



# Performance and Robustness Evaluation of Pathology Foundation Models

---

**Fredrik K. Gustafsson**

**University of Oxford**

Department of Engineering Science

[www.fregu856.com](http://www.fregu856.com)

April 29, 2026

Postdoc in the group of [David Clifton](#) at the University of Oxford, since July 2025.

Background:

- Dec 2023 - Jun 2025: Postdoc at Karolinska Institutet in Stockholm, ML/CV for *computational pathology*, in the group of [Mattias Rantalainen](#).
- 2018 - 2023: PhD in *Machine Learning*, Uppsala University.
  - Thesis: *Towards Accurate and Reliable Deep Regression Models*.  
Supervisors: [Thomas Schön](#) & [Martin Danelljan](#).
- 2013 - 2018: BSc & MSc in *Electrical Engineering*, Linköping University.
  - 2016 - 2017: Graduate exchange student, Stanford University.

*Machine learning for healthcare* with a particular focus on biosignals and wearables.

More broadly, I am interested in how to build and evaluate *reliable machine learning* models, for applications within *data-driven medicine and healthcare*.

I will briefly describe three papers where we evaluate the performance and *robustness* of a range of recent pathology foundation models:

## **I: Benchmarking Pathology Foundation Models for Breast Cancer Survival Prediction**

*Fredrik K. Gustafsson, Constance Boissin, Johan Vallon-Christersson, David A. Clifton, Mattias Rantalainen*  
Preprint, 2026

## **II: Evaluating Computational Pathology Foundation Models for Prostate Cancer Grading under Distribution Shifts**

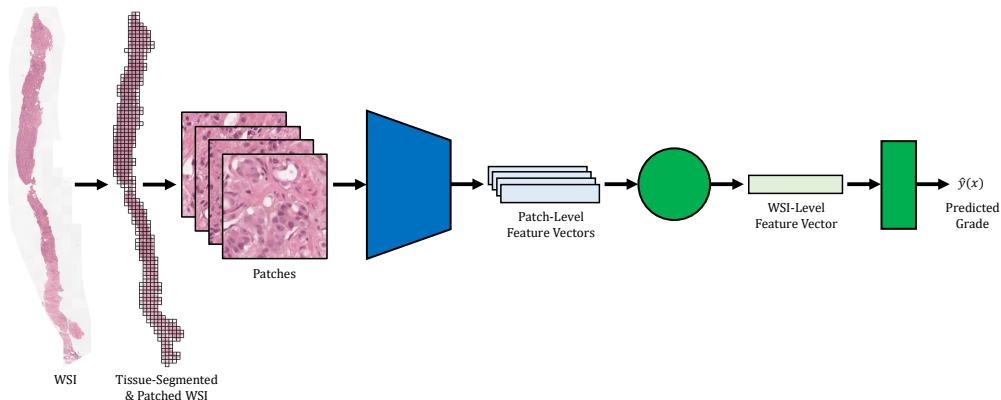
*Fredrik K. Gustafsson, Mattias Rantalainen*  
Preprint, 2026

## **III: Scanner-Induced Domain Shifts Undermine the Robustness of Pathology Foundation Models**

*Erik Thiringer, Fredrik K. Gustafsson, Kajsa Ledesma Eriksson, Mattias Rantalainen*  
Preprint, 2026

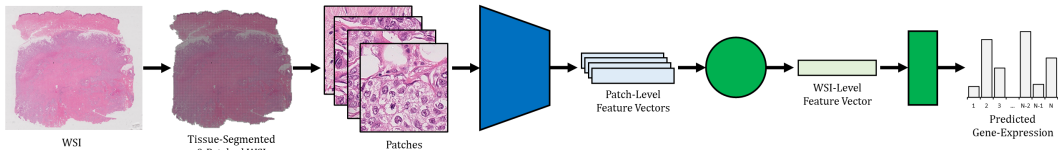
Computational pathology uses machine learning and computer vision to automatically extract useful information from histopathology whole-slide images (WSIs).

Given datasets of (WSI, label) pairs, models can be trained for applications such as *histological grading*, patient outcome prediction, and prediction of various biomarkers.



Computational pathology uses machine learning and computer vision to automatically extract useful information from histopathology whole-slide images (WSIs).

Given datasets of (WSI, label) pairs, models can be trained for applications such as histological grading, patient outcome prediction, and *prediction of various biomarkers*.



Foundation models are large models trained on *large amounts of unlabeled data* using *self-supervised learning*. They are intended to be general-purpose feature extractors.

Self-supervised learning enables models to be trained on “raw” unlabeled data. Large collections of unlabeled WSIs – *WSIs without known clinical info, patient outcomes or any other type of annotations* – can thus be directly utilized in model training.

Has recently become a popular research direction within computational pathology:

**UNI: Towards a General-Purpose Foundation Model for Computational Pathology**

Nature Medicine, 2024

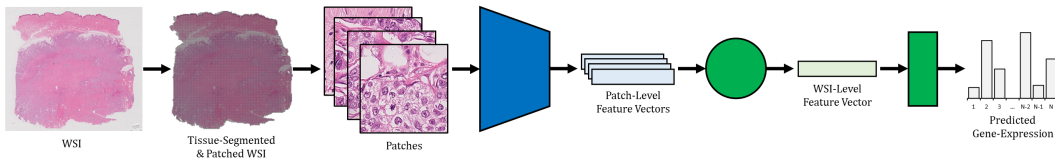
**Prov-GigaPath: A Whole-Slide Foundation Model for Digital Pathology from Real-World Data**

Nature, 2024

**Virchow: A Foundation Model for Clinical-Grade Computational Pathology and Rare Cancers Detection**

Nature Medicine, 2024

•  
•  
•



- Tissue-segment each WSI and divide it into image patches (e.g.  $224 \times 224$  pixels).
- Use a *frozen foundation model* to extract feature vectors for all images patches in each WSI (*typical range: 5,000 - 25,000 image patches per WSI*).
- Train a *small model* that, for each WSI, takes the extracted patch-level feature vectors as input and outputs a WSI-level prediction (standard supervised training).

## Benchmarking Pathology Foundation Models for Breast Cancer Survival Prediction

*Fredrik K. Gustafsson, Constance Boissin, Johan Vallon-Christersson, David A. Clifton, Mattias Rantalainen*

Preprint, 2026

Model Name	Architecture	Size	Feature Dimension	Pretraining Data
Resnet-IN [7]	ResNet-50	25M	1024	1.3M natural images
CTransPath [22]	CNN + Swin-T	22M	768	30K WSIs
RetCCL [23]	ResNet-50	25M	2048	32K WSIs
UNI [2]	ViT-L	307M	1024	100K WSIs
UNI2-h [13]	ViT-H	682M	1536	350K WSIs
H-optimus-0 [18]	ViT-G	1.1B	1536	500K WSIs
H-optimus-1 [1]	ViT-G	1.1B	1536	1M WSIs
Prov-GigaPath [25]	ViT-G	1.1B	1536	170K WSIs
Virchow [21]	ViT-H	632M	2560	1.5M WSIs
Virchow2 [26]	ViT-H	632M	2560	3.1M WSIs
H0-mini [5]	ViT-B	86M	768	500K + 6K WSIs
CONCH [11]	ViT-B	86M	512	21K WSIs + 1.1M image-text pairs
CONCHv1.5 [3]	ViT-L	307M	768	N/A

Rank	Model Name	Model Ranks	Mean Model Rank (↓)
1	H-optimus-1	1,1,1,3	1.5
2	H0-mini	2,2,4,5	3.25
3	H-optimus-0	5,5,5,3	4.5
4	CONCHv1.5	8,8,2,1	4.75
5	UNI2-h	4,4,6,6	5
6	Virchow2	3,3,7,8	5.25
7	Virchow	10,10,2,2	6
8	CONCH	6,6,8,7	6.75
9	Prov-GigaPath	6,6,9,9	7.5
10	UNI	9,9,10,10	9.5
11	CTransPath	12,12,11,11	11.5
12	RetCCL	11,11,13,13	12
13	Resnet-IN	13,13,12,12	12.5

We evaluate 13 models: a natural-image baseline, two early pathology-specific models, seven state-of-the-art PFMs, a compact distilled PFM, and two vision-language PFMs.

Downstream survival models are trained on a dataset of *2,315 breast cancer patients* (SöS-BC-4) and evaluated on two independent external datasets (KS-Solna and SCAN-B-Lund) comprising *3,119 patients* in total.

Models are evaluated under four settings: RFS and PFS, each assessed both for the full cohort ('All Patients') and for the clinically relevant 'ER+ & HER2-' patient subgroup.

The combined evaluation set of 3,119 patients contains 615 RFS events and 233 PFS events, with 2,524 patients (80.9% of the full set), 475 RFS events (77.2%) and 157 PFS events (67.4%) in the 'ER+ & HER2-' patient subgroup.

# I: Benchmarking PFMs for Survival Prediction - Main Results



Rank	Model Name	C-index ( $\uparrow$ )
1	H-optimus-1	0.678 (0.656 – 0.700)
2	H0-mini	0.676 (0.653 – 0.698)
3	Virchow2	0.675 (0.653 – 0.698)
4	UNI2-h	0.667 (0.644 – 0.689)
5	H-optimus-0	0.664 (0.642 – 0.686)
6	CONCH	0.663 (0.641 – 0.685)
6	Prov-GigaPath	0.663 (0.640 – 0.686)
8	CONCHv1.5	0.660 (0.637 – 0.683)
9	UNI	0.657 (0.635 – 0.680)
10	Virchow	0.656 (0.633 – 0.679)
11	RetCCL	0.645 (0.622 – 0.668)
12	CTransPath	0.632 (0.609 – 0.655)
13	Resnet-IN	0.612 (0.588 – 0.635)

(a) RFS – All Patients.

Rank	Model Name	C-index ( $\uparrow$ )
1	H-optimus-1	0.702 (0.666 – 0.737)
2	Virchow	0.695 (0.660 – 0.729)
2	CONCHv1.5	0.695 (0.660 – 0.729)
4	H0-mini	0.692 (0.655 – 0.728)
5	H-optimus-0	0.686 (0.651 – 0.722)
6	UNI2-h	0.683 (0.647 – 0.719)
7	Virchow2	0.681 (0.645 – 0.718)
8	CONCH	0.675 (0.639 – 0.710)
9	Prov-GigaPath	0.673 (0.636 – 0.710)
10	UNI	0.671 (0.634 – 0.708)
11	CTransPath	0.647 (0.609 – 0.686)
12	Resnet-IN	0.628 (0.590 – 0.666)
13	RetCCL	0.627 (0.589 – 0.666)

(c) PFS – All Patients.

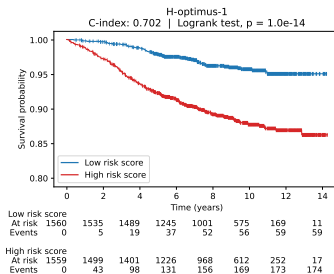
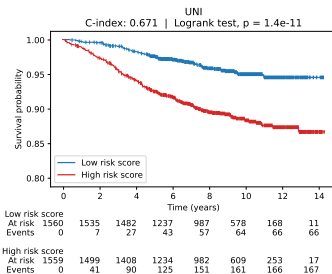
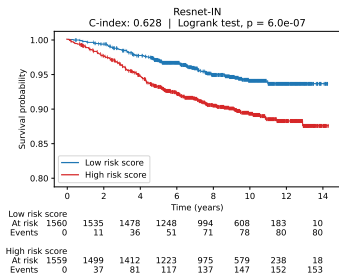
Rank	Model Name	C-index ( $\uparrow$ )
1	H-optimus-1	0.670 (0.645 – 0.695)
2	H0-mini	0.668 (0.642 – 0.693)
3	Virchow2	0.665 (0.639 – 0.690)
4	UNI2-h	0.663 (0.637 – 0.688)
5	H-optimus-0	0.660 (0.633 – 0.685)
6	CONCH	0.657 (0.631 – 0.681)
6	Prov-GigaPath	0.657 (0.631 – 0.682)
8	CONCHv1.5	0.651 (0.625 – 0.676)
9	UNI	0.648 (0.622 – 0.674)
10	Virchow	0.639 (0.613 – 0.665)
11	RetCCL	0.636 (0.609 – 0.662)
12	CTransPath	0.627 (0.601 – 0.652)
13	Resnet-IN	0.600 (0.572 – 0.627)

(b) RFS – Patient Subgroup (ER+ & HER2-).

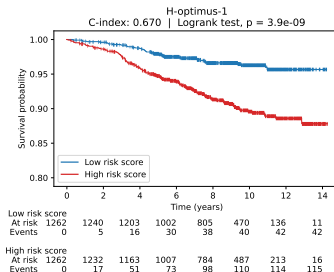
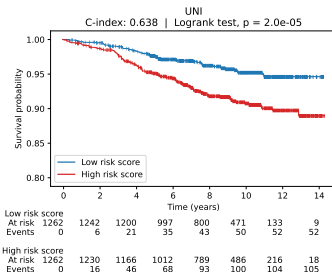
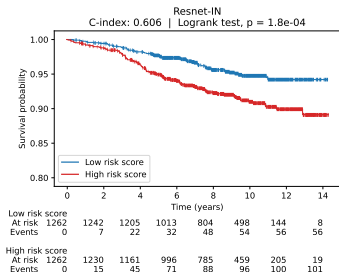
Rank	Model Name	C-index ( $\uparrow$ )
1	CONCHv1.5	0.680 (0.636 – 0.721)
2	Virchow	0.678 (0.634 – 0.720)
3	H-optimus-1	0.670 (0.626 – 0.715)
3	H-optimus-0	0.670 (0.624 – 0.714)
5	H0-mini	0.666 (0.619 – 0.711)
6	UNI2-h	0.657 (0.611 – 0.702)
7	CONCH	0.656 (0.612 – 0.700)
8	Virchow2	0.655 (0.608 – 0.700)
9	Prov-GigaPath	0.641 (0.594 – 0.687)
10	UNI	0.638 (0.592 – 0.683)
11	CTransPath	0.622 (0.573 – 0.670)
12	Resnet-IN	0.606 (0.557 – 0.655)
13	RetCCL	0.604 (0.555 – 0.653)

(d) PFS – Patient Subgroup (ER+ & HER2-).

# I: Benchmarking PFMs for Survival Prediction - Risk Stratification

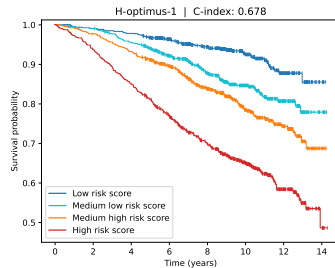
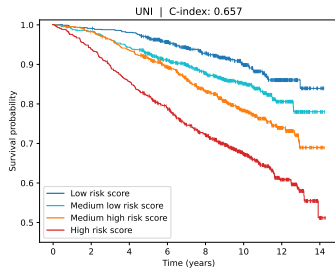
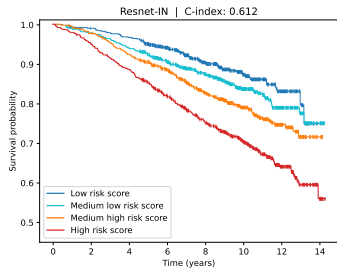


(c) PFS – All Patients.

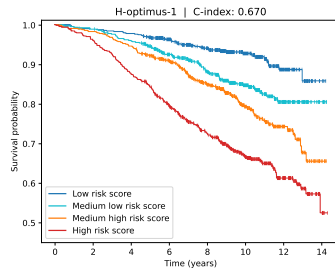
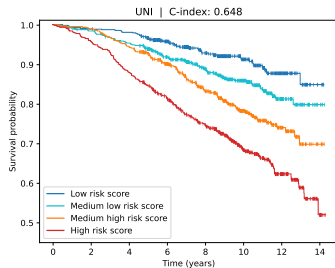
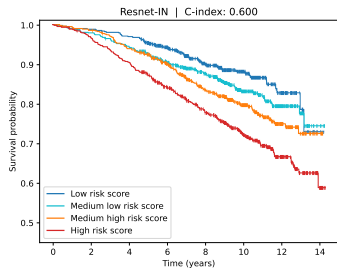


(d) PFS – Patient Subgroup (ER+ & HER2-).

# I: Benchmarking PFMs for Survival Prediction - Risk Stratification

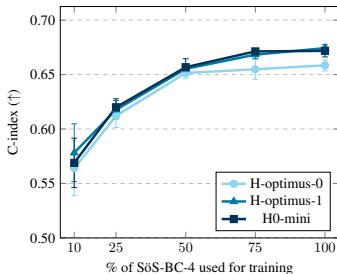


(a) RFS – All Patients.

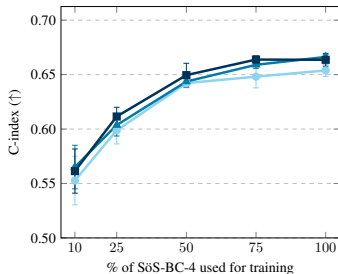


(b) RFS – Patient Subgroup (ER+ & HER2-).

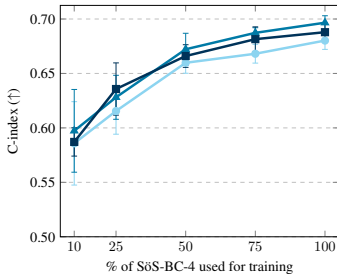
# I: Benchmarking PFMs for Survival Prediction - Detailed Comparison



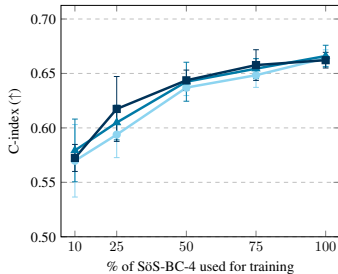
(a) RFS - All Patients.



(b) RFS - Patient Subgroup (ER+ & HER2-).

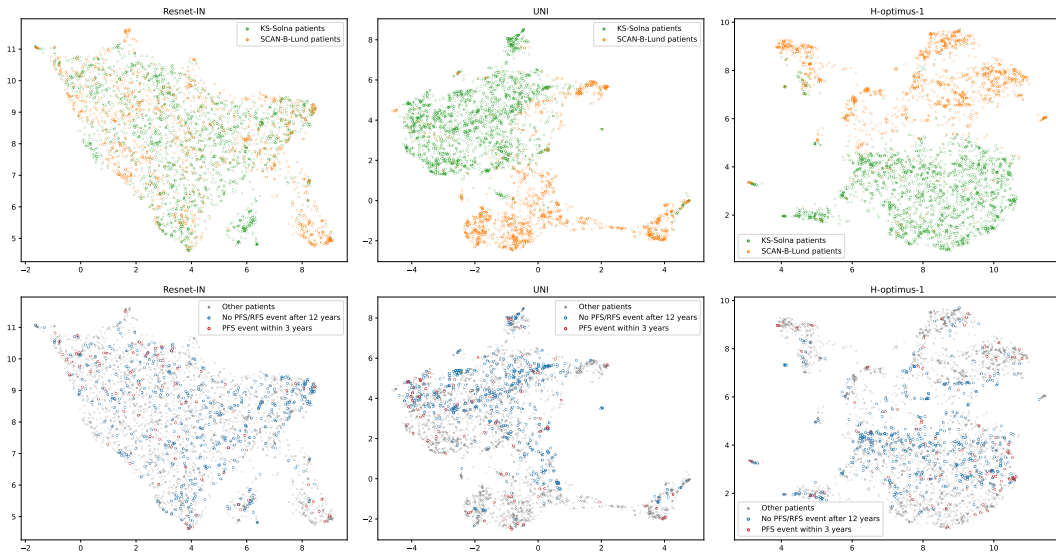


(c) PFS - All Patients.



(d) PFS - Patient Subgroup (ER+ & HER2-).

# I: Benchmarking PFMs for Survival Prediction - Feature Space Analysis





**(1/4)** H-optimus-1 achieves the strongest overall performance, but absolute differences between top-performing PFMs are small and confidence intervals substantially overlap.

**(2/4)** Across model families, consistent generational improvements are observed, with second-generation PFMs (H-optimus-1, CONCHv1.5, UNI2-h, Virchow2) outperforming their first-generation counterparts. However, these gains remain modest despite substantial increases in pretraining data scale, suggesting diminishing returns from scaling alone.

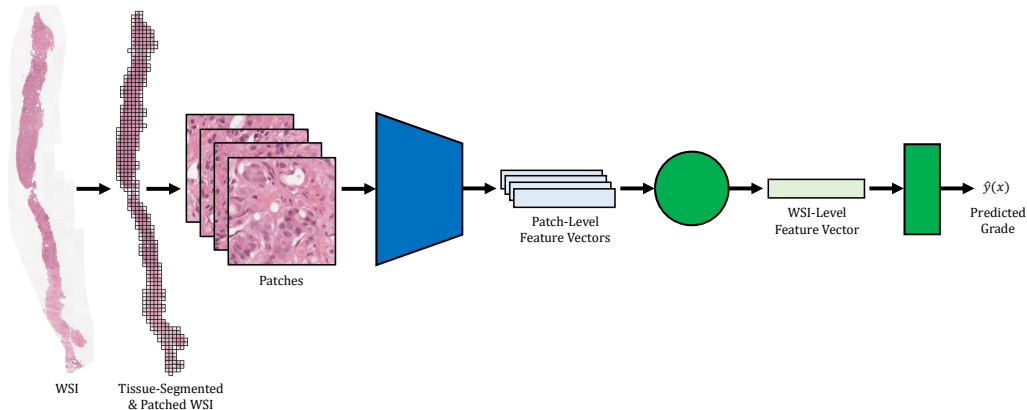
**(3/4)** Model size is not a reliable predictor of performance, as smaller and more efficient models can match or exceed much larger architectures, emphasizing the importance of training strategy and pretraining data quality over model scaling.

**(4/4)** The distilled H0-mini achieves the second-best overall ranking and slightly outperforms its teacher model H-optimus-0, while using less than 8% of the parameters and enabling much faster feature extraction. Knowledge distillation can thus yield highly efficient PFMs without sacrificing performance, making it a particularly promising approach for practical deployment.

## Evaluating Computational Pathology Foundation Models for Prostate Cancer Grading under Distribution Shifts

*Fredrik K. Gustafsson, Mattias Rantalainen*

Preprint, 2026



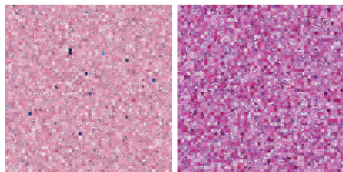
## II: Evaluating PFMs under Distribution Shifts



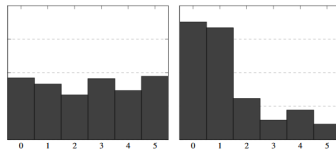
The PANDA dataset was collected from two different sites: *Radboud* University Medical Center in the Netherlands, and *Karolinska* Institutet in Sweden.

Radboud and Karolinska differ in terms of both the pathology lab procedures and utilized scanners, creating a clear distribution shift for the WSI image data.

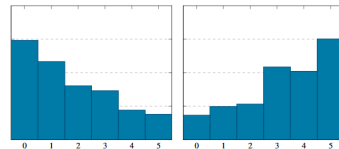
By creating further subsets of the PANDA dataset, we are also able to evaluate robustness in terms of shifts in the label distribution over the ISUP grades 0 - 5.



(a) WSI image data shift, *Radboud*  $\rightarrow$  *Karolinska*.

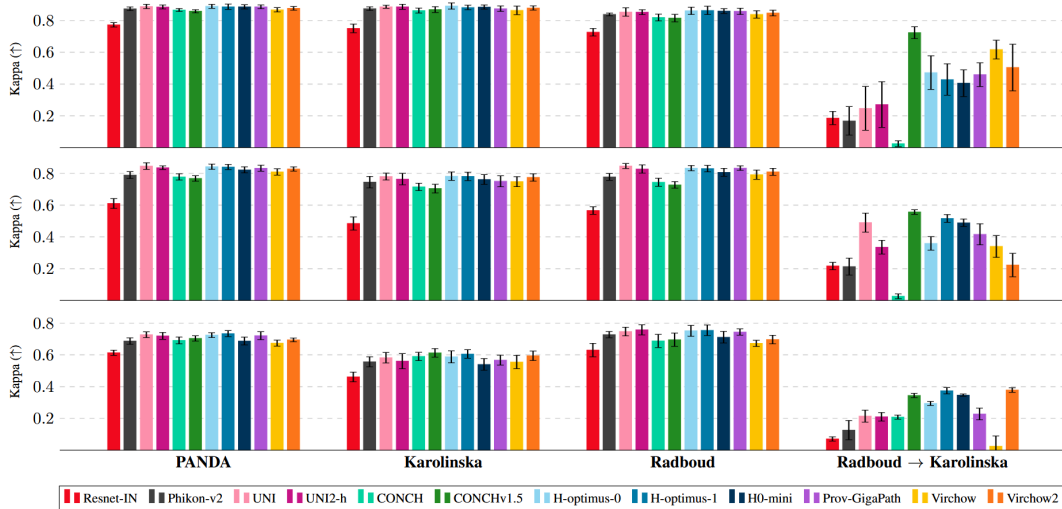


(b) Grade label shift, *Radboud*  $\rightarrow$  *Karolinska*.

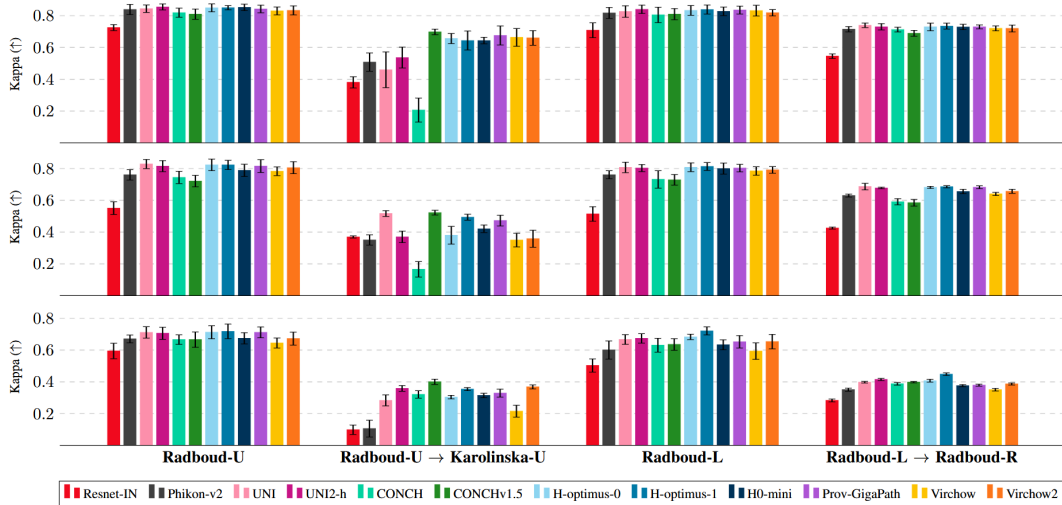


(c) Grade label shift, *Radboud-L*  $\rightarrow$  *Radboud-R*.

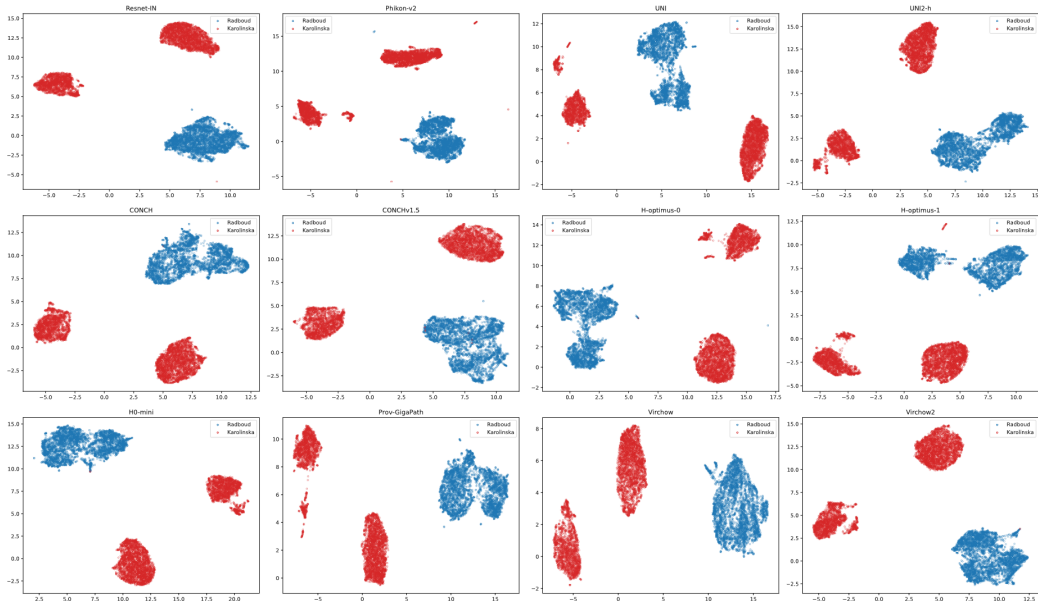
## II: Evaluating PFMs under Distribution Shifts - Results



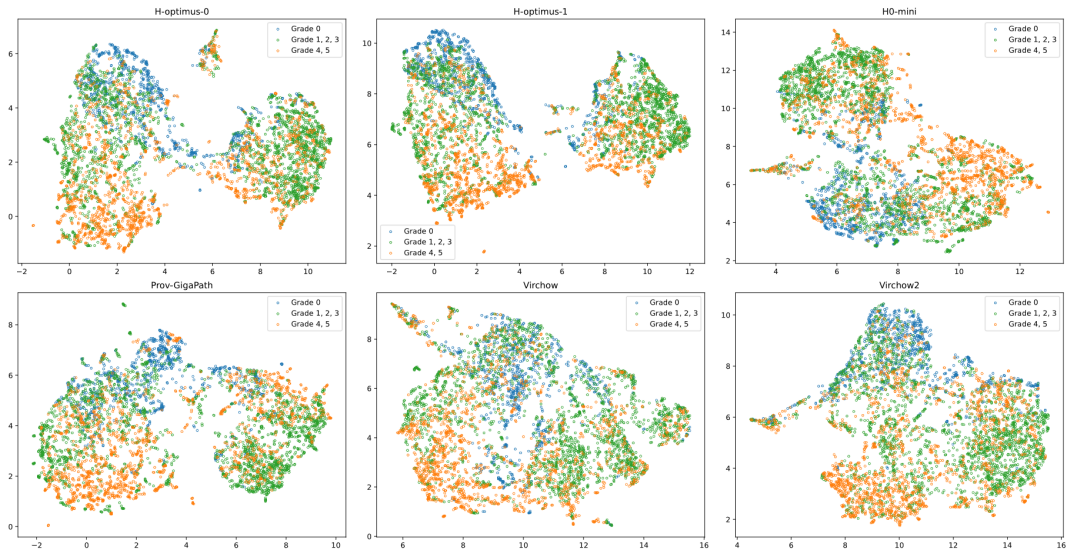
## II: Evaluating PFM under Distribution Shifts - Results



## II: Evaluating PFM's under Distribution Shifts - Feature Space Analysis



## II: Evaluating PFMs under Distribution Shifts - Feature Space Analysis





**(1/5)** PFMs substantially outperform a natural-image baseline for prostate cancer grading, but absolute performance under cross-site shifts can still deteriorate severely.

**(2/5)** Large-scale pretraining of PFMs on diverse datasets does not guarantee downstream cross-site generalization, and increasing model size or pretraining data alone does not reliably improve robustness.

**(3/5)** PFMs are significantly more sensitive to image-domain shifts than to label distribution shifts; acquisition-related variability seems to be the dominant challenge.

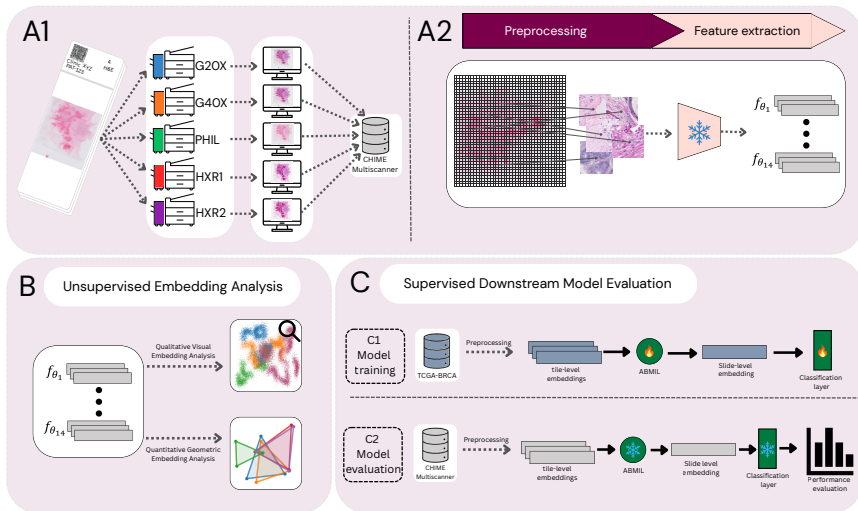
**(4/5)** Feature space analysis reveals persistent domain separation, suggesting that current PFMs encode site- and scanner-specific variation more strongly than pathology-relevant signals.

**(5/5)** Strong PFMs and high-quality downstream data complement rather than replace each other: even in the foundation-model era, the quality and diversity of the data used to train downstream prediction models remains critical for reliable deployment.

## Scanner-Induced Domain Shifts Undermine the Robustness of Pathology Foundation Models

Erik Thiringer, Fredrik K. Gustafsson, Kajsa Ledesma Eriksson, Mattias Rantalainen

Preprint, 2026



### III: PFM Scanner-Variability Robustness



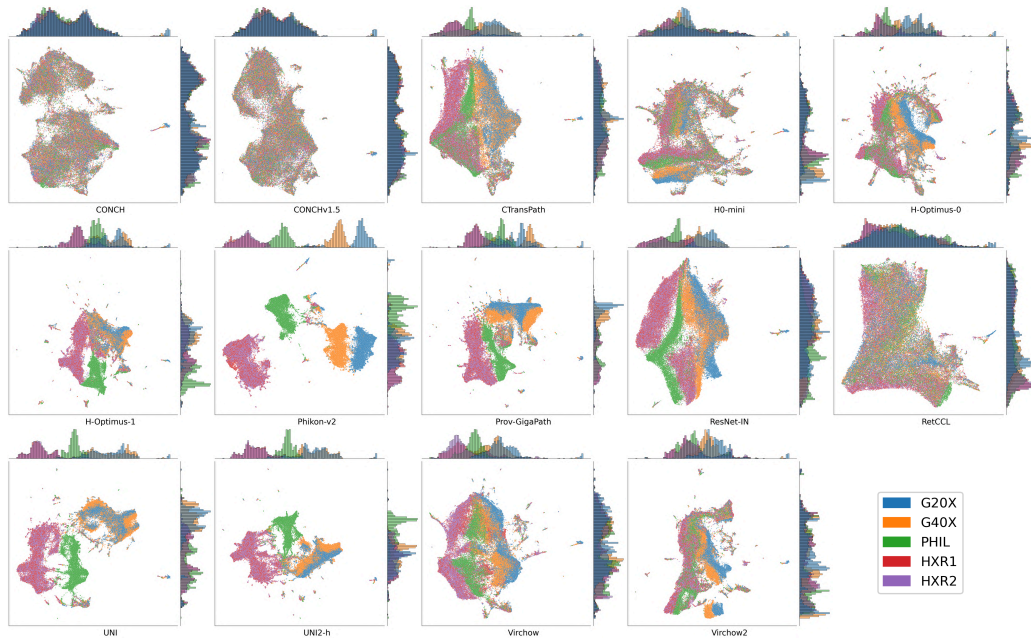
We systematically evaluate the robustness of 13 PFMs to scanner-induced variability.

Using a controlled multiscanner dataset comprising 384 breast cancer WSIs from as many patients, *scanned on five different whole-slide scanners*, we isolate scanner-induced effects from all other sources of variation.

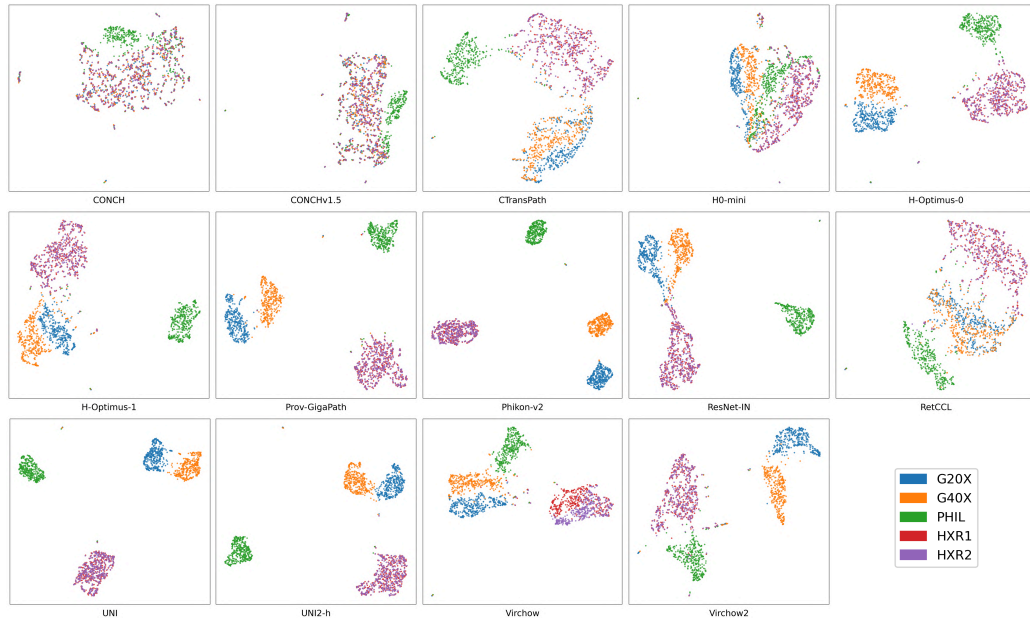
Robustness is assessed through complementary unsupervised analyses of the feature embedding space and a set of clinicopathological supervised prediction tasks.



# III: PFM Scanner-Variability Robustness - Patch-level Feature Space



### III: PFM Scanner-Variability Robustness - Slide-level Feature Space



# III: PFM Scanner-Variability Robustness - Quantitative Feature Analysis



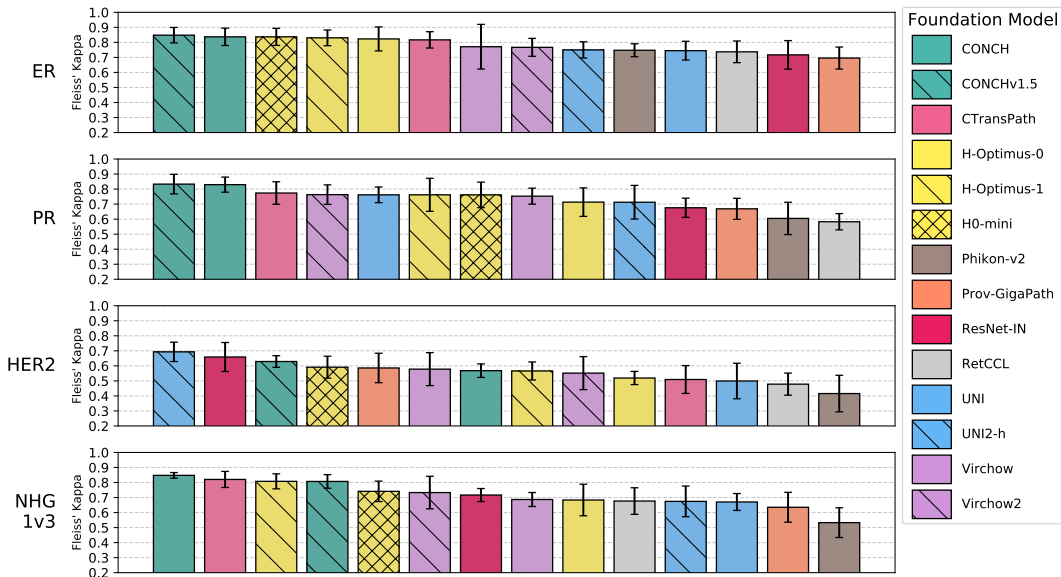
**A**

ResNet-IN	0.015	0.048	0.021	0.021	0.029	0.017	0.016	0.027	0.024	0.001
CONChv1.5	0.022	0.066	0.035	0.035	0.045	0.024	0.024	0.031	0.030	0.002
CONCH	0.034	0.092	0.050	0.049	0.057	0.027	0.025	0.042	0.038	0.002
CTransPath	0.055	0.152	0.092	0.099	0.077	0.081	0.074	0.085	0.081	0.004
RetCCL	0.056	0.141	0.099	0.095	0.070	0.085	0.078	0.086	0.080	0.004
Virchow	0.058	0.132	0.135	0.192	0.064	0.103	0.158	0.069	0.092	0.040
HD-mini	0.093	0.193	0.173	0.165	0.086	0.160	0.146	0.120	0.109	0.010
Virchow2	0.125	0.269	0.221	0.223	0.136	0.194	0.184	0.124	0.116	0.019
H-Optimus-1	0.125	0.387	0.246	0.242	0.257	0.194	0.182	0.259	0.257	0.013
H-Optimus-0	0.128	0.353	0.325	0.312	0.209	0.261	0.241	0.225	0.215	0.015
Phikon-v2	0.173	0.456	0.326	0.341	0.264	0.284	0.292	0.310	0.289	0.055
Prov-GigaPath	0.173	0.607	0.379	0.366	0.402	0.301	0.283	0.362	0.356	0.014
UNI	0.180	0.598	0.353	0.379	0.435	0.331	0.328	0.394	0.368	0.077
UNI2-h	0.202	0.601	0.343	0.332	0.404	0.394	0.374	0.469	0.444	0.035
	G20X & G40X	G20X & PHIL	G20X & HXR1	G20X & HXR2	G40X & PHIL	G40X & HXR1	G40X & HXR2	PHIL & HXR1	PHIL & HXR2	HXR1 & HXR2

**B**

CONChv1.5	372 (96.9%)	314 (81.8%)	367 (95.6%)	364 (94.8%)	346 (90.1%)	376 (97.9%)	371 (96.6%)	335 (87.2%)	336 (87.5%)	384 (100.0%)
CONCH	372 (96.9%)	292 (76.0%)	363 (94.5%)	363 (94.5%)	335 (87.2%)	375 (97.7%)	376 (97.9%)	330 (85.9%)	341 (88.8%)	384 (100.0%)
UNI	355 (92.4%)	288 (75.0%)	297 (77.3%)	284 (74.0%)	351 (91.4%)	328 (85.4%)	347 (90.4%)	310 (80.7%)	329 (85.7%)	384 (100.0%)
HD-mini	354 (92.2%)	301 (78.4%)	308 (80.2%)	308 (80.2%)	366 (95.3%)	265 (69.0%)	274 (71.4%)	273 (71.1%)	275 (71.6%)	383 (99.7%)
H-Optimus-1	356 (92.7%)	223 (58.1%)	345 (89.8%)	344 (89.6%)	252 (65.6%)	352 (91.7%)	358 (93.2%)	239 (62.2%)	218 (56.8%)	384 (100.0%)
H-Optimus-0	362 (94.3%)	258 (67.2%)	278 (72.4%)	279 (72.7%)	312 (81.2%)	273 (71.1%)	288 (75.0%)	283 (73.7%)	292 (76.0%)	383 (99.7%)
CTransPath	351 (91.4%)	236 (66.2%)	317 (82.6%)	326 (84.9%)	319 (83.1%)	291 (75.8%)	303 (78.9%)	202 (52.6%)	213 (55.5%)	384 (100.0%)
UNI2-h	353 (91.9%)	306 (79.7%)	275 (71.6%)	260 (67.7%)	345 (89.8%)	220 (57.3%)	249 (64.6%)	181 (47.1%)	213 (55.5%)	383 (99.7%)
Prov-GigaPath	305 (79.4%)	192 (50.0%)	143 (37.2%)	147 (38.3%)	300 (78.1%)	291 (75.8%)	293 (76.3%)	292 (76.0%)	289 (75.3%)	384 (100.0%)
Virchow2	322 (83.9%)	228 (59.4%)	250 (65.1%)	248 (64.6%)	268 (69.8%)	189 (49.2%)	215 (56.0%)	203 (52.9%)	236 (61.5%)	382 (99.5%)
RetCCL	336 (87.5%)	147 (38.2%)	196 (51.0%)	206 (53.6%)	235 (61.2%)	219 (57.0%)	220 (57.3%)	180 (46.9%)	190 (49.5%)	382 (99.5%)
Virchow	335 (87.2%)	150 (39.1%)	132 (34.4%)	75 (19.5%)	211 (54.9%)	131 (34.1%)	77 (20.1%)	221 (57.6%)	127 (33.1%)	356 (92.7%)
ResNet-IN	280 (72.9%)	56 (14.6%)	205 (53.4%)	217 (56.5%)	71 (18.5%)	196 (51.0%)	214 (55.7%)	58 (15.1%)	80 (20.8%)	383 (99.7%)
Phikon-v2	249 (64.0%)	99 (25.8%)	170 (44.3%)	154 (40.1%)	192 (50.0%)	147 (38.3%)	180 (46.9%)	81 (21.1%)	89 (23.2%)	384 (100.0%)
	G20X & G40X	G20X & PHIL	G20X & HXR1	G20X & HXR2	G40X & PHIL	G40X & HXR1	G40X & HXR2	PHIL & HXR1	PHIL & HXR2	HXR1 & HXR2

### III: PFM Scanner-Variability Robustness - Prediction Consistency



### III: PFM Scanner-Variability Robustness - Main Takeaways



**(1/5)** Current state-of-the-art PFMs are not invariant to scanner-induced domain shifts. Most models encode pronounced scanner-specific variability in their feature embedding space.

**(2/5)** Although predictive performance largely remains similar when measured by AUC, this masks an important failure mode: scanner variability systematically alters the embedding space and impacts calibration of downstream model predictions, resulting in scanner-dependent bias.

**(3/5)** Robustness is not a simple function of training data scale, model size, or model recency.

**(4/5)** The models trained on the most diverse data, vision-language models, appear to have an advantage with respect to robustness. More systematic analysis is however needed, comparing vision-only and multimodal pretraining under controlled settings.

**(5/5)** Targeted robustness-oriented training strategies appear promising, as the distilled model H0-mini consistently outperforms its larger teacher model H-optimus-0 across multiple aspects of embedding stability and downstream prediction consistency.

**Fredrik K. Gustafsson**

University of Oxford

*fredrik.gustafsson@eng.ox.ac.uk*

[www.fregu856.com](http://www.fregu856.com)