

Abstract 7441 Valentin Bernu<sup>1</sup>, Guillaume Appé<sup>1</sup>, Abdelkader Behdenna<sup>1</sup>, Antoine Gaston<sup>1</sup>, Julien Haziza<sup>1</sup>, Léa Meunier<sup>1</sup>, Akpéli Nordor<sup>1</sup>. <sup>1</sup>Epigene Labs, Paris, France

# Can a single generalist pipeline infer multiple clinical variables from transcriptomic data?

profiles.





![](_page_0_Figure_10.jpeg)

**Results for the test set** 

# Machine-learning-based inference of clinical metadata from gene expression data

	Test							
	Samples	Classes	Performance					
	3187	2	99% accuracy					
	4233	11	<b>81%</b> accuracy					
	5868	4	91% accuracy					
	3938	20	<b>89%</b> accuracy					
	3458	7	<b>94%</b> accuracy					
	440	2	49% accuracy (work in progress)					
_								

(e.g. 100 datasets for histological type)

### Comparison with existing cross-dataset models trained on TCGA

genes reflecting

sex effects in

profiles

Genes in

cytogenetic

band chrYq11

transcriptomic

• Most previous studies focus on predicting the tissue of origin, often in metastatic cancers with unknown origin, to improve diagnosis and outcomes.

	Papers	Metadata	<b>Cross-validation</b>			Test		
			Samples	Classes	Performance	Samples	Classes	Performance
	Van et al., 2024	Tissue type	7192	14	Undisclosed	876 (ICGC & 4 GEO datasets)	6	0.80 weighted F1-score
	Chen et al., 2021	Cancer type	7715	21	96.38% R2 score	42 (1 GEO dataset)	5	83.3% R2 score
	He et al., 2023	Cancer type	9911	32	97.50% accuracy	1988 (ICGC)	10	82.67% accuracy
,	7bao at al. 2020	Cancer type	18217 (+ ICGC)	32	98.54% accuracy	23 & 69 (JAX & Melbourne)	16 & 18	86.96% & 72.46% accuracy
nt	21140 et al., 2020	Molecular subtype	3367	11	4 classifiers 60% - 83.5% (acc.)	1784 & 215 (1 EGA & 1 GEO dataset)	4&4	84.19% ovarian subtype 79.88% breast subtype
	Elmahy et al., 2021	Stage	TCGA-KIRC	4	82% accuracy	_	-	-
	$\mathbf{\downarrow}$							
	To our knowledge, <b>no existing</b> approach performs pan-cancer inference across multiple clinical variables.			Effective within databases	Limited representativeness of test datasets compared to the TCGA distribution → Difficulty to properly evaluate cross-database performance			

Contact Akpéli Nordor, PharmD, PhD <u>akpeli@epigenelabs.com</u>

The authors have no conflict of interest to declare.

![](_page_0_Picture_27.jpeg)

![](_page_0_Picture_29.jpeg)

### Applications

- Clinical metadata quality control. • Labeling samples with missing
- clinical metadata.

### Limitations

The lower performance on rare visible in overall classes is not and accuracy can be accuracy for variables with fewer inflated classes.

## **Future Work**

- Explore batch effect correction methods without covariates or using predicted primary site as a covariate to improve proxy performance.
- Develop cancer-specific models, particularly to improve Stage inference where general models seem to underperform.
- Employ stricter criteria for gene selection to **enhance model** applicability.