

*tS* *tranSkriptorium*

**Beyond automatic transcription:**

**A better source of information on which to apply NLP with success**

# Index

- The problem
- Solutions
- Beyond Automatic Transcription
- Searching
- Transcription + Alternatives
- Metadata Extraction
- Anonymization
- Classification & Segmentation
- Big Data Analysis
- Demonstration



# The Problem

- Large collections of physical documents
- Some are still growing
- Handwritten, Typed, printed and mixed
- Structured, Semi-structured, Unstructured
- Changing Templates
- Different access needs depending on the collection
- Depending on the country some laws are in place

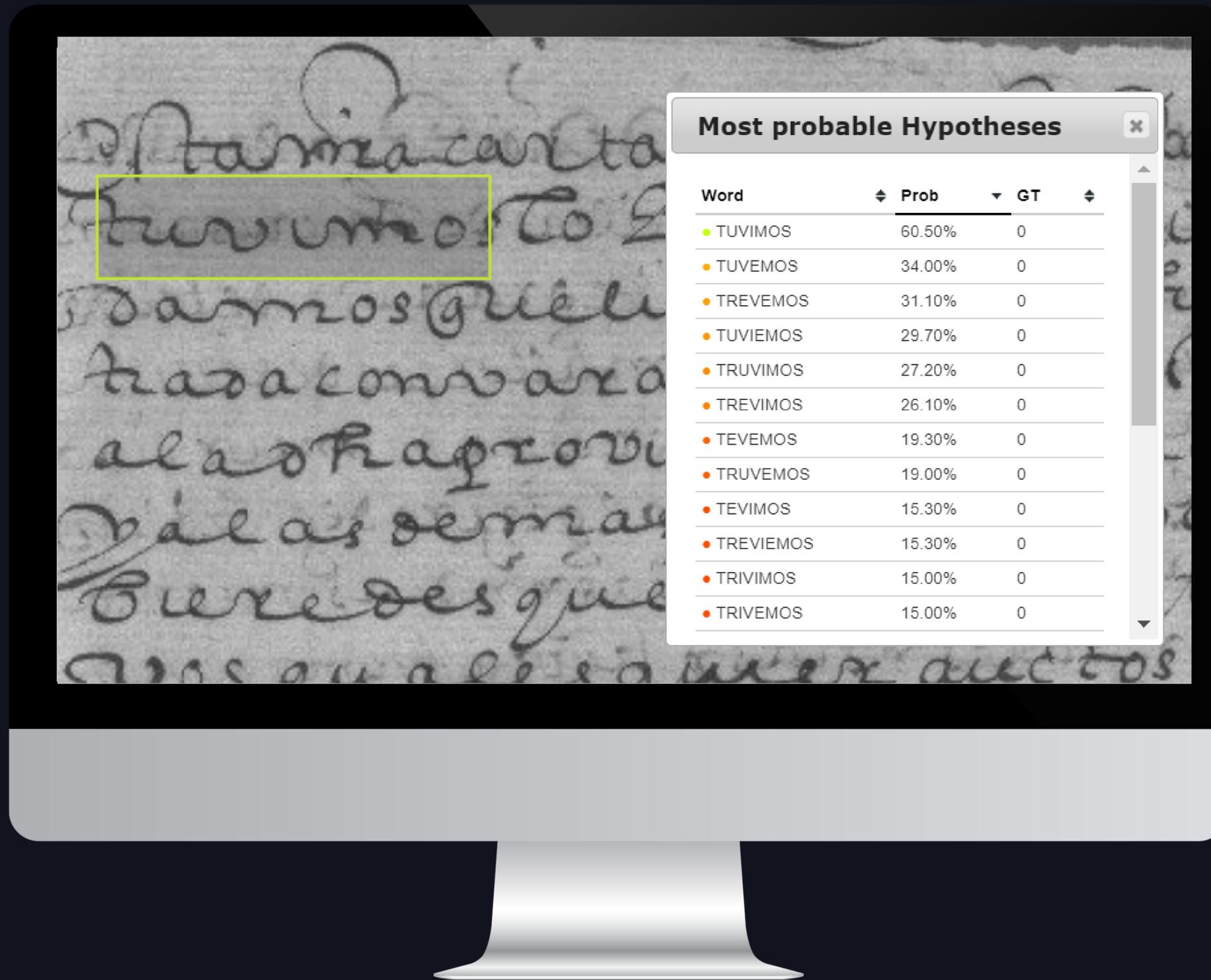


# Solutions

- Manual Process
- Automatic Transcription
- You may not need it
- Hybrid Approach



# Beyond Automatic Transcription

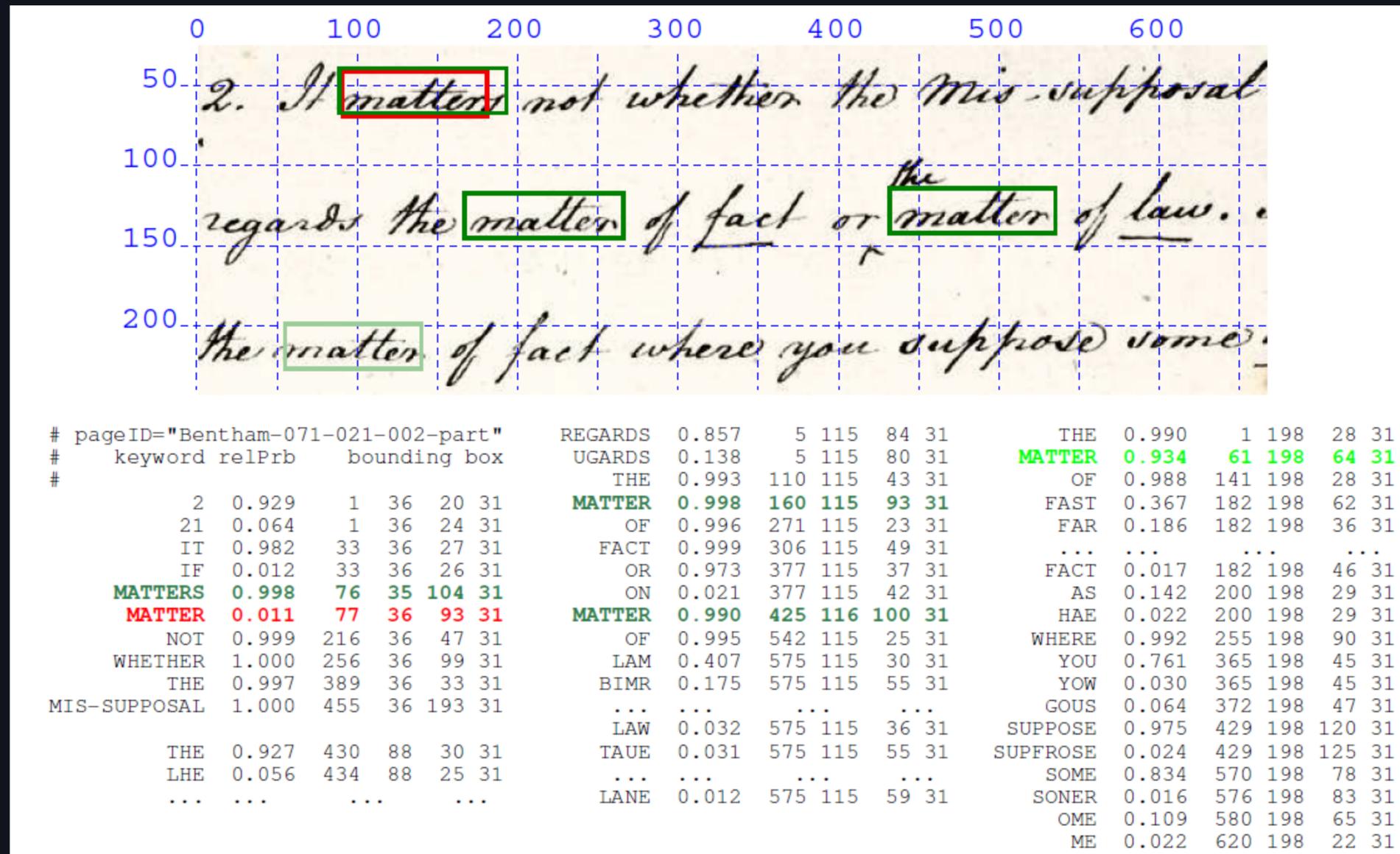


The image shows a computer monitor displaying a handwritten document. A yellow box highlights the word "Tuvimos" in the second line of the text. A window titled "Most probable Hypotheses" is overlaid on the right side of the screen, showing a list of words and their probabilities.

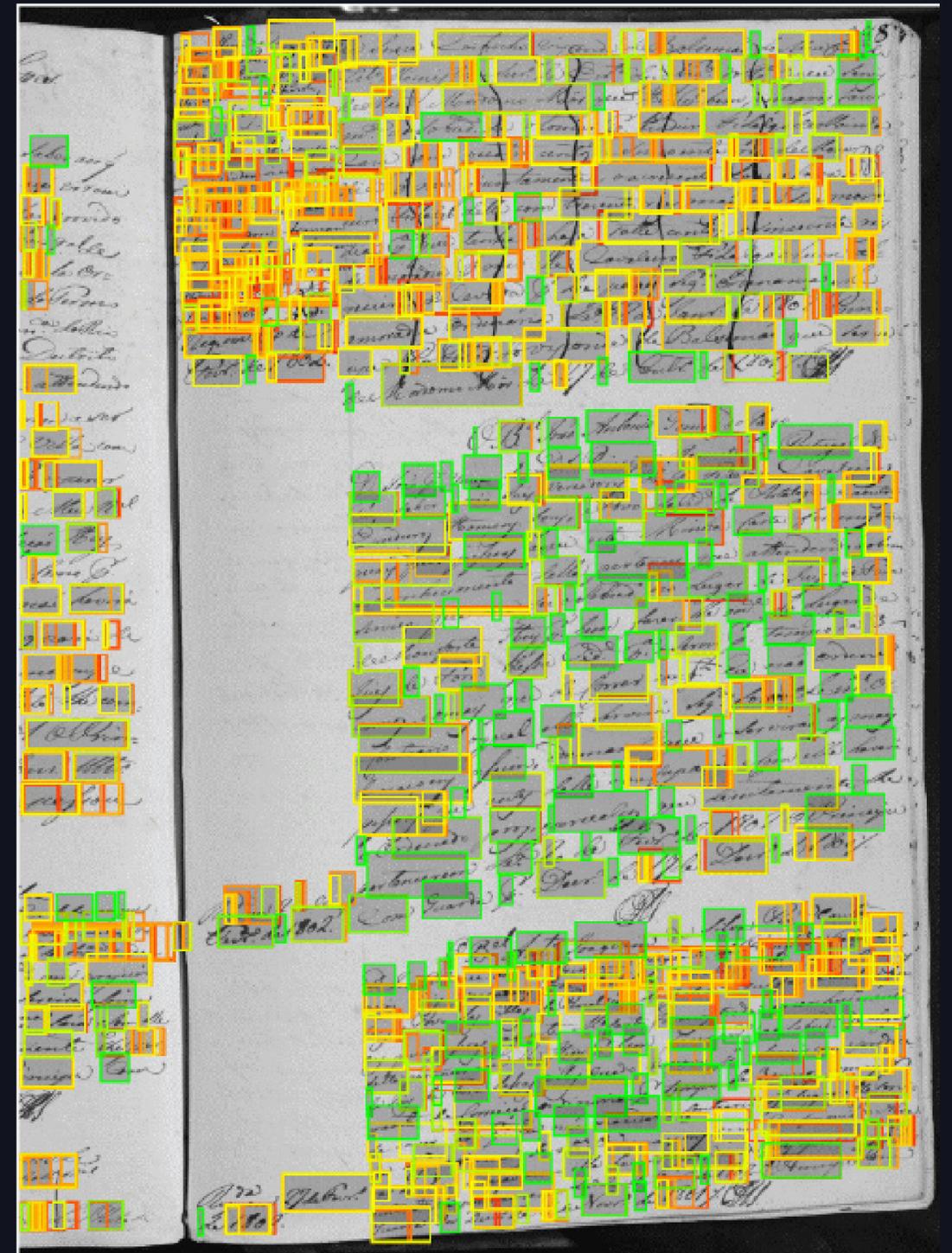
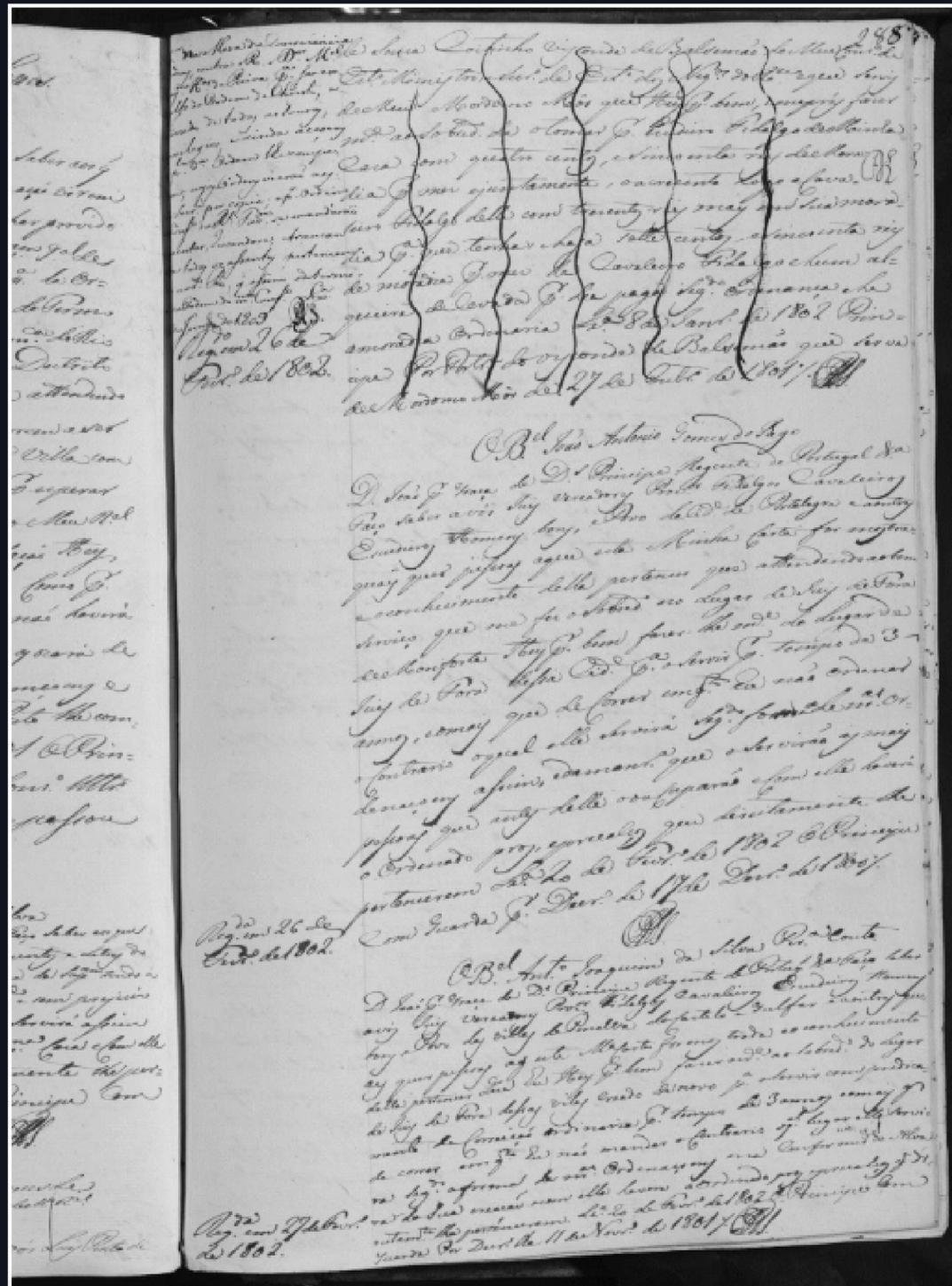
Word	Prob	GT
TUVIMOS	60.50%	0
TUVEMOS	34.00%	0
TREVEMOS	31.10%	0
TUVIEMOS	29.70%	0
TRUVIMOS	27.20%	0
TREVIMOS	26.10%	0
TEVEMOS	19.30%	0
TRUVEMOS	19.00%	0
TEVIMOS	15.30%	0
TREVIEMOS	15.30%	0
TRIVIMOS	15.00%	0
TRIVEMOS	15.00%	0



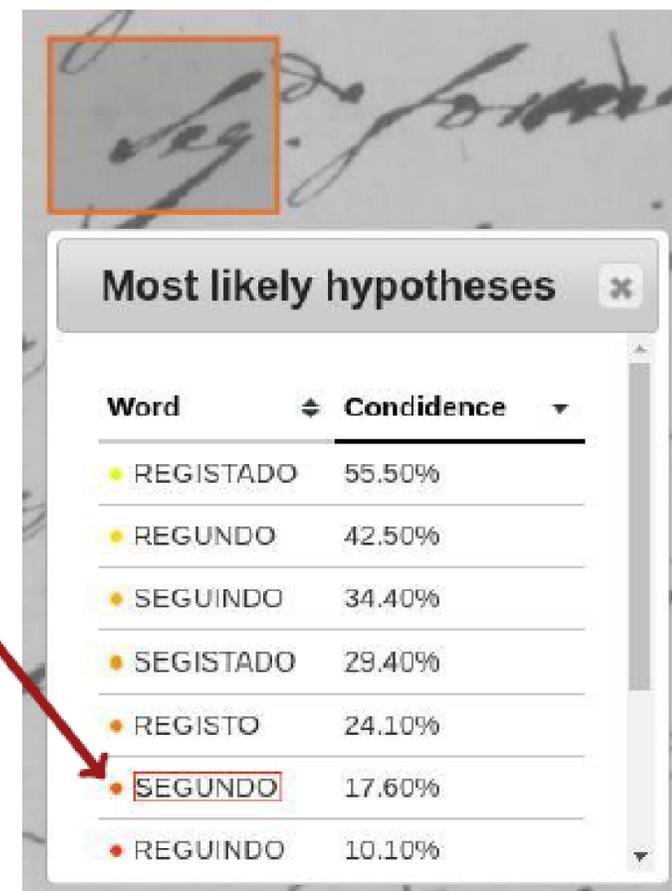
# Probabilistic Indexing



# Probabilistic Indexing



# Searching



# Transcription + Alternatives

del artículo sesenta y ocho de la ley vigente, se declara ~~se~~  
 admisible dicha excepción.  
 Se presentó Domingo Ferreiro no ocho de la primera  
 vez por el cupo de Santa Comba y reemplazó de mil ochocientos cincuenta  
 y uno, que tallado y resultó corto y en atención a lo que resulta del  
 expediente de prófugo se confirma la declaración hecha por el Ayuntamiento  
 y se condena al expresado Ferreiro en cincuenta días de malla y  
 por su insolencia en cien días de Carcel que había de sufrir en la  
 de la Capital del partido remitiéndose certificación del Alcalde Car-  
 celero que acredite haber cumplido la condena con expresión del día en

# Metadata Extraction

Nombre y apellidos *José Rafael*

Most probable Hypotheses

Rafael!nombre

Add

Word	Confidence	GT	
Rafael!nombre	53.60%	0	X

Most probable Hypotheses

José!nombre

Add

Word	Confidence	GT	
José!nombre	97.90%	0	X

Provincia de *Coruña*

Most probable Hypotheses

Coruña!provincia

Add

Word	Confidence	GT	
Coruña!provincia	100.00%	0	X

De *34* años Estado *casado* Profesión *empleado*

Most probable Hypotheses

empleado!oficio

Add

Word	Confidence	GT	
empleado!oficio	100.00%	0	X

Most probable Hypotheses

casado!estadocivil

Add

Word	Confidence	GT	
casado!estadocivil	100.00%	0	X

Most probable Hypotheses

34!edad

Add

Word	Confidence	GT	
34!edad	100.00%	0	X

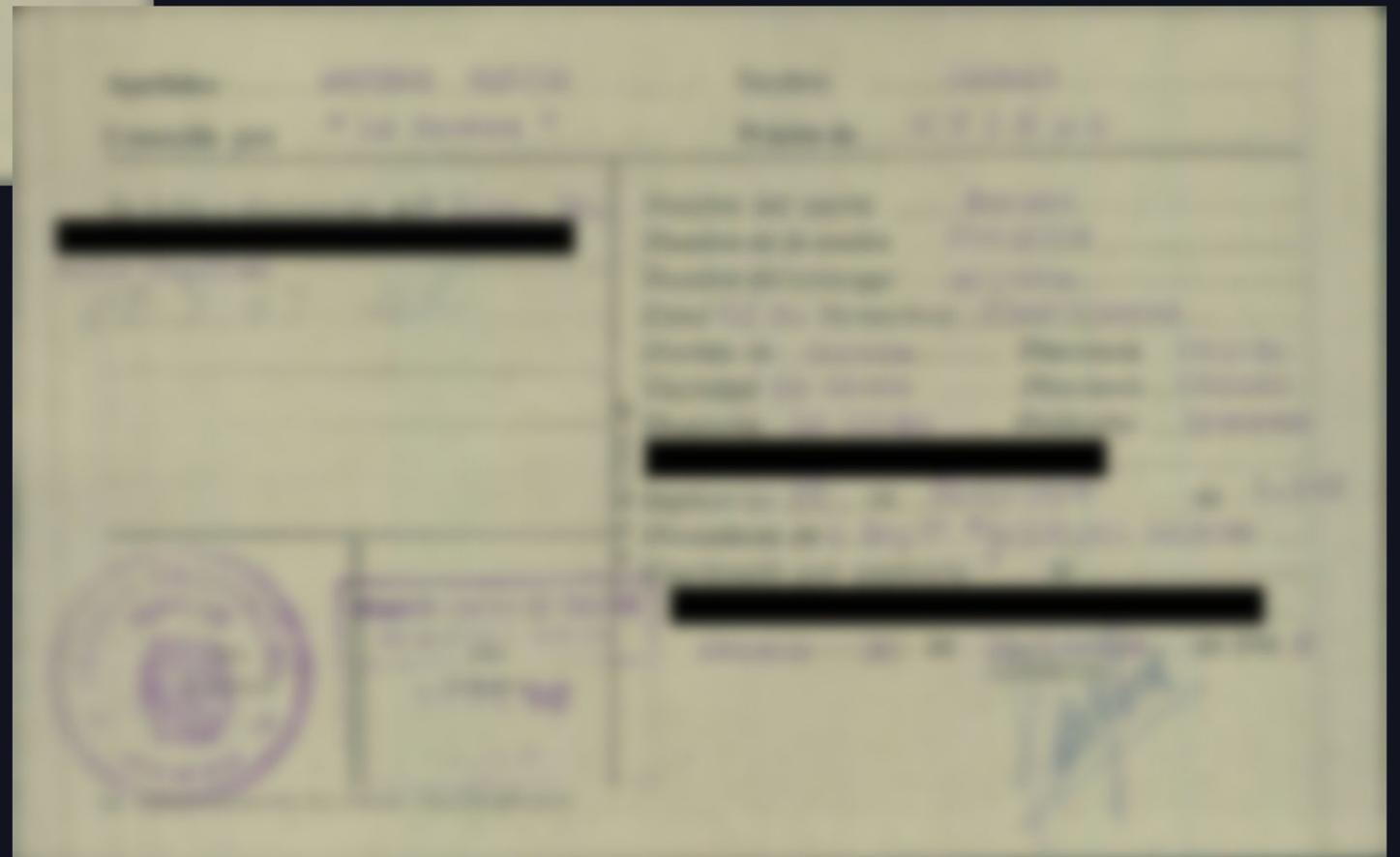
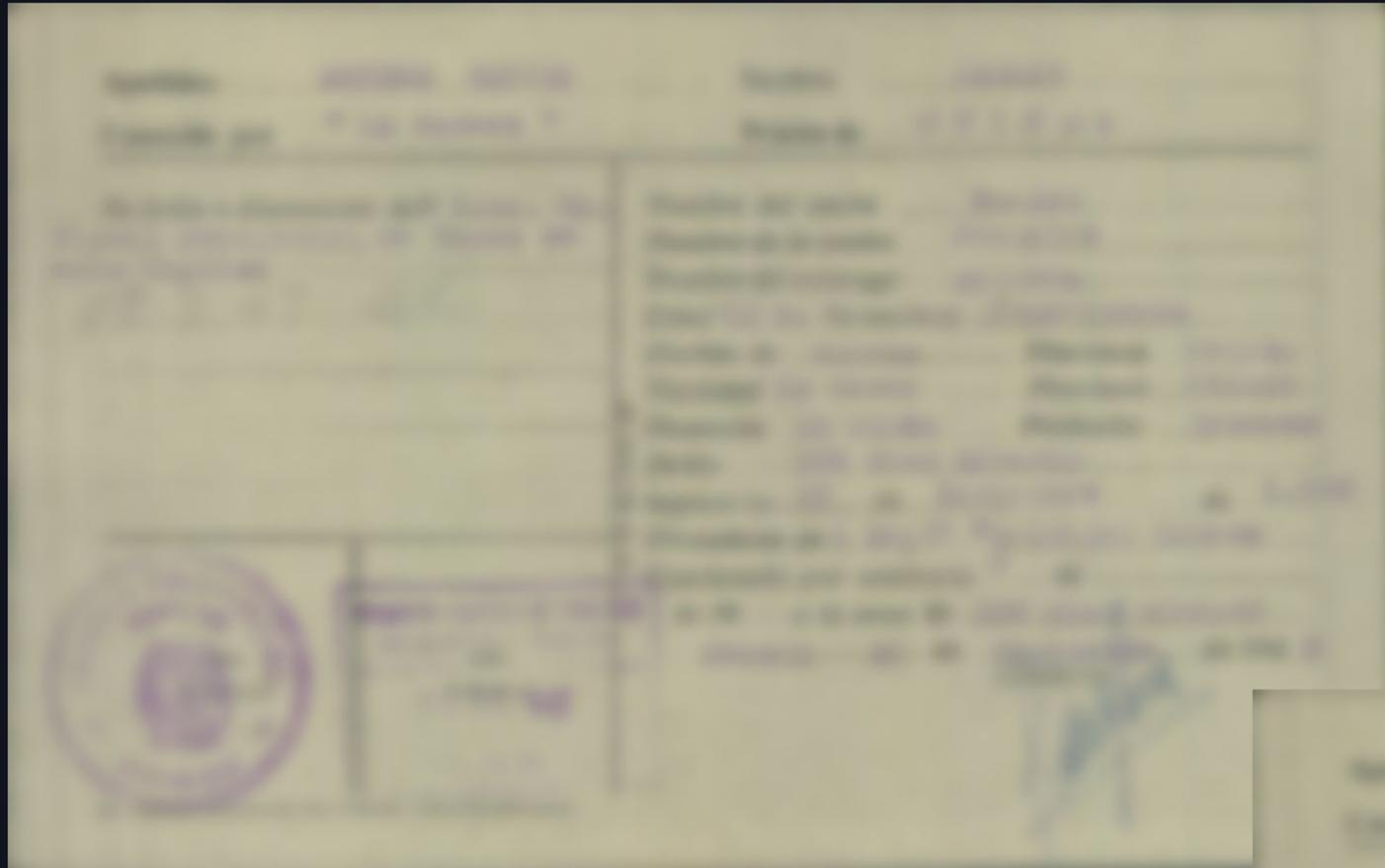
Most probable Hypotheses

Nombre!impr

Add

Word	Confidence	GT	
Nombre!impr	100.00%	0	X

# Anonymization







# Big Data Analysis

## Civil State

-----			0%
1096.8	soltero	36%	
1015.7	casado	34%	
-----			70%
437.0	soltera	14%	
209.0	casada	7%	
-----			91%
120.5	viuda	4%	
70.3	viudo	2%	
21.7	célibe célibero	1%	
-----			98%
62.2	[OTHER and ERRORS]	2%	
-----			100%

## Genre

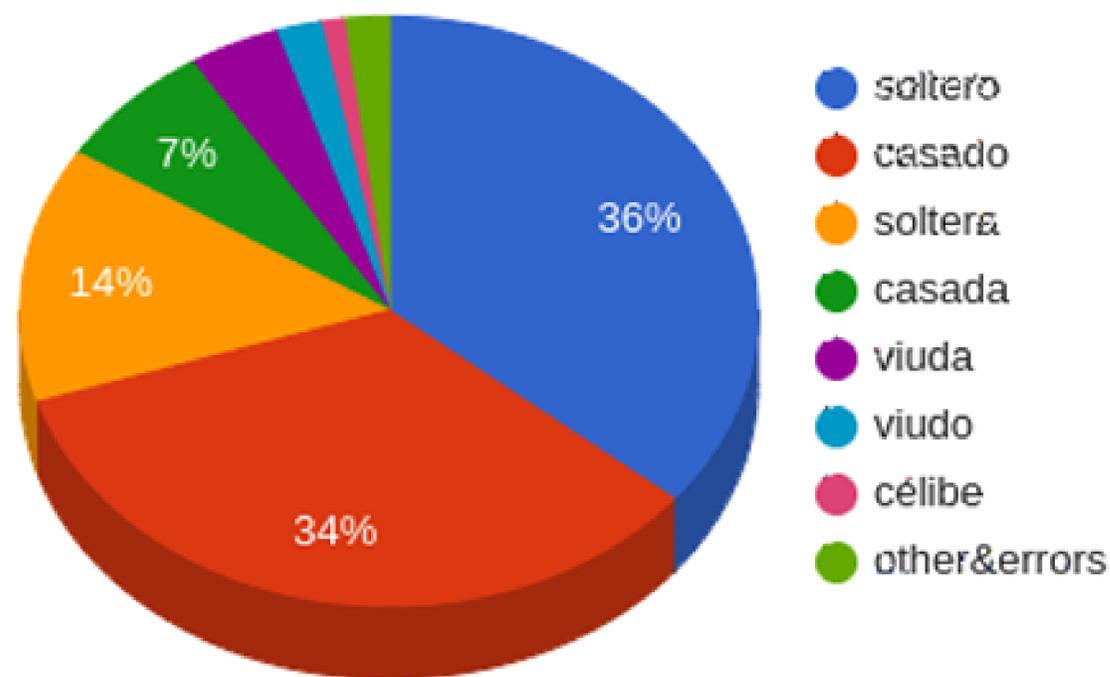
-----		
2182.8	Men:	72%
766.5	Women:	25%
73.9	Unknown:	3%
-----		

## Age (years old)

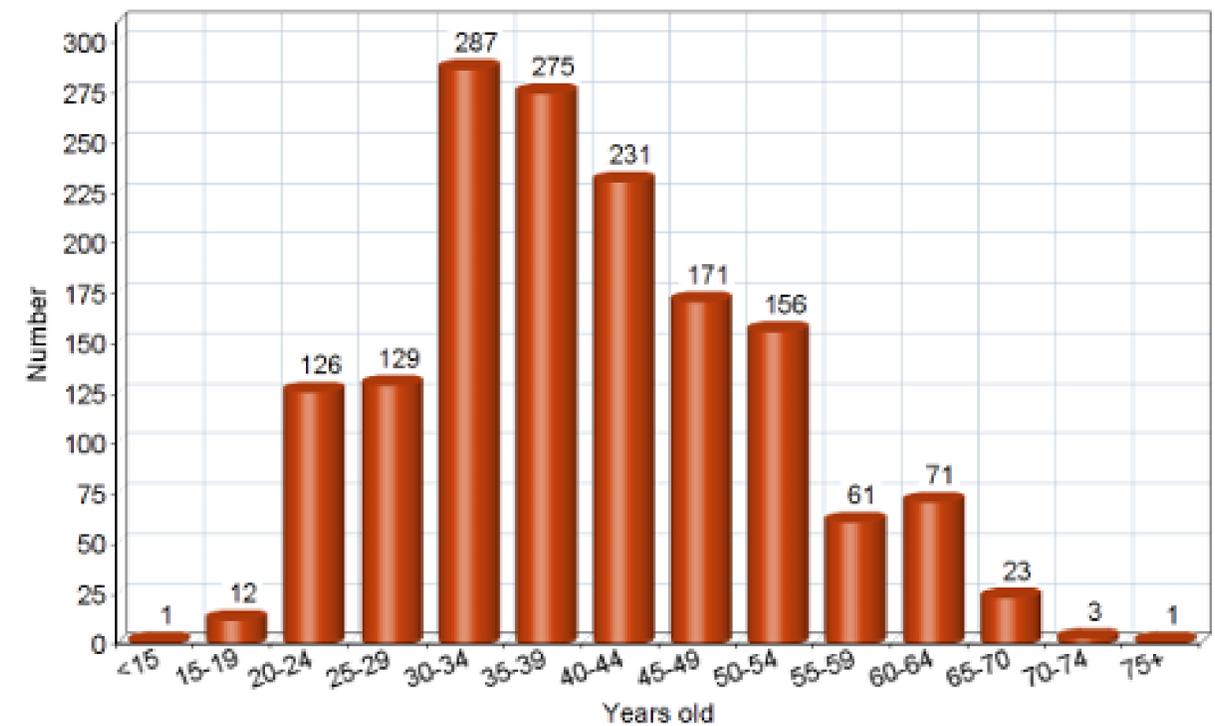
-----		0%
1.4	<15	
11.9	15-19	
-----		10%
126.3	20-24	
129.0	25-29	
286.6	30-34	
-----		50%
275.0	35-39	
230.5	40-44	
170.7	45-49	
156.2	50-54	
-----		90%
61.2	55-59	
71.4	60-64	
-----		99%
23.1	65-70	
3.2	70-74	
1.1	75+	
-----		100%

# Big Data Analysis

Civil state



Age



# Demonstration





November 2023

# Thank you!