

Building Early Colonial Corpora for Digital Scholarship



Hannah Alpert-Abrams
CLIR Postdoctoral Fellow
University of Texas at Austin

LLILAS BENSON

<https://goo.gl/atVK2Q>



Hannah Alpert-Abrams

with

Dan Garrette

Taylor Berg-Kirkpatrick

Bryan Tarpley

Maria Victoria Fernández

Albert Palacios

Stephanie Wood

Kent Norsworthy

...

LLILAS BENSON

<https://goo.gl/atVK2Q>

Corpus Creation Challenges

Curation
Creation
Encoding

Colonial Corpus Challenges

Language
Orthography
Epistemology

+ more?

How do we draw on our training as humanists to build better colonial corpora?

How do we draw on our training as humanists to build better colonial corpora?

Methods

Tools

Evaluation

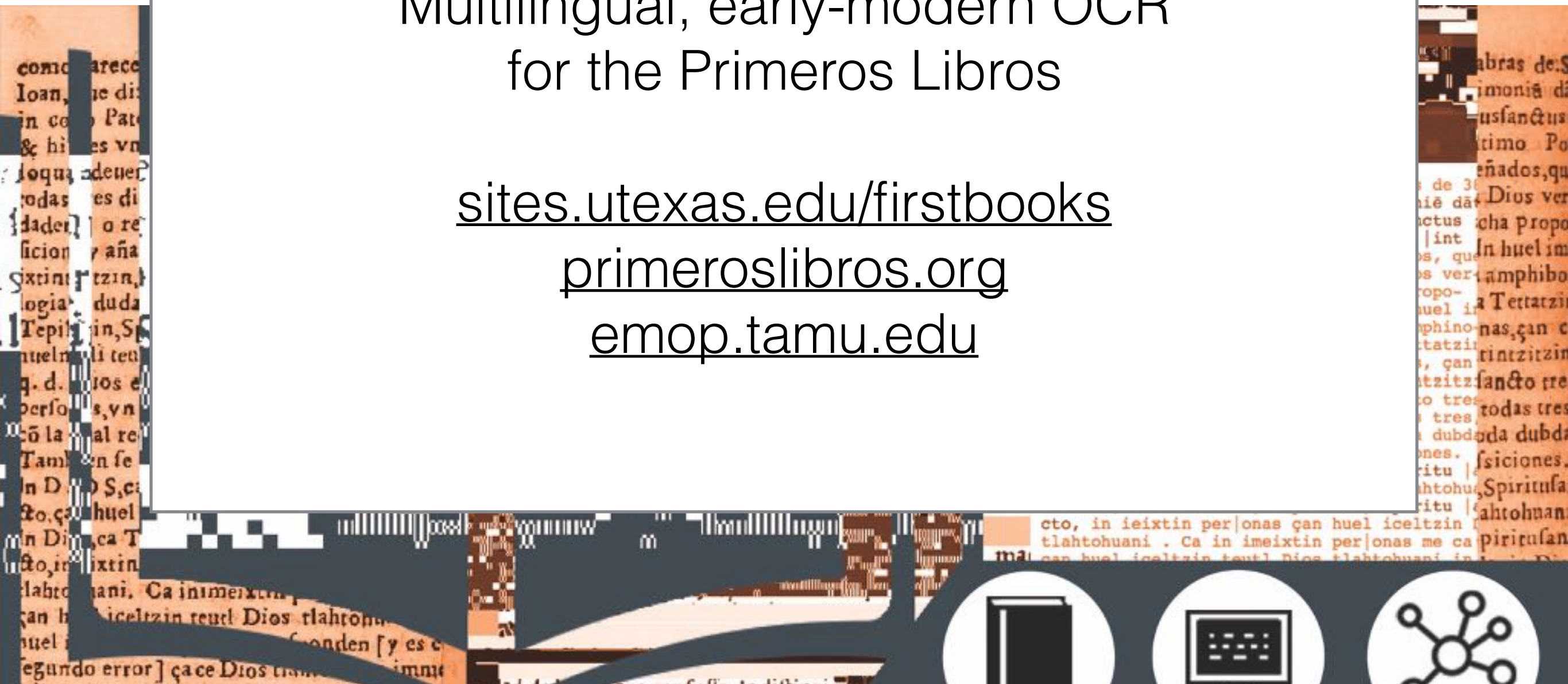
Reading the First Books

Multilingual, early-modern OCR
for the Primeros Libros

sites.utexas.edu/firstbooks

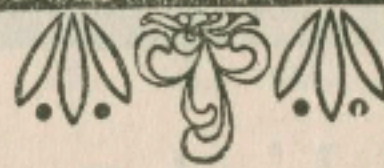
primeroslibros.org

emop.tamu.edu



Digitization methods shape corpora.

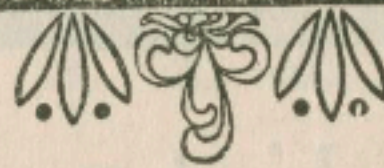
PIANOS THAT SERVE



Smithsonian Transcription Center:

National Baptist Metoka and Galeda Bible Class Magazine, September 1917
Hannah Alpert-Abrams <https://goo.gl/atVK2Q>

PIANOS THAT SERVE



[[boxed advertisement]]

PIANOS THAT SERVE

[[double line]]

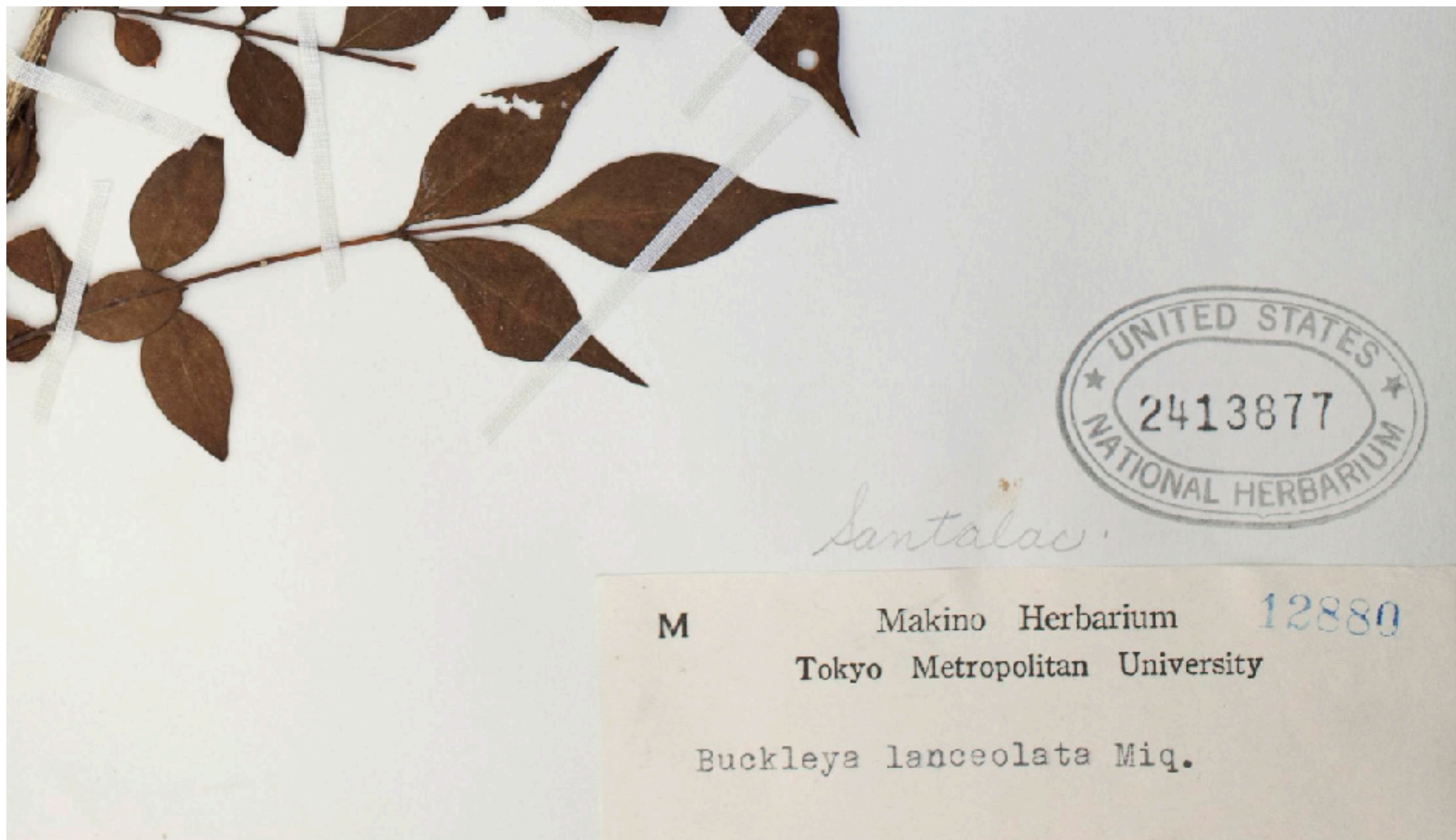
[[image - small floral sketch]]

Smithsonian Transcription Center:

National Baptist Metoka and Galeda Bible Class Magazine, September 1917

Hannah Alpert-Abrams

<https://goo.gl/atVK2Q>



TRANSCRIPTION FORM

INSTRUCTIONS

1) General

2) Collector Details

3) Location

Collector #

12880

Collection Date

1933

single date

Collector

Makino, T. (Makino)

Notes on Transcribing this page (optional)

UNITED STATES

2413877

NATIONAL HERBARIUM

M

Makino Herbarium

12880

Tokyo Metropolitan University

Buckleya lanceolata Miq.



Elizabeth Hill Boone, *The Red and the Black*

Hannah Alpert-Abrams

<https://goo.gl/atVK2Q>



Viceroy Mendoza
(maguey + tozan [gopher])

Elizabeth Hill Boone, *The Red and the Black*

Transcription is interpretation.

Transcription is
specialized interpretive labor.

Digitization methods shape corpora.

We need to think about the ways we employ
labor in colonial corpus creation.

Transcription is
specialized interpretive labor even when it is
done by a machine.

Digitization tools shape corpora.

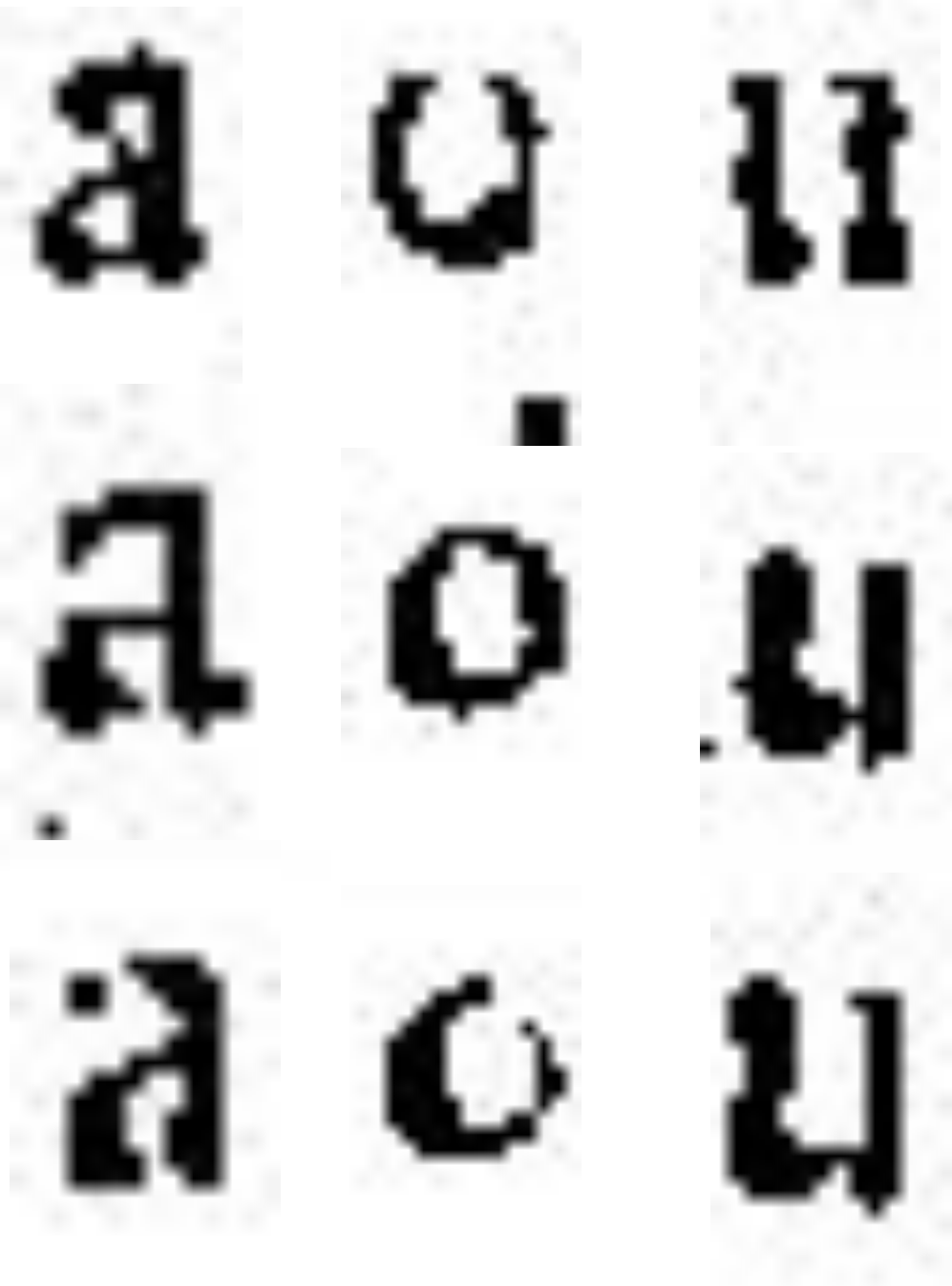


အ
အ
အ

ဝ
ဝ
ဝ

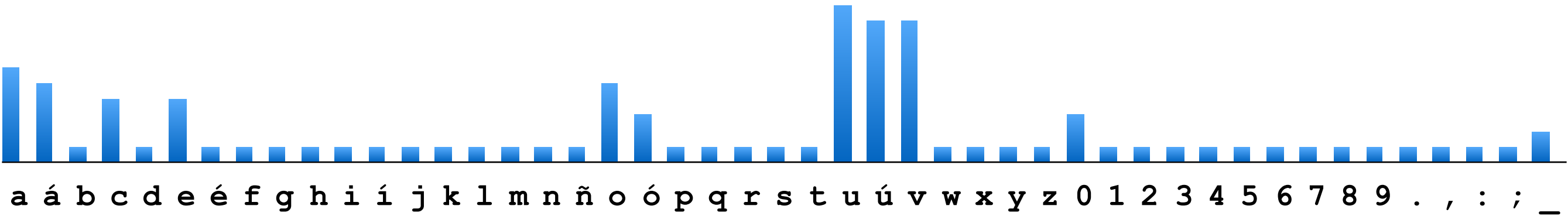
၁၁
၁၁
၁၁

၁၁

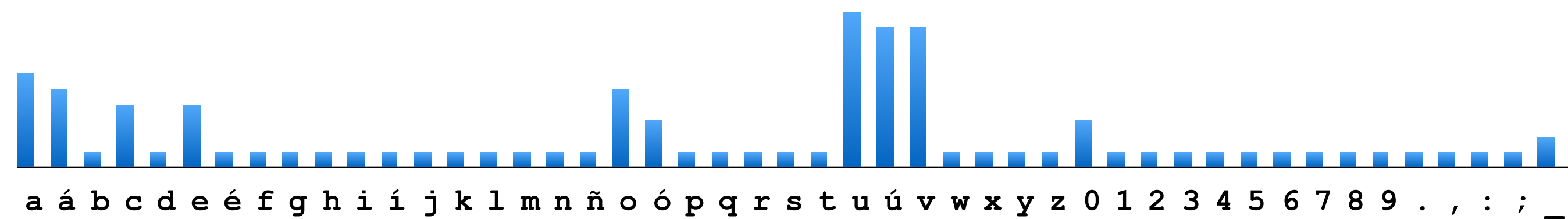


a á b c d e é f g h i í j k l m n ñ o ó p q r s t u ú v w x y z 0 1 2 3 4 5 6 7 8 9 . , : ;

Font Model



conque fe



con que fe

con que fe

resp. por la mayor parte. Ca yehuatzin in

resp pool a may or parte. Ca yeftnat win fit

responsortia mayor Paris Chairman analysts

refpi por la mayor parte. Ca yehuatzin in

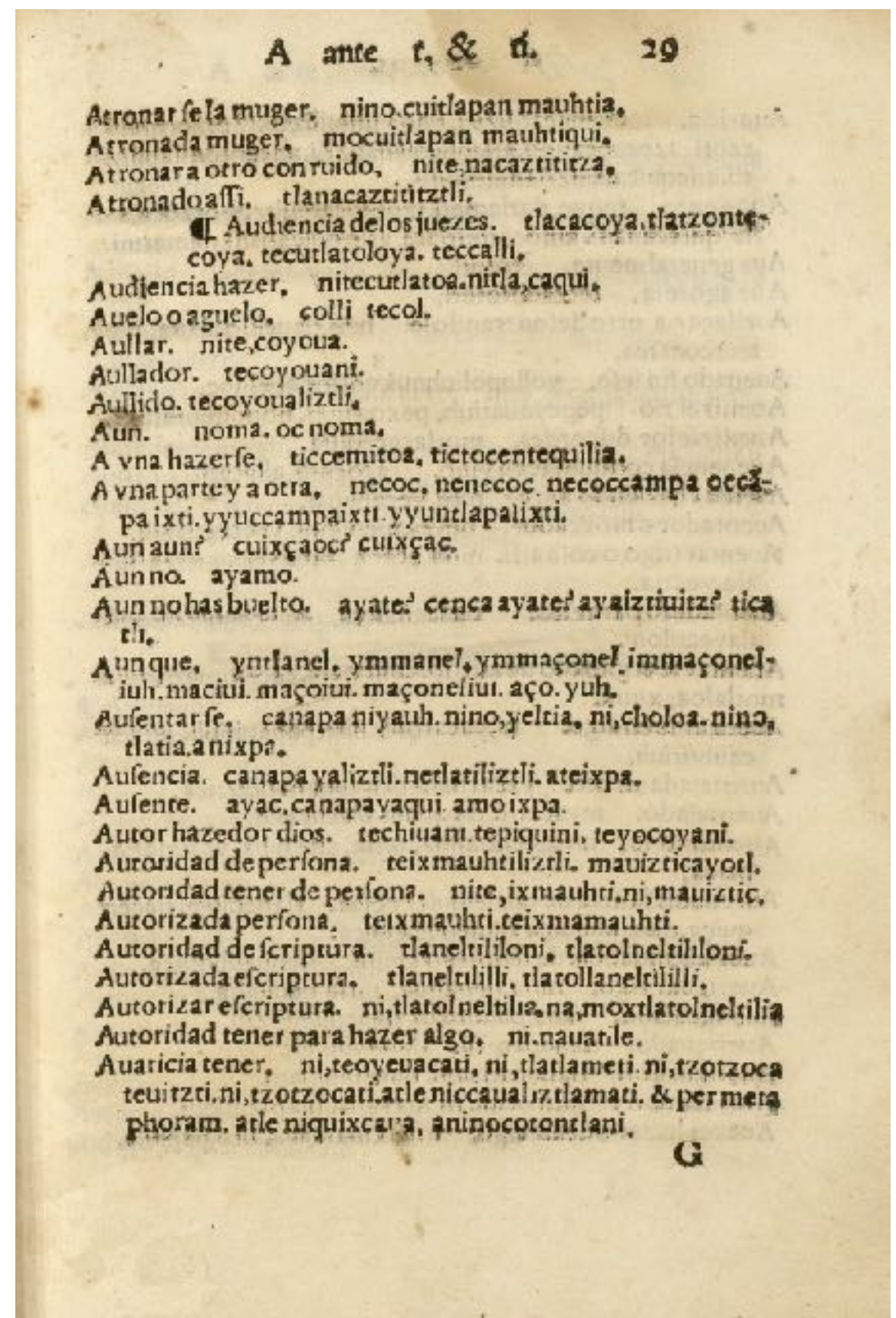
Digitization tools shape corpora.

We need tools that are sensitive to the
needs of our collections.

Elvia Arroyo Ramírez, Invisible Defaults and Perceived
Limitations: Processing the Juan Gelman Files.

Evaluation systems must be
context-specific

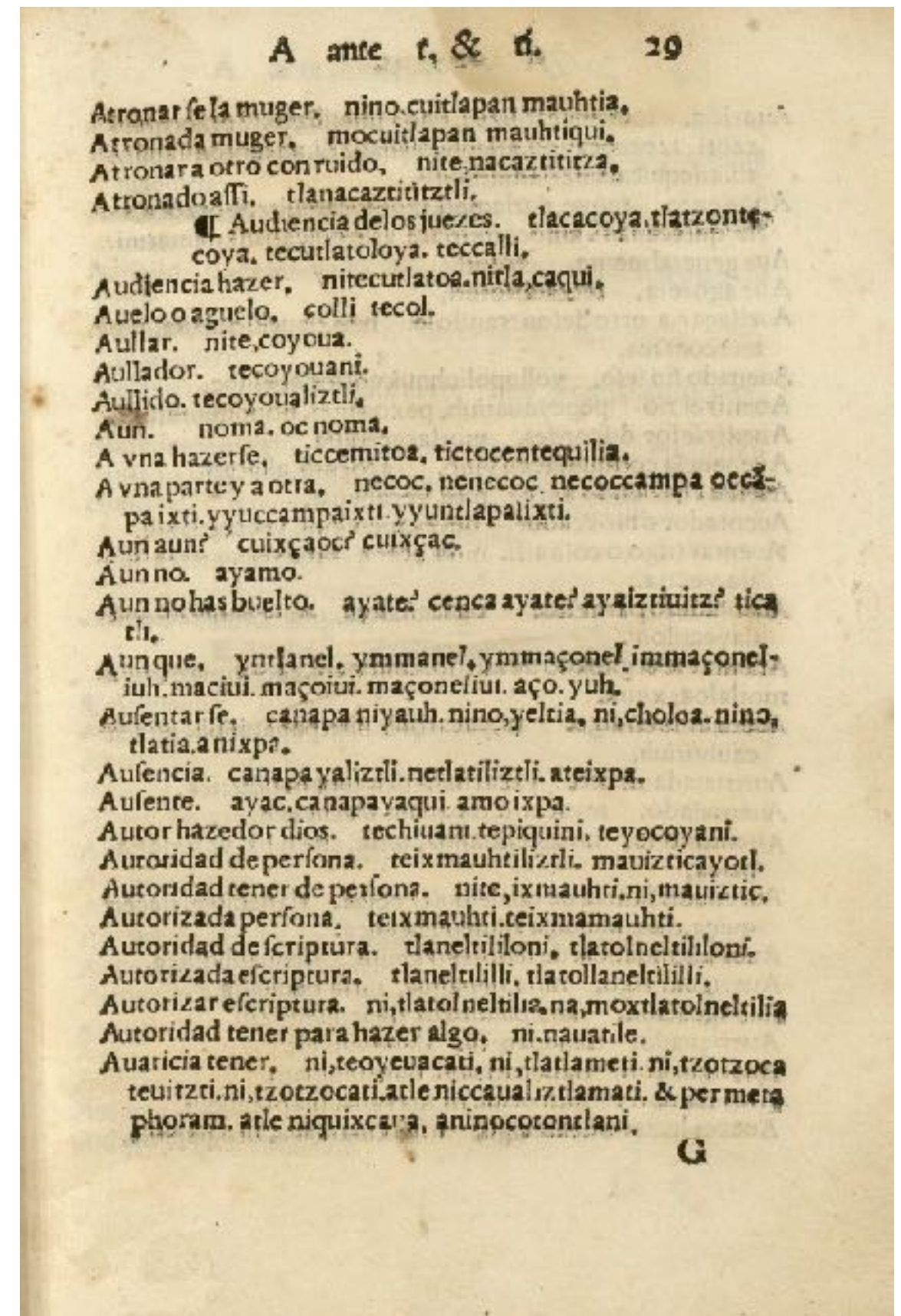
Evaluating Colonial Corpora



Evaluating Colonial Corpora

?? ??

Arronarse la muger, ni no, cuitlapan mauhtia.
 Atró nada muger, mocuitlapan mauhtiqui.
 Atronar á otro con ruido, ni te nacaztititza,
 á tronado allí, tlanacaztititzli,
 ā[Audiencia de los Juezes. tlacaco ya
 coya. tecutlatoloya, teccalli,
 Audiencia hazer, ni tecutlatoa. nitla-caqui.
 Aue io o aguelo. colli tecol.
 Aul lar, nite, coycoa.
 At illa dor. tecoyouatit-
 Auilido. tecoyoualiztli,
 Aun. Jsoma, oc noma.
 A yn a hazerfe, ticcemitoa, tictocentequilia,
 Ayna parte y a otra, necoc, nenecoc, necocca
 pa ixti yyuccampaixti yyuntlapalixti.
 Aun aun̄cuix ça oc̄c97xçac.
 Aun no. ayamo.
 Aun no ha abuelto, ayatēccen casayatēcayaíz
 tli.
 Aun que, yntla nel, ymmanel, ymmaçonel, ii
 iuh maciui, maçoiui, maçoneltin, aço yuh.
 Aufentat fe. canapa niyauh. nino-yeltia, ni-cl
 tlatia. an̄xpa,
 Aufencia. canapa yaliztli, netlatiliztli, atēixpa,
 Aufente. ayac. canapa yaqui amo ixpa.
 Autor hazedor dios. techiuani. tepi quini, te yc
 Autoridad de persona, teixmauhtiliztli. mau
 Autoridad tener de persona, ni tesix mauhti, ni
 Autorizada persona, teixmauhti, teixmamauj
 Autoridad de feriptura. tlaneltililoni, tlatolti
 Autorizada escriptura, tlaneltillili, tlatollanel
 Autoriza referiptura, ni-tlatol neltilia. na, moxi
 Autoridad tener para hazer algo, ni nauatile,
 Auaticia tener, ni ste oyeuacatl, nist latjanteti,
 teuitzti-ni, tzotzocati. atle niccaualiztlamati.
 photam. atle niquixcatia. an̄nocotontlani.



Evaluating Colonial Corpora

Atronada muger, mocuitlapan mauhtiqui.

Atró nada muger, mocuitlapan mauhtiqui.

Atronar a otro con ruido, nite, nacaztititza.

Atronar á otro con ruido, ni te nacaztititza,

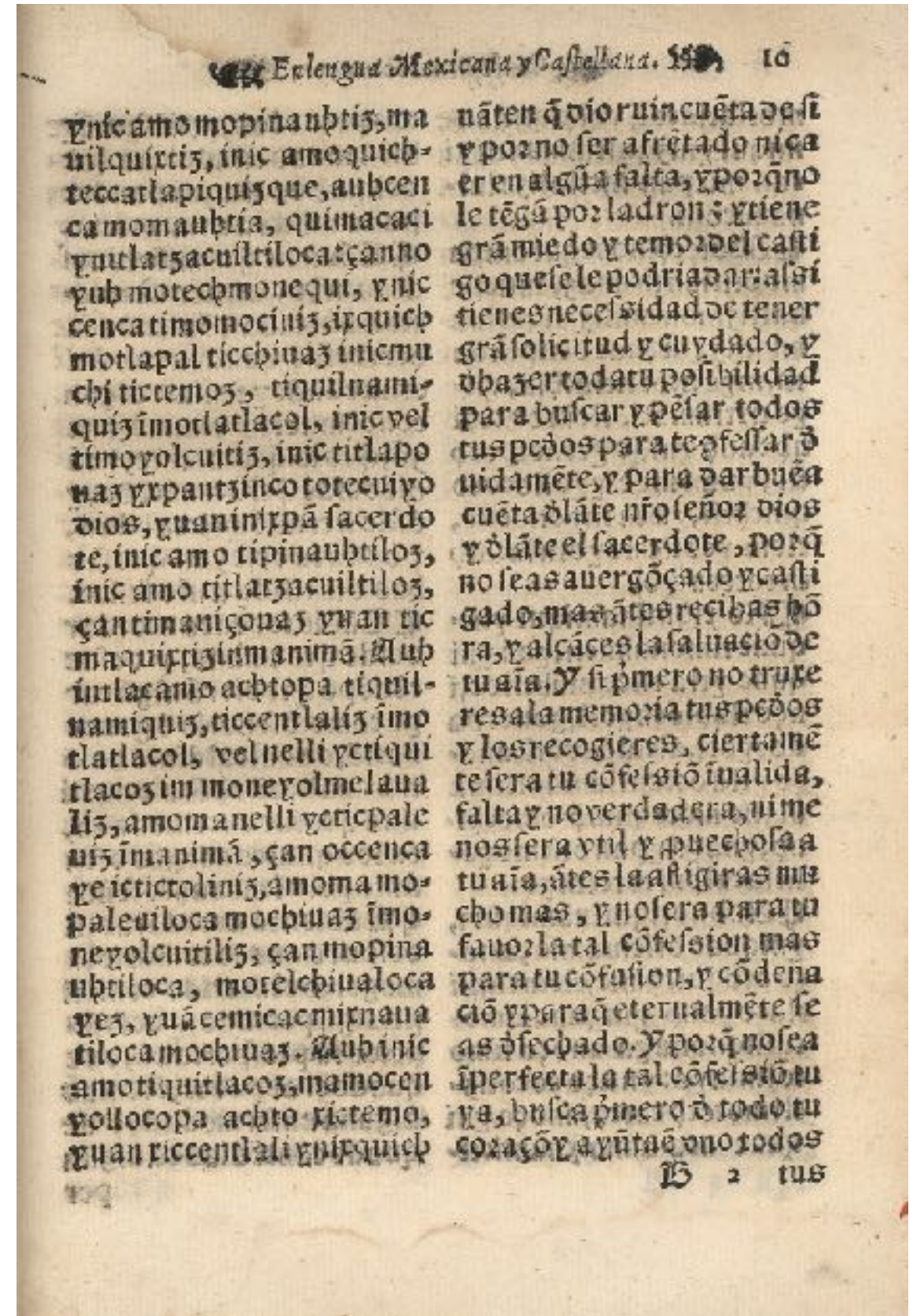
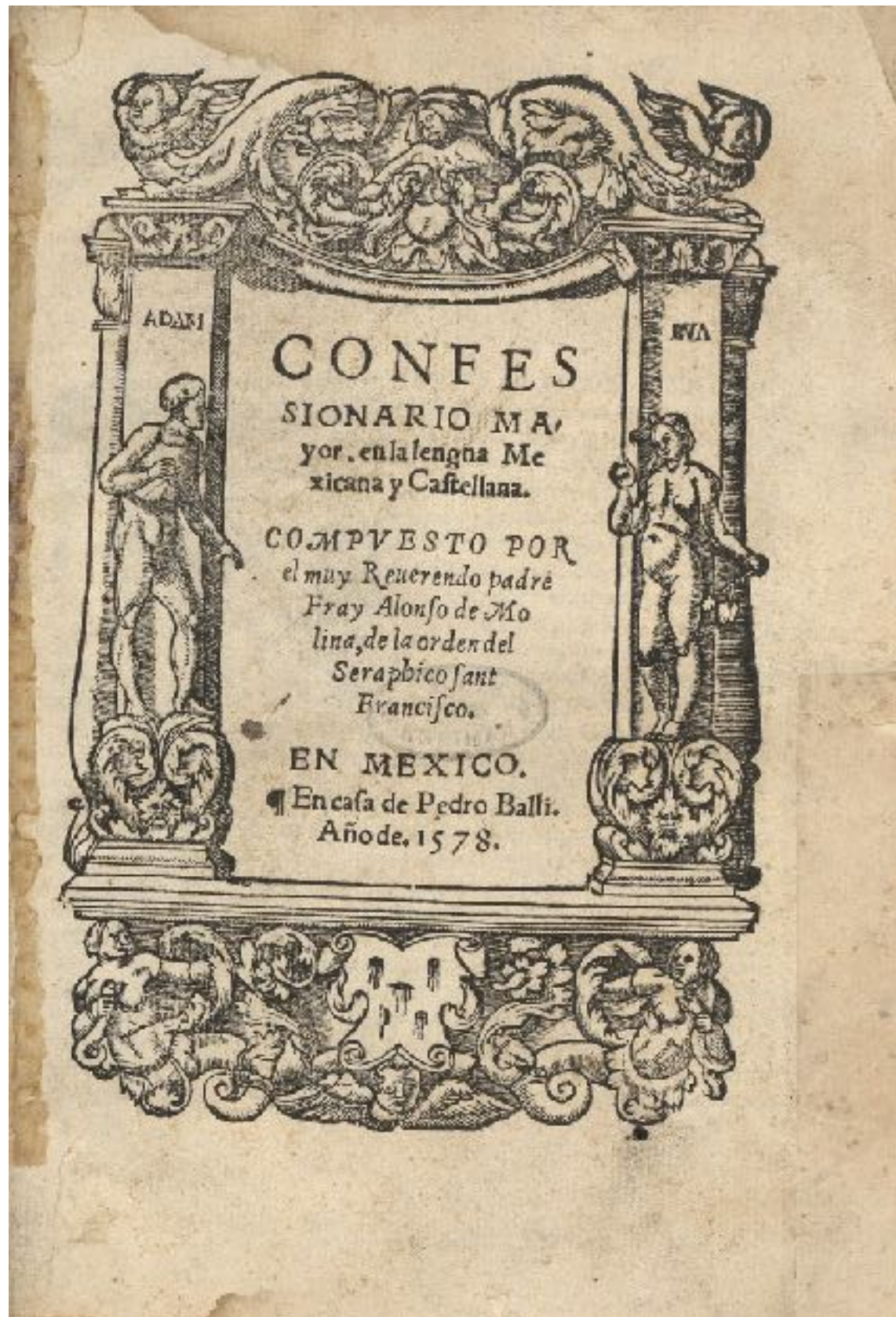
Atronado allí, tlanacaztititztli.

á tronado allí, tlanacaztititztli,

¶ Audiencia de los juezes. tlacacoya

ã[Audiencia de los Juezes. tlacaco ya

Evaluating Colonial Corpora



Evaluating Colonial Corpora

ynic amo mopinaubtiz, ma nāten q̃dior uincuētaoe-ĩ

ynic amo mopĩnāuhtij, ma nāten. q̃Dior uincuētaoe-ĩ

uilquixtiz, inic amo quicb - y poz no fer atrétado nica

uilquirtij, inic amo quicb - y **poz no fer atrétado** nica

teccatlapiquique, aubcen er enalgua falta, y poz q̃no

teccatlapiquique, auh cen **eren** **algua falta**, **y poz** q̃n.o

camomaubtia, quimacaci le tēgā poz ladron ; y tiene

ca momaubtia, quimacaci **le tēgā poz ladron ; y tiene**

Evaluation systems must be
context-specific

We need standards for evaluating
our colonial corpora.