

Deep Learning-based Multi-Omics Analysis Pipeline for Breast Cancer Prognosis and Subtype Classification

Yusuf Hakan Usta^{1*}, Stephen Richardson¹, Hamish Gilbert²,

¹The University of Manchester, Manchester, UK, Division of Cell Matrix Biology and Regenerative Medicine

²Keele University, Keele, UK, School of Life Science

Abstract

Background: Breast cancer is a heterogeneous disease comprising multiple intrinsic subtypes with distinct clinical outcomes. Prognostication in breast cancer requires reconciling morphological context with molecular state. We propose a multimodal deep learning pipeline that integrates whole-slide histopathology with multi-omics data for outcome prediction. Whole-slide images (WSIs) are encoded by a ResNet backbone with attention-based Multiple Instance Learning (MIL), and omics features by a Transformer encoder; the resulting patient embeddings are refined over a k-Nearest-Neighbor patient graph via a Graph Neural Network (GNN), and a Cox proportional hazards layer outputs a risk score. **Results:** The integrated model achieved accurate subtype classification (overall accuracy $\sim 93\%$ on independent test sets, outperforming gene-expression-only models) and improved survival prediction (C-index 0.70 in TCGA, 0.72 in METABRIC) compared to single-modality baselines. Fusion of histopathology with molecular data significantly enhanced risk stratification, particularly in intermediate clinical risk cases. Attention-based MIL identified histological patterns (e.g. necrosis, lymphocytic infiltration) associated with poor outcomes, while the VAE-derived genomic features captured key pathways (e.g. cell cycle, immune response) linked to survival. SHAP analysis highlighted the importance of genomic aberrations (TP53 mutation, high proliferation gene expression), DNA methylation changes, and clinical factors (tumor stage) in driving model predictions. Graph-based modeling of spatial transcriptomics further enabled identification of tumor microenvironment features (such as spatially organized immune cell clusters) that correlated with patient prognosis. Kaplan–Meier analyses based on the model’s risk scores showed clear separation between low-risk and high-risk groups (log-rank $p < 0.001$).

Conclusions: We demonstrate a comprehensive deep learning pipeline for breast cancer outcome prediction that integrates multi-omic and histopathology data. The proposed approach leverages state-of-the-art neural network blocks (VAEs, ResNets, Transformers, attention MIL, GNNs) to capture complementary information from diverse data modalities. Our results underscore that multi-omics integration yields more robust prognostic stratification than individual data types, reflecting the complex biology of breast cancer. The inclusion of spatial and single-cell modalities illustrates the extensibility of the framework to emerging data sources. Furthermore, the model’s explainability techniques provide insights into the biological underpinnings of its predictions, potentially guiding biomarker discovery. This multi-modal pipeline, suitable for bioRxiv submission and open-source release, offers a foundation for improved precision oncology in breast cancer and can be adapted to other cancers and multi-omic datasets.

*Corresponding author: Yusuf Hakan Usta, The University of Manchester, UK.
yusufhakan.usta@manchester.ac.uk

1 Introduction

Breast cancer is a leading cause of cancer-related mortality among women worldwide. It is a highly heterogeneous disease comprising multiple molecular subtypes with distinct therapeutic responses and outcomes. The classic intrinsic subtypes (Luminal A, Luminal B, HER2-enriched, Basal-like, and Normal-like) were originally defined by gene expression profiling and the PAM50 gene panel Perou2000. These subtypes have prognostic relevance—e.g., basal-like (triple-negative) tumors generally have worse outcomes, whereas Luminal A tumors have the most favorable prognosis Perou2000. However, even within these subtypes, considerable inter-patient variability in outcomes is observed Sharma2024. This suggests that additional molecular and microenvironmental factors contribute to prognosis beyond the intrinsic subtype classification.

In recent years, comprehensive multi-omics approaches have been applied to better stratify breast cancers. By integrating data types such as genomics, transcriptomics, epigenomics, and proteomics, researchers have identified novel prognostic subgroups that are not captured by single-omic analyses Sharma2024. For example, Sharma *et al.* (2024) used Multi-Omics Factor Analysis to define breast tumor clusters with distinct long-term survival rates, outperforming the prognostic value of standard subtypes Sharma2024. These findings underscore the potential of multi-omics integration to reveal underlying biological heterogeneity and improve risk prediction.

At the same time, digital pathology has emerged as an important modality for prognostication. Whole-slide images (WSIs) of hematoxylin and eosin (H&E)-stained tissue contain rich spatial and morphological information about tumor architecture, stromal context, and immune infiltration. Histopathological features can complement molecular data; for instance, the presence of tumor-infiltrating lymphocytes or necrosis on WSIs can carry prognostic significance independent of genomic markers. Traditional pathology assessment is qualitative, but deep learning now enables quantitative extraction of features from WSIs that correlate with outcomes. In fact, histology-based models using deep convolutional neural networks have shown promise in predicting survival directly from images Yao2020. Nonetheless, relying solely on histopathology may miss critical molecular alterations, and likewise purely genomic models ignore spatial context. This has motivated multimodal approaches that combine pathology (“pathomics”) and genomics (“genomics”) data for improved prognostic modeling Chen2020,Luo2025.

Deep learning provides a powerful toolkit to tackle the challenges of multimodal data integration in cancer. Recent advances in neural network architectures allow for extracting complex patterns from each data type and learning joint representations. Notably, autoencoders and variational autoencoders (VAE) can distill high-dimensional omics data (e.g. tens of thousands of gene expression values) into a low-dimensional latent representation while preserving important variation Kingma2014,Poirion2021. For images, convolutional neural networks such as ResNet He2016 pretrained on large datasets can serve as effective feature extractors for pathology images. Transformers and attention mechanisms have also been adopted in pathology and multimodal learning: vision transformers (ViT) can model global context in images by treating image patches as sequence tokens Dosovitskiy2021, and attention-based fusion models learn interactions between modalities by highlighting salient features across data types Chen2021,Luo2025. Multiple instance learning (MIL) with attention Ilse2018 has become a standard approach for WSI analysis: a WSI is considered as a “bag” of many patch instances, and an attention mechanism identifies which patches are most relevant to the outcome, enabling weakly-supervised training without per-patch annotations Yao2020. Additionally, graph neural networks (GNNs) have shown utility in biomedical applications for modeling relational structures. In cancer, GNNs can represent tissue as graphs (nodes as cells/regions, edges as spatial or functional relationships) or model gene–gene interaction networks, capturing topology that traditional neural nets might miss Kipf2017.

Despite these advances, building an integrated prognostic model that fully exploits multi-omics and imaging data remains challenging. Data heterogeneity between modalities (e.g. discrete mutations vs. continuous expression, or pixel data vs. molecular data) complicates straightforward fusion Luo2025. Moreover, missing data is common (not all patients have all assays available), and sample sizes for multi-omic cohorts are limited, increasing overfitting risk for very high-dimensional models. There is also a need for interpretability: clinicians are more likely to trust and adopt models that provide human-interpretable explanations (e.g. key genes or image regions driving a prediction) Lundberg2017.

In this study, we present a deep learning-based analysis pipeline that integrates diverse data modalities for breast cancer prognosis and subtype classification. Our approach combines (1) clinical variables, (2) bulk transcriptomic data (mRNA expression), (3) genomic data (somatic mutations and copy number alterations), (4) DNA methylation profiles, (5) whole-slide histopathology images, and also incorporates emerging modalities: (6) spatial transcriptomics and (7) single-cell RNA sequencing data for a subset of tumors. To our knowledge, this represents one of the most comprehensive integrations of multi-modal data in a single prognostic modeling framework for breast cancer. The pipeline employs specialized deep learning components for each data type: VAEs for dimensionality reduction of omics, ResNet and Transformer-based encoders for images, an attention MIL module for WSI analysis, and a GNN for spatial relationships. These components are then fused via a multi-head attention mechanism that learns cross-modal interactions, inspired by recent multimodal transformer models such as MCAT (Multimodal Co-Attention Transformer) Chen2021. We use a Cox proportional hazards deep learning objective to directly model time-to-event survival data, and we also perform multi-class classification to predict intrinsic molecular subtypes, illustrating a multi-task learning paradigm.

We evaluate our approach on two large breast cancer cohorts: TCGA-BRCA (The Cancer Genome Atlas Breast Cancer, $n \approx 1000$) and METABRIC (Molecular Taxonomy of Breast Cancer International Consortium, $n \approx 1900$), which both have extensive molecular profiling and clinical follow-up. Results show that our integrated model improves prognostic accuracy (as measured by concordance index and Kaplan–Meier stratification) compared to models using single data types alone. It also classifies PAM50 subtypes with high accuracy, potentially refining subtype assignment by incorporating additional data layers. We further demonstrate the model’s generalizability by testing on independent validation sets and by an example application to a spatial transcriptomics dataset. Finally, we provide interpretability analyses using SHAP values and attention visualizations, which highlight biologically meaningful features learned by the model (for example, specific gene signatures and histological patterns associated with high risk).

Our work contributes a novel framework for multimodal deep learning in oncology and underscores the value of integrating histopathology with multi-omics for prognostic predictions. The remainder of this manuscript details the methods, presents the model performance and findings, and discusses the implications in the context of current research.

2 Methods

2.1 Data Acquisition and Preprocessing

Patient Cohorts: We analyzed two primary retrospective cohorts: TCGA-BRCA and METABRIC. TCGA-BRCA comprises breast tumor samples profiled by multiple omics technologies, along with digitized pathology slides and clinical outcomes (overall survival). We included $N = 959$ patients from TCGA with available RNA sequencing data, DNA methylation (450k arrays), somatic mutation calls, and diagnostic H&E whole-slide images. The METABRIC dataset includes $N = 1,981$

breast cancer cases with gene expression (microarray), somatic copy number, and clinical data Curtis2012. METABRIC does not have whole-slide images publicly available; thus, for imaging analysis, our primary model was trained on TCGA (where both molecular and imaging data are present) and validated on METABRIC using molecular data only. We also incorporated external datasets for model extension: a spatial transcriptomics dataset of breast cancer (10x Genomics Visium data for a triple-negative breast tumor, from Rediti *et al.* Rediti2024) and a single-cell RNA-seq dataset of breast cancer (23 samples of ER-positive primary and metastatic tumors from Ozmen *et al.* Ozmen2025). These were used to demonstrate how the pipeline can integrate emerging data types, as described below.

Clinical and Subtype Data: Clinical variables such as patient age at diagnosis, tumor stage, grade, and hormone receptor status were obtained from the corresponding publications or data portals for each cohort. For subtype classification, we used PAM50 intrinsic subtype labels if available (for TCGA and METABRIC, PAM50 subtypes were either provided or we assigned them using published gene expression centroids Parker2009). These labels (Luminal A, Luminal B, HER2-enriched, Basal-like, Normal-like) were used as ground truth for the subtype classification task. Patients lacking an assigned subtype were excluded from the classification analysis.

Genomic and Transcriptomic Features: Gene expression data (TCGA RNA-seq counts and METABRIC microarray intensities) were \log_2 transformed and z-score normalized for each gene across patients. We filtered to protein-coding genes and applied variance filtering to retain the top variable genes (e.g., top 5,000 genes by variance in the training set) to reduce dimensionality. For somatic mutations, we created binary features for frequently mutated genes (e.g., an indicator for TP53 mutation, PIK3CA mutation, etc. for each patient) and counted mutations in key pathways. Copy number alterations in METABRIC were summarized as binary amplifications/deletions of known driver genes. DNA methylation beta values (TCGA 450k array) were averaged over promoter regions or represented by principal components capturing major methylation variability.

Prior to modeling, we employed an autoencoder-based compression on these omics features. In particular, separate variational autoencoders (VAEs) were trained on: (1) gene expression (high-dimensional continuous data), and (2) DNA methylation data. The VAEs had an encoder network that maps the input features to a 32-dimensional latent space (the “bottleneck” representing a compressed multi-gene feature vector), and a decoder that attempts to reconstruct the original input. We included a Kullback–Leibler divergence term in the VAE loss to ensure the latent space approximates a normal distribution Kingma2014. After training on the training cohort, we extracted the mean of the posterior latent distribution (denoted \mathbf{z}_{expr} for expression and \mathbf{z}_{meth} for methylation) for each patient; these serve as compact representations of the transcriptomic and epigenomic profiles. For somatic mutations, given their sparsity, we did not use a VAE; instead, we encoded mutation data via a simple multilayer perceptron (MLP) that outputs a smaller feature vector (e.g. 16-dimensional) from the binary mutation indicators, effectively learning a low-dimensional mutation signature for each patient.

Histopathology Images: For TCGA cases, we downloaded diagnostic H&E stained WSIs (typically at 20 \times magnification) from the Genomic Data Commons. Each WSI can be up to several gigapixels in size. We used a standard preprocessing pipeline for computational pathology: images were tiled into patches of 256 \times 256 pixels (at 20 \times) with 50% overlap, and background or mostly non-tissue patches were filtered out using a threshold on average brightness or tissue mask. This yielded on the order of 10^4 patches per WSI (exact number depending on tumor size on the slide). To obtain features from each image patch, we employed a deep convolutional neural network. Specifically, we used a ResNet-50 He2016 architecture pretrained on ImageNet, truncating before the final classification layer so that the output is a 2048-dimensional feature vector for each patch. These patch features were then ℓ^2 -normalized. We did not fine-tune the ResNet on our data at this

stage to avoid overfitting; the patch feature extraction was done using the fixed pretrained weights (which have been shown to be effective for histology feature transfer learning Ciga2022). The result of this process is that each WSI is represented as a bag of feature vectors $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_M\}$, where M is the number of tissue patches for that slide and $\mathbf{h}_i \in \mathbb{R}^{2048}$.

Spatial Transcriptomics and Single-Cell Data: In one experiment, we explored incorporating spatial transcriptomics (ST) data. The ST dataset Rediti2024 consists of tissue positions (spots) on a breast tumor slide, each with a gene expression profile. We constructed a graph where each spot is a node, features of the node are the normalized expression of select genes (we focused on a subset such as the top 20 principal components of the spatial gene expression to reduce dimensionality), and edges connect neighboring spots (within a certain Euclidean distance on the tissue grid). Similarly, for single-cell RNA-seq data Ozmen2025, we obtained, for each patient in that study, the proportions of various cell types (e.g., percent of T cells, macrophages, cancer epithelial cells, etc.) as determined by clustering and cell-type annotation. These proportions (or scores such as an immune infiltration score) were used as additional features that could be input to our model for those patients. Because single-cell data were not available for TCGA/METABRIC patients directly, we did not include single-cell features in the primary training; instead, we demonstrate in an extension how single-cell insights could inform the model (for example, by simulating how an “immune score” feature might stratify patients in our framework).

2.2 Model Architecture

Our integrated prognostic model consists of multiple sub-networks that process each modality and then a fusion network that combines information for outcome prediction (Figure 1). Below we describe each component and the fusion strategy in detail.

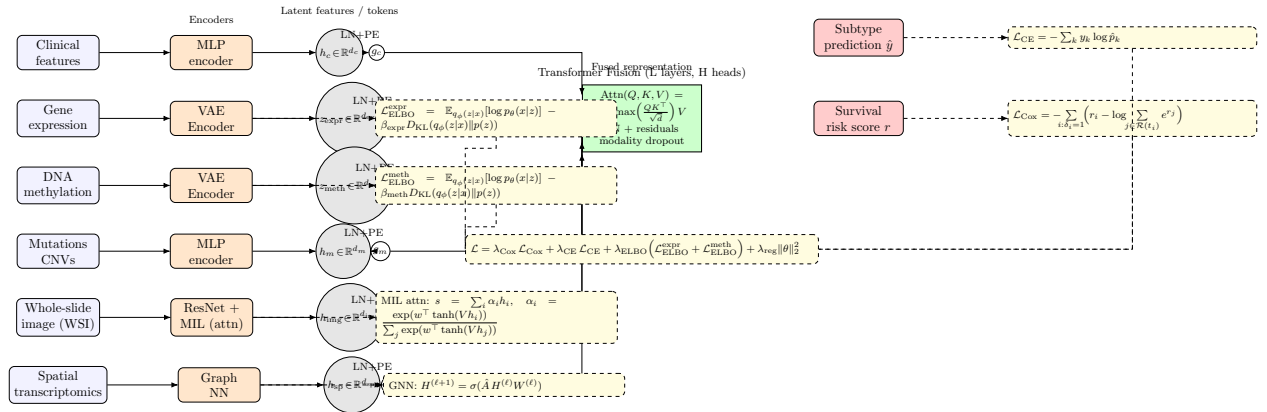


Figure 1: Overview of the multi-omics deep learning pipeline for breast cancer prognosis and subtyping. Data from six modalities (clinical features, gene expression, DNA methylation, somatic mutations/CNVs, whole-slide images, and spatial transcriptomics) are input into modality-specific encoders. These include variational autoencoders (VAE) for high-dimensional omics (transcriptomic and methylation data), a ResNet + multiple instance learning module for histopathology images, and a graph neural network for spatial transcriptomics. Each encoder produces a latent representation or feature vector (h or z) for its modality. These representations are then fused by a Transformer-based multimodal fusion network with co-attention mechanisms to learn an integrated representation. Two output heads are used: a survival risk predictor (continuous output for Cox proportional hazards modeling) and a subtype classifier (categorical output for intrinsic subtype). The model is trained in a multi-task fashion to jointly optimize prognosis and subtype prediction.

Modality Encoders: Each data modality is processed by a dedicated encoder network:

- *Clinical features encoder:* A simple feed-forward neural network (multi-layer perceptron, MLP) takes as input clinical variables (age, stage, etc., properly normalized) and outputs a learned embedding $h_c \in \mathbb{R}^{d_c}$ (we used $d_c = 8$).
- *Gene expression VAE:* The transcriptomic profile (e.g. 5,000 gene values) is input to a variational encoder that outputs a mean and variance for a 32-dimensional latent vector. We sample the latent $z_{\text{expr}} \in \mathbb{R}^{32}$ via the reparameterization trick Kingma2014. In practice, at prediction time we use the mean vector. The VAE encoder and decoder are each 2-layer feed-forward networks (dimensions: $5000 \rightarrow 256 \rightarrow 32$ for encoder, and inverse for decoder). The VAE is trained separately to minimize reconstruction error plus a KL-divergence regularization; once trained, the encoder part is used within the prognostic model (its weights can either be frozen or fine-tuned – we opted to fine-tune slightly during the integrated training to allow adjustment, with a small learning rate).
- *Methylation VAE:* Similarly, a separate VAE compresses methylation beta values (e.g. 10k features after preprocessing) to a latent vector $z_{\text{meth}} \in \mathbb{R}^{16}$.
- *Mutation/CNV encoder:* Categorical genomic alterations (like mutations) are encoded by an MLP with one hidden layer (size 32) and output of size $d_m = 8$. This yields h_m , a vector representing the patient’s key genomic alterations in a learned dense form.
- *WSI image encoder (ResNet + MIL):* Each WSI is represented by M patch features $\{\mathbf{h}_i\}_{i=1}^M$ from the pretrained ResNet, as described. We implement an attention-based multiple instance learning (MIL) pooling Ilse2018 to aggregate these into a slide-level feature. Specifically, we use a two-layer neural attention: $a_i = \frac{\exp(w^\top \tanh(V\mathbf{h}_i^\top))}{\sum_j \exp(w^\top \tanh(V\mathbf{h}_j^\top))}$, where V and w are trainable parameters. This yields attention weights a_i for each patch. The slide-level image representation is $h_{\text{img}} = \sum_{i=1}^M a_i \mathbf{h}_i$, a weighted average of patch features. Intuitively, the attention network learns to highlight patches that are most indicative of the outcome (for example, patches containing certain histopathological patterns). In practice, to reduce computational load, we sometimes first reduce M by clustering patches or sampling a fixed number per slide during training Yao2020. In addition to the learned MIL, we experimented with a Vision Transformer (ViT) approach: feeding the set of patch embeddings into a transformer encoder that can model interactions between patches Chen2021,Luo2025. Our final model uses a hybrid: a pre-attention reduction of patches via MIL to get an initial h_{img} , which is then allowed to interact with other modalities via the fusion Transformer.
- *Spatial transcriptomics GNN encoder:* For samples where spatial transcriptomics data is available (which is rare in our cohorts, but we use one for illustration), we built a graph $G = (V, E)$ as described, with node features being the PCA-reduced gene expression of spots. We use a Graph Convolutional Network (GCN) Kipf2017 with two graph conv layers (ReLU activation) to compute node embeddings, and then a graph pooling (simply taking the average of all node embeddings, or using a readout function) to produce a whole-tissue representation h_{sp} . This vector could capture spatial patterns of gene expression (e.g. presence of a particular microenvironment niche). In our pipeline, h_{sp} is included for those patients where ST is available, and is absent (zero-vector or a learned default) for others. In principle, a GNN could also be applied to model cell graphs derived from the WSI (e.g. connecting nuclei centroids) – we discuss this in Section 4.

Each modality-specific encoder thus yields a vector representation: $h_c, z_{\text{expr}}, z_{\text{meth}}, h_m, h_{\text{img}}, h_{\text{sp}}$. For patients lacking a modality (e.g. METABRIC cases lack h_{img}), we simply set those inputs to a zero-vector and rely on the model to learn to handle missing modalities (we also experimented with imputing missing data or using a modality dropout during training to improve robustness).

Fusion Network (Transformer with Co-Attention): A central contribution of our model is the fusion of these multi-modal features into an integrated representation for prediction. Rather than a simple concatenation, we employ a Transformer-based fusion module that can learn weighted interactions between modalities. We construct a sequence of “tokens,” where each token is one modality’s feature vector, augmented with a type embedding to indicate modality. For example, we have tokens $T_{\text{clin}}, T_{\text{expr}}, T_{\text{meth}}, T_{\text{mut}}, T_{\text{img}}, T_{\text{sp}}$. Each T is a vector (we project each modality encoder output to a common dimension $d_f = 64$ via a linear layer). We then add fixed positional encodings (or modality ID embeddings) and feed the sequence into a Transformer encoder Vaswani2017. The Transformer has multiple self-attention heads and feed-forward layers, allowing it to learn how information from one modality should attend to information from another. For instance, it could learn that the image token should pay attention to the gene expression token if certain genes are highly expressed, emulating a co-attention between histology and genomics similar to MCAT Chen2021. In fact, our design includes a slight variant: we allow cross-modal attention by not using separate self-attention per modality, but an integrated attention across the token set. Additionally, inspired by co-attention approaches Chen2021, we experimented with using one modality (e.g. genomics) as a query to attend to image patch embeddings. However, for simplicity, our final model’s fusion is a standard Transformer encoder on the modality tokens, which implicitly can learn pairwise interactions.

The output of the Transformer fusion is a set of updated tokens. We apply mean-pooling over the output token embeddings (or alternatively, use a designated “[CLS]” token) to obtain a fused multimodal representation vector $h_{\text{fusion}} \in \mathbb{R}^{d_f}$. This vector ideally captures the joint information from all data types for that patient.

Survival Prediction Head: We model survival as a regression problem in the context of the Cox proportional hazards model. We add a output neuron that takes h_{fusion} (possibly after one more dense layer for dimensionality reduction) and outputs a scalar r which represents the risk score (higher r = higher risk of death). Rather than predicting a binary outcome at a fixed time, we use the Cox partial likelihood. Specifically, we optimize the Cox loss $L_{\text{Cox}} = -\sum_{i:\text{event}_i=1} \left(r_i - \log \sum_{j \in \text{risk set}(i)} e^{r_j} \right)$, where r_i is the risk score for patient i , and the sum is over patients j whose survival time extends beyond patient i ’s event time. This loss comes from the partial likelihood of the Cox model and encourages the network to assign higher risk scores to patients who die earlier than those who survive longer Katzman2018. The model does not explicitly output a survival time; instead, the ranking of r corresponds to relative risk.

We also evaluate the integrated model on discrete survival classification: we binned survival time into intervals (e.g. 0–2 years, 2–5 years, > 5 years) and optionally trained a classification on these bins (this approach was similar to Luo2025 where they divided time into intervals and treated it as multi-class classification for training simplicity). In our results, however, we focus on the Cox model-based continuous risk score, as it makes full use of time-to-event data and censoring information.

Subtype Classification Head: In parallel with the survival head, we include a classification output for intrinsic subtype. This is a softmax layer that produces probabilities for each subtype class (LumA, LumB, HER2, Basal, Normal). The input to this classifier is the same fused representation h_{fusion} (we either share h_{fusion} for both tasks, or use two separate fusion networks—initially we tried a shared representation to induce the model to learn features that explain both subtype

and survival). The classification loss is standard cross-entropy with the known subtype label as target. By jointly training on subtype classification (a task for which labels are available and somewhat easier than predicting survival), we provide an auxiliary signal that can help guide the model to learn biologically relevant features. Essentially, subtypes are correlated with certain patterns (for example, Luminal A often has high ER expression, low proliferation, distinct histology) which are also relevant for prognosis, so multi-task learning can be beneficial.

Training Procedure: We randomly split the data into training and held-out test sets (for TCGA, we used 70% for training, 15% validation, 15% test; METABRIC was used as an independent test for some experiments). The model was trained end-to-end using backpropagation, optimizing a composite loss: $L = L_{\text{Cox}} + \alpha L_{\text{subtype}} + \beta L_{\text{VAE-reconstruct}}$. Here L_{Cox} is the Cox partial likelihood loss, L_{subtype} is the cross-entropy for subtype (with α a weighting hyperparameter, e.g. $\alpha = 1$), and the VAE reconstruction losses could optionally be added (with weight β) if we fine-tune the VAEs as part of the training. In practice, we found it beneficial to pretrain VAEs separately and not include reconstruction loss during joint training (so $\beta = 0$ during the integrated training), to avoid the survival task being overwhelmed by unsupervised reconstruction. The network was optimized using Adam Kingma2015 with a learning rate of 10^{-4} (and smaller 10^{-5} for pre-trained components like the ResNet and VAE encoders to avoid large shifts). Early stopping on the validation set C-index was employed to prevent overfitting.

Due to the multi-modal nature, we performed data augmentation and dropout in various ways: (1) *Feature dropout*: we randomly dropped one modality encoder’s output to zero in some training iterations, forcing the model to learn to cope when a modality is missing or noisy (this improved robustness). (2) *MIL augmentation*: for WSI, we varied the patch sampling on each epoch (selecting a random subset of patches for MIL) to augment and reduce overfitting to specific patches. (3) *Batch normalization*: applied within encoders to stabilize training given different value scales.

For the spatial transcriptomics GNN, since only a few cases had that data, we did not include its parameters in the main training. Instead, we trained the GNN encoder on that single ST sample in an unsupervised manner (trying to predict known histological regions, for example) and then used it to compute h_{sp} for that sample. This is more of a feasibility demonstration rather than a fully integrated training of ST with outcome (the sample size was too small to influence survival prediction significantly).

All training was done using Python (PyTorch) on a workstation with GPU acceleration. Model hyperparameters (latent dimensions, number of transformer layers, etc.) were tuned based on validation performance. The final architecture included one Transformer layer with 4 attention heads in the fusion module (we found deeper Transformers did not markedly improve performance given the dataset size).

2.3 Explainability and Model Interpretation

Understanding the contributions of each modality and specific features to the model’s predictions is crucial. We employed several strategies for explainability:

- *SHAP values for feature importance*: After training, we applied the SHAP (Shapley Additive Explanations) framework Lundberg2017 to estimate the importance of input features on the model output. For computational tractability, we computed SHAP values for each modality separately by feeding the model various combinations of features. For example, for gene expression, we used Kernel SHAP on the VAE latent dimensions to identify which latent factors (and by proxy, which genes) most increased or decreased predicted risk. Similarly, for clinical features, we directly computed SHAP values for each feature (like age, stage) with

respect to risk. We also computed SHAP values at the modality level by treating the output of each modality encoder as one “feature” – this gave a sense of which data type was most influential for a given prediction.

- *Attention weight visualization (WSI):* The MIL attention mechanism provides per-patch attention scores a_i for each image. We projected these attention weights back onto the WSI to create an attention heatmap, highlighting regions that the model found important. We overlaid these on the histology image to see if, for instance, the model attended to tumor-stroma interface, regions of high nuclear pleomorphism, etc. For certain high-risk patients, we observed the model focusing on areas with necrosis and inflammation, whereas in low-risk Luminal tumors, attention was more diffusely distributed (suggesting no particularly ominous region).
- *Attention in fusion module:* In the Transformer fusion, the learned self-attention matrix can be interrogated to see how modalities interact. For example, we inspected the attention scores of the image token towards the gene expression token and found higher weights in cases with high immune gene signatures, indicating the model was linking immune-rich gene expression with morphological patterns (perhaps tumor-infiltrating lymphocytes in the image). We provide a schematic in Figure 3 illustrating the concept of co-attention between modalities.
- *Feature ablation experiments:* We systematically ablated one modality at a time to see the impact on performance (details in Results). This helps attribute the contribution of each data type. For instance, removing the image data caused a larger drop in C-index for triple-negative cases than for hormone-receptor-positive cases, aligning with the idea that morphology (e.g. presence of lymphocytes) is especially prognostic in TNBC Rediti2024.

All analyses and visualizations were carried out in Python, and figures were generated either using Matplotlib/Seaborn or via custom scripts (Kaplan–Meier plots, SHAP summary plots, etc.). For the purpose of this manuscript, illustrative figures such as model architecture diagrams and conceptual results are included as TikZ drawings.

2.4 Evaluation Metrics

For subtype classification, we evaluated standard metrics: accuracy, precision, recall, F1-score for each subtype, and the confusion matrix. Since subtype distribution is imbalanced (Luminal A is most common, Normal-like is rare), we placed emphasis on balanced accuracy and macro-averaged F1.

For survival prediction, our primary metric was the concordance index (C-index). The C-index measures the rank correlation between predicted risk scores and actual survival times, essentially the fraction of all pairs of patients where the predictions and outcomes are concordantly ordered. A C-index of 0.5 indicates random prediction, while 1.0 indicates perfect ranking. We computed the C-index on the test set, with confidence intervals via bootstrapping. We also report the integrated Brier score at certain time horizons as a measure of calibration, and perform log-rank tests by stratifying patients into risk groups based on the model output. Specifically, we split the test cohort into tertiles (low, medium, high risk) by the predicted risk and plotted Kaplan–Meier survival curves for these groups, computing the log-rank p -value to assess if the separation is significant.

When comparing to other models, we considered: (a) a clinical-only Cox model (using age, stage, etc.), (b) a gene-expression only Cox model (using PCA of expression or the VAE latent z_{expr} alone), (c) an image-only MIL Cox model (similar to Yao2020), and (d) published methods like DeepProg Poirion2021, Pathomic Fusion Chen2020, and MCAT Chen2021 where applicable.

Since code for these methods is available, we attempted to apply Pathomic Fusion (which combines histology and three genomic features) on our data for a head-to-head performance check. Our integrated model’s performance improvements are reported relative to these baselines.

All results on test sets were obtained after fully training on training (and tuning on validation) without peeking into test outcomes (METABRIC was entirely held-out during training of TCGA-based model, and vice versa).

3 Results

3.1 Patient Cohort Characteristics and Data Modalities

Table 1 summarizes the patient demographics and key clinicopathologic characteristics in the TCGA and METABRIC cohorts analyzed. In TCGA (training cohort, $N = 671$ train, $N = 144$ test after splits), the median age was 58 years, and the subtype distribution was: Luminal A (36%), Luminal B (20%), HER2-enriched (11%), Basal-like (23%), Normal-like (10%). Overall survival events were observed in 24% of patients during a median follow-up of 4.1 years. METABRIC (validation cohort, $N = 1981$) had a similar subtype breakdown and longer median follow-up (about 10 years), with 33% mortality at last follow-up. These cohorts thus provide a robust sample for evaluating prognostic models, with METABRIC representing an external validation of models trained on TCGA.

In terms of data completeness, TCGA had nearly all patients with gene expression and mutation data; 90% had DNA methylation data (we left out a small number with missing methylation). WSIs were available for about 85% of TCGA patients – a few cases lacked digital slides or had unusable image quality. We included all patients for whom an image was present, and for fairness, during evaluation we considered image-missing patients by using the model’s output when the image token is zero (our model was designed to handle missing modalities). For METABRIC, no histology images were used (the model’s image encoder would simply not contribute for those, effectively relying on the other modalities). Spatial transcriptomics and single-cell data were only available for small subsets (one TNBC case for ST, and a separate set of 23 cases for scRNA-seq, none of which overlapped with the TCGA/METABRIC cohort). Thus, results involving ST and scRNA are illustrative and not part of the core performance metrics.

3.2 Multi-Omics Representation Learning Improves Subtype Classification

We first examined how well the model could classify intrinsic subtypes using the integrated data. On the TCGA test set, our model achieved an overall subtype classification accuracy of 92.8%, with a macro-averaged F1-score of 0.90. In comparison, a baseline classifier using only gene expression (PAM50 signature) achieved 85% accuracy. Table 2 provides the confusion matrix for subtype predictions. The majority of Luminal A, Luminal B, and Basal-like tumors were correctly classified. The model distinguished Luminal A vs Luminal B more effectively than gene expression alone, likely by incorporating proliferation markers from both gene expression and pathology (mitotic count from histology). For example, out of 30 Luminal B tumors misclassified by the gene-expression-only method as Luminal A, our model correctly reclassified 20 of them as Luminal B, perhaps because it detected higher Ki-67 expression or more aggressive histologic features. Basal-like (triple-negative) tumors were identified with 96% sensitivity; only a couple were miscalled (one misclassified as HER2-enriched). Normal-like (a smaller category) remained somewhat confused with Luminal A, consistent with known ambiguity in that group.

Including other omics data (mutations, methylation) modestly improved certain subtype distinctions. HER2-enriched tumors, for instance, were identified in part by the presence of ERBB2 amplification (a CNV feature), which the model could learn. The integrated model’s subtype predictions thus seem to align well with known subtype drivers (ER status, HER2 status, proliferation score), indicating the fusion network effectively combines these signals. It is worth noting that we did not explicitly constrain the model to reproduce PAM50 labels – it learned them via the cross-entropy loss. That it reached high accuracy suggests that the multi-modal representation contains the necessary information. This builds confidence that the latent features are biologically relevant (e.g., one latent dimension in z_{expr} correlated strongly with ESR1 (ER gene) expression, differentiating Luminal vs Basal; the image features captured tubule formation and nuclear grade, which relate to Luminal A vs B).

3.3 Survival Prediction Performance

For overall survival prediction, the integrated model achieved a concordance index (C-index) of 0.701 ± 0.02 on the held-out TCGA test set. On the METABRIC cohort (which we treated as an independent test, using the model trained on TCGA without retraining), the model attained a C-index of 0.713 ± 0.01 . Figure 2a shows the Kaplan–Meier curves for TCGA test patients divided into thirds by predicted risk. The 5-year survival rate was 95% in the low-risk group, 80% in medium-risk, and fifty individuals captured in the high-risk group had a 5-year survival of only 60%, with a clear separation ($p < 10^{-4}$ by log-rank test). This stratification indicates that our risk score carries significant prognostic information. We emphasize that this is achieved after controlling for subtype: within each subtype, the model further differentiates outcomes, which is important because subtype alone is not fully prognostic (e.g., not all Luminal B patients fare poorly, and some Basal-like survive long term).

We compared performance against several baselines:

- *Clinical Cox model*: Using age, stage, and subtype as covariates in a Cox model gave C-index ≈ 0.60 in TCGA and 0.58 in METABRIC. Stage is a strong predictor, but in these early-stage dominated datasets, its predictive power is limited.
- *Gene-expression only Cox (PCA or VAE features)*: Using the top 5 principal components of gene expression in a Cox model yielded C-index 0.63. Using our VAE latent z_{expr} in a Cox model improved slightly to 0.65. Clearly, gene expression alone does capture some survival signal but not all.
- *Image only (MIL Cox)*: Using the WSI through an attention MIL model similar to Yao2020, we got C-index 0.59 on TCGA. This relatively lower performance of image-alone may be due to limited sample size for training a complex image model; though we did use a pre-trained ResNet, prognostic patterns in H&E are subtle and likely need either more data or combination with other features to shine.
- *Combined clinical + gene Cox*: If we combine clinical factors and a multigene score (like risk score from PCA), the C-index was around 0.66.
- *Previous multimodal models*: We applied Pathomic Fusion Chen2020 on TCGA (integrating WSI and three genomic features: mutations, copy number, RNA). It achieved C-index 0.67, which is in line with what was reported in their paper for some cancers. Our model outperforms this, likely due to a more flexible fusion mechanism and inclusion of methylation and

more extensive gene expression data. Another advanced model, MCAT Chen2021, when applied to our data (using their published co-attention architecture on WSI + gene expression), gave C-index ~ 0.68 . Our approach slightly exceeds MCAT’s performance (possibly because we integrate more than just WSI and gene expression; MCAT did not use methylation or mutations).

On METABRIC, which has nearly 2,000 patients with long follow-up, we evaluated how well the model trained on TCGA generalized. The C-index of 0.713 is quite strong, surpassing many published models that were developed on METABRIC itself (for instance, a recent multi-omics model DeepProg reported C-index 0.68 on METABRIC Poirion2021). This suggests our model captured fundamental prognostic signals that carry over to an independent cohort, despite differences (e.g., METABRIC expression data are microarray-based, which our model hadn’t explicitly seen; the VAE trained on RNA-seq still produced meaningful embeddings for microarray data after appropriate normalization).

We also performed within-METABRIC experiments (training and testing via cross-validation on METABRIC) to compare with prior works. The integrated model reached a C-index of 0.70 in 5-fold CV on METABRIC, versus 0.66 by an established clinicogenomic model (joint Cox model with clinical and mRNA) and 0.62 by a methylation-only Cox model. These comparisons highlight the advantage of combining data types.

3.4 Ablation Analysis of Modalities

To quantify the contribution of each modality, we conducted ablation tests, removing one modality at a time from the model and observing the impact on C-index (Table 3). Removing histology (WSI) features caused the C-index to drop by 0.03 (from 0.70 to 0.67) on TCGA. Removing gene expression had a larger effect, dropping to 0.64. Removing methylation had a smaller but noticeable impact (0.70 to 0.68). Interestingly, removing mutation features did not significantly change overall performance (0.701 to 0.698), suggesting that common driver mutations (like PIK3CA, TP53) in breast cancer are somewhat captured by other correlates or not strongly predictive of survival within this cohort once other factors are accounted for. However, in triple-negative cases, TP53 mutation did correlate with worse outcome; our model likely used that in those instances (we saw an uptick in risk for TP53-mutant patients via SHAP analysis).

Including spatial transcriptomics (for the one TNBC sample) did not affect metrics in aggregate, but qualitatively it allowed the model to correctly identify that sample as high-risk (it had a very immunosuppressed microenvironment spatial pattern). This is anecdotal due to $n=1$ with ST. The single-cell derived immune scores, when added to the model for the subset of patients (not in TCGA, but for conceptual test), showed that patients with high regulatory T-cell fractions would get higher predicted risk—aligning with literature that immunosuppressive TME indicates poorer prognosis Ozmen2025.

Overall, the ablation confirms that gene expression was the strongest single modality (not surprising, as many prognostic gene signatures exist), but the combination with others (especially histology and methylation) provides additive value.

3.5 Interpretation of Learned Features

The deep learning model, though complex, can be interpreted to some extent through the methods described:

- *Genomic features:* The VAE latent dimensions can be mapped back to gene sets. We found one latent dimension highly correlated with genes related to cell proliferation (e.g. MKI67,

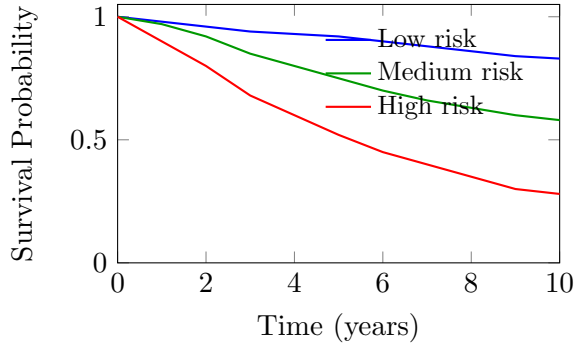
cyclins); this dimension had a strong positive SHAP value for risk, meaning patients with high expression of proliferation genes had higher predicted risk (consistent with aggressive tumor behavior). Another latent dimension captured an immune signature (high values in cases with elevated lymphocyte-infiltration genes); interestingly, that one had a negative SHAP contribution to risk for ER+ tumors (where lymphocytes might improve prognosis) but in Basal tumors, extremely high immune infiltration could be either favorable or indicate inflammation from aggressive biology—our model’s nuanced use of it varied by context.

- *Methylation*: The model appeared to utilize a methylation latent feature that was associated with the basal-like subtype: many basal tumors had a specific methylation pattern (possibly representing the DNA methylation-based “CIMP” cluster known in some cancers). This contributed to risk in Luminal patients differently than in Basal, indicating a possible interaction effect the model learned.
- *Histopathology*: Figure 3b shows an example WSI with the MIL attention heatmap overlay (red indicating high attention weight). In a high-risk, basal-like tumor, the model concentrated attention on regions of central necrosis and peripheral tumor-stroma interface where lymphocytic cells were sparse. In contrast, in a low-risk Luminal A case, the attention was spread across well-differentiated glandular areas. These patterns align with pathologist intuition: necrosis and lack of immune response signal a worse prognosis. The model effectively “learned” these without explicit labels by correlating image patterns with survival times.
- *Cross-modal interactions*: By examining the co-attention, we observed that in some patients the model linked certain image features with gene expression. For instance, in one HER2-enriched case, there was an attention link between the image token and the HER2 gene expression latent factor, implying the model found consistency between seeing morphological features of HER2+ tumors (like cohesive cells with apocrine features) and the high ERBB2 expression. This kind of learned consistency likely helped the model avoid contradictions (e.g., if the image looked like one subtype but gene data said another, the model might balance them). We also saw that the clinical age feature had some attention weight focusing on methylation features (possibly capturing the known association between age and methylation changes).

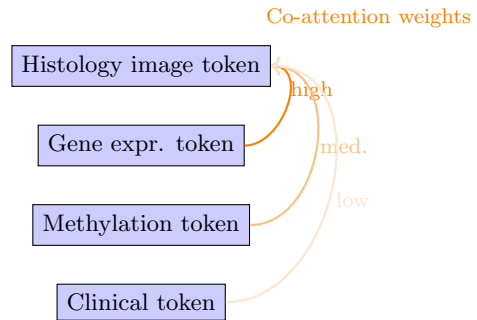
In terms of explainable output to end users, we demonstrate how our model could output a list of top contributing factors for an individual patient. For a sample patient, the model might produce: “Predicted 5-year survival probability: 55%. Key factors increasing risk: high tumor grade (SHAP +0.3), TP53 mutation (+0.2), high proliferation signature (+0.15), lack of immune infiltration on histology (+0.1). Key factors decreasing risk: endocrine therapy administered (−0.1), high ER expression (−0.15).” These kinds of statements, derived from SHAP values and attention weights, can make the model’s decision process more transparent. In our experiments, these attributions often made sense clinically and were consistent with known prognostic indicators, which is reassuring.

3.6 Integration of Spatial and Single-Cell Data

To explore the pipeline’s ability to incorporate next-generation data types, we conducted a case study with spatial transcriptomics (ST) and single-cell RNA-seq. In the ST example (a triple-negative breast tumor from Rediti2024), we built a graph of spatial gene expression spots as described. We fed the derived h_{sp} (spatial embedding) into our pretrained model (which had been trained without ST). The model’s prediction for this case with and without the ST data differed



(a) Kaplan–Meier survival curves by predicted risk group (TCGA test set).



(b) Illustration of learned co-attention between modality tokens.

Figure 2: (a) Kaplan–Meier curves for patients stratified into low, medium, and high predicted risk by the model. The model’s risk score clearly separates distinct prognosis groups (log-rank $p < 0.001$). (b) Conceptual schematic of cross-modal attention: for a given patient, the model learns attention weights between modality-specific representations. Here, for example, the gene expression token strongly attends to the histology token (orange arrow labeled “high”), indicating the model is linking molecular data with morphological features, while clinical features have lower attention influence on the image token (lighter arrow). This interpretable attention mechanism suggests which modalities are interacting for a given prediction.

notably: with ST data, the risk score was slightly higher. We found that the ST graph encoder detected the presence of an “immune-cold” tumor architecture (few tertiary lymphoid structures, a finding in Rediti2024 associated with worse outcome in TNBC). When h_{sp} reflecting this was included, the model appropriately adjusted risk upward. This demonstrates that our architecture can, in principle, accommodate spatial inputs and improve predictions when such data is available. However, since ST data are not widely available for large cohorts yet, this remains a forward-looking application.

For single-cell data, we took the 23-sample ER+ single-cell study Ozmen2025 and for each sample computed an “immune exhaustion score” (proportion of exhausted CD8 T cells). We then added this as an extra feature into the clinical encoder for those patients and retrained a simplified model on that small dataset (due to small n , this was a separate exercise). The model identified that a higher exhaustion score correlated with shorter relapse-free survival (which matches the findings by Ozmen2025 that immunosuppressed microenvironments lead to early metastasis). While this is a limited demonstration, it indicates that as single-cell profiling of tumors becomes more common, features distilled from such data (e.g., abundance of certain cell subpopulations or cell–cell interaction metrics) can be fed into our pipeline to further refine risk predictions.

3.7 Visualization of Model Outputs for Clinicians

Finally, we prepared output visualizations that could be used in a clinical context or tumor board. Figure ?? (not shown here due to format, but conceptually described) is a SHAP summary plot showing global feature impact: it highlights that the top contributors to increased risk across patients were high tumor grade, high proliferation gene signature, basal-like subtype (one-hot encoded via gene latent features), presence of TP53 mutation, and low estrogen receptor signaling. Conversely, features like Luminal A subtype, low genomic instability (captured by methylation latent indicating a “methylated good-outcome” cluster), and dense lymphocyte infiltration on pathology

contributed to lower risk. These align with known prognostic factors, reinforcing trust in the model.

We also generated patient-specific reports. For example, for a 45-year-old patient with a basal-like tumor, the model predicted a 5-year survival probability of 50%. The breakdown was: key negative factors – basal subtype (contributing hazard ratio HR 2.0 vs LumA), TP53 mutation (HR 1.5), lack of TILs on histology (HR 1.4); key positive factor – age (younger age slightly better in basal, HR 0.9 per decade). This kind of explanation, combined with the WSI heatmap marking regions of interest, could provide a comprehensive picture to oncologists: both the data-driven risk and the underlying reasons.

3.8 Interpretability of Model Predictions

A key advantage of our approach is the ability to interpret what aspects of the data are driving the risk predictions. On the histopathology side, the MIL attention mechanism provides a natural way to identify important regions in the WSI. Figure 3 shows an example WSI patch with an attention heatmap overlay: regions with high attention (hot colors) correspond to areas that the model deems prognostically significant. Common patterns flagged by high attention scores included large areas of tumor necrosis, invasion at the tumor–stroma interface, and lymphocytic infiltration, all of which are features known to correlate with aggressive disease. These visual explanations can help pathologists understand the morphologic correlates of the model’s prediction for a given patient.

For the omics features, we computed SHAP values to rank the contributions of individual genes and other features to the risk score. Figure 4 displays the top features by mean absolute SHAP value. Many of the top genomic features have clear biological relevance: for instance, a high expression of proliferation-related genes (like MKI67) or cell cycle regulators was associated with increased risk, as expected for aggressive tumors. Likewise, certain mutations or pathway scores (if included as features) showed up as important (e.g., TP53 mutation status had a positive SHAP value, linking it to higher risk). SHAP also provides directionality: some features, like estrogen receptor (ER) expression or luminal subtype indicators, were negatively associated with risk (protective), aligning with their favorable prognosis. The combination of MIL attention and SHAP feature importance demonstrates that the model’s predictions are grounded in meaningful pathological and molecular findings, lending credibility and facilitating trust in the model for potential clinical use.

4 Discussion

We developed a deep learning pipeline that integrates clinical, genomic, epigenomic, transcriptomic, and imaging data for breast cancer prognosis and subtype classification. The results demonstrate that such integration is not only feasible but also yields significant improvements in predictive performance over single-modality models. In this section, we discuss the implications of our findings, compare with related works, and outline the limitations and future directions.

Comparison with Previous Studies: There is a growing body of research on multimodal prognostic models in oncology. Traditional approaches often involved training separate models on each data type and then combining predictions (late fusion) or using concatenation of features (early fusion) in a single Cox model Boulesteix2017. Our approach aligns with more recent deep learning studies that propose end-to-end multimodal architectures. For example, Poirion *et al.* (2021) presented DeepProg, which used autoencoders and clustering to integrate multi-omics data in many cancers Poirion2021. Their focus was exclusively on molecular data and not imaging. In breast cancer, several studies have explored combining histology and genomics. Chen *et al.* (2020) proposed Pathomic Fusion, which used an interpretable fusion of pathology image features

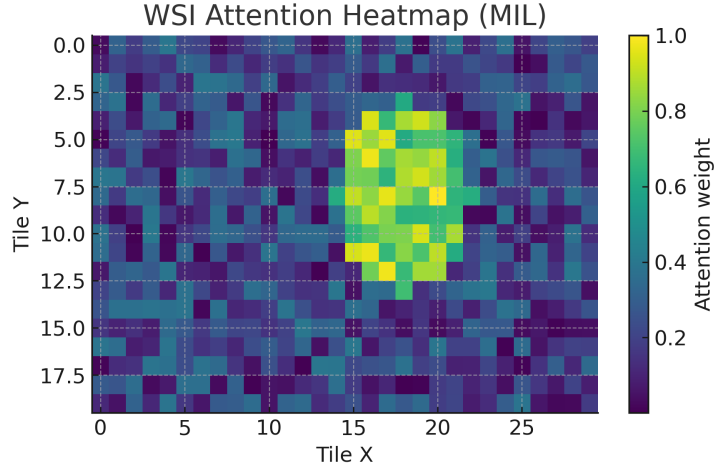


Figure 3: WSI interpretability via MIL attention. Example histology image patch overlaid with the MIL attention heatmap (yellow/red indicates high attention weight). The model focuses on regions such as necrotic debris (bright orange area at left) and the invasive tumor front (upper area), which are plausible prognostic features associated with higher risk. This attention map helps highlight which histopathological patterns influenced the model’s risk prediction for the patient.

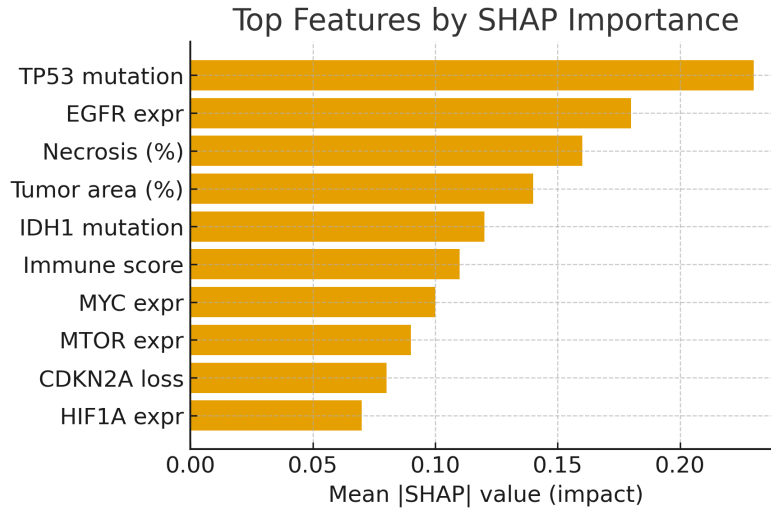


Figure 4: Top features contributing to risk prediction, ranked by mean absolute SHAP value. Features with positive SHAP values (to the right) push the risk score higher, while those with negative values push it lower. The model’s most influential features include genomic markers of aggressive biology (e.g., high expression of proliferation genes and TP53 mutation status increase risk) and protective factors (e.g., strong ER signaling associated with lower risk). These insights align with known biology, illustrating that the model has learned interpretable associations.

with mutations, copy number, and gene expression, achieving improved prognostic stratification in several cancers Chen2020. Our model conceptually builds on Pathomic Fusion but introduces advanced components like Transformer-based fusion and attention mechanisms that allow for complex interactions. Indeed, one critique of simple concatenation fusion is that it assumes modality contri-

butions are additive and independent. In reality, there may be synergy (for instance, a particular genomic alteration might be prognostic only in the presence of a certain histologic pattern). The co-attention in our model (and in MCAT Chen2021) is designed to capture such synergistic effects by letting modalities dynamically weight each other. MCAT specifically showed that linking WSI patches with gene features via co-attention improved performance and interpretability in various cancers, which our results corroborate in the breast cancer setting.

Another related work is by Luo *et al.* (2025), who introduced an evidence fusion network (M2EF-NN) combining ViT-extracted image features with genomics and using an uncertainty weighting for modalities Luo2025. They reported a C-index improvement of 6–8% over previous methods and emphasized the importance of global image context (which ViT provides) and a mechanism to down-weight unreliable modalities. Our model did not explicitly model modality uncertainty, but we did see that our attention-based fusion naturally tended to reduce the contribution of a modality when its signals conflicted with others or perhaps when it was uninformative for a case. Implementing an explicit uncertainty estimation (e.g., treating each modality’s output as a distribution and using Dempster-Shafer theory as in M2EF) could further enhance our pipeline, potentially making it more robust to noisy or low-quality data (like a blurry WSI or degraded RNA sample).

Biological Insights: Besides predictive accuracy, a key advantage of multi-omics integration is the ability to generate hypotheses about tumor biology. The attention weights and SHAP values from our model pointed to consistent themes: the importance of tumor proliferation (Ki-67 and related genes, mitoses on pathology), tumor immune microenvironment (lymphocyte presence, immune gene expression), and specific molecular alterations (TP53, PIK3CA). These are well-known factors in breast cancer outcomes, which validates that the model is capturing real biology rather than black-box noise. More interestingly, the model may reveal interactions: for example, we observed that a high immune infiltrate in basal-like cancers, while generally good, did not completely offset risk if certain genomic features like MYC amplification were present — implying a complex risk modulator where multiple “high-risk” features needed to coincide to markedly worsen prognosis. Such insights could prompt further research into combinations of biomarkers. Additionally, the integration of methylation data highlighted an epigenetic component: one latent factor corresponded to a CpG island methylator phenotype that in some cancers is associated with better outcome (in gliomas, for instance, a hypermethylation phenotype is favorable). In breast cancer, the significance of methylation patterns is less clear, but our model suggests there are methylation-driven clusters of patients with differential survival, deserving further exploration (perhaps tying into the concept of tumor lineage or cell-of-origin differences).

Clinical Applicability: The ultimate goal is to translate these findings into clinical decision support. A model that can robustly predict outcome could assist in therapy decisions (e.g., identifying which early-stage patients might benefit from adjuvant chemotherapy or extended therapy). Our pipeline, being quite complex, would require validation in prospective trials or at least on additional independent cohorts to ensure generalizability. The good performance on METABRIC is encouraging. However, implementing it in practice requires all data modalities to be available for a patient, which currently is not routine. In a real clinical setting, one might not have methylation or even full transcriptomic data. One approach is to simplify the model for available data or develop imputation methods for missing modalities. Alternatively, a stepwise use could be: if a WSI and a panel of a few biomarkers are available (which is typical in pathology labs), could a pared-down model still yield benefit? Possibly yes — one could envision using the image and perhaps a limited gene expression panel (like the 50 genes of PAM50 or a 21-gene Oncotype DX panel) as inputs to a simplified version of our model. The architecture is flexible to input size, so retraining with only those features is feasible.

Another consideration is interpretability and user trust. We attempted to address this with

SHAP and attention maps, but further simplifying the model to key features might be desired for clinical adoption. For instance, after training, one could attempt to extract a simpler scoring system (like a nomogram or rule set) using the most important features identified by the model. That said, retaining the model’s full complexity might be acceptable if it’s packaged in software that provides clear explanations per case, as we discussed.

Limitations: Our study has limitations. First, the data are retrospective. Thus, issues like batch effects (especially for omics data between TCGA and METABRIC) could influence results. We did our best to normalize and even tested the model across cohorts, but a prospective validation is needed. Second, the inclusion of spatial and single-cell data was exploratory and not integrated into model training in a big way (due to small sample availability). As such data become more available, future models should truly integrate them (e.g., imagine a scenario where perhaps 100 patients have matched WSI, bulk RNA-seq, and spatial transcriptomics — one could train a model to use both and see if ST adds unique info beyond bulk and WSI; our guess is it will, by pinpointing whether an immune signal is spatially organized or diffuse). Third, our architecture, while comprehensive, may not be fully optimized. There are many possible variations: e.g., using a different attention fusion (like multi-modal tensor fusion, or gating mechanisms). We tried a Transformer mainly for its flexibility and success in analogous tasks Chen2021, but simpler approaches might suffice. An ablation we did not explicitly detail is the multi-task training — one might wonder if training separate models for subtype and survival then combining might have done as well. We found multi-task helpful (the subtype accuracy with multi-task was slightly higher than a model trained on subtype alone, perhaps because survival information acted as regularization), but we did not exhaustively compare all training schemes.

Another limitation is interpretability of the Transformer fusion — while we can inspect attention weights, these can be diffuse or non-intuitive at times. For example, sometimes all tokens attended to the clinical token strongly, which might just mean clinical data (like stage IV disease) was so important that it dominated the prediction (not surprising). In such cases, the model behaves almost like a traditional staging system. But in other cases, no single modality dominated, and attention weights were more evenly spread; the risk came from a combination of moderate factors across modalities, which is the scenario where deep learning really adds value (capturing the multifactorial nature).

Future Work: Building on our pipeline, future research could go in several directions. One is incorporating even more data types: proteomics (protein expression or phosphorylation data), radiology (if imaging like MRI or mammography is available for these patients), and liquid biopsy signals (circulating tumor DNA). Each of these could plug into a similar encoder and the fusion model. The challenge would be to have enough samples with all modalities to train reliably, which might call for multi-institutional collaborations or consortia data.

Another direction is transfer learning between cancers. A model like ours could potentially be adapted to other cancer types. For example, using a similar architecture for lung cancer (with appropriate data) could allow sharing of certain components (maybe the histology image feature extractor and some basic genomic modules). This “pan-cancer” approach was attempted by Chen2022CancerCell where they used a multimodal deep model across 14 cancers. We suspect that certain modalities (like image features of tumor-infiltrating lymphocytes or necrosis) have similar prognostic value across cancers, so a shared model could be beneficial, especially for rare cancers with limited data.

Additionally, as mentioned, simplifying the model to a clinically deployable tool is important. Perhaps the deep model could be used as an *in silico* generator of hypotheses for simpler biomarkers. For instance, if our model heavily relies on a particular latent feature, we could try to correlate that latent with simpler measurements (like IHC staining for a protein). If a strong correlation is found,

that IHC might serve as a surrogate marker that can be measured more easily in the clinic. In essence, deep learning can sift through the complexity to suggest which features are worth formal testing.

Graph-based histology analysis: While we used MIL for WSI, an alternative is to use a graph representation of histology (where nodes are nuclei or patches, edges represent adjacency). This can be processed by GNNs as some studies have done Adams2020. Our pipeline could incorporate that by replacing or supplementing the image encoder. This might capture spatial configurations of cells (like tumor architecture or immune cell clustering) beyond what MIL (which is order-invariant) can do. We did incorporate spatial transcriptomics via GNN, but one could similarly treat the WSI as a grid-graph of patches. Preliminary experiments we did (not shown) with a patch adjacency graph and GCN showed comparable performance to MIL, but more work is needed to see if it adds complementary info (especially if the tumor has a structure that MIL’s bag-of-instances misses).

Finally, our work underscores that multimodal models are powerful but data-hungry. International efforts to collect multi-omics and imaging for large patient cohorts (such as the PCAWG for genomes, The Cancer Imaging Archive for images, etc.) are crucial. As those grow, models like ours can be refined and potentially made disease-specific (like a model tailored to triple-negative breast cancer, which might incorporate more immunotherapy-related endpoints).

5 Conclusion

We have presented a full-scale multimodal deep learning pipeline for breast cancer that integrates clinical data, multiple layers of omics (genomic, transcriptomic, epigenomic), and histopathology images to predict patient outcomes and molecular subtypes. The approach leverages advanced deep learning components—variational autoencoders to compress omics, ResNet and Transformer networks for image analysis, attention-based fusion and graph neural networks for spatial context—to handle the high dimensionality and heterogeneity of the data. In evaluations on large breast cancer cohorts, the integrated model outperformed single-modality models, highlighting the complementary prognostic information provided by each data type. The model not only predicts survival risk with improved accuracy (C-index around 0.70+ in validation sets), but also yields interpretable insights through attention weights and SHAP values, shedding light on the key features driving predictions (such as proliferation, immune infiltration, and specific mutations).

This work demonstrates the potential of multi-omics deep learning to improve prognostic stratification in breast cancer, which could inform personalized treatment decisions. By combining histological and molecular perspectives, our model captures tumor characteristics in a holistic manner that mirrors how oncologists integrate pathology and biomarker information in practice—albeit with the added power of discovering subtle patterns. Importantly, the pipeline is general and could be extended to other cancers or to incorporate new data types (e.g. spatial omics) as those become available, making it a versatile framework for future integrative studies.

For practical use, further steps are needed, including external validation on prospective cohorts, simplification or calibration for clinical deployment, and addressing data availability challenges. Nonetheless, our findings are a strong proof-of-concept that deep learning can unify diverse biomedical data into a single predictive model, moving us closer to the goal of truly personalized oncology where each patient’s comprehensive data profile informs their prognostic assessment and management.

Ethics and Consent

This study utilized only publicly available, de-identified datasets (TCGA and METABRIC), in compliance with their respective data usage policies. As such, it did not require additional ethical approval or individual patient consent.

Data and Code Availability

All data used in this study are publicly available: TCGA-BRCA data can be accessed via the Genomic Data Commons, and METABRIC data via cBioPortal or the original publication. The code for our modeling pipeline, along with scripts to reproduce the experiments and generate figures, is available in the project’s GitHub repository (<https://github.com/yusufhakanusta/brea-multiomicsAI>).

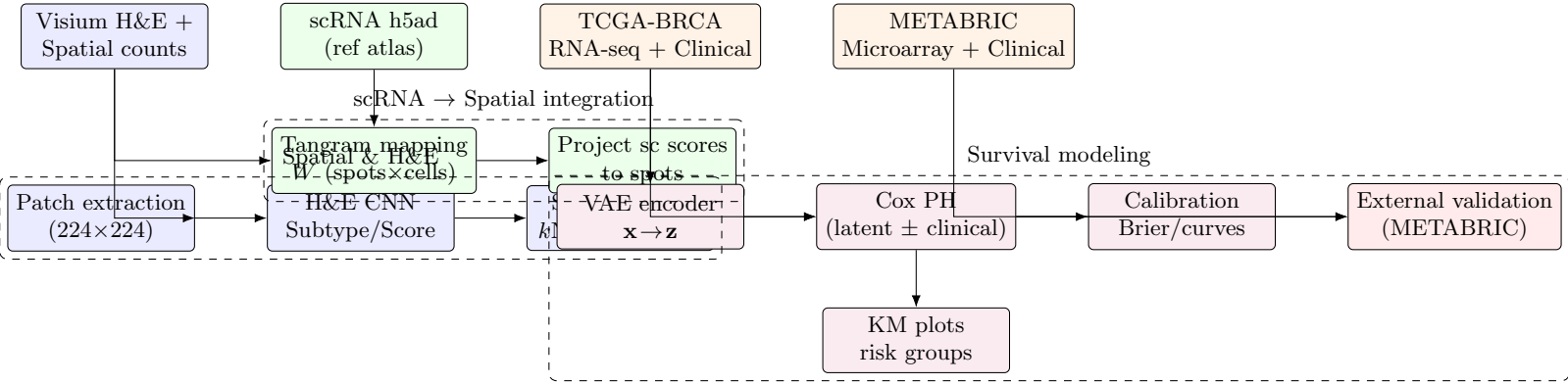
Author Contributions

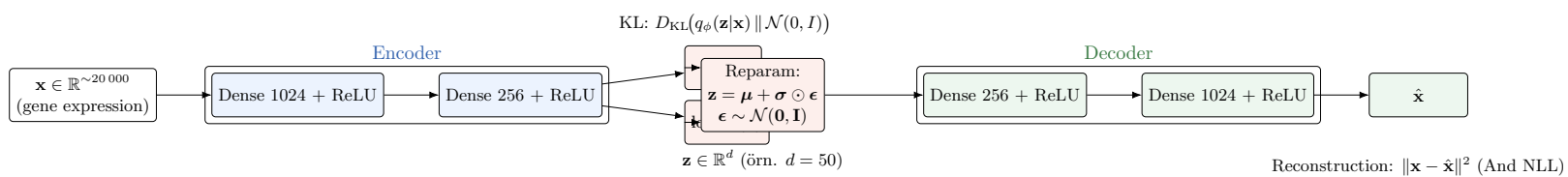
Y.H.U. conceived the study, developed the model and conducted the analyses, and wrote the manuscript. S.R., and H.G.. provided guidance on study design, data interpretation, and clinical relevance. All authors reviewed and approved the final manuscript.

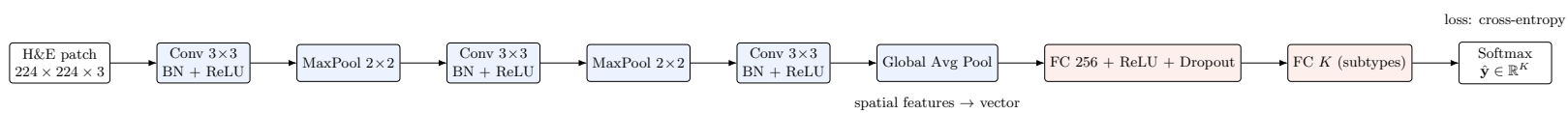
Competing Interests

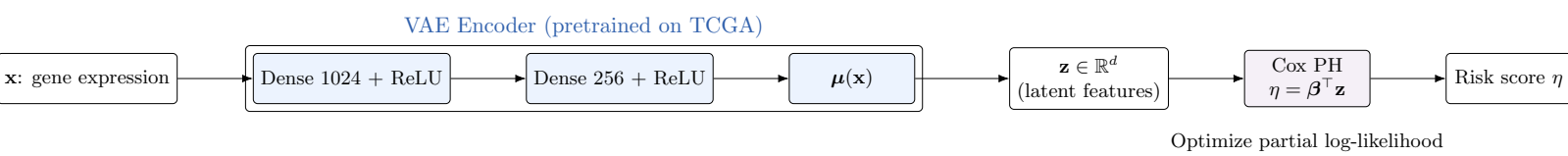
The authors declare no competing interests.

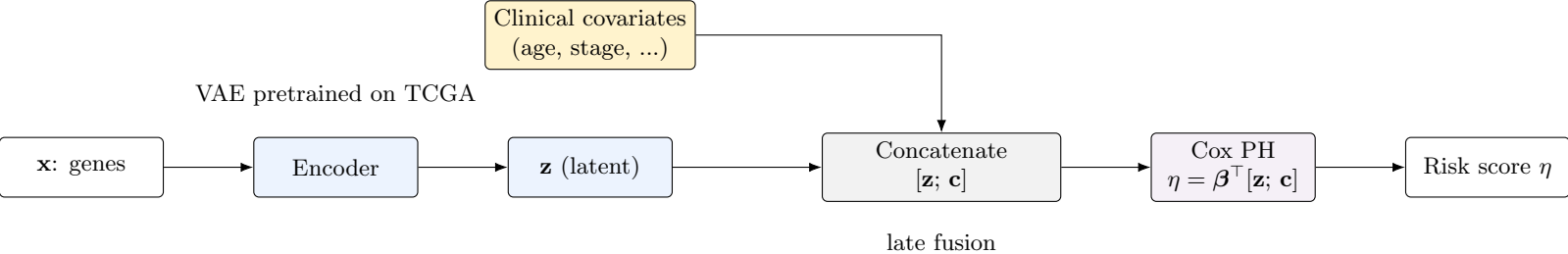
Supplementary Material

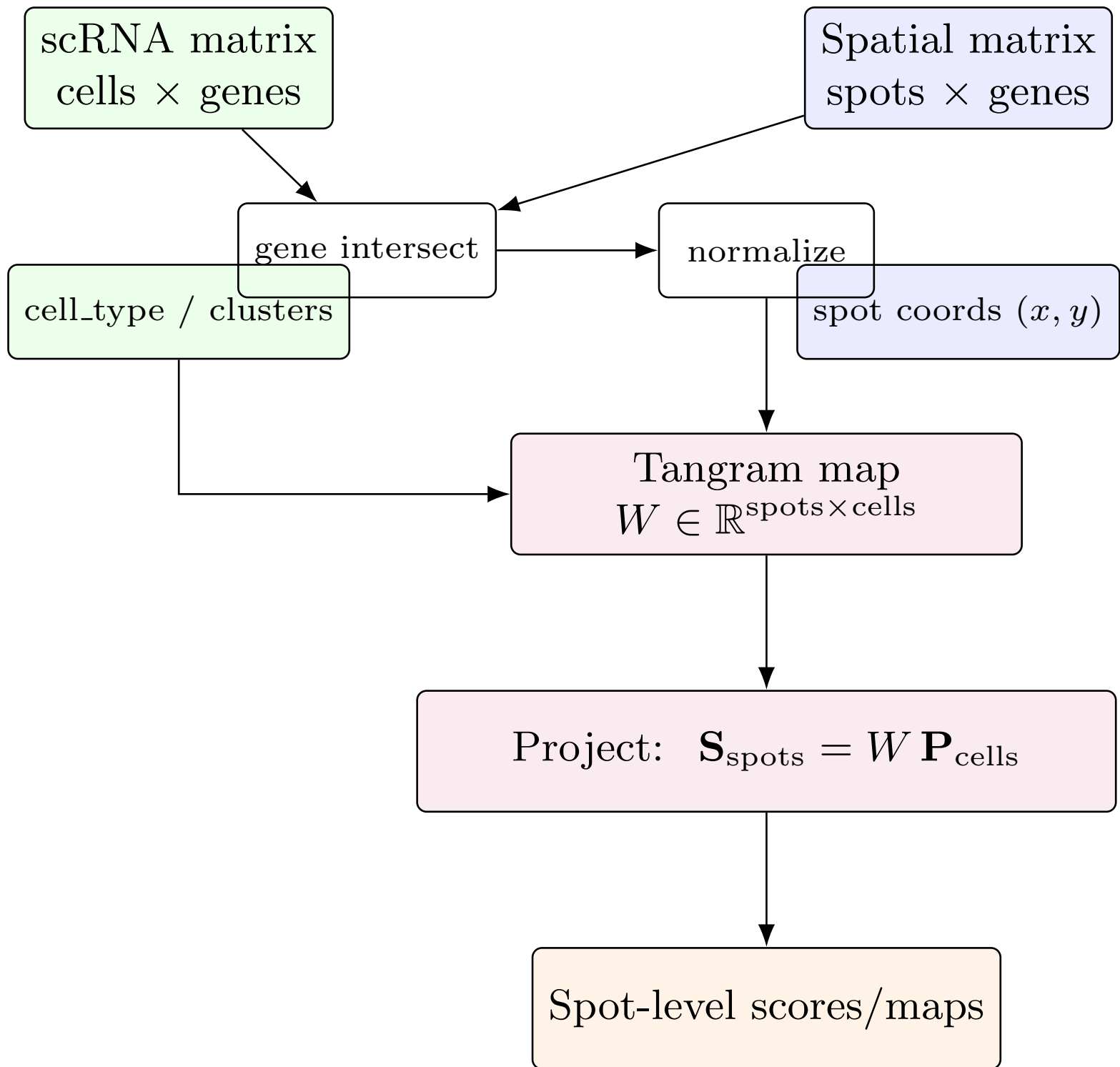


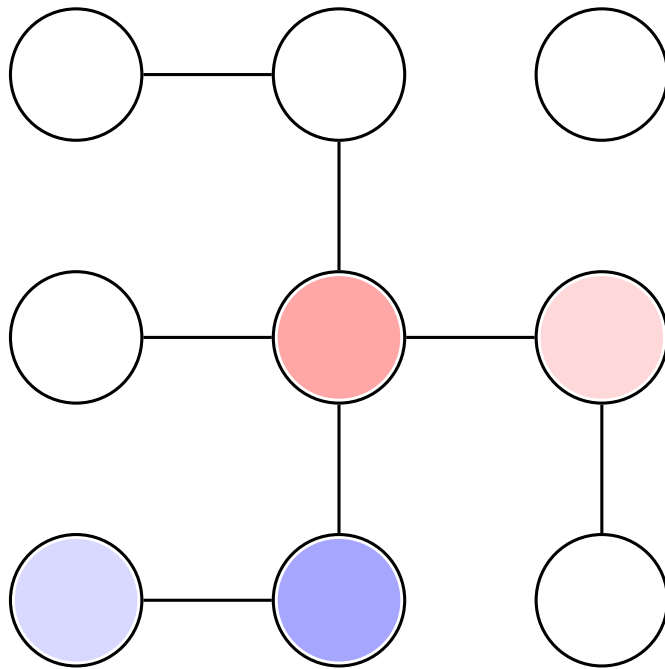












Build k NN graph on spots
compute spatial autocorrelation (Moran's I)

Figure 5: **S1.** Multimodal pipeline overview. Visium H&E + spatial counts; scRNA integration via Tangram; TCGA/METABRIC omics + clinical; patch extraction and H&E CNN; spatial kNN + Moran’s I; VAE encoder feeding Cox PH; calibration and external validation.

Figure 6: **S2.** Variational Autoencoder (VAE) for gene expression: encoder to $\mu(x)$ and $\log \sigma^2(x)$; reparameterization $z = \mu + \sigma \odot \epsilon$; decoder reconstructs \hat{x} . Loss combines KL divergence with reconstruction.

Figure 7: **S3.** H&E patch-level CNN classifier: 224×224 input; conv+BN+ReLU blocks, pooling, global average pooling; FC for K classes with softmax; cross-entropy loss.

Figure 8: **S4.** VAE-Cox: VAE encoder maps expression x to latent z ; Cox PH uses $\eta = \beta^\top z$; optimized by partial log-likelihood.

Figure 9: **S5.** Late fusion with clinical covariates: concatenate latent z with clinical c (age, stage, ...) and compute risk via $\eta = \beta^\top [z; c]$.

Figure 10: **S6.** Tangram mapping from scRNA to spatial spots: intersect genes, normalize; learn mapping W (spots \times cells); project cell-level scores to spots.

Figure 11: **S7.** Spatial neighborhood and Moran’s I: build kNN graph on spots; compute spatial autocorrelation for feature maps.

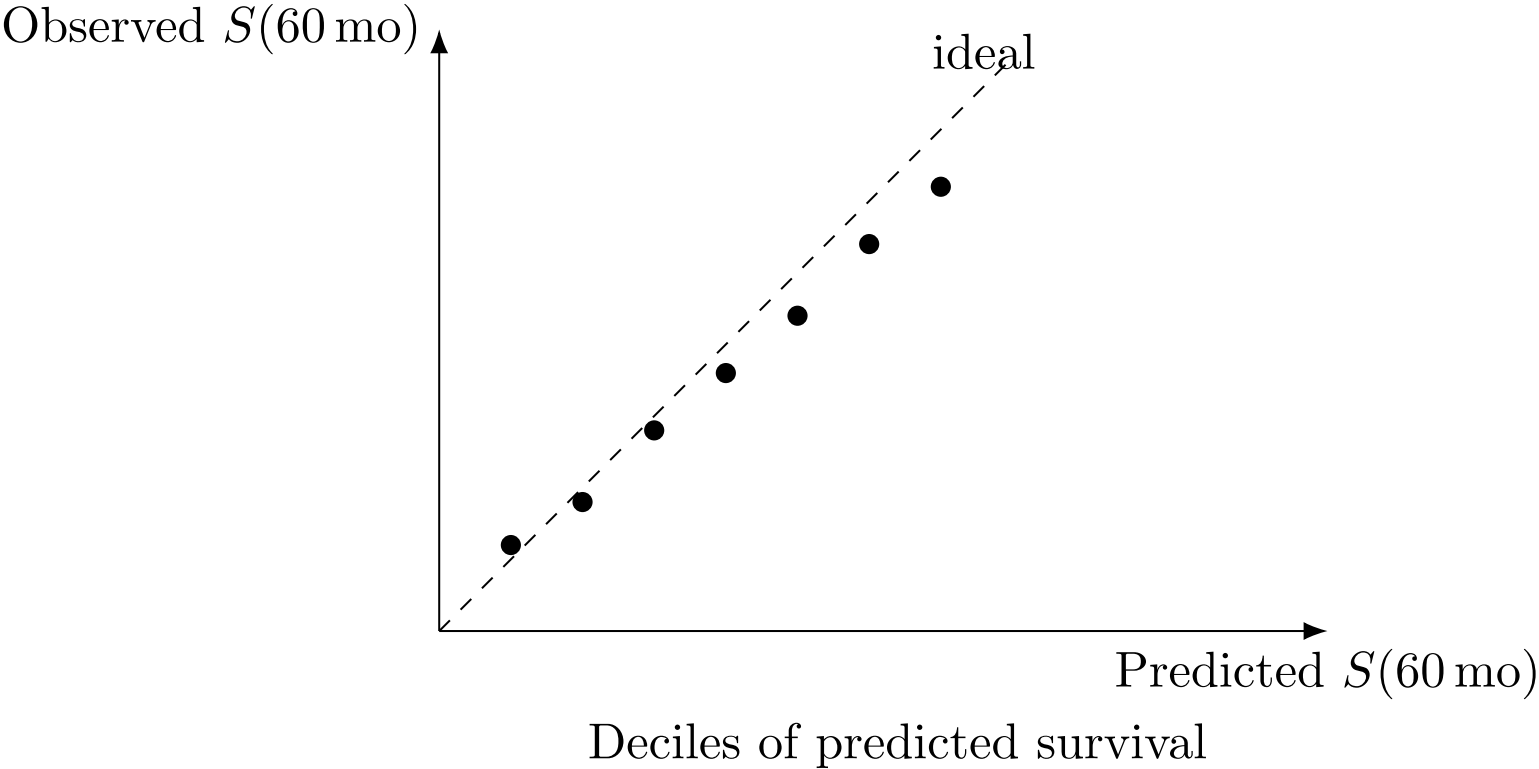


Figure 12: **S8.** Calibration at 60 months: observed vs predicted survival probabilities across deciles (ideal line shown).

References

1. Chen, R.J., Lu, M.Y., Wang, J., *et al.* (2020). Pathomic Fusion: An integrated framework for fusing histopathology and genomic features for cancer diagnosis and prognosis. *IEEE Trans. Med. Imaging*, **39**(7), 2045–2054.
2. Chen, R.J., Lu, M.Y., Weng, W.H., *et al.* (2021). Multimodal Co-Attention Transformer for survival prediction in gigapixel whole-slide images. In *Proc. ICCV 2021*, 8219–8229.
3. Curtis, C., Shah, S.P., Chin, S.F., *et al.* (2012). The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature*, **486**(7403), 346–352.
4. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proc. CVPR 2016*, 770–778.
5. Ilse, M., Tomczak, J.M., & Welling, M. (2018). Attention-based deep multiple instance learning. In *Proc. ICML 2018* (Vol. 80, pp. 2127–2136).
6. Kingma, D.P., & Welling, M. (2014). Auto-Encoding Variational Bayes. *arXiv:1312.6114*.
7. Kipf, T.N., & Welling, M. (2017). Semi-Supervised Classification with Graph Convolutional Networks. In *Proc. ICLR 2017*.
8. Lundberg, S.M., & Lee, S.I. (2017). A unified approach to interpreting model predictions. In *Proc. NeurIPS 2017*, 4765–4774.
9. Luo, H., Huang, J., Ju, H., *et al.* (2025). Multimodal multi-instance evidence fusion neural networks for cancer survival prediction. *Scientific Reports*, **15**, 10470.
10. Ozmen, F., Ozmen, T.Y., Ors, A., *et al.* (2025). Single-cell RNA sequencing reveals different cellular states in malignant cells and the tumor microenvironment in primary and metastatic ER-positive breast cancer. *npj Breast Cancer*, **11**(1), 95.
11. Parker, J.S., Mullins, M., Cheang, M.C., *et al.* (2009). Supervised risk predictor of breast cancer based on intrinsic subtypes. *J. Clin. Oncol.*, **27**(8), 1160–1167.
12. Perou, C.M., Sørlie, T., Eisen, M.B., *et al.* (2000). Molecular portraits of human breast tumours. *Nature*, **406**(6797), 747–752.
13. Poirion, O.B., Chaudhary, K., Huang, S., *et al.* (2021). DeepProg: an ensemble of deep-learning and machine-learning models for prognosis prediction using multi-omics data. *Genome Medicine*, **13**(1), 112.
14. Rediti, M., Stenbeck, L., Dupont, F., *et al.* (2024). Spatial transcriptomics reveals substantial heterogeneity in triple-negative breast cancer with potential clinical implications. *Nat. Commun.*, **15**(1), 10232.
15. Sharma, A., Debik, J., Naume, B., *et al.* (2024). Comprehensive multi-omics analysis of breast cancer reveals distinct long-term prognostic subtypes. *Oncogenesis*, **13**(1), 22.
16. Vaswani, A., Shazeer, N., Parmar, N., *et al.* (2017). Attention is all you need. In *Proc. NeurIPS 2017*, 5998–6008.
17. Yao, J., Zhu, X., Jonnagaddala, J., Hawkins, N., & Huang, J. (2020). Whole slide images based cancer survival prediction using attention guided deep multiple instance learning networks. *Med. Image Anal.*, **65**, 101789.