



DEEPHEALTH

Hackaton - Course 2

De-identification of Radiological reports

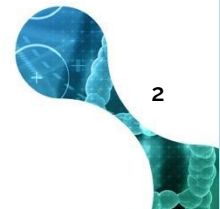


The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.



Index

1. Radiology Reports
 - a. Electronic Health Records (EHR)
 - b. Medical Texts
 - c. Radiology Reports
2. De-identification methodologies
 - a. Target Detection
 - b. Regular Expressions and List Matching
 - c. Named Entity Recognition
 - d. De-identification
3. Named Entity Recognition - NLP
 - a. NER - Machine Learning
 - b. Corpus construction
 - c. NER metrics
 - d. Applications
4. Dismed
 - a. The dataset
 - b. Corpus
 - c. Pipeline
 - d. Results





DEEPHEALTH

Radiology reports



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.



Electronic Health Records (EHR)

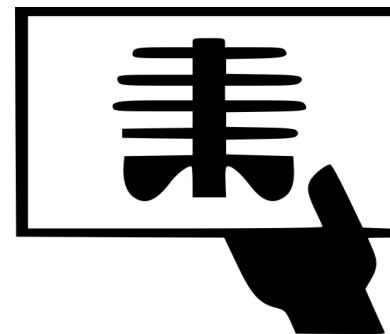
“Systematized collection of patient and population electronically stored health information in a digital format”
Structured and (mostly) Unstructured Data



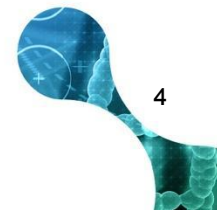
Tests results



Omics Data



Medical Imaging





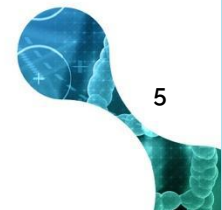
Medical Texts

Clinical data in biomedical research

- Clinical Trial Protocols
- Physician Progress Notes
- Surgery Reports
- Radiology Reports

A 12-year old girl with known hyperagglutinability, presented to the emergency department with a 2-week history of headaches and facial weakness. Neurologic examination indicated sensorineural hearing loss on the right side with Weber's test lateralizing to the left, and the Rinne's test demonstrating bone conduction greater than air conduction on the right. Magnetic resonance imaging of the head revealed severe structural defects of the right petrous temporal bone. No indication of cerebral infarction.

Figure 1, Starlinger, J. et al. (2016)



Radiology Reports

Clinical data in biomedical research

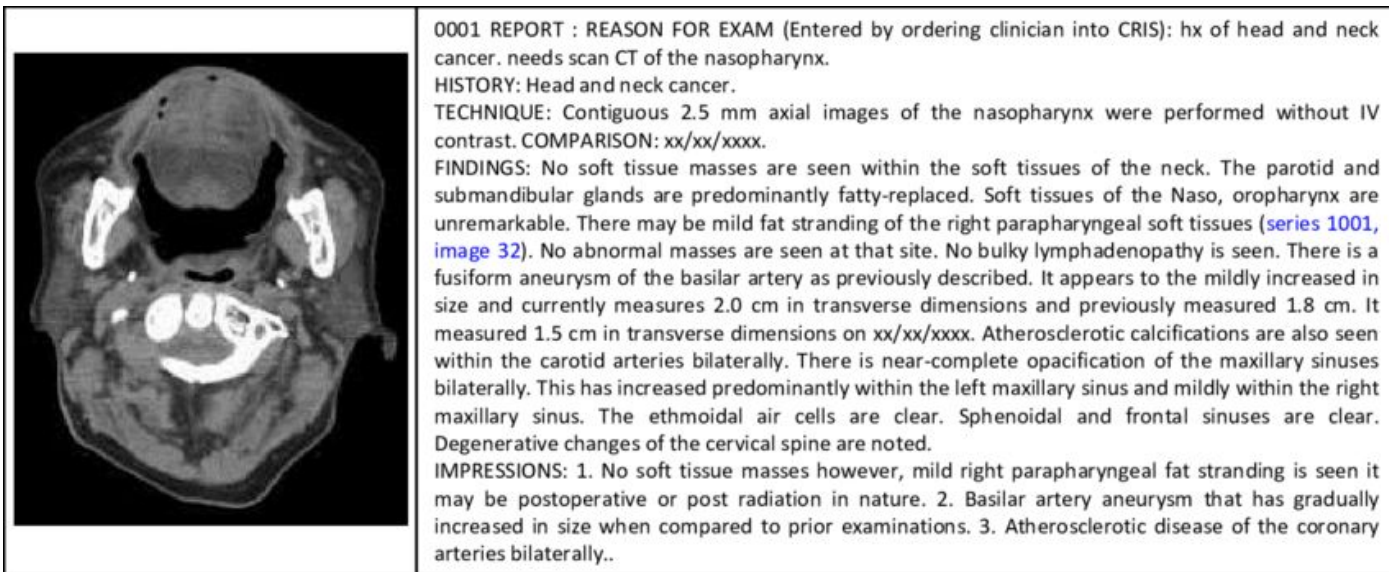


Figure 1, Shin, H. et al. (2016)



DEEPHEALTH

De-identification methodology



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.



Target Detection

Methods

HOSPITAL INST CLÍNICO INST DE INST VALENCIA INST
NÚM. D'HISTORIA 5942938 NUM CLÍNICA 541521 NUM
15/9/1959 FECHA
JOSE NAME MANUEL NAME AMADO NAME
JESUS NAME
ANA NAME BELEN NAME
Informe Radiología
488645 NAME NÚM.DE EPISODIO
CONCLUSIÓN
Malformación de Chiari tipo 1. Ocupación completa de celdillas mastoideas izquierdas.
NOMBRE / NOM: DANIEL NAME HERREIZ NAME GARCIA NAME
INFORME RADIOLOGÍA de / INFORME RADIOLOGIA de MARIANA NAME
SOBRON NAME PALOMARES NAME

Figure 3b, Pérez-Díez, I. et al. (2021)

Manual Detection

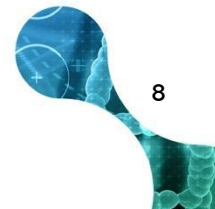
Automatic Detection

Regex, lists of terms,...

dd/mm/yyyy, dd.mm.yyyy
mm/dd/yyyy, mm.dd.yyyy

Machine Learning
Named Entity Recognition

Name, Date, Number,...





Target Detection

Regular expression and lists of terms matching

Matching dates

`^([1-9] |1[0-9]| 2[0-9]|3[0-1])(\.-)([1-9] |1[0-2])(\.-)20[0-9][0-9]$`

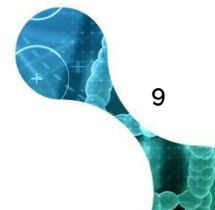


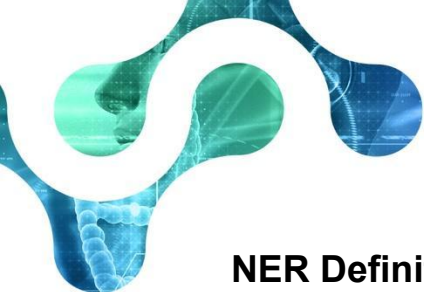
More Specific pattern

`^(?:(?:31(V|-|\.)(?:0?[13578]|1[02])|(?:Jan|Mar|May|Jul|Aug|Oct|Dec)))\1|(?:(?:29|30)(V|-|\.)(?:0?[1,3-9]|1[0-2])|(?:Jan|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec)))2)(?:1[6-9]||[2-9]\d)?\d{2}$|^(?:29(V|-|\.)(?:0?2|(?Feb))\3(?:1[6-9]||[2-9]\d)?(?:0[48]||[2468]||[048]||[13579]||[26])|(?:(?:16||[2468]||[048]||[3579]||[26])00)))$|^(?:0?[1-9]|1\d|2[0-8])(V|-|\.)(?:0?[1-9]|(?Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep))|(?:1[0-2]|(?Oct|Nov|Dec)))4(?:1[6-9]||[2-9]\d)?\d{2}$`

List of existing names

Nombre	Frecuencia
MARIA CARMEN	647,877
MARIA	589,055
CARMEN	375,835
ANA MARIA	271,616
JOSEFA	262,571
MARIA PILAR	260,302
ISABEL	258,873
LAURA	257,680
MARIA DOLORES	255,040
MARIA TERESA	248,131
ANA	246,892
CRISTINA	228,022
MARTA	226,974
MARIA ANGELES	223,560
LUCIA	208,878
FRANCISCA	204,718
MARIA ISABEL	203,839
MARIA JOSE	203,236





Target Detection

Named Entity Recognition (NER)

NER Definition - *The task of identifying and categorizing key information (entities) in text. NER is a form of Natural Language Processing (NLP).*

Albert Einstein
Person

is a physicist who was born in

German
Place

Heat
Action

a

frying pan
Cooking Tool

,

add
Action

pork
Food

miso paste
Food

, and

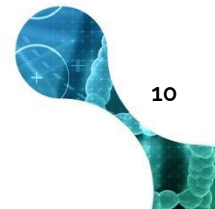
fry
Action

them

Examples of NERs

- Names
- Organizations
- Places
- Phone Numbers
- Genes
- Diseases
- Compounds

<https://medium.com/optuna/nerman-named-entity-recognition-system-built-on-allennlp-and-optuna-c044c319b955>





De-identification

MRI Brain
Date of service: 24/11/2020
Patient: Alberto Pérez
Findings: The ventricles, cisterns...

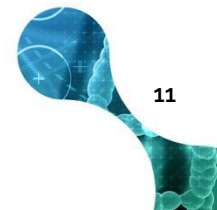
Substitution



Elimination

MRI Brain
Date of service: 13/04/2002
Patient: Juan Martínez
Findings: The ventricles, cisterns...

MRI Brain
Date of service: XX/XX/XXXX
Patient: XXXXX
Findings: The ventricles, cisterns...





DEEPHEALTH

Named Entity Recognition



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.



NER - Machine Learning

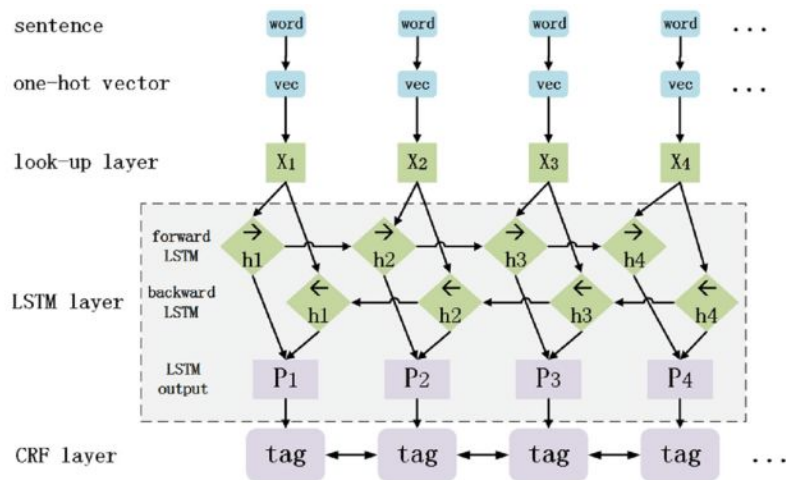
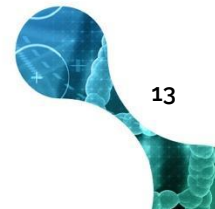


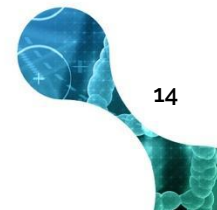
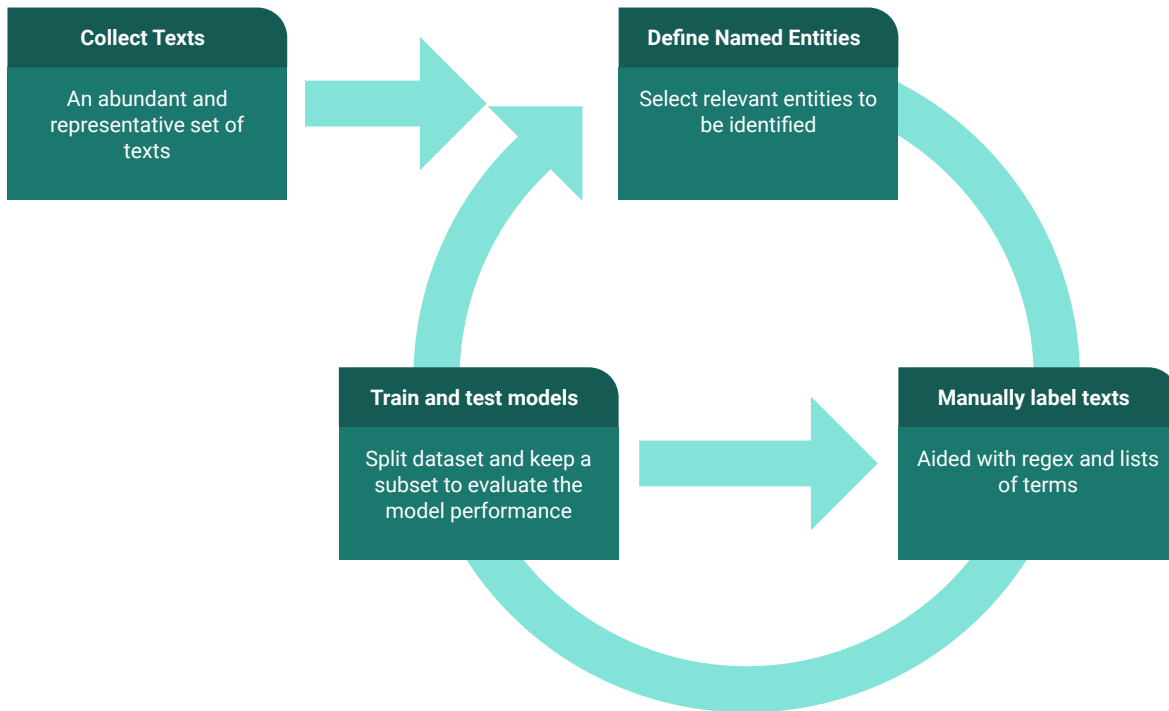
Figure 3, Ji, B. et al. (2019)

- Automatic processing of texts
- Fast and Accurate
- Adaptable to custom entities
- Requires a ML model trained with a suitable annotated dataset (corpus)
- Generation of corpus is expensive in terms of time and work





Corpus construction





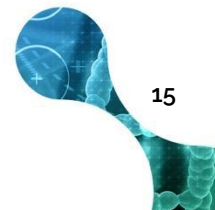
NER metrics

	Predicted Positives	Predicted Negatives
Positives	True Positives (TP)	False Negatives (FN)
Negatives	False Positives (FP)	True Negatives (TN)

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

$$F1score = \frac{2 \cdot precision \cdot recall}{precision + recall}$$





Applications of NER

Texts Anonymization

MRI Brain
Date of service: XX/XX/XXXX
Patient: XXXXX
Findings: The ventricles, cisterns...

Structure Data

MRI Brain
Date of service: 24/11/2020
Patient: Alberto Pérez
Findings: The ventricles,
cisterns...



Patient	Alberto Pérez
Date	24/11/2020
Disease	...
Treatment	...

Information Extraction

BioNER

BRCA1 gene causes predisposition
to breast cancer and ovarian cancer

gene
disease

Inferring
Relations

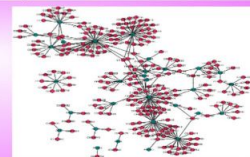
BRCA1 gene causes predisposition
to breast cancer and ovarian cancer.

association verb (entity-verb-entity)

Analyzing Polarity
& Strength of
Relations

Word Distance = 4 OR
Shortest Path Distance = 5
Negative Polarity

Visualization &
Integrating
Relations





DEEPHEALTH

Our de-identification method

Dismed



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.



The dataset

A

DIAGNÒSTIC / DIAGNÓSTICO

TORRES NAME

NOM FACULTATIU / NOMBRE FACULTATIVO

RAFAEL NAME TORRES NAME NAVAJAS NAME

METGE PETICIONARI / MÉDICO PETICIONARIO

JOSE NAME ALONSO NAME GIL NAME

R.M. Cerebral, 26/07/2017 FECHA

R.M. Columna cervical, 08/04/2011 FECHA

VALORACIÓ CAB / VALORACIÓN CAB

protocolo de volumetría cerebral y cervical, gadovist 10 ml iv.
ggc-38139

Tras valoración comparativa con previos, ausencia de incremento de carga lesional o de actividad inflamatoria, no existiendo focos con realce tras el contraste.

Volumen C2-C6: 6.788cm3

área C2. 94.67m2

Hospital INST Clínico INST Universitario INST de INST Granada INST

SIP 294321 NUM NÚM.D'HISTÒRIA CLÍNICA

NÚM. DE HISTORIA CLÍNICA 594637 NUM

DATA NAIXEMENT

FECHA NACIMIENTO 13 de enero de 1972 FECHA

DIRECCIÓ

DIRECCION CALLE DIR PAJARO DIR VERDE DIR - 21 18299 DIR

VALENCIA LOC

B

Datos del paciente

Nombre: Rocío NAME

Apellidos: Pérez NAME Ontiveros NAME

NHC: 22 75689632 36 NUM

Domicilio: Av DIR de DIR Leon DIR 66 DIR 1H DIR

Localidad/Provincia: Lleida LOC

CP: 06233 NUM

Fecha de nacimiento: 05/04/1937 FECHA

País de nacimiento: España LOC

Edad: 10 años Sexo: Varón

Fecha de ingreso: 15/08/2016 FECHA

Servicio: Oftalmología

Médico: Ender NAME Goñi NAME Moreno NAME NºCol: 15 15 31525 NUM

Consulta por dolor abdominal, observándose en la ecografía una masa renal. Se realizó biopsia tru-cut diagnosticada de tumor mesenquimal benigno. Fue intervenido quirúrgicamente.

Hallazgos CAB histológicos CAB

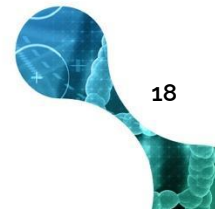
Se trataba de una proliferación de células fusocelulares con zonas de diferentes densidades sin atípias ni mitosis. Las células del estroma eran positivas para CD-34 y vimentina.

Ultraestructuralmente el tumor presenta células mesenquimales inmaduras

Dirección para correspondencia: Irene NAME Amat NAME Villegas NAME

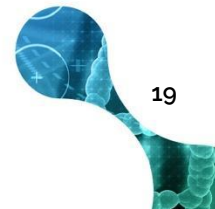
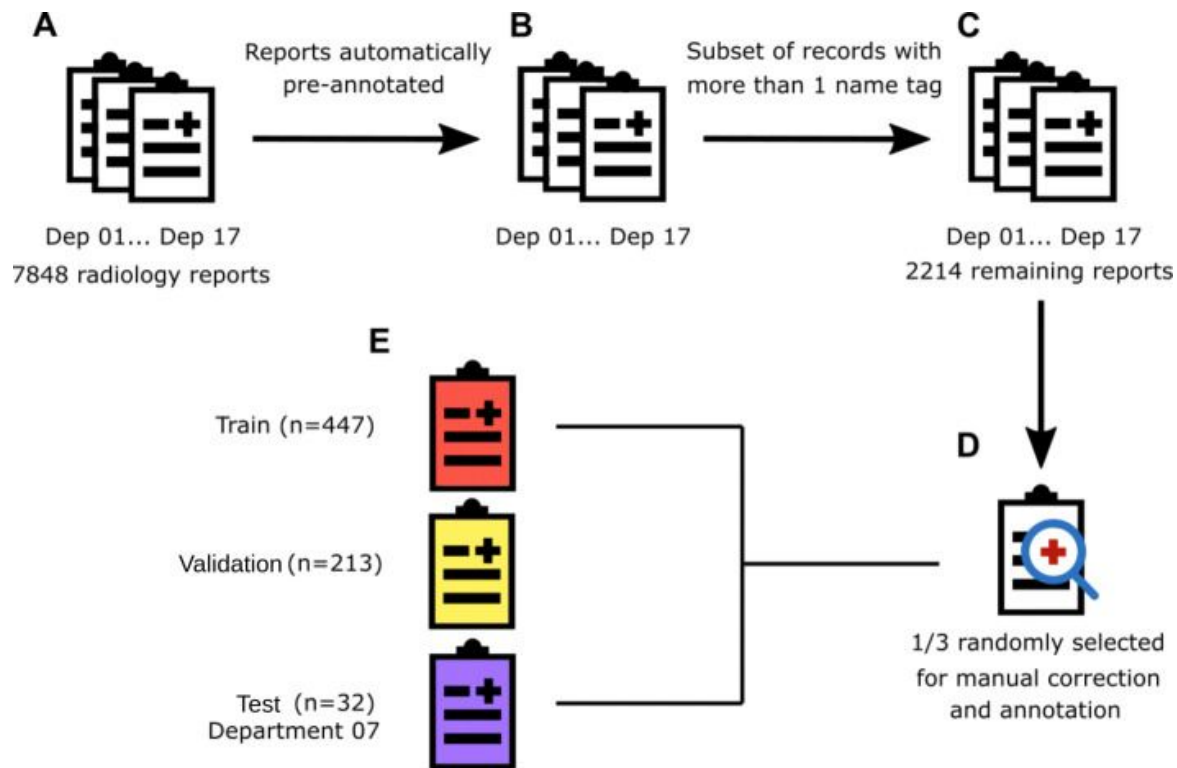
Pedro DIR de DIR Alejandría DIR Nº 1 DIR , 31014 DIR Pamplona LOC

Navarra LOC

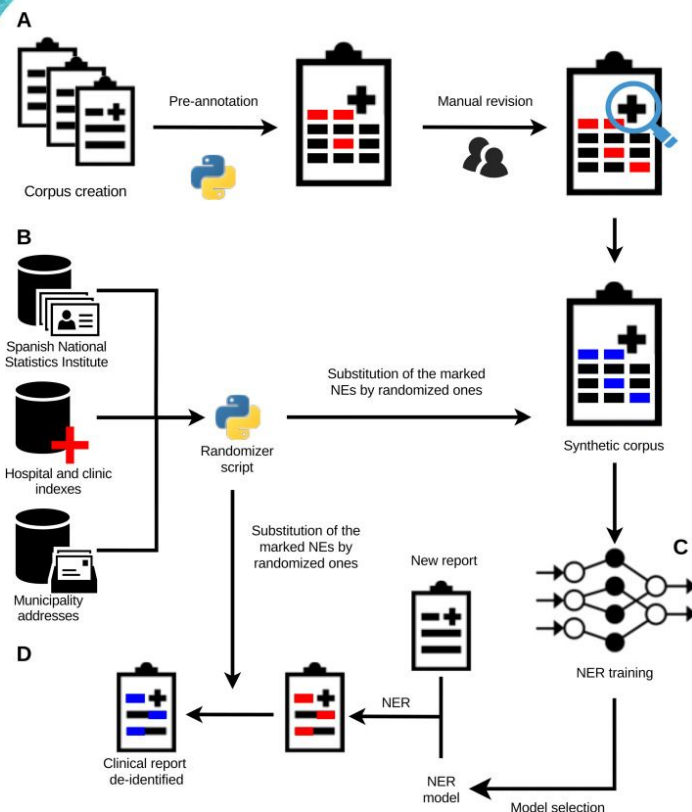




Corpus Construction



Dismed pipeline

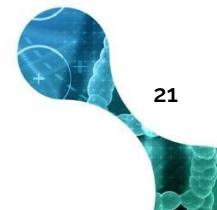


NEs	Description
CAB	Section headers
NAME	Names and surnames (patient and others)
DIR	Full addresses, including streets, numbers and zip codes
LOC	Cities, inside and outside addresses
NUM	Numbers or alphanumeric strings that might identify someone, including digital signatures, patient numbers, medical numbers, medical license numbers and others
FECHA	Dates
INST	Hospitals, healthcare centres or other institutions that might point to someone's location



Results

	Training set			Validation set			Test set			MEDDOCAN		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
CAB	99.43	96.54	97.96	98.28	93.98	96.08	92.54	74.49	82.52	4.76	33.33	8.33
DIR	100	100	100	94.28	63.96	76.01	87.79	74.77	61.46	43.15	4.47	8.01
FECHA	100	100	100	98.54	99.04	98.78	98.20	97.53	97.86	51.39	89.41	65.13
INST	99.97	99.96	99.98	98.19	97.24	97.71	93.50	98.00	95.69	45.72	12.28	19.27
LOC	100	100	100	76.64	54.66	63.80	61.04	26.85	36.79	7.19	0.32	0.59
NAME	100	99.99	99.99	98.34	98.28	98.31	88.78	94.29	93.19	75.62	83.91	79.23
NUM	100	100	100	97.81	95.65	96.72	95.11	87.56	91.18	68.50	60.32	63.99
	99.87	99.28	99.58	98.06	96.10	97.08	93.23	89.39	91.31	65.63	55.37	59.98





Publication

Pé et al. *Journal of Biomedical Semantics* (2021) 12:6
<https://doi.org/10.1186/s13326-021-00236-2>

Journal of
Biomedical Semantics

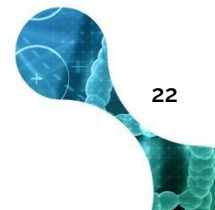
RESEARCH

Open Access

De-identifying Spanish medical texts - named entity recognition applied to radiology reports



Irene Pérez-Díez^{1,2†}, Raúl Pérez-Moraga^{1,3†}, Adolfo López-Cerdán^{1,2}, Jose-Maria Salinas-Serrano⁴ and
María de la Iglesia-Vayá^{1,5,6*} 





DEEPHEALTH

**Thank you for your
attention!**



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.