

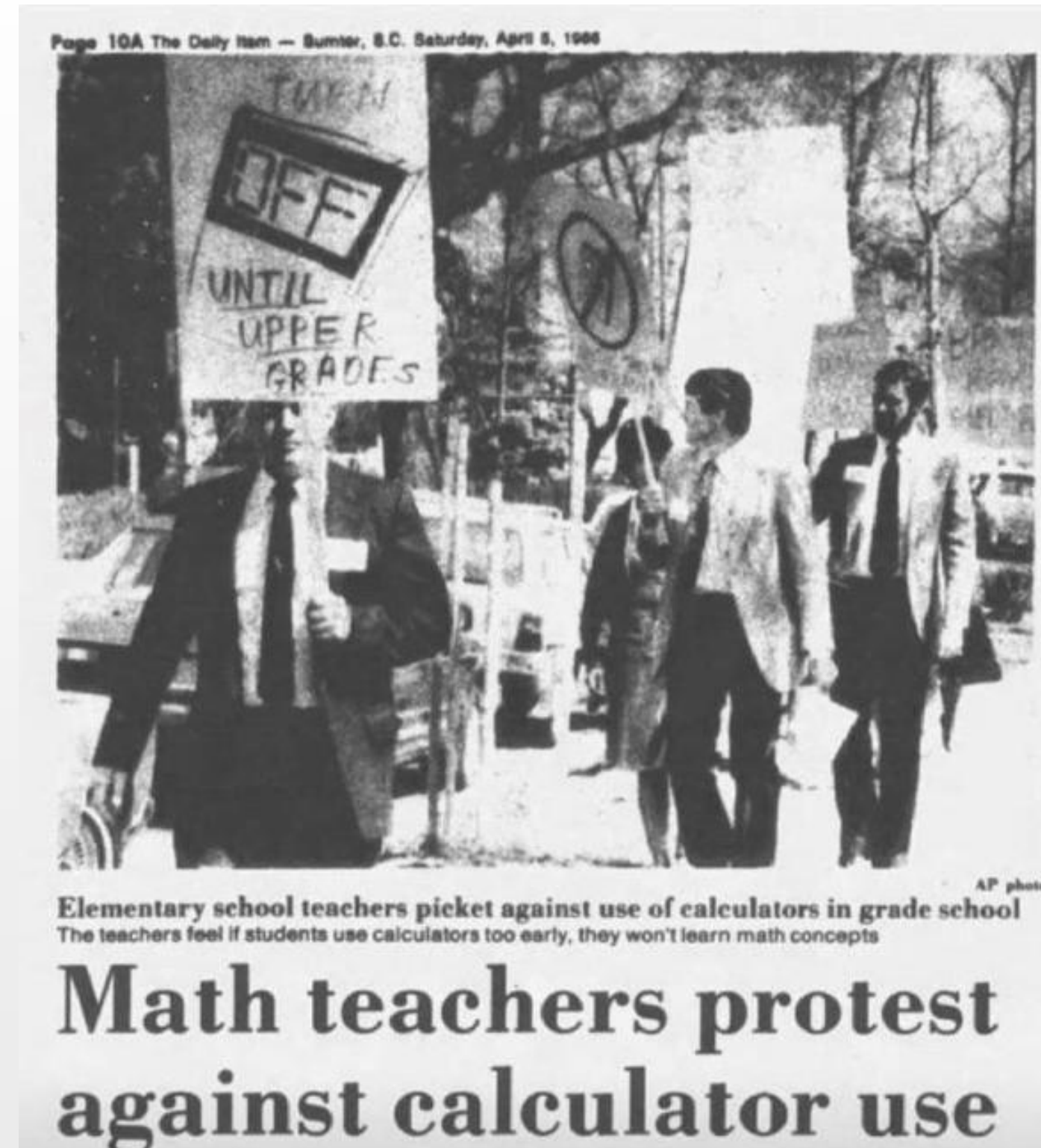


Integración de la inteligencia artificial en Radioterapia: avances tecnológicos y retos clínicos

Eva Ambroa Rey

*Especialista en Radiofísica Hospitalaria
Hospital del Mar (Barcelona)
evamaria.ambroa.rey@hmar.cat*

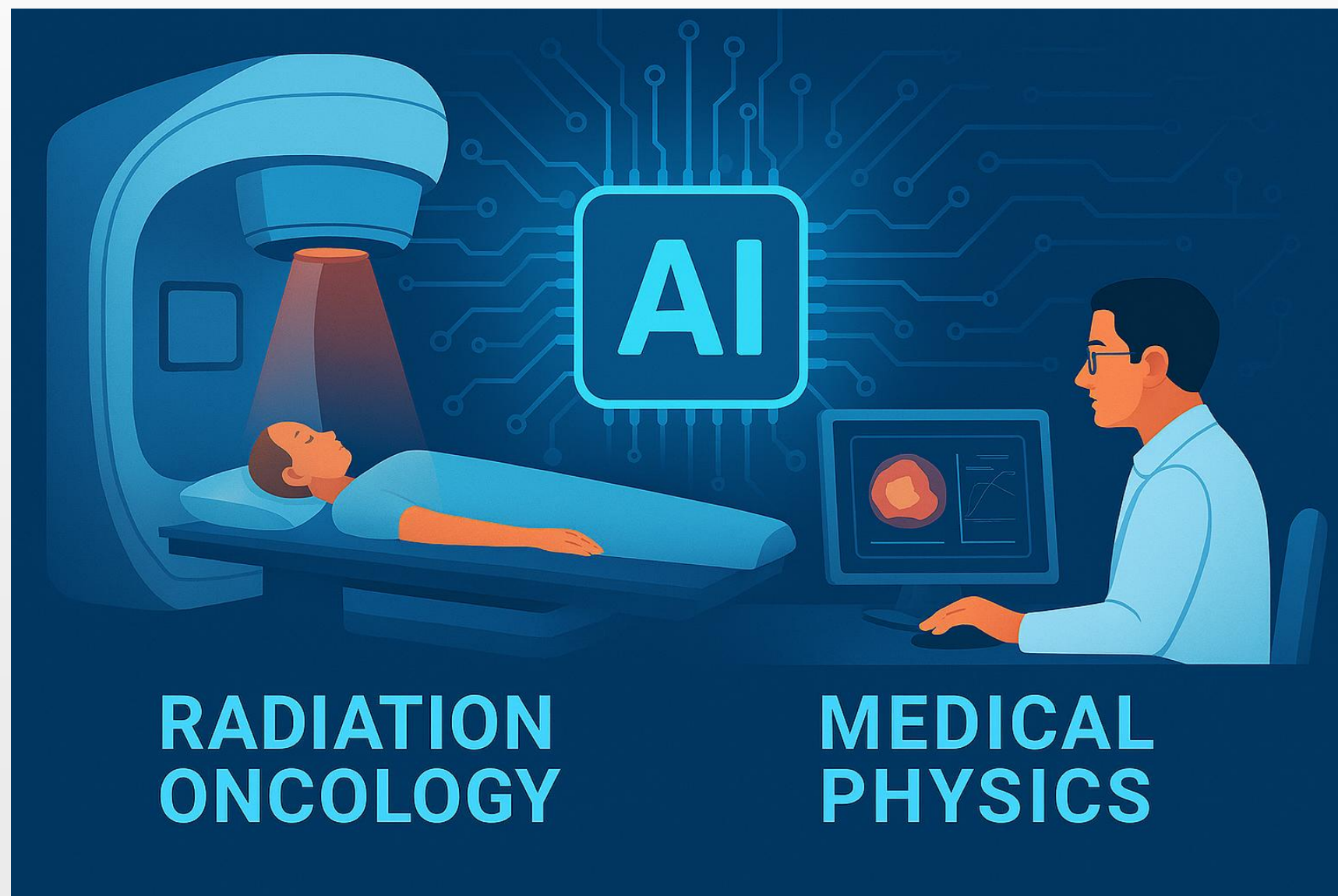
¿LA IA REEMPLAZARÁ NUESTRO TRABAJO?



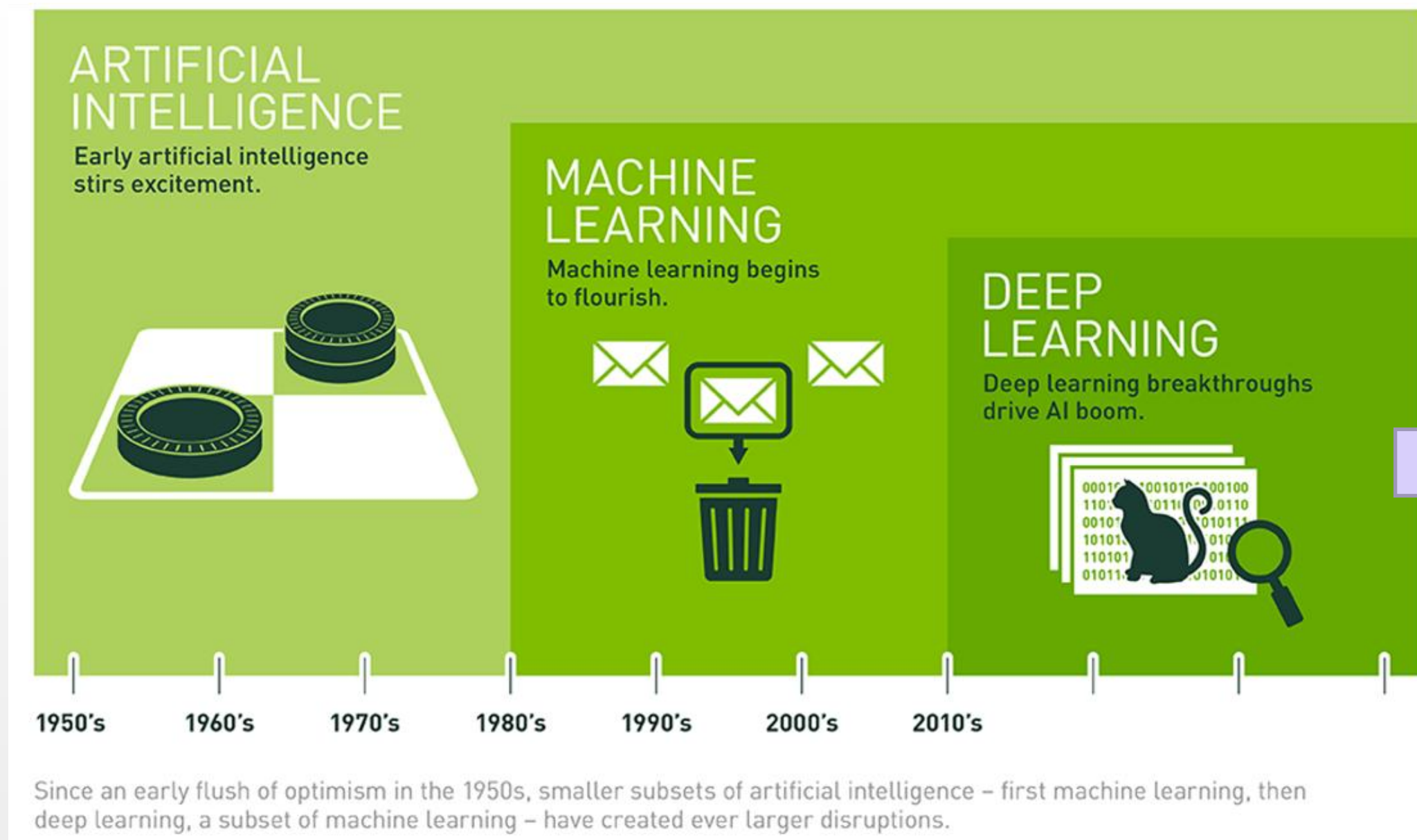
1966

FUTURO DE LA IA EN RT Y FÍSICA MÉDICA

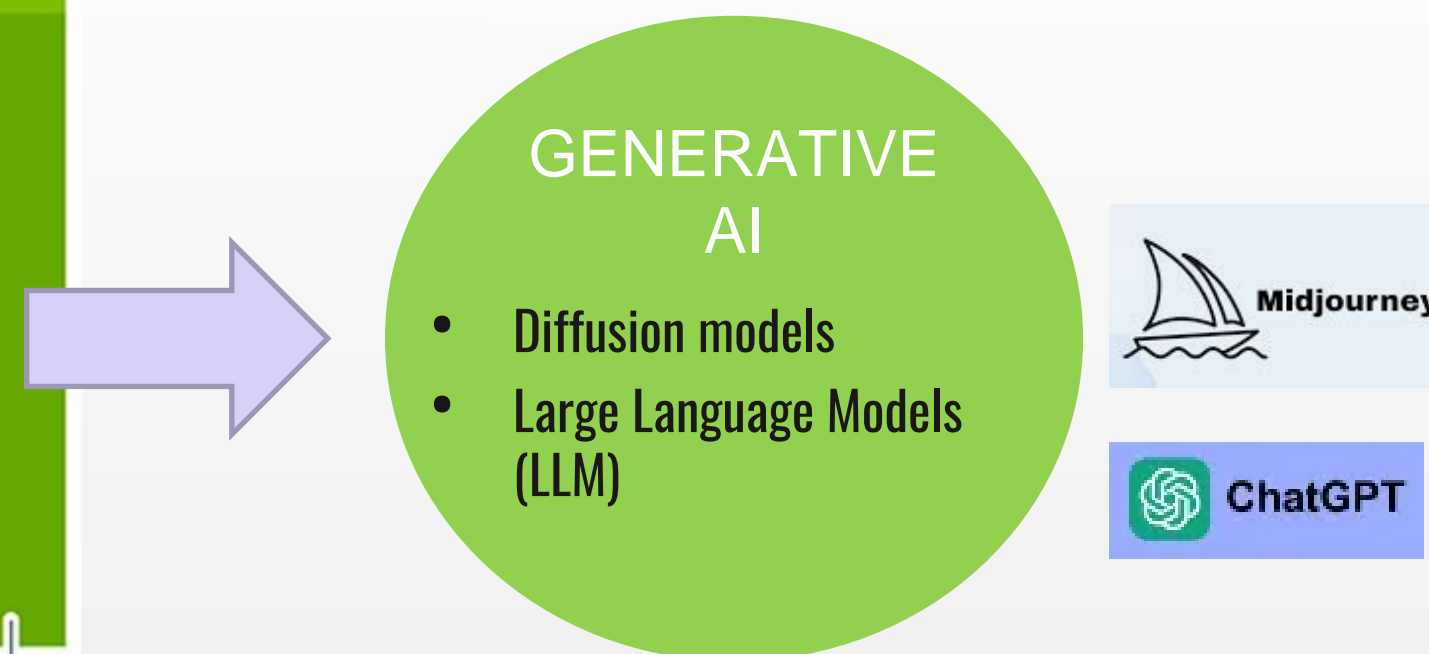
- Prompt Co-pilot: puedes crearme una imagen de cómo influirá la IA en el ámbito de la oncología radioterápica y la física médica?
- Prompt Co-pilot: puedes crearme una imagen de cuál será la promesa de la IA en sanidad?



Evolución histórica



Fuente: Michael Copeland What's The Difference Between AI, Machine Learning And Deep Learning, NVIDIA Blogs



DATOS Y ESTRUCTURA

- 97% de los datos clínicos son 'dark data' perdidos
- Cuando se estructuran los datos y se comparten:

Benedict S, et al. *Introduction to Big Data in Radiation Oncology*. (2016)

Big data initiatives

OncoSpace – AI-based predictive planning

Momentum – MRgRT

UniCancer – CANTO-RT

AAPM Data Science Committee – Big Data Subcommittee

TCIA - RTStruct data

MROCQ – registry for benchmarking and practice improvement

NSIR-RT – nation-wide incident reporting

E²-RADIatE – EORTC/ESTRO pan-European platform for prospective registries

Grand Challenge datasets

IAEA DIRAC – comprehensive database on radiotherapy resources



APLICACIONES IA EN RT

DECISIÓN TTO

01.

Herramientas que combinan datos clínicos con genómicos e imagen para dar soporte a la decisión.



IMAGEN

02.

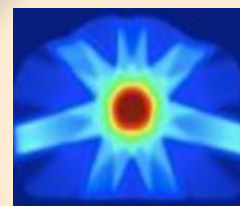
Reconstrucción, mejora de imagen, generación de imágenes sintéticas, auto-contorneo...



PLANIFICACIÓN DE TTO

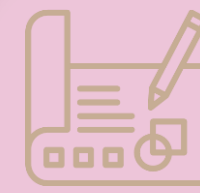
03.

Predicción dosis, DVH, técnica, ...



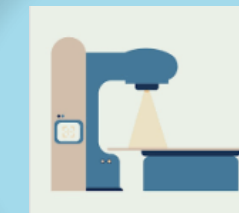
05. CONTROL DE CALIDAD

QA paciente, errores MLC, output, simetría, ...

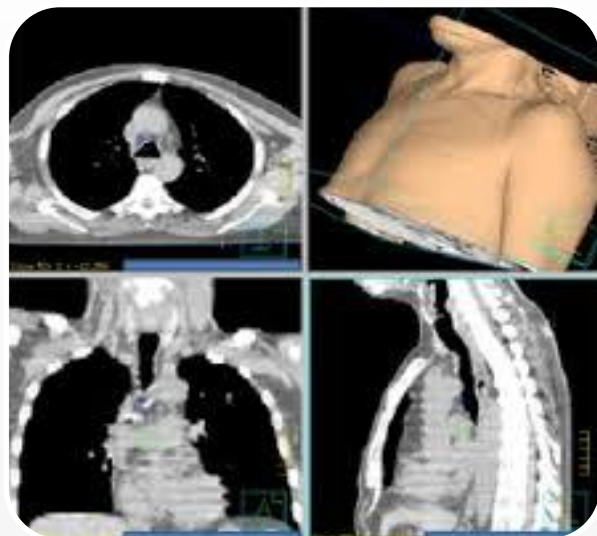


04. ADMINISTRACIÓN TTO

Posicionamiento, monitorización durante el tto. predecir cambios anatómicos, adaptativa sobre CBCT, ...



IMAGEN



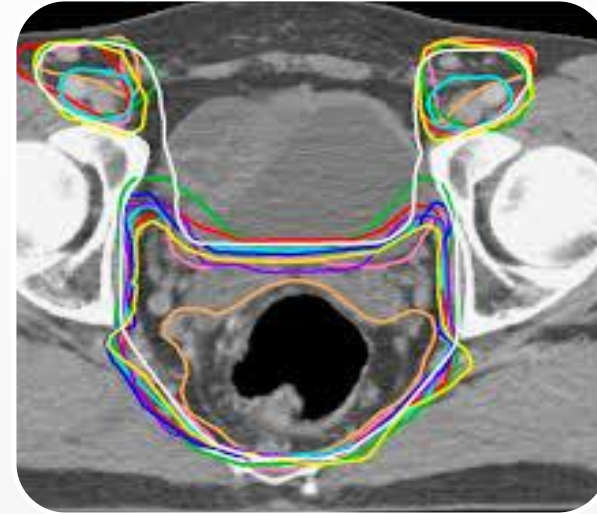
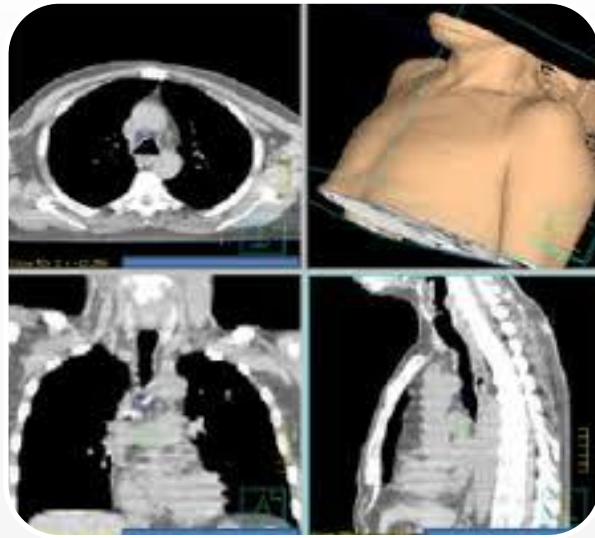
Generación/Reconstrucción CT/RM

- ↑ calidad imagen
- ↓ ruido
- ↓ artefactos
- Supervisados
- U-Net, GAN

Table 1. Available commercial solution for simulation CT generation in October 2023.

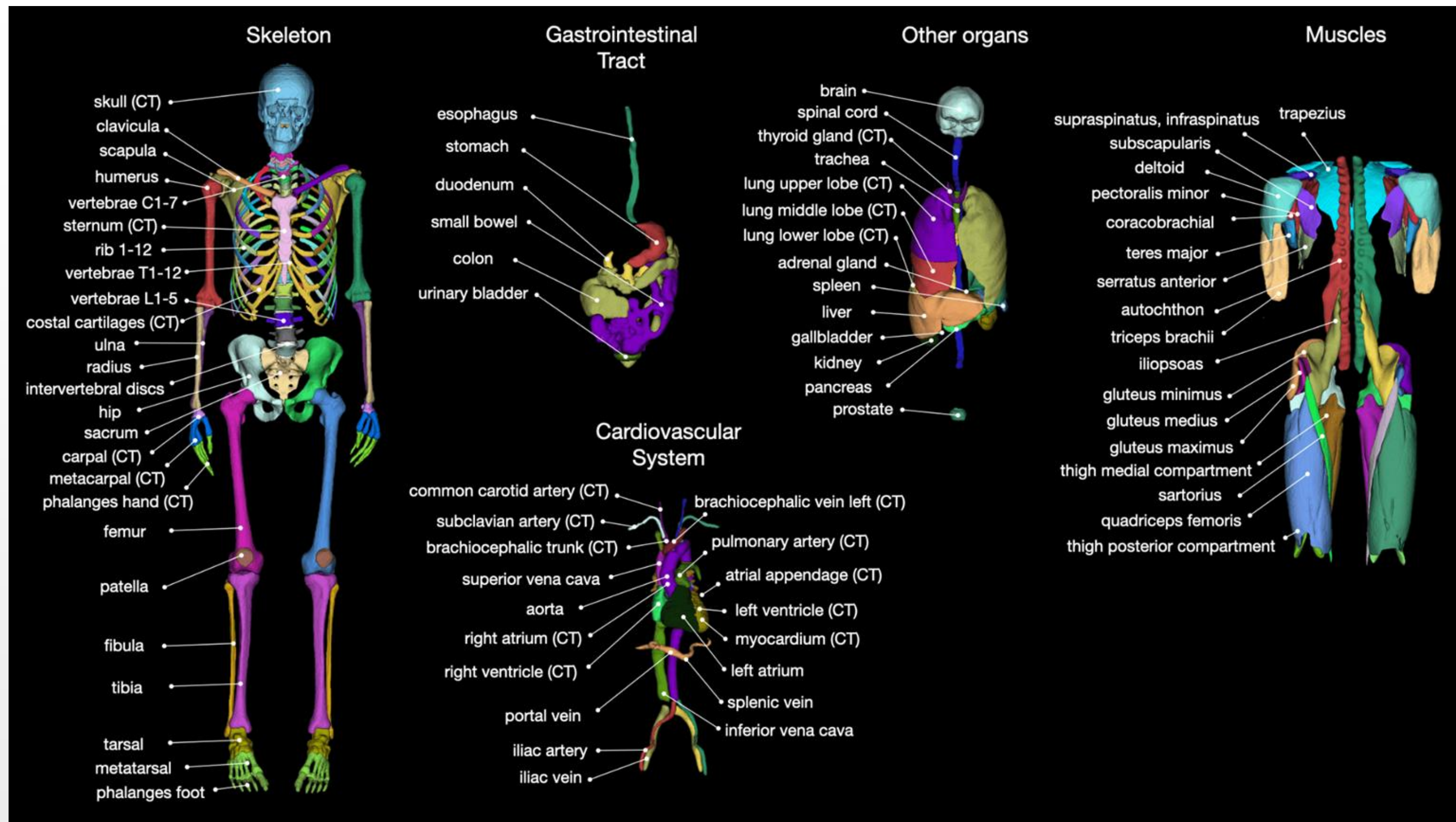
Brand name (Company)	Imaging technique	Anatomical site	Certifications
MRCAT® (Philips Healthcare)	T1-weighted GRE + Dixon	Pelvis, brain, head & neck	CE, FDA March 2016
MRI Planner® (Spectronics Medical)	T2-weighted	Pelvis	FDA, CE June 2016
	T1-weighted GRE	Brain, head & neck	
SyngoVia.syntheticCT® (Siemens Healthcare)	Multiple MRI sequences	Brain	FDA, CE January 2018
	T1-weighted GRE + Dixon	Pelvis	
MR-Box® (Therapanacea)	T1-weighted GRE	Brain	CE, FDA 2021–2022
	T1-weighted GRE, T2-weighted	Pelvis	
Adapt-Box® (Therapanacea)	Cone beam CT	Pelvis	CE 2023, FDA pending

IMAGEN

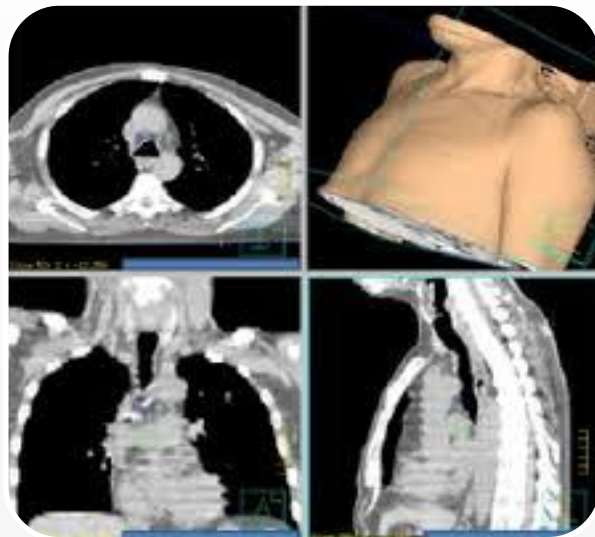


Reconstrucción CT/RM

- ↑ calidad imagen
- ↓ ruido
- ↓ artefactos
- Supervisados
- U-Net, GAN



IMAGEN



Generación/Reconstrucción CT/RM

- ↑ calidad imagen
- ↓ ruido
- ↓ artefactos
- Supervisados
- U-Net, GAN

Segmentación

- CT, PET, RM
- ↓ tiempo
- ↑ estandarización
- Supervisados
- U-Net, ResNet, DenseNet

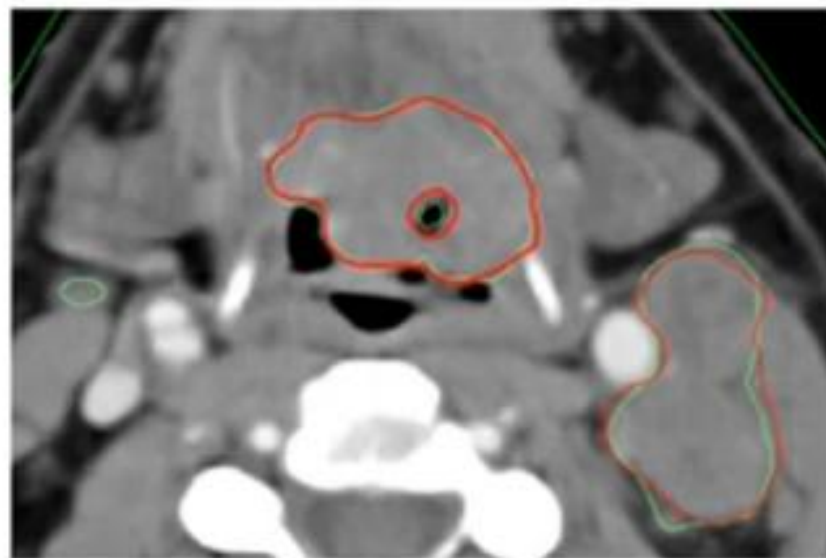
NO SOLO OARS!

Original Article

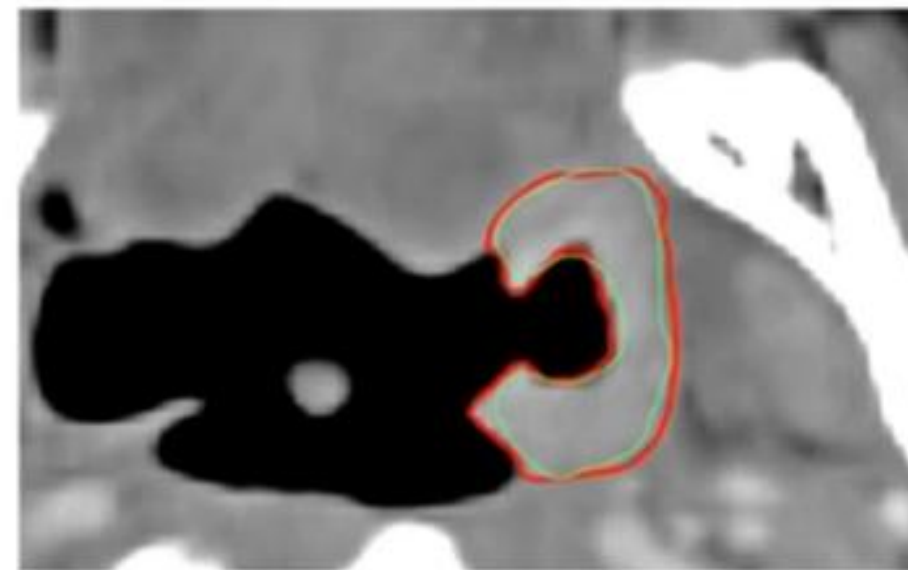
Benefits of automated gross tumor volume segmentation in head and neck cancer using multi-modality information

Heleen Bollen^{a,*,1}, Siri Willems^{b,1}, Marilyn Wegge^a, Frederik Maes^b, Sandra Nuyts^a

^aKU Leuven, Dept. Oncology, Laboratory of Experimental Radiotherapy, & UZ Leuven, Radiation Oncology; and ^bKU Leuven, Dept. ESAT, Processing Speech and Images (PSI), & UZ Leuven, Medical Imaging Research Center, B-3000 Leuven, Belgium

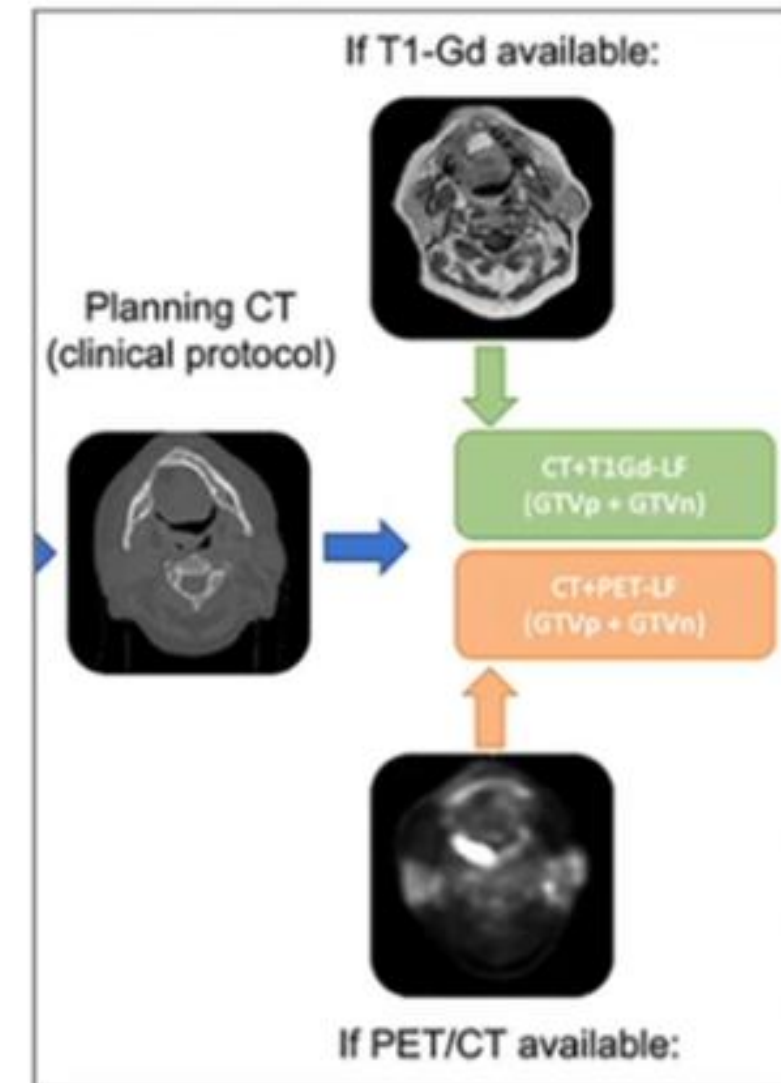


(g)



(h)

Fig. 2. Examples of automated GTV delineations by the LF CNN for CT + PET (a-d) and CT + MRI (e-h). Red, manual ground truth; green, automated delineation.

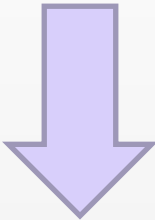


- 150 casos clínicos H&N
- CT, CT+PET o CT+MRI
- CNN con multimodalidad mejor que CNN solo con CT

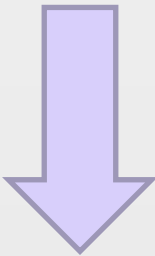
The time necessary to manually delineate the whole GTV was on average 15,5 min per case, while correcting the automated delineations by applying only necessary corrections took 8min on average per case, i.e. an increase in efficiency of 48%.

DL auto-segmentation (1 año uso clínico)

- 1. Evaluar el impacto del autocontorneo en el tiempo
- 2. Evaluar el impacto sobre el tiempo total para completar la planificación del tto (desde el CT hasta el inicio del TTO)



- 1. Se ahorró un 70% aprox
- 2. El tto no se inició antes a pesar de reducirse el tiempo de contorneo



Necesidad de revisar el workflow y adaptarlo a las nuevas necesidades



ELSEVIER

Contents lists available at ScienceDirect

Radiotherapy and Oncology

journal homepage: www.thegreenjournal.com



Original Article

Real world AI-driven segmentation: Efficiency gains and workflow challenges in radiotherapy

Ciaran Malone^{a,*}, Jill Nicholson^{a,b}, Samantha Ryan^a, Pierre Thirion^{a,c}, Ruth Woods^a, Peter McBride^a, Orla McArdle^{a,d}, Frances Duane^{a,b}, Gerard G. Hanna^{a,b}, Brendan McClean^{a,e}, Sinead Brennan^{a,b}

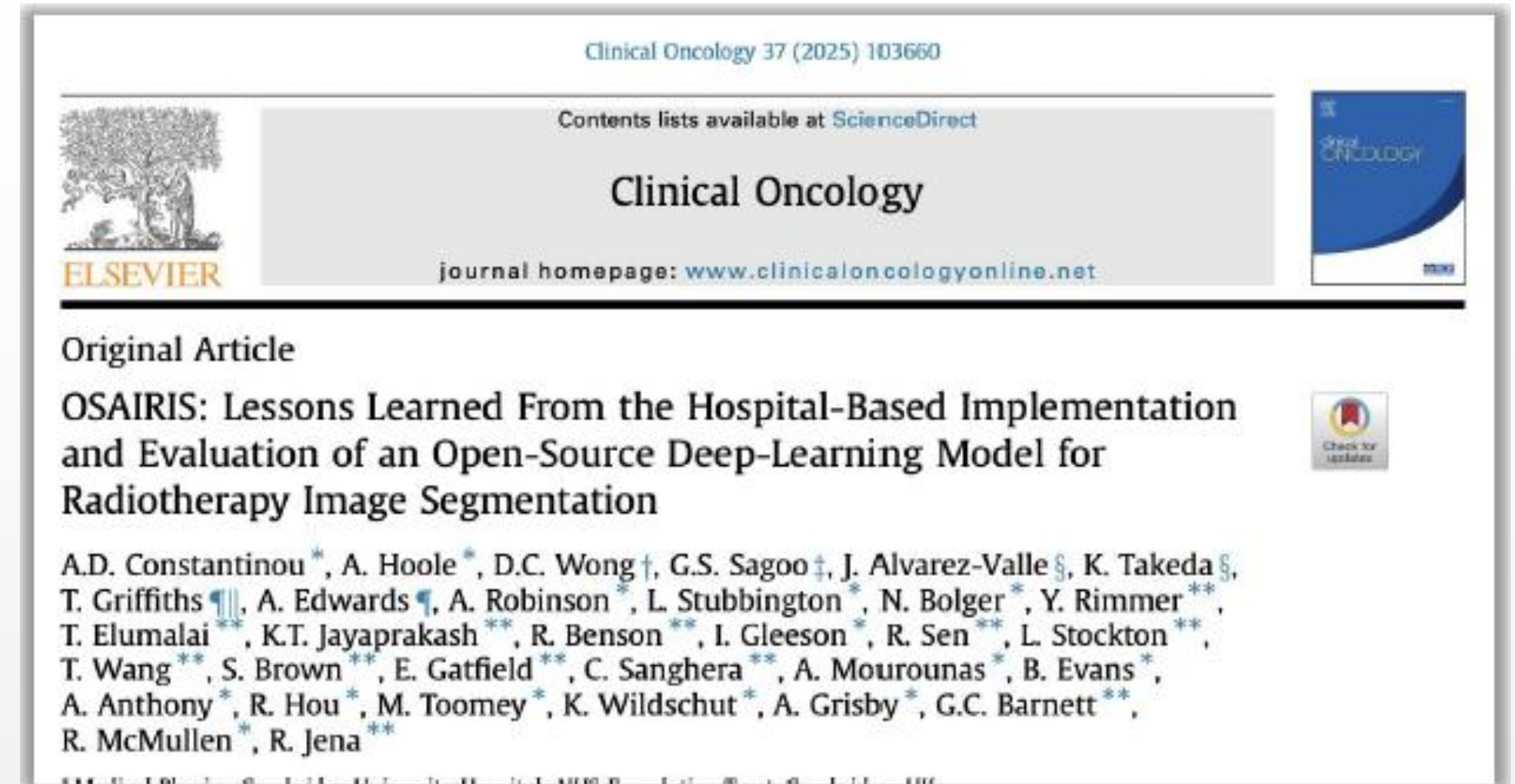
^a St. Luke's Radiation Oncology Network, Dublin, Ireland
^b Applied Radiation Therapy Trinity, Discipline of Radiation Therapy & Trinity St James's Cancer Institute, Trinity College Dublin, Dublin, Ireland
^c School of Medicine, Discipline of Radiation Therapy & Trinity St James's Cancer Institute, Trinity College Dublin, Dublin, Ireland
^d Centre for Physics in Health and Medicine, School of Physics, UCD, Dublin, Ireland
^e Beaumont RSCG Cancer Centre, Royal College of Surgeons, Ireland

Site	N		Median (Hrs)			Significance p
	Pre	Post	Pre	Post	%	
All	3472	3820	3.3	1.6	-51.5%	<0.001
Thorax	594	489	4	1.3	-67.5%	<0.001
Pelvis	1021	1422	3.8	2.4	-36.8%	<0.001
H&N	821	912	3.4	1.1	-67.6%	<0.001
Breast	220	310	2.7	2.6	-3.7%	0.152
Abdomen	68	51	4	1.1	-72.5%	0.004
Rescan	369	292	1.8	1.4	-22.2%	0.022
Other	379	344	1.9	0.7	-63.2%	<0.001

A diferencia de otros trabajos, no se centra solo en métricas de precisión (como Dice), sino en **tiempo real de edición, utilidad práctica y seguridad clínica**.

Principales hallazgos

- **Ahorro de tiempo significativo:**
 - 36 % en próstata
 - 67 % en cabeza y cuello
- **Hausdorff distance** fue la métrica que **mejor correlacionó con el tiempo de edición** ($\rho = 0.70$), superando al Dice, que no refleja bien el esfuerzo clínico.
- Se detectó **“automation bias”**: los clínicos tienden a confiar en los contornos de la IA y a anclarse en ellos, aunque la mayoría de errores fueron corregidos (72 % en H&N, 81 % en próstata).
- En una evaluación ciega, los clínicos **prefirieron los contornos de la IA** frente a los de colegas humanos (“gold standard”).



LLM-driven multimodal target volume contouring in radiation oncology

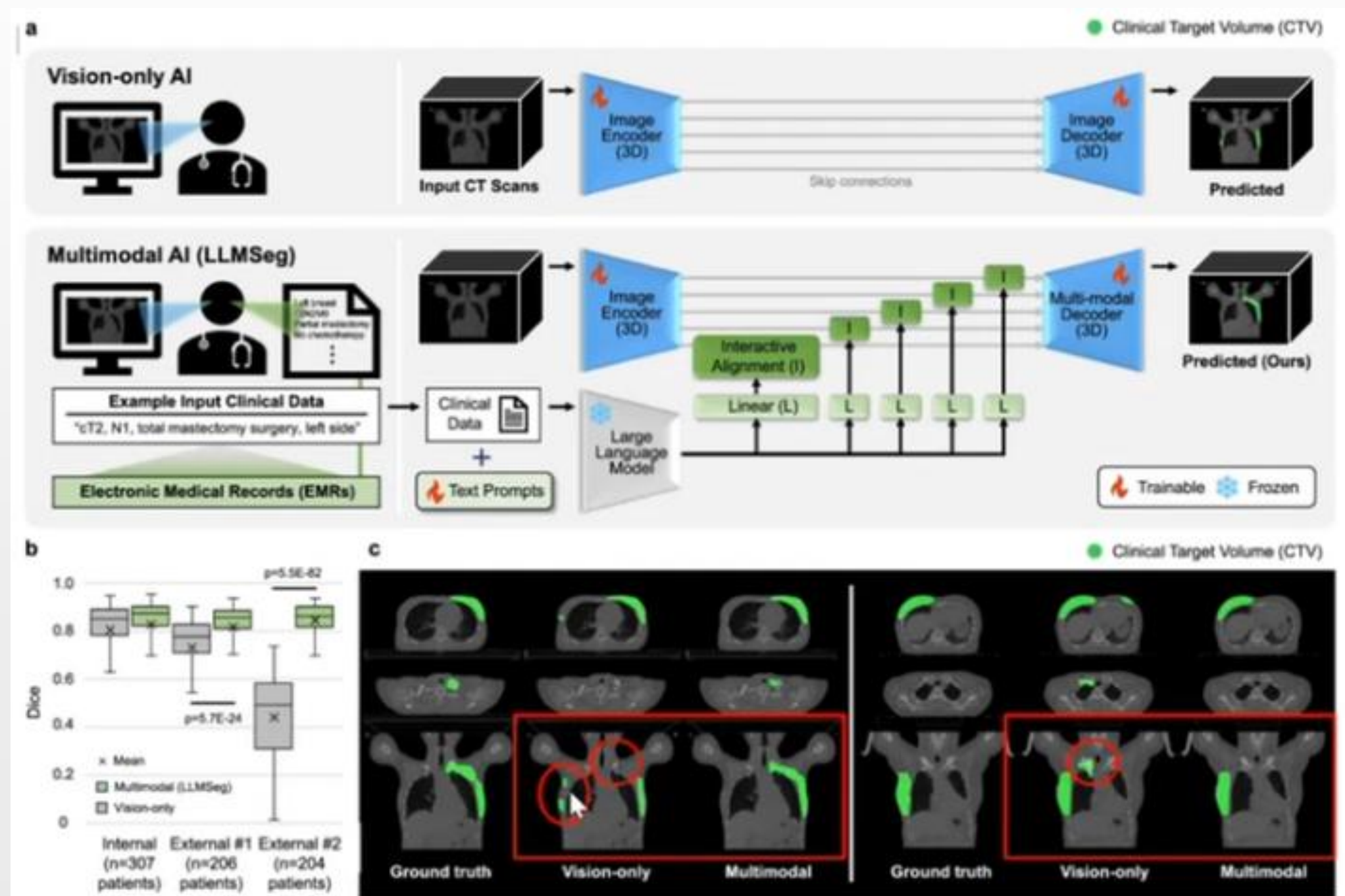
[Yujin Oh](#), [Sangjoon Park](#), [Hwa Kyung Byun](#), [Yeona Cho](#), [Ik Jae Lee](#), [Jin Sung Kim](#)  & [Jong Chul Ye](#) 

[Nature Communications](#) **15**, Article number: 9186 (2024) | [Cite this article](#)

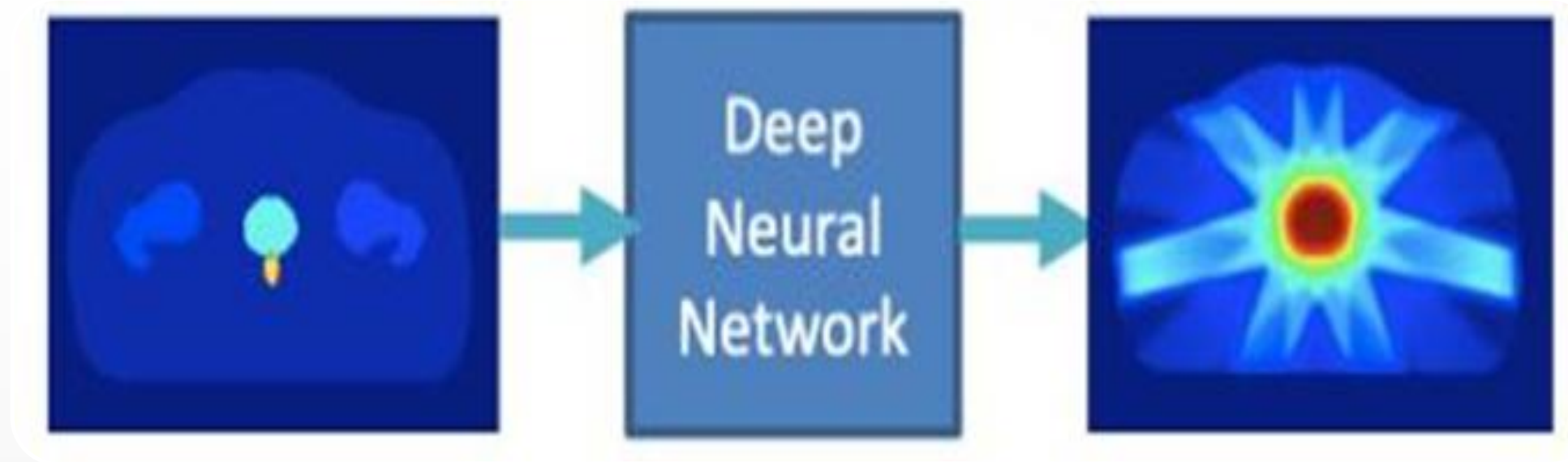
- Muestra que **la incorporación de información clínica textual** junto con imagen mejora la automatización de tareas de contorno,
- Refuerza la idea de que **menos datos pueden bastar** si se usan de forma multimodal — algo interesante si en tu entorno tienes limitaciones de datos de entrenamiento.
- Plantea un camino hacia modelos explicables que integran “razones clínicas” en la predicción, lo cual puede mejorar la confianza del clínico en la IA.

La validación se limita a ciertos tipos de cáncer (mama y próstata) con volúmenes de contorno relativamente estandarizados; no está validado aún para tumores más complejos o menos estructurados.

Lo validan sobre un conjunto de datos de cáncer de mama (y también aplican a próstata) con imagen de simulación de radioterapia + información clínica.



PREDICCIÓN DE DOSIS



- Guía para planificación de ttos
- Planificación automática
- Predicción de DVH, distribución dosis
- Comparación de modalidades (IMRT/VMAT/3D...)
- Cálculo de dosis
- Selección de ángulos
- Planificación interactiva con lenguaje natural

- Supervisados
- 3D
- 2D

Breast radiotherapy planning: A decision-making framework using deep learning

Pedro Gallego^{1 2}, Eva Ambros³, Jaime Pérez-Alíja¹, Nuria Iornet¹, Cristina Anson¹,
 Natalia Tejedor¹, H...
 Victor Riu¹, Javier R...

Affiliations + expar

PMID: 39625151 PI

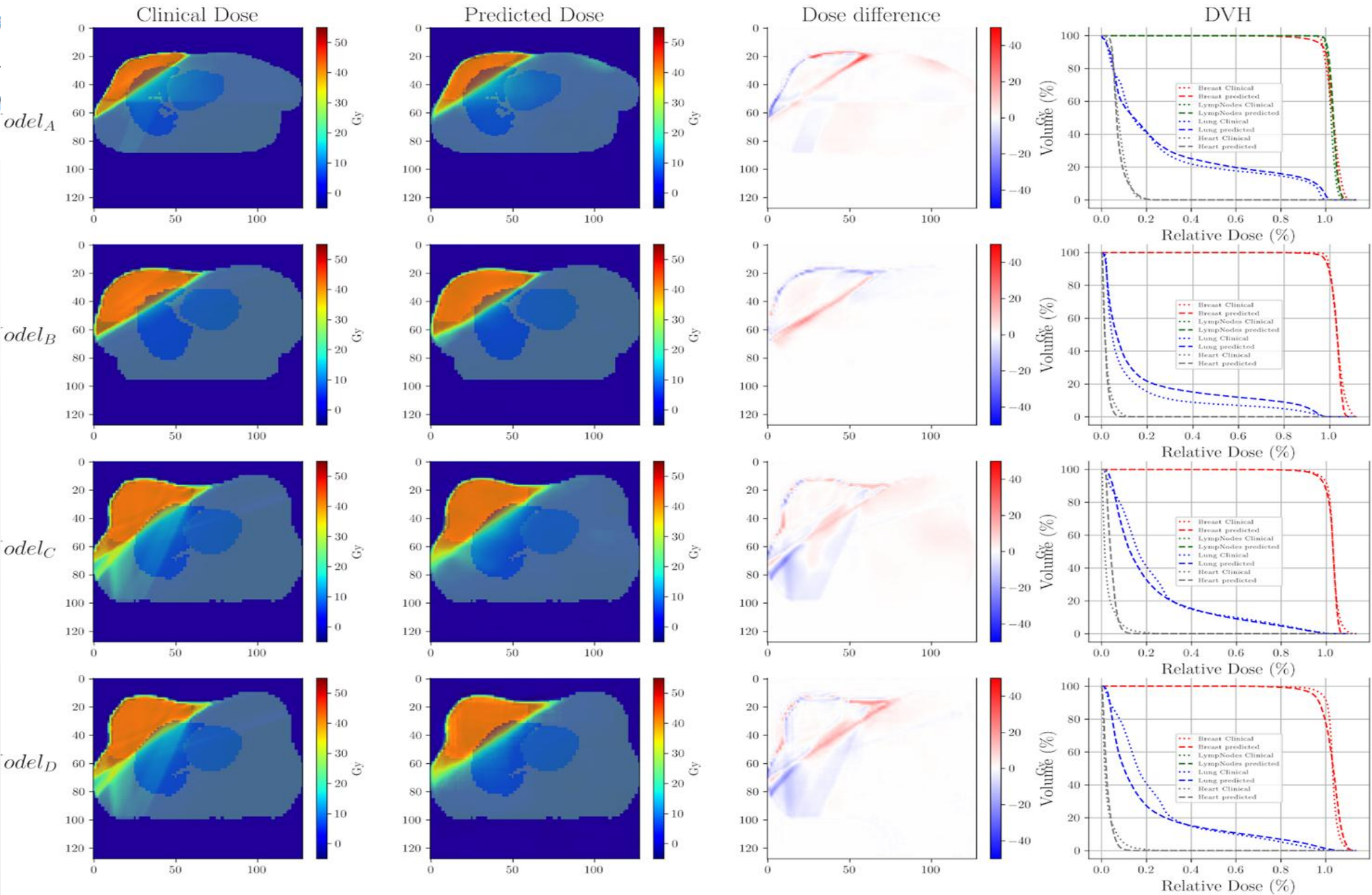


TABLE 3 Confusion matrix between the decisions made by the independent observer (ground truth) against the clinical decisions made historically.

	Actually 3D-CRT	Actually IMRT
Historical 3D-CRT	11	3
Historical IMRT	12	4

Abbreviations: 3DCRT, three-dimensional conformal radiation therapy; IMRT, intensity-modulated radiation therapy.

To test the decision-making framework, *Model_B* and *Model_C* were used, with the latter chosen due to its better performance compared to *Model_D* (Table 2, Figure 6).In comparing the historical decision and the decision-making framework, as shown in Table 3, the results exhibit a mix of true positives, true negatives, false positives, and false negatives. The recall rate for this comparison was 47.8%, precision was 78.6%, accuracy was 50%, and the F1 score was 59.5%.

TABLE 4 Confusion matrix between the decisions made by the independent observer (ground truth) against the decision making framework.

	Actually 3D-CRT	Actually IMRT
Predicted 3D-CRT	22	2
Predicted IMRT	1	5

Abbreviations: 3DCRT, three-dimensional conformal radiation therapy; IMRT, intensity-modulated radiation therapy.

In the second comparison between the independent observer's decision (ground truth) and the decision-making framework, represented in Table 4, the outcomes were more favorable. The recall rate in this comparison was 95.7%, precision was 91.7%, accuracy was 90.0%, and the F1 score was 93.6%.



Clinical integration of machine learning for curative-intent radiation treatment of patients with prostate cancer

Chris McIntosh^{1,2,3,4,5,6,8}, Leigh Conroy^{1,2,7,8}, Michael C. Tjong^{1,7}, Tim Craig^{1,2,7}, Andrew Bayley^{1,7}, Charles Catton^{1,7}, Mary Gospodarowicz^{1,7}, Joelle Helou^{1,7}, Naghmeh Isfahanian^{1,7}, Vickie Kong^{1,7}, Tony Lam^{1,7}, Srinivas Raman^{1,7}, Padraig Warde^{1,7}, Peter Chung^{1,7}, Alejandro Berlin^{1,2,7,8} and Thomas G. Purdie^{1,2,6,7,8} ✉

- **Aceptabilidad clínica:** 89% de los planes ML fueron considerados aptos para tratar pacientes.
- **Selección plan tto:** el 72% ML vs humanos (83% en simulación, 61% en despliegue real).
- **Eficiencia:** el tiempo total de planificación se redujo en un 60% (de 118 h a 47 h).
- **Calidad dosimétrica:** los planes ML mantuvieron calidad constante entre fases y se ajustaron a las guías clínicas.
- La menor selección de planes ML en la fase clínica no se debió a una pérdida de calidad, sino a **factores humanos** (percepción y confianza del médico tratante).

- **Random Forest** con 99 planes clínicamente aprobados.
- Se realizaron tres fases:
 1. Factibilidad técnica (n=17)
 2. Simulación retrospectiva (n=50)
 3. Despliegue clínico prospectivo (n=50), con selección ciega entre plan ML y plan humano; el plan elegido se utilizó realmente en el tratamiento.

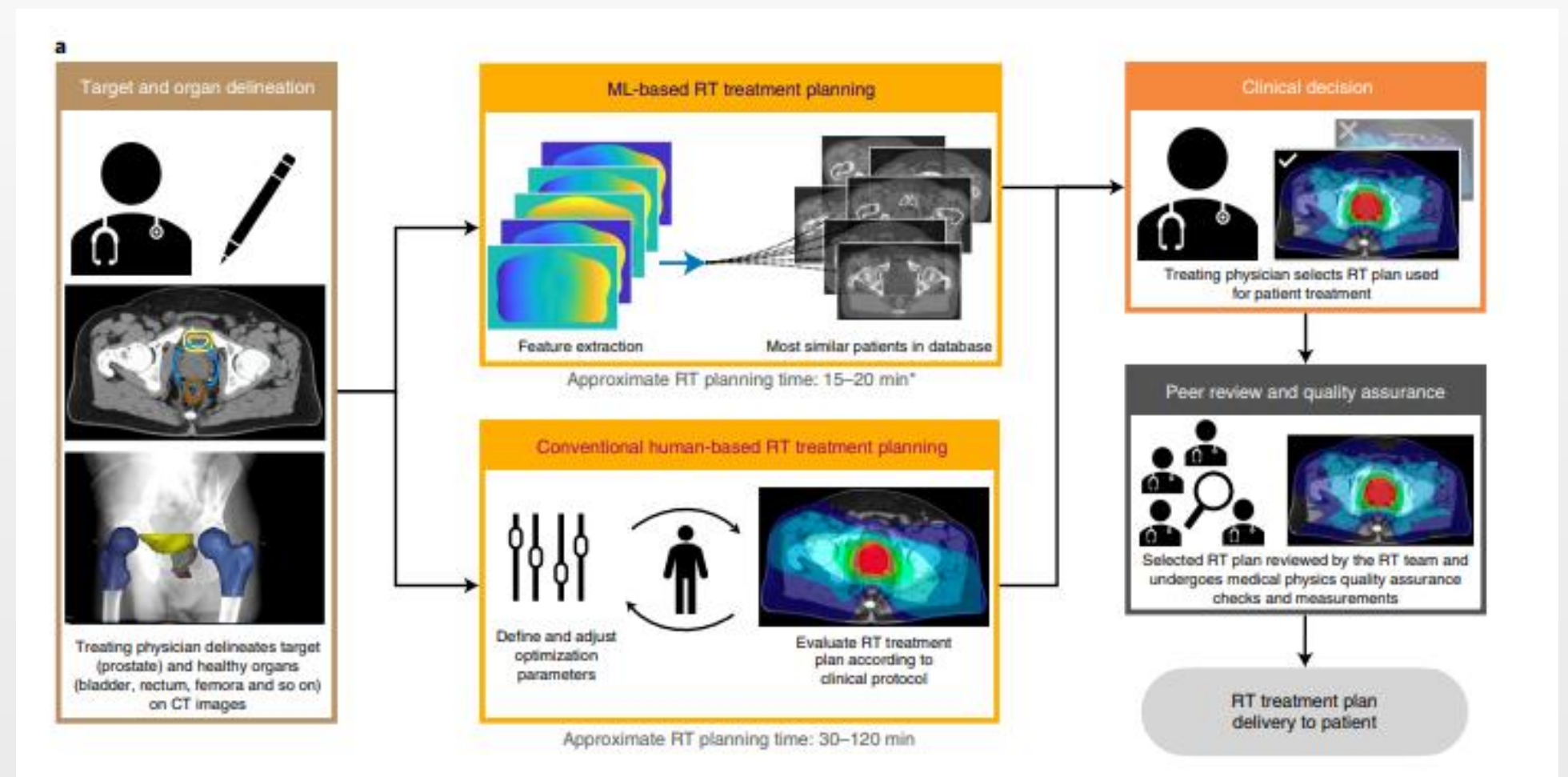
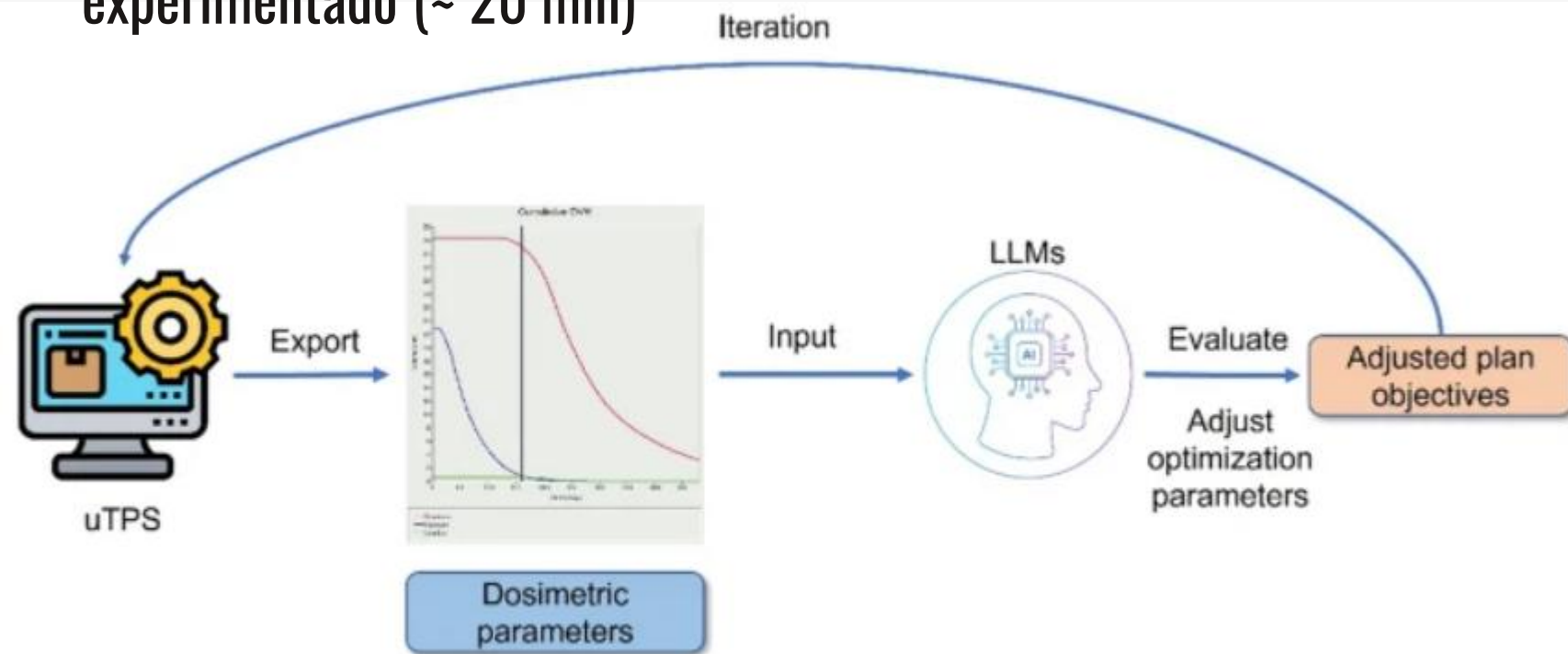




Fig. 1. Overview of the DLO implementation workflow.

- Deep learning planning automates robust proton therapy planning for oropharyngeal cancer.
 - Plan quality was evaluated in blinded retrospective and prospective studies.
 - Deep learning plans were preferred over or comparable to manual plans in 92% of the patients.
-
- Aunque los resultados son prometedores, todavía se requiere supervisión humana para revisar y aprobar los planes generados automáticamente.
 - La generalización a otros sitios anatómicos, otros centros con distintos sistemas de protonterapia o diferentes protocolos aún necesita validación adicional.
 - Como siempre en IA aplicada al entorno clínico, la gestión del cambio, la aceptación del usuario y la integración en el flujo de trabajo son factores clave que deben considerarse más allá de la técnica.

- 35 pacientes cérvix
- Conformidad PTV y OARs
- Varios modelos LLM (2.5-max, LLama-3.2, Gemini 1.5 flash)
- Un modelo generó halucinaciones
- Dos de los LLMs (Qwen-2.5-max y Llama-3.2) generaron planes aceptables en tiempos de $\approx 16.3 \pm 5.0$ min y $\approx 9.8 \pm 2.1$ min, respectivamente, mejorando el tiempo estimado de un físico experimentado (≈ 20 min)



The workflow of using LLM for radiotherapy treatment planning

Parameters fed into the LLMs comprised current doses, dose goals, constraint target values, ideal ranges of constraint target values for OAR optimization objectives, and PTV D95 values.

Research | [Open access](#) | Published: 15 May 2025

Feasibility study of automatic radiotherapy treatment planning for cervical cancer using a large language model

[Shuoyang Wei](#), [Ankang Hu](#), [Yongguang Liang](#), [Jingru Yang](#), [Lang Yu](#), [Wenbo Li](#), [Bo Yang](#) & [Jie Qiu](#)

Radiation Oncology **20**, Article number: 77 (2025) | [Cite this article](#)

2022 Accesses | 2 Citations | 1 Altmetric | [Metrics](#)

A feasibility study of automating radiotherapy planning with large language model agents

Qingxin Wang, Zhongqiu Wang, Minghua Li, Xinye Ni, Rong Tan, Wenwen Zhang, Maitudi Wubulaishan, Wei Wang, Zhiyong Yuan, Zhen Zhang [Show full author list](#)

Published 21 March 2025 • © 2025 Institute of Physics and Engineering in Medicine. All rights, including for text and data mining, AI training, and similar technologies, are reserved.

[Physics in Medicine & Biology](#), Volume 70, Number 7

3764 • Volume 123, Issue 1, Supplement , E799, September 01, 2025

Development and Performance Evaluation of an LLM-Based Automated Recommendation System for Radiotherapy Treatment Planning

[Y.X. Yang](#)¹ · [L. Jia](#)² · [H. Li](#)² · ... · [G.Y. Wang](#)¹ · [G. Zhou](#)¹ · [Y. Sun](#)⁵ ... [Show more](#)

[Affiliations & Notes](#) [Article Info](#)

169 • Volume 123, Issue 1, Supplement , S68, September 01, 2025

Automating Head-and-Neck Cancer Intensity Modulated Radiation Therapy Treatment Planning with a ReAct Large Language Model Agent

[D. Yang](#)¹ · [X. Wu](#)¹ · [Y. Xie](#)¹ · ... · [Q. Wu](#)¹ · [Q.J.J. Wu](#)² · [Y. Sheng](#)¹ ... [Show more](#)

Aplicaciones clínicas desarrolladas

<http://www.dlinrt.eu>



Treatment Planning RT Plan

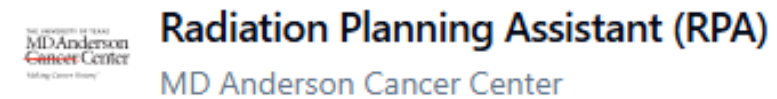
AI-powered treatment planning assistant that optimizes radiation therapy plans using machine learning algorithms for improved plan quality and...

Key Features:

- DL-driven plan optimization
- Full workflow automatin from dose prediction to final plan
- Auto-planning models validated on standard protocols

+1 more features

[Visit](#)



Treatment Planning CT

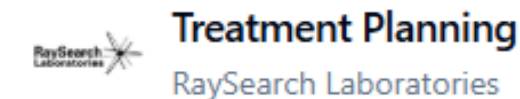
The Radiation Planning Assistant offers a suite of fully automated contouring and radiotherapy planning tools for various anatomical sites including cervix,...

Key Features:

- Fully automated contouring tools
- Comprehensive radiotherapy planning suite
- Multi-anatomical site coverage

+2 more features

[Visit](#)



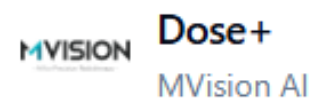
Treatment Planning CT

Advanced treatment planning system with AI-enhanced optimization algorithms (i.e. deep learning dose prediction followed by dose mimicking) for...

Key Features:

- Auto-planning with deep learning driven dose prediction
- Customizable post-processing to match your protocols, clinical goals, and machines
- Models validated against common radiotherapy protocols

[Visit](#)

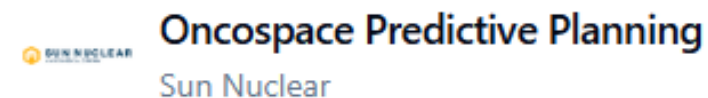


Treatment Planning CT

Dose+ is an AI-based treatment planning solution that automates VMAT and IMRT radiotherapy plans creation, reducing planning time from hours to...

Key Features:

- Fully automated treatment planning



Treatment Planning Treatment Planning System Data

Cloud-based, AI-powered solution that uses machine learning models to derive achievable, best-practice dosimetric goals for plan optimization and evaluatio...

Key Features:

- Machine learning-derived dosimetric goals

[Visit](#)

+2 more features



Machine learning applications in radiation oncology: Current use and needs to support clinical implementation

Charlotte L. Brouwer^{a,*}, Anna M. Dinkla^b, Liesbeth Vandewinckele^{c,d}, Wouter Crijns^{c,d}, Michaël Claessens^{e,f}, Dirk Verellen^{e,f}, Wouter van Elmpt^g

- Survey
 - 213 responders in 202 radiotherapy clinics
 - 37% of responders implemented at least 1 ML application clinically
- Main reasons introducing AI:
 - Time saving
 - Quality improvement (also consistency)
- Main reasons **not** introducing AI:
 - Lack of knowledge how to implement
 - No software available & lack of time
 - Resistance or scepticism against AI

Inconvenientes principales de la adopción de la IA en la clínica



- Ética
- Transparencia de los modelos de IA (models cards)
- De-skilling (pérdida de conocimiento)
- Opiniones de los pacientes ante el uso de IA
- Educación y formación del personal

Intelligent Machines

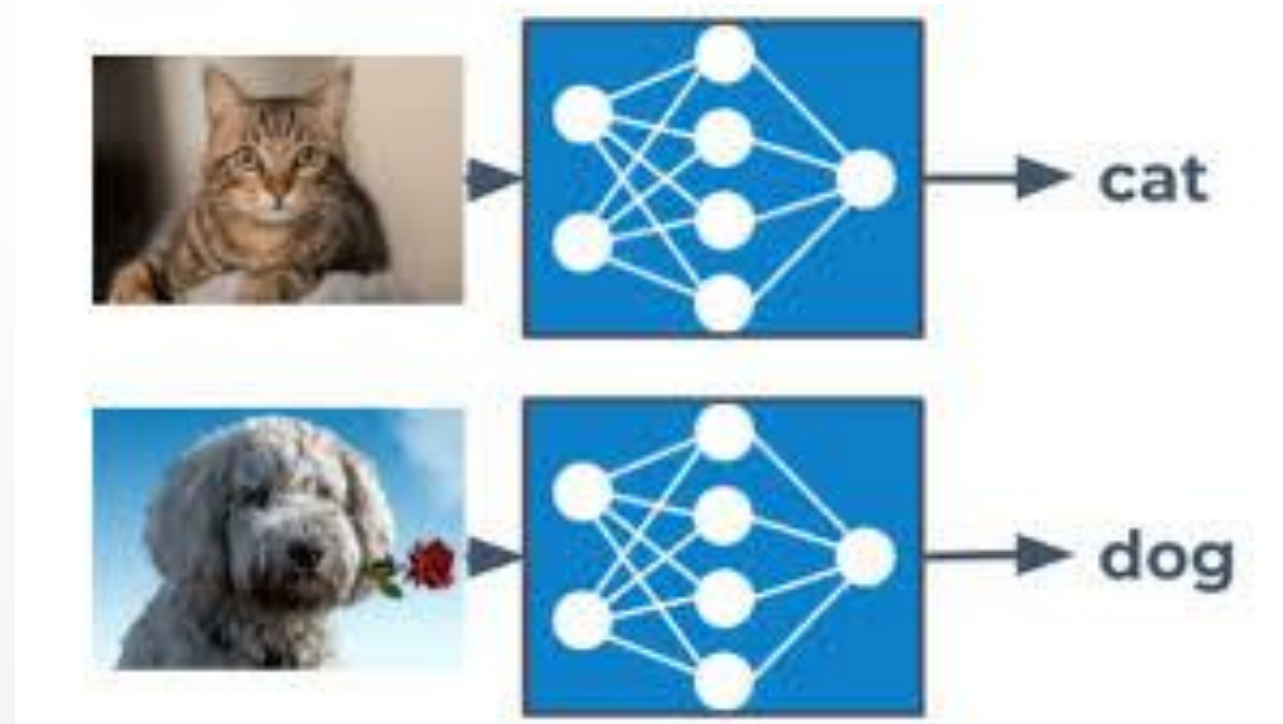
The Dark Secret at the Heart of AI

No one really knows how the most advanced algorithms do what they do. That could be a problem.

by Will Knight April 11, 2017



Ética



- toma de decisiones sobre pacientes
- necesidad de saber las razones y generar confianza sobre predicciones
- herramientas para evaluar las predicciones



Official Journal
of the European Union

EN
L series

2024/1689

12.7.2024

REGULATION (EU) 2024/1689 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

of 13 June 2024

laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act)

(Text with EEA relevance)

THE EUROPEAN PARLIAMENT AND THE COUNCIL OF THE EUROPEAN UNION,



Transparencia: MODEL CARDS

- Son como el prospecto de un medicamento, pero para un modelo de IA:
- qué hace,
- con qué datos se entrenó,
- para qué sirve
- para qué no.



Articles

Standardization of Artificial Intelligence Development in Radiotherapy

Alessia de Biase ^{† 1}, Nikos Sourlos ^{‡ 1}, Peter M.A. van Ooijen ^{§ ||}  

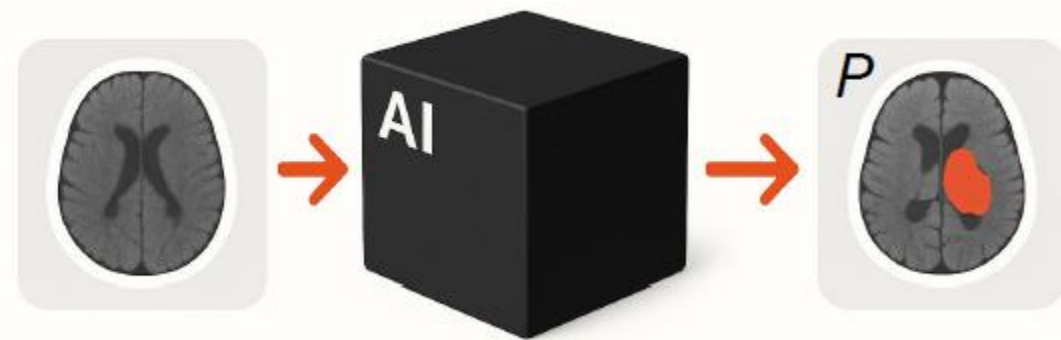
Model Card	Model Fact Labels	Fact Sheets
Model Details	Summary/ Mechanism	Statement of purpose
<p>Person or organization developing model: It was developed in 2021 by a group of PhD students in Deep Learning in Medical Imaging from the University Medical Hospital of Groningen. With the implemented models the group participated in the HECKTOR 2021 challenge.</p> <p>Model type: Co-learning method and a 3D Skip Spatial and Channel Squeeze and Excitation Multi-Scale Attention method (Skip-scSE-M), both based on Convolutional Neural Networks</p> <p>Paper or other resource for more information:</p> <p>De Biase, A., et al.: Skip-SCSE multi-scale attention and co-learning method for oropharyngeal tumor segmentation on multi-modal PET-CT images. In: Andrearczyk, V., Orellier, V., Hatt, M., Depeursinge, A. (eds.) HECKTOR 2021. LNCS, vol. 13209, pp. 109–120. Springer, Cham (2022)</p>		
Intended Use/ Factors/Ethical considerations	Uses and directions	Statement of purpose
<p>Intended Use: intended to be used by radiation oncologists before radiotherapy treatment planning for the automatic delineation of the gross tumor volume (GTVt)</p> <p>Target Population: oropharyngeal cancer patients</p> <p>Benefits: manual delineation of contours is time consuming and prone to errors</p> <p>Use case: after PET and CT images are collected for a new patient, the primary tumor contour is immediately available and the treatment planning can start</p> <p>Appropriate decision support: the model identifies a primary tumor for the patient X and the oncologist discusses the possible treatment options</p> <p>Before using this model: test the model on internal data to establish the model validity on a different dataset</p> <p>Safety and efficacy evaluation: results are promising for this task when comparing with oncologists' manual contours</p>		
Metrics/ Quantitative Analysis	Validation and Performance	Basic Performance
<p>The metrics used are the Dice Similarity Coefficient (DSC) and the Hausdorf distance at 95% (HD95) between the ground-truth annotations and the predictions of the implemented algorithms.</p> <p>All metrics are reported at a threshold of 0.5</p> <p>The best results we obtained with the proposed methods on the test set were a mean DSC of 0.762 and median HD of 3.143</p>		
Evaluation Data	< no separate section >	< no separate section >
<p>3D PET/CT scans of 101 from 2 centers (CHUP - which is also present in the training set - and CHUV)</p> <p>Data was provided by the organizers of the challenge in August 2021 as 3D bounding boxes of size 144x144x144 containing the tumor region.</p> <p>As a pre-processing technique we used a z score normalization.</p>		
Training Data	Mechanism	Lineage
<p>3D PET/CT scans and GTVt manual contours of 224 from 5 centers (CHGJ, CHMR, CHUM, CHUS and CHUP)</p> <p>Data was provided by the organizers of the challenge in June 2021 as 3D bounding boxes of size 144x144x144 containing the tumor region.</p> <p>As a pre-processing technique we used a z score normalization.</p>		
Caveats and Recommendations	Warnings	Safety & Security
<p>Risks: The segmentation model can output false positive or false negative. A missed primary tumor could lead to mortality. A patient without tumor receiving radiation therapy leads to undesired side-effects and toxicities.</p> <p>Inappropriate Settings: this model was trained on imaging data used for radiotherapy planning. Do not use on other image modalities or images collected for other purposes.</p> <p>Clinical Rationale: The model can not be interpret and does not provide rationale for good or bad segmentation results. The model outputs need to be approved by a doctor before they can be used.</p> <p>Inappropriate decision support: This model was trained and tested on patients from 5 different cohorts. It may not be accurate outside the target population.</p>		

De-skilling

- ❖ *La IA: aprende rápido, comete errores y necesita supervisión (constante).*
- podría ser necesaria la **intervención humana**, pero esto se convierte en un problema si las personas ya no son competentes en las tareas que deben realizar.
- es fundamental **mantener la capacidad de supervisar las actividades de la IA** e intervenir cuando sea necesario, lo que exige una comprensión profunda de sus propias áreas de competencia.
- ❖ *La clave para un futuro exitoso en la era de la IA reside en entender su naturaleza dual y en saber navegar con prudencia hacia una convivencia armoniosa entre la inteligencia humana y la artificial.*



Retos clínicos



EU AI Act, Article 13:

“High-risk AI systems shall be designed and developed in such a way to ensure that their operation is sufficiently transparent to enable users to interpret the system’s output and use it appropriately. An appropriate type and degree of transparency shall be ensured ...”

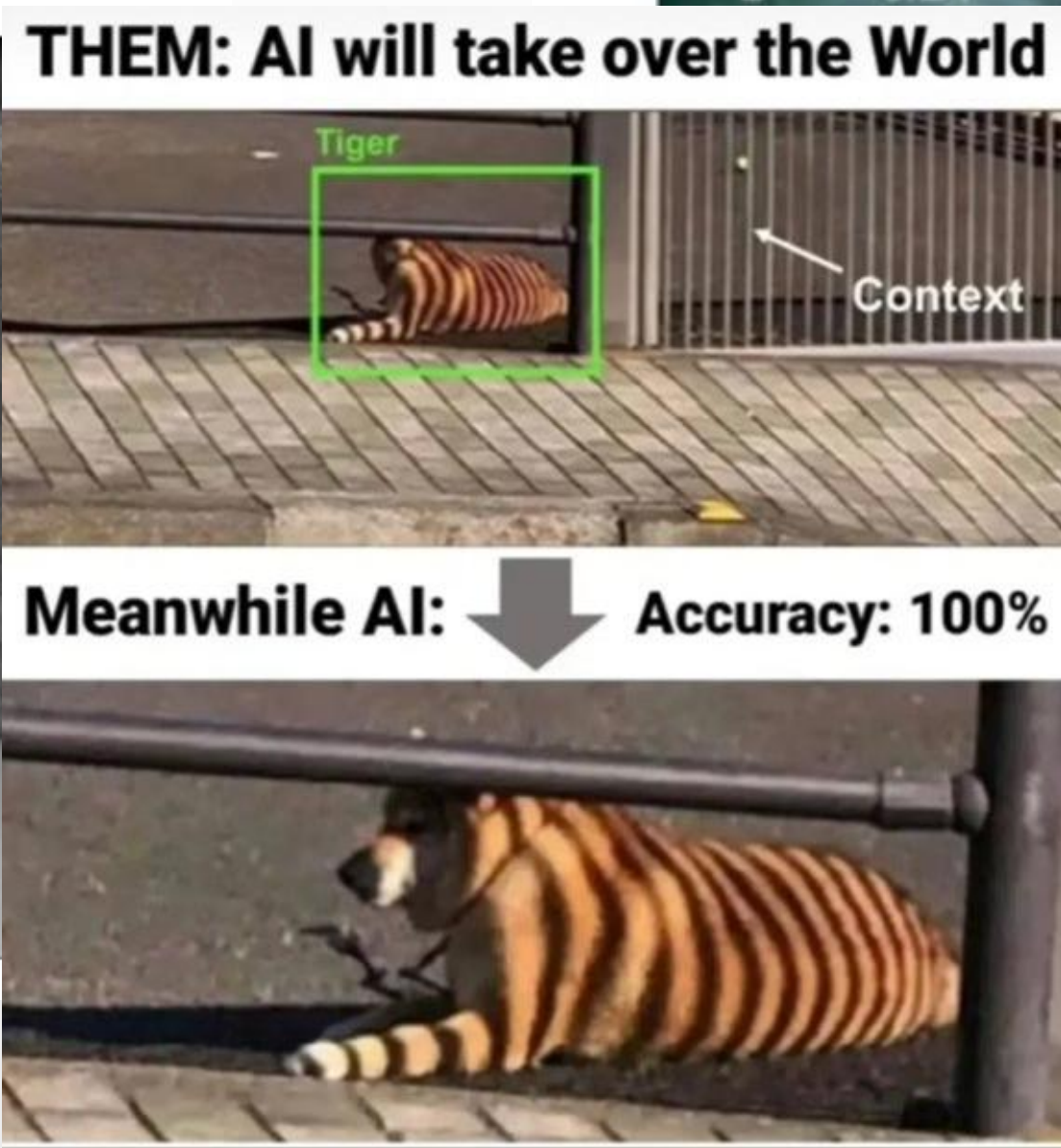
- Interpretabilidad: el propio modelo (como funciona el modelo desde dentro)
- Explicabilidad: las razones de un predicción específica
- Incertidumbre: las limitaciones del modelo o de las predicciones

IA Robusta y accesible

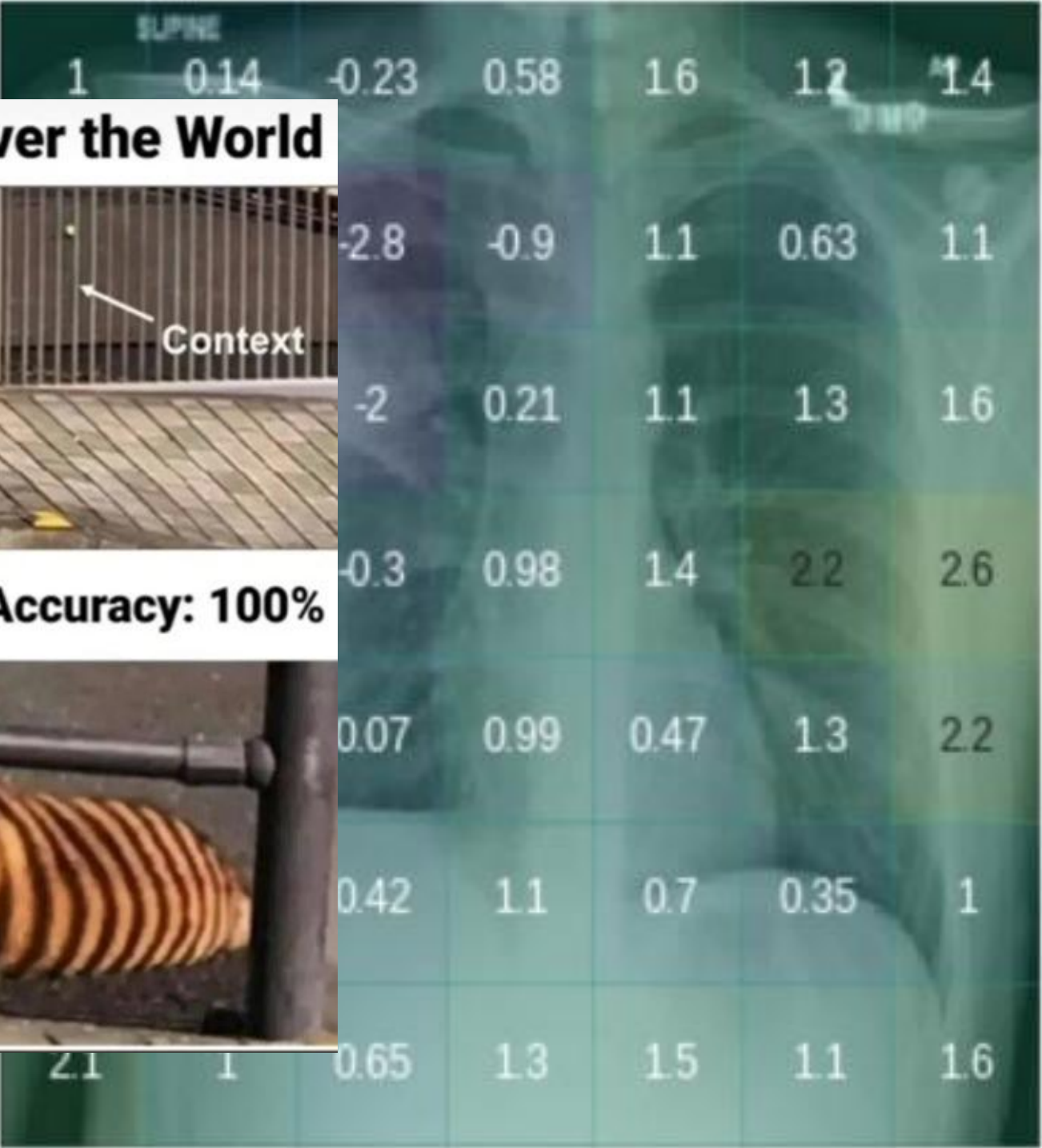
Neumonía Accuracy 88%



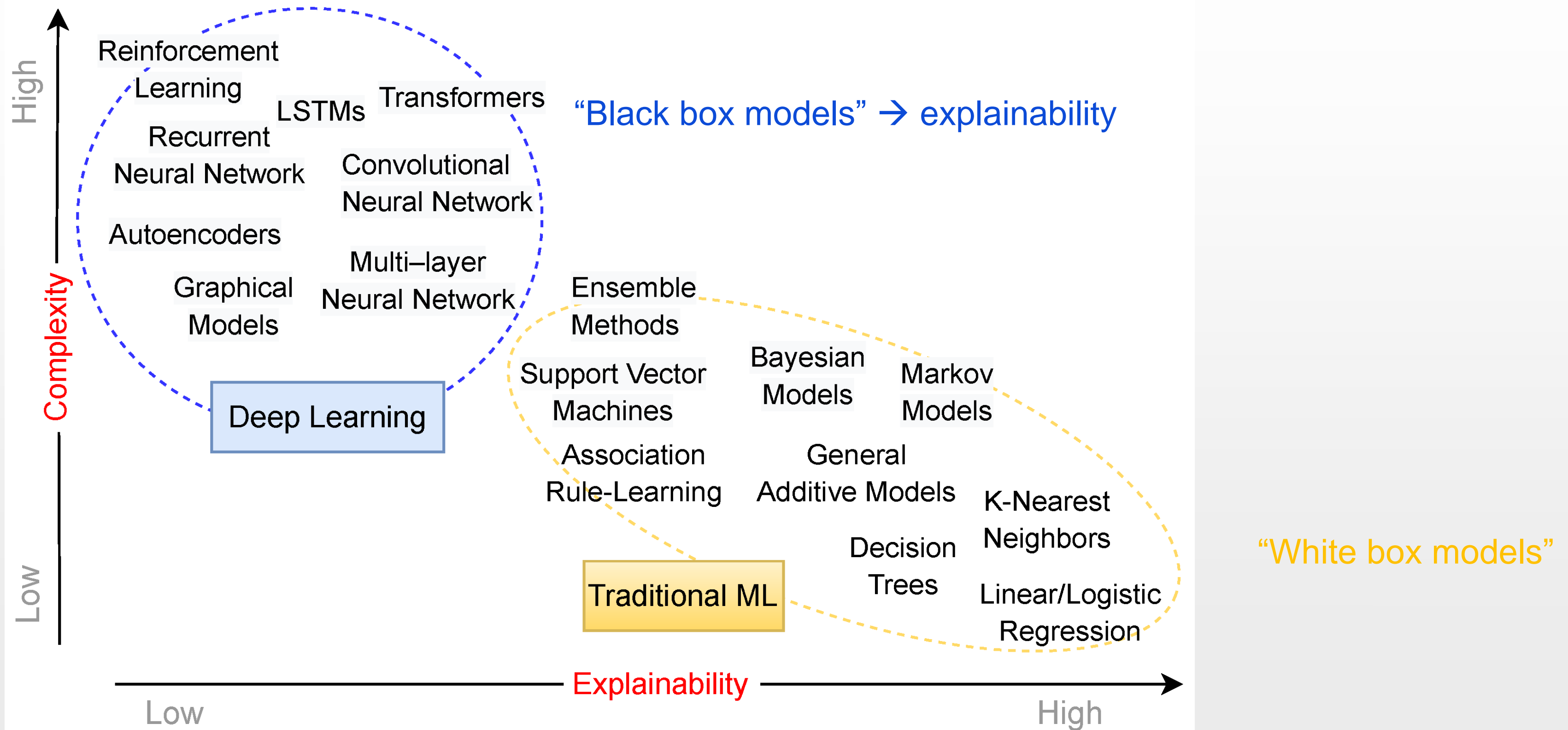
CNN on left, original image on right



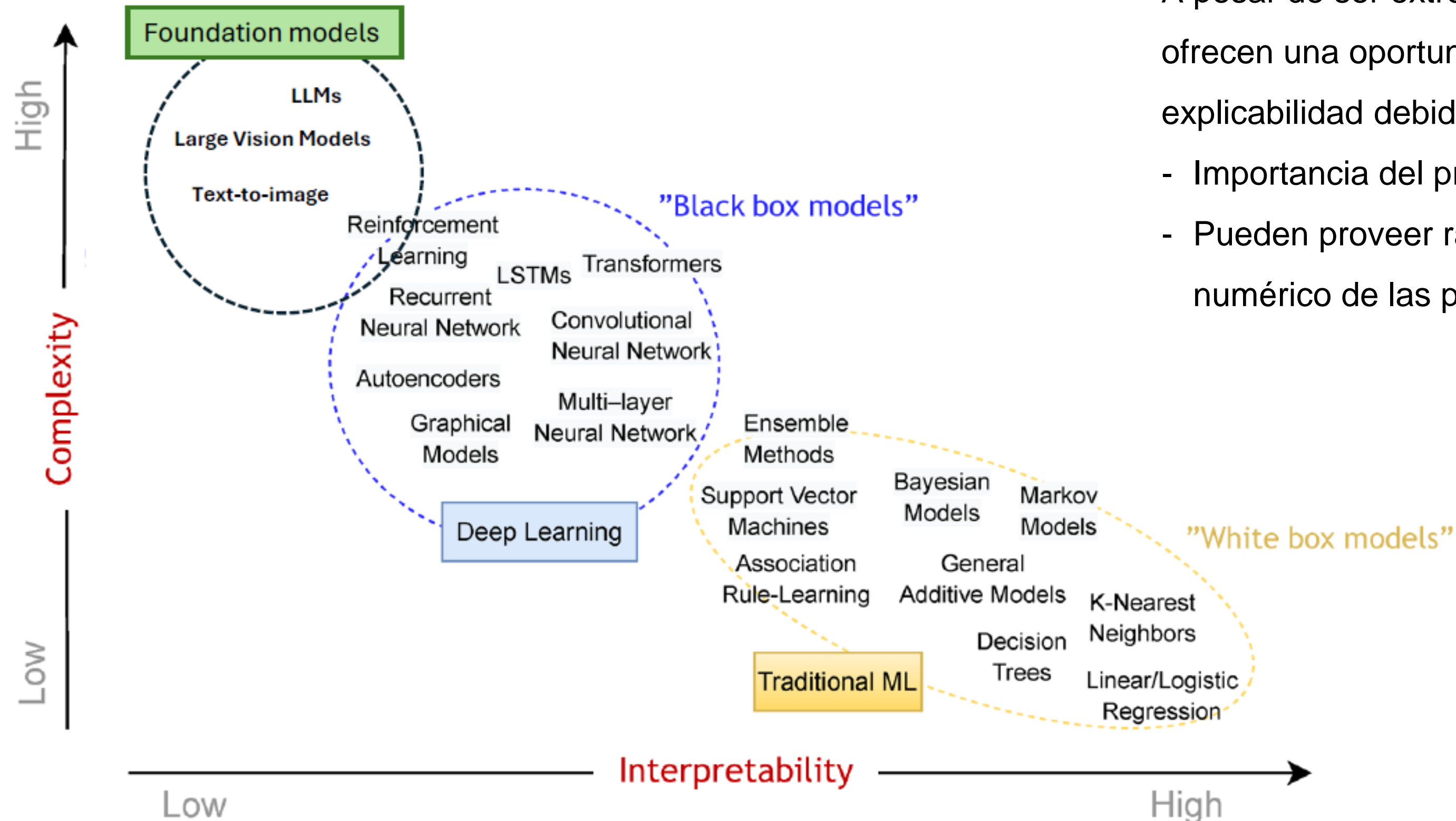
$P(\text{Pneumonia})=0.024$



La interpretación depende del modelo



La interpretación depende del modelo



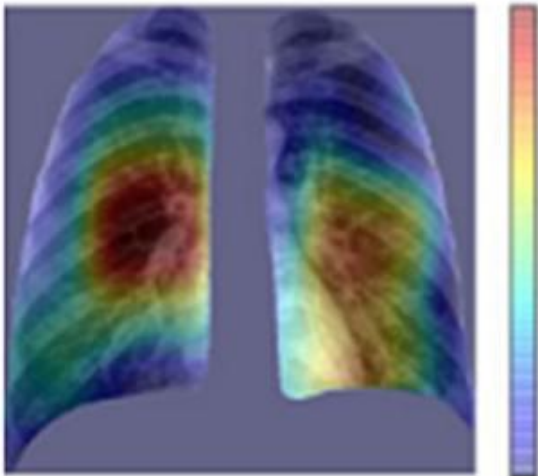
Foundation models

A pesar de ser extremadamente complejo, ofrecen una oportunidad para la explicabilidad debido a:

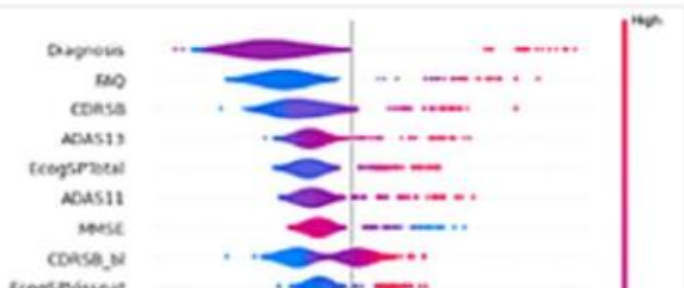
- Importancia del prompt usado
- Pueden proveer raciocinio textual/visual y numérico de las predicciones

Métodos de explicabilidad

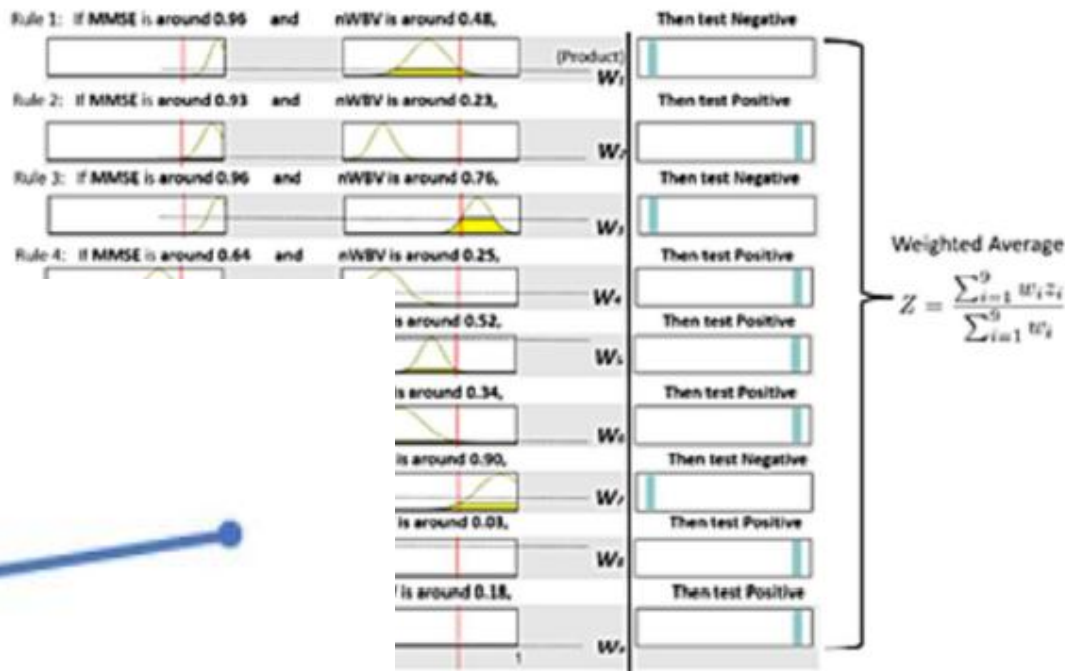
a. Visual



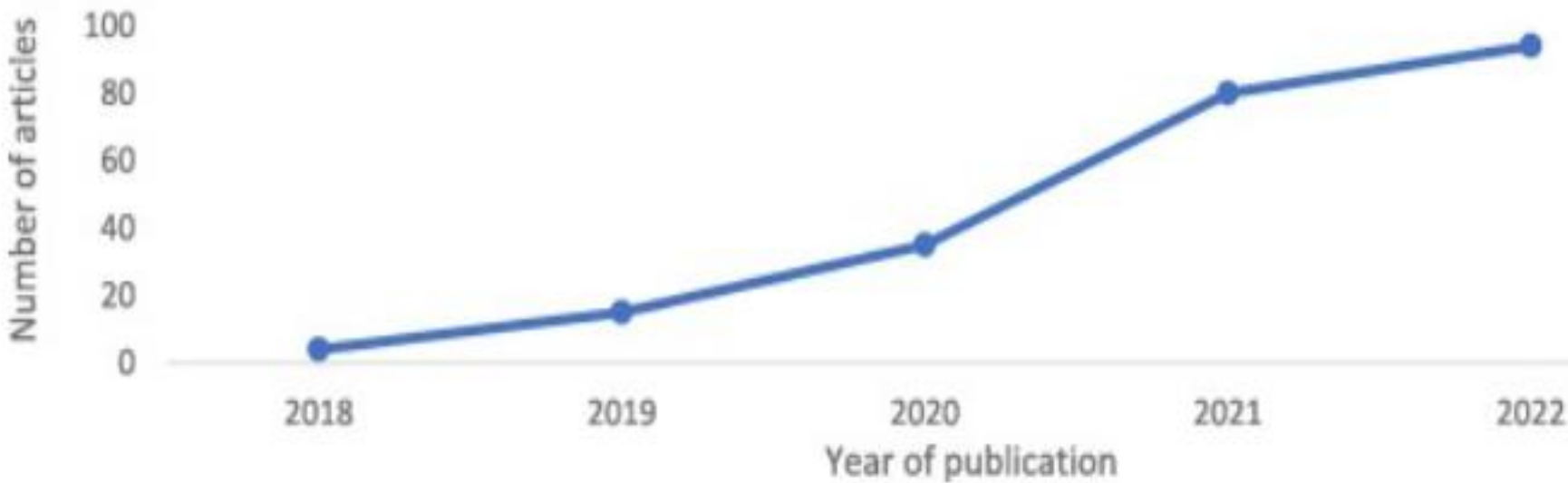
b. Numerical



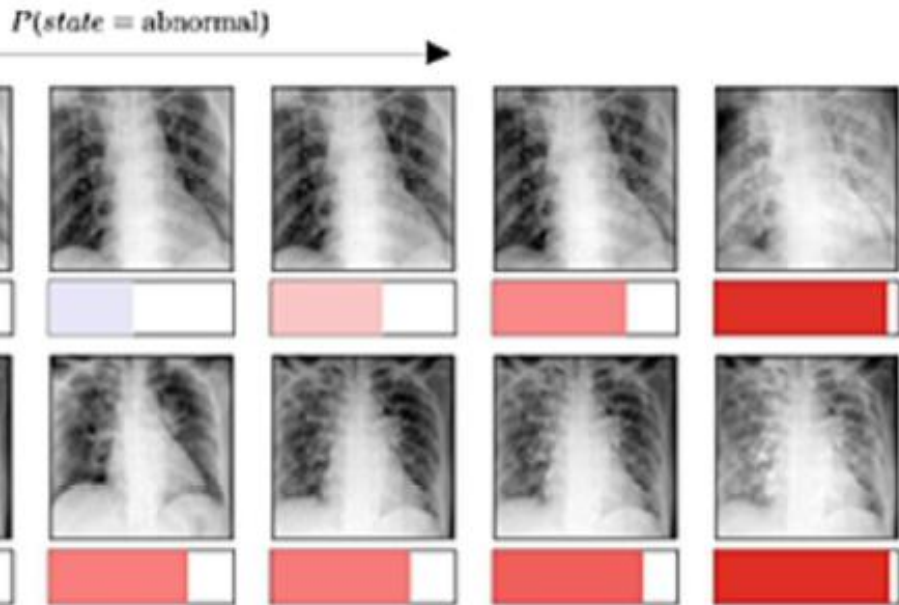
c. Rule-based



d.



e. Example-based



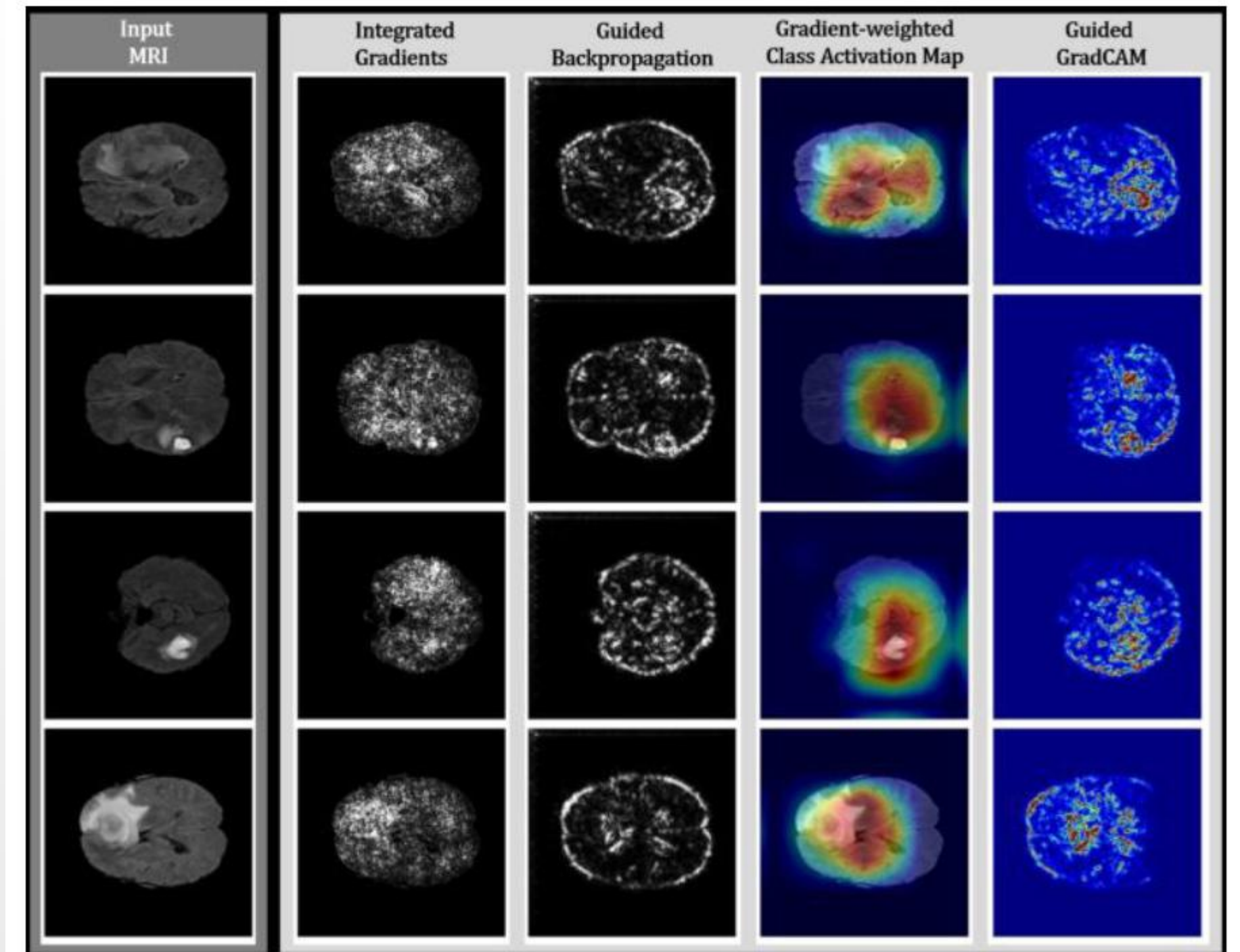
Radiology output	Diagnosis: Alzheimer's (ALZ)
	Criteria in favour AD:
	1. Hippocampal Volume is very Low.
	2. Cluster Prominence is very Low in combination with Low Volume.
	Criteria in favour NC:
	1. Entropy is High, Variance is Nominal and Contrast is not Very High
	The criteria supporting AD are in general stronger than those in favour for NC and therefore the predicted diagnosis is AD.

Explicabilidad a través de la visualización

Mapas de atención (attention maps)

Señalan que voxeles o regiones de una imagen en la que se ha enfocado el modelo usando:

- Gradientes de imagen
- Pesos
- Bordes
- Texturas
- Patrones

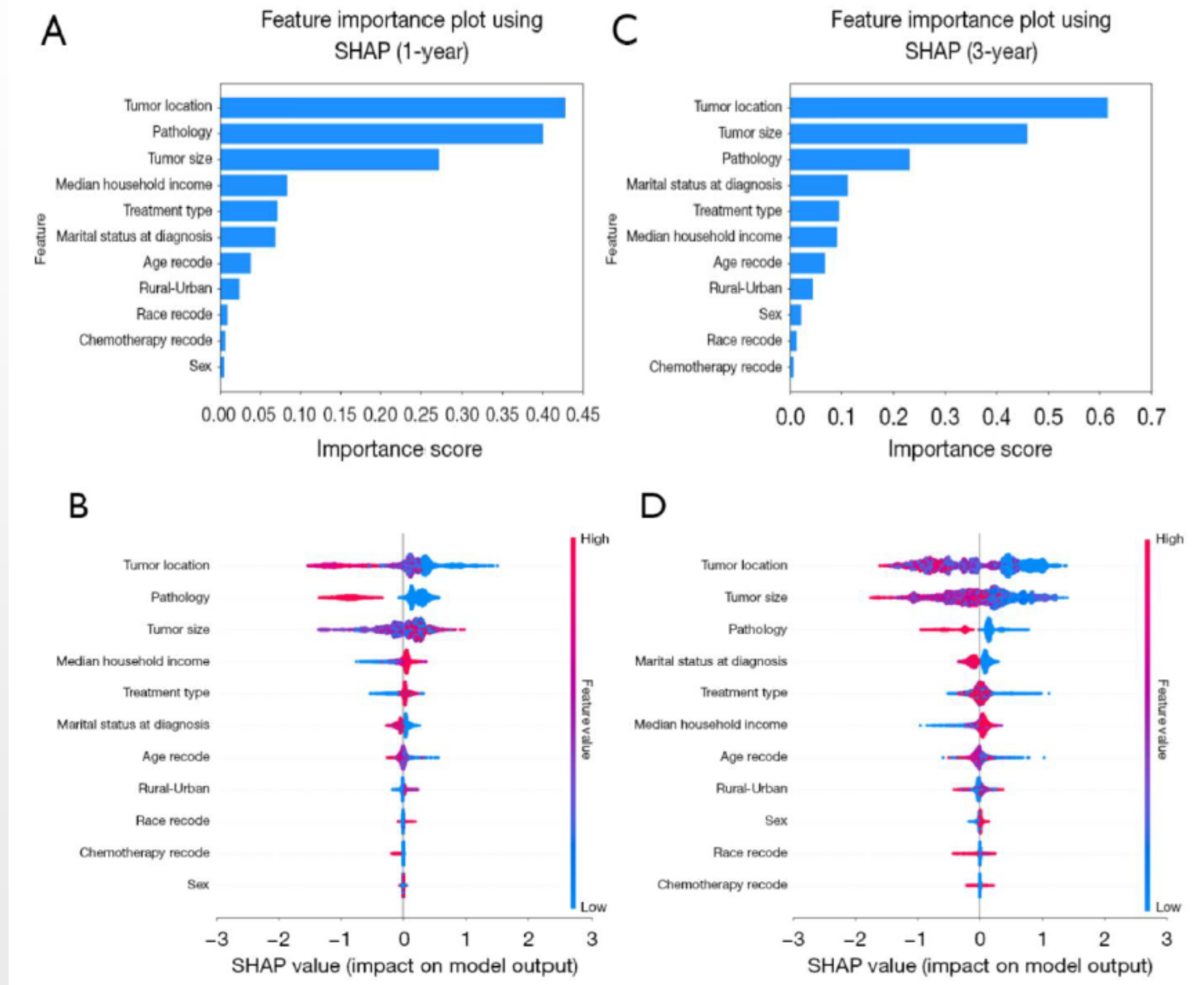


Explicabilidad a través de rangos numéricos (feature weights)

SHAP (Shapley Additive exPlanations)

Señala la importancia y el impacto de las diferentes features:

- Datos clínicos
- Biomarcadores
- Características de pacientes
- Modalidades
- ...

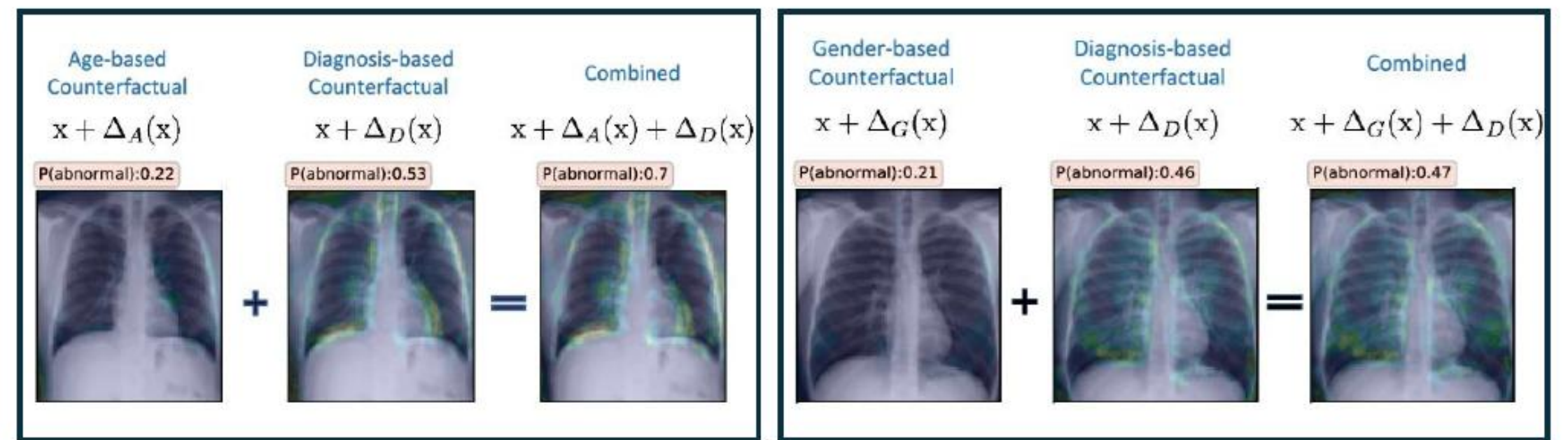
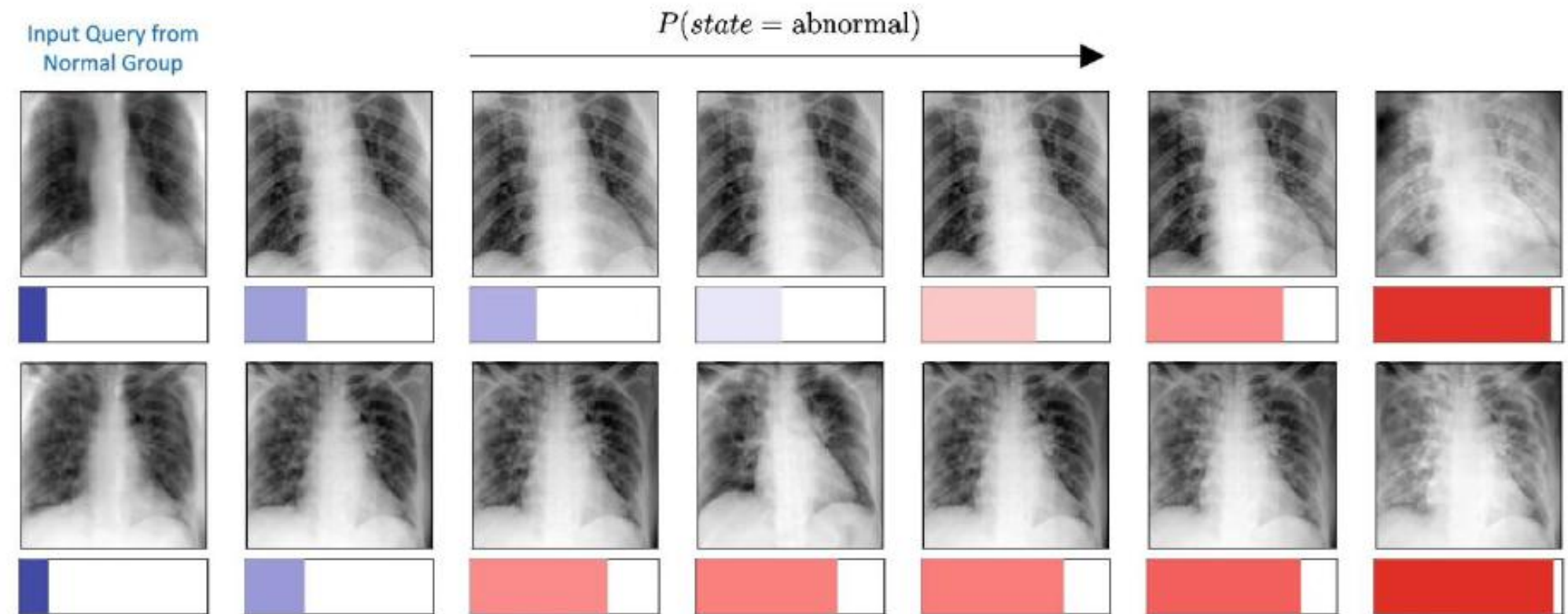


Explicabilidad a través de la comparación

Genera imágenes sintéticas modificadas (“contrafactuales”) que muestran visualmente qué cambios harían que el modelo cambiara su predicción.

- Diagnóstico
- Género
- ...

Por ejemplo, si el modelo dice “*Neumonía*”, el contrafactual busca generar una imagen muy parecida, pero que el modelo clasifique como “*Normal*”.



Article | [Open access](#) | Published: 12 January 2022

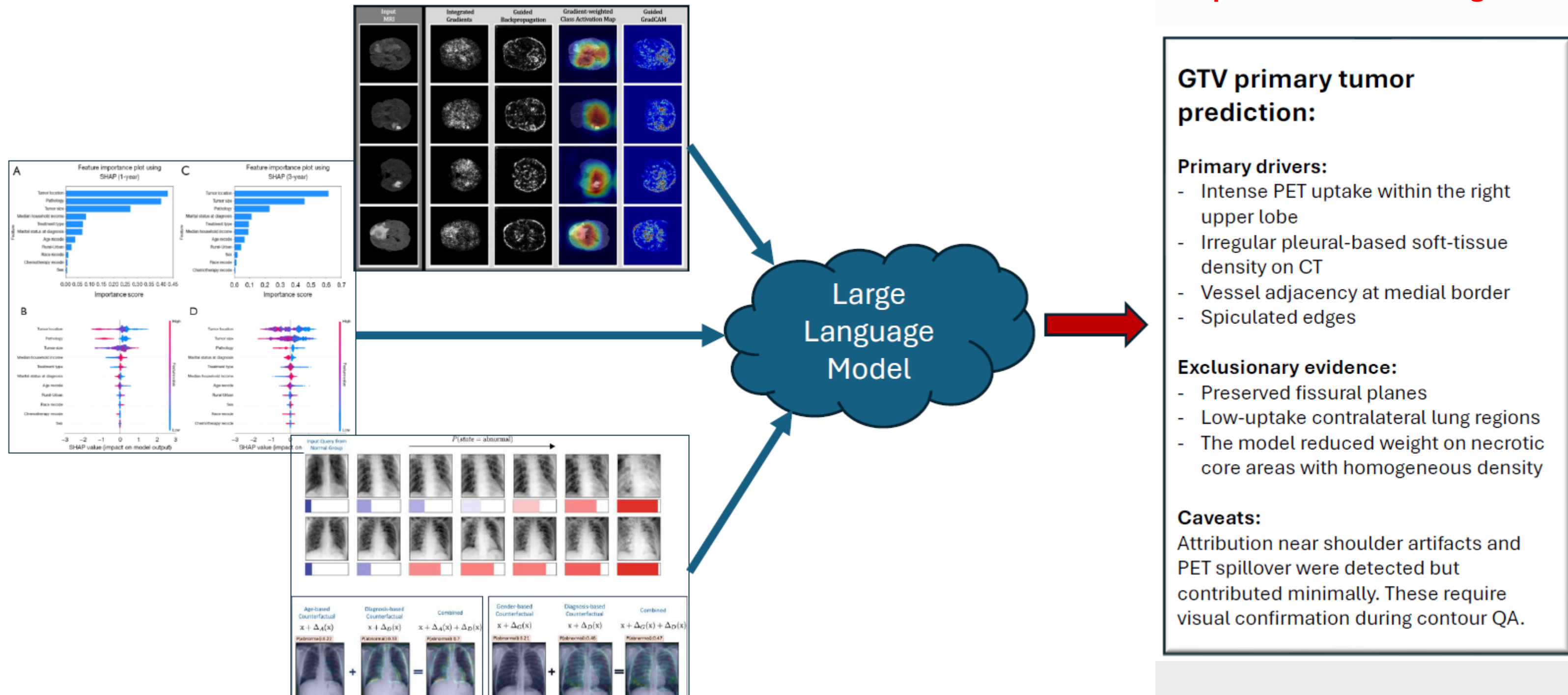
Training calibration-based counterfactual explainers for deep learning models in medical image analysis

Jayaraman J. Thiagarajan , Kowshik Thopalli, Deepta Rajan & Pavan Turaga

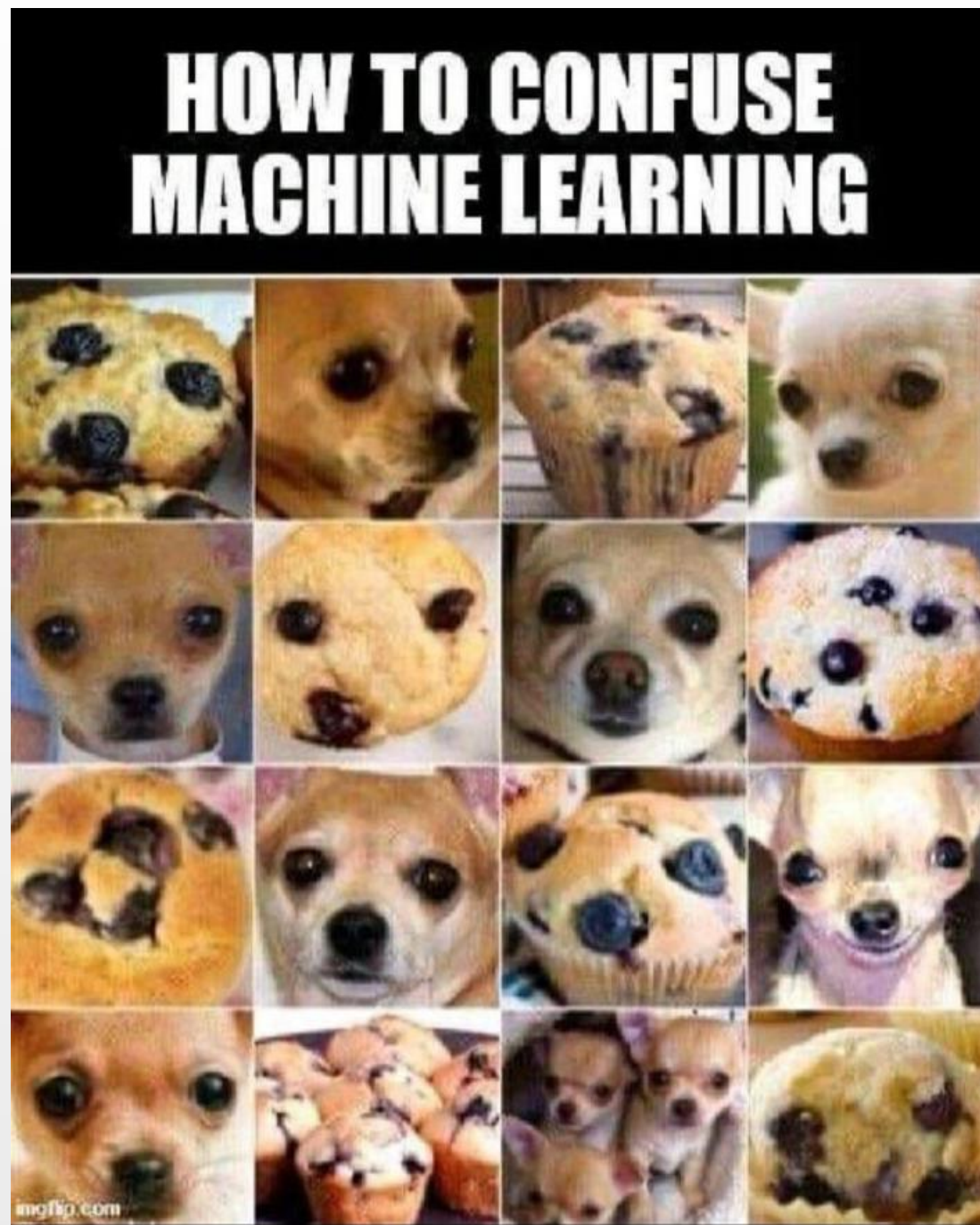
Explicabilidad a través de información textual

- Usa información textual para explicar las predicciones

¡Cuidado con los sesgos!

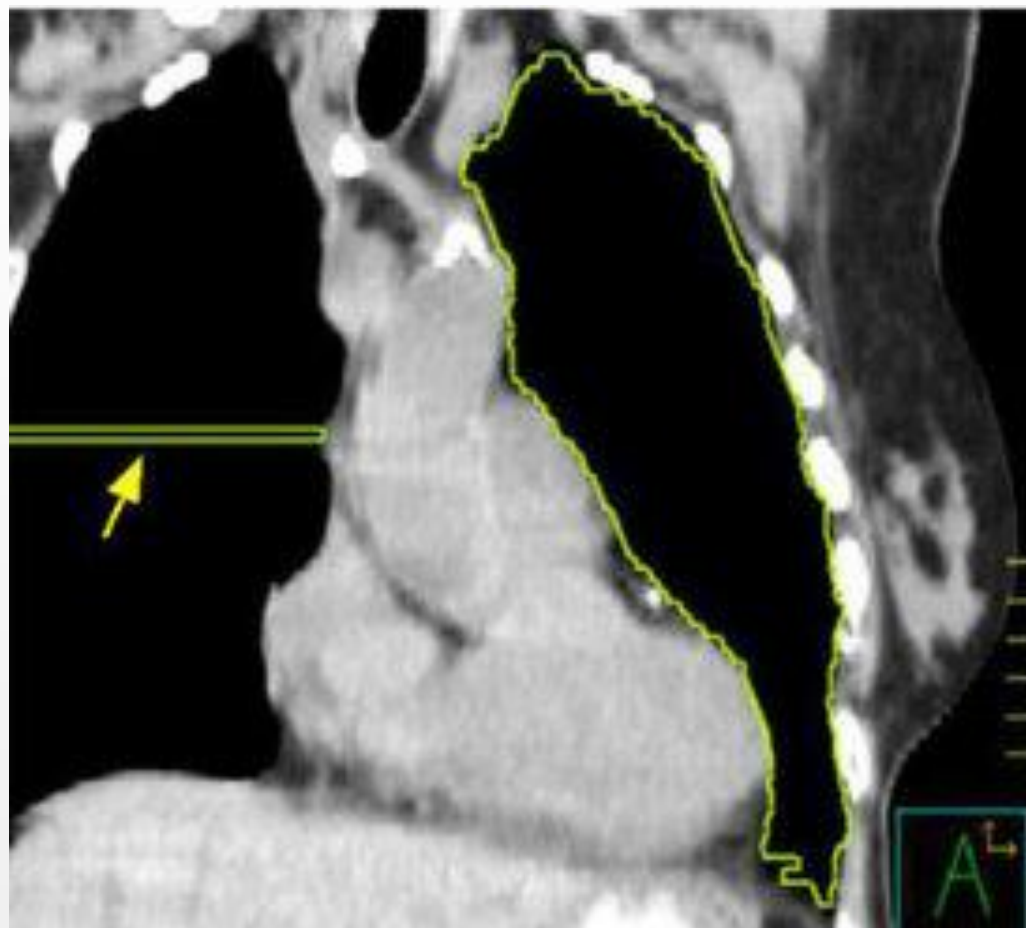


Interpretabilidad sí, pero...



- Importancia de la calidad y cantidad de datos (positivos y negativos)
- Buena en reconocimiento de patrones, pero no entiende lo que está mirando: aprende por asociación de características.
- El contexto importa.

Errores IA vs humano



HUMANO



IA

El cerebro humano es extremadamente complejo

y tiene muy poca interpretabilidad!!



Muchas de las preocupaciones realmente son sobre las interacciones humanas con las IA, en vez de la IA por si sola...

Muchas gracias por vuestra atención

