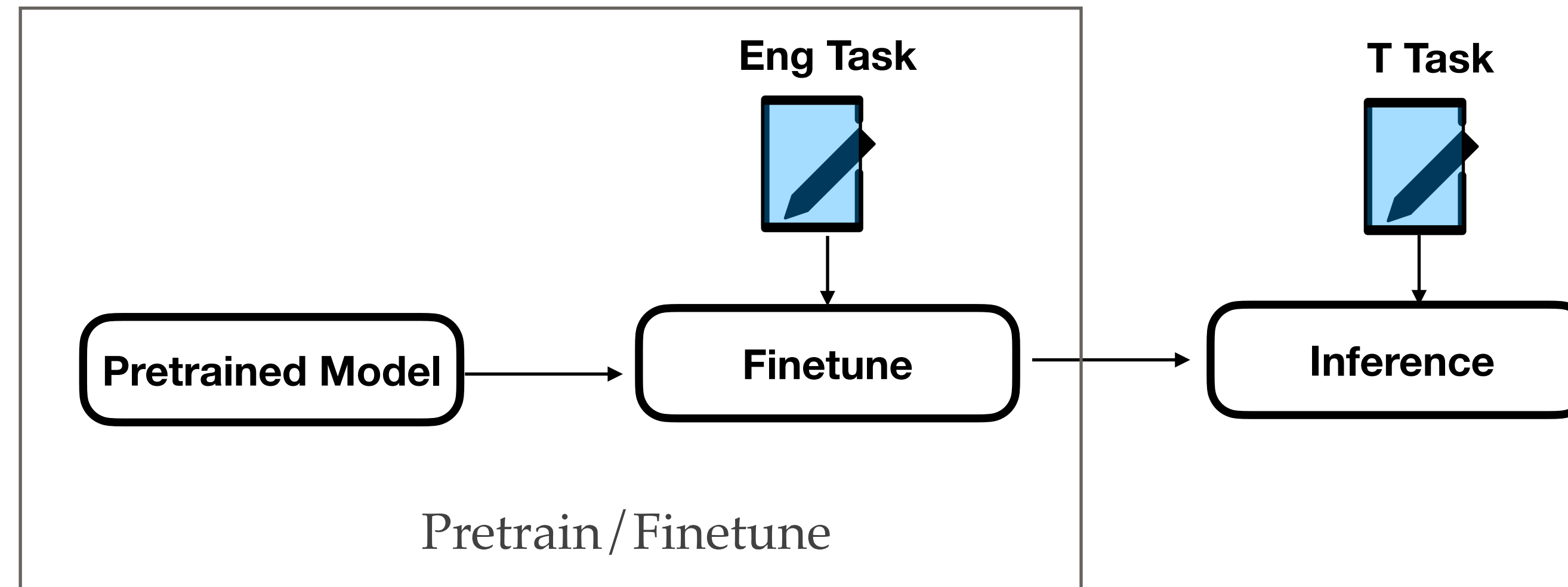
Print dictionaries in Igbo are high-coverage, but not digitized

Unlocking Bi-lingual Lexicons

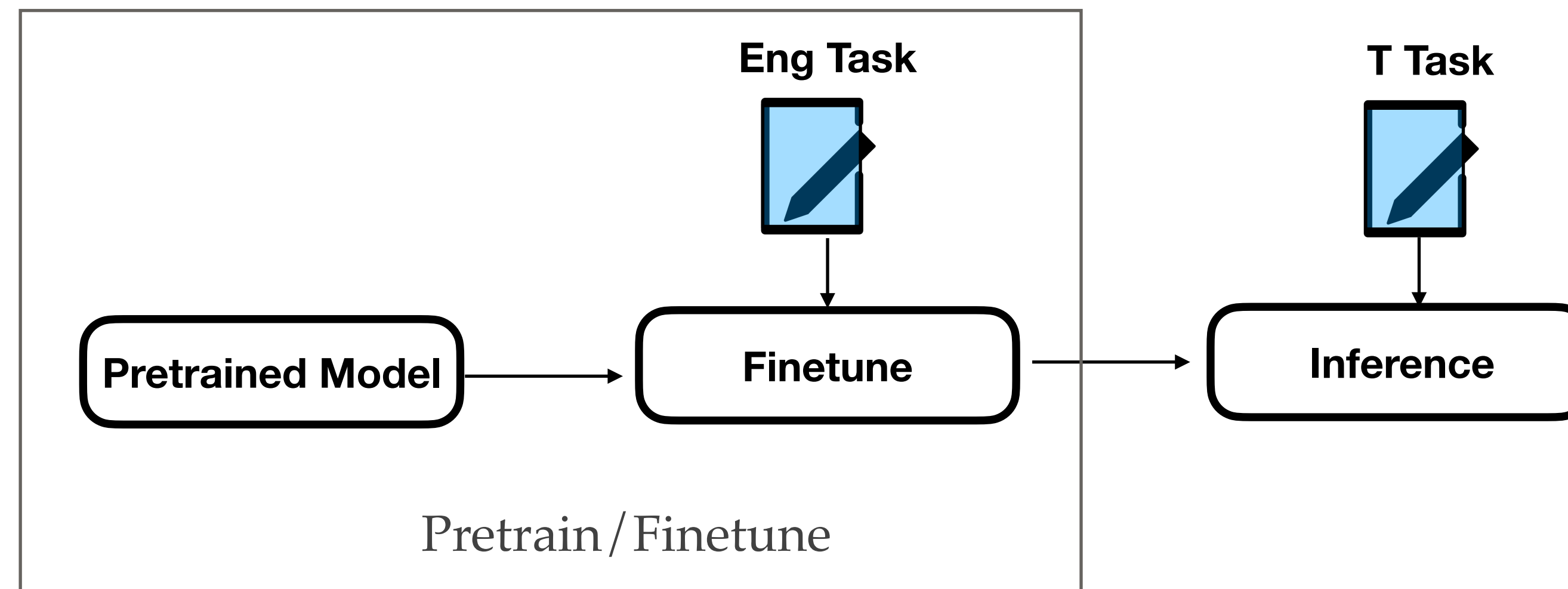


Xinyi Wang, Sebastian Ruder, Graham Neubig.
Expanding Pretrained Models to Thousands More Languages via Lexicon-based Adaptation.
ACL 2022.

Multilingual Pretrained Models

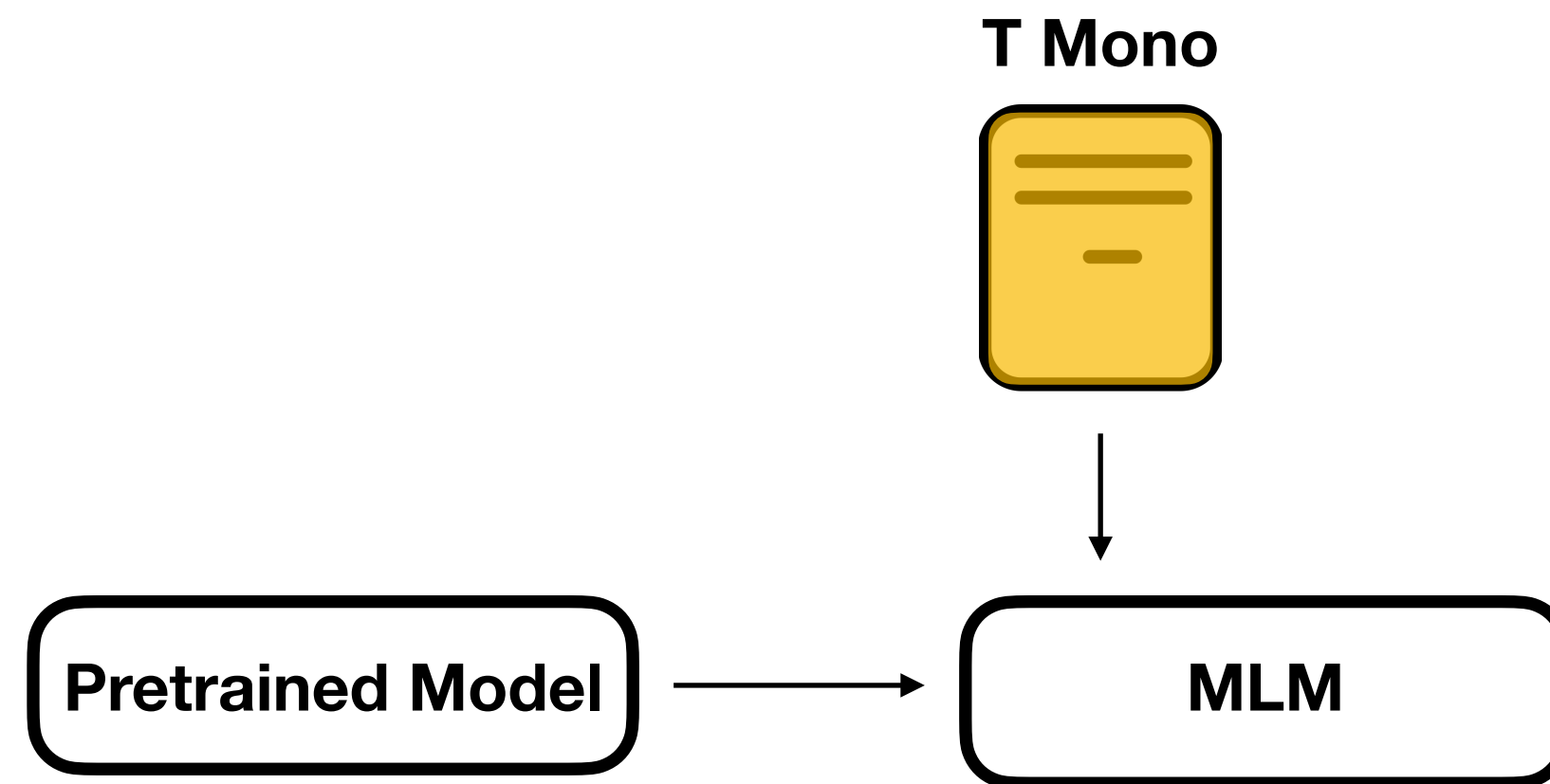


Multilingual Pretrained Models



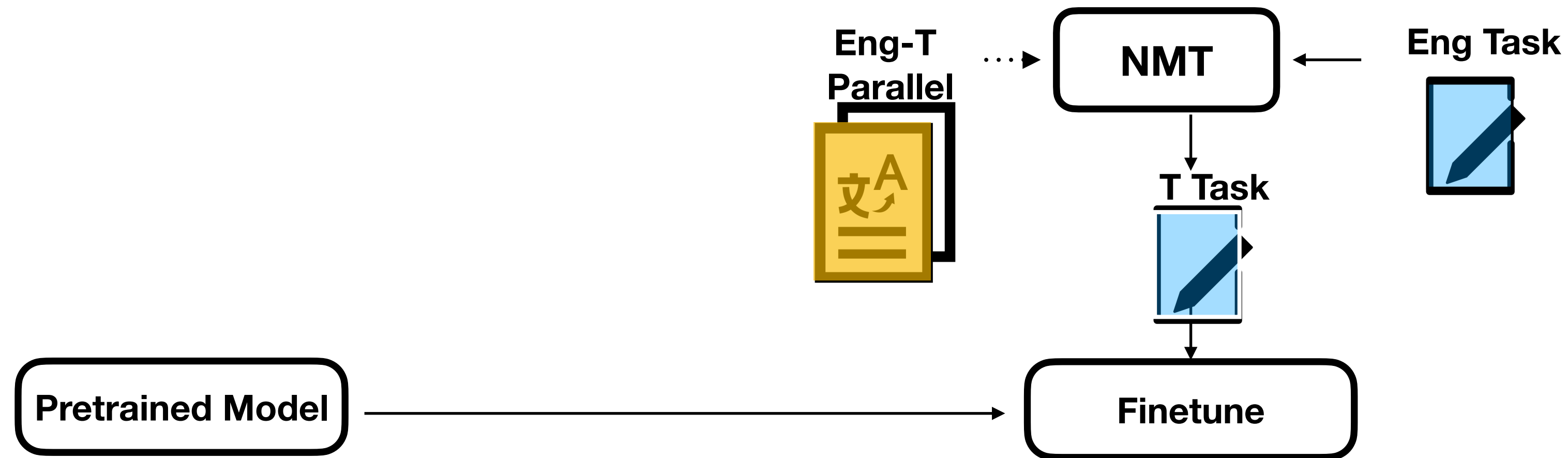
How to adapt the model for the language T?

Adaptation: Monolingual Data



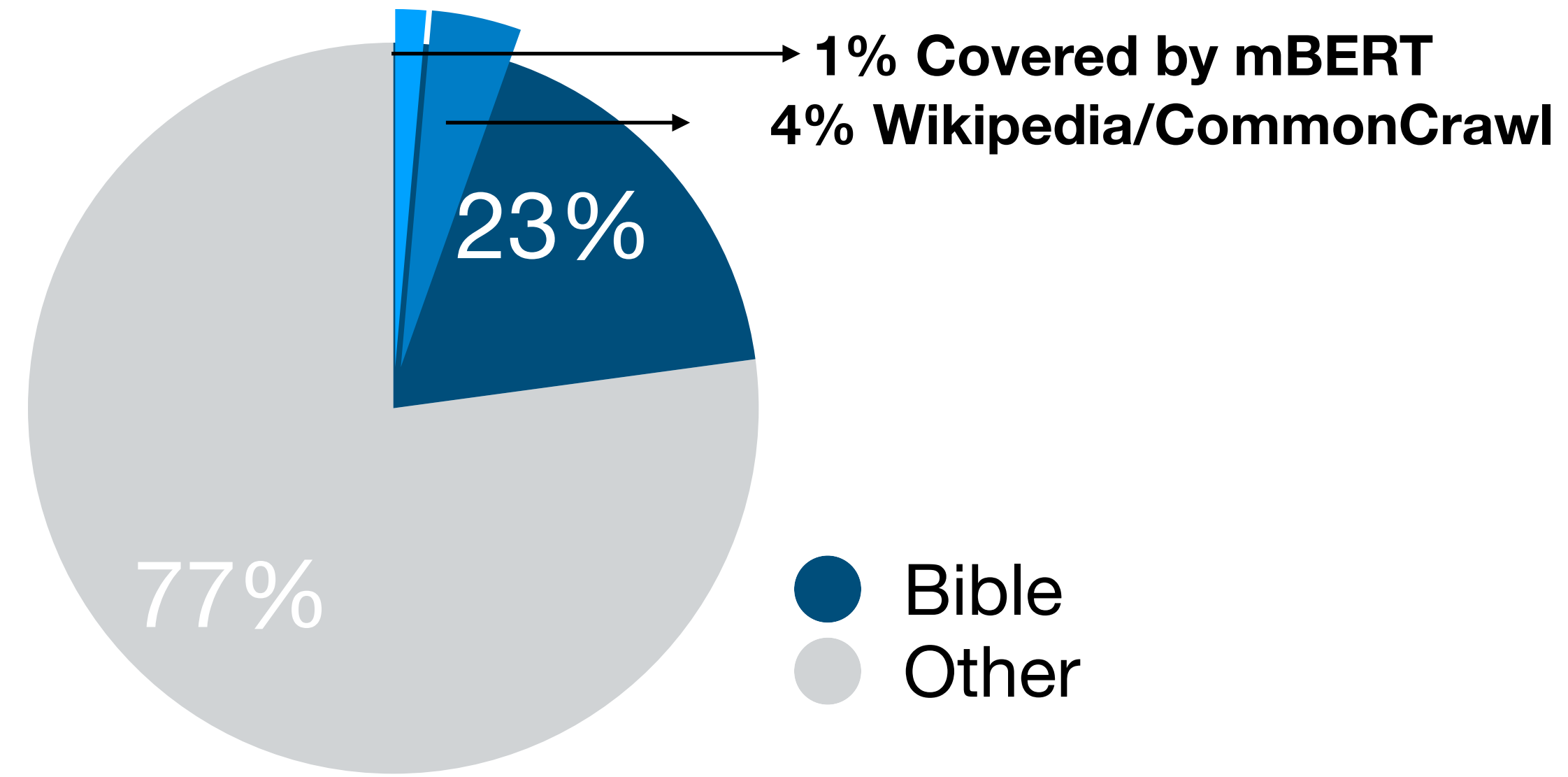
- e.g. Continued Masked Language Modeling(MLM) using monolingual data in the target language T

Adaptation: Parallel Data

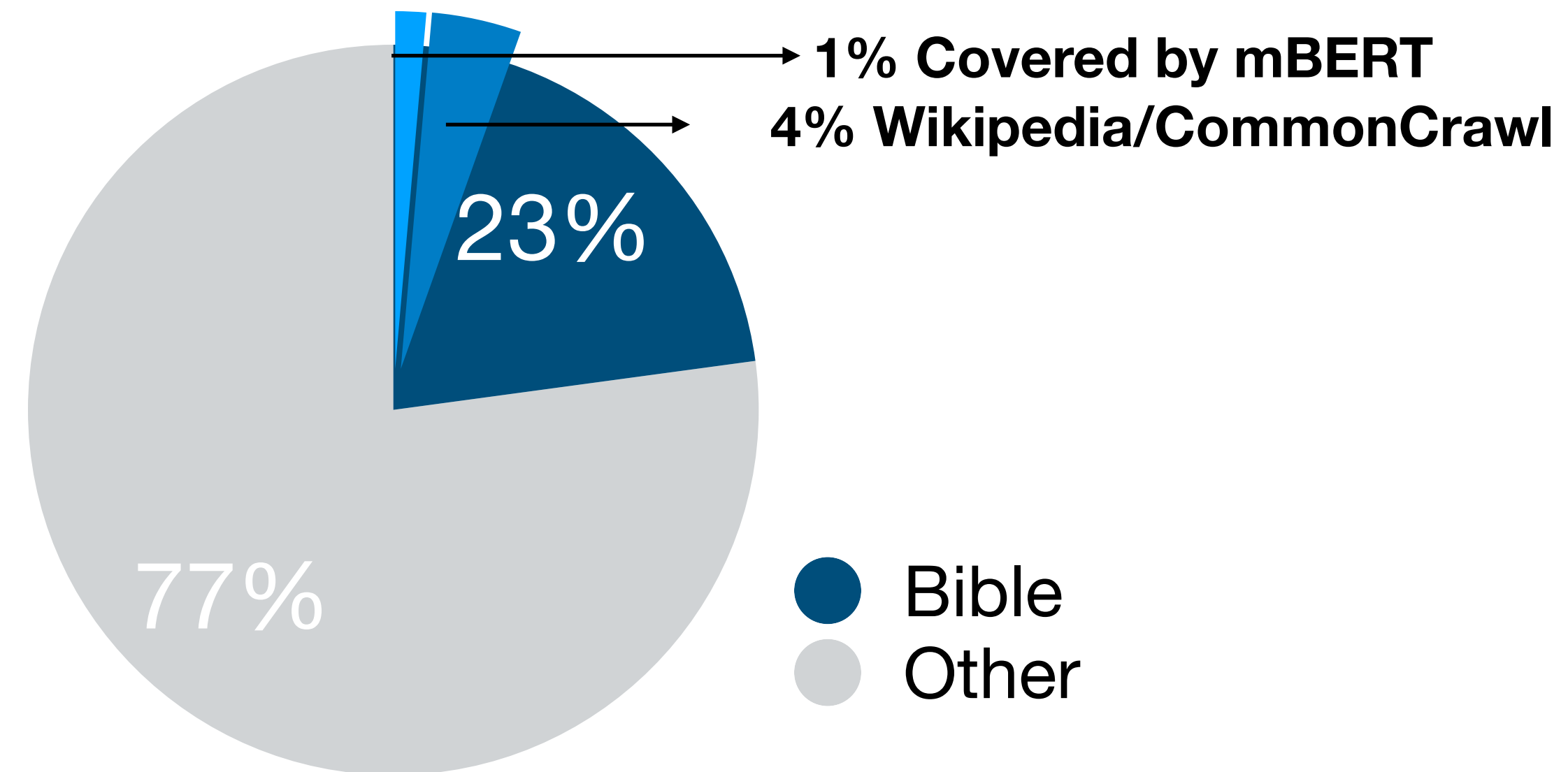


- e.g. Parallel Data: use best NMT system available to translate English task data into the target language T

Languages without Conventional Data

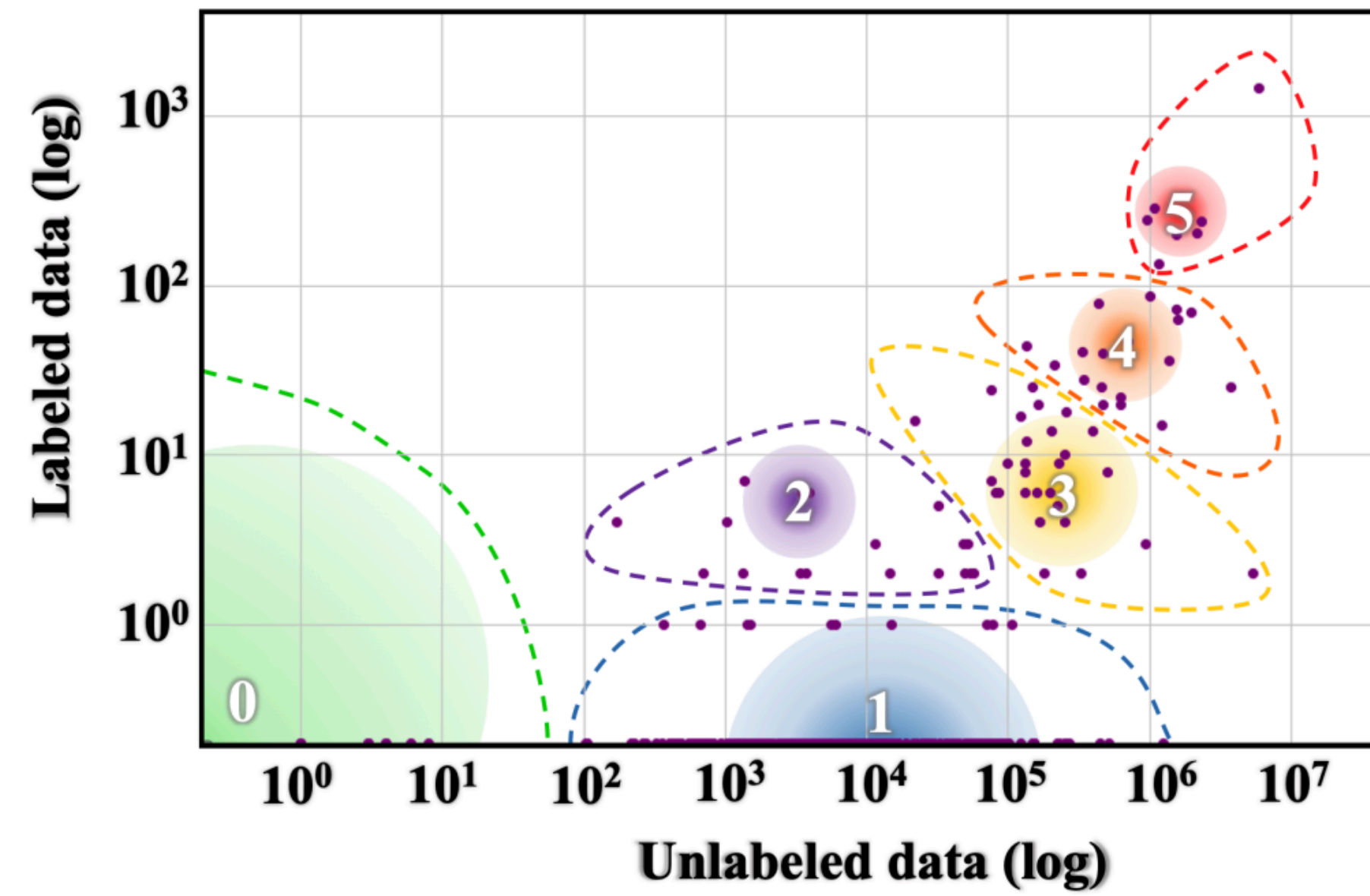


Languages without Conventional Data

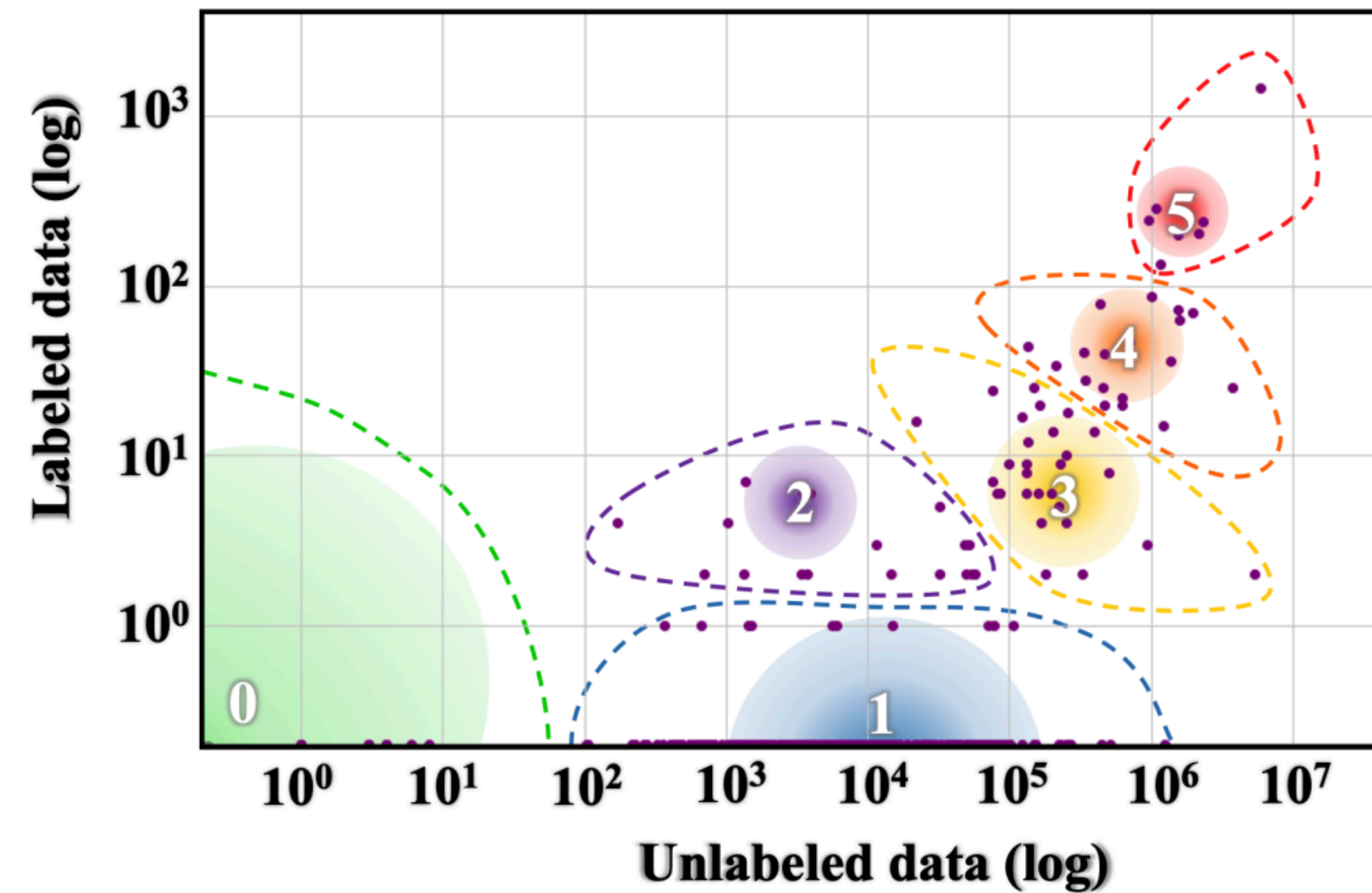


- Majority of languages in the world cannot benefit from progress in NLP due to lack of data

Two Groups of Low-resource Languages

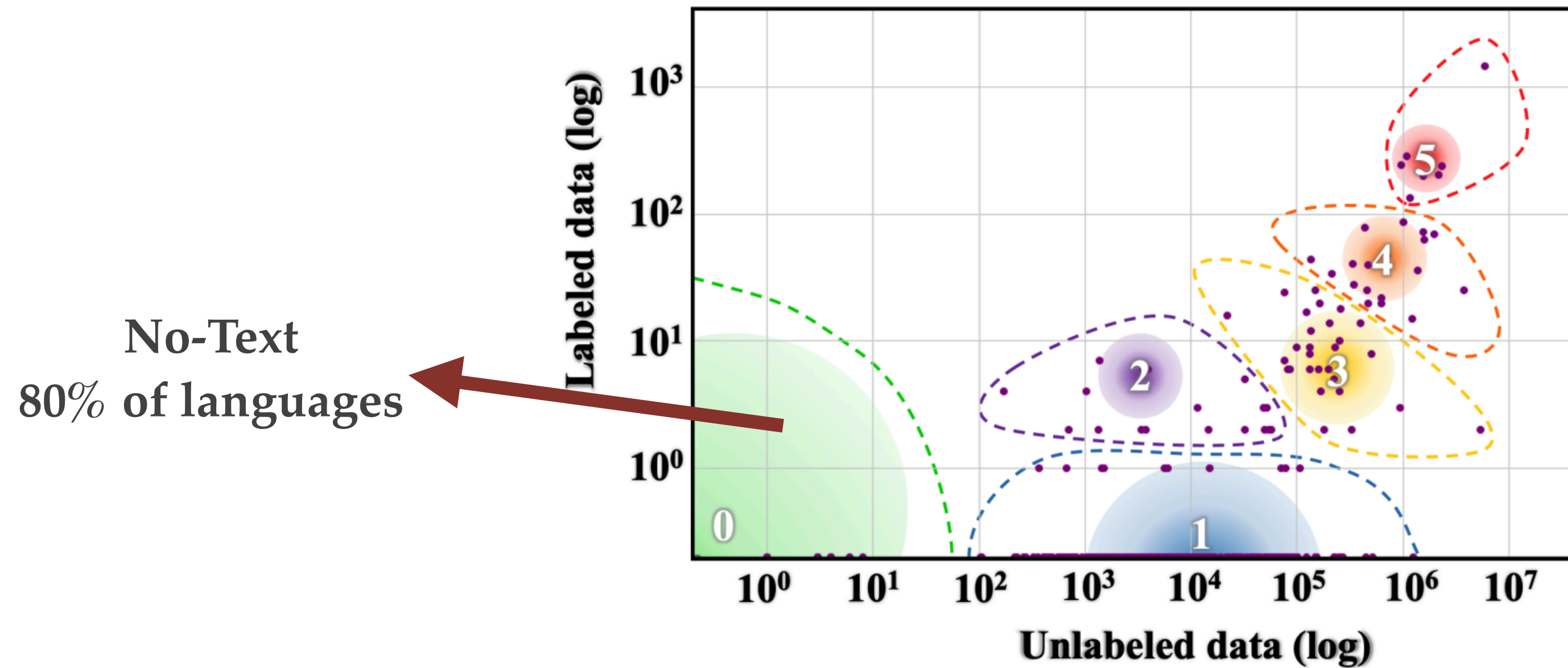


Two Groups of Low-resource Languages



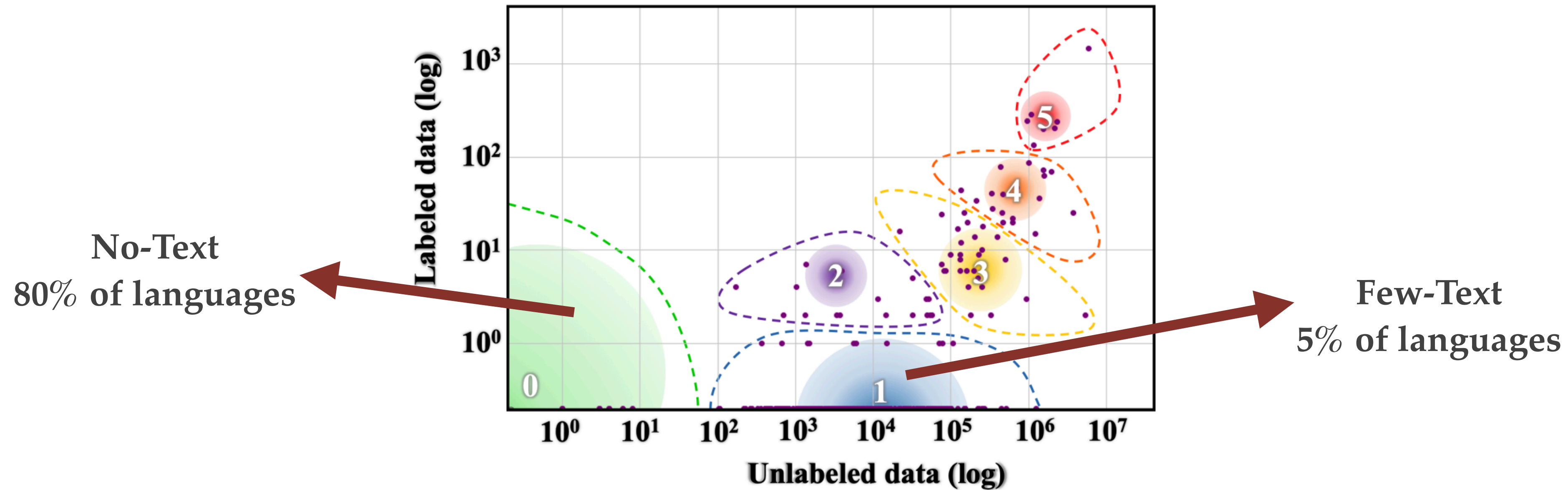
- Majority of World's languages cannot benefit from progress in NLP (Joshi et al. 2020)

Two Groups of Low-resource Languages



- Majority of World's languages cannot benefit from progress in NLP (Joshi et al. 2020)
 - No-Text: virtually no resource

Two Groups of Low-resource Languages



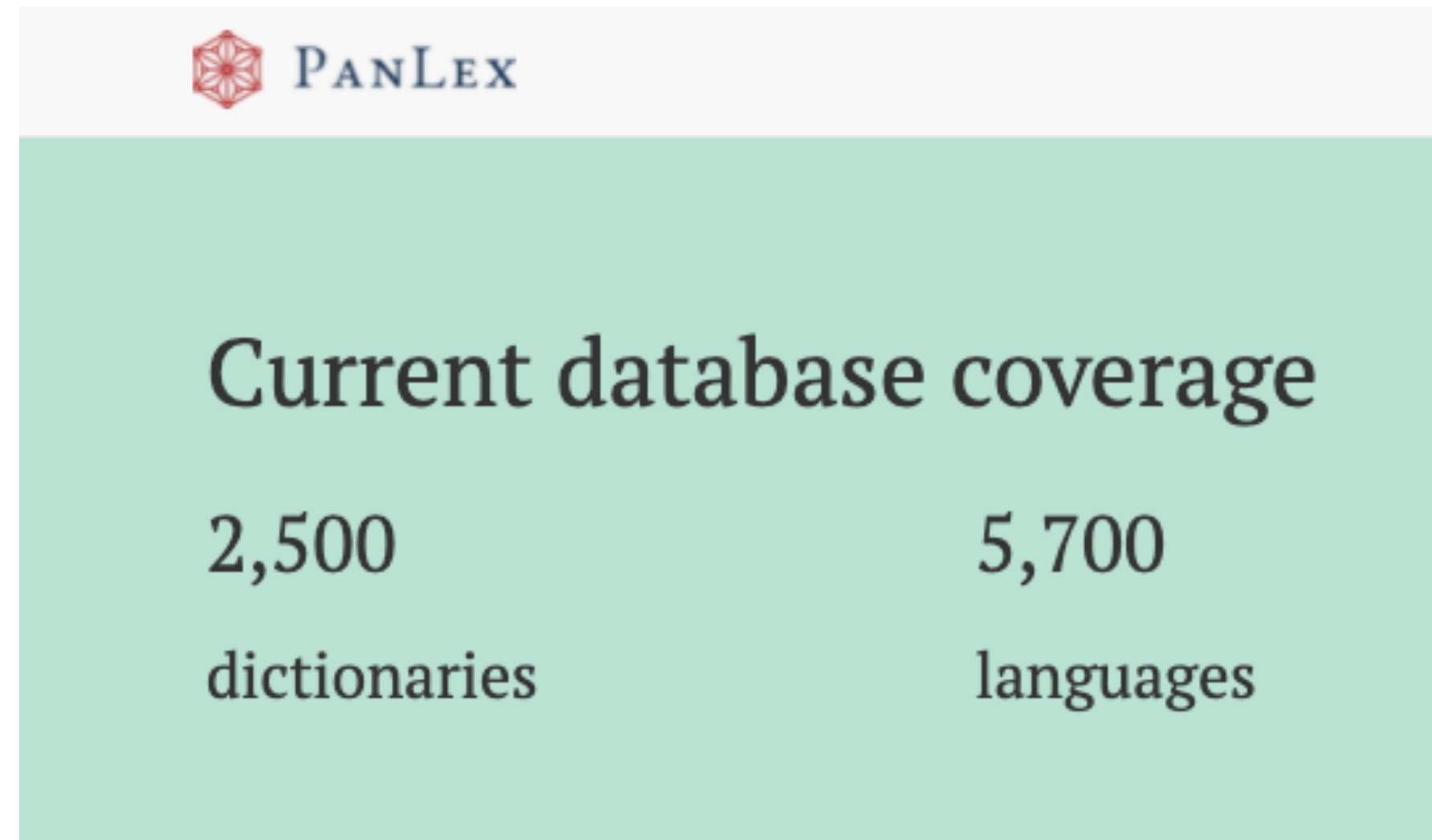
- Majority of World's languages cannot benefit from progress in NLP (Joshi et al. 2020)
 - No-Text: virtually no resource
 - Few-Text: very limited resource

Alternative Data Source

Alternative Data Source

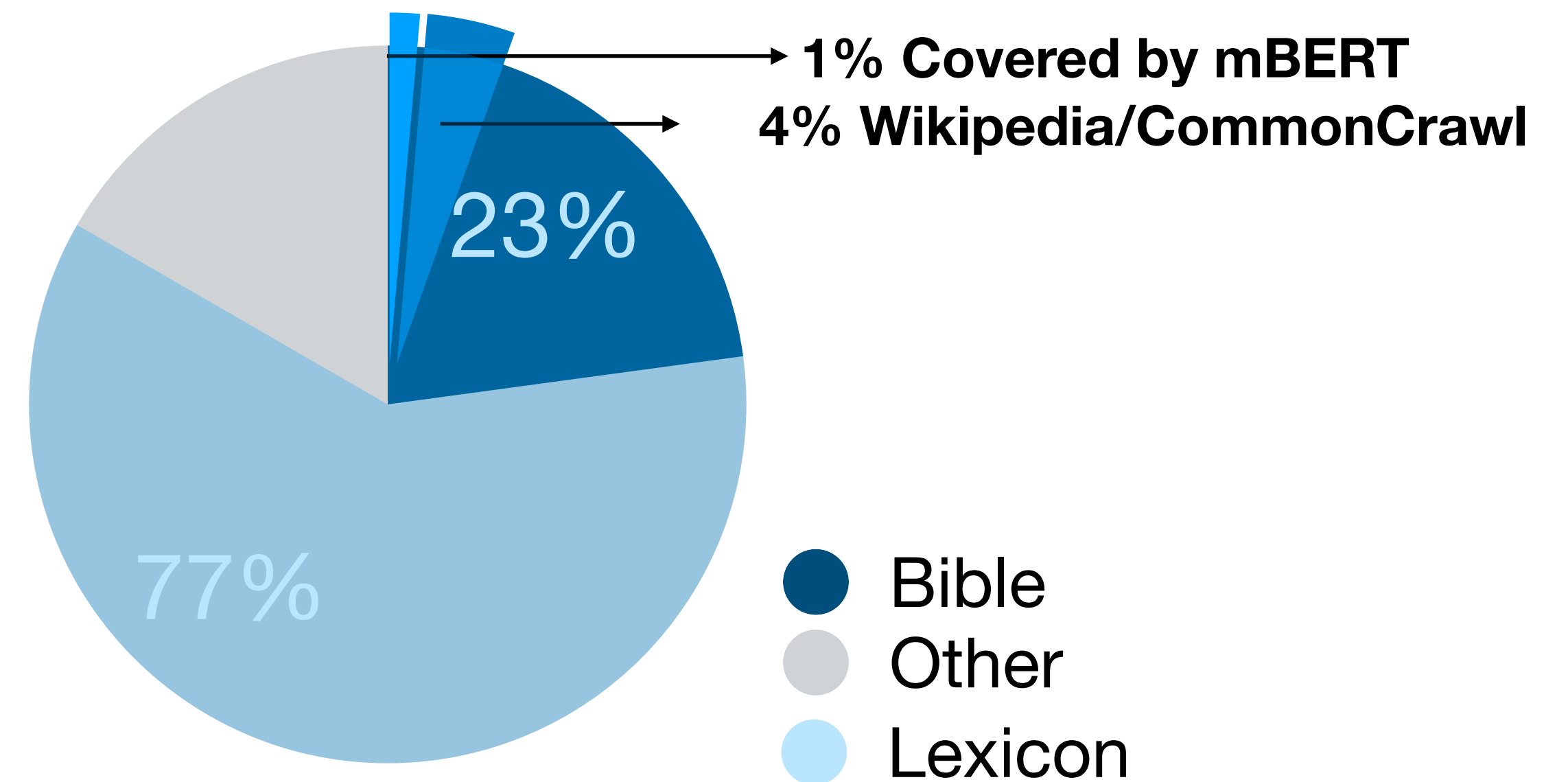
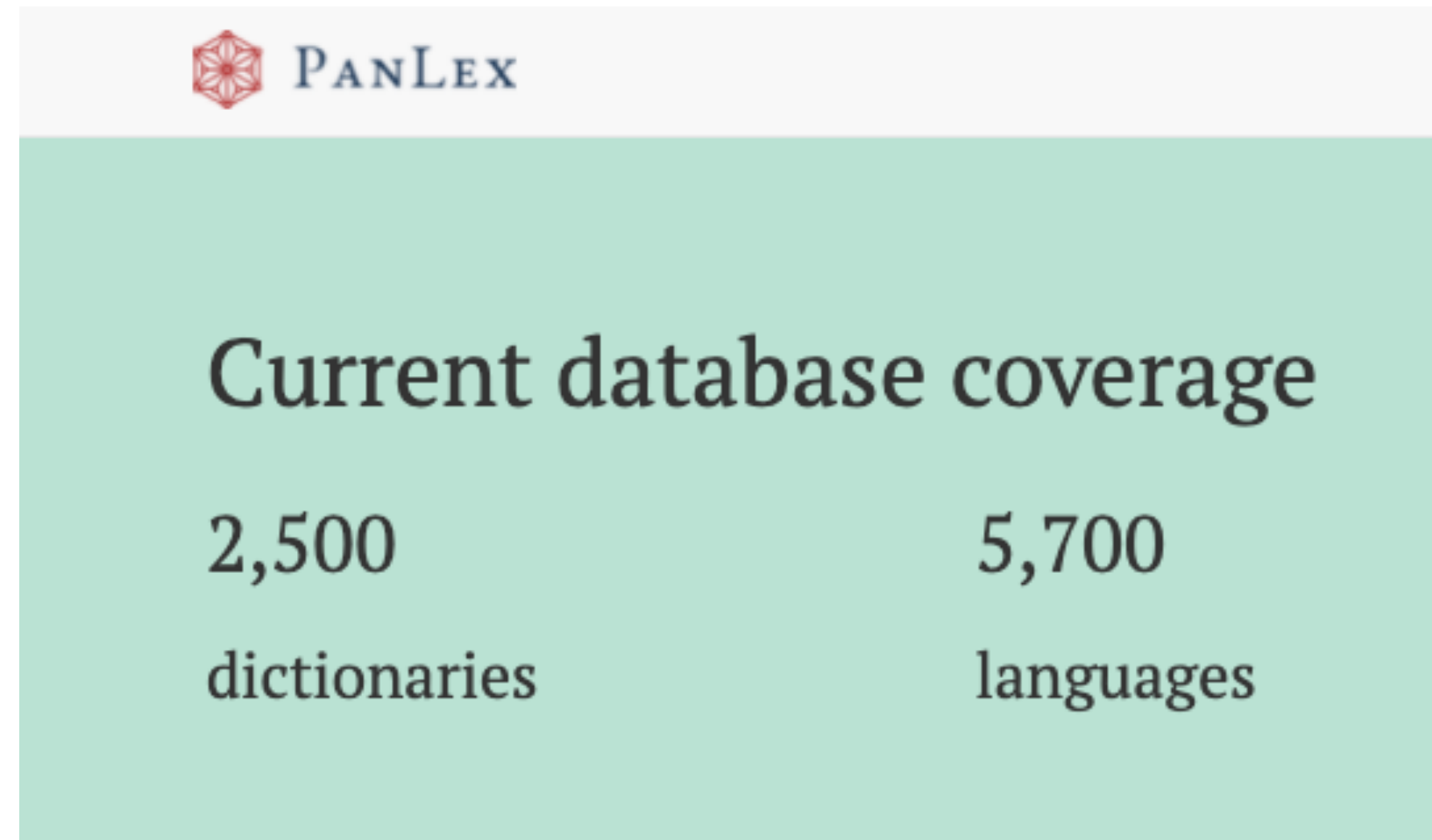
- Linguists have been documenting languages for years in formats such as lexicons

Alternative Data Source



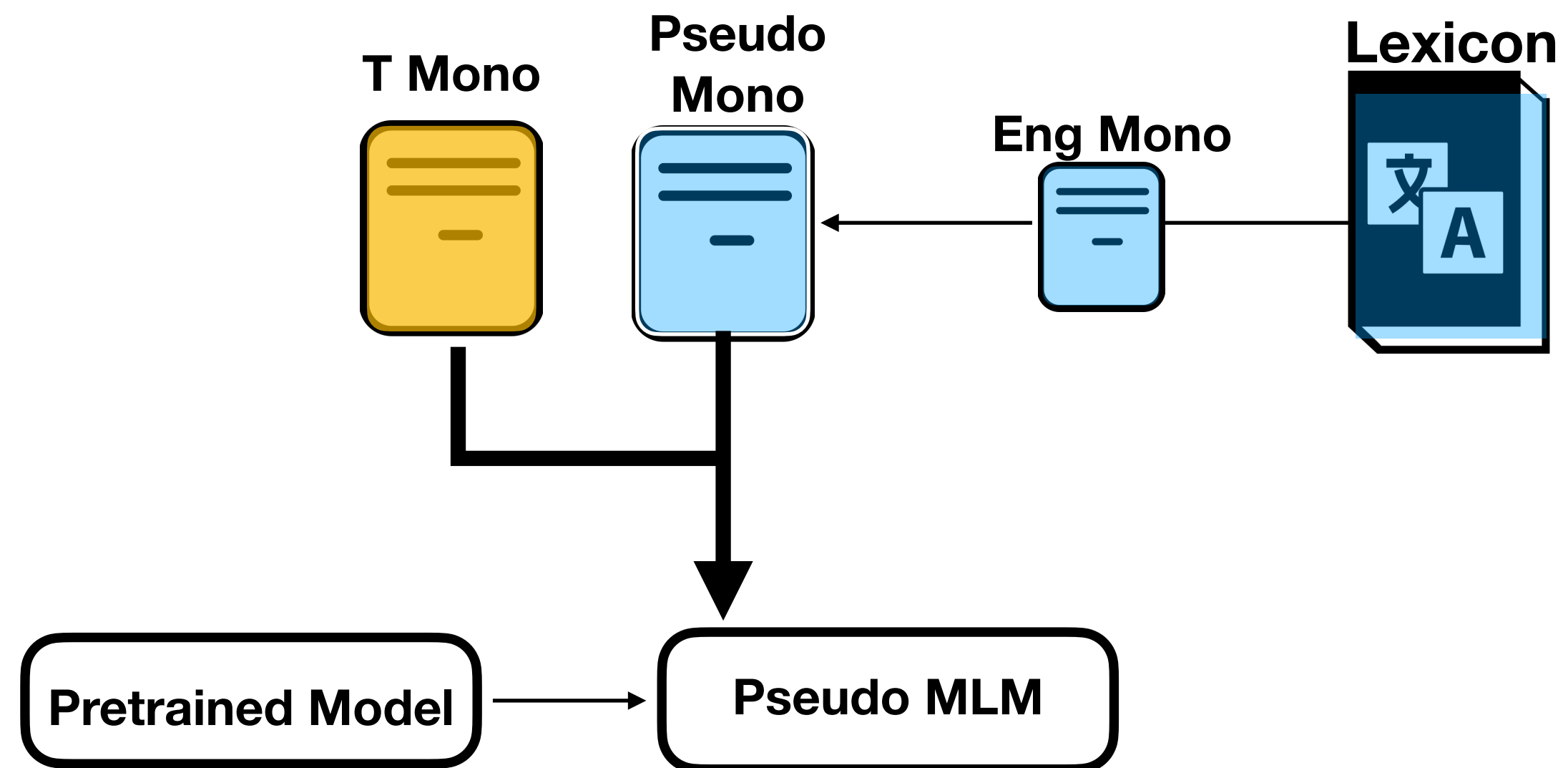
- Linguists have been documenting languages for years in formats such as lexicons
- PanLex: open-sourced database of lexicons with much better language coverage

Alternative Data Source

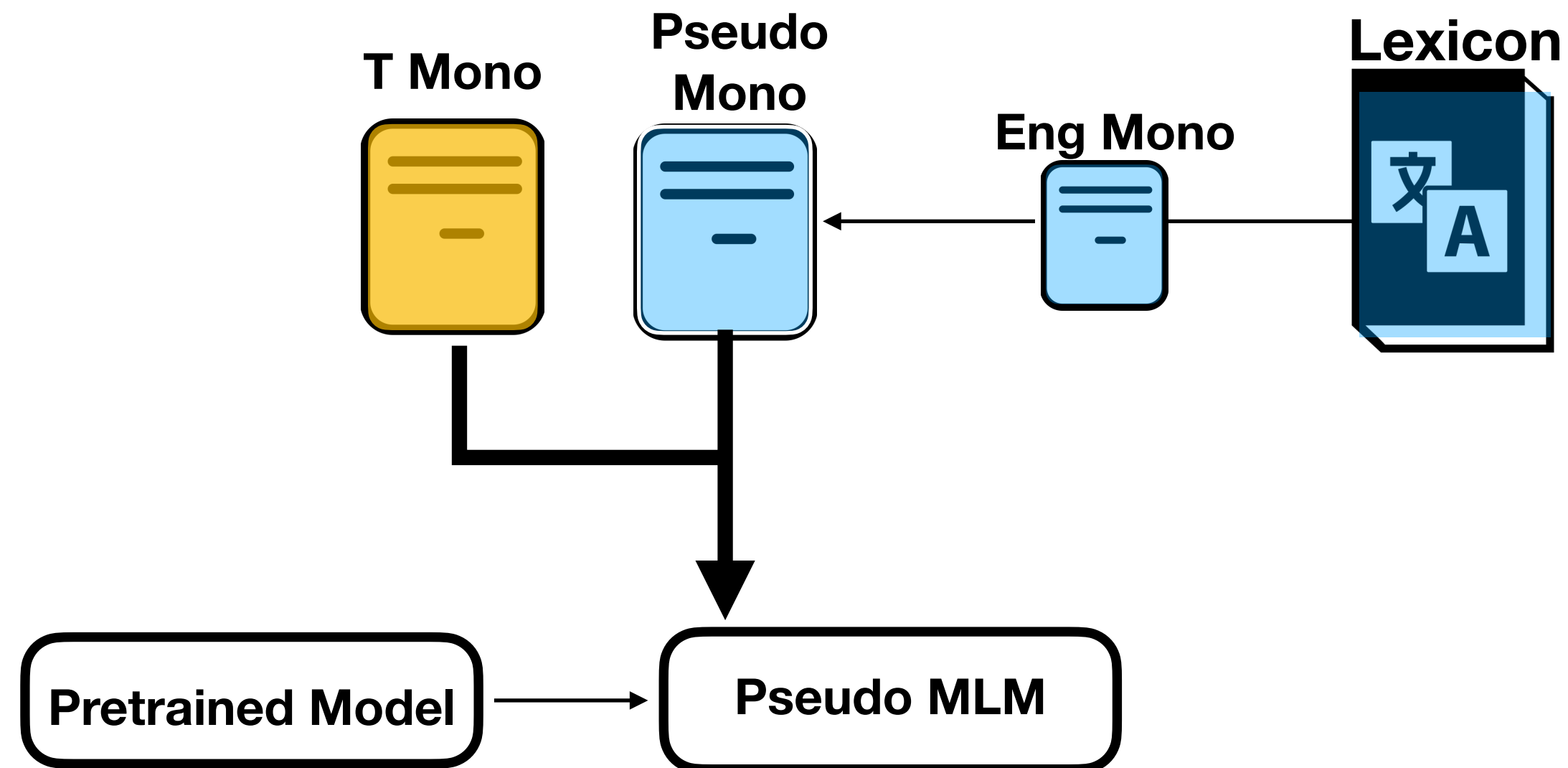


- Linguists have been documenting languages for years in formats such as lexicons
- PanLex: open-sourced database of lexicons with much better language coverage

Synthesizing Data Using Lexicons

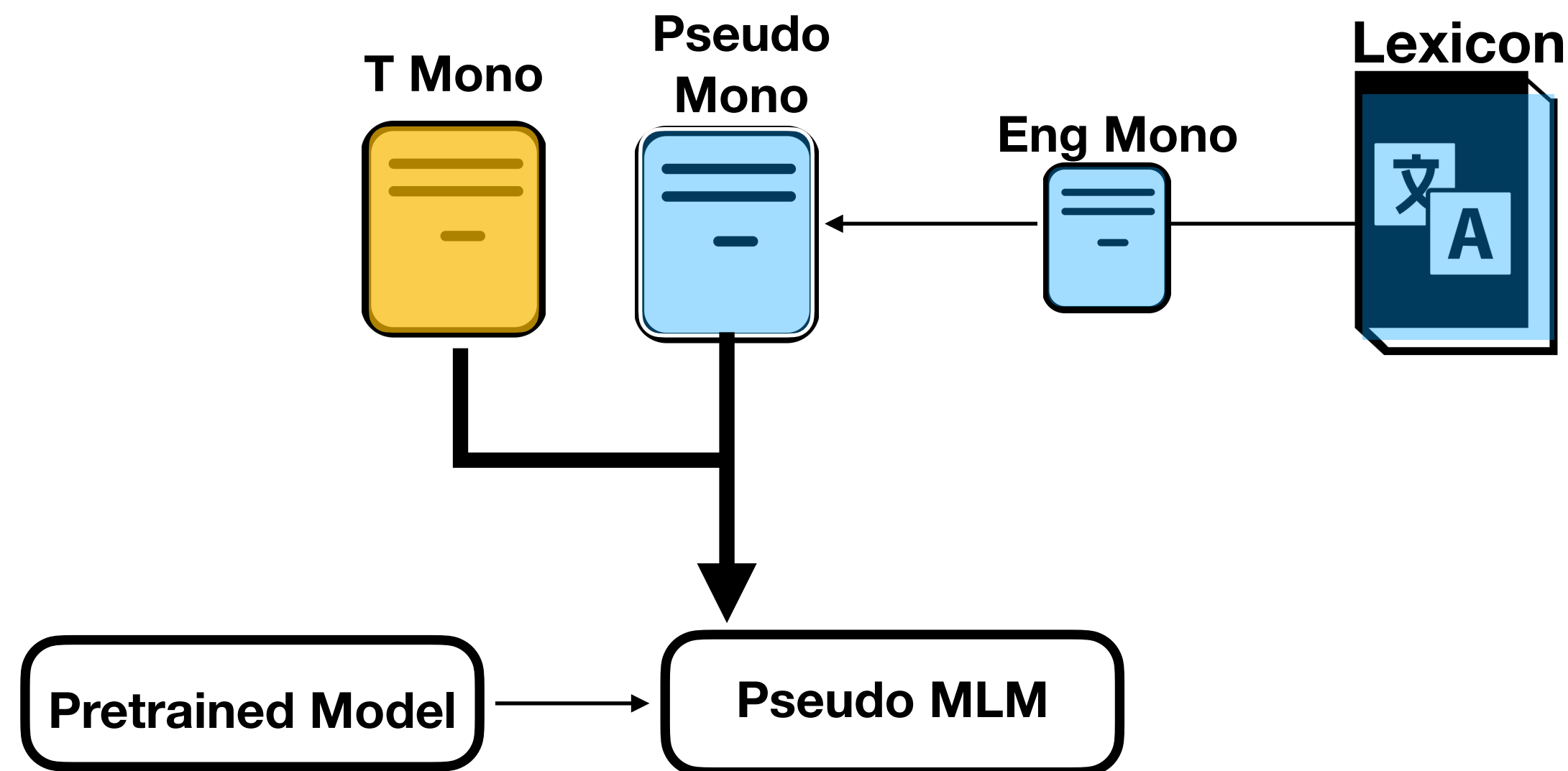


Synthesizing Data Using Lexicons



- Pseudo Mono Data: replace words in **English monolingual data** to its corresponding translation in the target language T

Synthesizing Data Using Lexicons

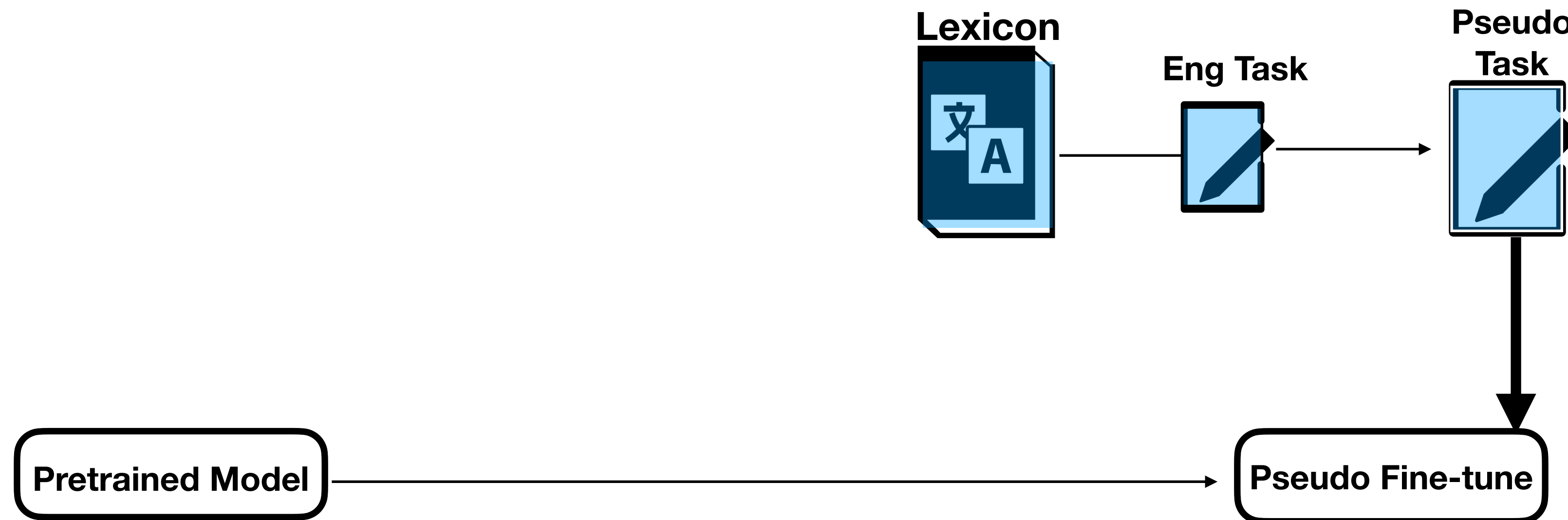


Eng Mono Anarchism calls for the abolition of the state , which it holds to be undesirable , unnecessary , and harmful .

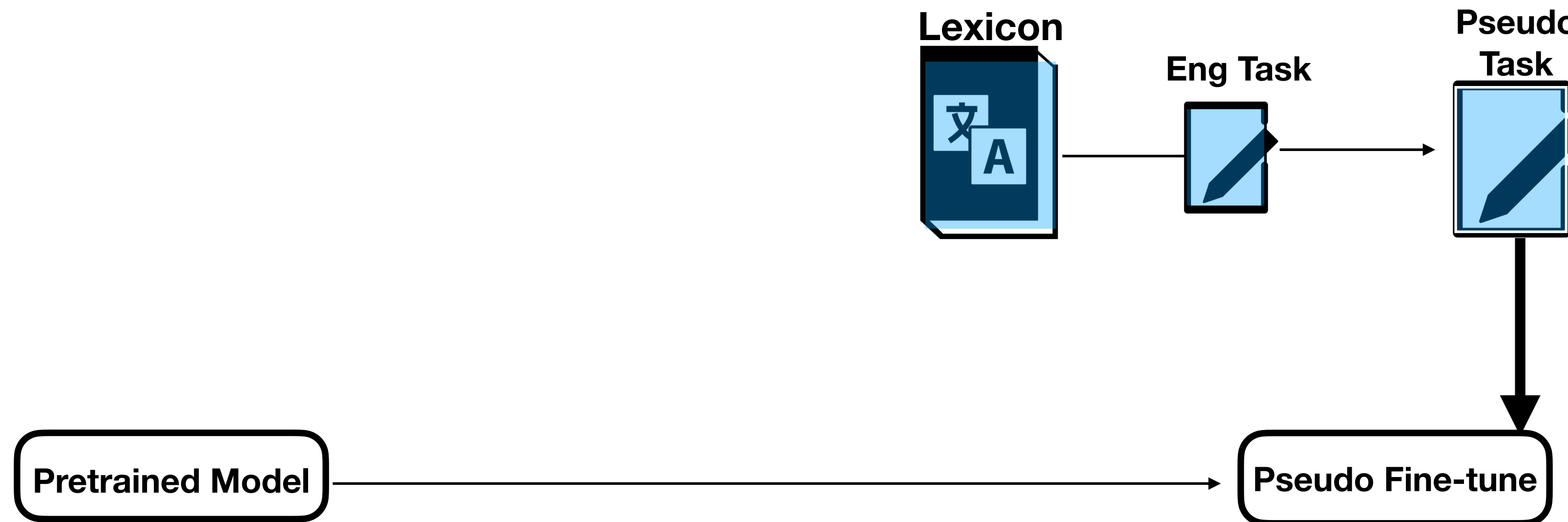
Pseudo Mono Anarchism calls *ghal il* abolition *ta' il stat* , *lima hi* holds *ghal tkun* undesirable , *bla bzonn* , *u* harmful .

- Pseudo Mono Data: replace words in **English monolingual data** to its corresponding translation in the target language T

Synthesizing Data Using Lexicons

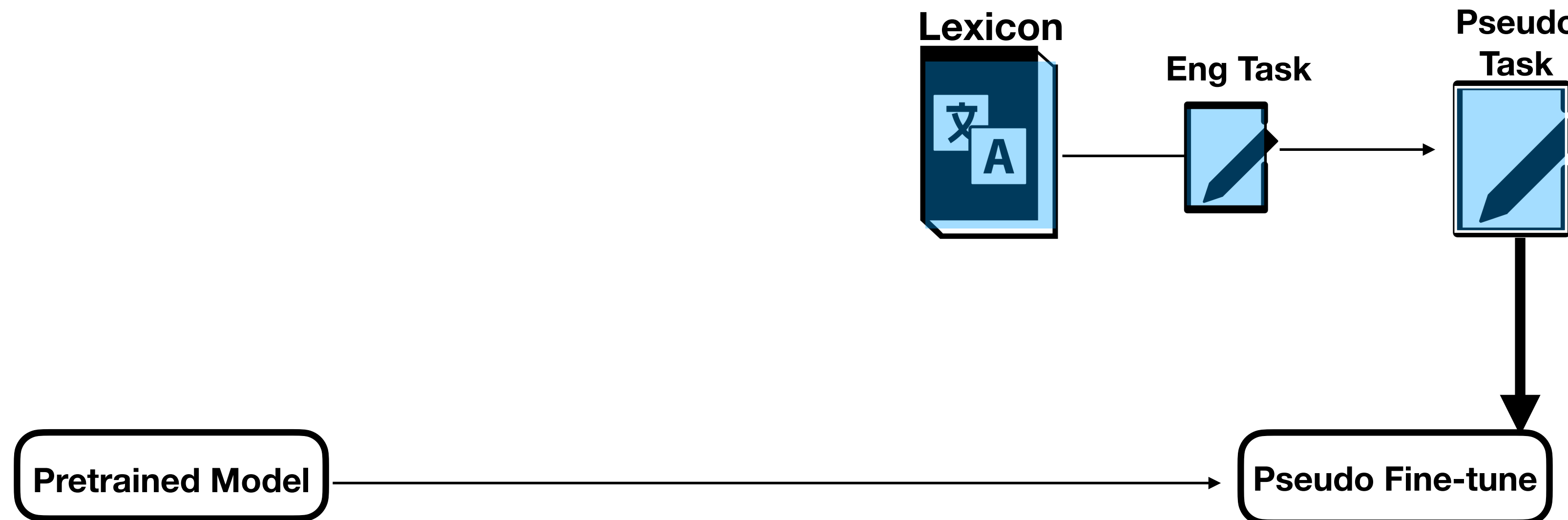


Synthesizing Data Using Lexicons



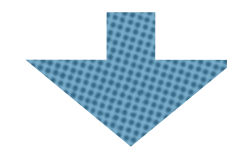
- Pseudo Task Data: replace words in **English task data** to its corresponding translation in the target language T

Synthesizing Data Using Lexicons



Eng Task

I suspect the streets of Baghdad will look as if a war is looming this week .
 PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

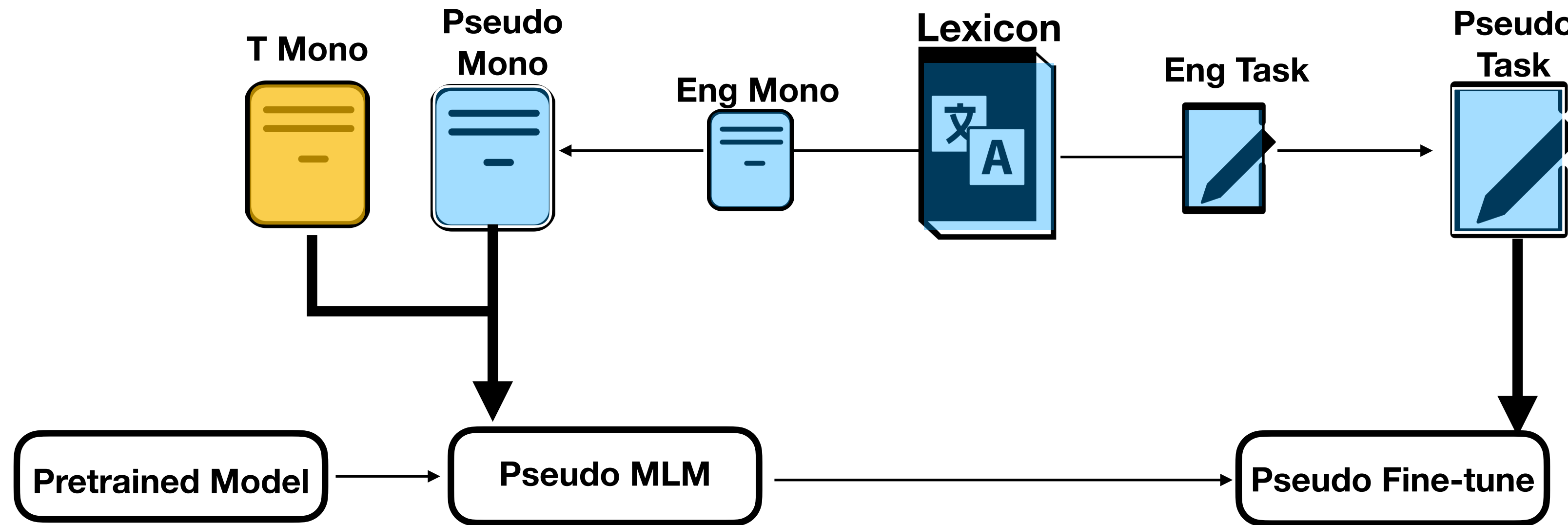


Pseudo Task

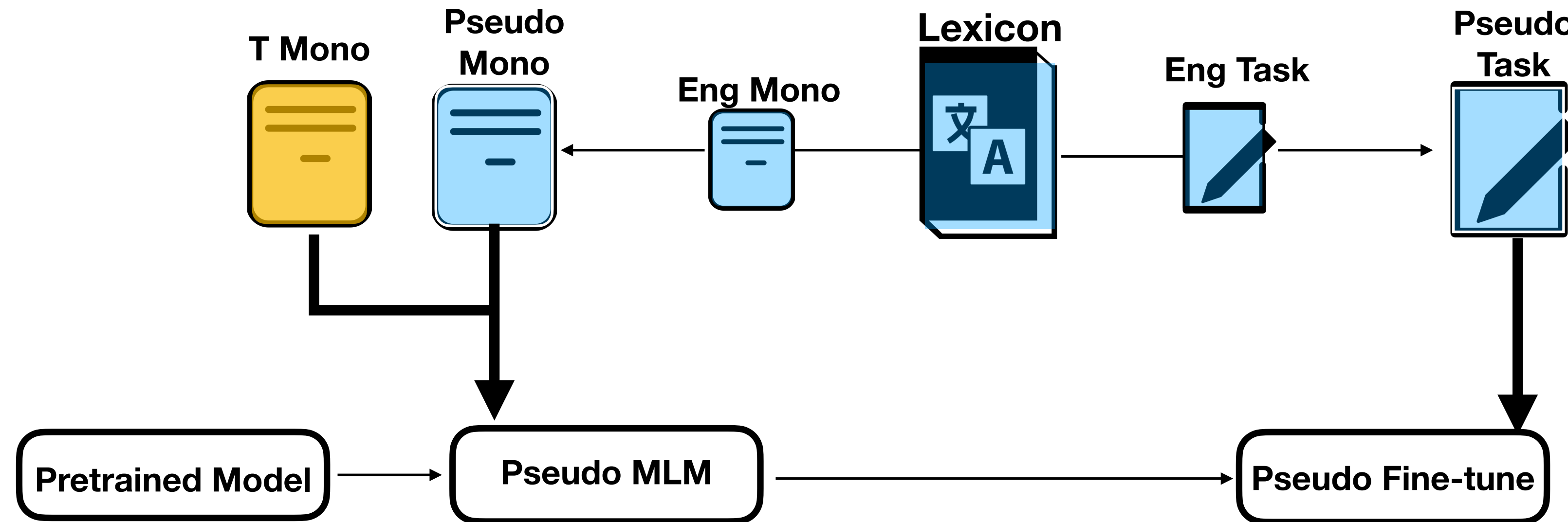
jien iddubita il streets ta' Bagdad xewqa hares kif jekk a gwerra is looming dan gimgha .
 PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

- Pseudo Task Data: replace words in **English task data** to its corresponding translation in the target language T

Synthesizing Data Using Lexicons



Synthesizing Data Using Lexicons



- Use either pseudo MLM or Pseudo Fine-tune, or both

Experiments

Experiments

- Model: mBERT

Experiments

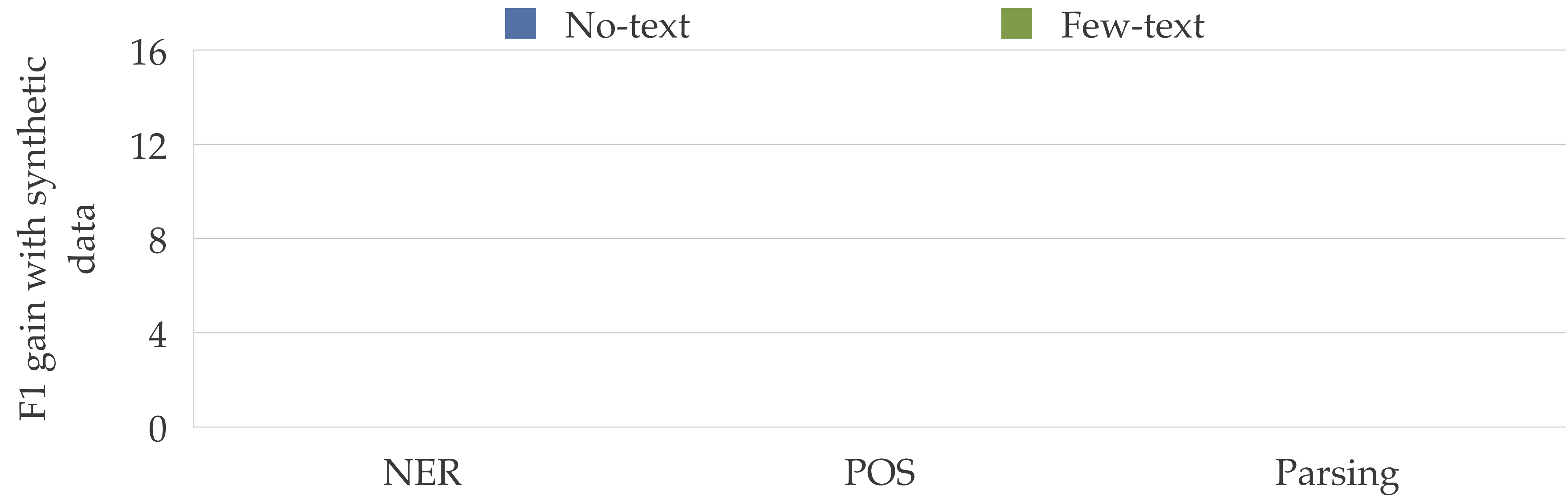
- Model: mBERT
- Tasks:
 - NER
 - POS tagging
 - Dependency Parsing

Experiments

- Model: mBERT
- Tasks:
 - NER
 - POS tagging
 - Dependency Parsing
- Languages: 19 languages not covered by mBERT pretraining

Results

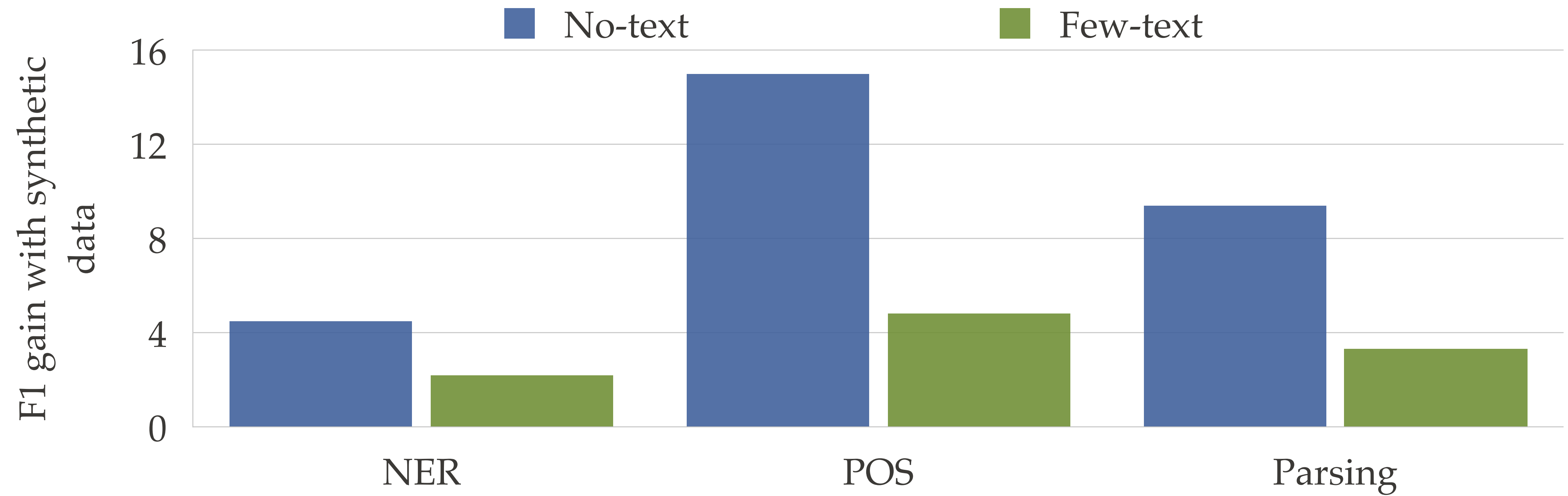
Results



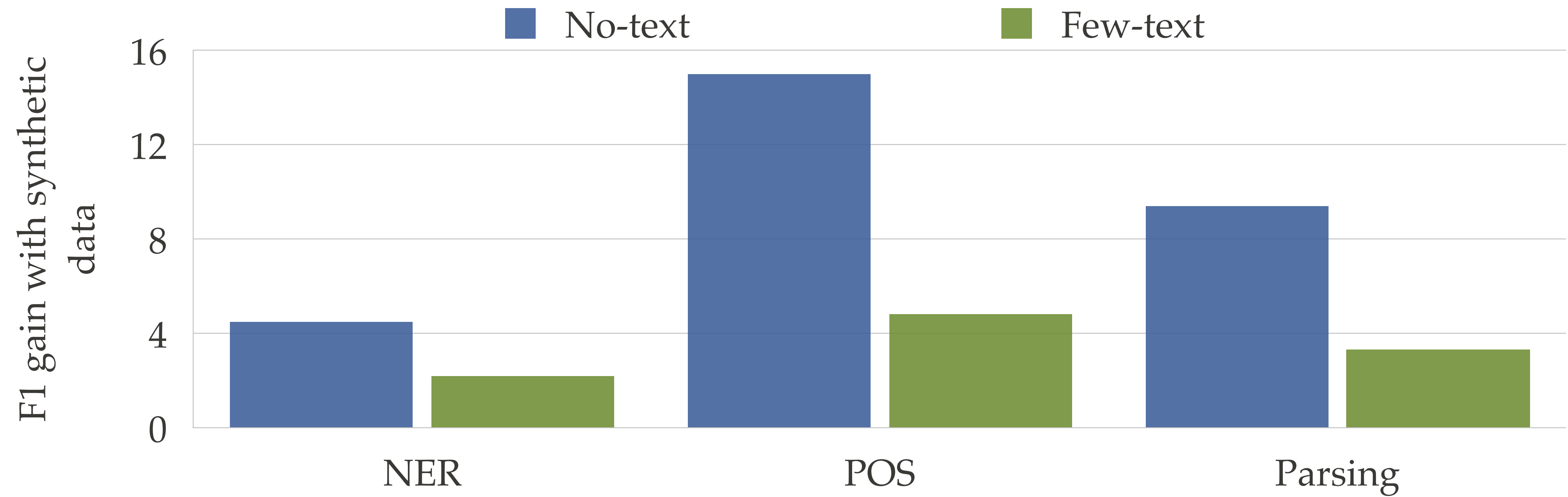
Results



Results



Results



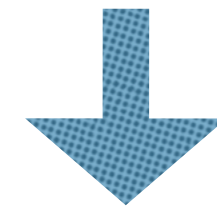
- Using synthetic data leads to significant improvements for both no-text and few-text setting

Label Noise

Eng Task

I suspect the streets of Baghdad **will** look as if a war is looming this week .

PRON VERB DET NOUN ADP PROPN **AUX** VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

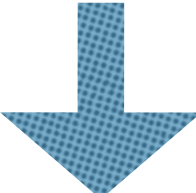


Pseudo Task

jien iddubita il streets *ta'* Bagdad **xewqa** hares kif jekk a gwerra is looming dan gimgha .

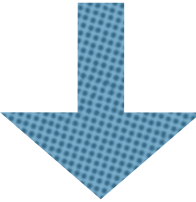
PRON VERB DET NOUN ADP PROPN **AUX** VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

Label Noise

Eng Task	<p>I suspect the streets of Baghdad will look as if a war is looming this week .</p> <p>PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT</p>
	
Pseudo Task	<p>jien iddubita il streets ta' Bagdad xewqa hares kif jekk a gwerra is looming dan gimgha .</p> <p>PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT</p>

- “**xewqa**” is a noun meaning “desire,will”

Label Noise

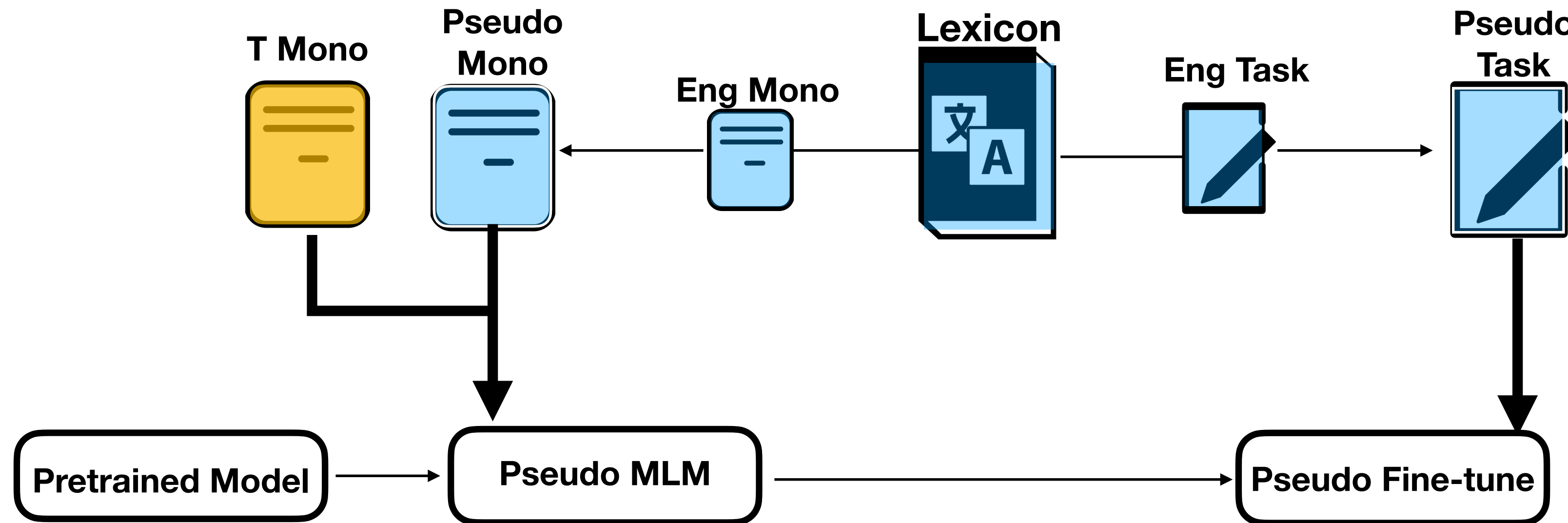
Eng Task	I suspect the streets of Baghdad will look as if a war is looming this week . PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT
	
Pseudo Task	jien iddubita il streets ta' Bagdad xewqa hares kif jekk a gwerra is looming dan gimgha . PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

- “**xewqa**” is a noun meaning “desire,will”
- But the original English POS tag is inconsistent with the replaced word

Label Noise

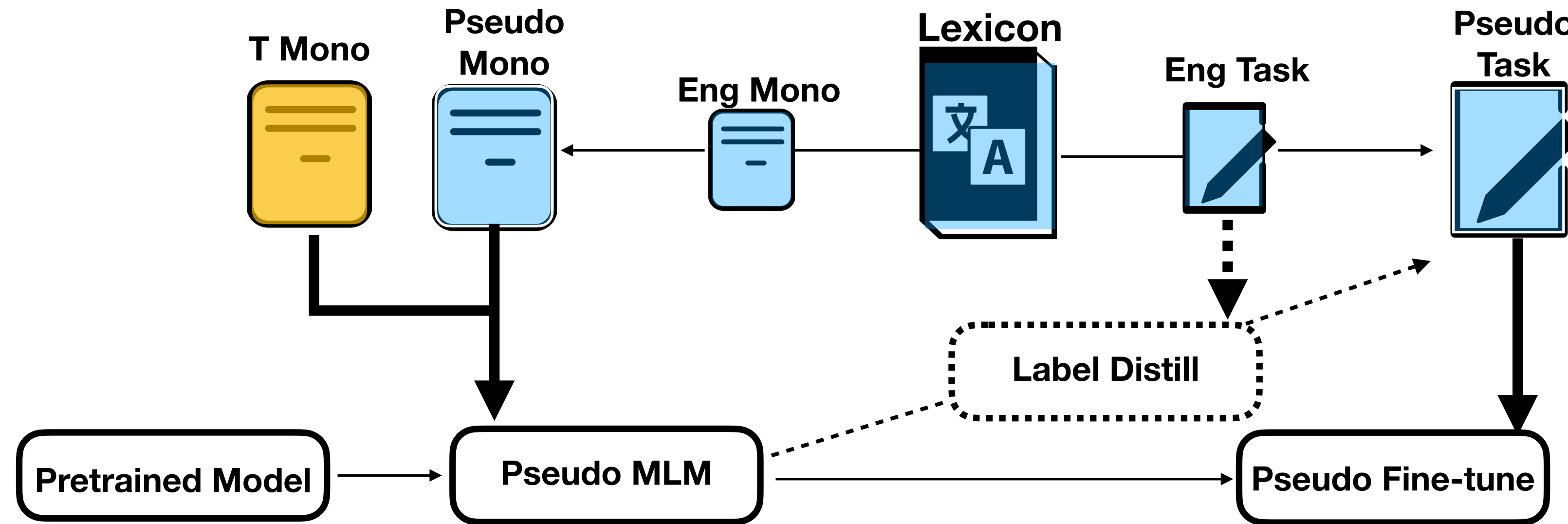
- Use the fine-tuned model to “correct” the labels for the Pseudo task data

Label Noise



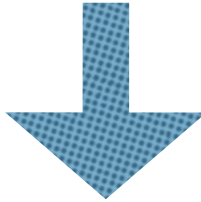
- Use the fine-tuned model to “correct” the labels for the Pseudo task data

Label Noise



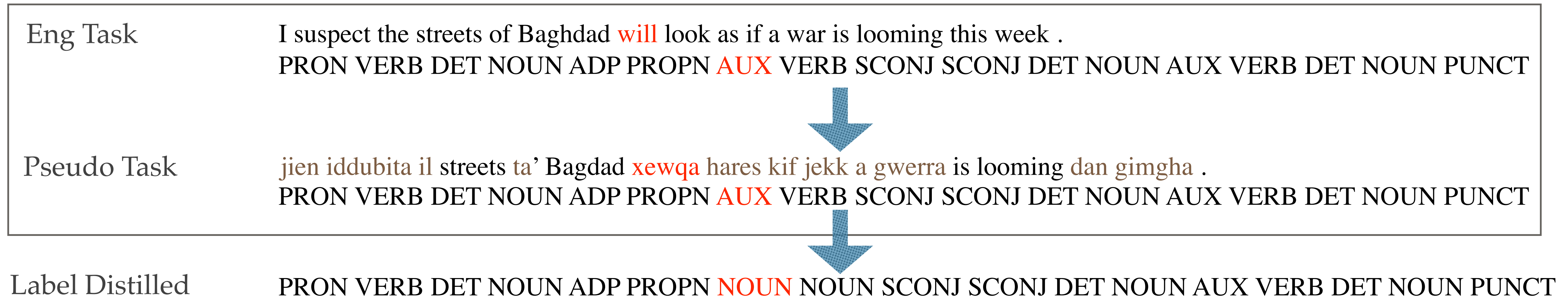
- Use the fine-tuned model to “correct” the labels for the Pseudo task data

Label Noise

Eng Task	I suspect the streets of Baghdad will look as if a war is looming this week . PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT
	
Pseudo Task	jien iddubita il streets ta' Bagdad xewqa hares kif jekk a gwerra is looming dan gimgha . PRON VERB DET NOUN ADP PROPN AUX VERB SCONJ SCONJ DET NOUN AUX VERB DET NOUN PUNCT

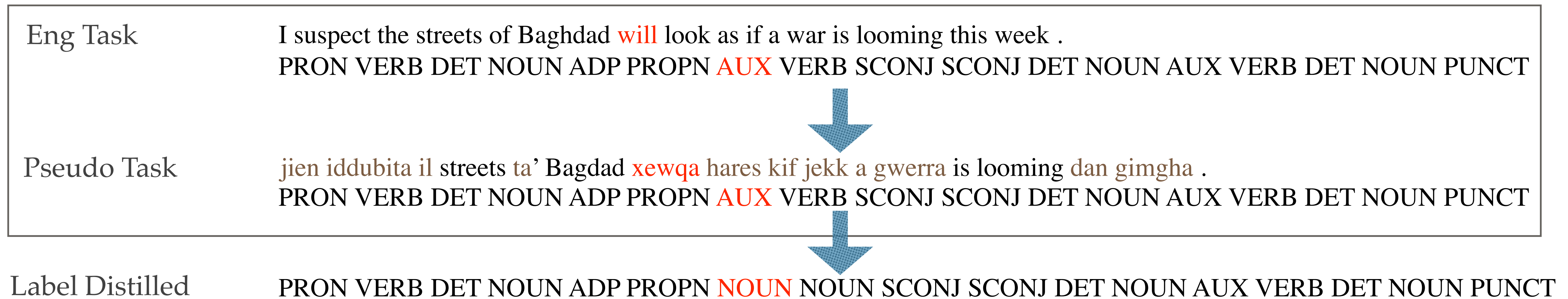
- “**xewqa**” is a noun meaning “desire,will”

Label Noise



- “**xewqa**” is a noun meaning “desire,will”

Label Noise



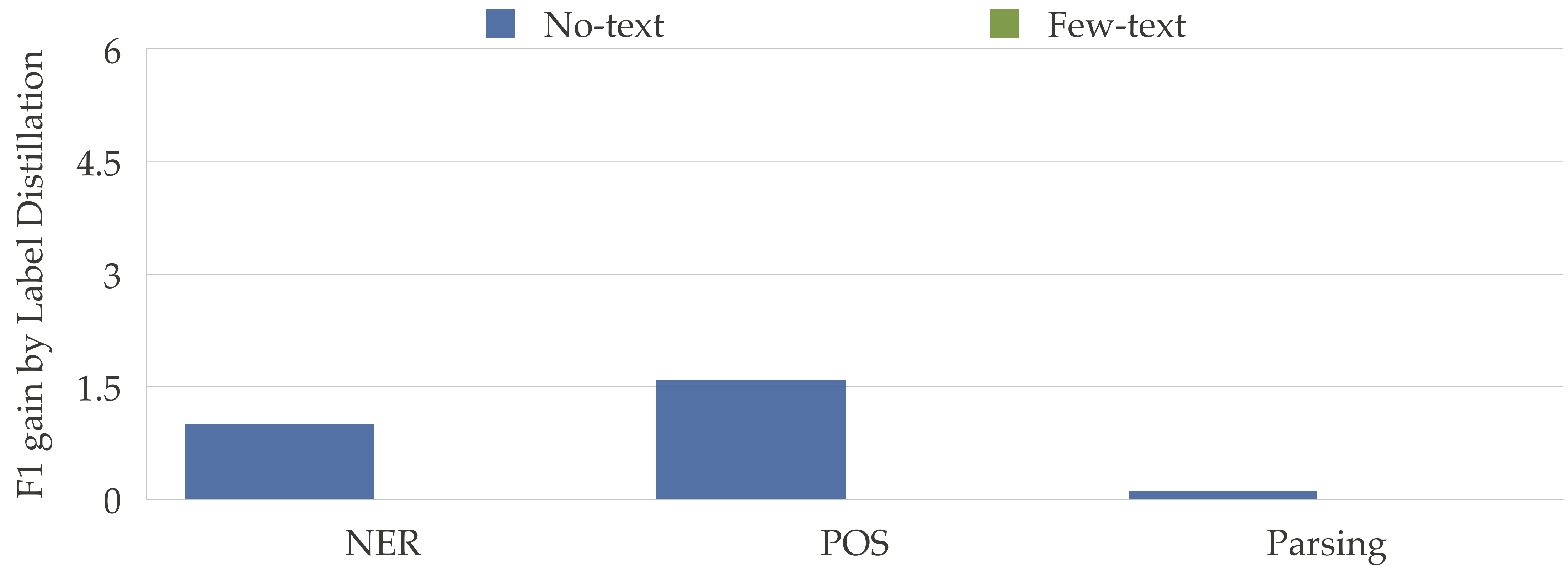
- “**xewqa**” is a noun meaning “desire,will”
- The model is able to assign the correct label of noun

Label Noise

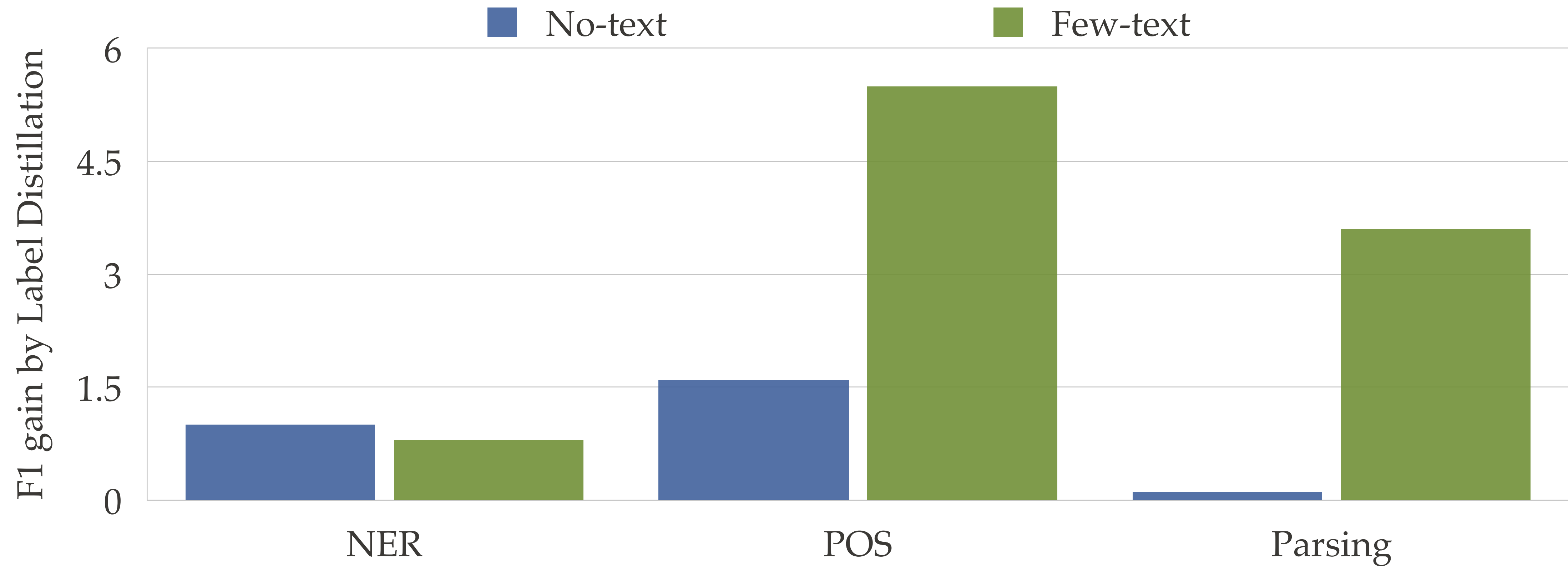
Label Noise



Label Noise



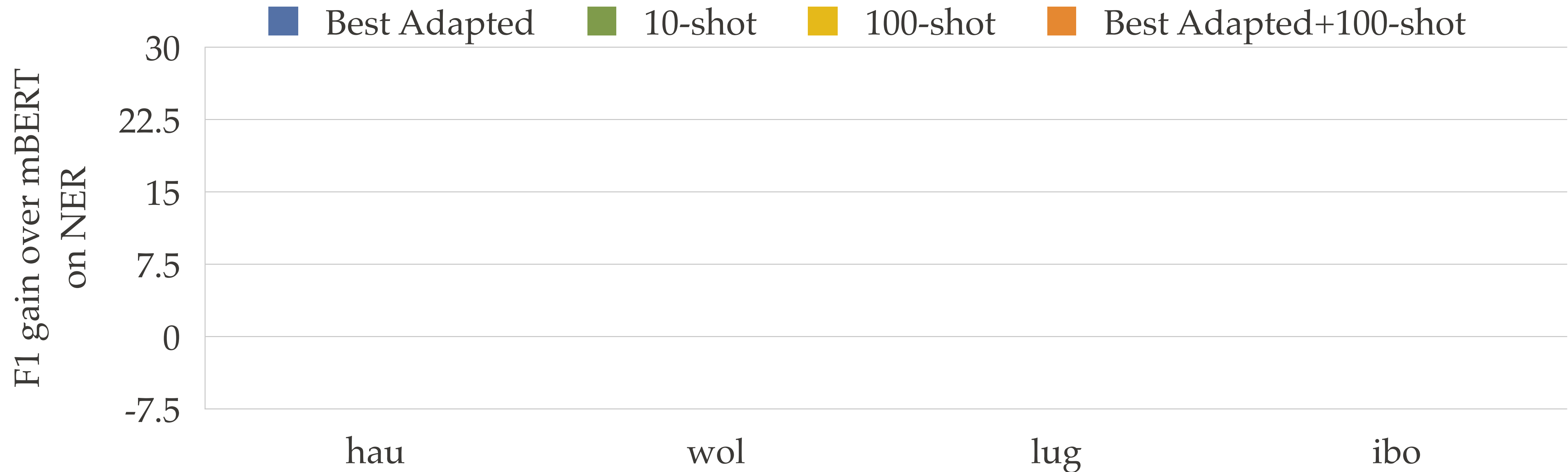
Label Noise



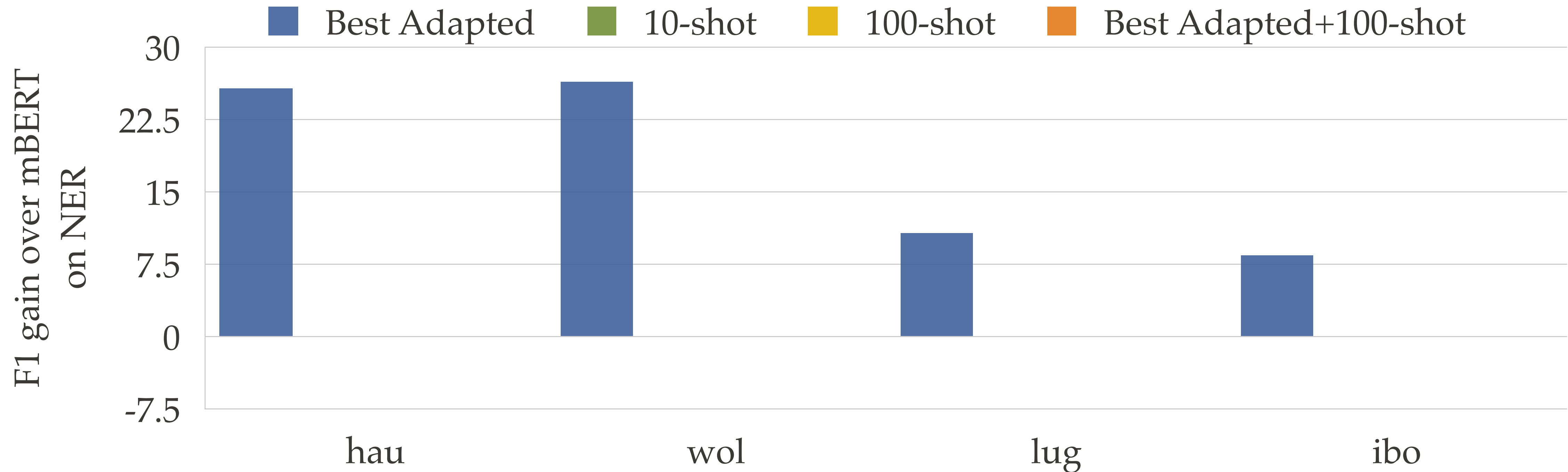
- Label Distillation is especially helpful for syntactic tasks

Comparison to Few-shot Learning

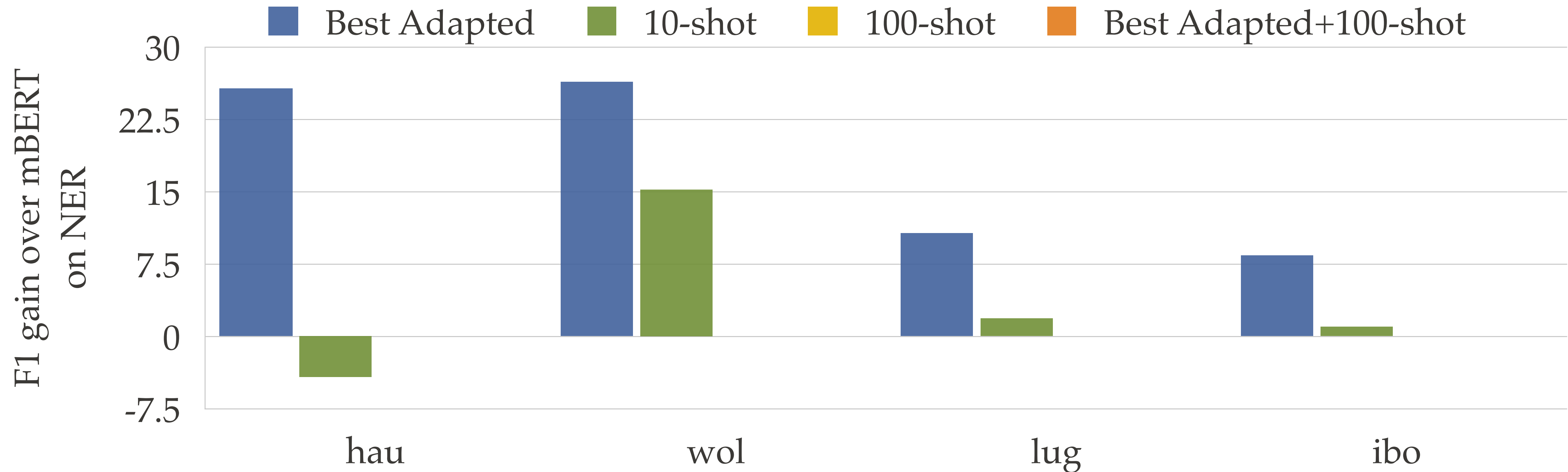
Comparison to Few-shot Learning



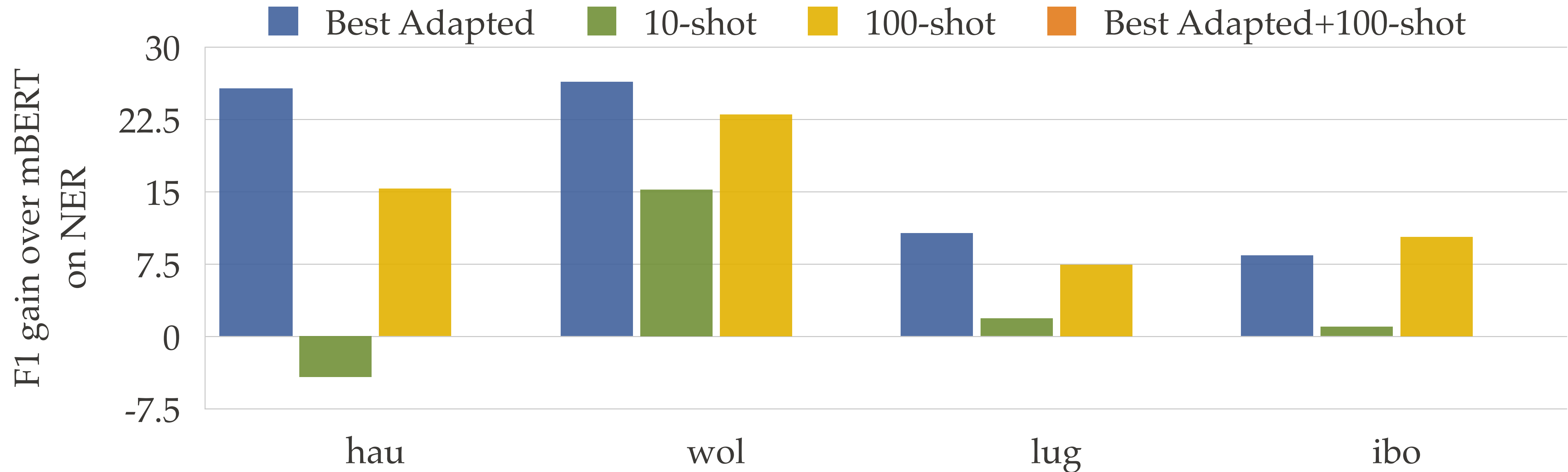
Comparison to Few-shot Learning



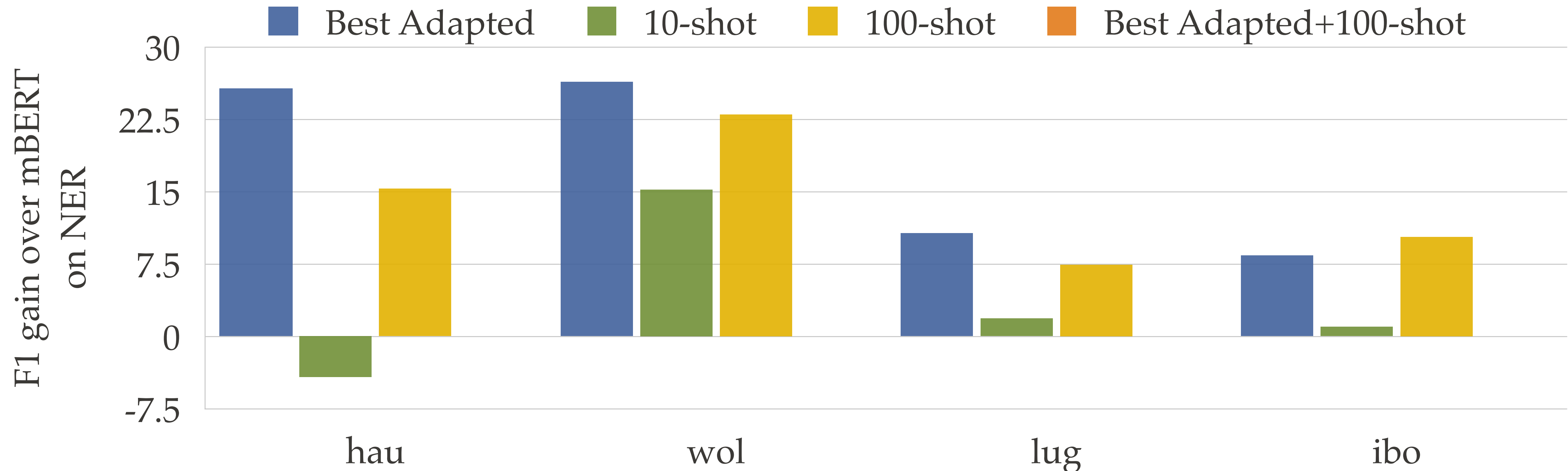
Comparison to Few-shot Learning



Comparison to Few-shot Learning

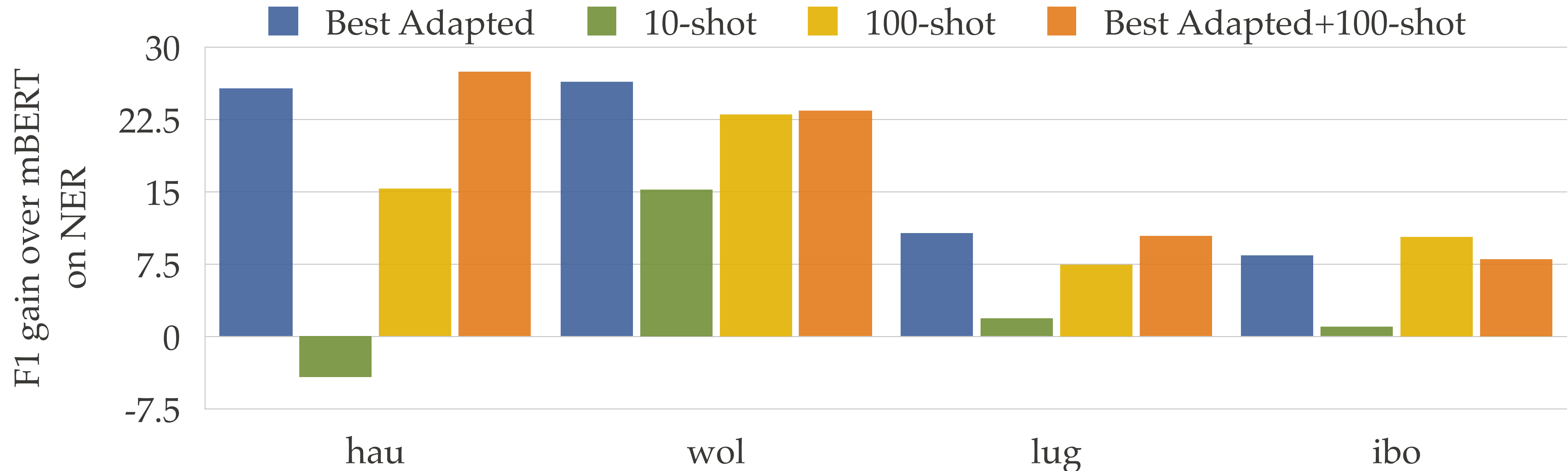


Comparison to Few-shot Learning



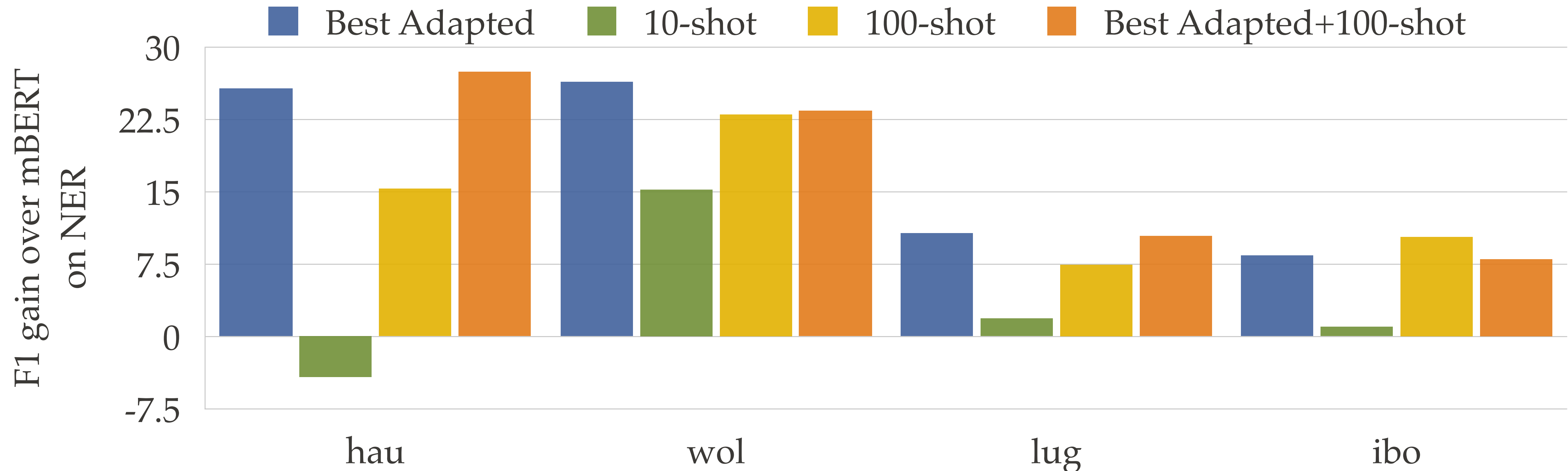
- Few-shot learning needs more annotated data for languages with limited text

Comparison to Few-shot Learning



- Few-shot learning needs more annotated data for languages with limited text

Comparison to Few-shot Learning

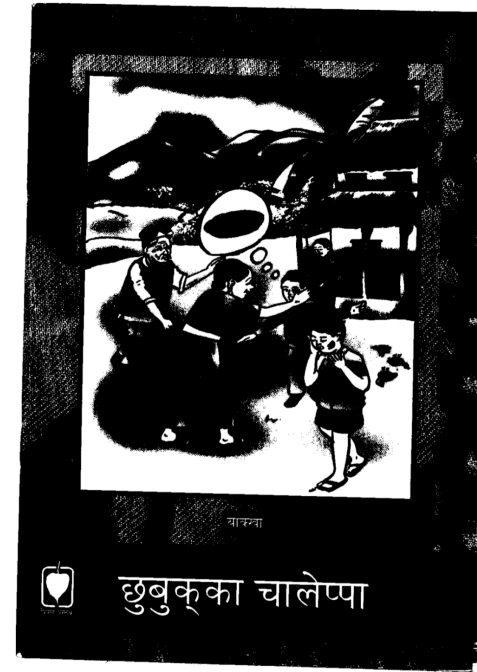
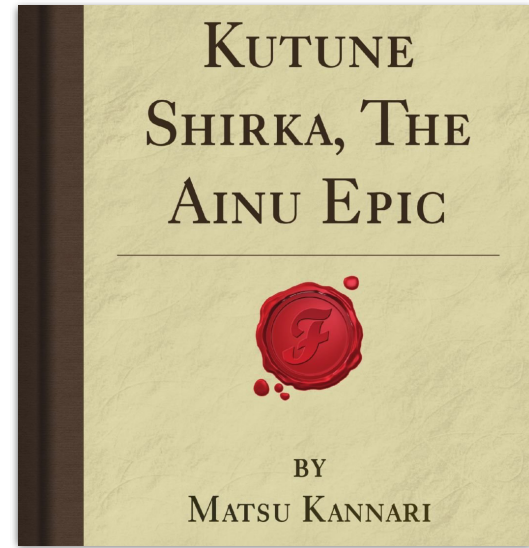


- Few-shot learning needs more annotated data for languages with limited text
- Combining adaptation and few-shot doesn't bring consistent improvements

Conclusion

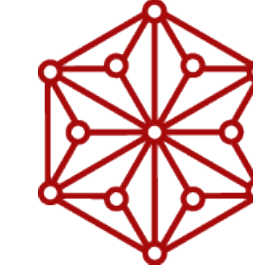
Conclusion

Conclusion



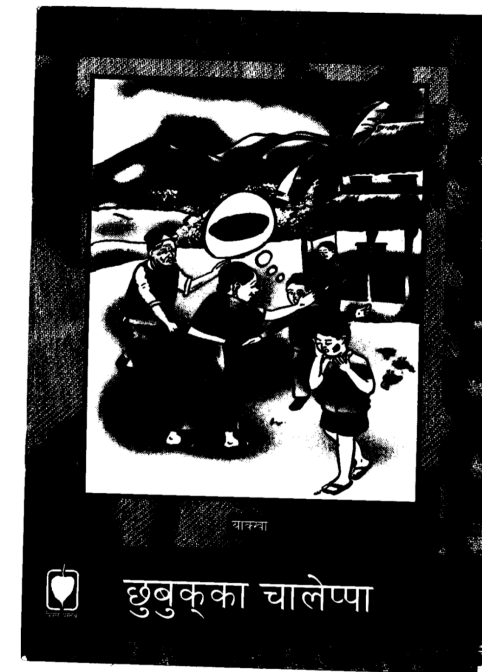
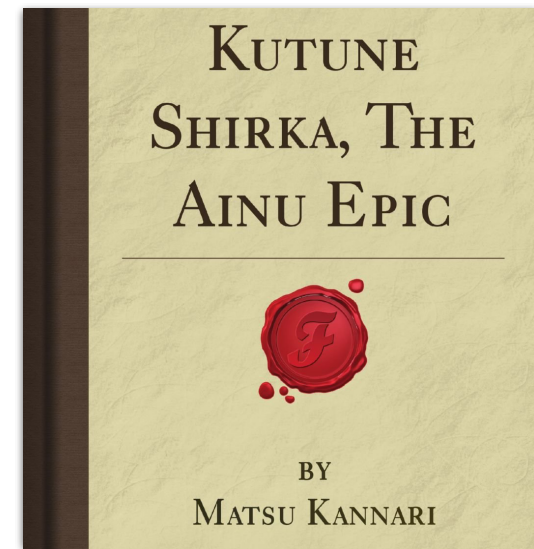
naúuzi wa ná'túkwa
mí kittóonaxipilikáan
naúuzi wa ná'túkwa

masēxa ts!ēx·inaxs
 Wä, g·il^εmēsē ^εwilg
 laē äx^εēdxēs gālay
 ts!ēx·mēsē. Wä,



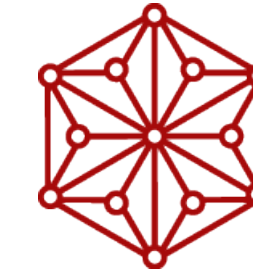
PANLEX

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

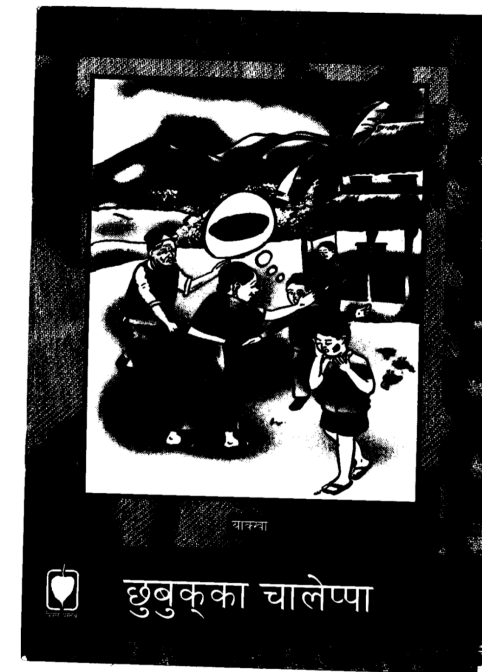
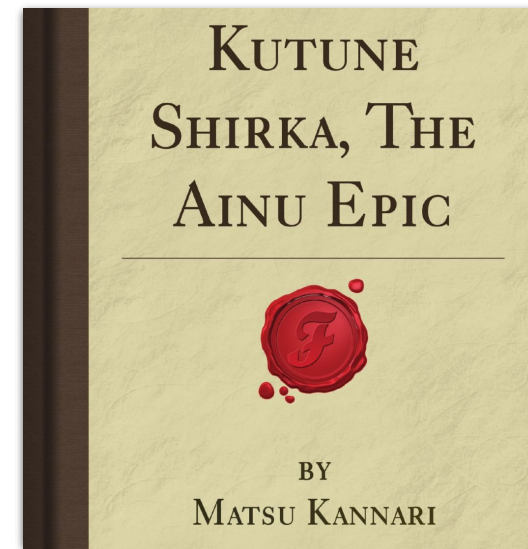
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

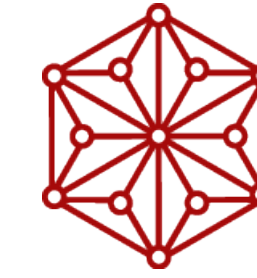
- Methods to **unlock new resources** for human or machine use in under-resourced languages

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

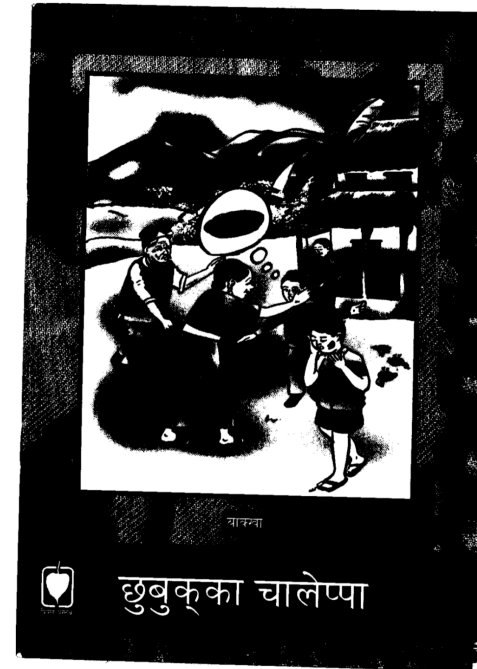
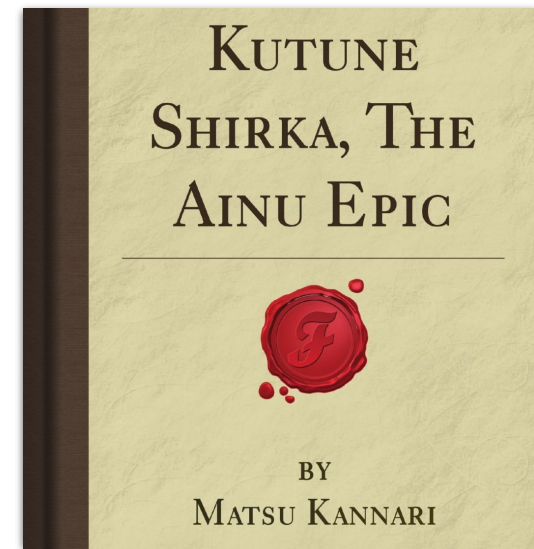
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?

Conclusion



*naúzi wa ná'túkwá
mí kittóonaxipilikáan
naúzi wa ná'túkwá*

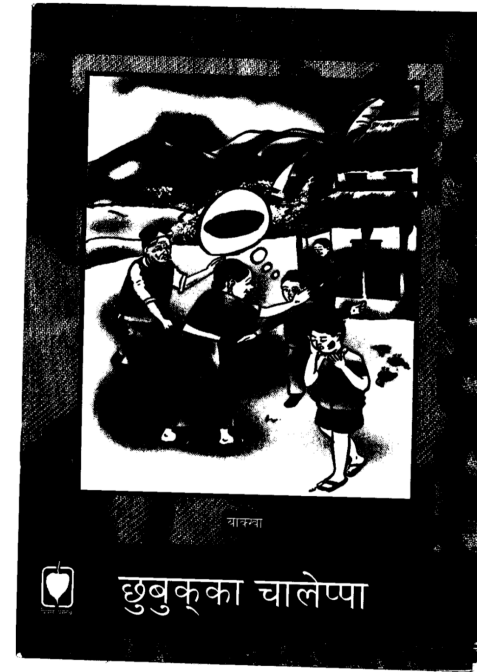
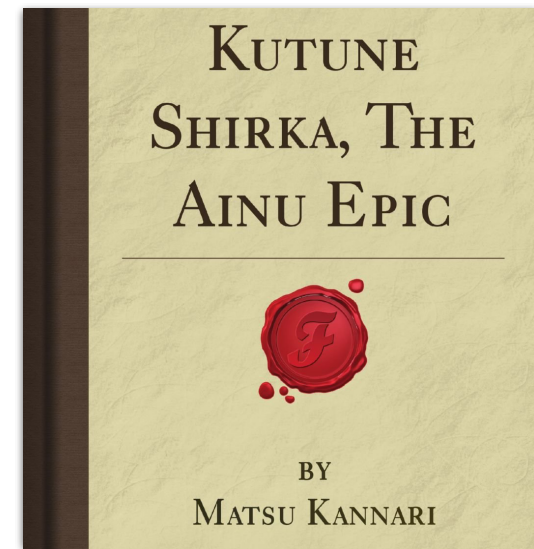
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

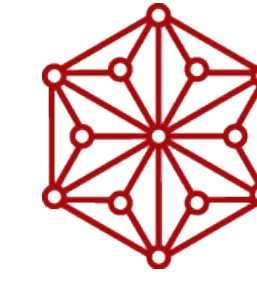
- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?
- Should we **put some linguistics in the models?**

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

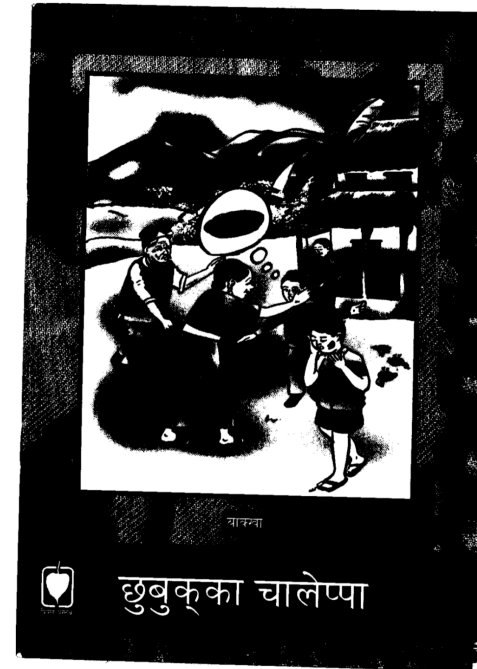
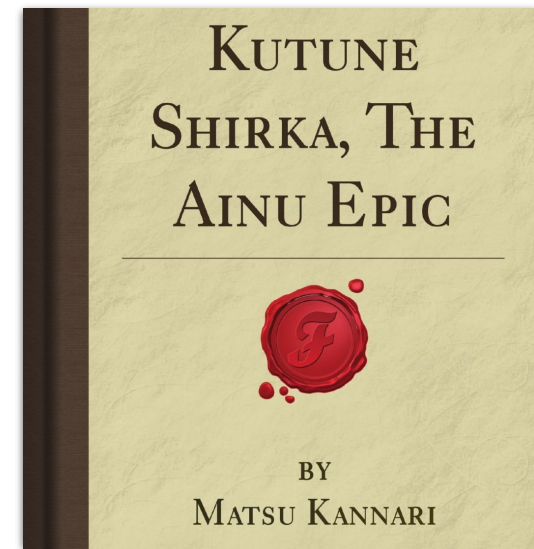
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

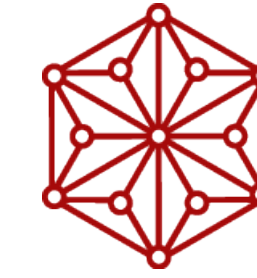
- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?
- Should we **put some linguistics in the models?**
 - Morphologically aware soft constraints for OCR?

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

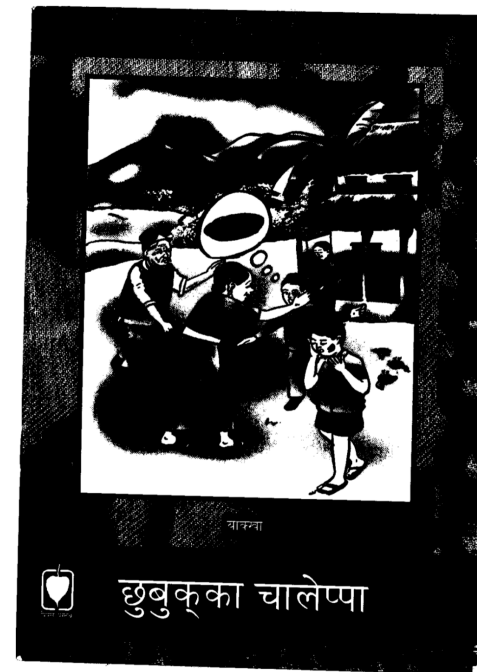
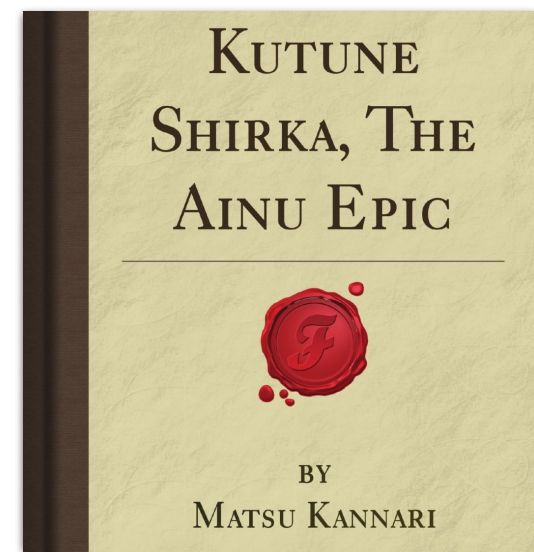
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

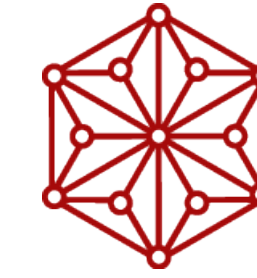
- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?
- Should we **put some linguistics in the models?**
 - Morphologically aware soft constraints for OCR?
 - Morphologically/syntactically aware data synthesis using lexicons?

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

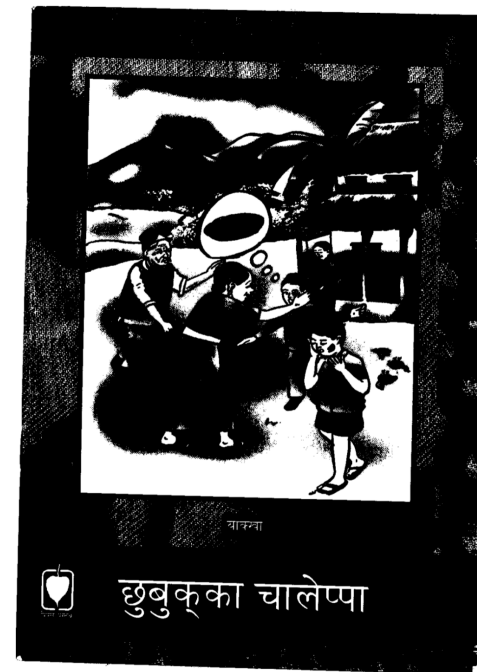
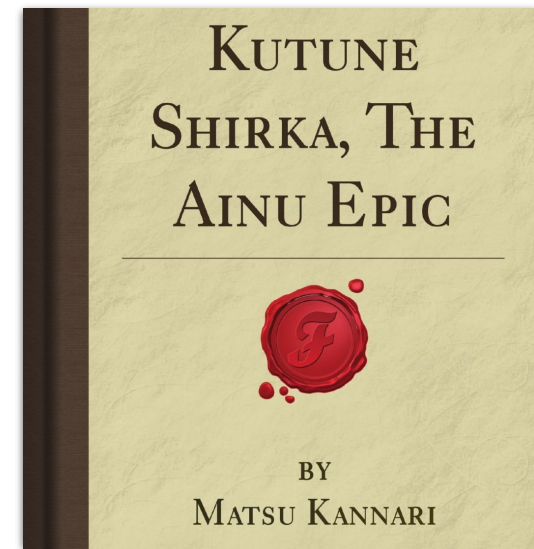
masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?
- Should we **put some linguistics in the models?**
 - Morphologically aware soft constraints for OCR?
 - Morphologically/syntactically aware data synthesis using lexicons?
- Should we **use the models in language learning or linguistics?**

Conclusion



*naúzi wa ná'túkwa
mí kittóonaxipilikáan
naúzi wa ná'túkwa*

masēxa ts!ēx·inaxs
Wä, g·il^εmēsē ^εwilg
laē äx^εēdxēs gālay
ts!ēx·mēsē. Wä,



PANLEX

- Methods to **unlock new resources** for human or machine use in under-resourced languages
- What's next?
- Should we **put some linguistics in the models?**
 - Morphologically aware soft constraints for OCR?
 - Morphologically/syntactically aware data synthesis using lexicons?
- Should we **use the models in language learning or linguistics?**
 - Large-scale extraction of text or inter-linear glosses for use in developing language materials?