



NewsEye International Conference
March, 17th 2021

Two Examples of Analysis of Textual Document in Oriental and Under-Resourced Languages

Chahan Vidal-Gorène

Calfa, École Nationale des Chartes-PSL

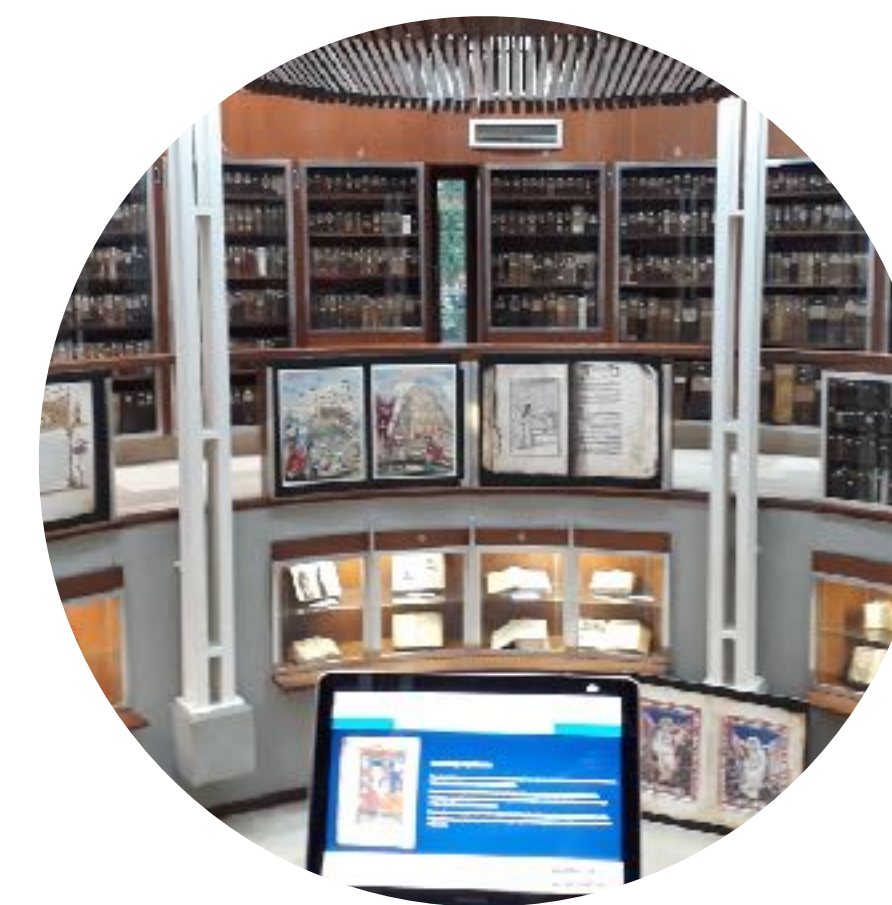




Open databases for public and
NLP Projects



Document Analysis for Under-
Resourced Languages

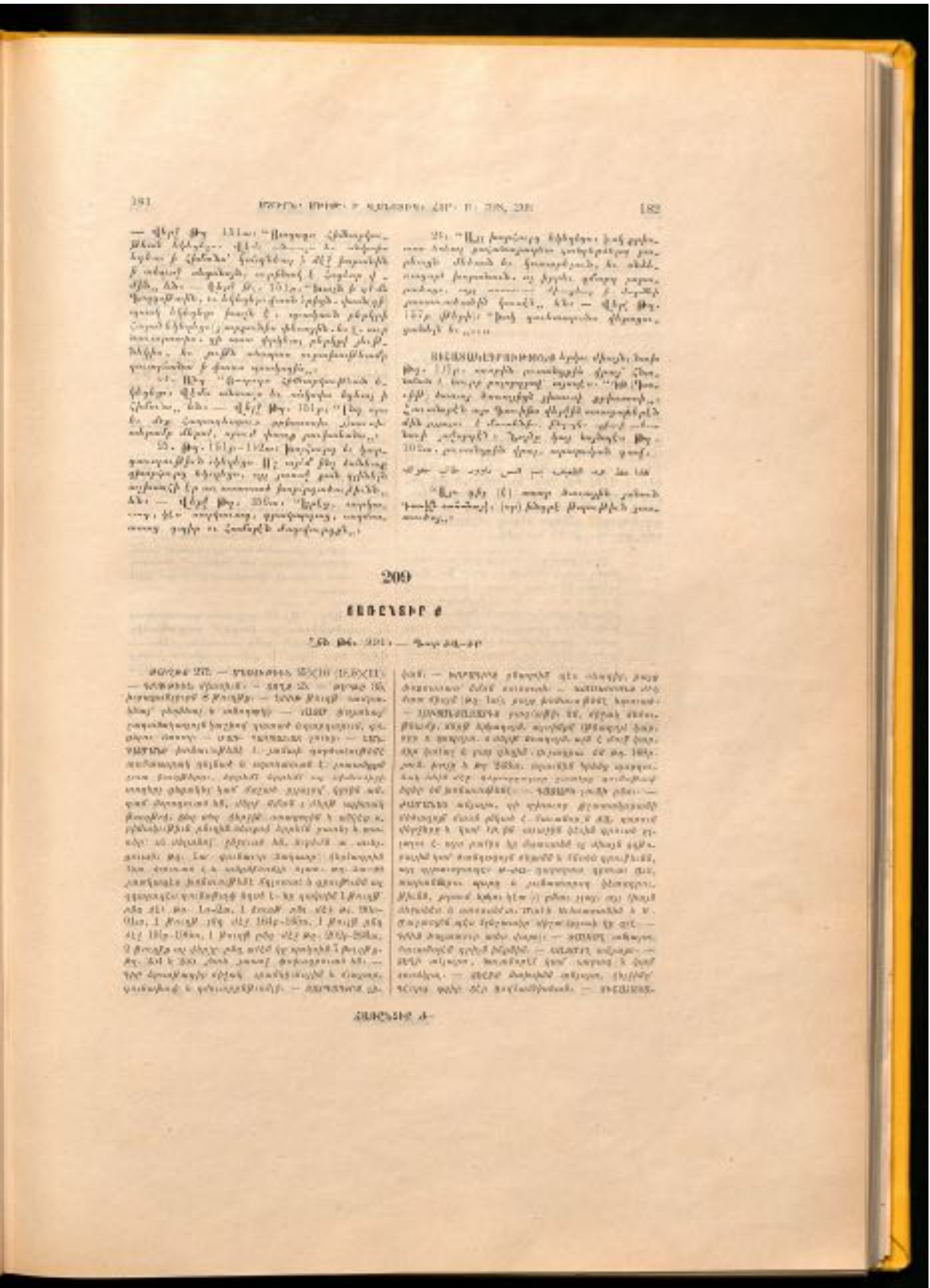


Customized IT Projects for
Libraries

Specialized in **Automatic Document Analysis**

➡ **Handwritten** and Printed **Text Recognition** (HTR / OCR) for Oriental Languages and **structure understanding**

Ongoing projects for Armenian Newspapers

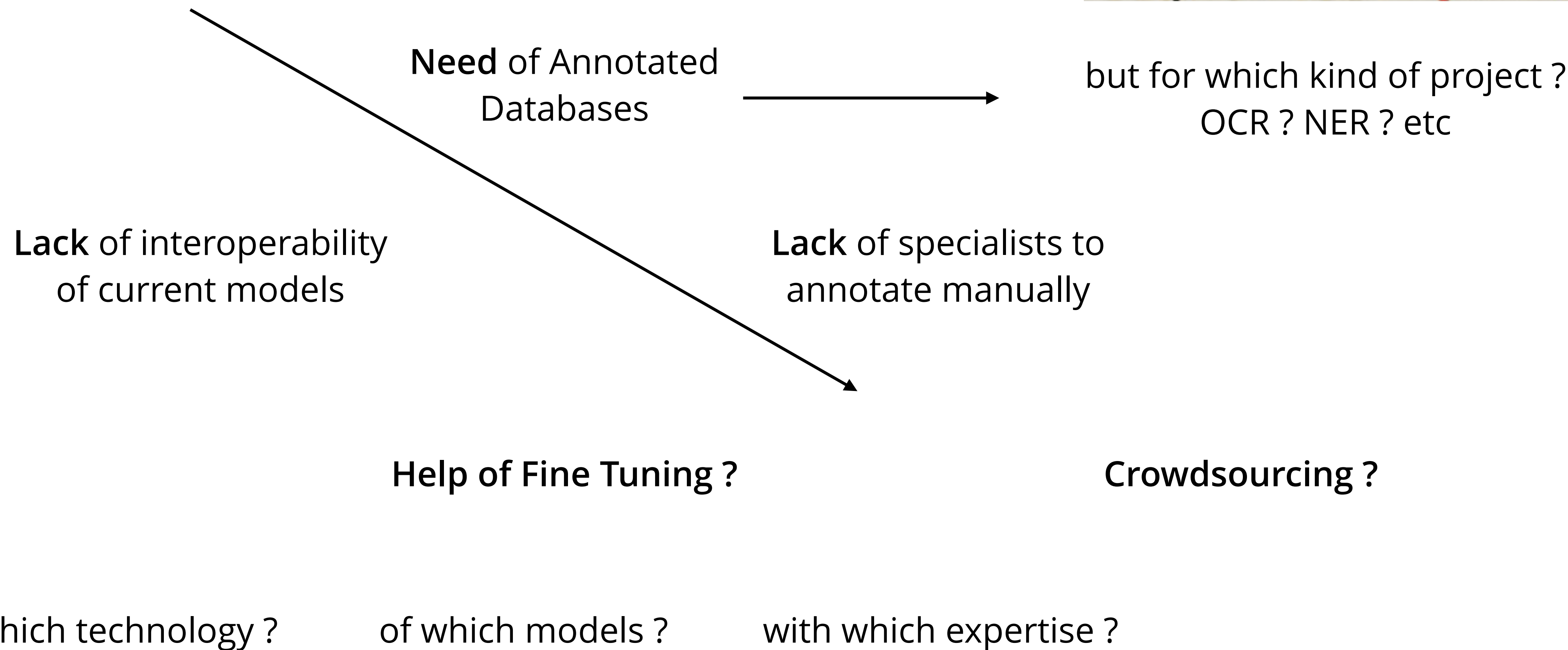
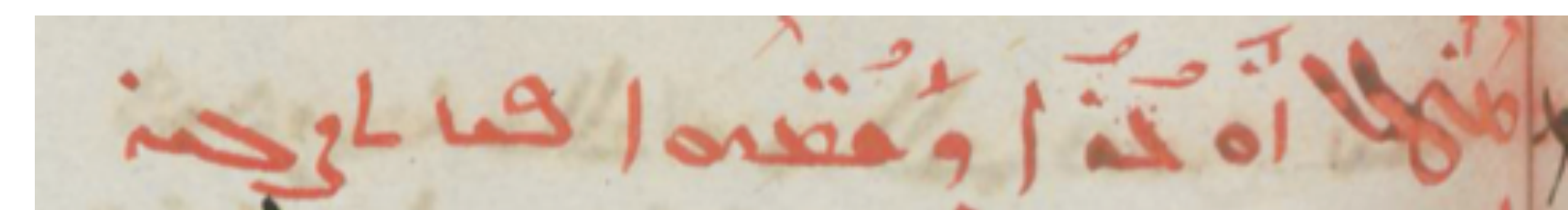
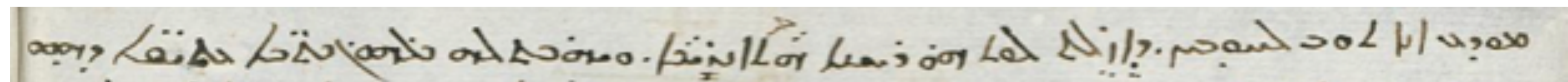
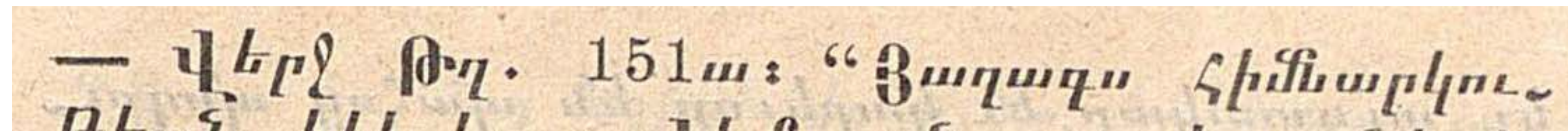


<https://webaram.com/biblio/presse/haratch-յառաջ>
Digitized by BULAC and ARAM association
with the support of the Gulbenkian Foundation

<https://arar.sci.am/dlibra/collectiondescription/10>
Digitized by the Fundamental Scientific Library
of the NAS RA

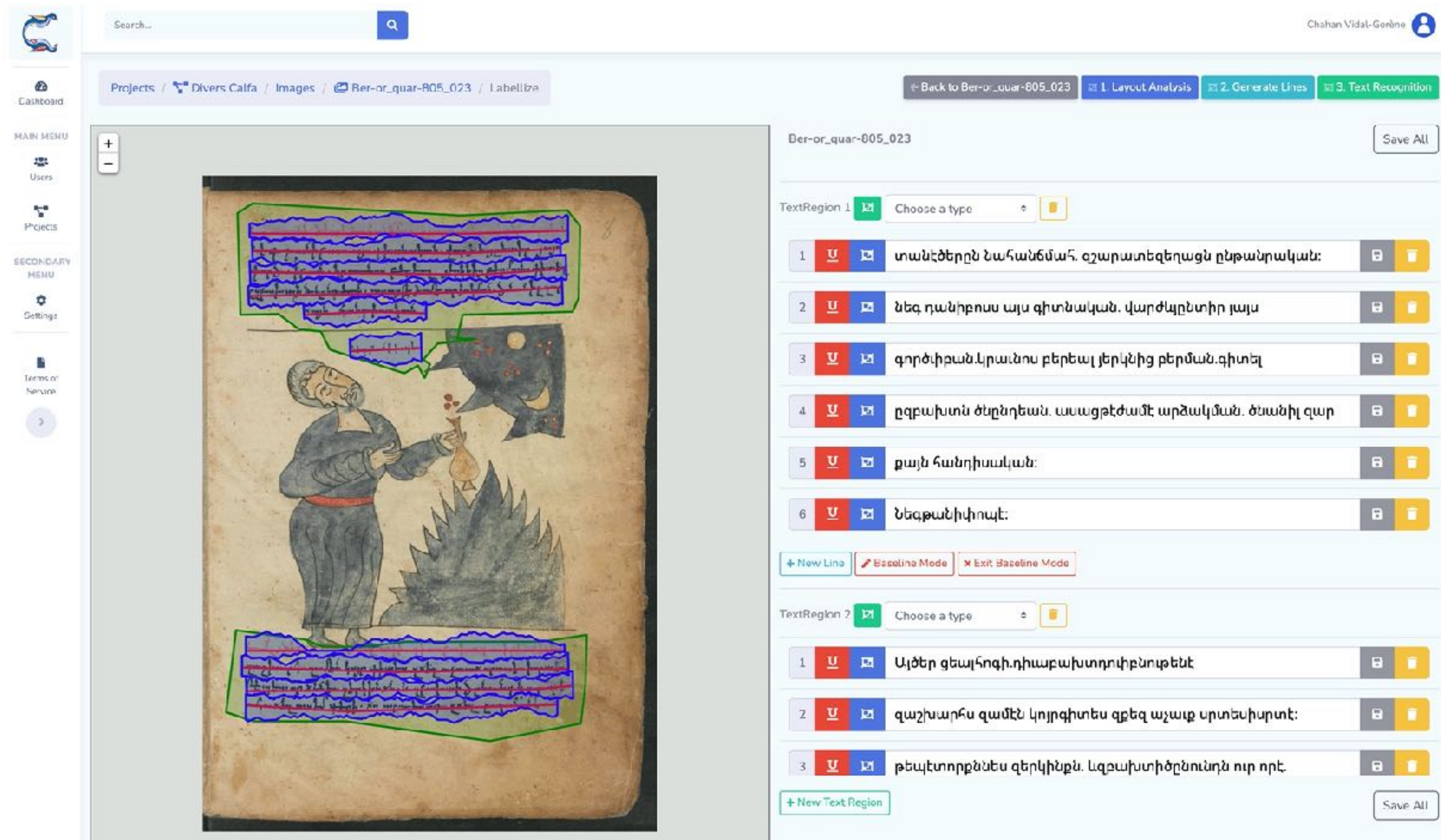
Previously OCRized, with noise in Layout extraction and CER ~ 85% (Tesseract)

Issues with Under-Resourced Languages (or very specific layouts)



<https://calfa.fr/vision>

Calfa Vision is a web-based annotation tool for documents and images, collaborative and free

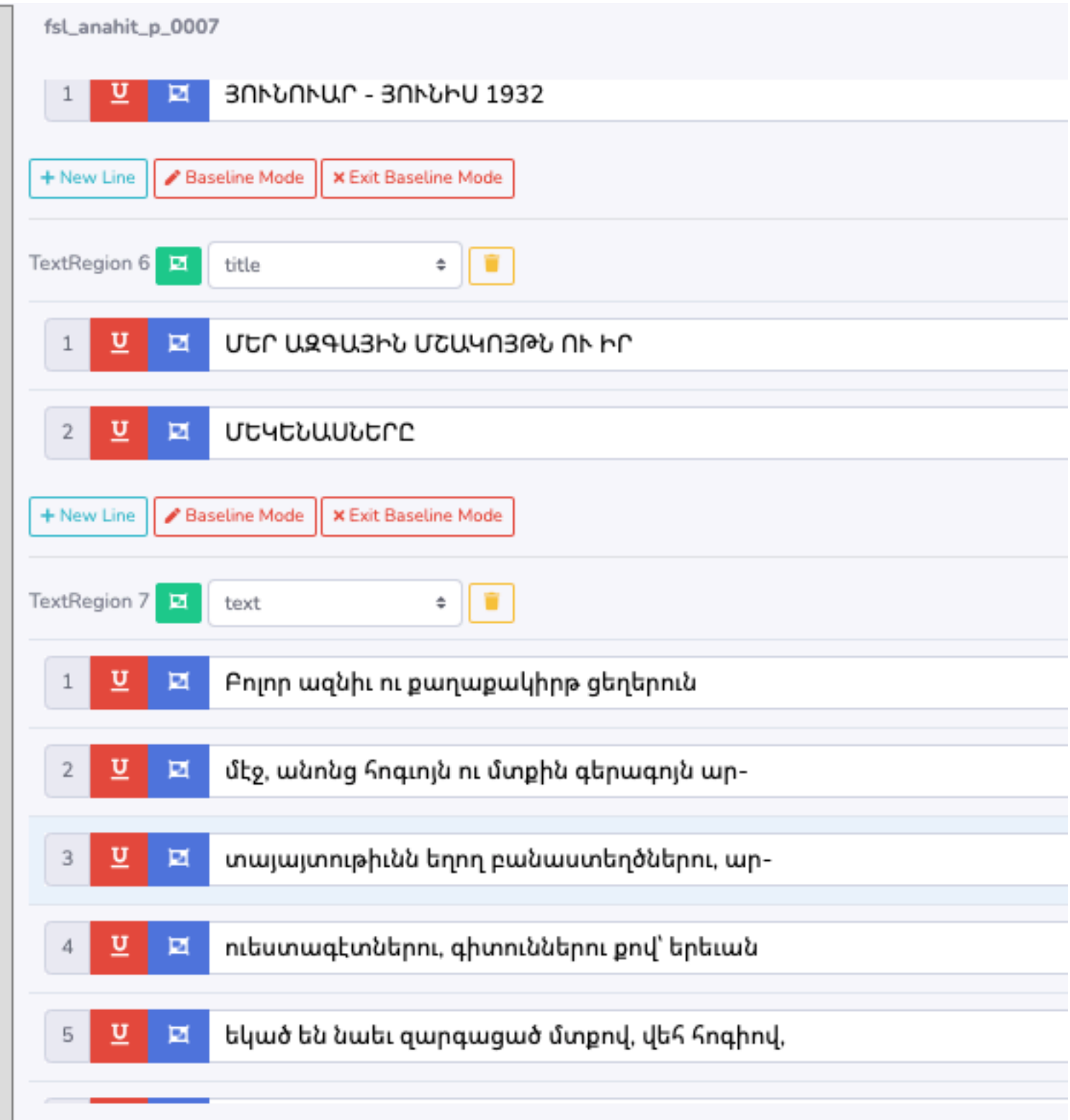
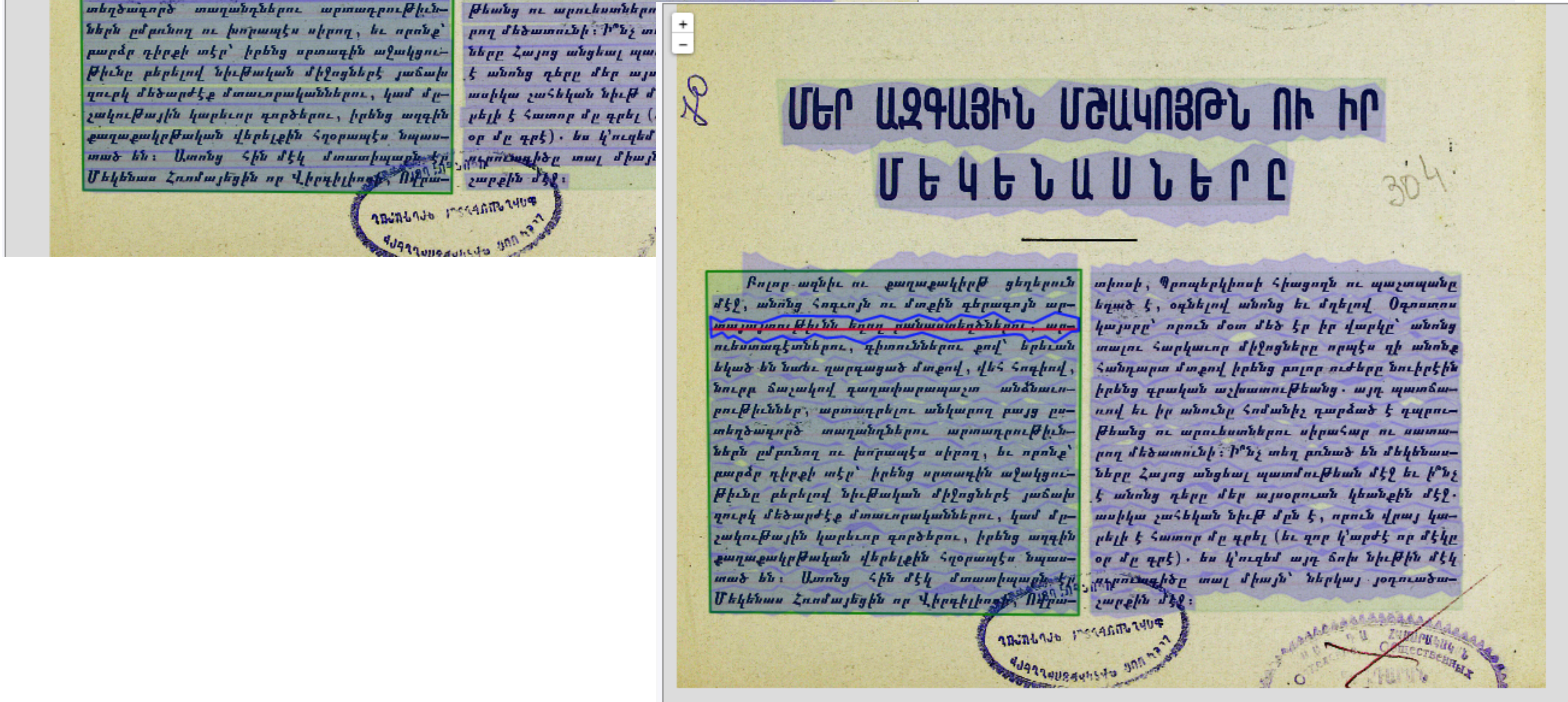


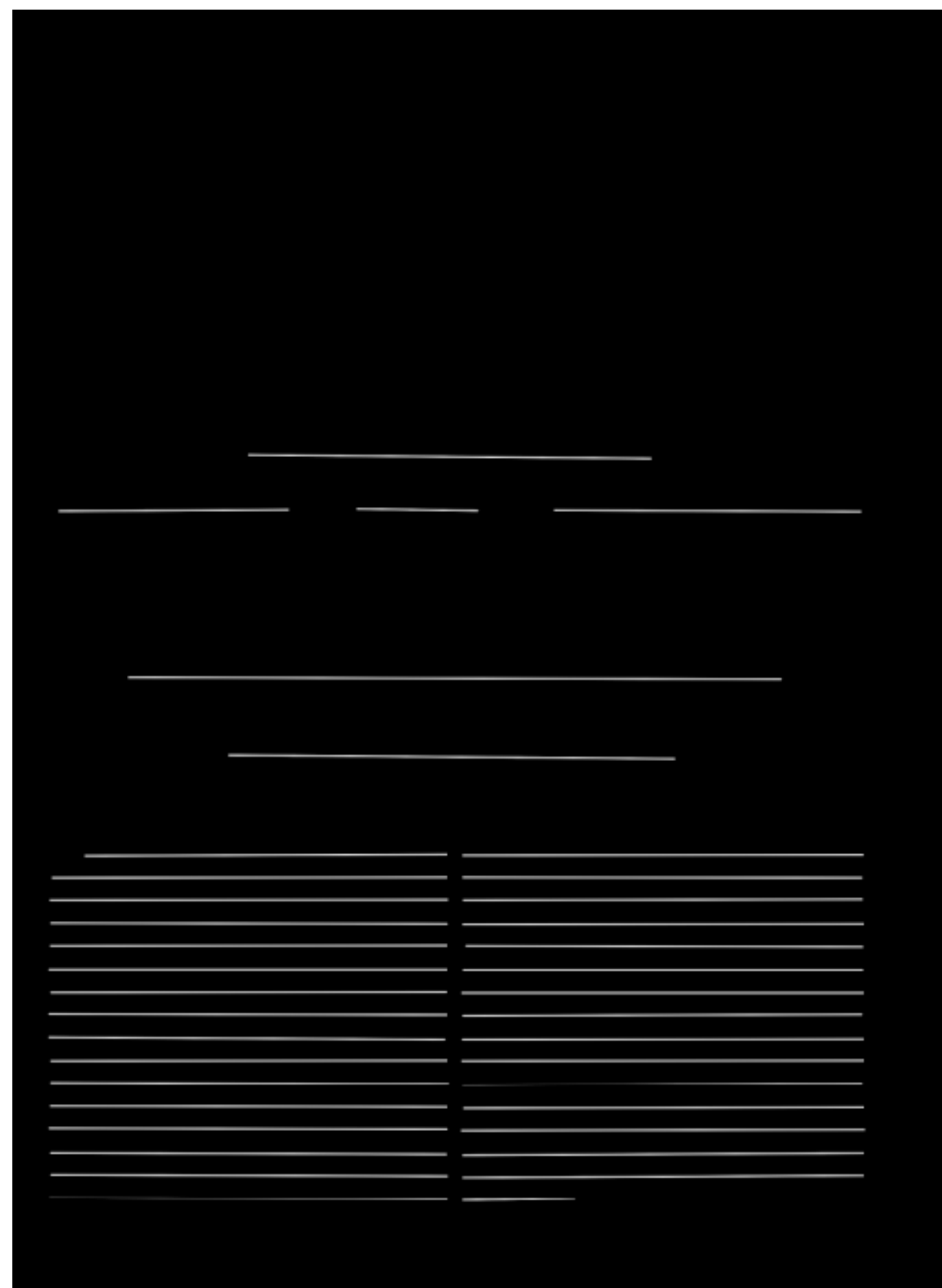
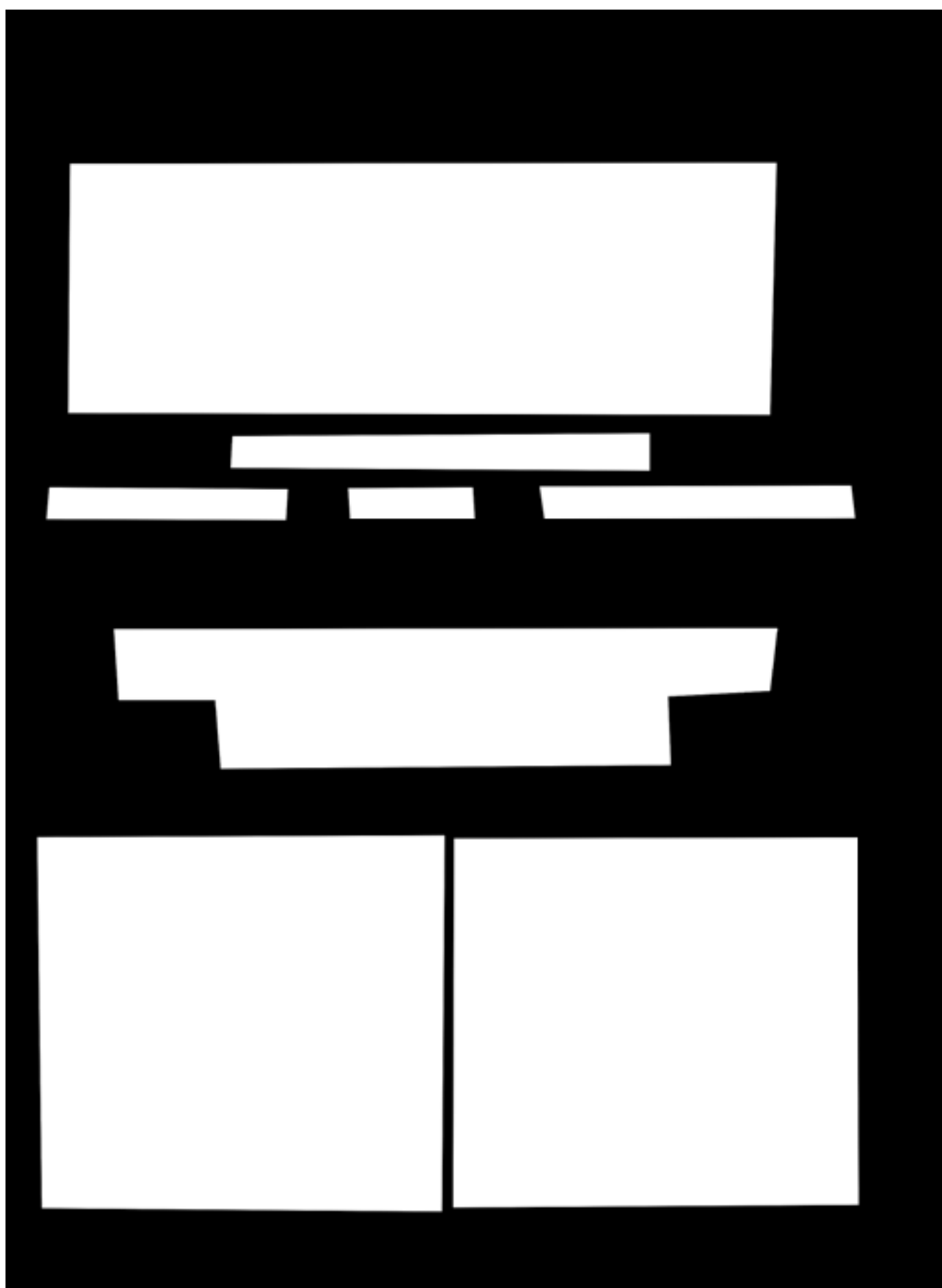
- Originally developed for Historical Manuscripts in Oriental Languages
- Includes automated tasks and fine-tunable models in real time

=> applied to process newspapers



Automatic analysis (layout analysis and OCR) of Anahit Journal, Fundamental Scientific Library (Armenia)

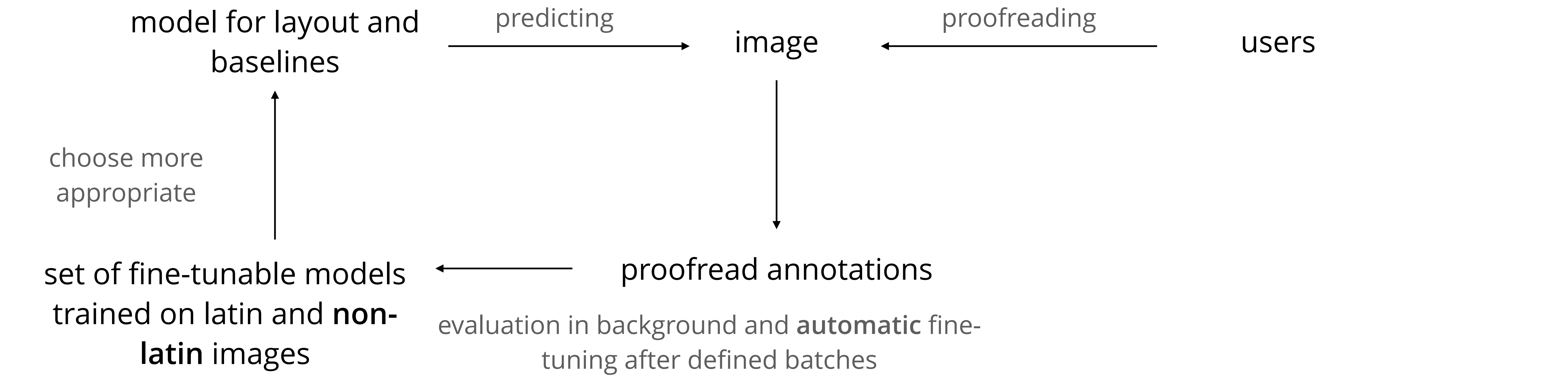




```
<?xml:version="1.0"?>
<PcGts xmlns="http://schema.primaresearch.org/PAGE/gts/pagecontent/2013"
  <<Metadata>
    <<Creator>Calfa</Creator>
    <<Created>2021-03-17T10:33:01+00:00</Created>
  </Metadata>
  <<Page imageFilename="fsl_anahit_p_0007.jpg" imageHeight="2938" imageWidth="1920">
    <<TextRegion id="8398" custom="structure-{type:heading;}">
      <<Coords points="149,353-1766,351-1752,927-145,923-149,353"/>
      <<TextEquiv>
        <<Unicode></Unicode>
      </TextEquiv>
    </TextRegion>
    <<TextRegion id="8399" custom="structure-{type:title;}">
      <<Coords points="521,977-1476,971-1476,1053-517,1047-521,977"/>
      <<TextLine id="99443">
        <<Coords points="539,965-576,961-620,971-683,964-709,951-744,972"
          <<Baseline points="539,1020-1468,1028"/>
        <<TextEquiv>
          <<Unicode>ՀԱՆԴԵՍ ՄՏԱԾՄԱՆ ԵՒ ԱՐՈՒԵՄՏԻ</Unicode>
        </TextEquiv>
      </TextLine>
      <<TextEquiv>
        <<Unicode>ՀԱՆԴԵՍ ՄՏԱԾՄԱՆ ԵՒ ԱՐՈՒԵՄՏԻ</Unicode>
      </TextEquiv>
    </TextRegion>
    <<TextRegion id="8400" custom="structure-{type:numeration;}">
      <<Coords points="99,1097-643,1101-639,1169-93,1167-99,1097"/>
      <<TextLine id="99444">
        <<Coords points="104,1104-157,1091-194,1094-207,1104-231,1095-26"
          <<Baseline points="104,1147-632,1144"/>
        <<TextEquiv>
          <<Unicode>Գ. ՏԱՐԻ, ԹԻՒՆ 5-6</Unicode>
        </TextEquiv>
      </TextLine>
      <<TextEquiv>
        <<Unicode>Գ. ՏԱՐԻ, ԹԻՒՆ 5-6</Unicode>
      </TextEquiv>
    </TextRegion>
    <<TextRegion id="8401" custom="structure-{type:footnote;}">
```

Output : text, ALTO, pageXML, customized json, masks, pairs GT

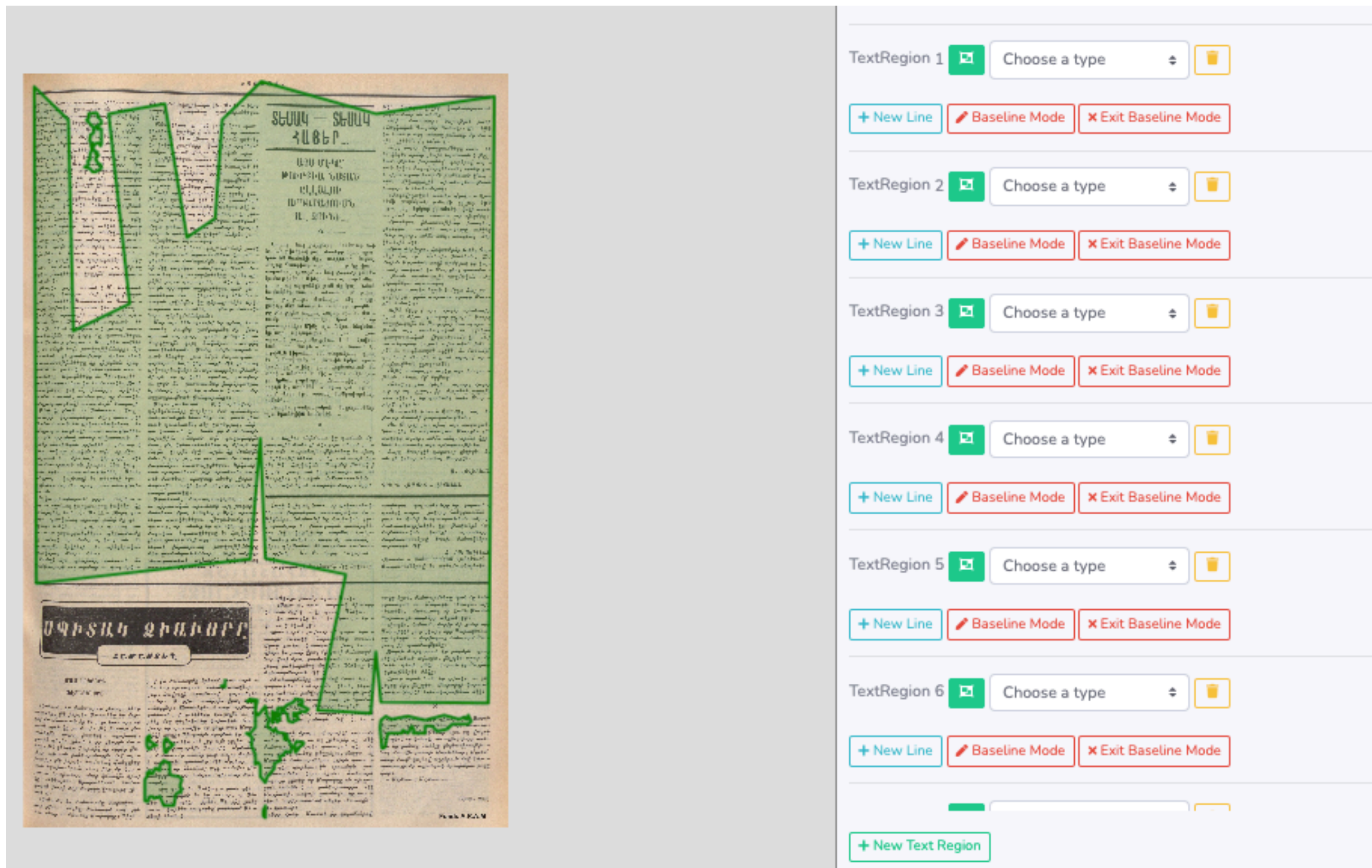
Pipeline to reduce OCR/HTR correction and proofreading



Basically, we could also considered that damaged documents consist in themself a under-resourced variation.

Dataset	Precision (%)	Recall (%)	F1-score (%)
BADAM	0.8253	0.8541	0.8395
cBAD complex track (ICDAR 2017)	0.9071	0.9052	0.9061
cBAD simple track (ICDAR 2017)	0.9511	0.9538	0.9525
cBAD (ICDAR 2019)	0.9312	0.9310	0.9311
READ (ICFHR 2016)	0.9590	0.9888	0.9737
Training Dataset	0.9531	0.9878	0.9701

Vidal-Gorène et al, 2021, *submitted*

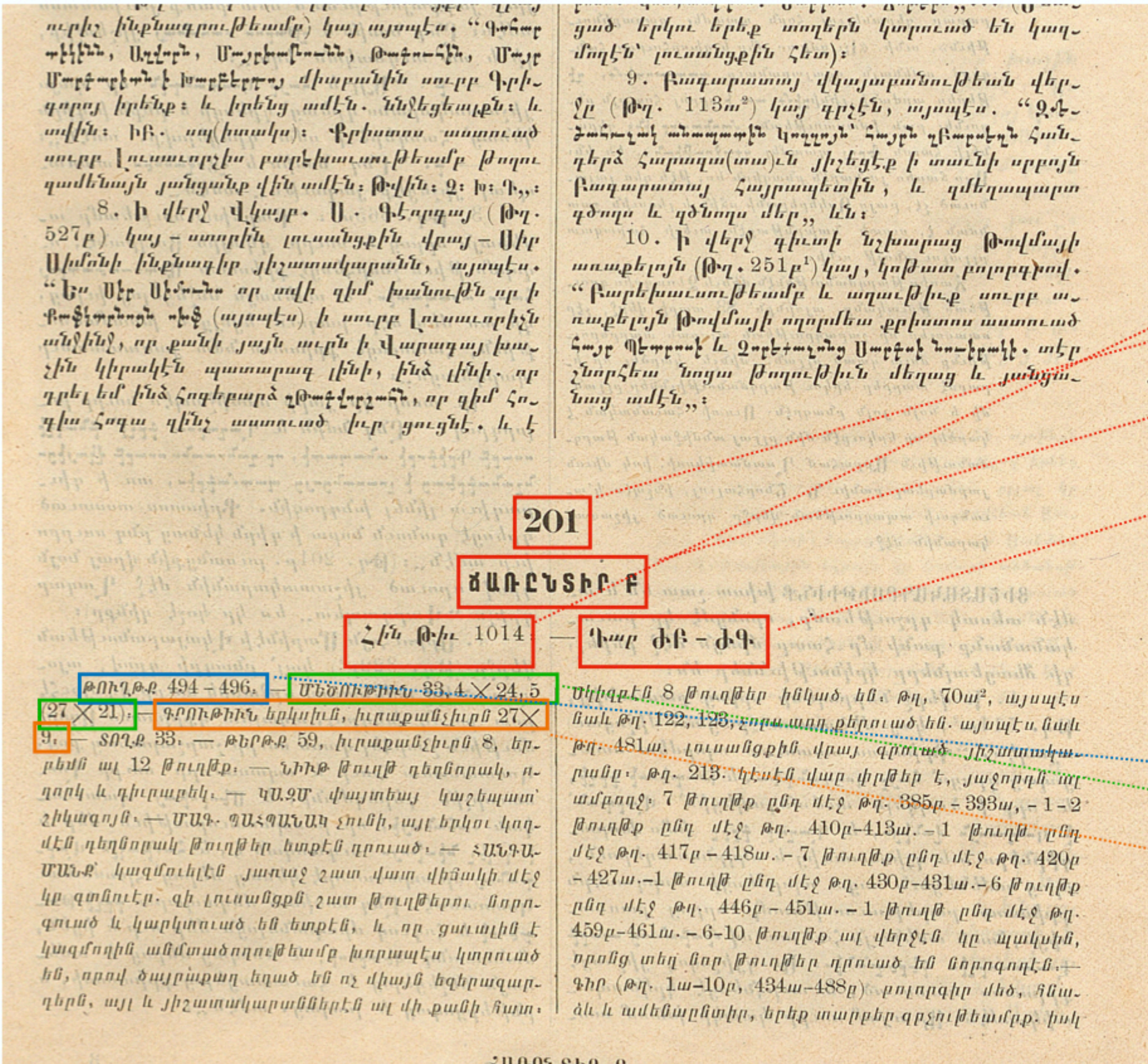


default model, 1st image



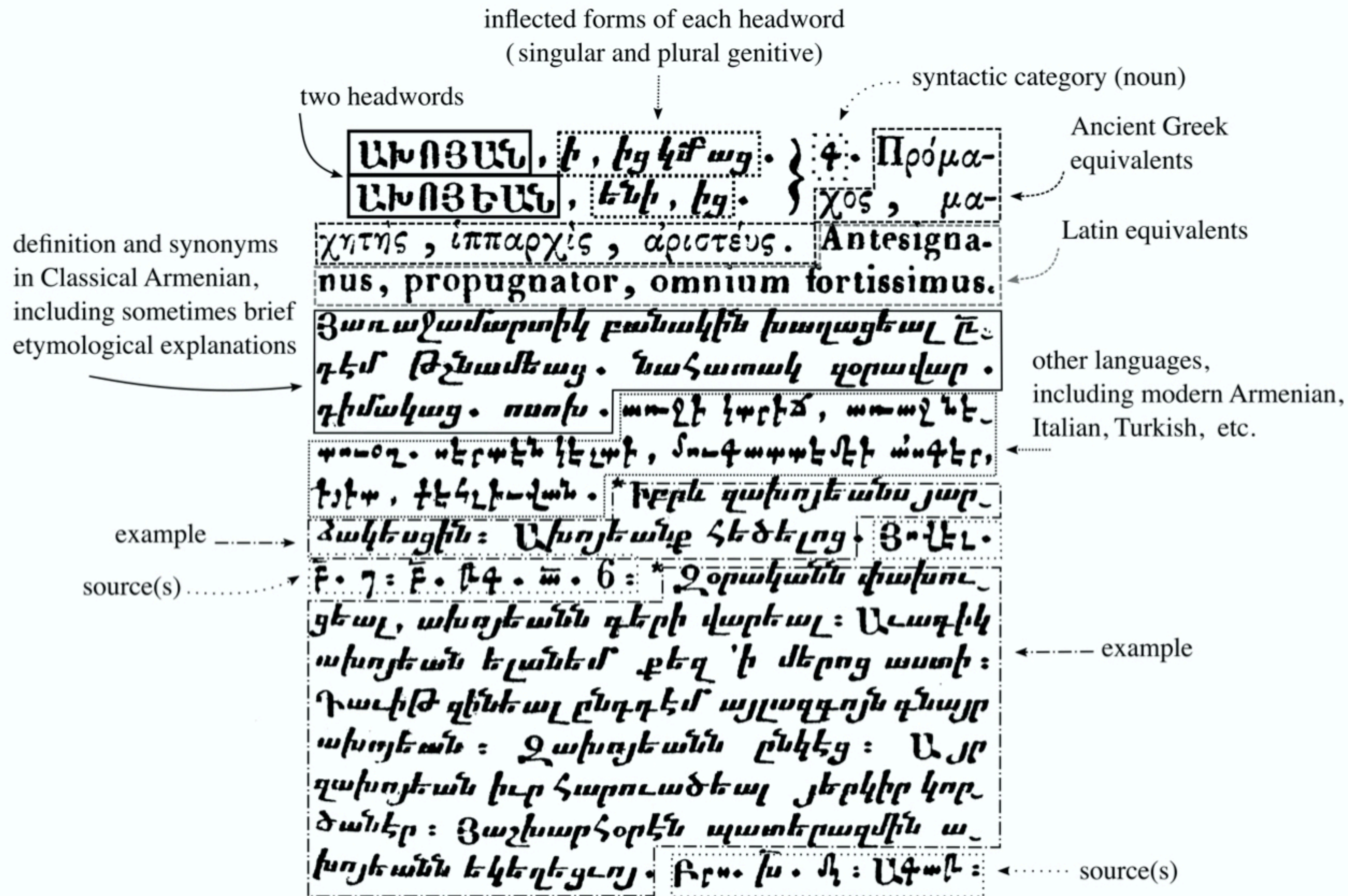
fine-tuned model (batch 20)

On average, text class pixel accuracy 87,03 % - results of fine-tuning on Calfa Vision after 20 pages manually proofread



IDENTIFIANT	V1014
NUMÉRO	1014
NUMÉRO DE NOTICE	201
TITRE	ՃԱՌԸՆՏԻՐ Բ
DATE DE DÉBUT	XII
DATE DE FIN	XIII
DÉTAILS DE DATE	անյայտ, սակայն դատելով ներքին հանգամանքներէն հաւանաբար առաջին կեսին, և կամ ԺԳ. թղին առաջին քառորդին. վասն զի յիշատակարաններն (տես Թղ. 123թ, 114թն). կը կրեն առաջին -1377) թուականները, և որոնք համեմատութեամբ երկաթա-
LIEU DE COPIE	անյայտ, հաւանօրէն Ս. Աստուածածին կամ Ս. Յովհաննէս վ.
RÉFÉRENCE	t. II, 1924, c. 35-66
GENRE	
NOMBRE DE PAGES	494-496
DIMENSIONS	33,4X24,5 (27X21)
MISE EN PAGE	երկսիւն, իւրաքանչիւրն 27X9
TYPE D'ÉCRITURE 1	

Catalog of Armenian Manuscripts of Venice
5 000 images processed with a model of layout analysis (fine-tuned from newspaper’s model) and OCR content-based segmenter

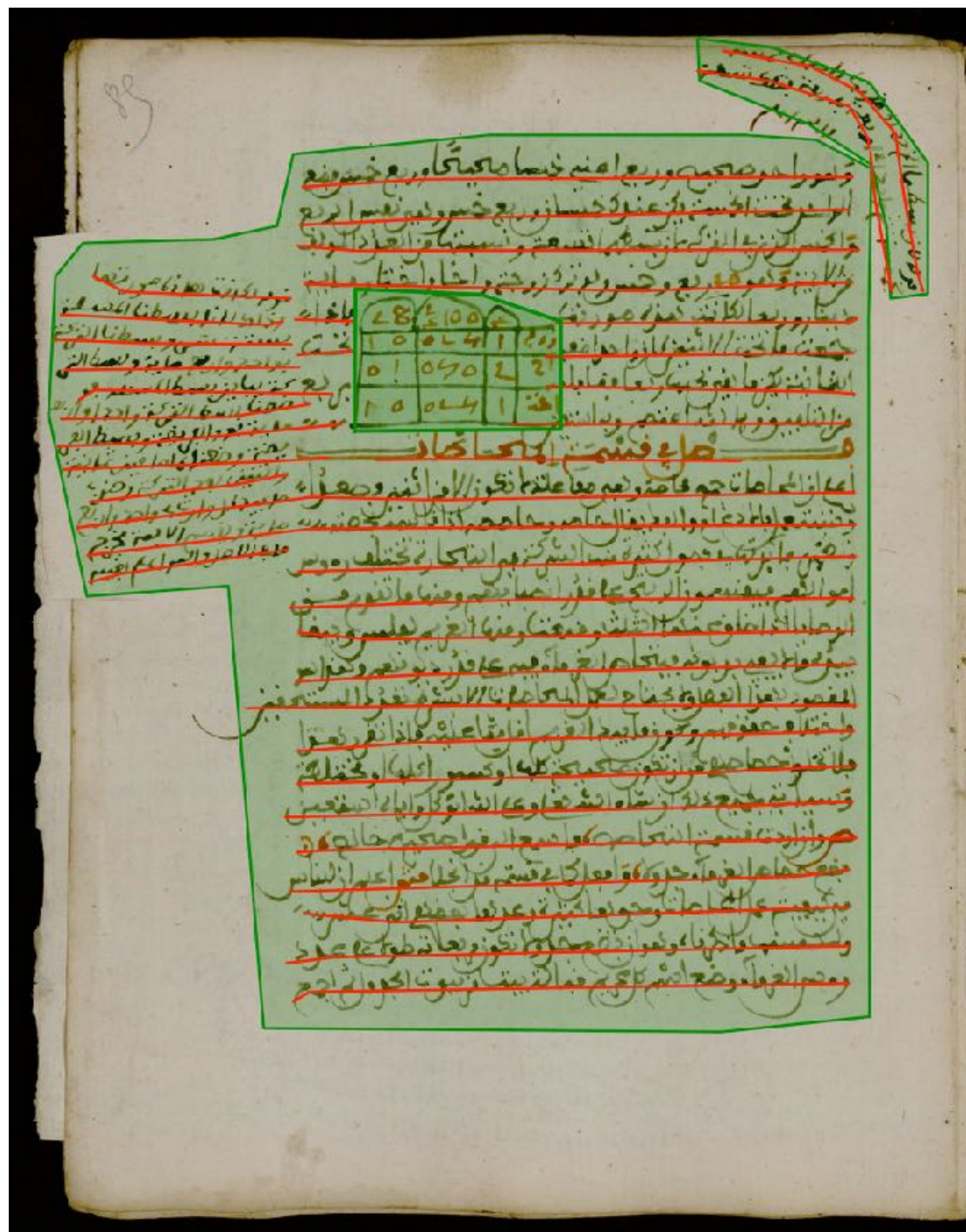


Vidal-Gorène et al, JSAS 27, 2021

Ախոյան, ի, ից, աց

Article NBHL Étymologie Mots dérivés

- ♦ ԱԽՈՅԱՆ. **πρόμαχος, προμαχητής, ἰππαρχίς, ἀριστεύς**. Յառաջամարտիկ բանակին խաղացեալ ընդդէմ թշնամի կտրիճ, առաջ նետուող. սերտէն կէշտի, մուգատտէմէի.
- ❖ «Իբրեւ զախոյեանս հարձակեսցին: Ախոյեանք հարձակեցին:» (Յար.)
- ❖ «Զօրականն փախուցեալ, ախոյեանն գերի վարեցաւ զինեալ ընդդէմ այլազգոյն գնայր ախոյեան: Զախոյեան կործանէր: Յաշխարհօրէն պատերազմին ախոյեան կը լայնէր: Շար.:»
- ❖ «Ընդդէմ խրոխտ ախոյեանին (կամ ախոյանին): Որքան չործանին վատանուն յախոյանէ անտի: Ի փախուցեալ ախոյանից: Նախամարտիկ եւ առաջին ախոյանից: Եզնիկ.: Պիտ.: Սարկ. քհն.:
- ❖ «Մատնեցաւ ի ձեռս ախոյեան թշնամոյն.» (Խոսր.)
- ❖ «Առ զյաղթութիւն ի վերայ ախոյանին իւրոյ ... Յախոյեանն եկեղեցւոյ. Բր. խ. Պի: Ագ.:
- ♦ Նմանութեամբ ասի.
- ❖ «Զօրեղ ախոյեանս աստուածականս: Ախոյեան ներհակաց.» (Նար.):
- ♦ «Լեռնայն ախոյեանն իւր զօրս ախոյեանն: Որքան չործանին վատանուն յախոյանէ անտի: Ի փախուցեալ ախոյանից: Նախամարտիկ եւ առաջին ախոյանից: Եզնիկ.: Պիտ.: Սարկ. քհն.:



Hackathon with the BULAC library and GIS MOMM from January to April 2021 with 15 volunteers

Database of 250 images and annotations will be released in June 2021

Fine-tuning of a model of newspapers with three classes : text, tables and marginalia. Shown example corresponds to the batch 75

From 16/02/2021 to 16/03/2021 for Armenian Newspapers and Maghrebi Arabic manuscripts :

- 1,081 images have been annotated and proofread
- 1,875 text-regions
- 12,551 lines of text

Thank you



To know more about the project of Arabic manuscripts : <https://www.bulac.fr/espace-recherche/la-bulac-et-la-recherche/>

About the platform : <https://calfa.fr/vision> and process : <https://www.aclweb.org/anthology/2020.lrec-1.385.pdf>



FONDATION
CALOUSTE
GULBENKIAN

Fondation des Frères
Ghoukassiantz



Moyen-Orient et
Mondes Musulmans
Groupement d'Intérêt Scientifique

