

Unsupervised Code-Switching for Multilingual Historical Document Transcription

Dan Garrette

UT-Austin Computer Science

Hannah Alpert-Abrams

UT-Austin Comparative Literature

Taylor Berg-Kirkpatrick

UC Berkeley Computer Science

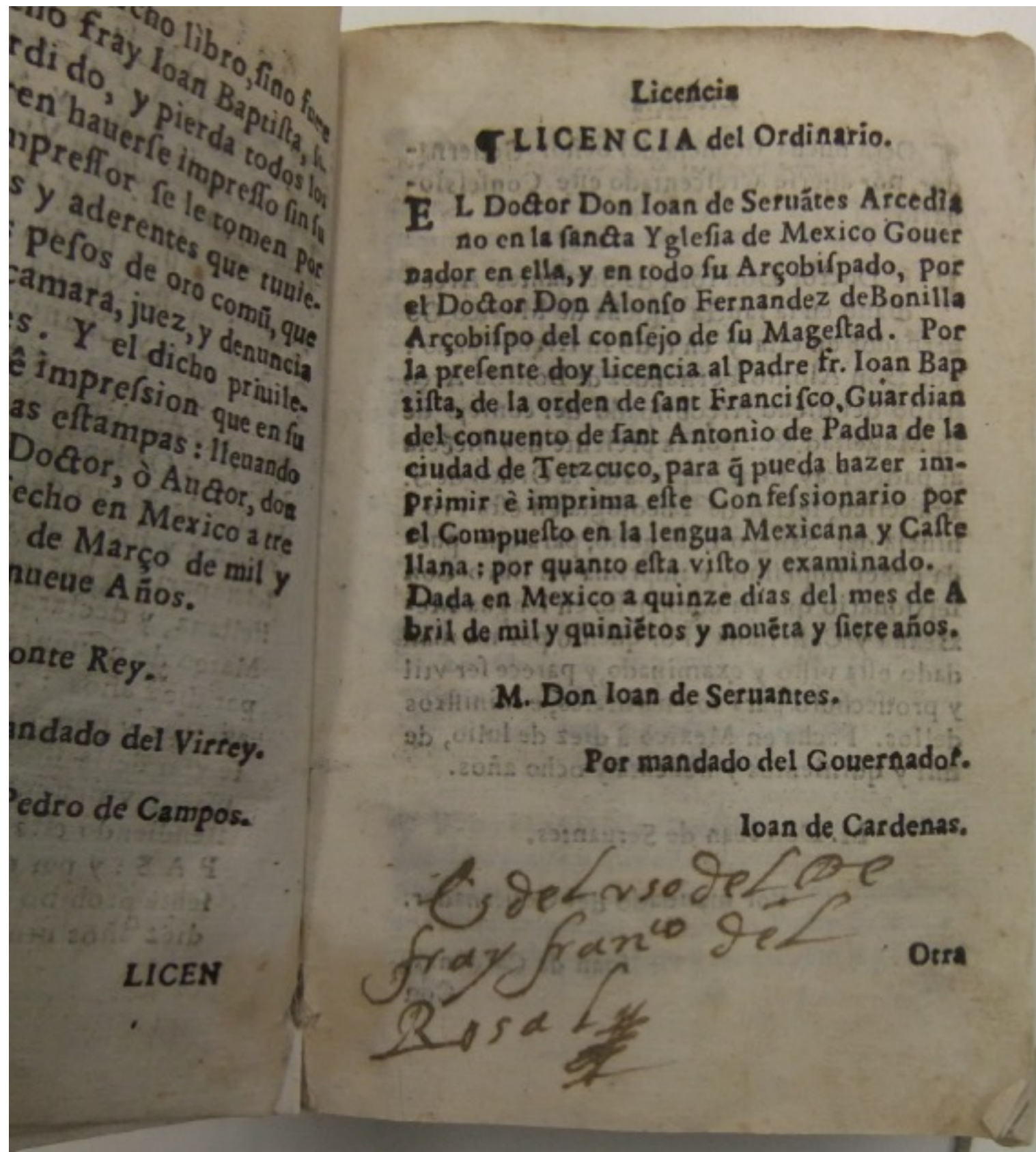
Dan Klein

UC Berkeley Computer Science

Historical Document Transcription

Working with scholars in humanities who want to study texts from the 1500s.

Standard OCR systems don't work well on printing-press books.



State-of-the-Art: *Ocular*

- Berg-Kirkpatrick, Durrett, and Klein 2013

prisoner

role along in silence

the Death of the Deceased,

Multilingual Texts

- But many historical documents are written in, and switch readily between, multiple languages.

Aduertencias para

como parece manifesto en las palabras de S. Ioan, que dize. Tres sunt qui testimoniū dāt in cælo Pater, Verbum, & Spiritus sanctus : & hi tres vnum sunt. i. Ioann. vltimo. Por lo qual deuen ser instruydos y enseñados, que todas tres diuinas personas son vn Dios verdadero; o reformando la sobre dicha proposicion, y añadiendo esta palabra. In huel imeixtintzin, con que se quita toda amphibologia y duda diziendo. In Dios, ca Tettatzin Tepiltzin, Spiritu sancto, ei personas, çan ce huelnelli teutl Dios in huel imeixtintzin, q. d. Dios es Padre, Hijo, y Spū sancto tres personas, vn solo Dios verdadero todas tres, cō la qual reduplicacion se quita toda dubda. Tambien se quita con estas proposiciones. In DIOS, ca Tettatzin, Tepiltzin, Spiritus sancto, çan huel iceltzin teutl Dios tlahtohuani. In Dios, ca Tettatzin, Tepiltzin, Spiritus sancto, imeixtin personas çan huel iceltzin Dios tlahtohuani. Ca in imeixtin personas me ca çan huel iceltzin teutl Dios tlahtohuani in huel imeixtin. ¶ Otros responden [y es e segundo error] ça ce Dios tlahtohuani, immetihttotica, y a algūos de sus ministros les ha parecido el metehttotica, vn vocablo en si d

uino

Aduertencias para

como parece manifesto en las palabras de S. Ioan, que dize.

[Redacted]

Por lo qual deuen ser instruydos y enseñados, que todas tres diuinas personas son vn Dios verdadero; o reformando la sobre dicha proposicion, y añadiendo esta palabra.

[Redacted] con que se quita toda amphibologia y duda diziendo. In Dios,

[Redacted] personas, Dios

q. d. Dios es Padre, Hijo, y [Redacted] tres personas, vn solo Dios verdadero todas tres, cō la qual reduplicacion se quita toda dubda. Tambien se quita con estas proposiciones.

DIOS, [Redacted] Dios

Dios, [Redacted] personas Dios

[Redacted] personas me Dios

Otros responden [y es e segundo error Dios

[Redacted] y a algños de sus ministros les ha parecido el [Redacted] vn vocablo en si d

[Redacted] uino

Spanish

Latin

Nahuatl

Tres sunt qui testimoniā dāt
in cælo Pater, Verbum, & Spiritus sanctus :
& hi tres vnum sunt. i. Ioann. vltimo.

Spiritu sancto, ei

Spū sancto

Spiritulan

cto,

Spiritusan.

cto,

uno

Spanish

Latin

Nahuatl

[Redacted text block]

In huel ime

ixtintzitzin,

ca Tettatzin

Tepiltzin,

çan ce

huelnelli teutl

in huel imeixtintzitzin,

[Redacted text block]

In ca Tettatzin, Tepiltzin,

çan huel iceltzin teutl tlahtohuani,

In ca Tettatzin, Tepiltzin,

imeixtin çan huel iceltzin

tlahtohuani. Ca inimeixtin ca

çan huel iceltzin teutl tlahtohuani in

huel imeixtin.

çace tlahtohuani, imne

teihhtotica,

meteihhtotica,

Spanish

Latin

Nahuatl

[Redacted text block]

[Redacted text block]

Spanish

Latin

Nahuatl

[Redacted text block]

Starting Point: *Ocular*

Generative Model in 3 parts:

1. Language model
2. Typesetting model
3. Rendering model

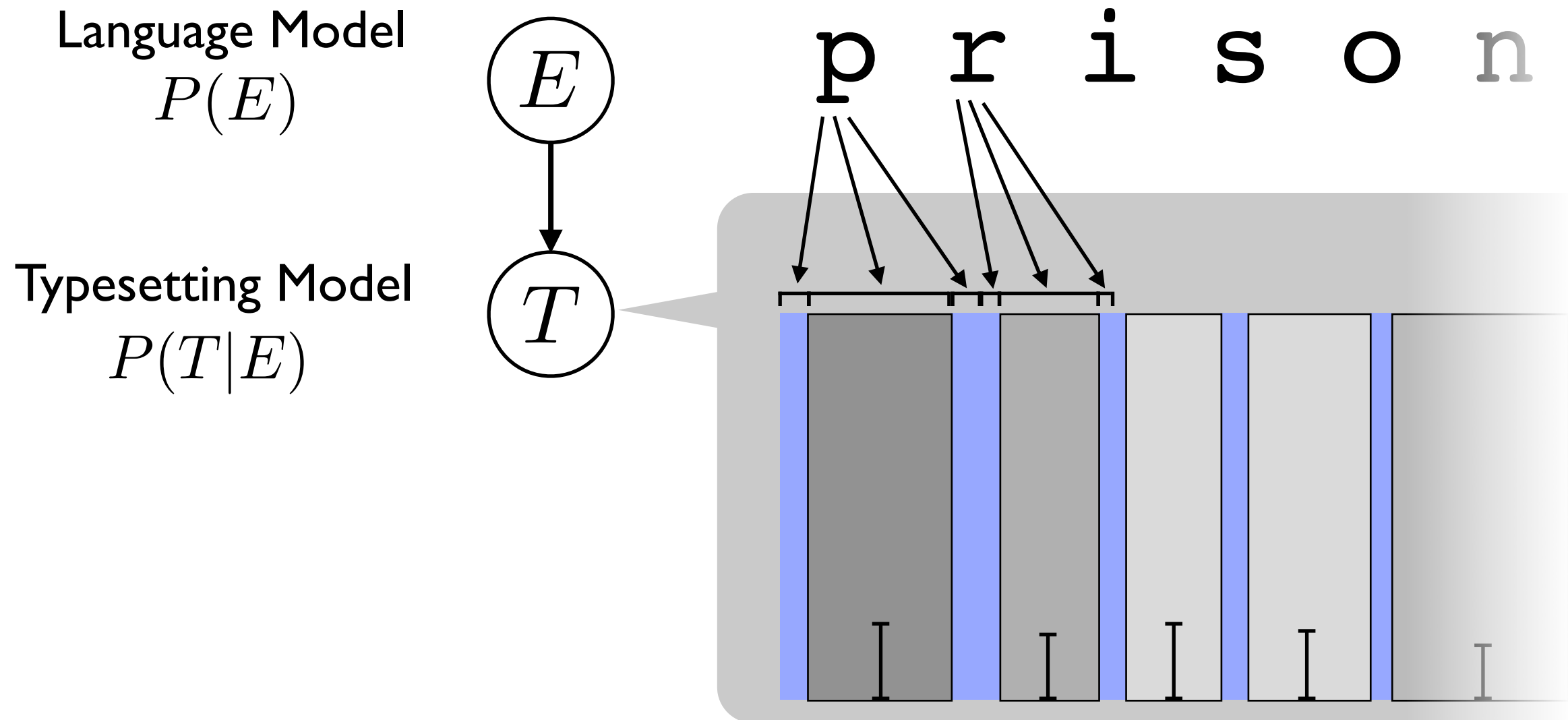
Ocular's Generative Model

Language Model
 $P(E)$

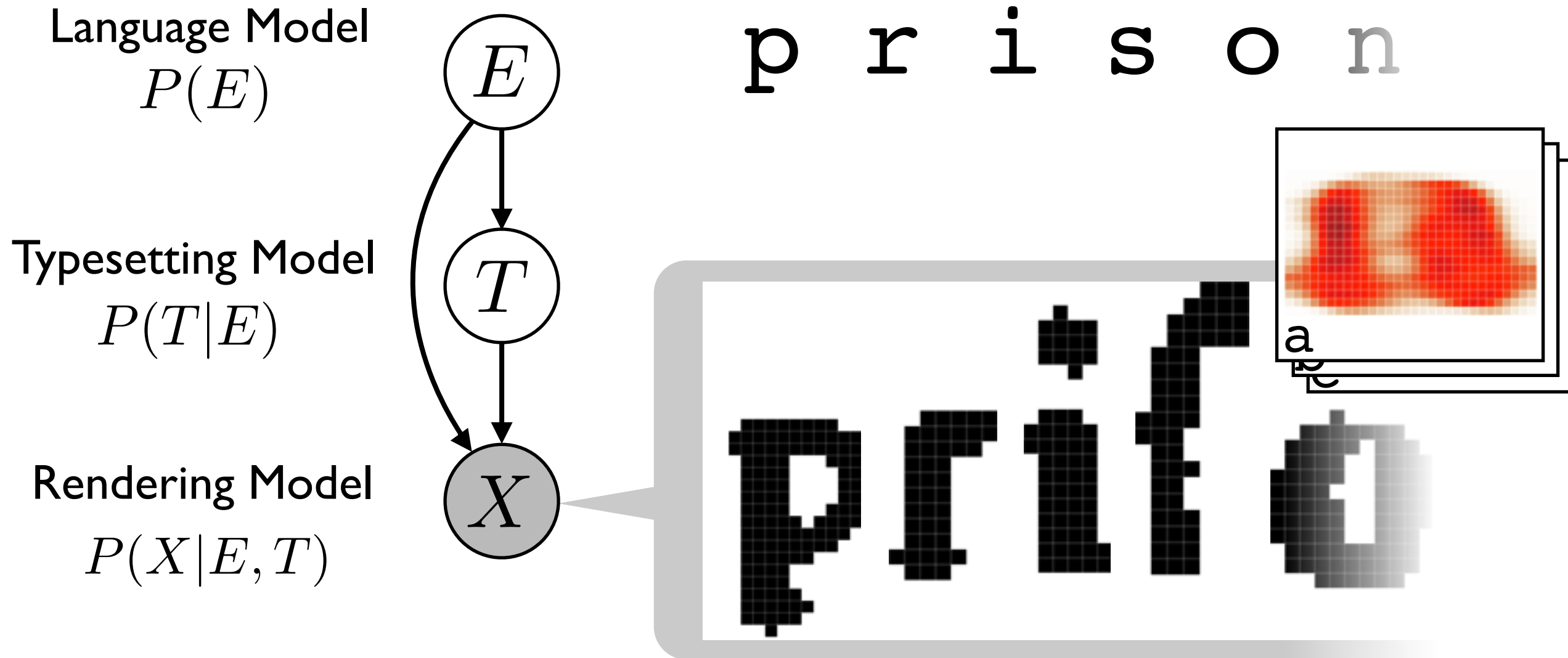
E

p r i s o n

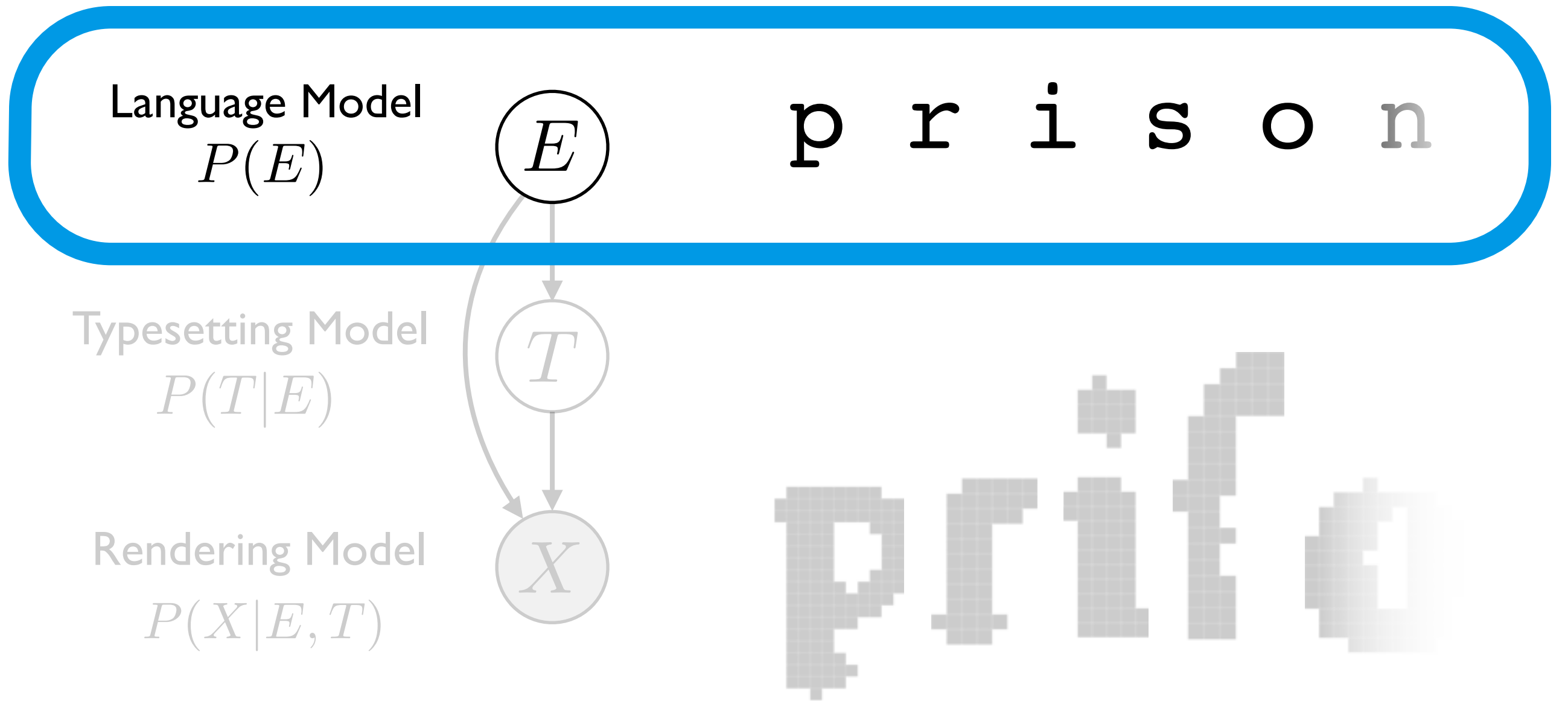
Ocular's Generative Model



Ocular's Generative Model



Our Focus



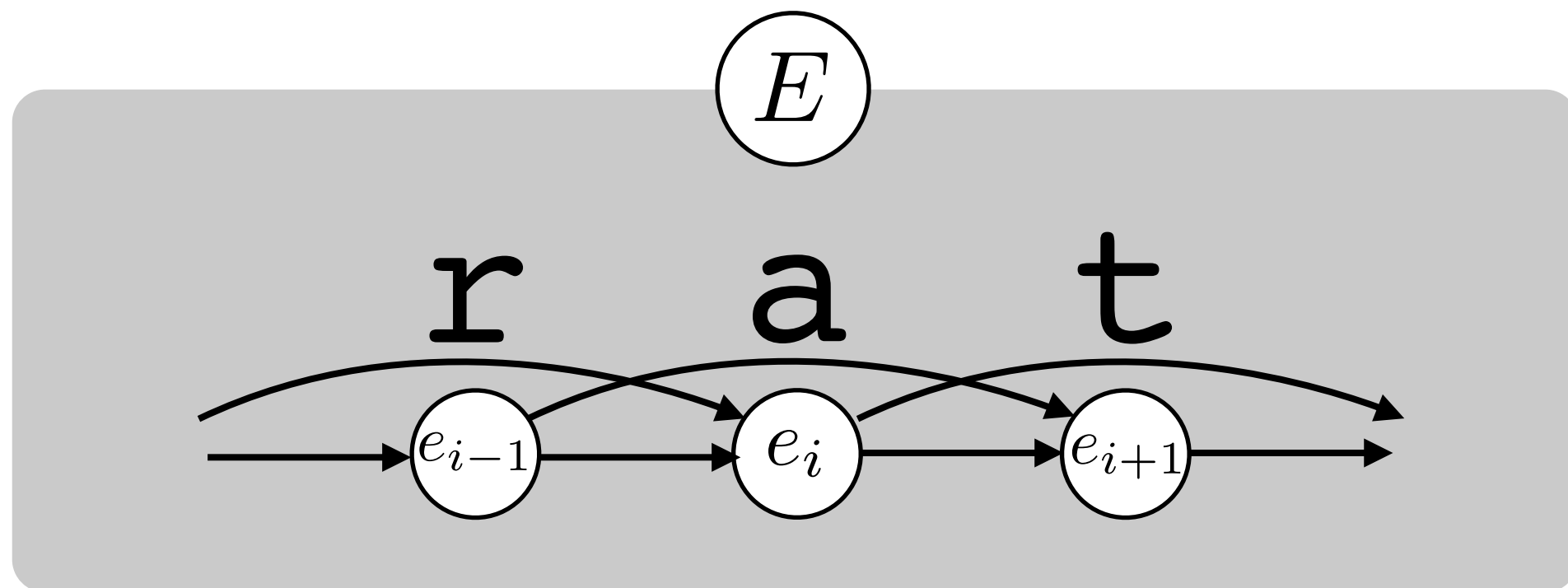
Starting Point: *Ocular*

- The language model helps Ocular work well, but creates additional challenges for many documents.
- Our work helps to overcome those challenges.

Our Focus

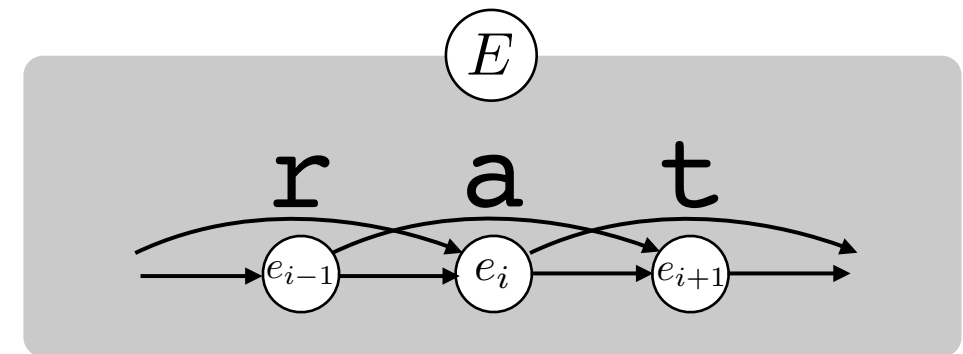
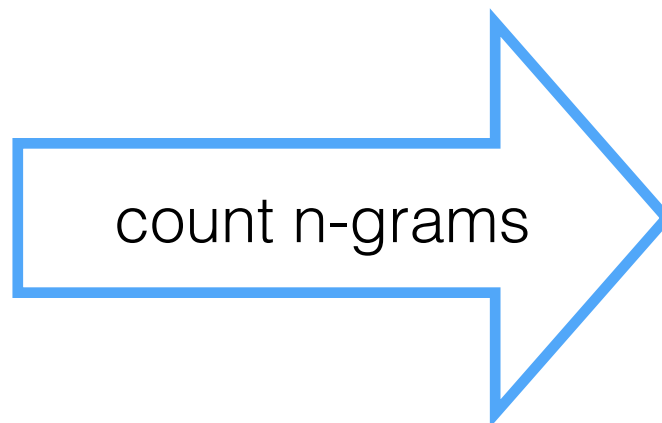
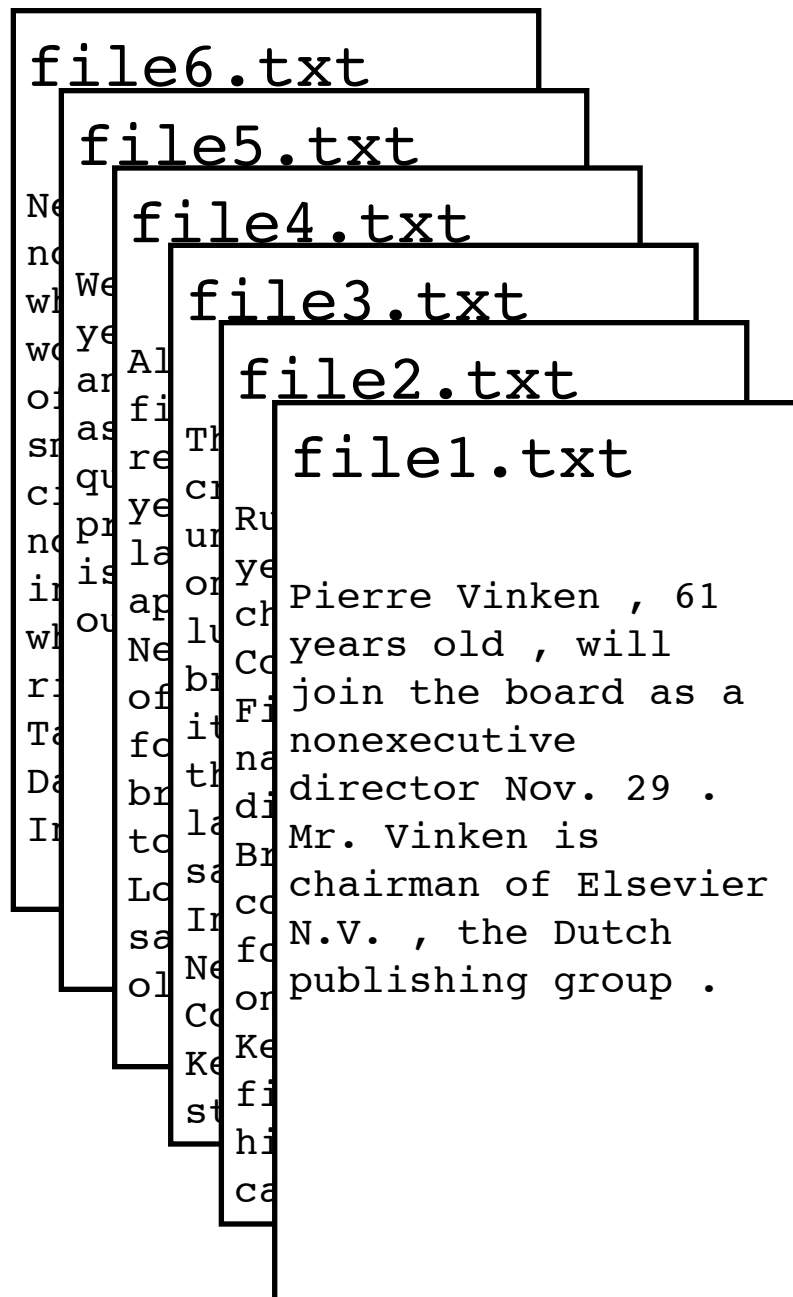
1. Multilingual code-switching
2. Inconsistent/outdated orthography

Ocular's Language Model



Kneser-Ney smoothed **character** 6-gram

Ocular's Language Model

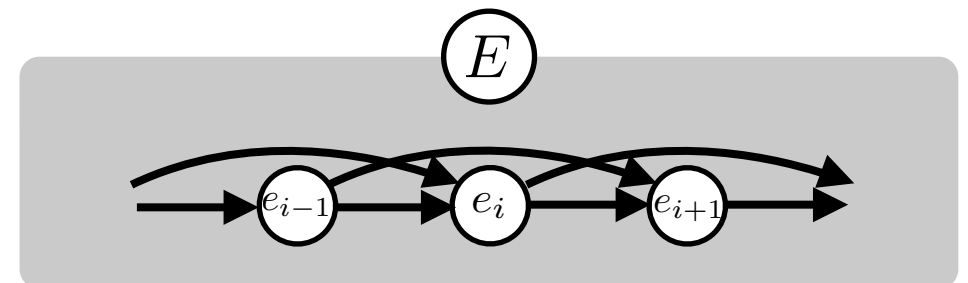


Baseline Multilingual Model

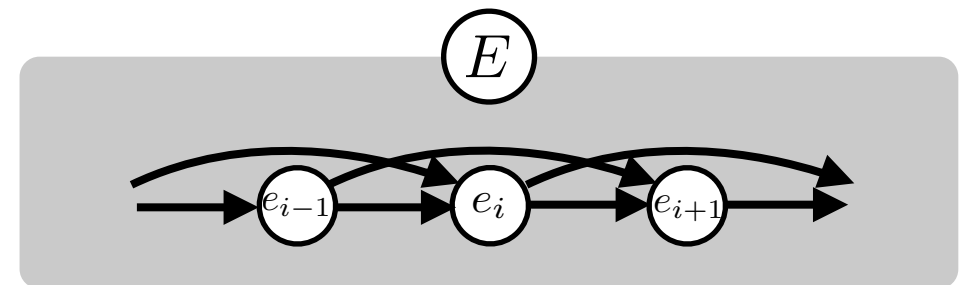
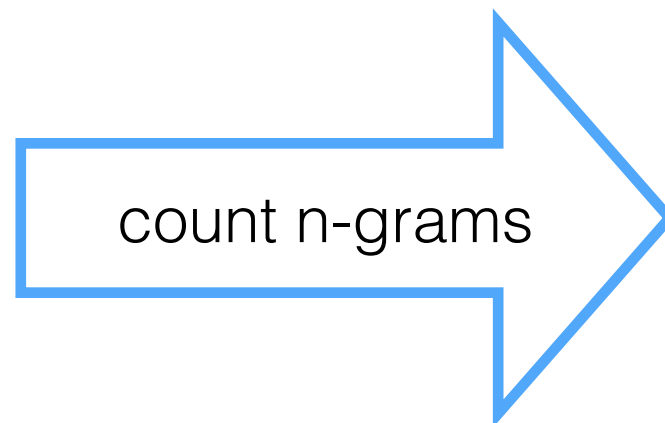
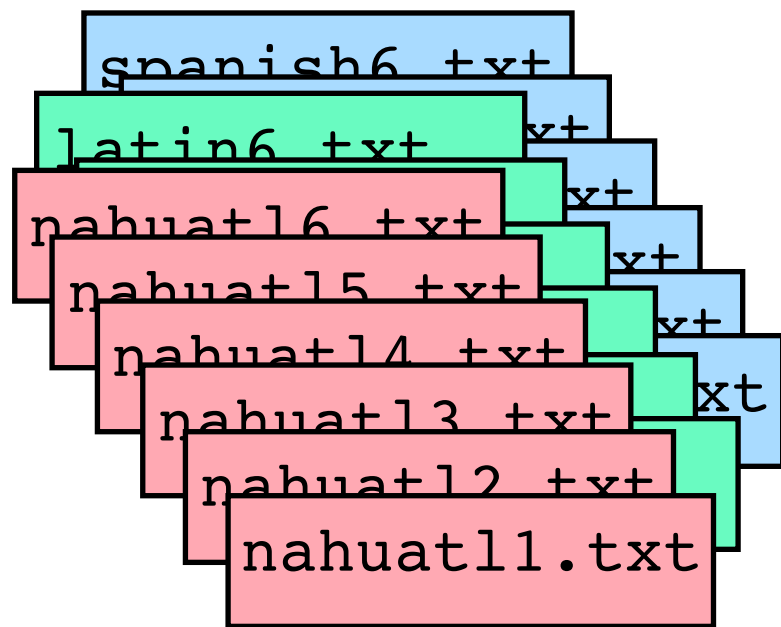
spanish6.txt
spanish5.txt
spanish4.txt
spanish3.txt
spanish2.txt
spanish1.txt

latin6.txt
latin5.txt
latin4.txt
latin3.txt
latin2.txt
latin1.txt

nahuatl6.txt
nahuatl5.txt
nahuatl4.txt
nahuatl3.txt
nahuatl2.txt
nahuatl1.txt



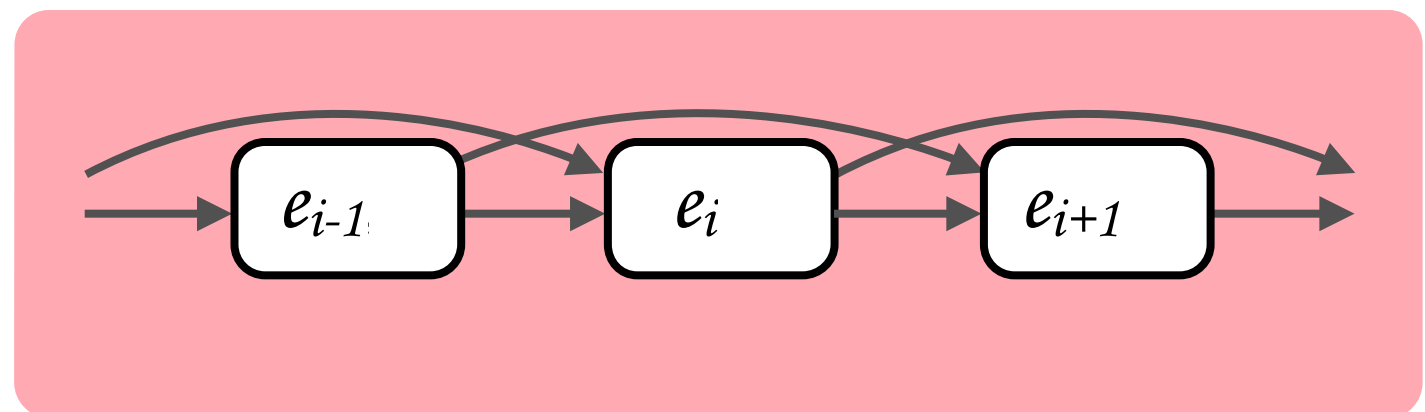
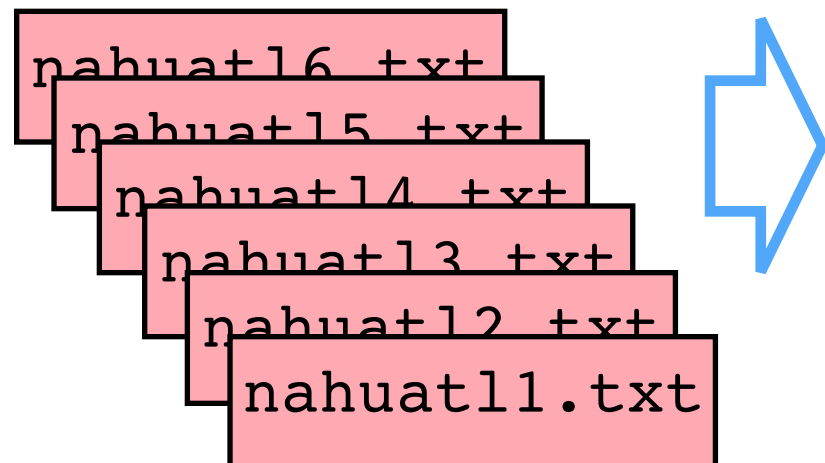
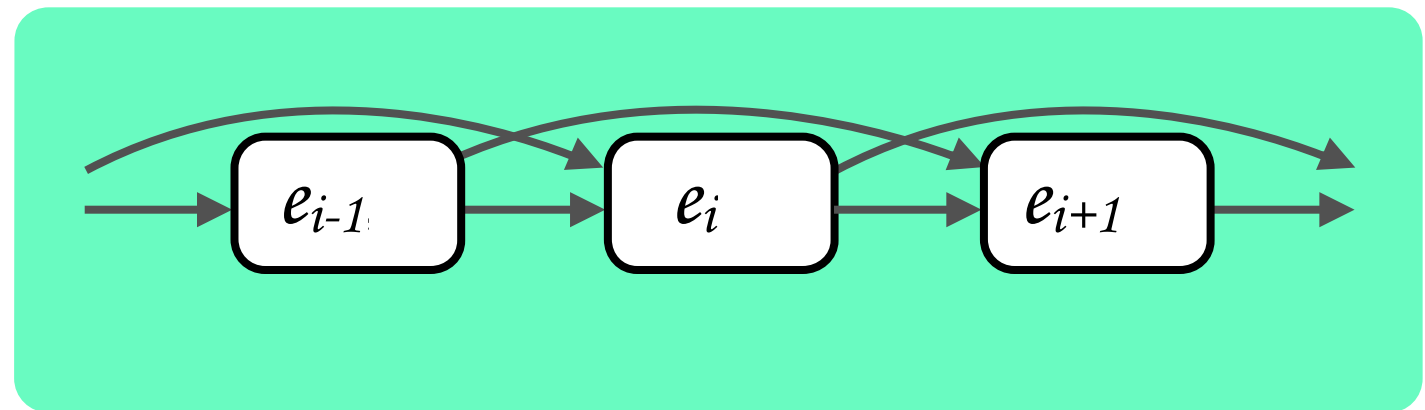
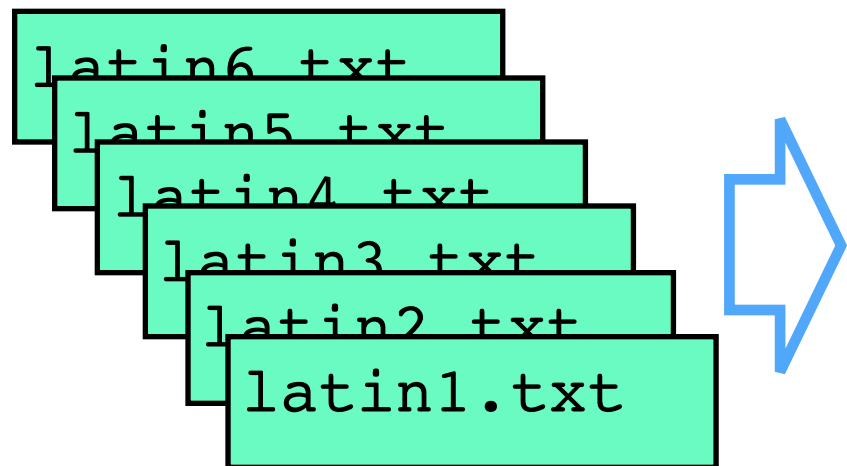
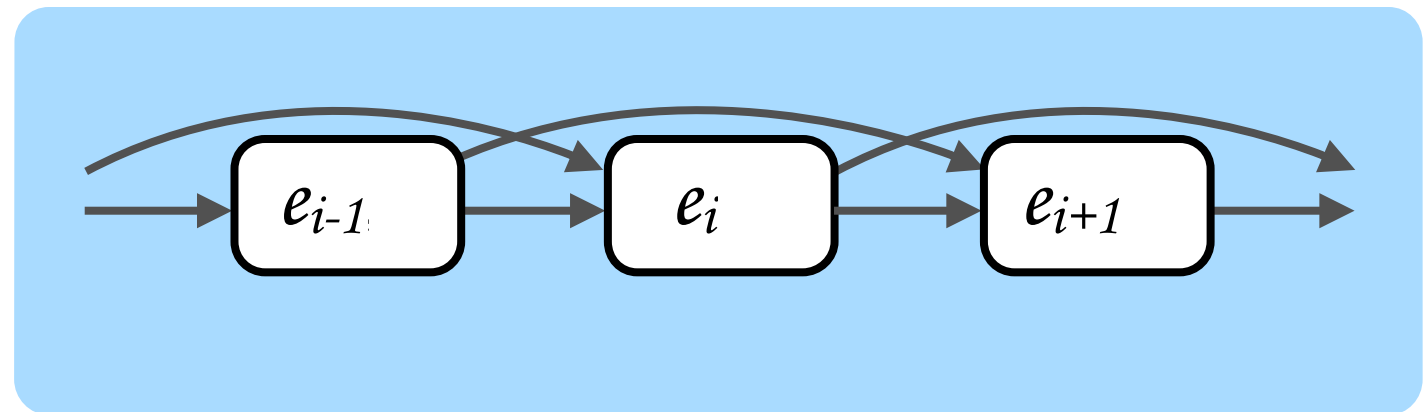
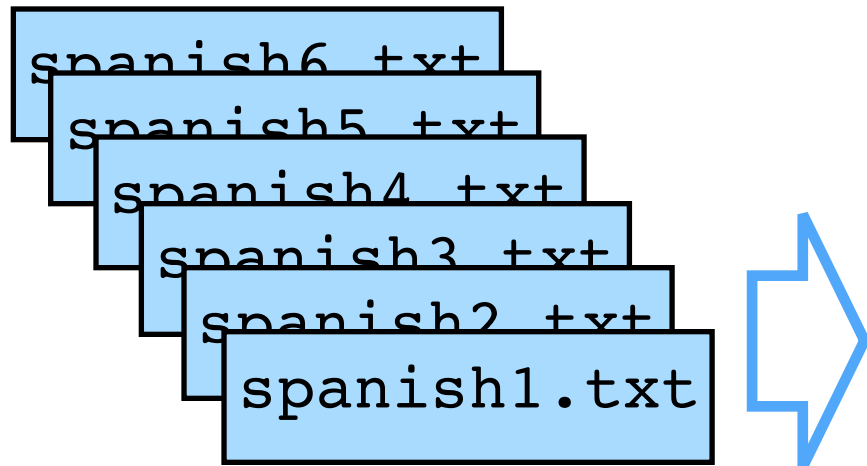
Baseline Multilingual Model



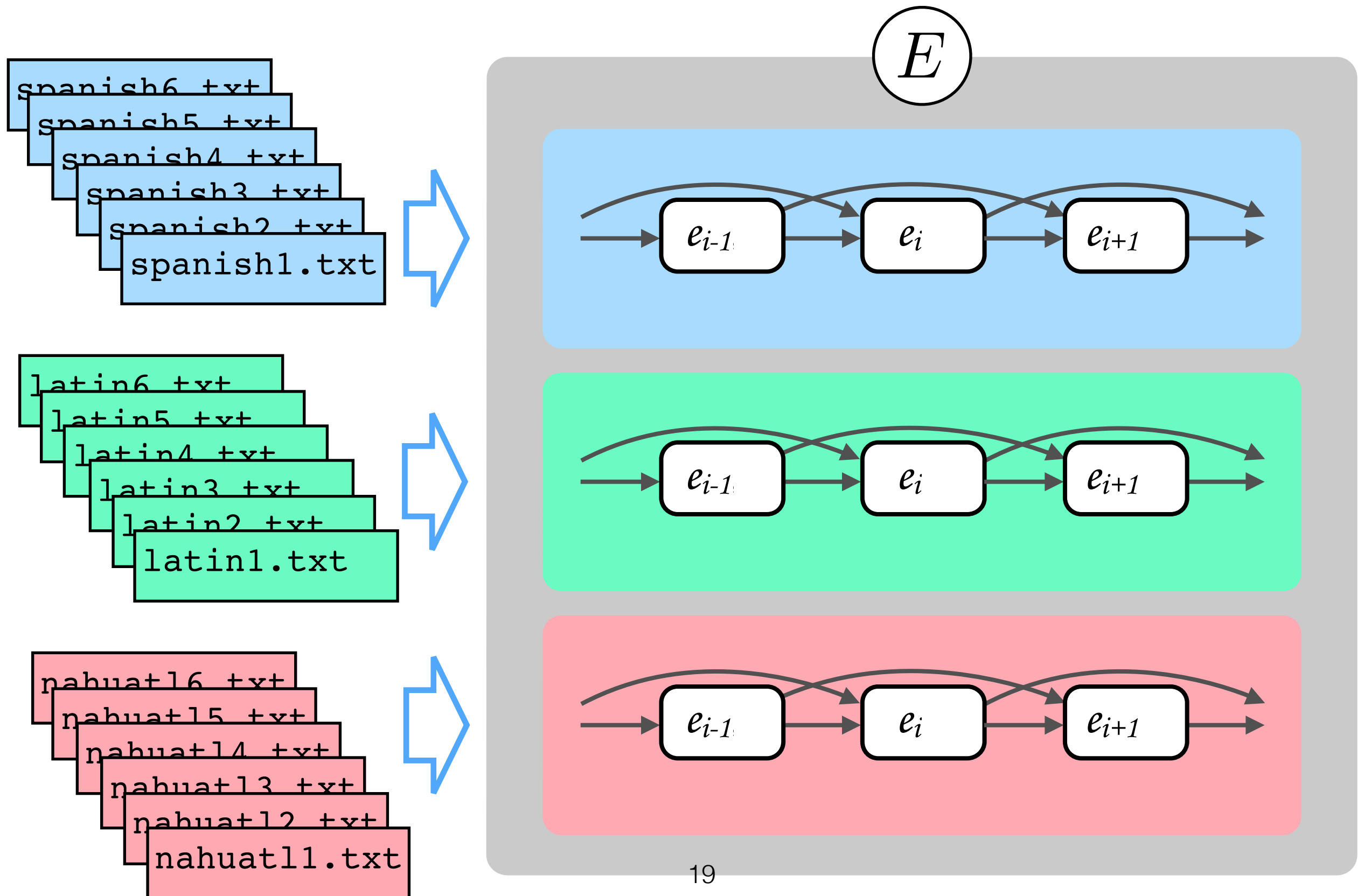
Baseline Multilingual Model

- Poor results
- “Multilingual blur”

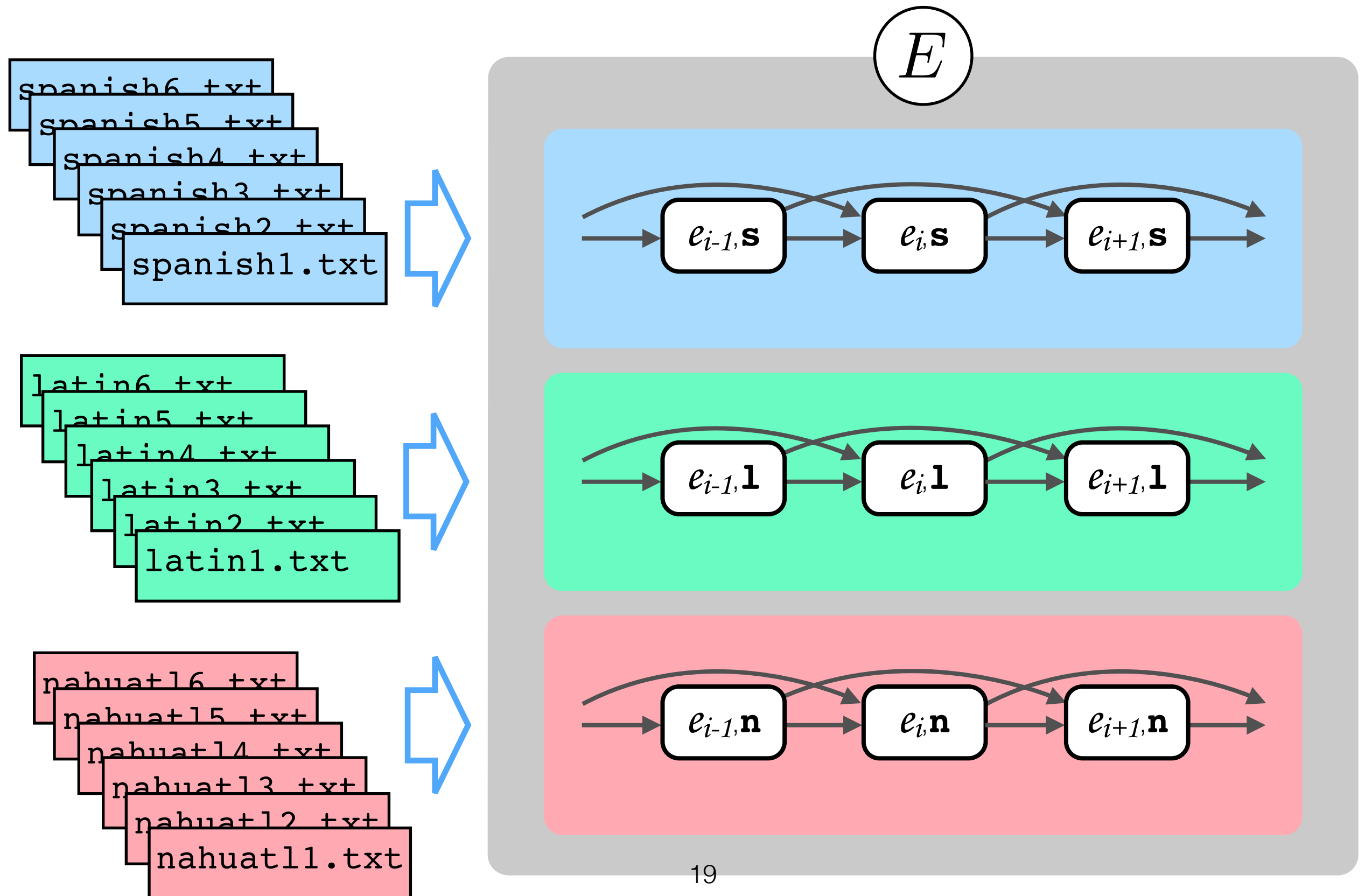
Code-Switching Language Model



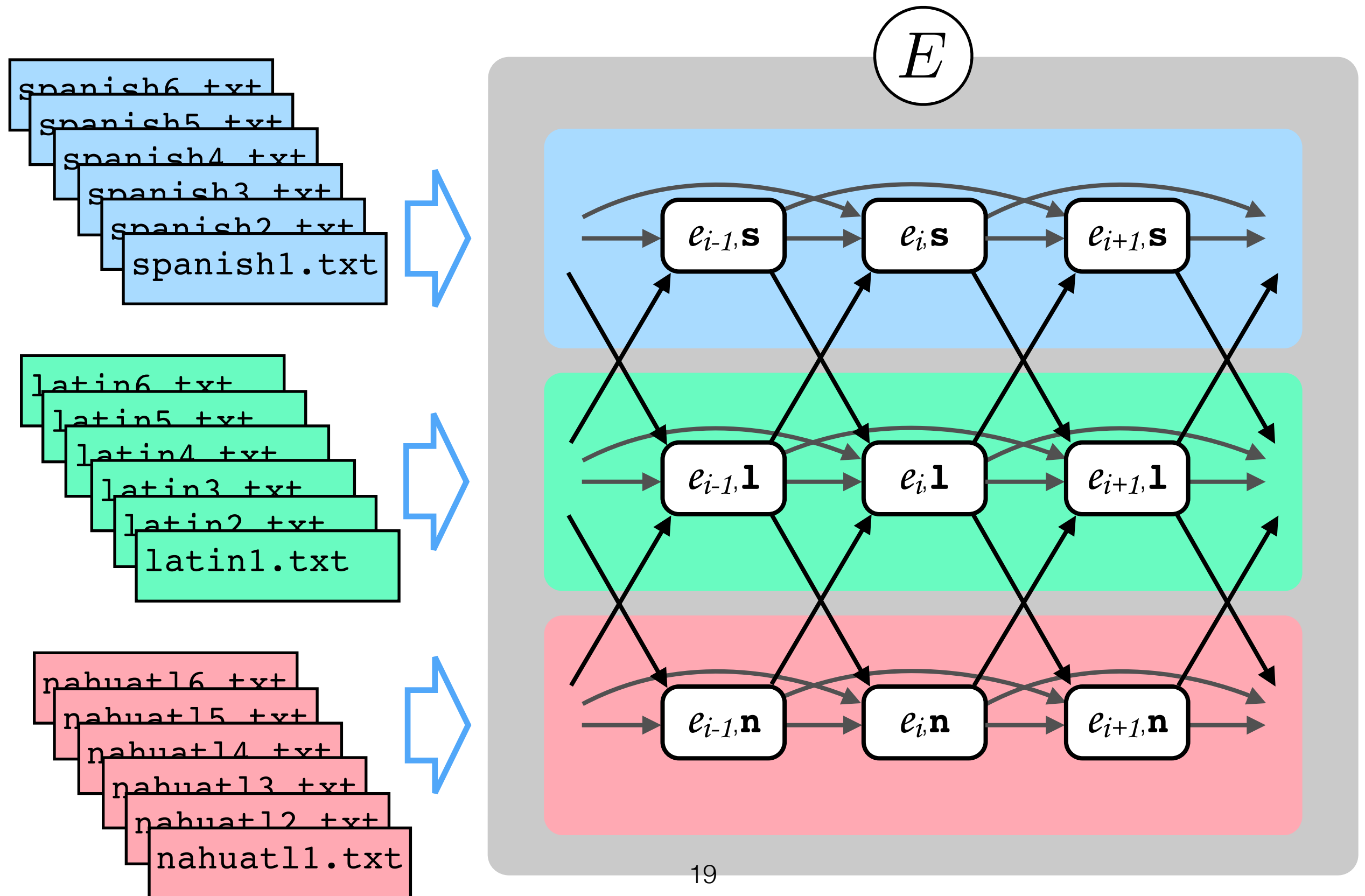
Code-Switching Language Model



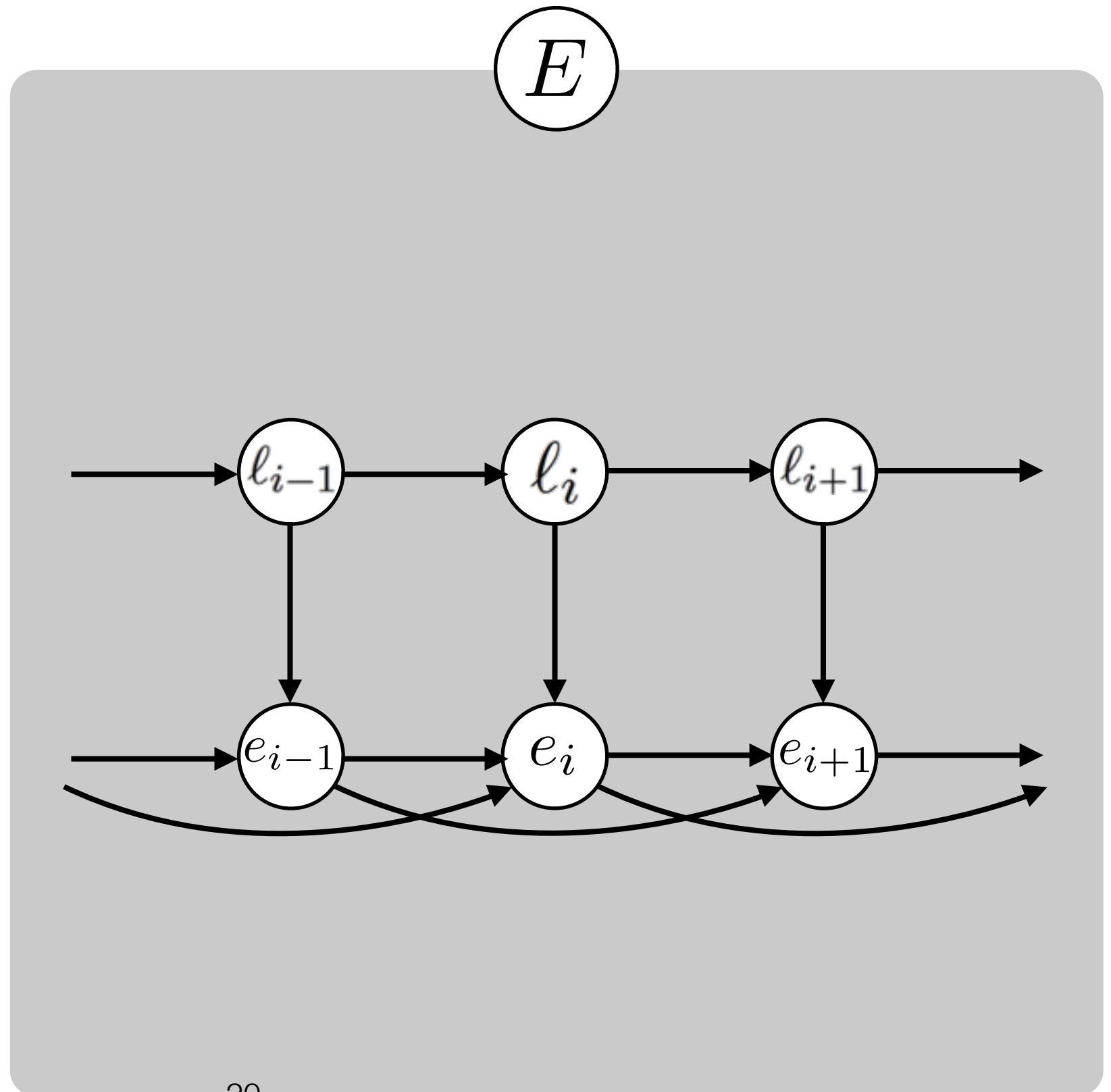
Code-Switching Language Model



Code-Switching Language Model

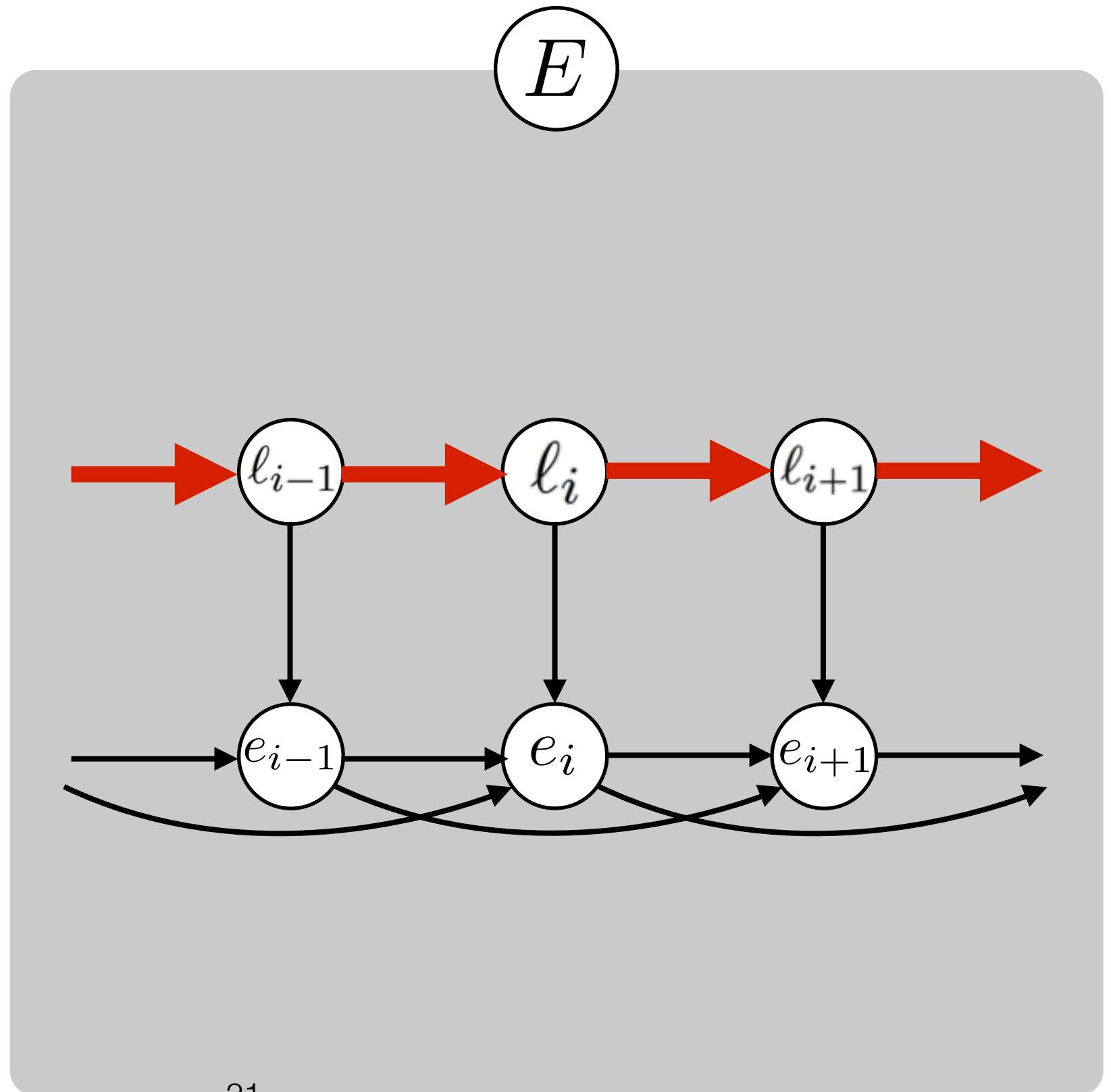


Code-Switching Language Model



Code-Switching Language Model

$P(l_i | l_{i-1})$ is learned unsupervised via EM, with a hyperparameter biasing the model toward **not** switching (long language spans)



Aduertencias para

como parece manifesto en las palabras de S. Ioan, que dize. Tres sunt qui testimoniū dāt in cælo Pater, Verbum, & Spiritus sanctus : & hi tres vnum sunt. i. Ioann. vltimo. Por lo qual deuen ser instruydos y enseñados, que todas tres diuinas personas son vn Dios verdadero; o reformando la sobre dicha proposicion, y añadiendo esta palabra. In huel imeixtintzin, con que se quita toda amphibologia y duda diziendo. In Dios, ca Tettatzin Tepiltzin, Spiritu sancto, ei personas, çan ce huelnelli teutl Dios in huel imeixtintzin, q. d. Dios es Padre, Hijo, y Spū sancto tres personas, vn solo Dios verdadero todas tres, cō la qual reduplicacion se quita toda dubda. Tambien se quita con estas proposiciones. In DIOS, ca Tettatzin, Tepiltzin, Spiritus sancto, çan huel iceltzin teutl Dios tlahtohuani. In Dios, ca Tettatzin, Tepiltzin, Spiritus sancto, imeixtin personas çan huel iceltzin Dios tlahtohuani. Ca in imeixtin personas me ca çan huel iceltzin teutl Dios tlahtohuani in huel imeixtin. ¶ Otros responden [y es e segundo error] ça ce Dios tlahtohuani, immetihttotica, y a algūos de sus ministros les ha parecido el metehttotica, vn vocablo en si d

ficion, y añadiendo esta palabra. In huel ime
ixtintzitzin, con que se quita toda amphibo-
logia y duda diziendo. In Dios,ca Tettatzin

añadiendo esta palabra. In huel ~~ime~~ ixtintzitzin, con que se quita toda amphibologia y duda diziendo. In Dios,

ixtintzitzin, con que

Spanish

AÁBCDÉFGHIÍJKLMÑOÓPQRSTUÚVWXYZ
aábcdeéfgghiíjklmñoópqrstuúvwxyz
01234567890.,/\()?!"':;-

Latin

ABCDFGHIJKLMOPQRSTUVWXYZ
abcdfghijklmopqrstuvwxyz
01234567890.,/\()?!"':;-

Nahuatl

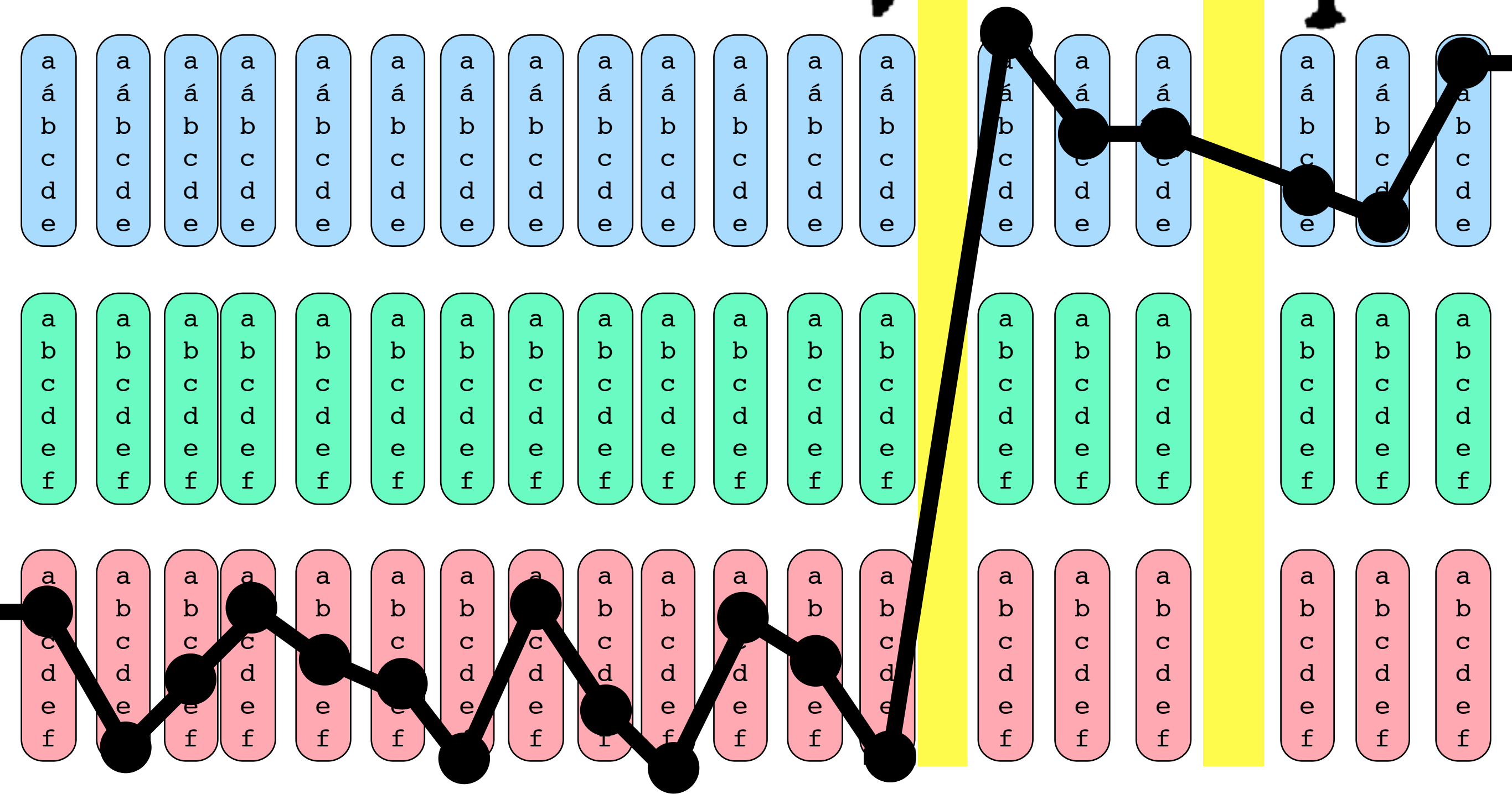
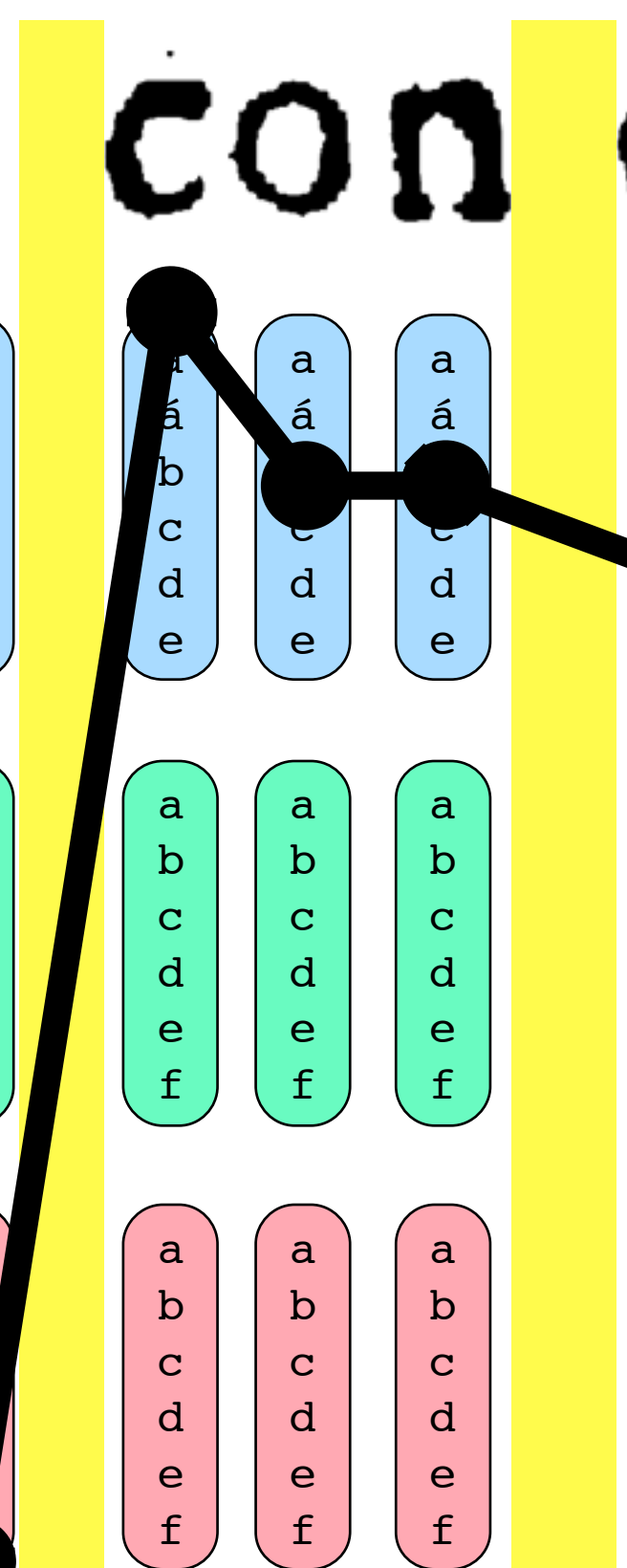
ABCDFGHIJKLMOPQRSTUVWXYZ
abcdfghijklmopqrstuvwxyz
01234567890.,/\()?!"':;-

ixtintzitzin, con que

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a |
| á | á | á | á | á | á | á | á | á | á | á | á | á | á | á | á | á | á | á |
| b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b |
| c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c |
| d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d |
| e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e |

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a |
| b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b |
| c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c |
| d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d |
| e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e |
| f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f |

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a | a |
| b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b | b |
| c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c | c |
| d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d | d |
| e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e | e |
| f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f | f |



Code-Switching Language Model

- Improves transcription quality, *and*
- Implicitly identifies language spans in text (metadata of the transcription)

Orthographic Variability

Orthographic Variability

- We train our language models from available text (e.g. Project Gutenberg)
- Modern transcribers use modern spellings, which often do not match the printed documents

Orthographic Variability

Transcription

Modern Form

dize

dize

dice

numero

numero

número

Dõde

Dõde

Donde

Orthographic Variability

Simple solution:

Modify the modern corpora to use old conventions.

Orthographic Variability

Modern Spanish

spanish6.txt
spanish5.txt
spanish4.txt
spanish3.txt
spanish2.txt
spanish1.txt



Replacement Rules

u → v
c → z
ú → u
on → ñ
que → ñ
...



Old Spanish

spanish6b.txt
spanish5b.txt
spanish4b.txt
spanish3b.txt
spanish2b.txt
spanish1b.txt

Experiments

Experiments

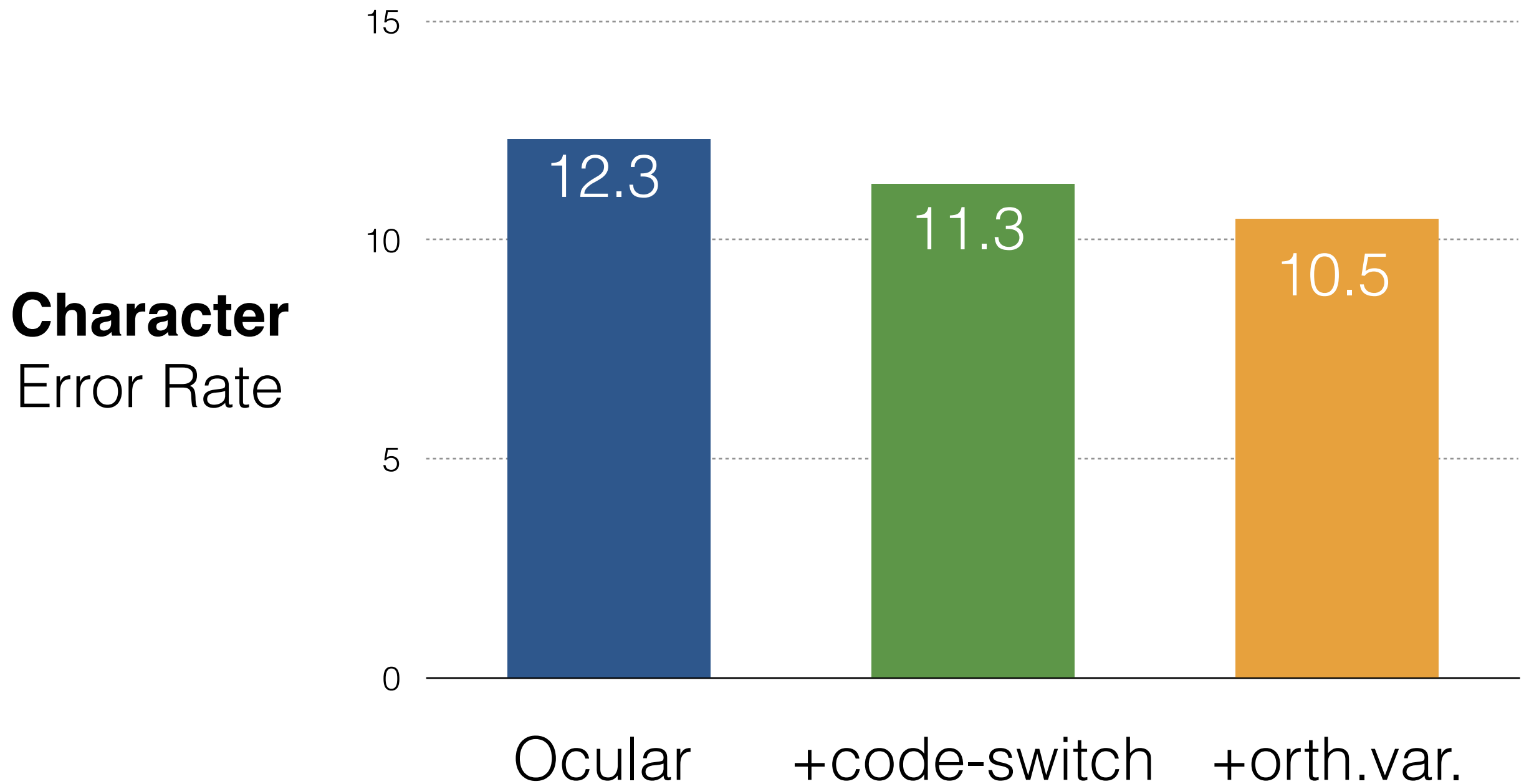
- Evaluated on five different books *Primeros Libros*
- Years 1553 to 1600
- Differing fonts, language proportions, clarity

Unknown Fonts

| | | |
|-------------|--------|---|
| Gante | (1553) | motlacatilia: ynica sacramento Baptis |
| Anunciación | (1565) | ¶ Rincóneltoquitia yndios |
| Sahagún | (1583) | Yoan oquihui in Emperador, in taca |
| Rincón | (1595) | etion.v.g.tetlaçotlaliztli.amatio, vel, |
| Bautista | (1600) | ¶ imo, hæc supra dictus doctor Medina. |

Experimental Results

(lower is better. ~90% characters are correct)



A thing we do well

mēcira

Without handling orth. variation: `merita`

With handling orth. variation: `mētira`

Modern form: `mentira`

A thing we do wrong

Los verbales en, li, o en tli, antepo-

A thing we do wrong

Los verbales en, li, o en tli, antepo-

Model output:

sí

tĩi

← Spanish

Gold transcription:

li

tli

← Nahuatl

Model avoids switching languages, but this is actually from a description of Nahuatl grammar.

A thing that's hard

tetechtlāmieccaquixtiliztli

tetechtla miec caquixtiliztli

- All letters are correct, but the model adds spaces
- No agreed-upon standards for Nahuatl spacing
- Hard to evaluate what is “correct”

Conclusion

Conclusion

- By accounting for **multilingual text** and **obsolete orthography**, we can improve the state-of-the-art for historical OCR.
- These are common characteristics of texts from all over the world and from all eras.
- Expansion of OCR abilities means a wider range of texts may be available for study.