# Improved Multimodal Fusion for Small Datasets with Auxiliary Supervision



Gregory Holste<sup>1,2</sup>, Douwe Van der Wal<sup>2</sup>, Hans Pinckaers<sup>2</sup>, Rikiya Yamashita<sup>2</sup>, Akinori Mitani<sup>2</sup>, Andre Esteva<sup>2</sup> <sup>1</sup>The University of Texas at Austin, <sup>2</sup> Artera Al



### MOTIVATION

- <u>Background</u>:
  - Prostate cancer is a common, deadly disease
     Diagnosis is multimodal: histopathology imaging + structured clinical risk factors
- <u>Problem</u>: Existing methods for fusing

## **EXPERIMENTAL SETUP**

- Train baseline *late joint fusion* approach that concatenates image and non-image features
- <u>Exp. #1</u>: Observe effect of auxiliary supervision on concatenation vs. Kronecker fusion
- <u>Exp. #2</u>: Isolate which auxiliary supervision

histopathology imaging + non-image data are
extremely expressive (ex: Kronecker fusion)
➢ Likely to overfit small/low-dimensional data

• <u>Question</u>: Can we develop a *parameter-efficient* method to learn from multimodal medical data?

## DATASET & TASK

- 4,581 patients from five phase III clinical trials w/ paired histopathology imaging + clinical data
  Kx128-dim "bag" of pre-extracted image features
  6-dim vector of numeric clinical features
  - Age, PSA, T-stage, Gleason scores

• <u>Goal</u>: Predict prostate cancer distant metastasis (DM)

## METHODS

methods improve upon the baseline (ablation)

 Use 5-fold cross-validation (CV) with identical base architecture and hyperparameters

• Evaluate with **mean AUC** across CV folds

## RESULTS

Fusion Operation	Auxiliary Supervision	# Params	AUC
Concatenation		17.1K	$0.781 \pm 0.024$
Kronecker		66.3K	$0.770\pm0.018$
Concatenation	$\checkmark$	43.1K	$\textbf{0.792} \pm \textbf{0.014}$
Kronecker	$\checkmark$	207.1K	$0.781\pm0.013$

#### Experiment #1:

• Auxiliary supervision improves performance

#### <u>Solution</u>: Use auxiliary sources of supervision



- Extra supervision: generate additional imageonly and clinical-only outcome predictions
   Minimize sum of cross-entropy (CE) losses
- Clinical prediction: use image-only features to predict/regress associate non-image inputs
   Minimize simple MSE loss

 Kronecker fusion increases parameter count by 4-5x and decreases performance

Extra Supervision	Clinical Prediction	<b>Dense Fusion</b>	AUC
			$0.781 \pm 0.024$
$\checkmark$			$0.778 \pm 0.021$
	$\checkmark$		$0.789 \pm 0.013$
		$\checkmark$	$0.776 \pm 0.025$
$\checkmark$	$\checkmark$		$0.787\pm0.015$
$\checkmark$		$\checkmark$	$0.785\pm0.016$
	$\checkmark$	$\checkmark$	$\underline{0.790 \pm 0.012}$
$\checkmark$	$\checkmark$	$\checkmark$	$\textbf{0.792} \pm \textbf{0.014}$

#### Experiment #2:

Combination of all three methods works best

Minimize simple MSE loss

- **Dense fusion**: encourage dense interaction of image-only, clinical-only, and fused features
- These approaches are all complementary
   Ex: a combination of all 3 methods can be optimized with 6 CE losses + 1 MSE loss



Gain in performance is not purely additive

Clinical prediction is the single most impactful
 > Highest mean AUCs and lowest std AUCs

PAPER

#### ACKNOWLEDGMENTS

This project was supported by U10CA180822 (NRG Oncology SDMC), U10CA180868 (NRG Oncology Operations), and U24CA196067 (NRG Specimen Bank) from the National Cancer Institute.

Special thanks to the entire AI team at Artera (<u>https://artera.ai/</u>)!