Dose
CT
Segmentations

C: Clinical variables
FC (n): Fully-connected layer with n units
NTCP: Normal Tissue Complication Probability ([0, 1])

96 × 96 × 96 @ 3
48 × 48 × 48 @ 8
24 × 24 × 24 @ 8
12 × 12 × 12 @ 16
6 × 6 × 6 @ 16
3 × 3 × 3 @ 32

FC (864)
FC (16)
C
NTCP

3D model input
Convolutional operations
Fully-connected operations
Output

FULLER LAB

Advances in AI-Based Prediction Models:
The Head and Neck Cancer Use-Case

DOI: 10.6084/m9.figshare.26817322

THE UNIVERSITY OF TEXAS
MD Anderson Cancer Center
Making Cancer History®

cdfuller@mdanderson.org

# Funding Acknowledgment/Disclosures

FULLERLAB

David Rosenthal,
MD
Professor/Section
Chief

Bill Morrison,
MD
Professor

Adam Garden,
MD
Professor

Steven Frank,
MD
Professor

Brandon Gunn
MD
Professor

Dave Fuller,
MD, PhD
Professor

THE UNIVERSITY OF TEXAS
MD Anderson
~~Cancer~~ Center

Making Cancer History®

Radiation
Oncology
Head and
Neck Section

Jack Phan,
MD, PhD
Assoc. Professor

Mike Spiotto
MD, PhD
Assoc. Professor

Jay Reddy
MD
Asst. Professor

Amy Moreno
MD,
Asst. Professor

Anna Lee
MD, MPH
Asst. Professor

# MDACC Head and Neck Team



Head and Neck Surgery

Thoracic/Head and Neck Medical Oncology

Neuroradiology

Radiation Oncology/Medical Physics

Pathology

Oncologic Dentistry

THE UNIVERSITY OF TEXAS
MD Anderson Cancer Center

Making Cancer History®

# SMART-ACT: Spatial Methodologic Approaches for Risk Assessment and Therapeutic Adaptation in Cancer Treatment

Liz Marai, PhD
Computer Science, UIC

Guadalupe Canahuate, PhD
Computer Science, UIowa

Xinhua Zhang, PhD
Computer Science, UIC

Dave Fuller, MD, PhD
Radiation Oncology
MDACC

# MD Anderson Multi-disciplinary Symptom Working Group



**Stephen Lai**
MD, PhD
Head
and Neck Surgery

**Kate Hutcheson**
PhD
Speech Pathology

**Amy Moreno, MD**
Radiation Oncology

**Abdallah Mohamed**
MD, MSc
Radiation Oncology

**Jihong Wang**
PhD
Radiation Oncology

**Dave Fuller**
MD, PhD
Radiation Oncology

John Christodouleas
MD
Elekta AB/UPenn

Dave Fuller
MD, PhD
MDACC

NIH Academic Industrial Partnership
(R01 DE028290-01)

Elekta

THE UNIVERSITY OF TEXAS
MD Anderson Cancer Center
Making Cancer History®

PHILIPS

Future state: Multi-sequence adaption

Inter-sequence registration & Auto-segmentation

Inter and intra-sequence registration & Contour propagation

ADC modified Plan re-optimization

Figure 8: Diagram illustrating the conceptualized Elekta software development pipeline.

# NSF-NIH Smart-Connected Health Program
# Rice-MDACC Operations Research in Oncology

Andrew Schaefer, PhD
Computational
and Applied Math
RiceU

Dave Fuller, MD, PhD
Radiation Oncology
MDACC

Figure 6. Graphical abstract of Aim 3, showing dose/NTCP and image implementation for adaptive radiotherapy decision support.

# Image Guided Cancer Therapy Program

## T32CA261856



Kristy Brock, PhD
**Director,**
**IGCTR Program**



Stephen Lai,
MD, PhD



Dave Fuller,
MD, PhD

Physicians

Post-docs

FULLER LAB

Computational Scientists

Lab Manager

Grad students

For radiation oncologists, *spatial* dose/response data is what separates us from other cancer paradigms

1954 PhD Berkeley (mathematics)

1960 -1967 UCLA (mathematics)

1969 -1982 Consultant

1982 - 1993 Berkeley (statistics)

1984 "Classification & Regression Trees"
    (with Friedman, Olshen, Stone)

1996 "Bagging"

2001 "Random Forests"

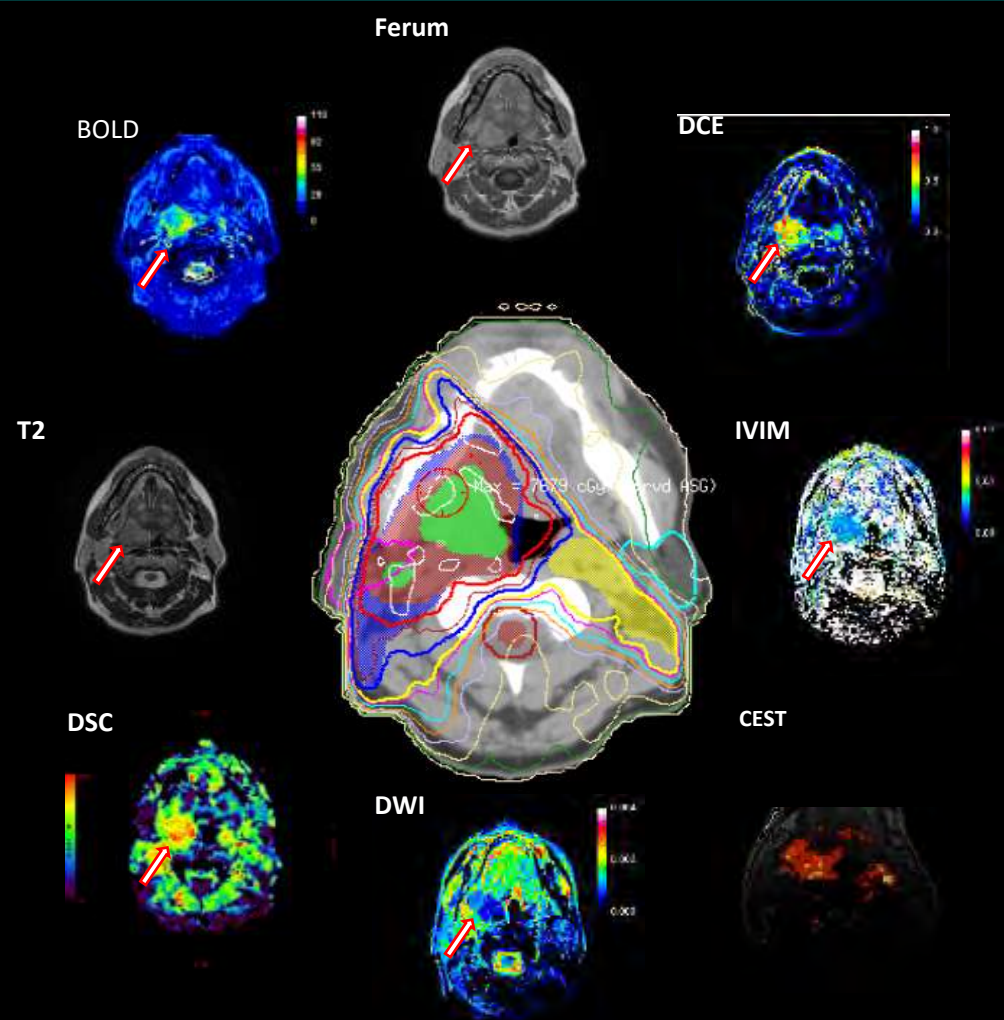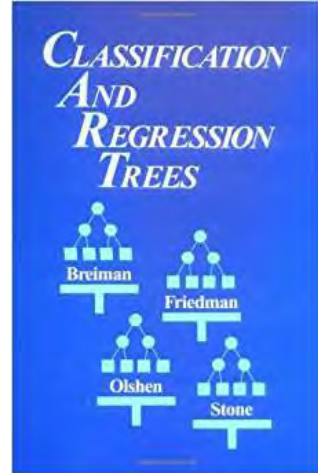# RECURSIVE PARTITIONING ANALYSIS (RPA) OF PROGNOSTIC FACTORS IN THREE RADIATION THERAPY ONCOLOGY GROUP (RTOG) BRAIN METASTASES TRIALS

LAURIE GASPAR, M.D.,* CHARLES SCOTT, M.S.,[†] MARVIN ROTMAN, M.D.,[‡]
SUCHA ASBELL, M.D.,[§] THEODORE PHILLIPS, M.D.,[¶] TODD WASSERMAN, M.D.,[#]
W. GILLIES MCKENNA, M.D., Ph.D.** AND ROGER BYHARDT, M.D.[††]

**(1) RTOG 79-16**

| | | |
|---|---|---|
| RANDOMIZE | 30 Gy | 10 fx / 10 days / 2 wks |
| | 30 Gy | 10 fx / 10 days / 2 wks + MISO |
| | 30 Gy | 6 fx / 6 days / 3 wks |
| | 30 Gy | 6 fx / 6 days / 3 wks + MISO |

**(3) RTOG 89-05**

| | | |
|---|---|---|
| RANDOMIZE | 37.5 Gy | 15 fx  3 wks |
| | 37.5 Gy | 15 fx  3 wks |

BUdR continuous 96 hr infusion,
0.8 g/m2/d starting 3-4 days
prior to weeks 1 and 2 of radiation

**(2) RTOG 85-28**

| | | |
|---|---|---|
| RANDOMIZE (DOSE ESCALATION) | 48 Gy | 1.6 Gy BID / 30 fx / 15 days |
| | 54.4 Gy | 1.6 Gy BID / 34 fx / 17 days |
| | 54.4 Gy | 1.6 Gy BID / 34 fx / 17 days |
| | 64 Gy | 1.6 Gy BID / 40 fx / 20 days |
| | 64 Gy | 1.6 Gy BID / 40 fx / 20 days |
| | 70.4 Gy | 1.6 Gy BID / 44 fx / 22 days |

Fig. 1. Protocol schemas.

## Recursive Tree



**Root** n=1200

**KPS ≥ 70** n=1011

**KPS < 70** n=175 → Class III

**Primary Controlled** n=588

**Primary uncontrolled** n=413 → Class II

**Age < 65 yrs** n=424

**Age ≥ 65 yrs** n=164 → Class II

**Metastases-Brain Only** n=236 → Class I

**Metastases-Brain & Other Sites** n=188 → Class II

Fig. 2. Recursive tree.

## BRAIN METASTASES SURVIVAL
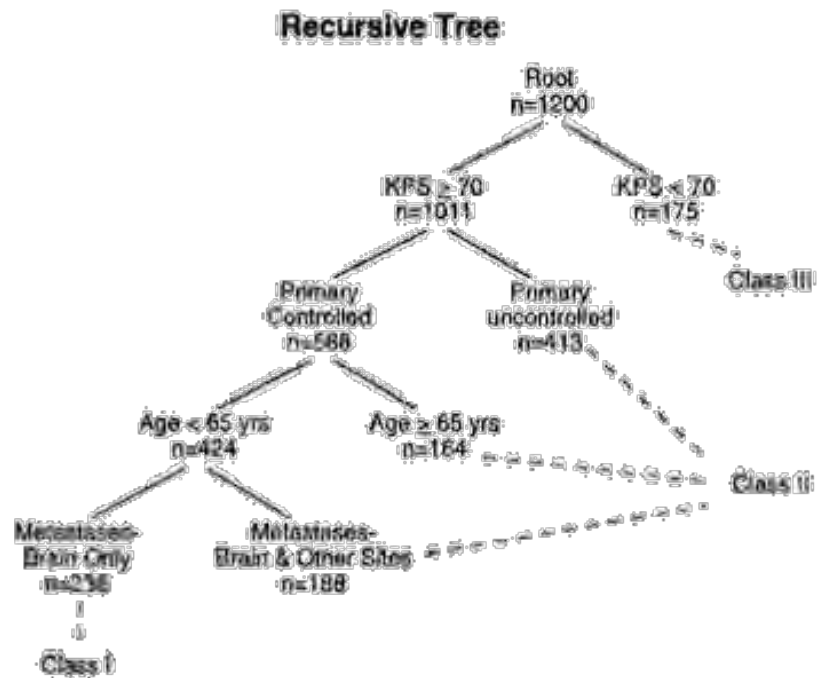


PERCENT ALIVE

MONTHS FROM ONSTUDY
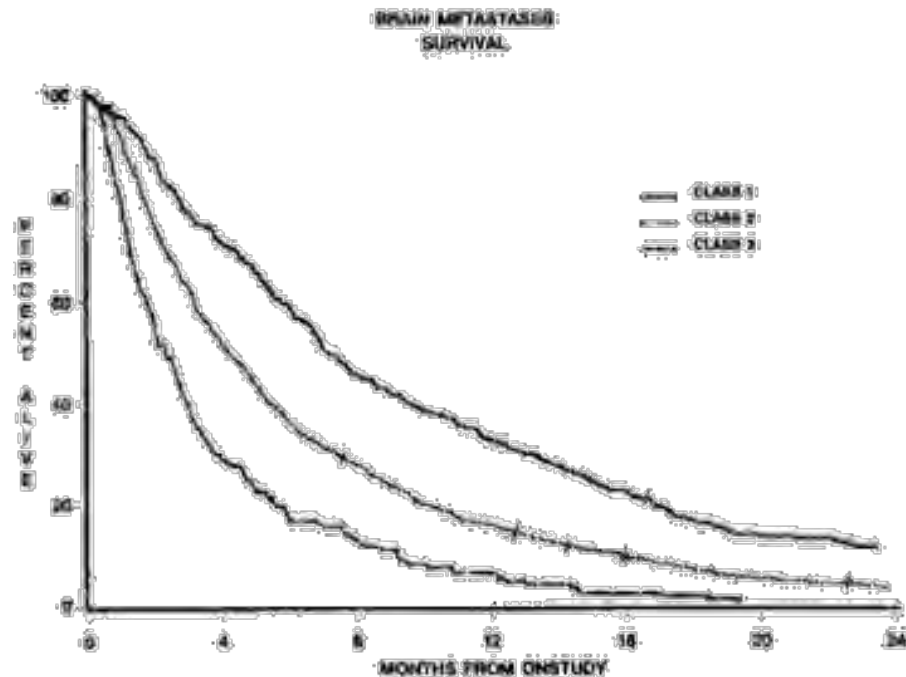
CLASS 1
CLASS 2
CLASS 3

Fig. 3. Survival curves for Class I, II, III.

# Statistical Modeling: The Two Cultures

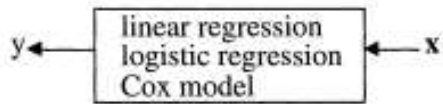**Leo Breiman**

Hypothesis Testers

AI Modelers

## The Data Modeling Culture

The analysis in this culture starts with assuming a stochastic data model for the inside of the black box. For example, a common data model is that data are generated by independent draws from

response variables = $f$(predictor variables, random noise, parameters)

The values of the parameters are estimated from the data and the model then used for information and/or prediction. Thus the black box is filled in like this:

y ← | linear regression / logistic regression / Cox model | ← x

*Model validation.* Yes–no using goodness-of-fit tests and residual examination.
*Estimated culture population.* 98% of all statisticians.

nature

goals in analy
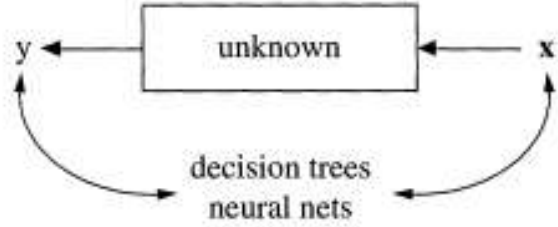able to predic
future input
extract som
ssociating the
ables.
different appr

## The Algorithmic Modeling Culture

The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$—an algorithm that operates on $\mathbf{x}$ to predict the responses $\mathbf{y}$. Their black box looks like this:

y ← | unknown | ← x

decision trees
neural nets

*Model validation.* Measured by predictive accuracy.
*Estimated culture population.* 2% of statisticians, many in other fields.

# Richard Bellman: "The curse of dimensionality"



Bellman, first editor of Mathematical Biosciences, was working in dynamic optimization

-Referred initially to issues that arise in higher-order analyses that are hard for humans to conceptualize as we move increase dimensions or add time-varying components

-Broadly, refers to typical increase in sparsity of data in high-dimensions and information reduction through dimensional summarization.
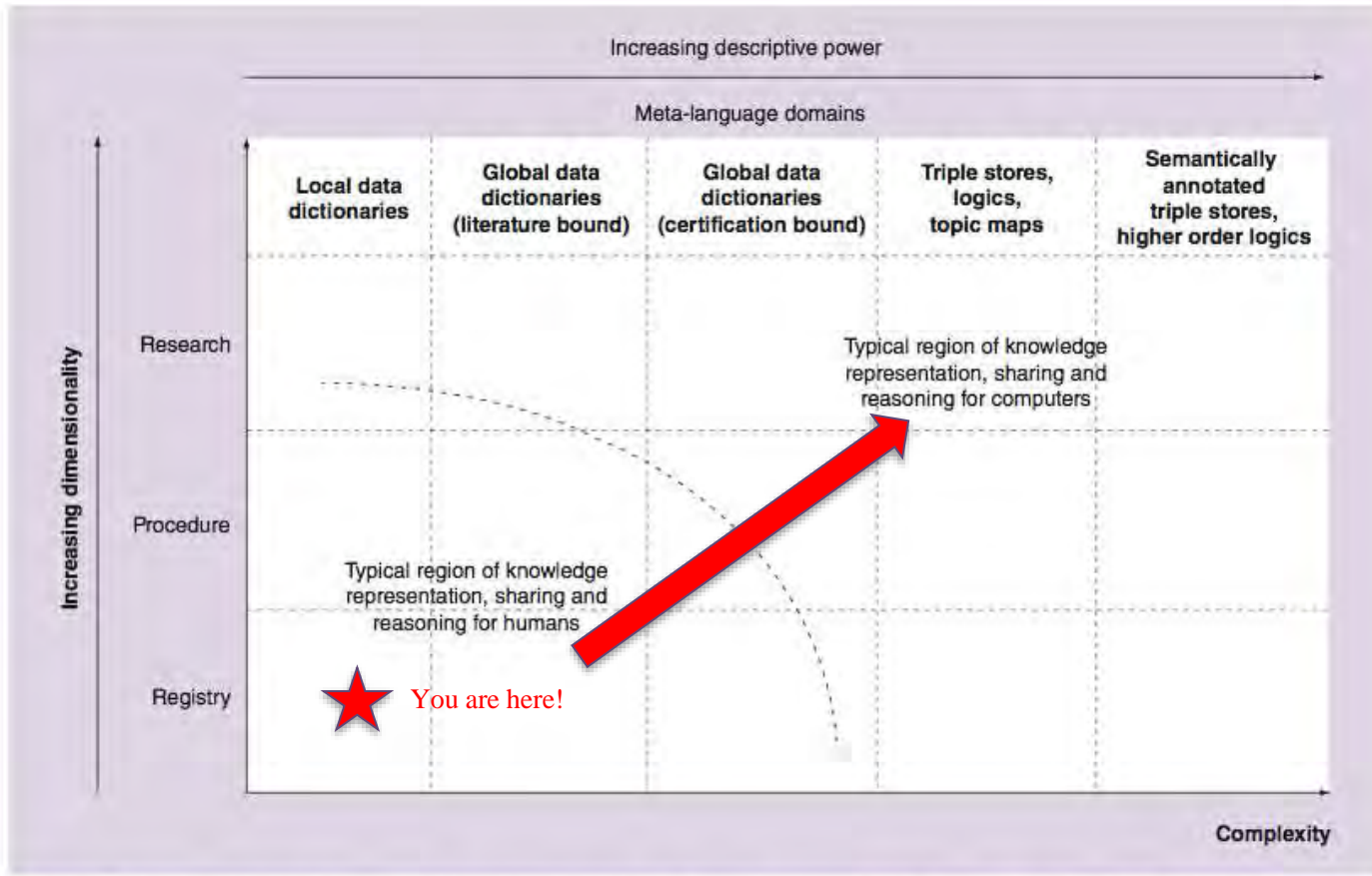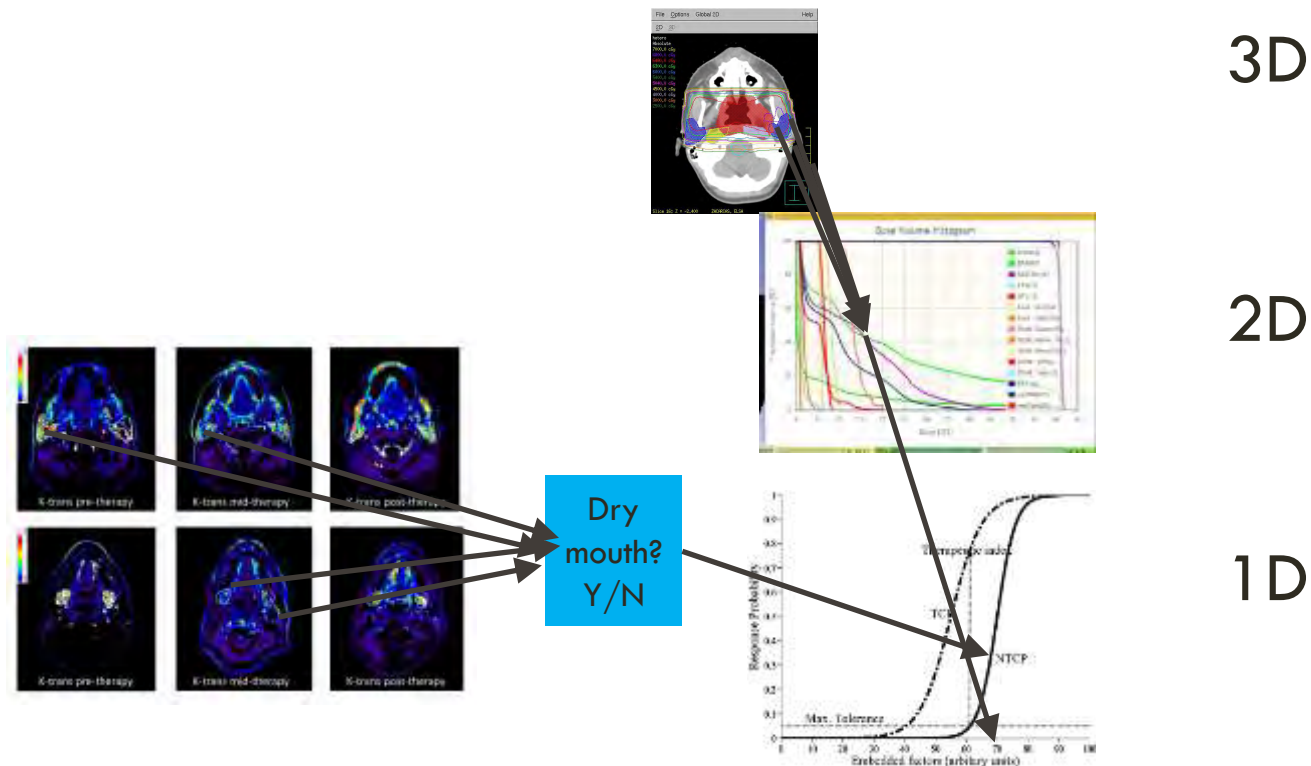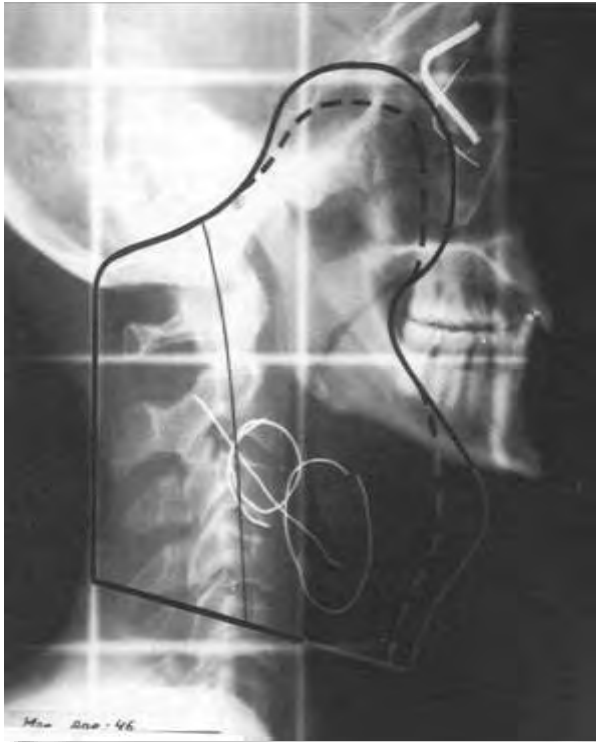
**Figure 1. Possible evolution in knowledge representation, seen from the perspective of computer science, under a qualitative point of view.**

at MD Anderson

# Example:
# Information loss through summarization by dimensionality reduction
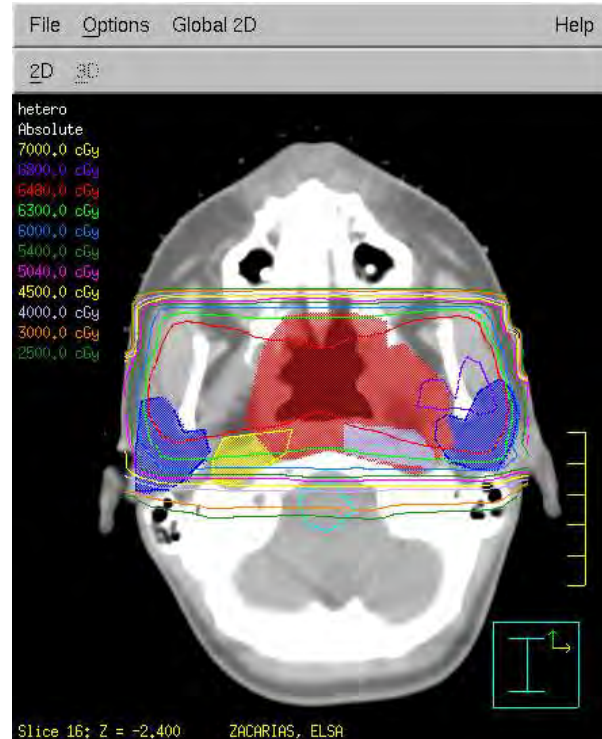


3D

2D

Dry
mouth?
Y/N

1D

# What NTCP models were built for…

**1990**

**2000**

# PAROTID GLAND FUNCTION AFTER RADIOTHERAPY: THE COMBINED MICHIGAN AND UTRECHT EXPERIENCE

Tim Dijkema, M.D.,* Cornelis P. J. Raaijmakers, Ph.D.,* Randall K. Ten Haken, Ph.D.,†
Judith M. Roesink, M.D., Ph.D.,* Pètra M. Braam, M.D., Ph.D.,* Anette C. Houweling, M.Sc.,*
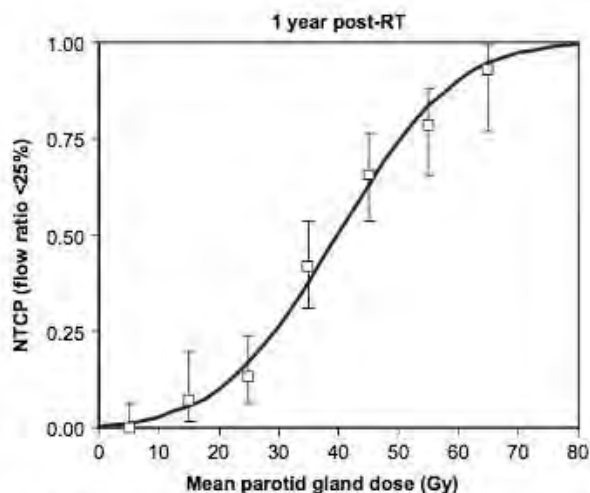Marinus A. Moerland, Ph.D.,* Avraham Eisbruch, M.D.,† and Chris H. J. Terhaard, M.D. Ph.D.*

Fig. 3. Combined Michigan and Utrecht normal tissue complication probability (NTCP) curve as a function of the mean parotid gland dose. Clinical NTCP values (using mean dose bins of 10 Gy) are shown, including 95% confidence intervals. RT = radiotherapy.
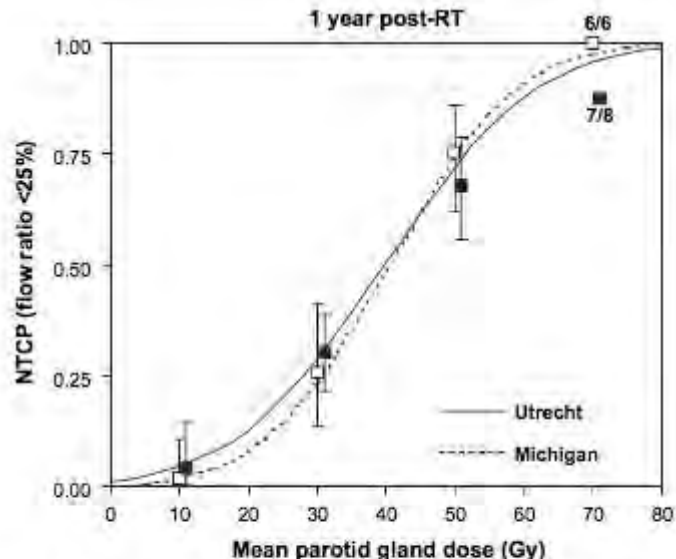


Fig. 2. Normal tissue complication probability (NTCP) curves as a function of the mean parotid gland dose for Michigan (dashed line) and Utrecht (solid line). Clinical NTCP values (using mean dose bins of 20 Gy) are shown for Michigan (open squares) and Utrecht (black squares), including 95% confidence intervals. RT = radiotherapy.
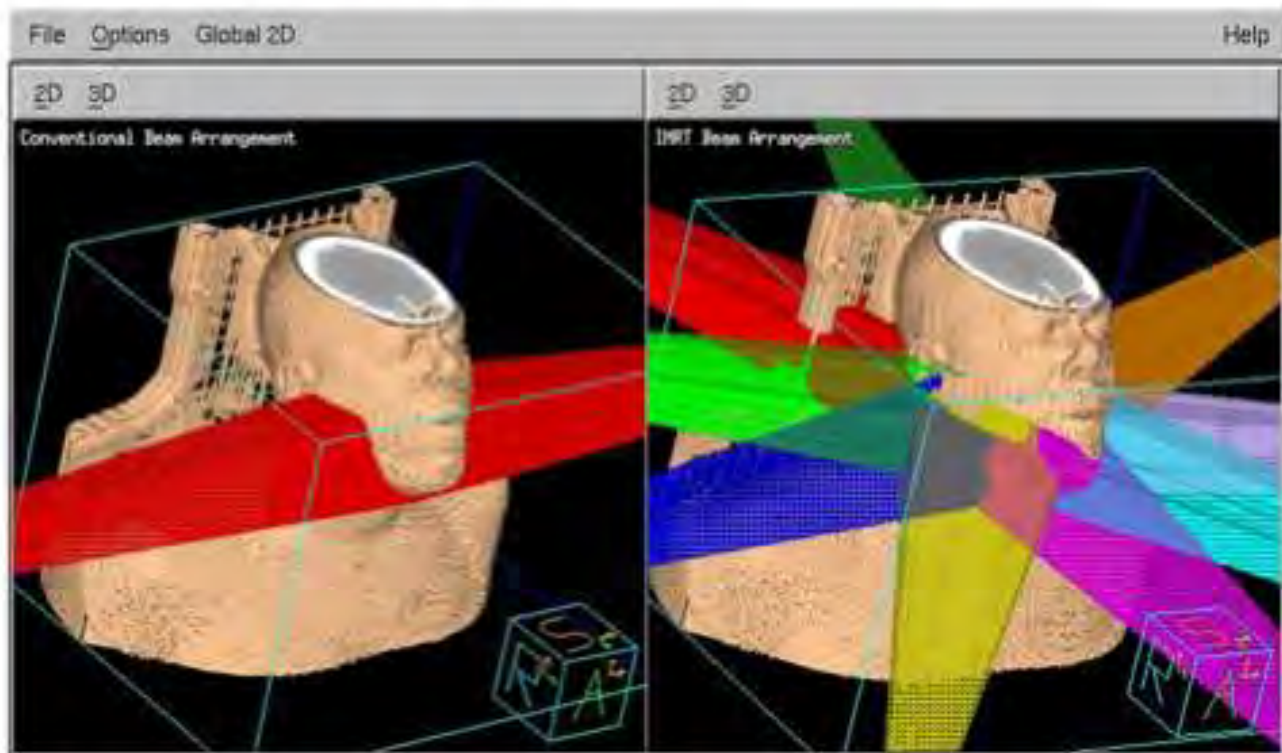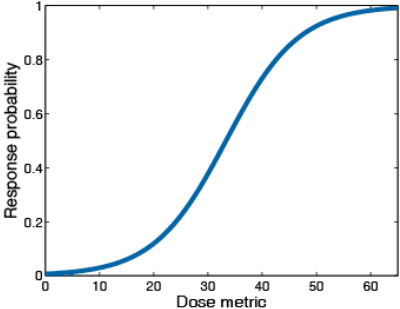
Fig. 1. Comparison of nontarget beam paths in intensity-modulated radiotherapy (top) vs. conventional three-dimensional technique (bottom).

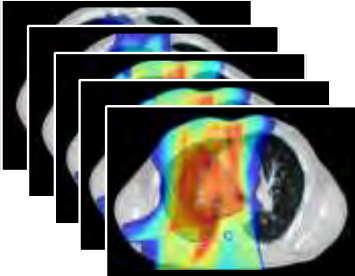# Standard phenomenological modelling methodology
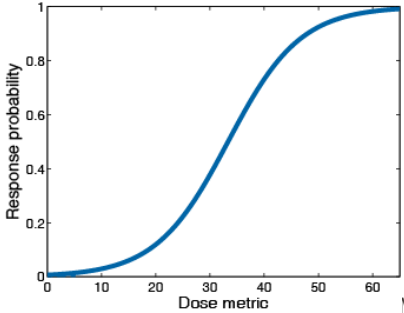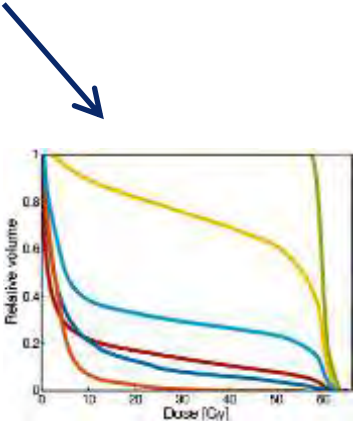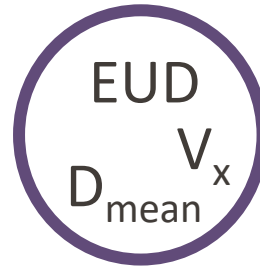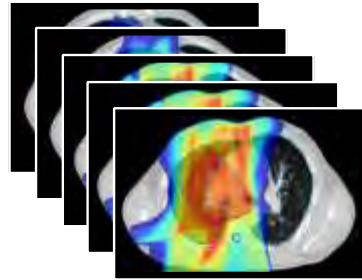


"Dose"

Generalised linear modelling

MD Anderson

# Standard phenomenological modelling methodology



"Dose"

Generalised linear modelling

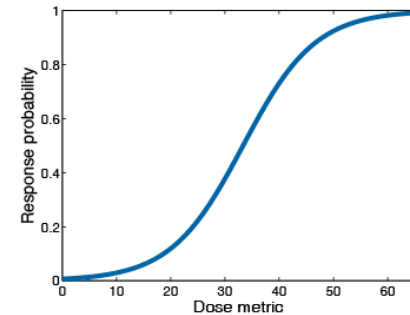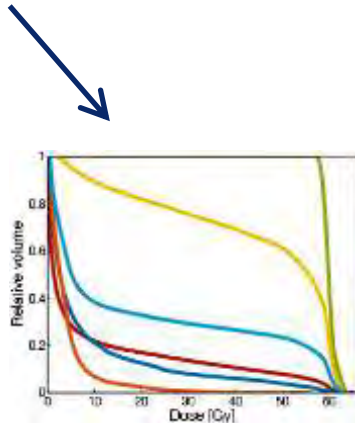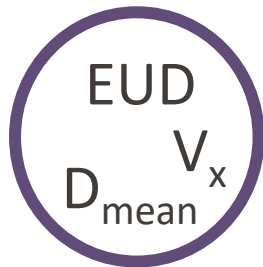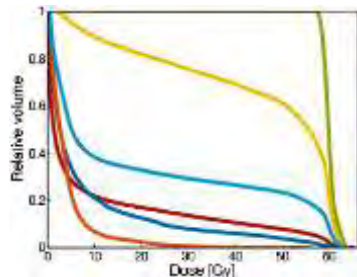MD Anderson

# Standard phenomenological modelling methodology



EUD
$V_x$
$D_{mean}$

Generalised linear modelling

MD Anderson

# Standard phenomenological modelling methodology



EUD
$V_x$
$D_{mean}$

Reduce DVH to one (or a limited number of) dose metrics

$$EUD = \left( \sum_k d_k^a \frac{v_k}{V_{tot}} \right)^{1/a}$$

$$V_x = \sum_k E(d_k) v_k \qquad E(d_k) = \begin{cases} 0 \; for \, d_k < x \, Gy \\ 1 \; for \, d_k \geq x \, Gy \end{cases}$$

Find the dose representation that best correlates with toxicity

# Potential problems with the standard dose-reduction approach
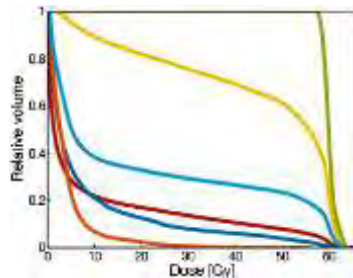


Reduce dose distribution to DVH
- Removes all spatial information
- Assumes equal sensitivity/response of all parts of OAR

Alternatives:
- Divide into anatomical substructures
- Dose surface histograms
- Consider (and/or explicitly model) local response on voxel-to-voxel basis

# Adding spatial information to (N)TCP models – general strategies



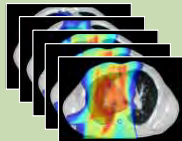| | Dose variable(s) | Response measure |
|---|---|---|
| Traditional NTCP modelling | One per patient | One per patient |
| Voxel-based analysis (VBA), convolutional neural networks (CNN) | Many per patient (2D or 3D data) | One per patient |
| Image-based response models | Many per patient (2D or 3D data) | Many per patient (2D or 3D data) |

Palma et al. Cancers 2021;13(14):3553. Palma et al. Phys Med 2020;69:192-204. Appelt et al. Clin Oncol 2022;34(2):e87-e96

# Adding spatial information to (N)TCP models



VBA

$p_1(O|D_{x1})$
$p_2(O|D_{x2})$
$p_3(O|D_{x3})$
$p_4(O|D_{x4})$

p-value map

Single patient-level prediction

$p(O|D)$

# New anatomical insights from voxel-based analysis of dose?



Image Based Data Mining Using Per-voxel Cox Regression

Original Research

Radiation dose to heart base linked with poorer survival in lung cancer patients

Generally for VBA based studies:

- How dependent are the results by structures in the dose data (e.g. dose gradients and robustness of planned relative to delivered dose)?
- Issues with statistical analysis in some parts of the published literature
  - Shortall et al. Flogging a Dead Salmon? IJROBP 2021
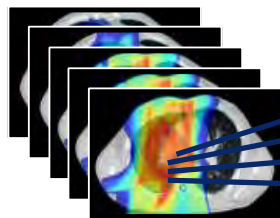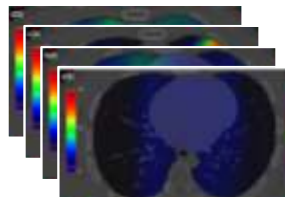
# Adding spatial information to (N)TCP models



VBA

$p_1(O|D_{x1})$
$p_2(O|D_{x2})$
$p_3(O|D_{x3})$
$p_4(O|D_{x4})$

p-value map

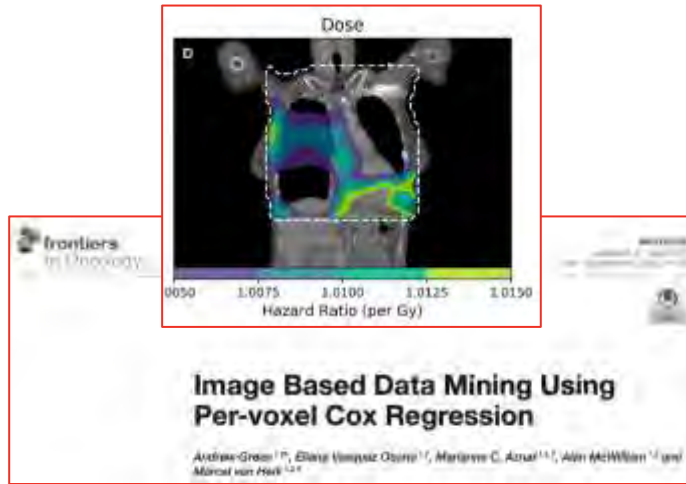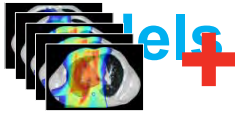Single patient-level prediction

$p(O|D)$

CNN

single patient-level prediction / classification

$p(O|D)$

Saliency map?
Sensitivity map?

Palma et al. Cancers 2021;13(14):3553. Palma et al. Phys Med 2020;69:192-204. Appelt et al. Clin Oncol 2022;34(2):e87-e96

# Improved toxicity prediction with voxel-based analysis?

| Patient number | Cancer site | Ref | Improvement over GLM | External validation |
|---|---|---|---|---|
| 42 | Cervical | Zhen 2017 | ✚ | ✖ |
| 125 | Liver | Ibragimov 2018 | ✚ | ✖ |
| 784 | Head and neck | Men 2019 | ✚ | ✖ |
| 120 | Liver | Ibragimov 2019 | ✚ | ✖ |
| 122 | Liver | Ibragimov 2020 | - | ✖ |
| 160 | Oropharyngeal | Welch 2020 | ✖ | ✖ |
| 70 | NSCLC | Liang 2019 | ✚ | ✖ |
| 66 | Oropharyngeal | Wang 2020 | - | ✖ |
| 52 | Post-prostatectomy | Yang 2021 | - | ✖ |
| 217 | Thoracic | Liang 2021 | ✚ | ✖ |

**Appelt et al. Deep Learning for Radiotherapy Outcome Prediction Using Dose Data - A Review. Clin Oncol 2022**

# Adding spatial information to (N)TCP models – general strategies

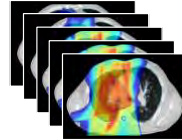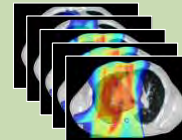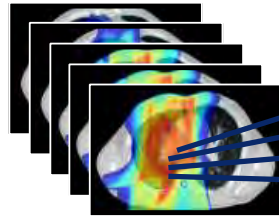| | **Dose variable(s)** | | **Response measure** | |
|---|---|---|---|---|
| Traditional NTCP modelling |  | One per patient |  | One per patient |
| Voxel-based analysis (VBA), convolutional neural networks (CNN) |  | Many per patient (2D or 3D data) |  | One per patient |
| Image-based response models |  | Many per patient (2D or 3D data) |  | Many per patient (2D or 3D data) |

Palma et al. Cancers 2021;13(14):3553. Palma et al. Phys Med 2020;69:192-204. Appelt et al. Clin Oncol 2022;34(2):e87-e96

# Adding spatial information to (N)TCP models

Image-based response models

single model linking dose &
local response

$$p(R_{x1}|D_{x1})$$
$$p(R_{x2}|D_{x2})$$
$$p(R_{x3}|D_{x3})$$
$$p(R_{x4}|D_{x4})$$

Multilevel mixed effect model

$$p(O|D)$$

# Better or novel biological insights from voxel-based analysis of dose?

Dose

LET

Imaging changes



**+**

Model of RBE dependence on dose and LET

A systematic review of clinical studies on proton Relative Biological Effectiveness (RBE)
- 13 studies using voxel-wise analyses of patient effects versus dose and LET
  - **3/13: No effect of LET on RBE**
  - **6/16: Maybe effect of LET on RBE**
  - **4/13: Effect of LET on RBE**
- Significant methodological modelling issues
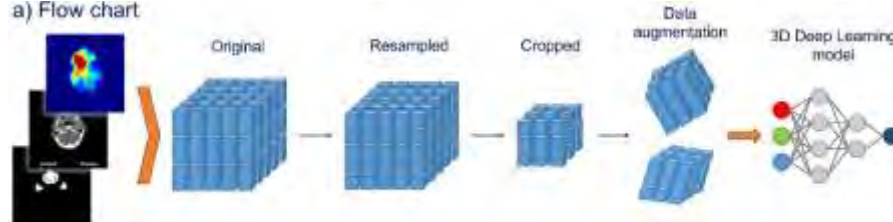  - E.g. no consideration of nested / multi-level data

Underwood et al. A systematic review of clinical studies on variable proton Relative Biological Effectiveness (RBE). Radiother Oncol. 2022

# 3D deep learning Normal Tissue Complication Probability model to predict late xerostomia in head and neck cancer patients
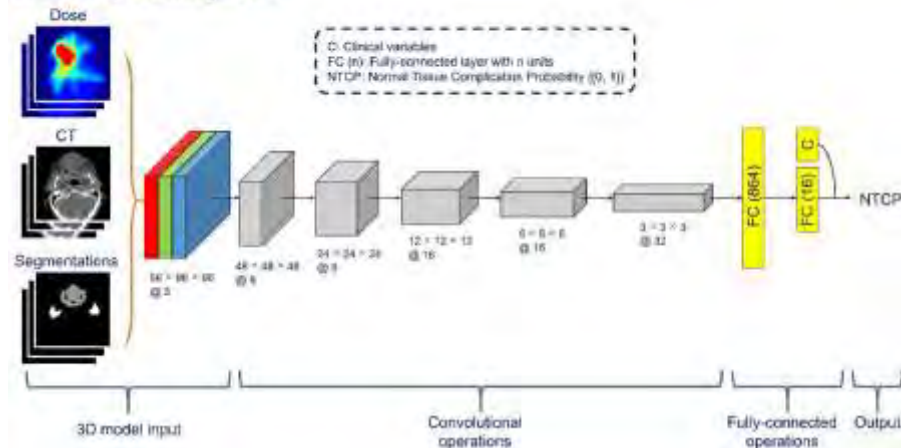
Hung Chu MSc [1] ✉, Suzanne P.M. de Vette MSc [1], Hendrike Neh MSc [1], Nanna M. Sijtsema PhD [1], Roel J.H.M. Steenbakkers MD, PhD [1], Amy Moreno MD [2], Johannes A. Langendijk MD, PhD [1], Peter M.A. van Ooijen PhD [1], Clifton D. Fuller MD, PhD [2], Lisanne V. van Dijk PhD [1] ✉
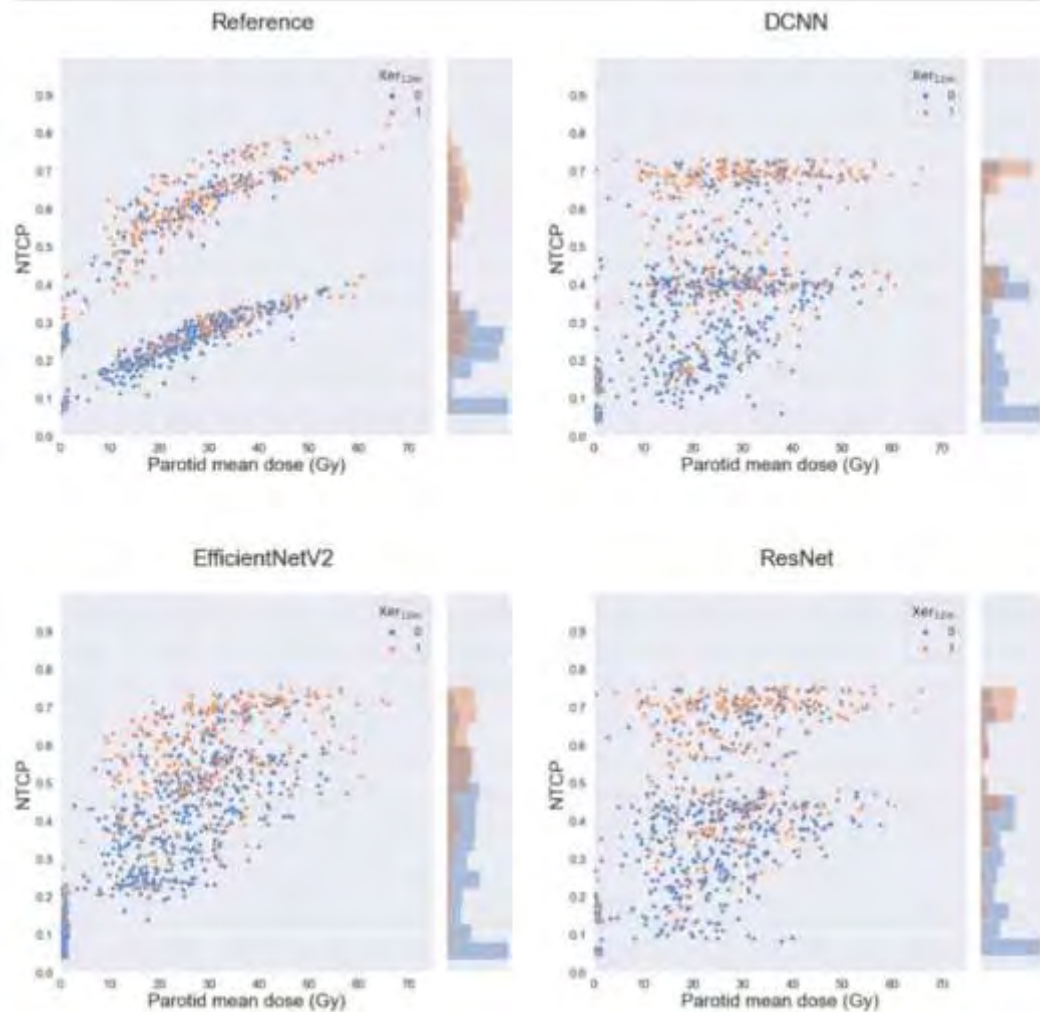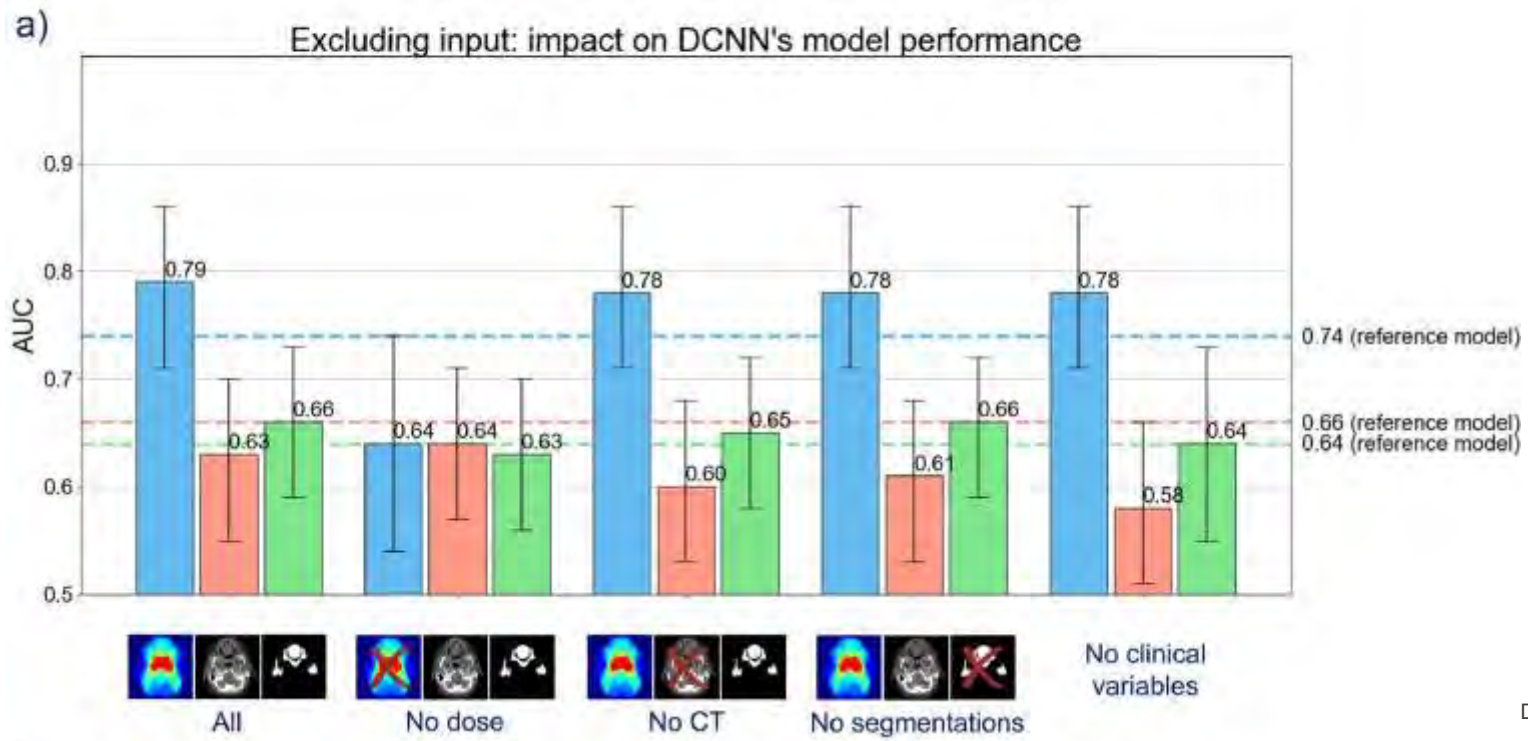
Figure 2. These scatterplots display the relationship between parotid mean dose (in Gy) and NTCP (Normal Tissue Complication Probability) value for all models. Patients who experienced moderate-to-severe xerostomia 12 months post-radiotherapy are represented by orange, while the remaining patients are represented by blue. The accompanying histogram illustrates the distribution of the NTCP values.

| | Training | Independent test | External validation |
|---|---|---|---|
| Total | 759 | 138 | 311 |

Legend:
- ■ Independent test
- ■ External validation
- ■ External validation (transfer learning)

a)

Excluding input: impact on DCNN's model performance



AUC values:

- All: 0.79 (Independent test), 0.63 (External validation), 0.66 (transfer learning)
- No dose: 0.64, 0.64, 0.63
- No CT: 0.78, 0.60, 0.65
- No segmentations: 0.78, 0.61, 0.66
- No clinical variables: 0.78, 0.58, 0.64

Reference lines:
- 0.74 (reference model)
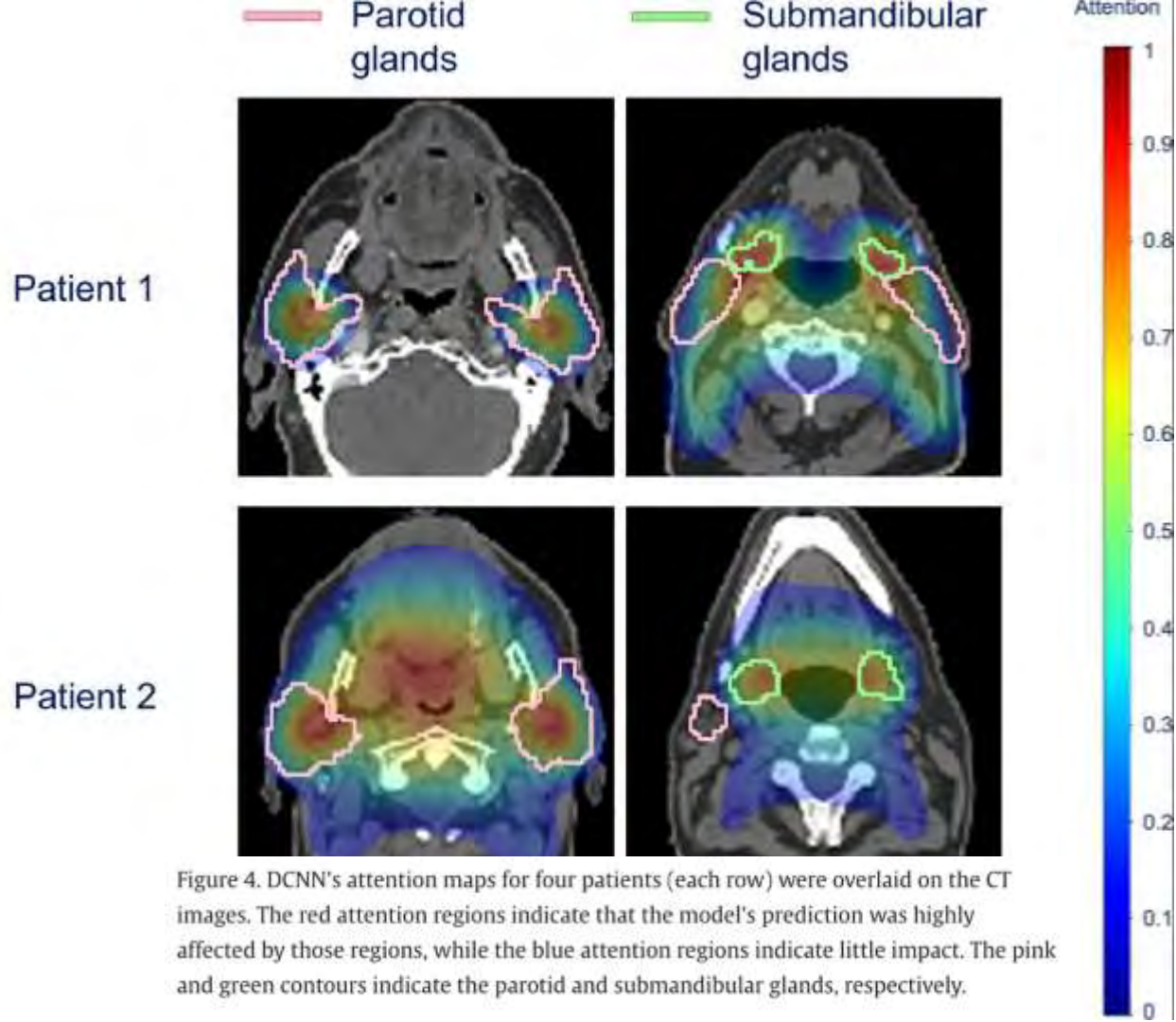- 0.66 (reference model)
- 0.64 (reference model)

Figure 4. DCNN's attention maps for four patients (each row) were overlaid on the CT images. The red attention regions indicate that the model's prediction was highly affected by those regions, while the blue attention regions indicate little impact. The pink and green contours indicate the parotid and submandibular glands, respectively.

# Beyond mean pharyngeal constrictor dose for beam path toxicity in non-target swallowing muscles: Dose–volume correlates of chronic radiation-associated dysphagia (RAD) after oropharyngeal intensity modulated radiotherapy ✶

MD Anderson Head and Neck Cancer Symptom Working Group (Contributing authors Timothy Dale[a,b,1], Katherine Hutcheson[b,1], Abdallah S.R. Mohamed[a,d], Jan S. Lewin[b], G. Brandon Gunn[a], Arvind U.K. Rao[c], Jayashree Kalpathy-Cramer[e], Steven J. Frank[a], Adam S. Garden[a], Jay A. Messer[a,d], Benjamin Warren[b], Stephen Y. Lai[b], Beth M. Beadle[a], William H. Morrison[a], Jack Phan[a], Heath Skinner[a], Neil Gross[b], Renata Ferrarotto[c], Randal S. Weber[b], David I. Rosenthal[a], Clifton D. Fuller[a,b,*])
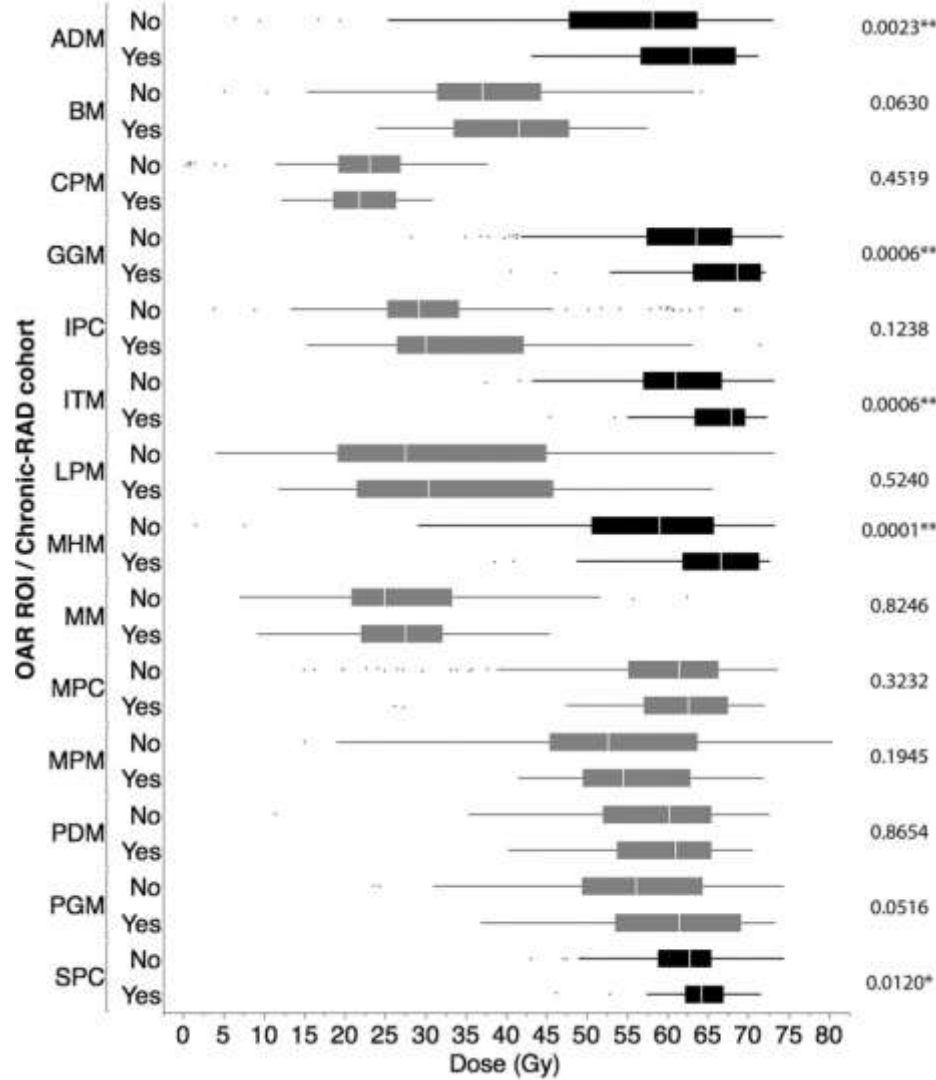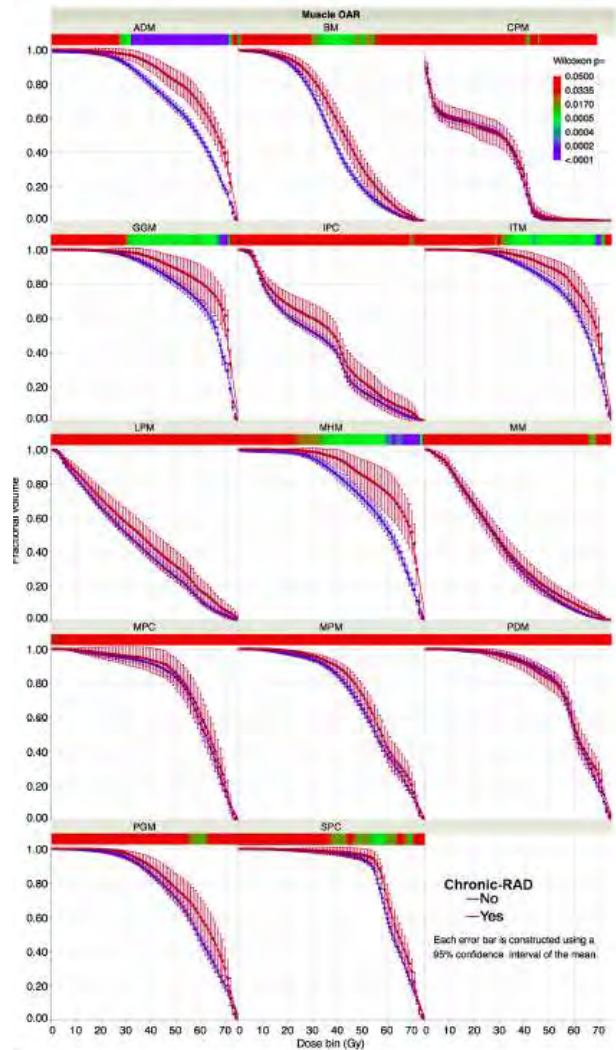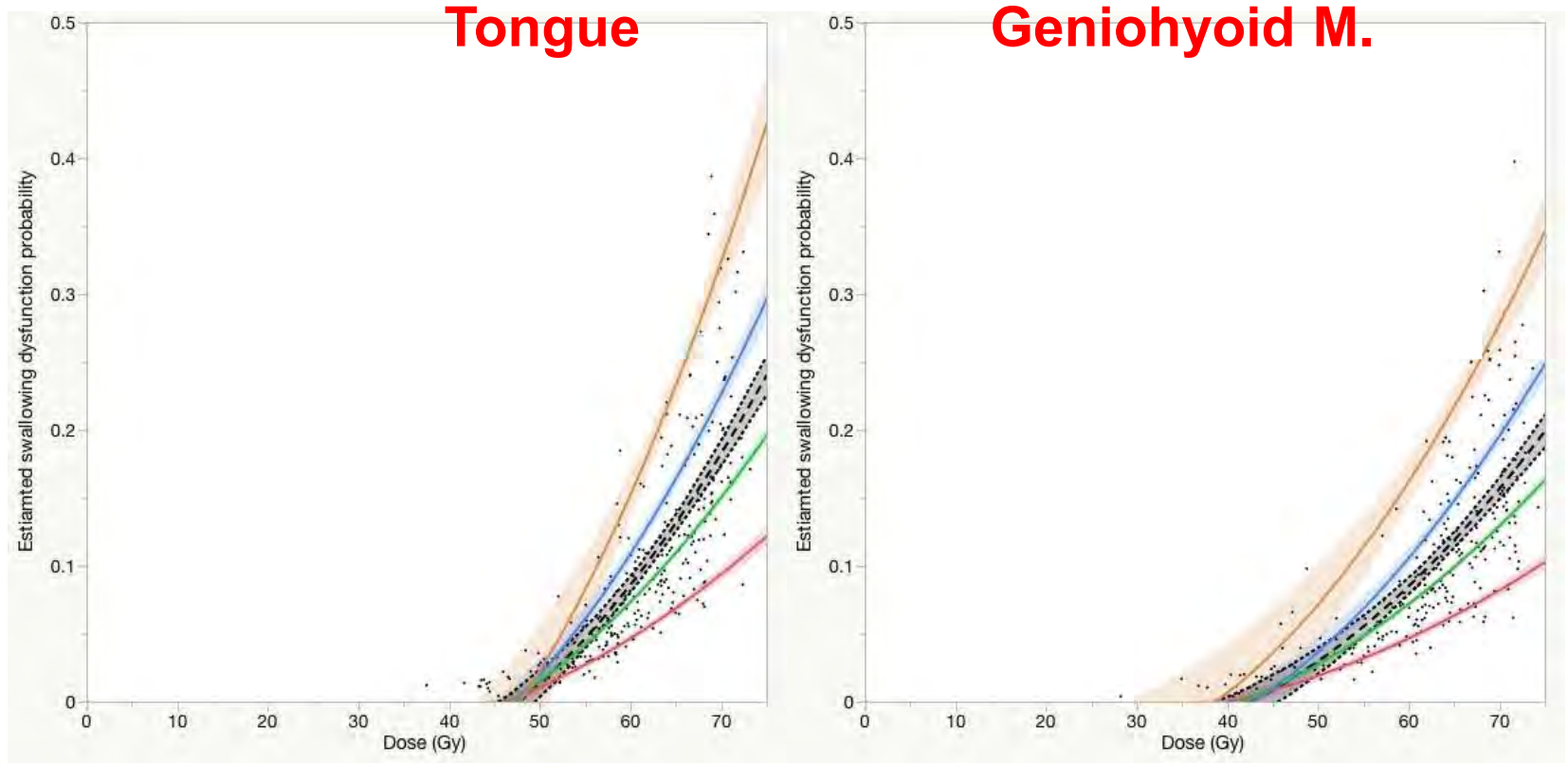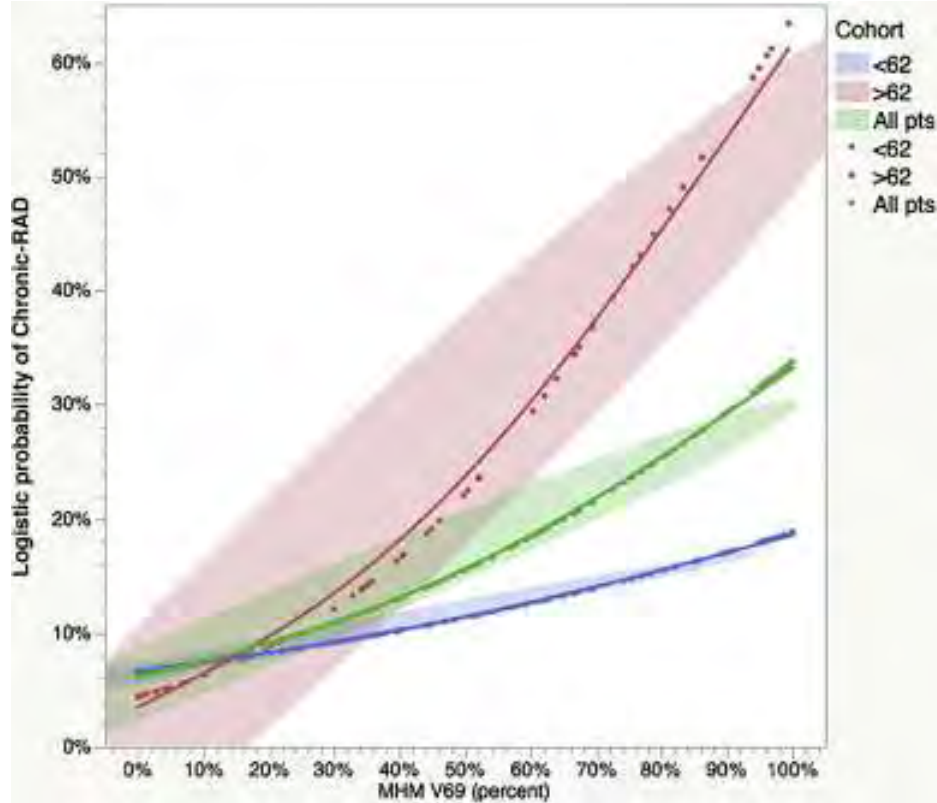


Fig. 1. Exemplar swallow-related ROI. Axial, coronal, and sagittal images of the contoured segments. *Abbreviations:* GGM – genioglossus muscle; HP – hard palate; IPC – inferior pharyngeal constrictor; ITM – intrinsic tongue muscles; LPM – lateral pterygoid muscle; MHM – mylo/geniohyoid complex; MM – masseter muscle; MPM – medial pterygoid muscle; PDM – posterior digastric muscle; SP – soft palate; SPC – superior pharyngeal constrictor, R.-right, L.-left.

# Example: Age and dysphagia

# Optimum OPC model includes mylohyoid/geniohyoid dose & age

# Adding spatial data…

Magnetic resonance imaging of swallowing-related structures in nasopharyngeal carcinoma patients receiving IMRT: Longitudinal dose–response characterization of quantitative signal kinetics

Jay A. Messer, Abdallah S.R. Mohamed, Katherine A. Hutcheson, Yao Ding, Jan S. Lewin, Jihong Wang, Stephen Y. Lai, Steven J. Frank, Adam S. Garden, Vlad Sandulache, Hillary Eichelberger, Chloe C. French, Rivka R. Colen, Jack Phan, Jayashree Kalpathy-Cramer, John D. Hazle, David I. Rosenthal, G. Brandon Gunn, Clifton D. Fuller



**Figure 1.** 1a) T1 Baseline. 1b) T2 Early Post-RT after 3 months 1c) T1 Late Post-RT after 29 months 1d) Radiation dose grid 1e) Co-registration of MRI and planning CT

Research at MD Anderson

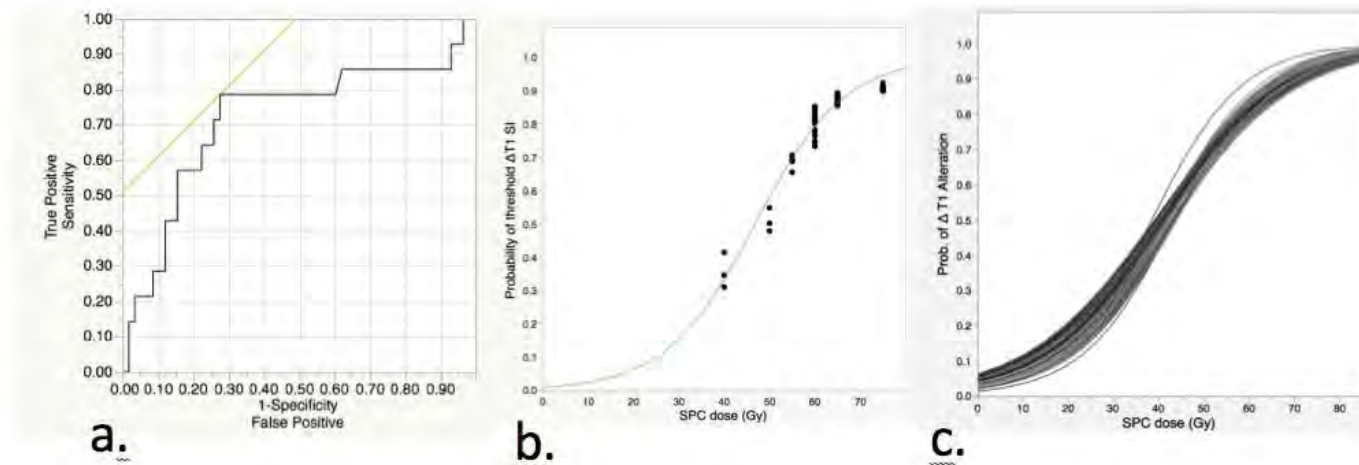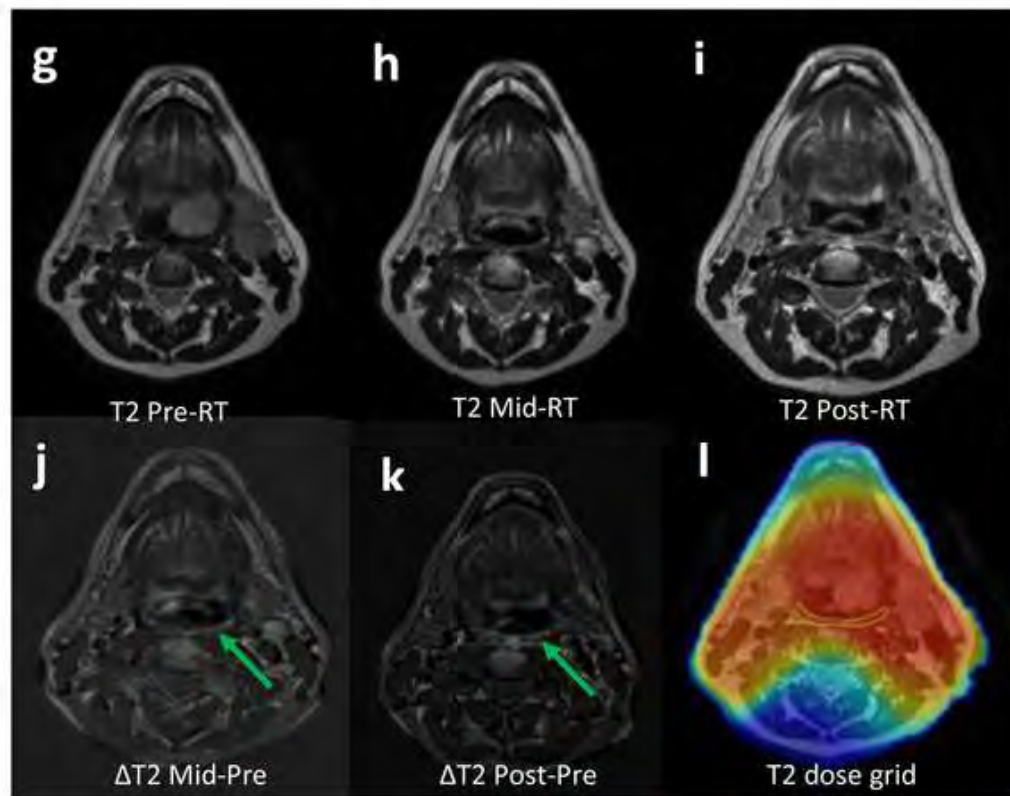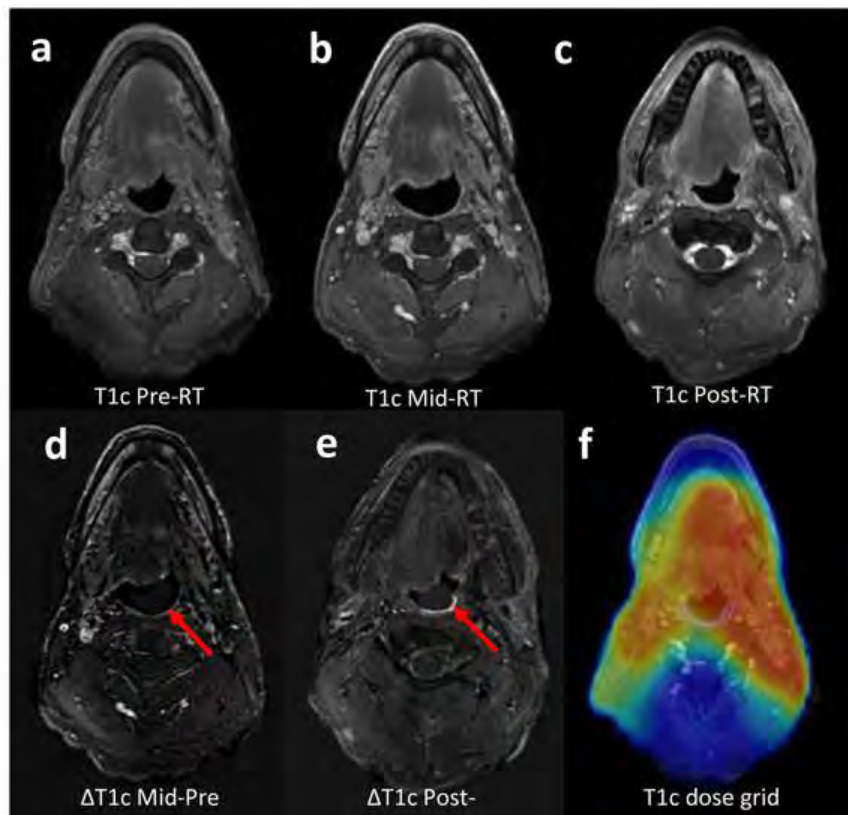# T1W Muscle damage/dose biomarker



**Figure 5**: Continuous (non-linear) dose response characterization of late T1 superior pharyngeal constrictor signal alteration from baseline. 5a. Confirmatory analysis of RPA-derived dose-threshold; Receiver operator characteristic curve (ROC), showing split performance for T1 signal intensity changes of greater than or less than 0.57 in the superior pharyngeal constrictors, as a function of $D_{mean}$, with area-under the curve (AUC) of 0.72 (P=0.013). 5b. Sigmoidal fit of observed probability of threshold T1 signal alteration as function of $D_{mean}$ to superior pharyngeal constrictor muscles (R2=0.93). 5c. Incidence–resampled bootstrap predicted probability of threshold T1 alteration as a function of dose; $10^4$ - independently-resampled distributions were individually fit using a maximum likelihood 2P-sigmoidal function, representing the range of possible dose-response normal tissue complication probability curves in order to best approximate a "true population incidence."
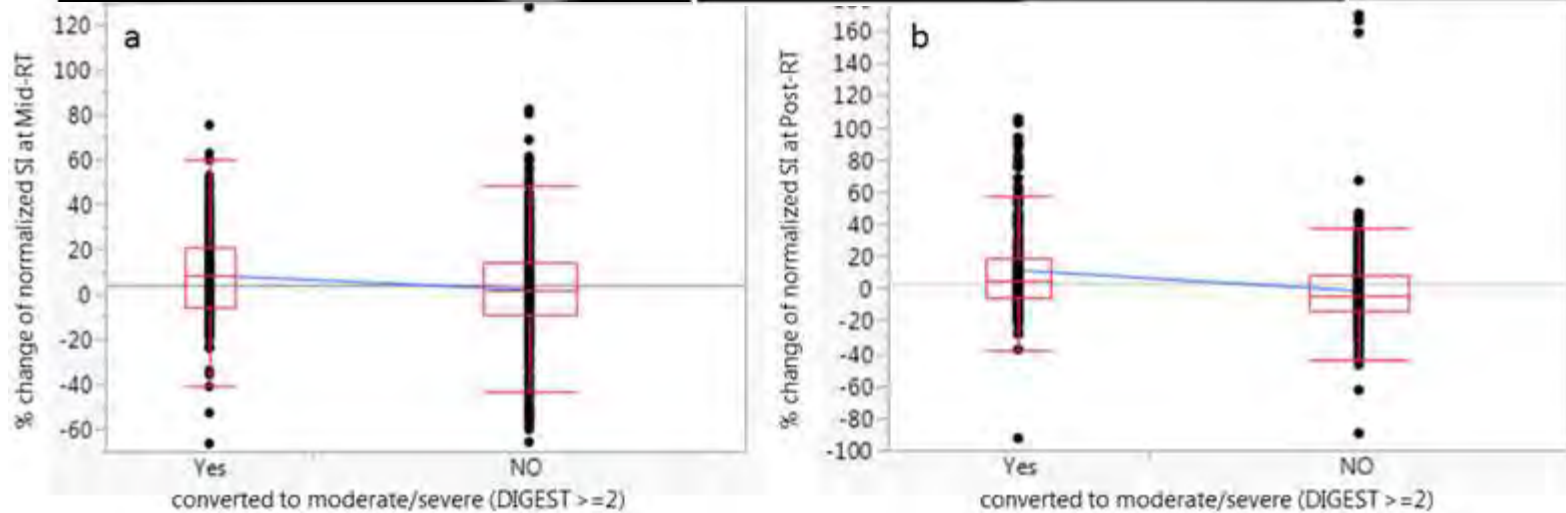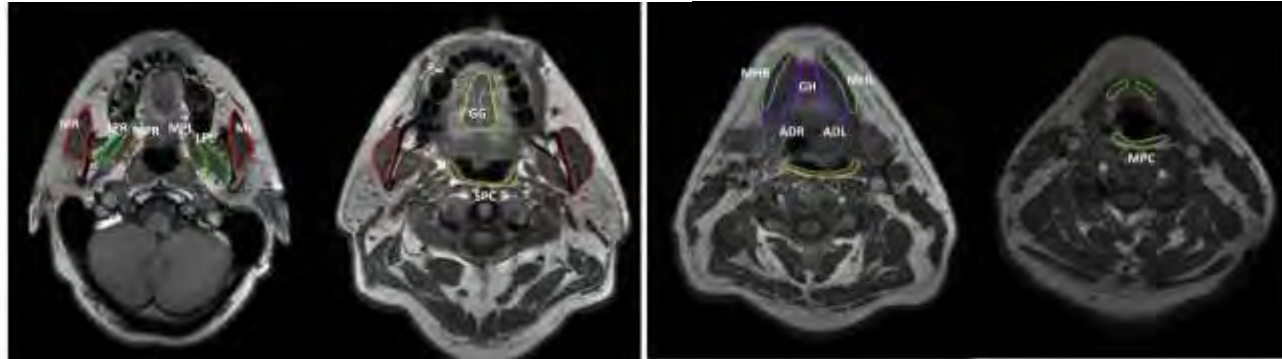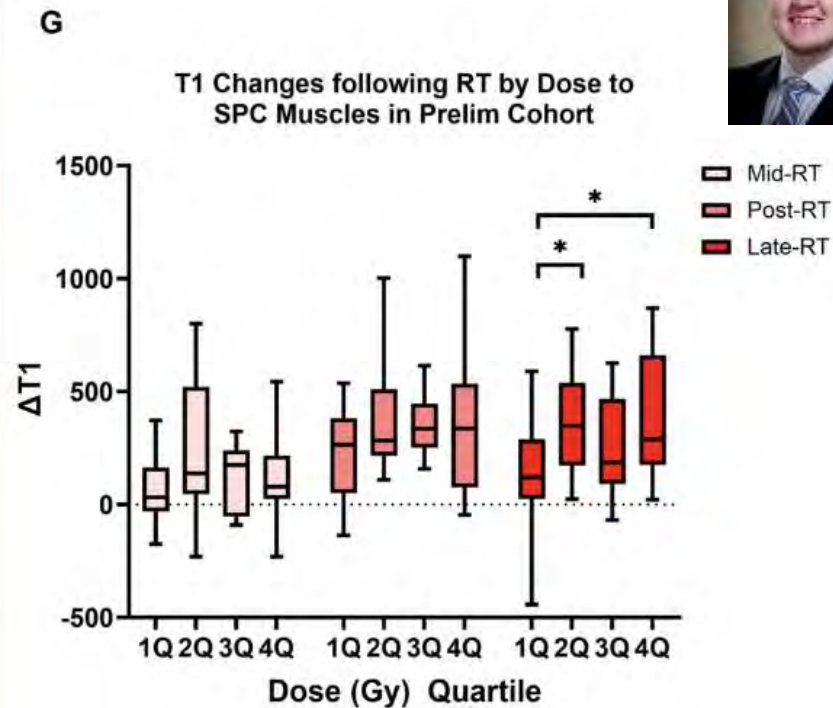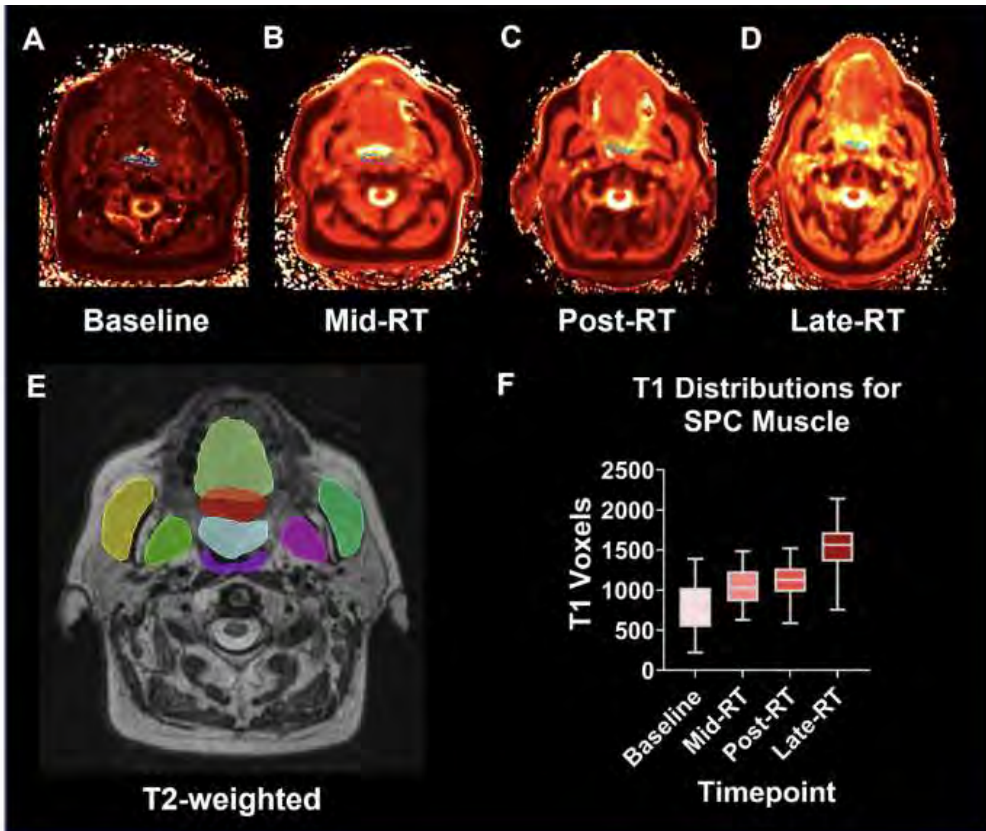
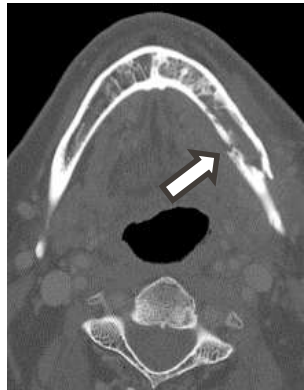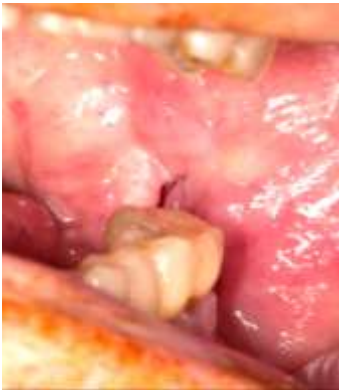# What if we just used standardized T1W/T2W MRI?

A prospective longitudinal assessment of MRI signal intensity kinetics of non-target muscles in patients with advanced stage oropharyngeal cancer in relationship to radiotherapy dose and post-treatment radiation-associated dysphagia: Preliminary findings from a randomized trial

A Baseline B Mid-RT C Post-RT D Late-RT

E T2-weighted

F **T1 Distributions for SPC Muscle**

G **T1 Changes following RT by Dose to SPC Muscles in Prelim Cohort**

Mid-RT
Post-RT
Late-RT

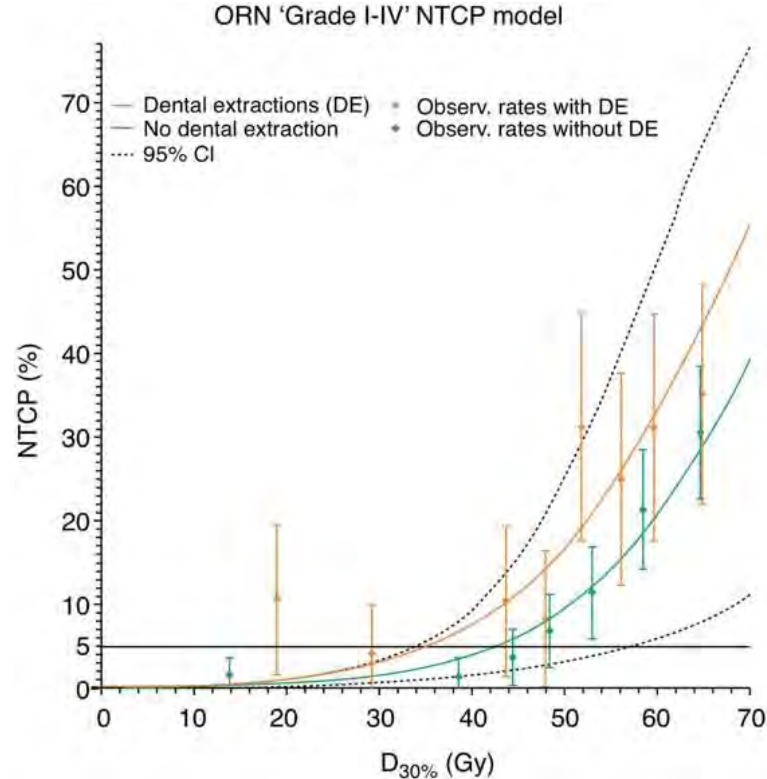# Osteoradionecrosis (ORN)



**"Exposed bone in a field of irradiation."**

**MDACC rate ~6-7%, which means about 65 cases/year**
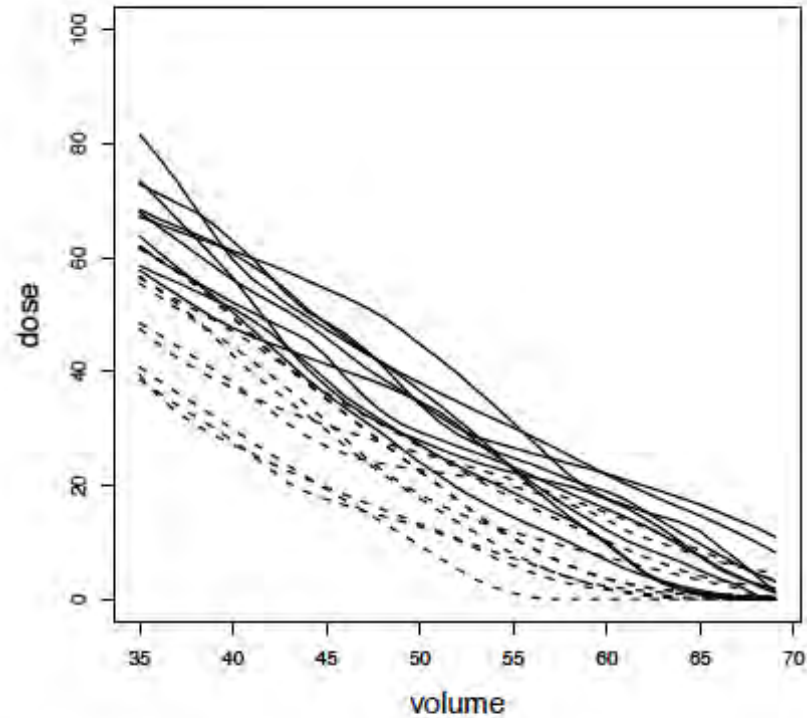
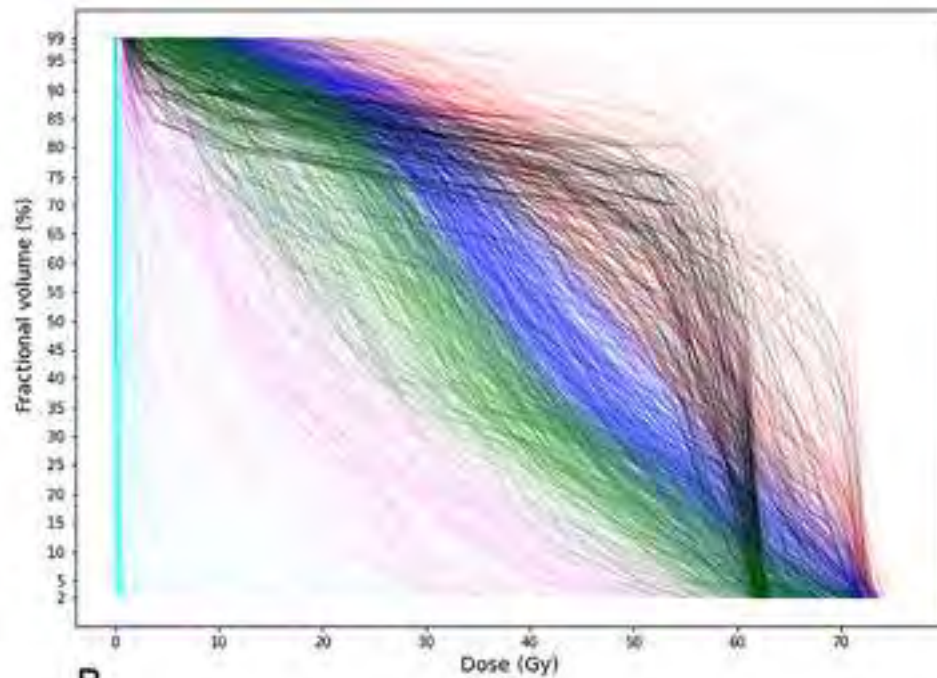# Normal Tissue Complication Probability (NTCP) For ORN



ORN 'Grade I-IV' NTCP model

adapted from van Dijk et al IJ*ROBP* 2021

Figure 1: The dose-volumes curves for patients with ORN (solid) versus controls (dotted).

Functional Principal Component Analysis for Dose-volume Correlates of
Mandibular Osteoradionecrosis

B

| Cluster no. | PDE = 0 | | PDE = 1 | |
|---|---|---|---|---|
| | ORN incidence | Risk index (95% CI) | ORN incidence | Risk index (95% CI) |
| 1 | 0 out of 58 | 0.0% | 0 out of 9 | 0.0% |
| 2 | 2 out of 68 | 2.9% (0.0%, 6.9%) | 0 out of 8 | 0.0% |
| 3 | 8 out of 273 | 2.9% (0.9%, 4.9%) | 10 out of 86 | 11.6% (4.8%, 18.4%) |
| 4 | 39 out of 318 | 12.3% (8.7%, 15.9%) | 34 out of 144 | 23.6% (16.7%, 30.5%) |
| 5 | 31 out of 118 | 26.3% (18.3%, 34.3%) | 14 out of 47 | 29.8% (16.7%, 42.9%) |
| 6 | 21 out of 82 | 25.6% (16.1%, 35.1%) | 14 out of 48 | 29.2% (16.3%, 42.1%) |

*Abbreviations:* ORN = osteoradionecrosis; PDE = preradiation dental extraction (0 = no/edentulous, 1 = dental extractions).

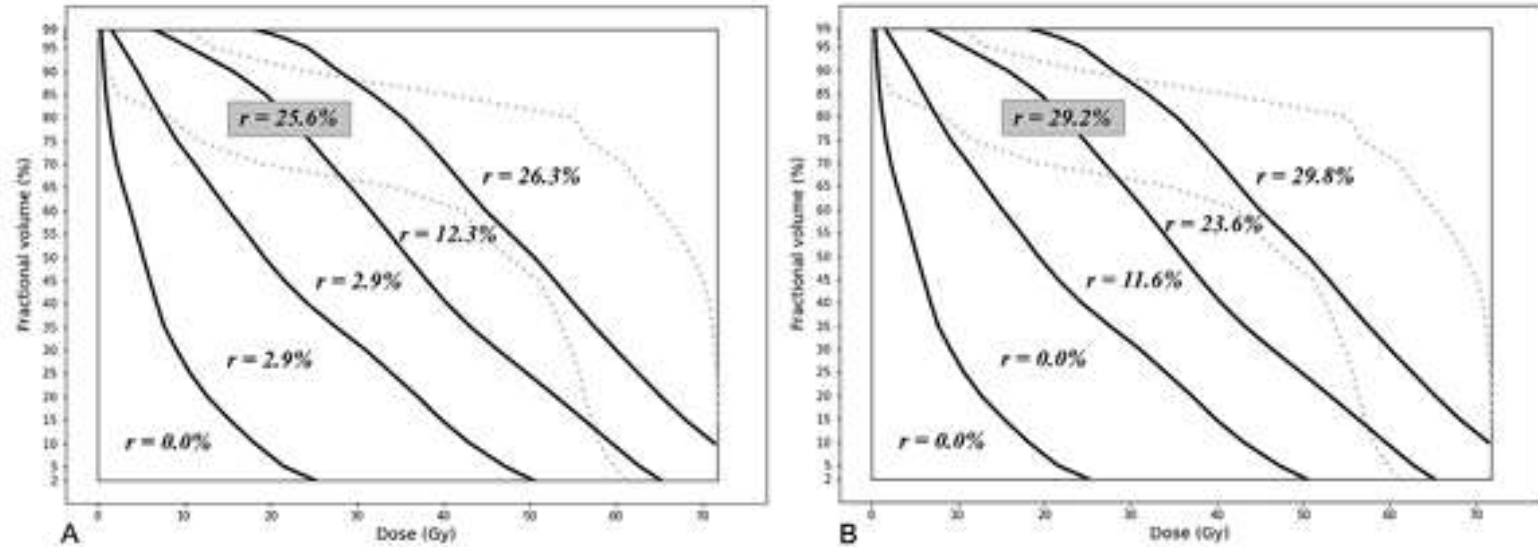**Fig. 4.** Risk indices of dose-volume regions for K = 6. (A) No/edentulous dental extractions (PDE = 0). (B) With dental extractions (PDE = 1).
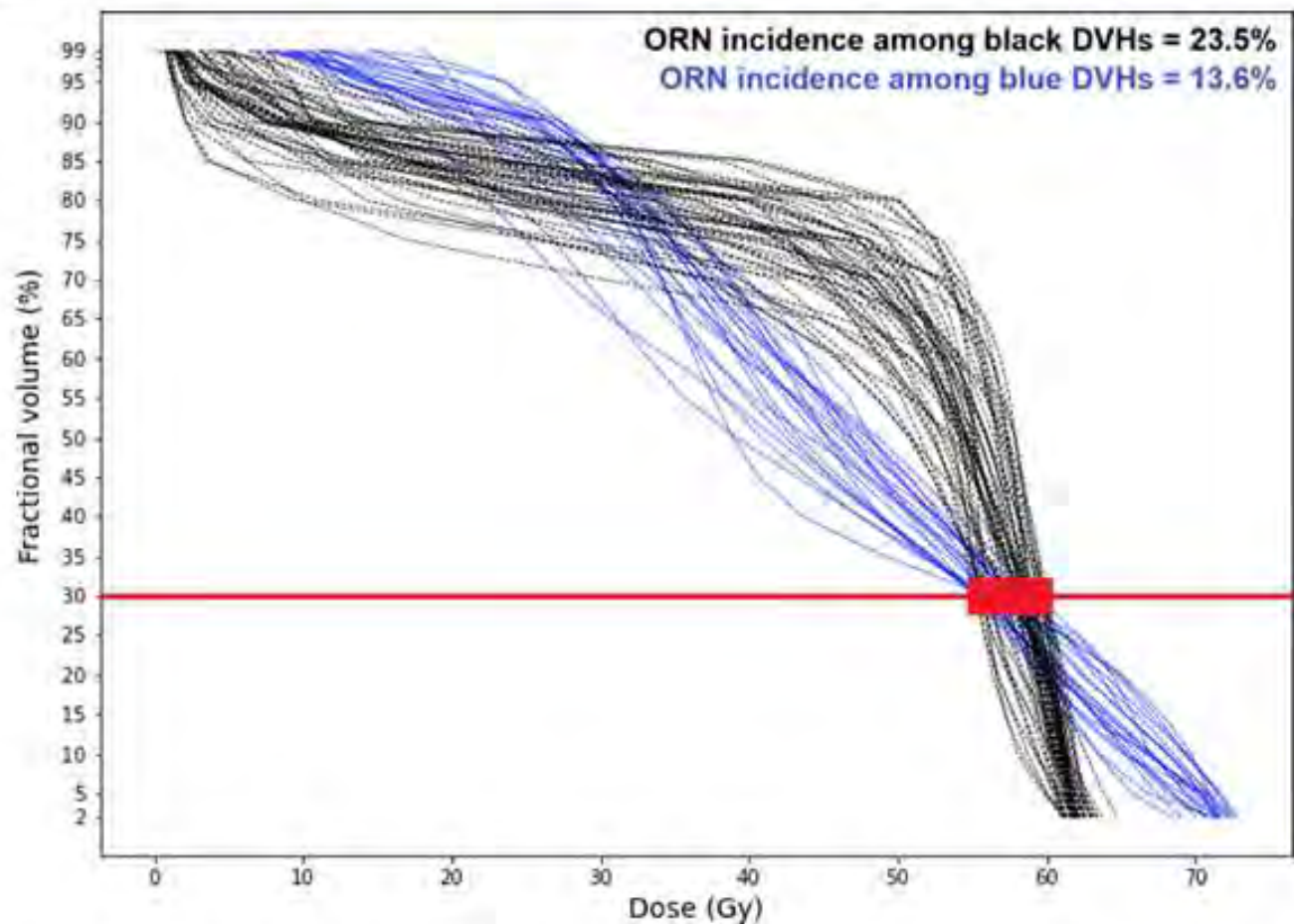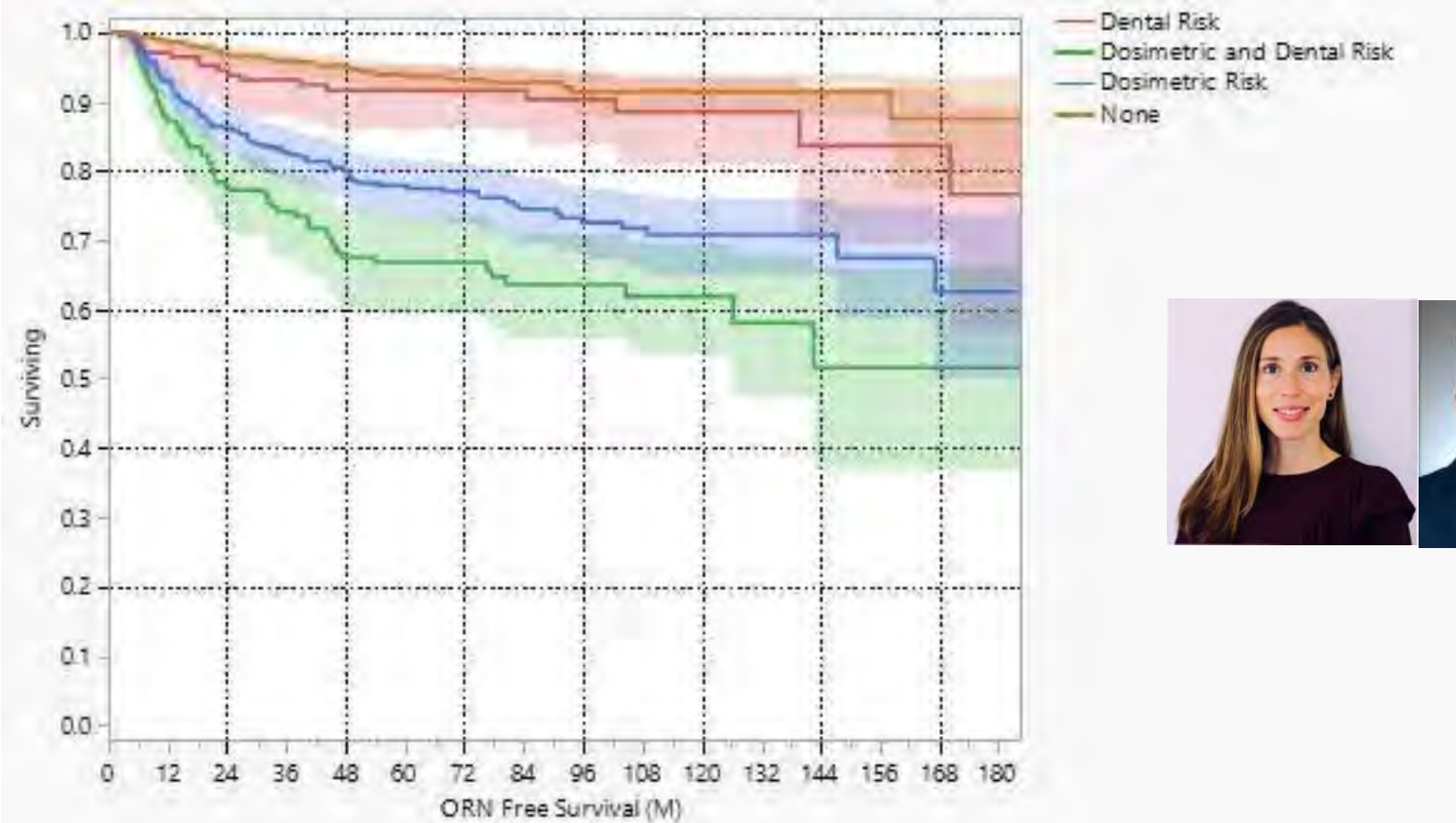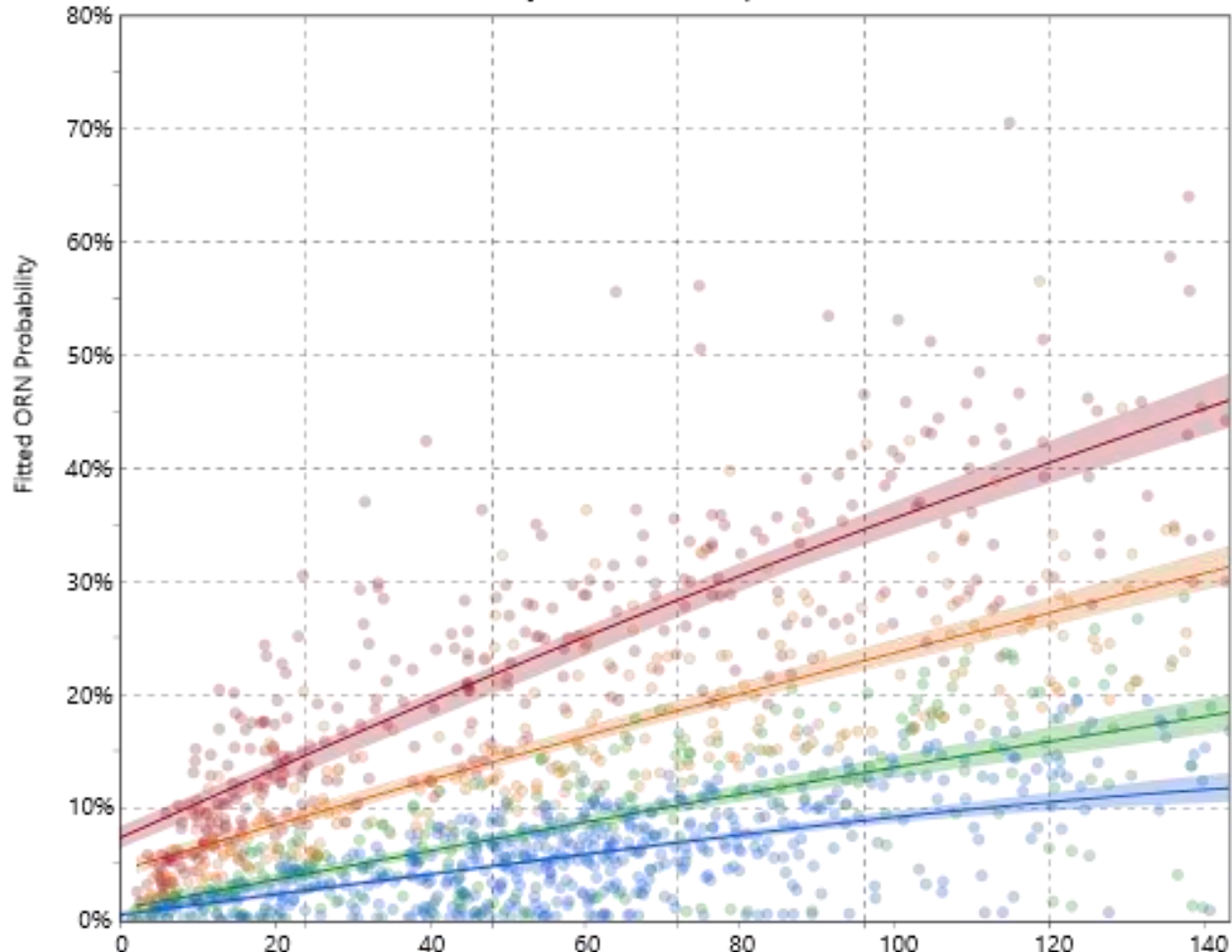
**Fig. 5.** Different osteoradionecrosis incidences among dose-volume histograms with the same D30% value.

# Temporal Awareness: Time to ORN

Fitted ORN Probability for Risk Groups Over Time in Months

# ORN Risk GUI

# Comparison of Machine-Learning and Deep-Learning Methods for the Prediction of Osteoradionecrosis Resulting From Head and Neck Cancer Radiation Therapy

Brandon Reber, BS,[a,e] Lisanne Van Dijk, PhD,[a,b] Brian Anderson, PhD,[a,c] Abdallah Sherif Radwan Mohamed, MD, PhD,[a] Clifton Fuller, MD, PhD,[a] Stephen Lai, MD, PhD,[a] and Kristy Brock, PhD[a]

**Table 1    Summary of subject demographics***

|  | ORN− | ORN+ |
|---|---|---|
| Number of subjects | 1086 | 173 |
| Age, y, median | 61 | 60 |
| Sex, male, n (%) | 894 (82%) | 150 (87%) |
| Smoking, current, n (%) | 153 (14%) | 27 (16%) |
| Smoking, pack-years, median | 7 | 8 |
| Postoperative RT | 172 (16%) | 44 (25%) |
| Dental extraction pre-RT | 270 (25%) | 72 (42%) |
| Tumor site |  |  |
| Oral cavity | 146 (13%) | 44 (25%) |
| Oropharynx | 703 (65%) | 123 (71%) |
| Hypopharynx/larynx/nasopharynx/unknown-primary | 237 (22%) | 6 (3%) |

*Abbreviations:* ORN = osteoradionecrosis; RT = radiation therapy.

\* Percent signs within cells indicate the percent of the subject cohort for the ORN− and ORN + cases separately that have each row attribute.

**Table 2    Mean (±SD) metric values for the cross-validation withheld folds for the ML models***

| Model | Accuracy | Balanced accuracy | Recall | Precision | F1 score | AUROC | AUPRC |
|---|---|---|---|---|---|---|---|
| Logistic regression | $0.69 \pm 0.05$ | $0.70 \pm 0.07$ | $0.72 \pm 0.14$ | $0.27 \pm 0.05$ | $0.39 \pm 0.07$ | $0.74 \pm 0.07$ | $0.28 \pm 0.08$ |
| Random forest | $0.65 \pm 0.05$ | $0.69 \pm 0.07$ | $0.74 \pm 0.14$ | $0.25 \pm 0.04$ | $0.37 \pm 0.06$ | $0.69 \pm 0.07$ | $0.23 \pm 0.04$ |
| Support vector machine | $0.69 \pm 0.04$ | $0.70 \pm 0.07$ | $0.71 \pm 0.13$ | $0.27 \pm 0.04$ | $0.39 \pm 0.06$ | $0.70 \pm 0.07$ | $0.24 \pm 0.04$ |
| Random classifier | $0.52 \pm 0.04$ | $0.49 \pm 0.08$ | $0.45 \pm 0.14$ | $0.14 \pm 0.04$ | $0.21 \pm 0.07$ | $0.50 \pm 0.00$ | $0.14 \pm 0.01$ |

*Abbreviations:* AUPRC = area under the precision recall curve; AUROC = area under the receiver operating characteristic curve; ML = machine learning.
*   Each cell shows the mean (±SD) of the metrics from the withheld folds of the stratified 10-fold cross-validation with 10 repeats.

**Table 3    Performance of the best DL models for each architecture type***

| Architecture | Accuracy | Balanced accuracy | Recall | Precision | F1 score | AUROC | AUPRC |
|---|---|---|---|---|---|---|---|
| ResNet | 0.87 | 0.69 | 0.04 | 0.50 | 0.07 | 0.57 | 0.23 |
| DenseNet | 0.83 | 0.54 | 0.10 | 0.21 | 0.14 | 0.58 | 0.17 |
| Autoencoder | 0.71 | 0.53 | 0.33 | 0.18 | 0.23 | 0.59 | 0.15 |
| Random | 0.49 | 0.46 | 0.46 | 0.11 | 0.17 | 0.49 | 0.13 |

*Abbreviations:* AUPRC = area under the precision recall curve; AUROC = area under the receiver operating characteristic curve; DL = deep learning.
*   The reported metrics are from the withheld test set not used during model training or selection. Metrics sensitive to data imbalance, such balanced accuracy, F1 score, and AUPRC, were lower than those for the logistic regression model using the test set.

**Figure 1**  Deep-learning model performance with increasing amounts of training data.

## Conclusion

In this work, we compared traditional ML algorithms to DL algorithms for the prediction of mandible ORN resulting from HNC RT. The traditional ML algorithms performed similarly to each other when using cross-validation and were successful at predicting ORN. The performance of the ML models shows promise in clinical integration for future studies. Despite our use of different architectures and model ensembles, the DL models continued to underperform compared to the best-performing ML algorithm identified by cross-validation, logistic regression, when evaluated on the test set. When we used additional training data, no performance improvement trends were evident, suggesting that more data are needed despite the relatively large HNC patient cohort. In further work, researchers could use more

# So how does AI model adoption practically occur?
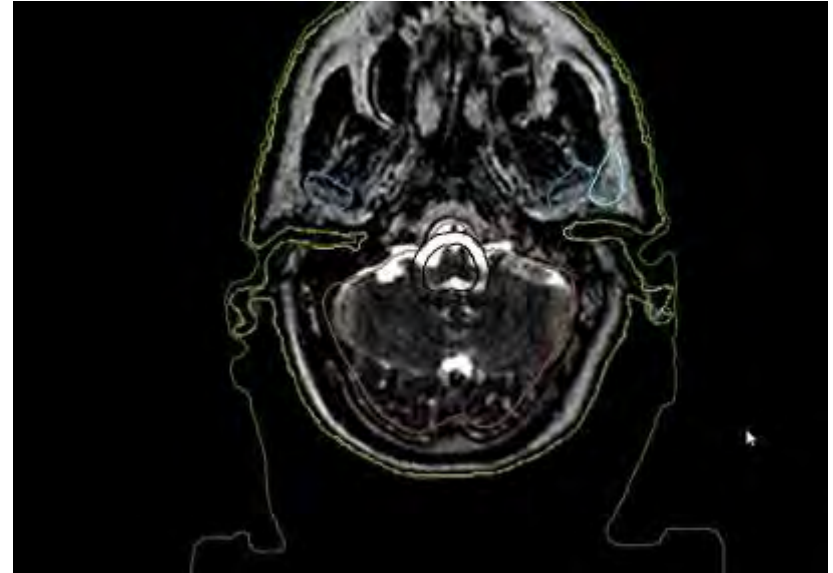
Research at MD Anderson

# Real Life: Use-case specific acceptance testing



1. Define acceptance criteria
2. Plan acceptance testing
3. Derive acceptance tests
4. Run acceptance tests
5. Negotiate test results
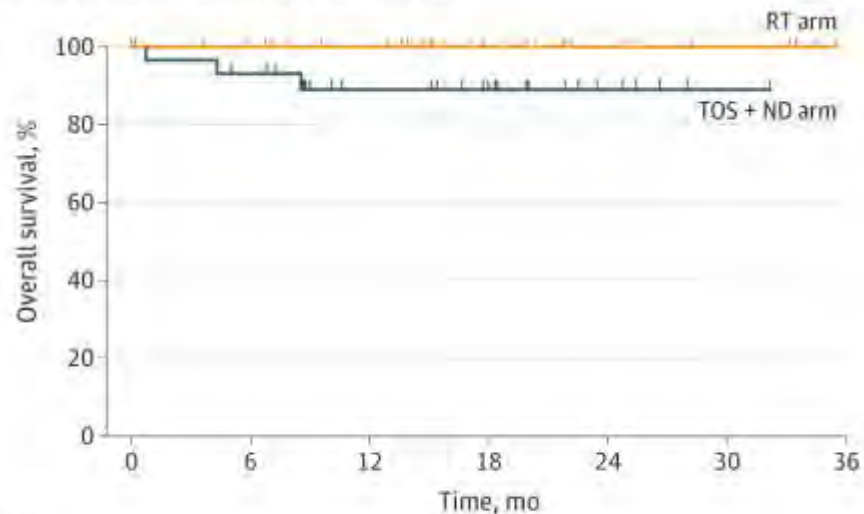6. Reject/accept system

# Example: Decision Support Tools for Surgical vs. Non-surgical therapy selection
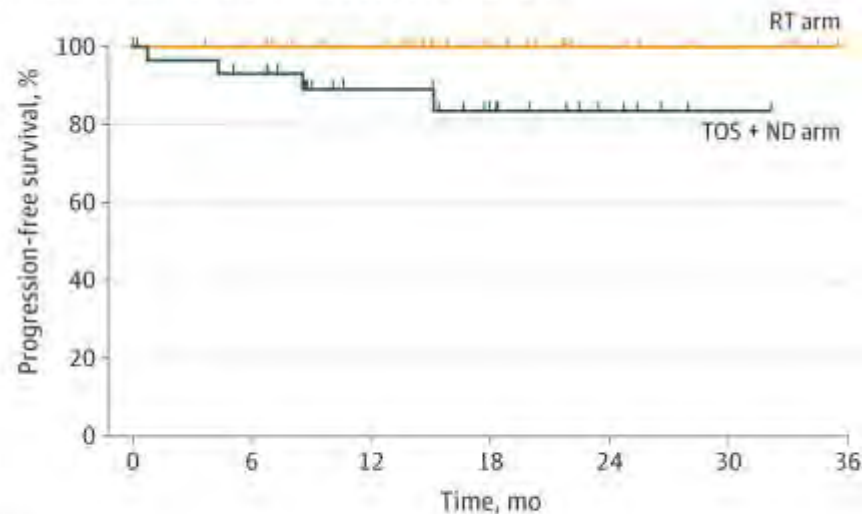
# Example: Decision Support Tools
## ORATOR2



A | Overall survival stratified by treatment arm

B | Progression-free survival stratified by treatment arm

Research at MD Anderson

# MDs/MDTs are bad at quantification of risk

**If I do TORS, there is no PM or ECE >> Best outcome**

• I have spared RT ☺

**If I do TORS, and there is low volume ENE or close margin**

• *Need adjuvant RT [bimodality]*

• MDADI is the same as RT alone,

• DIGEST is *worse* than RT alone 😐

**If I do TORS, and there is PM or >2mm ENE**

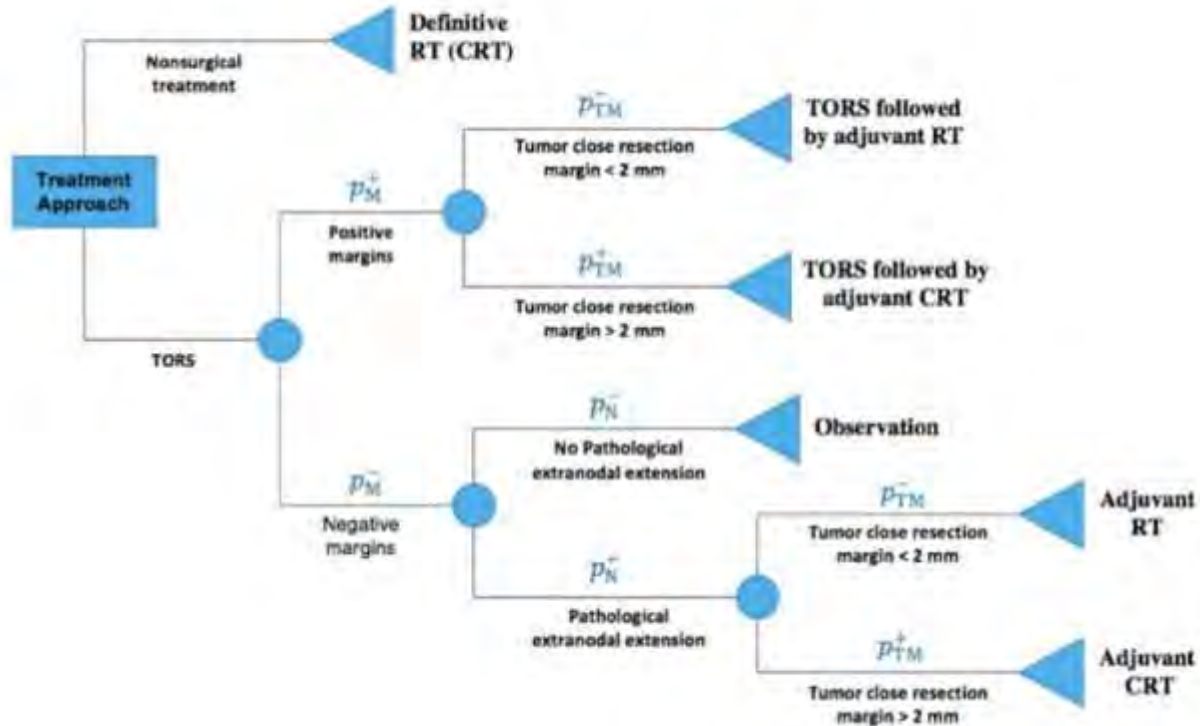• *Need adjuvant chemoRT*

• MDADI/DIGEST is worse than chemo(RT) ☹

We are bad at quantification of risk

**MDADI scores in HPV+ OPSCC**

McDowell L, et al (unpublished, 2023)

Research at MD Anderson

**Optimized decision support for selection of transoral robotic surgery or (chemo)radiation therapy based on posttreatment swallowing toxicity**
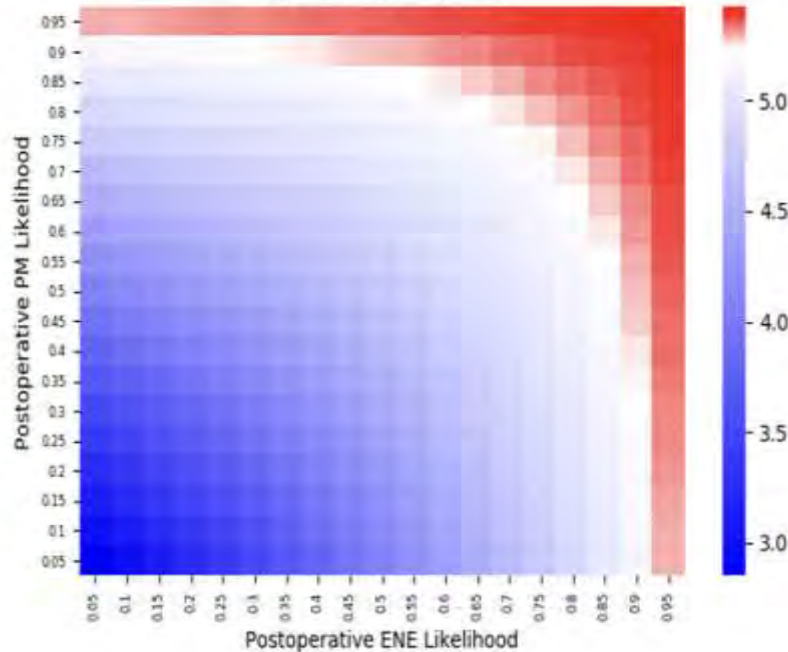


Research at MD Anderson

# Optimized decision support for selection of transoral robotic surgery or (chemo)radiation therapy based on posttreatment swallowing toxicity          DOI: 10.1002/cam4.5253

**TABLE 3** Range of likelihoods required for TORS and definitive therapies to become the optimal treatment under the second scenario

| Scenario I | | |
|---|---|---|
| **Instrument/measure** | **Confidence level of postoperative events for which TORS is optimal** | **Confidence level of postoperative events for which definitive RT is optimal** |
| MDADI | | |
| Short term (3–6 months) | — | Any likelihood associated with ENE and/or PM |
| Long term (18–24 months) | When both ENE and PM have likelihood <70% | If either of ENE or PM has a likelihood >90% |
| MDASI | | |
| Short term (3–6 months) | — | Any likelihood associated with ENE and/or PM |
| Long term (18–24 months) | Any likelihood associated with ENE and/or PM | — |
| DIGEST | | |
| Short term (3–6 months) | When both ENE and PM have likelihood <40% | If either of ENE or PM has a likelihood >75% |
| Long term (18–24 months) | When both ENE and PM have likelihood <10% | If either of ENE or PM has a likelihood >25% |
| **Scenario II** | | |
| **Instrument/measure** | **Confidence level of postoperative events for which TORS is optimal** | **Confidence level of postoperative events for which definitive CRT is optimal** |
| MDADI | | |
| Short term (3–6 months) | Any likelihood associated with ENE and/or PM | — |
| Long term (18–24 months) | Any likelihood associated with ENE and/or PM | — |
| MDASI | | |
| Short term (3–6 months) | Any likelihood associated with ENE and/or PM | — |
| Long term (18–24 months) | Any likelihood associated with ENE and/or PM | — |
| DIGEST | | |
| Short term (3–6 months) | When both ENE and PM have likelihood <55% | If either of ENE or PM has a likelihood >80% |
| Long term (18–24 months) | When both ENE and PM have likelihood <20% | If either of ENE or PM has a likelihood >40% |

Abbreviations: ENE, postoperative extranodal extension; PM, postoperative positive margin.
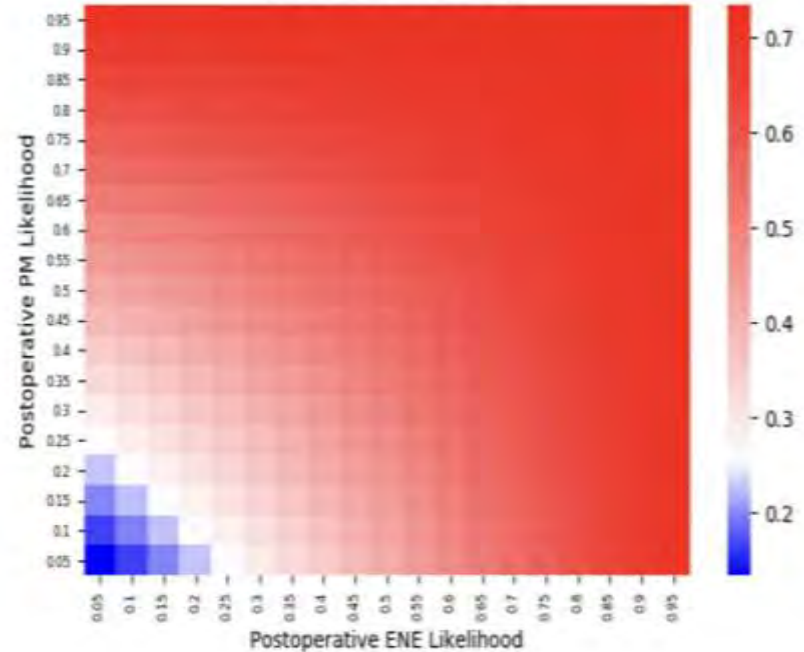
Research at MD Anderson

**Optimized decision support for selection of transoral robotic surgery or (chemo)radiation therapy based on posttreatment swallowing toxicity**

Red==RT better Blue==TORS better



**(B)**
Comparison of TORS vs. RT based on $\Delta_L^{MDADI}$-measure ($c_L = 5.246$, $r = 21.0\%$)
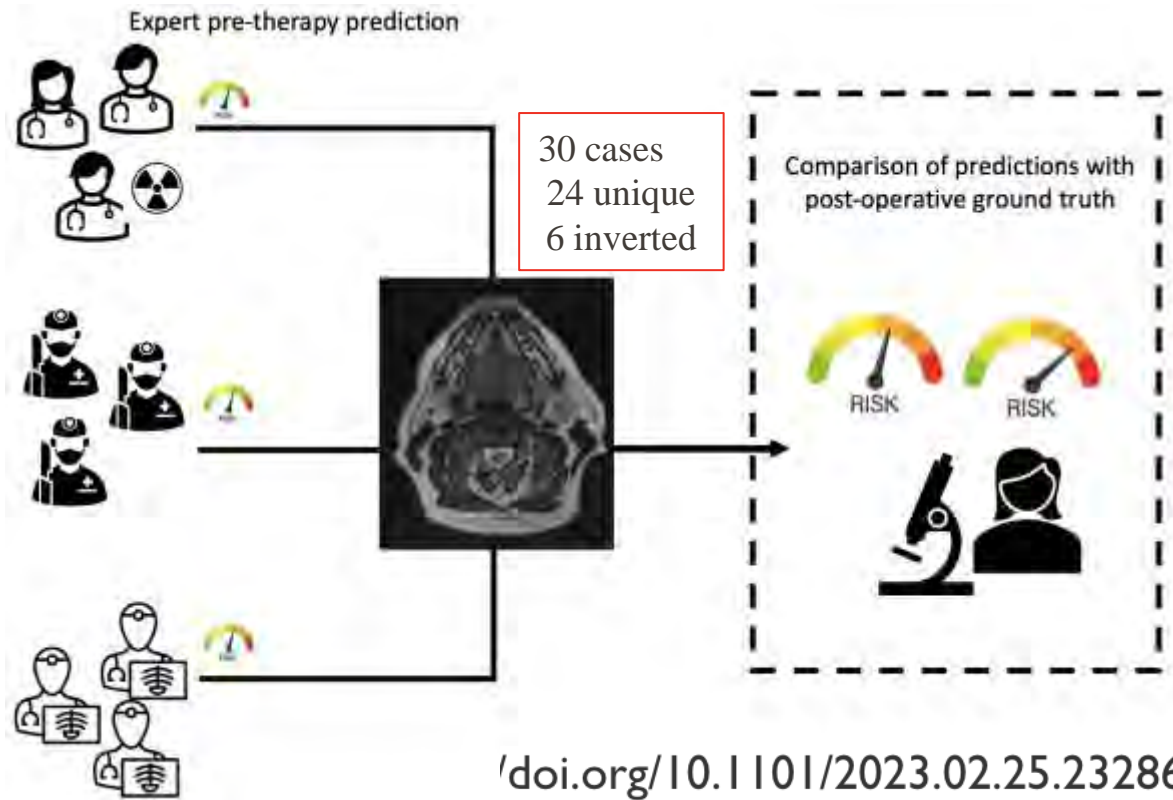
Comparison of TORS vs. RT based on $R^{DIGEST}$-measure ($c_L = 0.255$, $r = 97.0\%$)

# Multi-Specialty Expert Physician Identification of Extranodal Extension in Computed Tomography Scans of Oropharyngeal Cancer Patients: Prospective Blinded Human Inter-Observer Performance Evaluation

# Problem:
# Humans are crummy at pathologic ENE (pre)detection



95% CI of ROC Curves for Radiology

95% CI of ROC Curves for RadOnc

95% CI of ROC Curves for Surgery

# Pretreatment Identification of Head and Neck Cancer Nodal Metastasis and Extranodal Extension Using Deep Learning Neural Networks
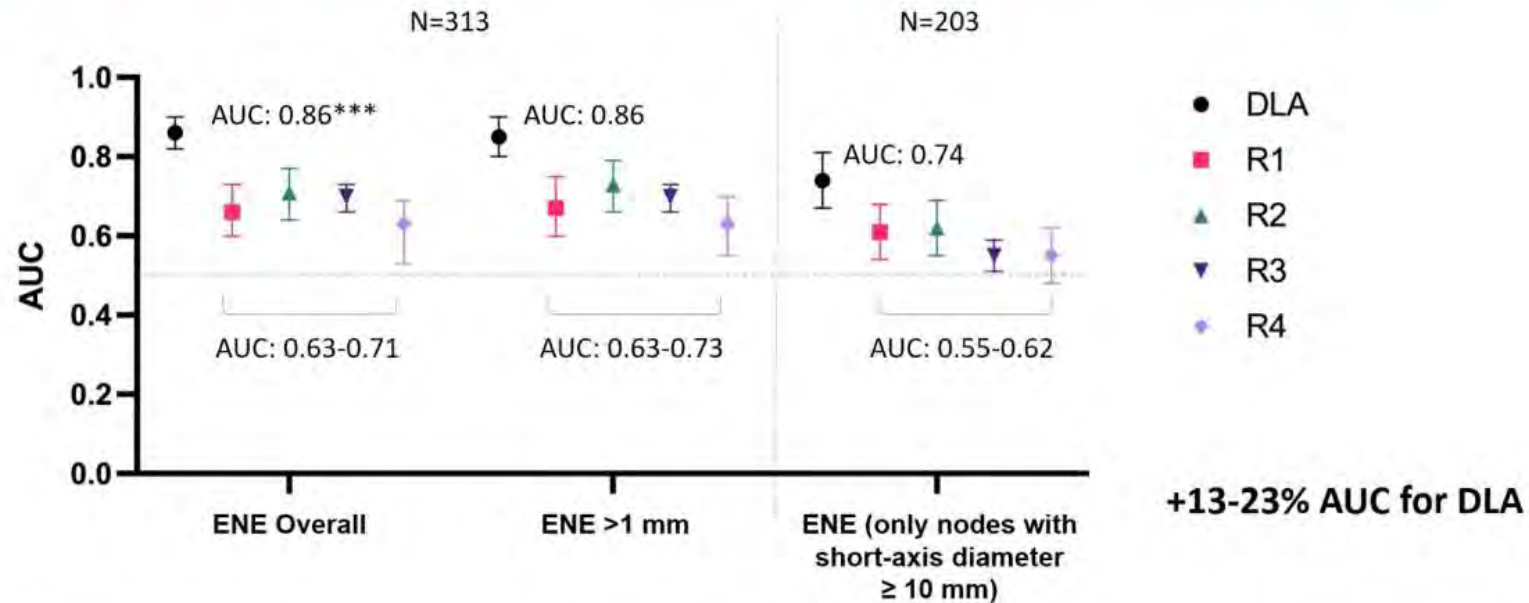
Benjamin H. Kann[1], Sanjay Aneja[1], Gokoulakrichenane V. Loganadane[1], Jacqueline R. Kelly[1], Stephen M. Smith[2], Roy H. Decker[1], James B. Yu[1], Henry S. Park[1], Wendell G. Yarbrough[3], Ajay Malhotra[4], Barbara A. Burtness[5] & Zain A. Husain[1]

**Figure 1.** (A,B) Lymph Node Region of Interest Preprocessing. (A) 2D representation of 3D lymph node segmentation preprocessing resulting in a dimension-preserving input (1) and a size-invariant, "zoomed-in" input (2). (B) Representation of actual 3D input arrays for dual-input deep learning neural network.

# Pretreatment Identification of Head and Neck Cancer Nodal Metastasis and Extranodal Extension Using Deep Learning Neural Networks

| Performance Metric | Extranodal Extension (ENE) | | | Nodal Metastasis (NM) | | |
|---|---|---|---|---|---|---|
| | ENE Test Set (n=98*) | | | Test Set (n=131) | | |
| | DualNet DLNN | Random Forest | Benchmark Logistic | DualNet DLNN | Random Forest | Benchmark Logistic |
| AUC | 0.91 | 0.88 | 0.81 | 0.91 | 0.91 | 0.86 |
| Accuracy | 85.7% | 82.6% | 77.7% | 85.5% | 84.7% | 76.1% |
| Sensitivity | 0.88 | 0.79 | 0.72 | 0.84 | 0.75 | 0.79 |
| Specificity | 0.85 | 0.84 | 0.80 | 0.87 | 0.92 | 0.74 |
| PPV | 0.66 | 0.61 | 0.54 | 0.88 | 0.87 | 0.69 |
| NPV | 0.95 | 0.93 | 0.89 | 0.82 | 0.83 | 0.83 |
| Youden Index | 0.73 | 0.63 | 0.51 | 0.71 | 0.67 | 0.53 |

Table 3. Model Performance and Benchmark Comparisons on Independent Test Set By Lymph Node Feature. *Test set for ENE includes lymph nodes with region of interest diameters $\geq$ 1 cm. Abbreviations: AUC = area under the curve; PPV = positive predictive value; NPV = negative predictive value. Youden index = Sensitivity + Specificity − 1.

MD Anderson

# Deep learning outperforms radiologists for ENE prediction in E3311

N=313  N=203

AUC (y-axis): 1.0, 0.8, 0.6, 0.4, 0.2, 0.0

AUC: 0.86*** AUC: 0.86 AUC: 0.74

AUC: 0.63-0.71 AUC: 0.63-0.73 AUC: 0.55-0.62

ENE Overall  ENE >1 mm  ENE (only nodes with short-axis diameter ≥ 10 mm)

Legend: ● DLA, ■ R1, ▲ R2, ▼ R3, ✦ R4

**+13-23% AUC for DLA**

***All P-values < 0.001 for DLA vs R1-4 comparisons. Error bars represent 95% CIs

*Manuscript under review*

**DLA outperformed expert radiologists and can have clinical utility in selection of patients appropriate for operative management and other de-escalation (or escalation) strategies for HPV + OPC**

MD Anderson

## So, why aren't we using these tools?

- **"I'm not sure about *this* case…"**
- **"What if it misses a node?"**
- **"I just don't trust it like I trust my colleagues…"**

# The current clinical problem: Trustworthiness/Uncertainty Estimation

Research at MD Anderson

# The current clinical problem: Trustworthiness/Uncertainty Estimation



Current AI Approaches

Data — Single **"Black box"** Estimator — Cohort Testing — Cohort Validation — 80% Model Performance — ? Unknown/ unquantified performance Certainty for individual patient

Uncertainty-quantified approaches

Data — Probabilistic Estimator — Cohort Testing — Cohort Validation — 80% Model Performance — Case-specific accuracy estimation

<45% certainty
<65% certainty
>85% certainty
>90% certainty

# Statement: Without uncertainty quantification, we cannot move forward

UQ consists of activities such as model verification, sensitivity analysis, calibration, surrogate modeling, validation, and uncertainty propagation. *Forward UQ* quantifies uncertainty in the model output given uncertainties in the inputs, model parameters, and model errors. *Inverse UQ* is related to model calibration which updates model parameter uncertainty using measurements (which are also uncertain).



**Sources of Uncertainty**

# The current **clinical** problem: Trustworthiness/Uncertainty Estimation

K. Zou et al. Meta-Radiology 1 (2023) 100003
https://doi.org/10.1016/j.metrad.2023.100003

# The current **clinical** problem:
Tl



**Fig. 2.** Visualization of the aleatoric (data) and the epistemic (model) uncertainty for the classification model.

Research at MD Anderson

# The current clinical problem: Trustworthiness/Uncertainty Estimation



Fig. 3. The different methods of uncertainty estimation.

# Application of simultaneous uncertainty quantification and segmentation for oropharyngeal cancer use-case with Bayesian deep learning

# Artificial Intelligence Uncertainty Quantification in Radiotherapy Applications – A Scoping Review

PMCID: PMC11118597

# Uncertainty estimation allows **direct** safety assessment

Risk Estimation flow charts from ISO 14971:2019

# Breiman's "Two Cultures" Revisited and Reconciled

Subhadeep Mukhopadhyay[1]

deep@unitedstatalgo.com

Kaijun Wang

kwang2@fredhutch.org

https://arxiv.org/abs/2005.13596



Figure 3: Integrated statistical learning framework at a glance; 'ML' stands for (an arbitrary) machine learning algorithm, and 'UPM' denotes uncertainty prediction machine.

# Oncologic prediction GUI



A) Serial prediction model design

B) Overview of datasets and splits for the clinical models

* Missing HPV, pack years and/or performance

Sanne van Dijk, PhD
UMC Gronigen

Research at MD Anderson

| | MDACC Train |
|---|---|
| **c-index** | **0.71** [0.65-0.77] |



**MDACC Training cohort**

p < 0.0001

Overall Survival

- risk<0.05
- risk>=0.05&<0.25
- risk>0.25

Survival time (years)

Number at risk

| | | | | | |
|---|---|---|---|---|---|
| 240 | 186 | 111 | 38 | 17 | 3 |
| 833 | 634 | 360 | 168 | 82 | 31 |
| 147 | 79 | 45 | 29 | 13 | 4 |

Survival time (years)

Sanne van Dijk, PhD
UMC Gronigen

# Web-based **individual** OS risk prediction in new patients

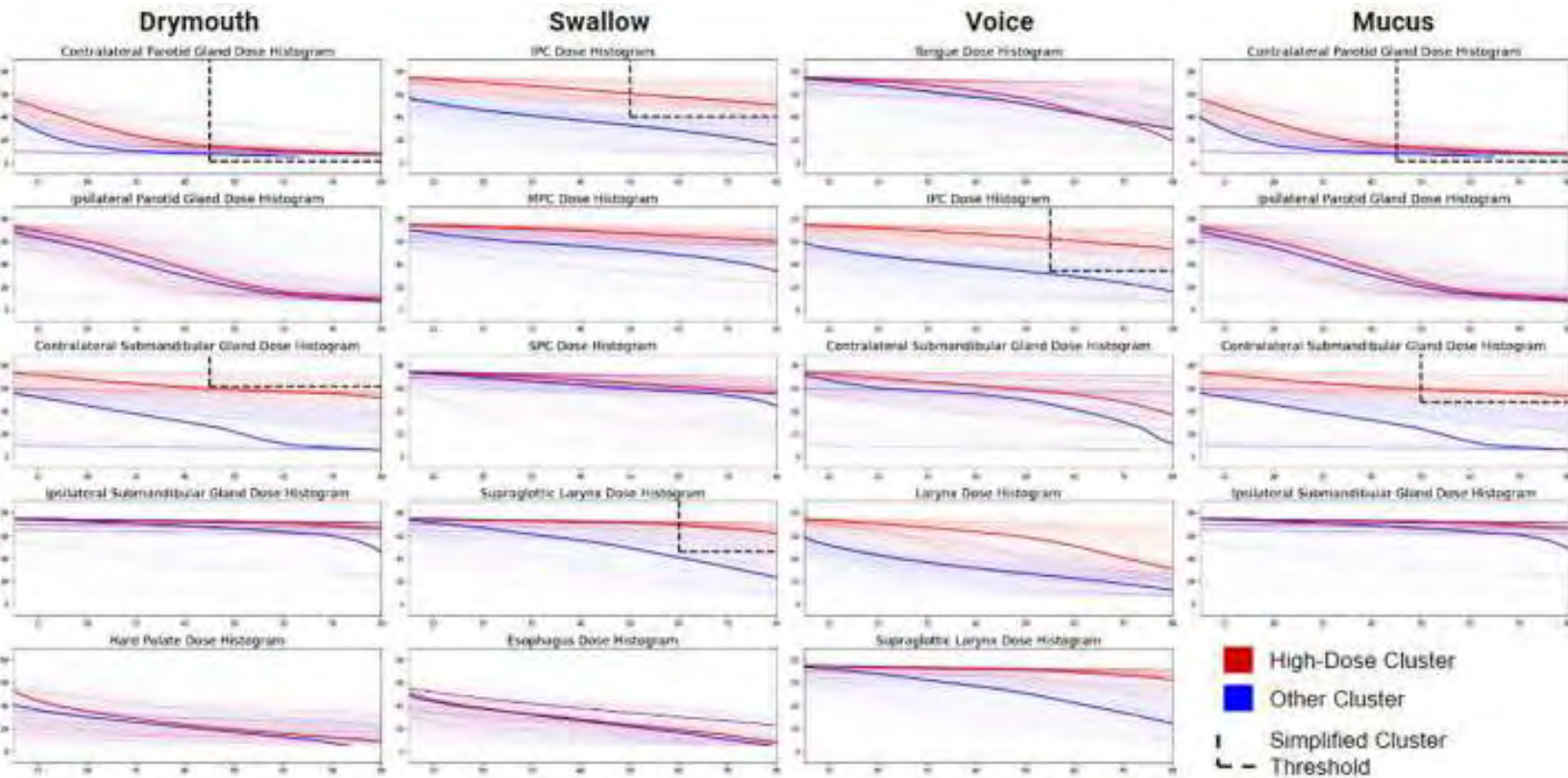# THALIS: Human–Machine Analysis of Longitudinal Symptoms in Cancer Therapy

# Multi-Organ Spatial Stratification of 3-D Dose Distributions Improves Risk Prediction of Long-Term Self-Reported Severe Symptoms in Oropharyngeal Cancer Patients Receiving Radiotherapy: Development of a Pre-Treatment Decision Support Tool.
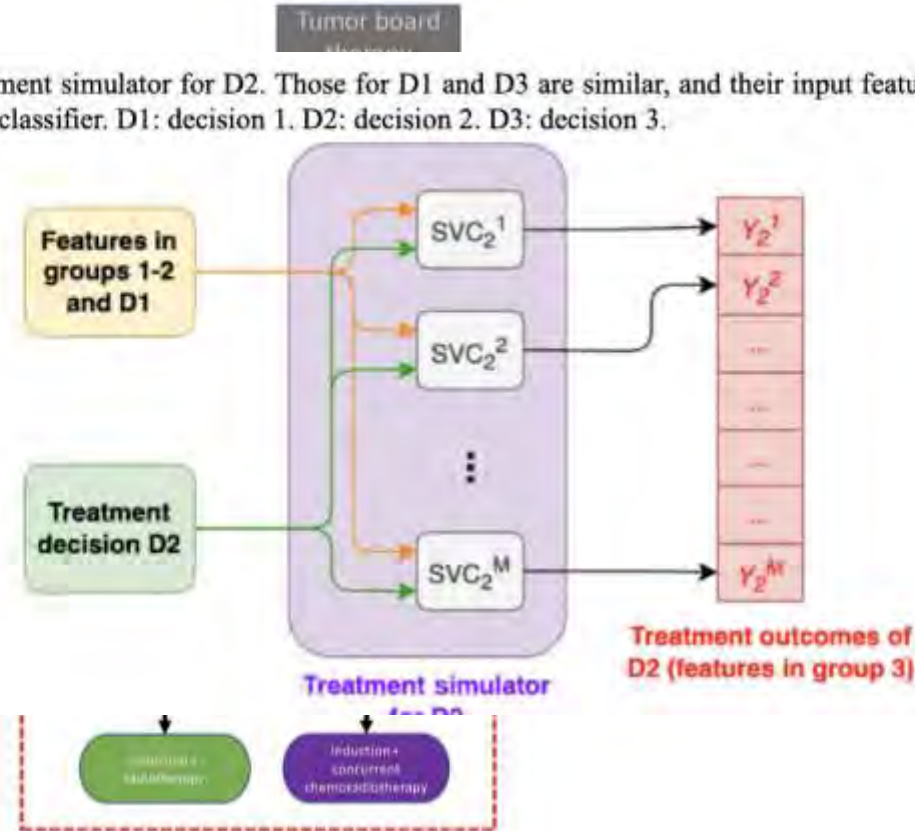
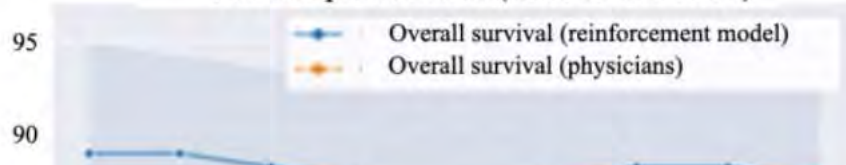# Predicting dynamic injury AND response kinetics

# Optimal Treatment Selection in Sequential Systemic and Locoregional Therapy of Oropharyngeal Squamous Carcinomas: Deep Q-Learning With a Patient-Physician Digital Twin Dyad
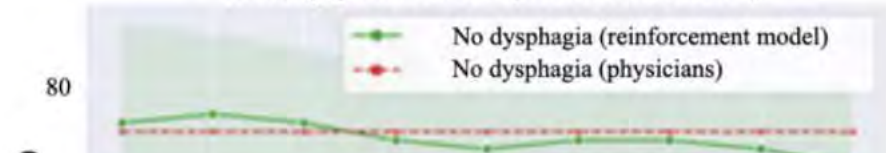


**Figure 4.** Illustration of the treatment simulator for D2. Those for D1 and D3 are similar, and their input features are from group 1 and groups 1-3, respectively. SVC: support vector classifier. D1: decision 1. D2: decision 2. D3: decision 3.

AI is good at survival prediction AND selecting therapy based on toxicity]

Optimal Treatment Selection in Sequential Systemic and
Locoregional Therapy of Oropharyngeal Squamous Carcinomas:
Deep Q-Learning With a Patient-Physician Digital Twin Dyad

Survival performance (without radiomics)
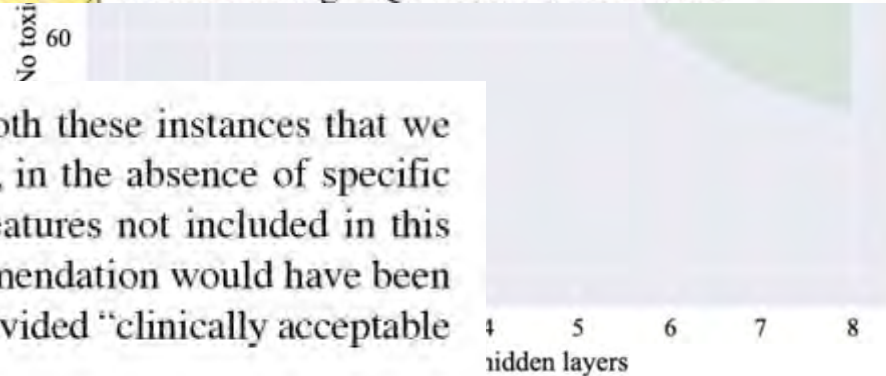
- Overall survival (reinforcement model)
- Overall survival (physicians)

Toxicity performance (without radiomics)

- No dysphagia (reinforcement model)
- No dysphagia (physicians)

**Results:** On the test set, we found mean 87.35% (SD 11.15%) and median 90.85% (IQR 13.56%) accuracies in treatment outcome prediction, matching the clinicians' outcomes and improving the (predicted) survival rate by +3.73% (95% CI −0.75% to 8.96%) and the dysphagia rate by +0.75% (95% CI −4.48% to 6.72%) when following DQL treatment decisions.

Overall, the physician review in both these instances that we investigated in detail suggests that, in the absence of specific *local* practices or occult clinical features not included in this decision platform, the DQL recommendation would have been a good strategy and that the dyad provided "clinically acceptable recommendations."]

# DITTO: A Visual Digital-twin for Interventions and Temporal Treatment Outcomes in Head and Neck Cancer

Andrew Wentzel, Serageldin Attia, Xinhua Zhang, Guadalupe Canahuate, Clifton David Fuller, and G.Elisabeta Marai
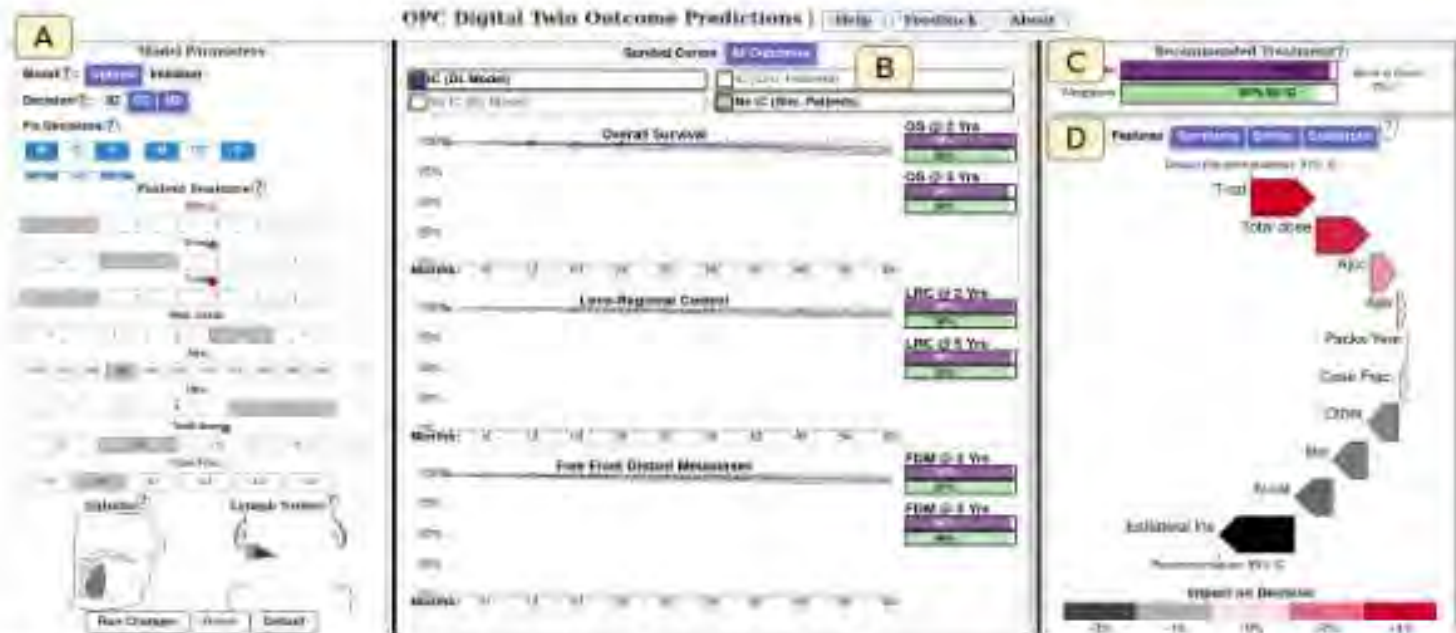
Fig. 1: Overview of DITTO. (A) Input panel to alter model parameters and input patient features. (B) Temporal outcome risk plots for the patient based on different models and treatment groups. (C) Treatment recommendation based on the twin model and similar patients. (D) Auxiliary data panel, currently showing a waterfall plot of how each feature cumulatively contributes to the model decision.

# But the view looks good for computational models in #RadOnc



**Please email/visit soon!**

**cdfuller@mdanderson.org**