# Information Extraction from Handwritten Tables in Historical Documents

José Andrés[1], Jose Ramón Prieto[1], Emilio Granell[1], Verónica Romero[2], Joan Andreu Sánchez[1], Enrique Vidal[1]

[1]Pattern Recognition and Human Language Technology Research Center
Universitat Politècnica de València, Spain

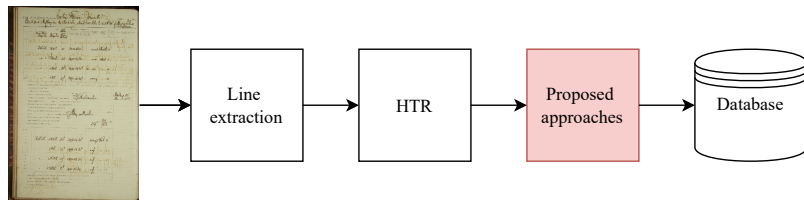[2]Departament d'Informàtica, Universitat de València, Spain.

May 23rd, 2022

PRHLT

# Introduction

# General pipeline

# Heuristic Geometric Information

# Log-linear model



Textlines with
HTR transcription

Grouping textlines → Replacing quotation marks → Performing information extraction → Database

# Log-linear model: grouping textlines



Grouping textlines into semantic cells

Merging semantic cells

# Log-linear model: replacing quotation marks



Substitute the quotation marks by the content of the precedent cell

Probability of being column header

# Log-linear model: performing IE



Probability of being aligned vertically

Probability of being a row header

Probability of being aligned horizontally

Graph of
textlines

| Identify columns, rows and headers | Performing information extraction | Database |

Created a graph by connecting the textlines by their line of sight

# Graph Neural Networks

# Graph Neural Networks

# Jeannette Statistics

Table: Basic statistics of the HisClima database for the three partitions.

| Number of: | Train | Validation | Test | Total |
|------------|------:|-----------:|-----:|------:|
| Pages | 143 | 15 | 50 | 208 |
| Lines | 23 617 | 2 284 | 7 838 | 33 739 |
| Running words | 46 599 | 4 604 | 15 611 | 66 814 |
| Lexicon | 1 287 | 491 | 924 | 1 483 |
| Character set size | 76 | 76 | 76 | 76 |

Table: Results of text recognition.

| Text type | Manuscript | Printed | Overall |
|-----------|:----------:|:-------:|:-------:|
| WER | 10.4% | 1.8% | 4.4% |

# Information Extraction

Table: Information extraction results. 95% confidence intervals are never larger than 0.01 when using GT lines and 0.02 when employing automatic lines.

| Lines | Ground Thruth | | | Automatic | | |
| Metric | P | R | $F_1$ | P | R | $F_1$ |
|---|---|---|---|---|---|---|
| Heuristic Geometric Information | 0.79 | 0.78 | 0.78 | 0.64 | 0.55 | 0.59 |
| Log-linear Model | 0.87 | 0.79 | 0.83 | 0.77 | **0.69** | **0.73** |
| Graph Neural Network | **0.88** | **0.83** | **0.85** | **0.78** | 0.67 | 0.72 |
| Oracle | 0.89 | 0.88 | 0.89 | 0.79 | 0.72 | 0.76 |

# Conclusions

- Machine learning approaches outperform the heuristic geometric information method.

- Log-linear model and graph neural networks have achieved a similar performance.

# Future works

- Improve automatic line extraction and HTR.

- Repeat this experiment employing PrIx instead of the HTR transcripts.

# Information Extraction from Handwritten Tables in Historical Documents

José Andrés[1], Jose Ramón Prieto[1], Emilio Granell[1], Verónica Romero[2], Joan Andreu Sánchez[1], Enrique Vidal[1]

[1]Pattern Recognition and Human Language Technology Research Center
Universitat Politècnica de València, Spain

[2]Departament d'Informàtica, Universitat de València, Spain.

May 23rd, 2022

**PRHLT**

# HTR performance

Table: Results of text recognition. 95% confidence intervals are never larger than 1.1% for manuscript text, 0.4% for printed text and 0.5 % overall.

| Text type | Manuscript | Printed | Overall |
|:---:|:---:|:---:|:---:|
| CER | 5.7% | 1.5% | 2.0% |
| WER | 10.4% | 1.8% | 4.4% |