

Multimodal Data Integration for Oncology in the Era of Deep Neural Networks: A Review

Asim Waqas 1,†,* , Aakash Tripathi 1,† , Ravi P. Ramachandran 2 , Paul Stewart 1 , and Ghulam Rasool 1

 1 Department of Machine Learning, Moffitt Cancer Center, Tampa, Florida, USA 2 Department of Electrical and Computer Engineering, Rowan University, Glassboro,

New Jersey, USA
Correspondence*:
Asim Waqas
asim.wagas@moffitt.org

†These authors contributed equally to this work and share first authorship

ABSTRACT

Cancer research encompasses data across various scales, modalities, and resolutions, 3 from screening and diagnostic imaging to digitized histopathology slides to various types of molecular data and clinical records. The integration of these diverse data types for personalized cancer care and predictive modeling holds the promise of enhancing the accuracy and reliability of cancer screening, diagnosis, and treatment. Traditional analytical methods, which often focus on isolated or unimodal information, fall short of capturing the complex and heterogeneous nature of cancer data. The advent of deep neural networks has spurred the development of sophisticated multimodal data fusion techniques capable of extracting and synthesizing information from disparate sources. 12 Among these, Graph Neural Networks (GNNs) and Transformers have emerged as powerful tools for multimodal learning, demonstrating significant success. This review presents the foundational principles of multimodal learning including oncology data modalities, taxonomy of multimodal learning, and fusion strategies. We delve into the recent advancements in GNNs and Transformers for the fusion of multimodal data in oncology, spotlighting key studies and their pivotal findings. We discuss the unique challenges of multimodal learning, such as data heterogeneity and integration complexities, alongside the opportunities it presents for a more nuanced and comprehensive understanding of cancer. Finally, 19 we present some of the latest comprehensive multimodal pan-cancer data sources. By 21 surveying the landscape of multimodal data integration in oncology, our goal is to underline the transformative potential of multimodal GNNs and Transformers. Through technological advancements and the methodological innovations presented in this review, we aim to 23 chart a course for future research in this promising field. This review may be the first that highlights the current state of multimodal modeling applications in cancer using GNNs 25 and transformers, presents comprehensive multimodal oncology data sources, and sets the stage for multimodal evolution, encouraging further exploration and development in 27 personalized cancer care.

Keywords: Multimodal, Graph Neural Networks, Transformers, Oncology, Deep Learning

62

63

64

65

1 INTRODUCTION

Cancer represents a significant global health challenge, characterized by the uncontrolled growth of abnormal cells, leading to millions of deaths annually. In 2023, the United States had around 1.9 31 million new cancer diagnoses, with cancer being the second leading cause of death and anticipated 32 to result in approximately 1670 deaths daily (Siegel et al., 2023). However, advancements in 33 oncology research hold the promise of preventing nearly 42% of these cases through early detection 34 and lifestyle modifications. The complexity of cancer, involving intricate changes at both the 35 microscopic and macroscopic levels, requires innovative approaches to its understanding and 36 management. In recent years, the application of machine learning (ML) techniques, especially 37 deep learning (DL), has emerged as a transformative force in oncology. DL employs deep neural 38 networks to analyze vast datasets, offering unprecedented insights into cancer's development 39 and progression (Çalışkan and Tazaki, 2023; Chen et al., 2023; Siam et al., 2023; Muhammad 40 41 et al., 2024; Talebi et al., 2024). This approach has led to the development of computer-aided diagnostic systems capable of detecting and classifying cancerous tissues in medical images, such 42 43 as mammograms and MRI scans, with increasing accuracy. Beyond imaging, DL also plays a crucial role in analyzing molecular data, aiding in the prediction of treatment responses, and the 44 identification of new biomarkers (Varlamova et al., 2024; Khan et al., 2023; Muhammad and Bria, 45 2023; Dera et al., 2021, 2019; Wagas et al., 2021; Barhoumi et al., 2023). DL methods can be 46 categorized based on the level of supervision involved. Supervised learning includes techniques 47 like Convolutional Neural Networks (CNNs) for tumor image classification and Recurrent Neural 48 49 Networks (RNNs) for predicting patient outcomes, both requiring labeled data (LeCun et al., 2015; 50 Iqbal et al., 2022, 2019). Unsupervised deep learning methods, such as Autoencoders and Generative 51 Adversarial Networks (GANs), learn from unlabeled data to perform tasks like clustering patients 52 based on gene expression profiles or generating synthetic medical images. Semi-supervised deep 53 learning methods, like Semi-Supervised GANs, leverage a mix of labeled and unlabeled data to enhance model performance when labeled medical data is limited. Self-supervised learning methods, 54 55 such as BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative 56 Pre-trained Transformer), use the structure of training data itself for supervision, enabling tasks like predicting patient outcomes or understanding the progression of cancer with limited labeled 57 examples. Reinforcement learning in cancer studies, exemplified by Deep Q-Networks (DQN) 58 59 and Proximal Policy Optimization (PPO), involves an agent learning optimal treatment strategies through rewards and penalties. 60

As the volume of oncology data continues to grow, DL stands at the forefront of this field, enhancing our understanding of cancer, improving diagnostic precision, predicting clinical outcomes, and paving the way for innovative treatments. This review explores the latest advancements in DL applications within oncology, highlighting its potential to revolutionize cancer care (Ghaffari Laleh et al., 2023; Chan et al., 2020; Tripathi et al., 2024a; Ibrahim et al., 2022).

Multimodal Learning (MML) enhances task accuracy and reliability by leveraging information from various data sources or modalities (Huang et al., 2021). This approach has witnessed a surge in popularity, as indicated by the growing body of MML-related publications (see Figure 1). By facilitating the fusion of multimodal data, such as radiological images, digitized pathology slides, molecular data, and electronic health records (EHR), MML offers a richer understanding

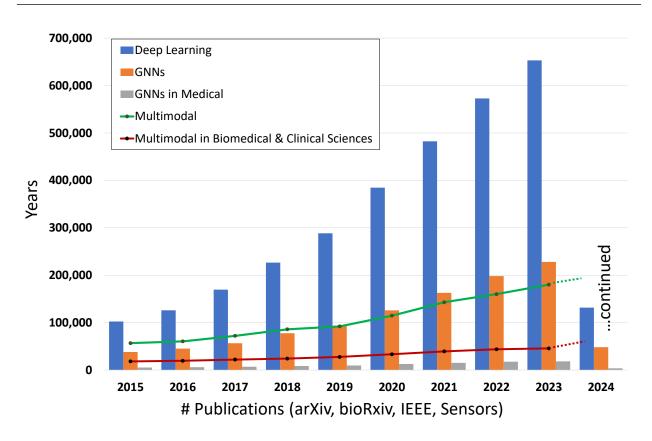


Figure 1. Number of publications involving DL, GNNs, GNNs in the medical domain, overall multimodal and multimodal in biomedical and clinical sciences in the period 2015-2024 (Hook et al., 2018).

71 of complex problems (Tripathi et al., 2024c). It enables the extraction and integration of relevant 72 features that might be overlooked when analyzing data modalities separately. Recent advancements in MML, powered by Deep Neural Networks (DNNs), have shown remarkable capability in learning 73 74 from diverse data sources, including computer vision (CV) and natural language processing (NLP) 75 (Achiam et al., 2023; Bommasani et al., 2022). Prominent multimodal foundation models such as Contrastive Language-Image Pretraining (CLIP) and Generative Pretraining Transformer (GPT-4) 76 by OpenAI have set new benchmarks in the field (Radford et al., 2021; Achiam et al., 2023). 77 Additionally, the Foundational Language And Vision Alignment Model (FLAVA) represents another 78 79 significant stride, merging vision and language representation learning to facilitate multimodal reasoning (Singh et al., 2022). Within the realm of oncology, innovative applications of MML 80 are emerging. The RadGenNets model, for instance, integrates clinical and genomics data with 81 Positron Emission Tomography (PET) scans and gene mutation data, employing a combination of 82 Convolutional Neural Networks (CNNs) and Dense Neural Networks to predict gene mutations 83 in Non-small cell lung cancer (NSCLC) patients (Tripathi et al., 2022). Moreover, GNNs and 84 Transformers are being explored for a variety of oncology-related tasks, such as tumor classification 85 (Khan et al., 2020), prognosis prediction (Schulz et al., 2021), and assessing treatment response 86 (Joo et al., 2021). 87

- 88 Recent literature has seen an influx of survey and review articles exploring MML (Boehm et al.,
- 89 2021; Xu et al., 2023; Baltrušaitis et al., 2018; Ektefaie et al., 2023b; Hartsock and Rasool, 2024).
- 90 These works have provided valuable insights into the evolving landscape of MML, charting key
- 91 trends and challenges within the field. Despite this growing body of knowledge, there remains a
- 92 notable gap in the literature regarding the application of advanced multimodal DL models, such as
- 93 Graph Neural Networks (GNNs) and Transformers, in the domain of oncology. Our article aims to
- 94 fill this gap by offering the following contributions:
- Identifying large-scale MML approaches in oncology. We provide an overview of the state-of-the-art MML with a special focus on GNNs and Transformers for multimodal data fusion in oncology.
- 98 2. Highlighting the challenges and limitations of MML in oncology data fusion. We discuss the challenges and limitations of implementing multimodal data-fusion models in oncology, including the need for large datasets, the complexity of integrating diverse data types, data alignment, and missing data modalities and samples.
- 3. Providing a taxonomy for describing multimodal architectures. We present a comprehensive taxonomy for describing MML architectures, including both traditional ML and DL, to facilitate future research in this area.
- 4. *Identifying future directions for multimodal data fusion in oncology*. We identify GNNs and
 Transformers as potential solutions for comprehensive multimodal integration and present the
 associated challenges.
- By addressing these aspects, our article seeks to advance the understanding of MML's potential in oncology, paving the way for innovative solutions that could revolutionize cancer diagnosis and treatment through comprehensive data integration.
- Our paper is organized as follows. Section 2 covers the fundamentals of MML, including data
- 112 modalities, taxonomy, data fusion stages, and neural network architectures. Section 3 focuses
- on GNNs in MML, explaining graph data, learning on graphs, architectures, and applications to
- unimodal and multimodal oncology datasets. Section 4 discusses Transformers in MML, including
- architecture, multimodal Transformers, applications to oncology datasets, and methods of fusing
- data modalities. Section 5 highlights challenges in MML, such as data availability, alignment,
- 117 generalization, missing data, explainability, and others. Section 6 provides information on data
- 118 sources. Finally, we conclude by emphasizing the promise of integrating data across modalities and
- 119 the need for scalable DL frameworks with desirable properties.

2 FUNDAMENTALS OF MULTIMODAL LEARNING (MML)

120 2.1 Data Modalities in Oncology

- A data *modality* represents the expression of an entity or a particular form of sensory perception,
- 122 such as the characters' visual actions, sounds of spoken dialogues, or the background music
- 123 (Sleeman IV et al., 2022). A collective notion of these modalities is called multi-modality
- 124 (Baltrušaitis et al., 2018). Traditional data analysis and ML methods to study cancer data use
- single data modalities (e.g., EHR (Miotto et al., 2016), radiology (Waqas et al., 2021), pathology

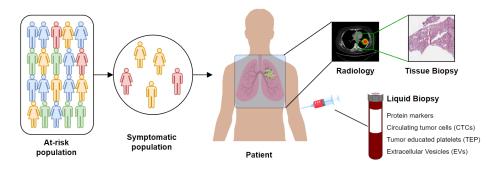


Figure 2a. An overview of data collected from population to a tissue.

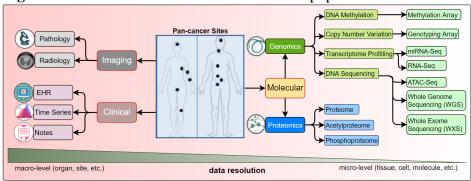


Figure 2b. Detailed look into data modalities acquired for cancer care.

Figure 2. We present various data modalities that capture specific aspects of cancer at different scales. For example, radiological images capture organ or sub-organ level abnormalities, while tissue analysis may provide changes in the cellular structure and morphology. On the other hand, various molecular data types may provide insights into genetic mutations and epigenetic changes.

(Litjens et al., 2017), or molecular, including genomics (Angermueller et al., 2017), transcriptomics 126 (Yousefi et al., 2017), proteomics (Wang et al., 2017), etc.). However, the data is inherently 127 multimodal, as it includes information from multiple sources or modalities that are related in 128 many ways. Figure 2a provides a view of multiple modalities of cancer at various scales, from the 129 population level to single-cell analysis. Oncology data can be broadly classified into 3 categories: 130 clinical, molecular, and imaging, where each category provides complementary information about 131 the patient's disease. Figure 2b highlights different clinical, molecular, and imaging modalities. 132 Multimodal analysis endeavors to gain holistic insights into the disease process using multimodal 133 data. 134

135 2.1.1 Molecular Data

Molecular data modalities provide information about the underlying genetic changes and alterations in the cancer cells (Liu et al., 2021). Efforts toward integrating molecular data resulted in the *multi-omics* research field (Waqas et al., 2024a). Two principal areas of molecular analysis in oncology are proteomics and genomics. *Proteomics* is the study of proteins and their changes in response to cancer, and it provides information about the biological processes taking place in cancer cells. *Genomics* is the study of the entire genome of cancer cells, including changes in DNA sequence, gene expression, and structural variations (Boehm et al., 2021). Other molecular

- 143 modalities include transcriptomics, pathomics, radiomics and their combinations, radiogenomics,
- and proteogenomics. Many publicly available datasets provide access to molecular data, including
- 145 the Proteomics Data Commons for proteomics data and the Genome Data Commons for genetic
- 146 data (Thangudu et al., 2020; Grossman et al., 2016).

147 2.1.2 Imaging Data

- 148 Imaging modalities play a crucial role in diagnosing and monitoring cancer. The imaging category
- 149 can be divided into 2 main categories: (1) radiological imaging and (2) digitized histopathology
- 150 slides, referred to as Whole Slide Imaging (WSI). Radiological imaging encompasses various
- 151 techniques such as X-rays, CT scans, MRI, PET, and others, which provide information about the
- 152 location and extent of cancer within the body. These images can be used to determine the size and
- shape of a tumor, monitor its growth, and assess the effectiveness of treatments. *Histopathological*
- imaging is the examination of tissue samples obtained through biopsy or surgery (Rowe and Pomper,
- 155 2022; Waqas et al., 2023). Digitized slides, saved as WSIs, provide detailed information about
- 156 the micro-structural changes in cancer cells and can be used to diagnose cancer and determine its
- 157 subtype.

158 2.1.3 Clinical Data

- 159 Clinical data provides information about the patient's medical history, physical examination,
- and laboratory results, saved in the patient's electronic health records (EHR) at the clinic. EHR
- 161 consists of digital records of a patient's health information stored in a centralized database. These
- 162 records provide a comprehensive view of a patient's medical history, past diagnoses, treatments,
- laboratory test results, and other information, which helps clinicians understand the disease (Asan
- 164 et al., 2018). Within EHR, time-series data may refer to the clinical data recorded over time, such as
- repeated blood tests, lab values, or physical attributes. Such data informs the changes in the patient's
- 166 condition and monitors the disease progression (Quinn et al., 2019).

167 **2.2 Taxonomy of MML**

- We follow the taxonomy proposed by Sleeman IV et al. (2022) (see Figure 3), which defines 5
- 169 main stages of multimodal classification: preprocessing, feature extraction, data fusion, primary
- 170 learner, and final classifier, as given below:

171 2.2.1 Pre-processing

- 172 Pre-processing involves modifying the input data to a suitable format before feeding it into the
- model for training. It includes data cleaning, normalization, class balancing, and augmentation. Data
- 174 cleaning removes unwanted noise or bias, errors, and missing data points (Al-jabery et al., 2020).
- 175 Normalization scales the input data within a specific range to ensure that each modality contributes
- equally to the training (Gonzalez Zelaya, 2019). Class balancing is done in cases where one class
- may have a significantly larger number of samples than another, resulting in a model bias toward the
- dominant class. Data augmentation artificially increases the size of the dataset by generating new
- 179 samples based on the existing data to improve the model's robustness and generalizability (Al-jabery
- 180 et al., 2020).

181 2.2.2 Feature Extraction

Different data modalities may have different features, and extracting relevant features may 182 improve model learning. Several manual and automated feature engineering techniques generate 183 representations (or *embeddings*) for each data modality. Feature engineering involves designing 184 features relevant to the task and extracting them from the input data. This can be time-consuming but 185 may allow the model to incorporate prior knowledge about the problem. Text encoding techniques, 186 such as bag-of-words, word embeddings, and topic models (Devlin et al., 2019; Zhuang et al., 2021), 187 transform textual data into a numerical representation, which can be used as input to an ML model 188 (Wang et al., 2020a). In DL, feature extraction is learned automatically during model training(Dara 189 and Tumma, 2018). 190

191 2.2.3 Data Fusion

Data fusion combines raw features, extracted features, or class prediction vectors from multiple modalities to create a single data representation. Fusion enables the model to use the complementary information provided by each modality and improve its learning. Data fusion can be done using early, late, or intermediate fusion. Section 2.3 discusses these fusion stages. The choice of fusion technique depends on the characteristics of the data and the specific problem being addressed (Jiang et al., 2022a).

198 2.2.4 Primary Learner

The primary learner stage is training the model on the pre-processed data or extracted features. 199 Depending on the problem and data, the primary learner can be implemented using various 200 ML techniques. DNNs are a popular choice for primary learners in MML because they can 201 automatically learn high-level representations from the input data and have demonstrated state-of-202 the-art performance in many applications. CNNs are often used for image and video data, while 203 recurrent neural networks (RNNs) and Transformers are commonly used for text and sequential 204 data. The primary learner can be implemented independently for each modality or shared between 205 modalities, depending on the problem and data. 206

207 2.2.5 Final Classifier

The final stage of MML is the classifier, which produces category labels or class scores and 208 can be trained on the output of the primary learner or the fused data. The final classifier can be 209 implemented using a shallow neural network, a decision tree, or an ensemble model (Sleeman IV 210 et al., 2022). Ensemble methods, such as stacking or boosting, are often used to improve and 211 robustify the performance of the final classifier. Stacking involves training multiple models and then 212 combining their predictions at the output stage, while boosting involves repeatedly training weak 213 learners and adjusting their weights based on the errors made by previous learners (Borisov et al., 214 215 2022).

216 **2.3 Data Fusion Strategies**

Fusion in MML can be performed at different levels, including early (feature level), intermediate (model level), or late (decision level) stages, as illustrated in figure 3. Each fusion stage has its

advantages and challenges, and the choice of fusion stage depends on the characteristics of the data and the task.

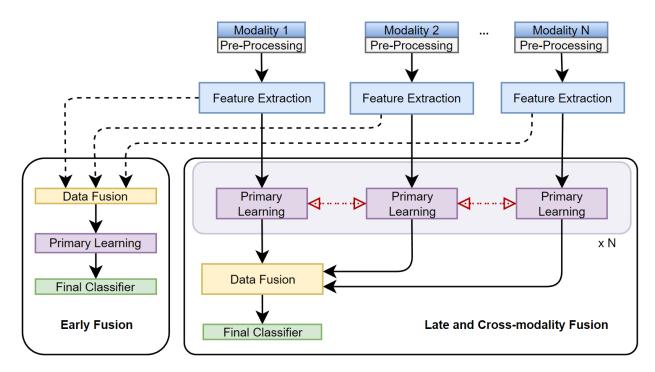


Figure 3. Taxonomy, stages, and techniques of multimodal data fusion are presented. *Early, late, cross-modality* fusion methods integrate individual data modalities (or extracted features) *before, after,* or *at* the primary learning step, respectively.

221 2.3.1 Early Fusion

222

223

224

225

226227

228

229

230

231

232

The early fusion involves merging features extracted from different data modalities into a single feature vector before model training. The feature vectors of the different modalities are combined into a single vector, which is used as the input to the ML model (Sleeman IV et al., 2022). This approach can be used when the modalities have complementary information and can be easily aligned, such as combining visual and audio features in a video analysis application. The main challenge with early fusion is ensuring that the features extracted from different modalities are compatible and provide complementary information.

2.3.2 Intermediate Fusion

Intermediate fusion involves training separate models for each data modality and then combining the outputs of these models for inference/prediction (Sleeman IV et al., 2022). This approach is suitable when the data modalities are independent of each other and cannot be easily combined at the feature level using average, weighted average, or other methods. The main challenge with intermediate fusion is selecting an appropriate method for combining the output of different models.

235 2.3.3 Late Fusion

- In late fusion, the output of each modality-specific model is used to make a decision independently.
- 237 All decisions are later combined to make a final decision. This approach is suitable when the
- 238 modalities provide complementary information but are not necessarily independent of each other.
- 239 The main challenge with late fusion is selecting an appropriate method for combining individual
- 240 predictions. This can be done using majority voting, weighted voting, or employing other ML
- 241 models.

261

262

263

264

265

266

267

268

269

270

242 2.4 MML for Oncology Datasets

Syed et al. (2021) used a Random Forest classifier to fuse radiology image representations learned 243 from the singular value decomposition method with the textual annotation representation learned 244 from the fastText algorithm for prostate and lung cancer patients. Liu et al. (2022) proposed a 245 hybrid DL framework for combining breast cancer patients' genomic and pathology data using 246 fully-connected (FC) network for genomic data, CNN for radiology data and a Simulated Annealing 247 algorithm for late fusion. Multiview multimodal network (MVMM-Net) (Song et al., 2021a) 248 249 combined 2 different modalities (low-energy and dual-energy subtracted) from contrast-enhanced spectral mammography images, each learned through CNN and late-fusion through FC network in 250 breast cancer detection task. Yap et al. (2018) used a late-fusion method to fuse image representations 251 from ResNet50 and clinical representations from a random forest model for a multimodal skin lesion 252 classification task. An award-winning work (Ma and Jia, 2020) on brain tumor grade classification 253 adopted the late-fusion method (concatenation) for fusing outputs from two CNNs (radiology and 254 pathology images). SeNMo, a self-normalizing deep learning model has shown that integrative 255 analysis on 33 cancers having five different molecular (multi-omics) data modalities can improve the 256 patient outcome predictions and primary cancer type classification (Waqas et al., 2024a). Recently, 257 GNNs-based pan-squamous cell carcinoma analysis on lung, bladder, cervicall, esophageal, and 258 head and neck cancers has outperformed different classical and deep learning models (Waqas et al., 259 2024b). 260

The single-cell unimodal data alignment is one technique in MML. Jansen et al. (2019) devised an approach (SOMatic) to combine ATAC-seq regions with RNA-seq genes using self-organizing maps. Single-Cell data Integration via Matching (SCIM) matched cells in multiple datasets in low-dimensional latent space using autoencoder (AEs) (Stark et al., 2020). Graph-linked unified embedding (GLUE) model learned regulatory interactions across omics layers and aligned the cells using variational AEs (Cao and Gao, 2022). These aforementioned methods cannot incorporate high-order interactions among cells or different modalities. Single-cell data integration using multiple modalities is mostly based on AEs (scDART (Zhang et al., 2022b), Cross-modal Autoencoders (Yang et al., 2021a), Mutual Information Learning for Integration of Single Cell Omics Data (SMILE) (Xu et al., 2022)).

3 GRAPH NEURAL NETWORKS (GNNs) IN MULTIMODAL LEARNING

Graphs are commonly used to represent the relational connectivity of any system that has interacting entities (Li et al., 2022a). Graphs have been used in various fields, such as to study brain

TO 1 1 1	D C	D' 1	•	α \cdot α
Table I	References	I hechieced	1n	Section 7
Table 1.	1 CICICIOCO	Discusseu	111	occuon 4.

Sections		References Discussed
Data Modalities in Oncology	Molecular	Liu et al. (2021), Waqas et al. (2024a), Boehm et al. (2021), Thangudu et al. (2020), Grossman et al. (2016)
	Imaging Clinical	Rowe and Pomper (2022), Waqas et al. (2023) Asan et al. (2018), Quinn et al. (2019)
Taxonomy of MML		Sleeman IV et al. (2022), Al-jabery et al. (2020), Gonzalez Zelaya (2019), Devlin et al. (2019), Zhuang et al. (2021), Wang et al. (2020a), Dara and Tumma (2018), Jiang et al. (2022a), Borisov et al. (2022)
Data Fusion Strategies		Sleeman IV et al. (2022)
MML for Oncology Datasets		Syed et al. (2021), Liu et al. (2022), Song et al. (2021a), Yap et al. (2018), Ma and Jia (2020), Waqas et al. (2024a), Waqas et al. (2024b), Jansen et al. (2019), Stark et al. (2020), Cao and Gao (2022), Zhang et al. (2022b), Yang et al. (2021a), Xu et al. (2022)

networks (Farooq et al., 2019), analyze driving maps (Derrow-Pinion et al., 2021), and explore the structure of DNNs themselves (Waqas et al., 2022). GNNs are specifically designed to process data represented as a graph (Waikhom and Patgiri, 2022), which makes them well-suited for analyzing multimodal oncology data as each data modality (or sub-modality) can be considered as a single node and the structures/patterns that exist between data modalities can be modeled as edges (Ektefaie et al., 2023b).

279 3.1 The Graph Data

280 A graph is represented as G=(V,E) having node-set $V=\{v_1,v_2,...,v_n\}$, where node v has feature 281 vector \mathbf{x}_v , and edge set $E=\{(v_i,v_j)\mid v_i,v_j\in V\}$. The neighborhood of node v is defined as 282 $N(v)=\{u\mid (u,v)\in E\}$.

283 3.1.1 Graph Types

284

285 286

287

288

289 290

291

292

293

294

As illustrated in figure 4(a), the common types of graphs include undirected, directed, homogeneous, heterogeneous, static, dynamic, unattributed, and attributed. *Undirected graphs* comprise undirected edges, i.e., the direction of relation is not important between any ordered pair of nodes. In the *directed graphs*, the nodes have a directional relationship(s). Homogeneous graphs have the same type of nodes, whereas heterogeneous graphs have different types of nodes within a single graph (Yang et al., 2021b). Static graphs do not change over time with respect to the existence of edges and nodes. In contrast, dynamic graphs change over time, resulting in changes in structure, attributes, and node relationships. *Unattributed graphs* have unweighted edges, indicating that the weighted value for all edges in a graph is the same, i.e., 1 if present, 0 if absent. *Attributed graphs* have different edge weights that capture the strength of relational importance (Waikhom and Patgiri, 2022).

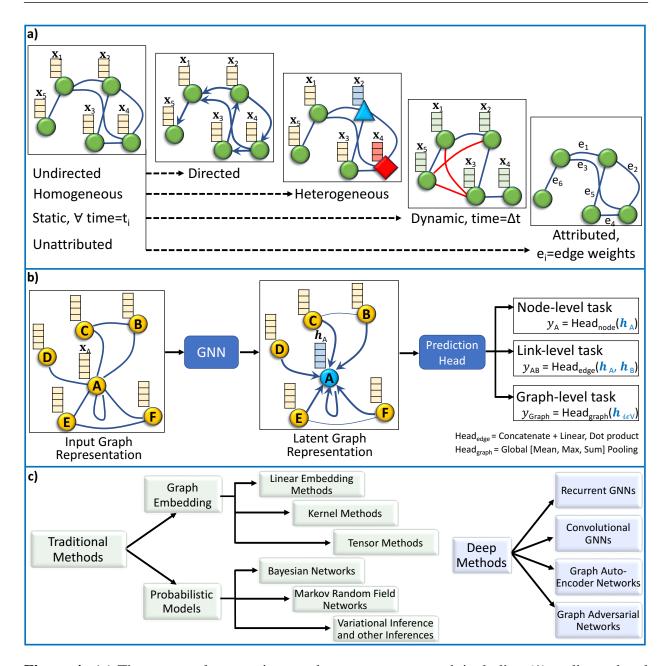


Figure 4. (a) The commonly occurring graph types are presented, including (1) undirected and directed, (2) homogeneous and heterogeneous, (3) dynamic and static, (4) attributed (edges) and unattributed. (b) Three different types of tasks performed using the graph data are presented and include (1) node-level, (2) link-level, and (3) graph-level analyses. (c) Various categories of representation learning for graphs are presented.

5 3.1.2 Tasks for Graph Data

296

297

298

In figure 4(b), we present 3 major types of tasks defined on graphs, including (1) *node-level tasks* - these may include node classification, regression, clustering, attributions, and generation, (2) *edge-level task* - edge classification and prediction (presence or absence) are 2 common edge-level

tasks, (3) *graph-level tasks* - these tasks involve predictions on the graph level, such as graph classification and generation.

301 3.2 ML for Graph Data

Representing data as graphs can enable capturing and encoding the relationships among entities of the samples (Wu et al., 2020). Based on the way the nodes are encoded, representation learning on graphs can be categorized into the traditional (or shallow) and DNN-based methods, as illustrated in Figure 4(c) (Jiao et al., 2022; Wu et al., 2020).

306 3.2.1 Traditional (Shallow) Methods

These methods usually employ classical ML methods, and their two categories commonly found 307 in the literature are graph embedding and probabilistic methods. Graph embedding methods 308 represent a graph with low-dimensional vectors (graph embedding and node embedding), preserving 309 the structural properties of the graph. The learning tasks in graph embedding usually involve 310 dimensionality reduction through linear (principal component or discriminant analysis), kernel 311 (nonlinear mapping), or tensor (higher-order structures) methods (Jiao et al., 2022). Probabilistic 312 graphical methods use graph data to represent probability distribution, where nodes are considered random variables, and edges depict the probability relations among nodes (Jiao et al., 2022). 314 Bayesian networks, Markov's networks, variational inference, variable elimination, and others are 315 used in probabilistic methods (Jiao et al., 2022).

317 3.2.2 DNN-based Methods - GNNs

327

328

329

330

331

318 GNNs are gaining popularity in the ML community, as evident from figure 1. In GNNs, the information aggregation from the neighborhood is fused into a node's representation. Traditional 319 DL methods such as CNNs and their variants have shown remarkable success in processing the 320 321 data in Euclidean space; however, they fail to perform well when faced with non-Euclidean or relational datasets. Compared to CNNs, where the locality of the nodes in the input is fixed, GNNs 322 have no canonical ordering of the neighborhood of a node. They can learn the given task for any 323 324 permutation of the input data, as depicted in figure 5. GNNs often employ a message-passing 325 mechanism in which a node's representation is derived from its neighbors' representations via a recursive computation. The message passing for a GNN is given as follows:

$$\mathbf{h}_v^{(l+1)} = \sigma \left(W_l \sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l)}}{|N(v)|} + B_l \mathbf{h}_v^{(l)} \right)$$

$$\tag{1}$$

where $h_v^{(l+1)}$ is the updated embedding of node v after l+1 layer, σ is the non-linear function (e.g., rectified linear unit or ReLU), $h_u^{(l)}$ and $h_v^{(l)}$ represent the embeddings of nodes u and v at layer l. W_l and B_l are the trainable weight matrices for neighborhood aggregation and (self)hidden vector transformation, respectively. The message passing can encode high-order structural information in node embedding through multiple aggregation layers. GNNs smooth the features by aggregating neighbors' embedding and filter eigenvalues of graph Laplacian, which provides an extra denoising mechanism (Ma et al., 2021b). GNNs comprise multiple permutation equivariant and invariant

functions, and they can handle heterogeneous data (Jin et al., 2022). As described earlier, traditional ML models deal with Euclidean data. In oncology data, the correlations may not exist in Euclidean space; instead, its features may be highly correlated in the non-Euclidean space (Yi et al., 2022). Based on the information fusion and aggregation methodology, GNNs-based deep methods are classified into the following:

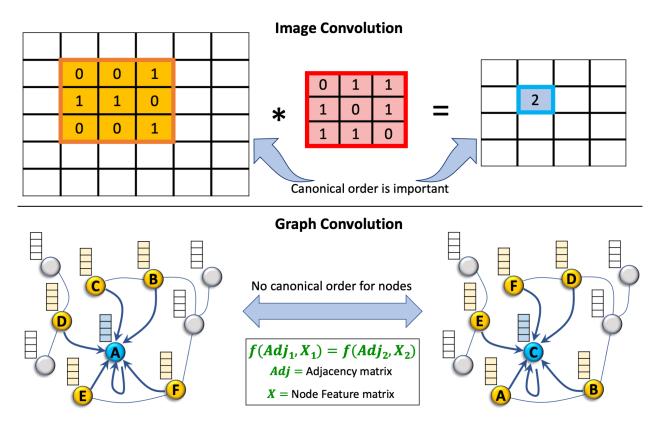


Figure 5. Convolution operation for graphs vs. image data. The canonical order of the input is important in CNNs, whereas in GNNs, the order of the input nodes is not important. From the convolution operation perspective, CNNs can be considered a subset of GNNs (Hamilton, 2020).

3.2.2.1 Recurrent GNNs

339

340

341342

343

344 345 RecGNNs are built on top of the standard Recurrent Neural Network (RNN) by combining with GNN. RecGNNs can operate on graphs with variable sizes and topologies. The recurrent component of the RecGNN captures temporal dependencies and learns latent states over time, whereas the GNN component captures the local structure. The information fusion process is repeated a fixed number of times until an equilibrium or the desired state is achieved (Hamilton et al., 2017). RecGNNs employ the model given by:

$$\mathbf{h}_{v}^{(l+1)} = \operatorname{RecNN}\left(\mathbf{h}_{u}^{(l)}, \mathbf{Msg}_{N(v)}^{(l)}\right), \tag{2}$$

where, RecNN is any RNN, and $Msg_{N(v)}^{(l)}$ is the neighborhood message-passing at layer l.

347 3.2.2.2 Convolutional GNNs

ConvGNNs undertake the convolution operation on graphs by aggregating neighboring nodes' embeddings through a stack of multiple layers. ConvGNNs use the symmetric and normalized summation of the neighborhood and self-loops for updating the node embeddings given by:

$$\mathbf{h}_{v}^{(l+1)} = \sigma \left(W_{l} \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_{v}}{\sqrt{|N(v)||N(u)|}} \right). \tag{3}$$

The ConvGNN can be spatial or spectral, depending on the type of convolution they 351 implement. Convolution in spatial ConvGNNs involves taking a weighted average of the 352 353 neighboring vertices. Examples of spatial ConvGNNs include GraphSAGE(Hamilton et al., 2017), Message Passing Neural Network (MPNN)(Gilmer et al., 2017), and Graph Attention Network 354 (GAT) (Veličković et al., 2017). The spectral ConvGNNs operate in the spectral domain by using 355 the eigendecomposition of the graph Laplacian matrix. The convolution operation is performed 356 on the eigenvalues, which can be high-dimensional. Popular spectral ConvGNNs are ChebNet 357 (Defferrard et al., 2016) and Graph Convolutional Network (GCN)(Kipf and Welling, 2016). 358 An interesting aspect of these approaches is representational containment, which is defined as: 359 convolution \subseteq attention \subseteq message passing. 360

3.2.2.3 Graph Auto-Encoder Networks (GAEs)

GAEs are unsupervised graph learning networks for dimensionality reduction, anomaly detection, and graph generation. They are built on top of the standard AEs to work with graph data. The encoder component of the GAE maps the input graph to a low-dimensional latent space, while the decoder component maps the latent space back to the original graph (Park et al., 2021).

366 3.2.2.4 Graph Adversarial Networks (GraphANs)

Based on Generative Adversarial Networks, GraphANs are designed to work with graph-structured data and can learn to generate new graphs with similar properties to the input data. The generator component of the GraphAN maps a random noise vector to a new graph, while the discriminator component tries to distinguish between the generated vs. the actual input. The generator generates graphs to fool the discriminator, while the discriminator tries to classify the given graph as real or generated.

373 3.2.2.5 Other GNNs

361

Other categories of GNNs may include scalable GNNs (Ma et al., 2019), dynamic GNNs (Sankar et al., 2018), hypergraph GNNs (Bai et al., 2021), heterogeneous GNNs (Wei et al., 2019), and many others (Ma and Tang, 2021).

377 3.2.3 Graph-based Reinforcement Learning

GNNs have been combined with Reinforcement Learning (RL) to solve complex problems involving graph-structured data (Jiang et al., 2018). GNNs enable RL agents to effectively process and reason about relational information in environments represented as graphs (Nie et al., 2023).

This combination has shown promise in various domains, including multi-agent systems, robotics, and combinatorial optimization (Fathinezhad et al., 2023; Almasan et al., 2022). However, the use of Graph-based RL on cancer data is still less-explored area of research.

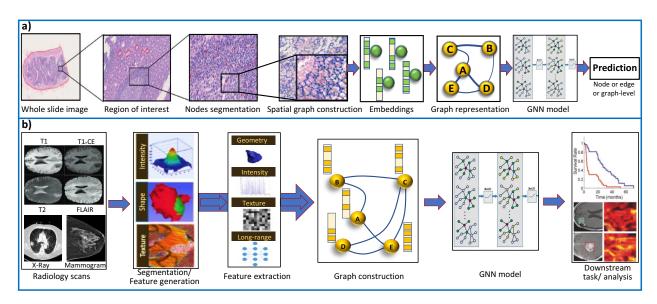


Figure 6. (a) Data processing pipeline for histopathology images using GNNs (Chen et al., 2020). (b) Graph processing pipeline on radiology data. Adapted from (Singh et al., 2021).

384 3.3 GNNs and ML using Unimodal Oncology Datasets

385 3.3.1 Pathology Datasets

386

387

388 389

390

391

392

393

394

395

396 397

398

399

400

401

402

403

Traditionally, CNN-based models are used to learn features from digital pathology data (Iqbal et al., 2022). However, unlike GNNs, CNNs fail to capture the global contextual information important in the tissue phenotypical and structural micro and macro environment (Ahmedt-Aristizabal et al., 2022). For using histology images in GNNs, the cells, tissue regions, or image patches are depicted as nodes. The relations and interactions among these nodes are represented as (un)weighted edges. Usually, a graph of the patient histology slide is used along with a patient-level label for training a GNN, as illustrated in Figure 6(a). Here, we review a few GNN-based pathology publications representative of a trove of works in this field. Histographs (Anand et al., 2020) used breast cancer histology data to distinguish cancerous and non-cancerous images. Pre-trained VGG-UNet was used for nuclei detection, micro-features of the nuclei were used as node features, and Euclidean distance among nuclei was incorporated as edge features. The resulting cell graphs were used to train the GCN-based robust spatial filtering (RSF) model, which performed superior to the CNN-based classification, citewang 2020 weakly analyzed grade classification in tissue micro-arrays of prostate cancer using the weakly-supervised technique on a variant of GraphSAGE with selfattention pooling (SAGPool). Cell-Graph Signature ($CG_{signature}$) (Wang et al., 2022) predicted patient survival in gastric cancer using cell-graphs of multiplexed immunohistochemistry images processed through two types of GNNs (GCNs and GINs) with two types of pooling (SAGPool, TopKPool). Besides the above-mentioned cell graphs, there is an elaborate review of GNN-based tissue graphs or patch-graphs methods implemented on unimodal pathology cancer data given in

405 (Ahmedt-Aristizabal et al., 2022). Instead of individual cell- and tissue-graphs, a combination of the 406 multilevel information in histology slides can help understand the intrinsic features of the disease.

407 3.3.2 Radiology Datasets

408 GNNs have been used in radiology-based cancer data for segmentation, classification, and prediction tasks, especially on X-rays, mammograms, MRI, PET, and CT scans. Figure 6(b) 409 illustrates a general pipeline of using radiology-based data to train GNNs. Here we give a non-410 exhaustive review of GNNs-based works on radiological oncology data as a single modality input. 411 Mo et al. (2020) proposed a framework that improved the liver cancer lesion segmentation in the 412 MRI-T1WI scans through guided learning of MRI-T2WI modality priors. Learned embeddings from 413 fully convolutional networks on separate MRI modalities are projected into the graph domain for 414 learning by GCNs through the co-attention mechanism and finally to get the refined segmentation 415 by re-projection. Radiologists usually review radiology images by zooming into the region of 416 interest (ROIs) on high-resolution monitors. Du et al. (2019) used a hierarchical GNN framework to 417 automatically zoom into the abnormal lesion region of the mammograms and classify breast cancer. 418 The pre-trained CNN model extracts image features, whereas a GAT model is used to classify the 419 nodes for deciding whether to zoom in or not based on whether it is benign or malignant. Based 420 on the established knowledge that lymph nodes (LNs) have connected lymphatic system and LNs 421 cancer cells spread on certain pathways, Chao et al. (2020) proposed a lymph node gross tumor 422 volume learning framework. The framework was able to delineate the LN appearance as well as 423 the inter-LN relationship. The end-to-end learning framework was superior to the state-of-the-art 424 on esophageal cancer radiotherapy dataset. Tian et al. (2020) suggested interactive segmentation 425 of MRI scans of prostate cancer patients through a combination of CNN and two GCNs. CNN 426 427 model outputs a segmentation feature map of MRI, and the GCNs predict the prostate contour from this feature map. Saueressig et al. (2021) used GNNs to segment brain tumors in 3D MRI images, 428 formed by stacking different modalities of MRI (T1, T2, T1-CE, FLAIR) and representing them 429 as supervoxel graph. The authors reported that GraphSAGE-pool was best for segmenting brain 430 tumors. Besides radiology, a parallel field of radiomics has recently gained attraction. Radiomics is 431 the automated extraction of quantitative features from radiology scans. A survey of radiomics and 432 radiogenomic analysis on brain tumors is presented by Singh et al. (2021). 433

434 3.3.3 Molecular Datasets

435 Graphs are a natural choice for representing molecular data such as omic-centric (DNA, RNA, or proteins) or single-cell centric. Individual modalities are processed separately to generate graph 436 representations that are then processed through GNNs followed by the classifier to predict the 437 438 downstream task, as illustrated in Figure 7. One method of representing proteins as graphs is to depict the amino acid residue in the protein as the node and the relationship between residues 439 denoted by edge (Fout et al., 2017). The residue information is depicted as node embedding, whereas 440 441 the relational information between two residues is represented as the edge feature vector. Fout et al. (2017) used spatial ConvGNNs to predict interfaces between proteins which is important in drug 442 discovery problems. Deep predictor of drug-drug interactions (DPDDI) predicted the drug-drug 443 interactions using GCN followed by a 5-layer classical neural network (Feng et al., 2020). Molecular 444 pre-training graph net (MPG) is a powerful framework based on GNN and Bidirectional Encoder 445

- 446 Representations from Transformers (BERT) to learn drug-drug and drug-target interactions (Li et al.,
- 447 2021b). Graph-based Attention Model (GRAM) handled the data inefficiency by supplementing
- 448 EHRs with hierarchical knowledge in the medical ontology (Choi et al., 2017). A few recent works
- 449 have applied GNNs to single-cell data. scGCN is a knowledge transfer framework in single-cell
- 450 omics data such as mRNA or DNA (Song et al., 2021b). scGNN processed cell-cell relations through
- 451 GNNs for the task of missing-data imputation and cell clustering on single-cell RNA sequencing
- 452 (scRNA-seq) data (Wang et al., 2021a).

453 3.4 MML - Data Fusion at the Pre-Learning Stage

The first and most primitive form of MML is the pre-learning fusion (see Figure 3), where 454 features extracted from individual modalities of data are merged, and the fused representations are 455 then used for training the multimodal primary learner model. In the context of GNNs being the 456 primary learning model, the extraction step of individual modality representations can be hand-457 engineered (e.g., dimensionality reduction) or learned by DL models (e.g., CNNs, Transformers). 458 Cui et al. (2021) proposed a GNN-based early fusion framework to learn latent representations from 459 radiological and clinical modalities for Lymph node metastasis (LNM) prediction in esophageal 460 squamous cell carcinoma (ESCC). The extracted features from the two modalities using UNet and 461 CNN-based encoders were fused together with category-wise attention as node representation. The 462 message passing from conventional GAT and correlation-based GAT learned the neighborhood 463 weights. The attention attributes were used to update the final node features before classification 464 by a 3-layer fully connected network. For Autism spectrum disorder, Alzheimer's disease, and 465 ocular diseases, a multimodal learning framework called Edge-Variational GCN (EV-GCN) fuses 466 the radiology features extracted from fMRI images with clinical feature vectors for each patient 467 468 (Huang and Chung, 2020). An MLP-based pairwise association encoder is used to fuse the input feature vectors and to generate the edge weights of the population graph. The partially labeled 469 population graph is then processed through GCN layers to generate the diagnostic graph of patients. 470

471 3.5 MML - Data Fusion using Cross-Modality Learning

472 Cross-MML involves intermediate fusion and/or cross-links among the models being trained on individual modalities (see Figure 3). For this survey, we consider the GNN-based hierarchical 473 474 learning mechanisms as the cross-MML methods. Hierarchical frameworks involve learning for one 475 modality and using the learned latent embeddings in tandem with other data modalities sequentially 476 to get the final desired low-dimensional representations. Lian et al. (2022) used a sequential learning framework where tumor features learned from CT images using the ViT model were used as node 477 features of the patient population graph for subsequent processing by the GraphSAGE model. The 478 479 hierarchical learning from radiological and clinical data using Transformer-GNN outperformed the ResNet-Graph framework in survival prediction of early-stage NSCLC. scMoGNN is the first 480 method to apply GNNs in multimodal single-cell data integration using a cross-learning fusion-481 482 based GNN framework (Wen et al., 2022). Officially winning first place in modality prediction task at the NeurIPS 2021 competition, scMoGNN showed superior performance on various tasks 483 by using paired data to generate cell-feature graphs. Hierarchical cell-to-tissue-graph network 484 (HACT-Net) combined the low-level cell-graph features with the high-level tissue-graph features 485 through two hierarchical GINs on breast cancer multi-class prediction (Pati et al., 2020). Data 486

498

499 500

501

502 503

504

505

506

507

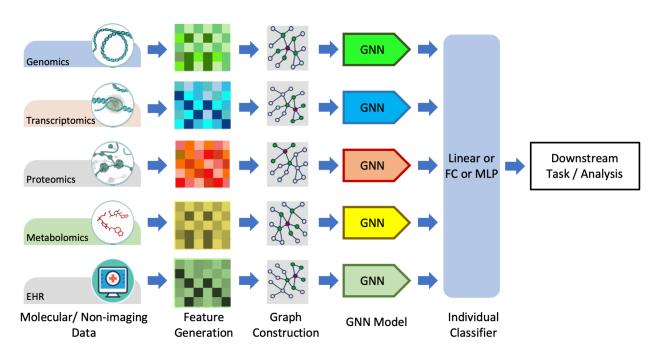


Figure 7. Graph data processing pipeline for non-imagery data, including molecular and textual data. Adapted from (Wang et al., 2021b). Abbreviations used: GNN - graph neural network, FC - Fully-Connected, MLP - Multi-Layer Perception.

imputation, a method of populating the missing values or false zero counts in single-cell data 487 mostly done using DL autoencoders (AE) architecture, has recently been accomplished using 488 489 GNNs. scGNN (Wang et al., 2021a) used imputation AE and graph AE in an iterative manner for imputation, and GraphSCI (Rao et al., 2021) used GCN with AE to impute the single-cell RNA-seq 490 data using the cross-learning fusion between the GCN and the AE networks. Clustering is a method 491 of characterizing cell types within a tissue sample. Graph-SCC clustered cells based on scRNA-seq 492 data through self-supervised cross-learning between GCN and a denoising AE network (Zeng et al., 493 2020). Recently, a multilayer GNN framework, Explainable Multilayer GNN (EMGNN), has been 494 proposed for cancer gene prediction tasks using multi-omics data from 16 different cancer types 495 (Chatzianastasis et al., 2023). 496

3.6 MML - Data Fusion in Post-Learning Regime

Post-learning fusion methods include processing individual data modalities and later fusing them for the downstream predictive task (Tortora et al., 2023). In the post-learning fusion paradigm, the hand-crafted features perform better than the deep features when the dimensionality of input data is low, and vice versa (Tortora et al., 2023). Many interesting GNN-based works involving the post-learning fusion mechanism have recently been published. Decagon used a multimodal approach on GCNs using proteins and drug interactions to predict exact side effects as a multirelational link prediction task (Zitnik et al., 2018). Drug—target affinity (DTA) experimented with four different flavors of GNNs (GCN, GAT, GIN, GAT-GCN) along with a CNN to fuse together molecular embeddings and protein sequences for predicting drug-target affinity (Nguyen et al., 2021). PathomicFusion combined the morphological features extracted from image patches (using CNNs),

cell-graph features from cell-graphs of histology images (GraphSAGE-based GCNs), and genomic features (using a feed-forward network) for survival prediction on glioma and clear cell renal cell 509 carcinoma (Chen et al., 2020). Shi et al. (2019) proposed a late-fusion technique to study screening 510 of cervical cancer at early stages by using CNNs to extract features from histology images, followed 511 by K-means clustering to generate graphs which are processed through two-layer GCN. BDR-CNN-512 GCN (batch normalized, dropout, rank-based pooling) used the same mammographic images to 513 extract image-level features using CNN and relation-aware features using GCN (Zhang et al., 2021). 514 The two feature sets are fused using a dot product followed by a trainable linear projection for breast 515 cancer classification. Under the umbrella of multi-omics data, many GNN-based frameworks have 516 been proposed recently. Molecular omics network(MOOMIN), a multi-modal heterogeneous GNN 517 to predict oncology drug combinations, processed molecular structure, protein features, and cell 518 lines through GCN-based encoders, followed by late-fusion using a bipartite drug-protein interaction 519 graph (Rozemberczki et al., 2022). Multi-omics graph convolutional networks (MOGONET) used 520 521 a GCN-GAN late fusion technique for the classification of four different diseases, including 522 three cancer types: breast, kidney, and glioma (Wang et al., 2021b). Leng et al. (2022) extended MOGONET to benchmark three multi-omics datasets on two different tasks using sixteen DL 523 networks and concluded that GAT-based GNN had the best classification performance. Multi-Omics 524 525 Graph Contrastive Learner (MOGCL) used graph structure and contrastive learning information to generate representations for improved downstream classification tasks on the breast cancer 526 527 multi-omics dataset using late-fusion (Rajadhyaksha and Chitkara, 2023). Similar to MOGCL, Park et al. (2022) developed a GNN-based multi-omics model that integrated mRNA expression, DNA 528 methylation, and DNA sequencing data for NSCLC diagnosis. 529

4 TRANSFORMERs IN MML

545

546

547

530 Transformers are attention-based DNN models originally proposed for NLP (Vaswani et al., 2017). 531 Transformers implement scaled dot-product of the input with itself and can process various types of data in parallel (Vaswani et al., 2017). Transformers can handle sequential data and learn long-range 532 533 dependencies, making them well-suited for tasks such as language translation, language modeling, 534 question answering, and many more (Otter et al., 2021). Unlike Recurrent Neural Networks (RNNs) and CNNs, Transformers use self-attention operations to weigh the importance of different input 535 tokens (or embeddings) at each time step. This allows them to handle sequences of arbitrary length 536 and to capture dependencies between input tokens that are far apart in the sequence (Vaswani et al., 537 2017). Transformers can be viewed as a type of GNN (Xu et al., 2023). Transformers are used 538 to process other data types, such as images (Dosovitskiy et al., 2020), audio (Zhang, 2020), and 539 time-series analysis (Ahmed et al., 2022b), resulting in a new wave of multi-modal applications. 540 Transformers can handle input sequences of different modalities in a unified way, using the same 541 self-attention mechanism, which processes the inputs as a fully connected graph (Xu et al., 2023). 542 This allows Transformers to capture complex dependencies between different modalities, such as 543 visual and textual information in visual question-answering (VQA) tasks (Ma et al., 2021a). 544

Pre-training Transformers on large amounts of data, using unsupervised or self-supervised learning, and then fine-tuning for specific downstream tasks, has led to the development of foundation models Boehm et al. (2021), such as BERT (Devlin et al., 2019), GPT (Radford et al., 2018),

Table 2. References Discussed in Section 3.

Sections		References Discussed
Graphs and GNNs		Li et al. (2022a), Farooq et al. (2019), Derrow-Pinion et al. (2021), Waqas et al. (2022), Waikhom and Patgiri (2022), Ektefaie et al. (2023b), Yang et al. (2021b), Waikhom and Patgiri (2022), Wu et al. (2020), Jiao et al. (2022), Wu et al. (2022), Ma et al. (2021b), Jin et al. (2022), Yi et al. (2022), Hamilton et al. (2017), Hamilton et al. (2017), Gilmer et al. (2017), Veličković et al. (2017), Defferrard et al. (2016), Kipf and Welling (2016), Park et al. (2021), Ma et al. (2019), Sankar et al. (2018), Bai et al. (2021), Wei et al. (2019), Ma and Tang (2021), Jiang et al. (2018), Nie et al. (2023), Fathinezhad et al. (2023), Almasan et al. (2022)
GNNs and ML using Unimodal Oncology Datasets	Pathology Radiology	Iqbal et al. (2022), Ahmedt-Aristizabal et al. (2022), Anand et al. (2020), Wang et al. (2020b), Wang et al. (2022), Ahmedt-Aristizabal et al. (2022) Mo et al. (2020), Du et al. (2019), Chao et al. (2020), Tian et al. (2020), Saueressig et al. (2021), Singh et al.
	Molecular	(2021) Fout et al. (2017), Feng et al. (2020), Li et al. (2021b), Choi et al. (2017), Song et al. (2021b), Wang et al. (2021a)
MML Data Fusion Stages		Cui et al. (2021), Huang and Chung (2020), Lian et al. (2022), Wen et al. (2022), Pati et al. (2020), Wang et al. (2021a), Rao et al. (2021), Zeng et al. (2020), Chatzianastasis et al. (2023), Tortora et al. (2023), Tortora et al. (2023), Tortora et al. (2021), Chen et al. (2020), Shi et al. (2019), Zhang et al. (2021), Rozemberczki et al. (2022), Wang et al. (2021b), Leng et al. (2022), Rajadhyaksha and Chitkara (2023), Park et al. (2022)

RoBERTa (Zhuang et al., 2021), CLIP (Radford et al., 2021), T5 (Raffel et al., 2020), BART (Lewis et al., 2019), BLOOM (Scao et al., 2022), ALIGN (Jia et al., 2021), CoCa (Yu et al., 2022) and more. Multimodal Transformers are a recent development in the field of MML, which extends the capabilities of traditional Transformers to handle multiple data modalities. The intermodality dependencies are captured by the cross-attention mechanism in multimodal Transformers, allowing the model to jointly reason and extract rich data representations. There are various types of multimodal Transformers, such as Unified Transformer (UniT) (Hu and Singh, 2021), Multiway Multimodal Transformer (MMT) (Tang et al., 2022), CLIP (Radford et al., 2021), Flamingo (Alayrac et al., 2022), CoCa (Yu et al., 2022), Perceiver IO (Jaegle et al., 2021), and GPT-4(Achiam et al., 2023).

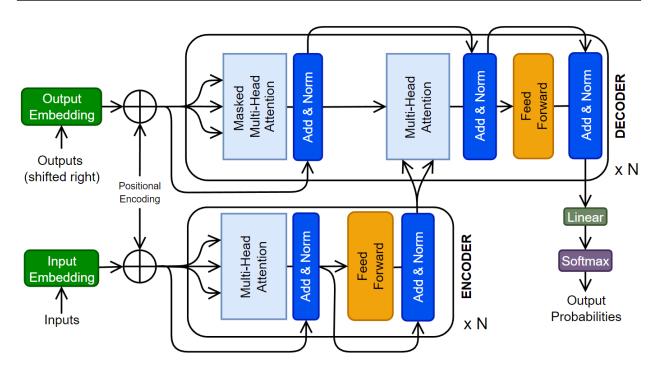


Figure 8. The original Transformer architecture is presented (Vaswani et al., 2017). A Transformer can have multiple encoder and decoder blocks, as well as some additional layers.

4.1 Model Architecture

The original Transformer (Figure 8) was composed of multiple encoder and decoder blocks, each made up of several layers of self-attention and feed-forward neural networks. The encoder takes the input sequence and generates hidden representations, which are then fed to the decoder. The decoder generates the output sequence by attending to the encoder's hidden representations and the previous tokens (i.e., auto-regressive). The self-attention operation (or scaled dot-product) is a crucial component of the Transformer. It determines the significance of each element in the input sequence with respect to the whole input. Self-attention operates by computing a weighted sum of the input sequence's hidden representations, where the weights are determined by the dot product between the *query* vector and the *key* vector, followed by a scaling operation to stabilize the gradients. The resulting weighted sum is multiplied by a *value* vector to obtain the output of the self-attention operation. There has been a tremendous amount of work on various facets of Transformer architecture. The readers are referred to relevant review papers (Otter et al., 2021; Xu et al., 2023; Han et al., 2023; Galassi et al., 2021).

4.2 Multimodal Transformers

Self-attention allows a Transformer model to process each input as a fully connected graph and attend to (or equivalently learn from) the global patterns present in the input. This makes Transformers compatible with various data modalities by treating each token (or its embedding) as a node in the graph. To use Transformers for a data modality, we need to tokenize the input and select an embedding space for the tokens. Tokenization and embedding selections are flexible and can be done at multiple granularity levels, such as using raw features, ML-extracted features, patches from

Data Modalities	Tokenization Level	Token Embeddings Model
Pathology images	Patch	CNNs (Chen et al., 2021)
Radiology images	Patch	CNNs (Xie et al., 2021)
EHR data	ICD code	GNNs (Shang et al., 2019), ML models (Rasmy et al., 2021)
-Omics	Graphs K-mers	GNNs (Kaczmarek et al., 2021) ML model (Ji et al., 2020)
Clinical notes	Word	BERT (Devlin et al., 2019), RoBERTa (Zhuang et al., 2021), BioBERT (Lee et al., 2019)

Table 3. Oncology data modalities and their respective tokenization and embeddings selection techniques

the input image, or graph nodes. Table 3 summarizes some common practices used for various types 579 of data in cancer data sets. Handling inter-modality interactions is the main challenge in developing 580 multimodal Transformer models. Usually, it is done through one of these fusion methods: early 581 fusion of data modalities, cross-attention, hierarchical attention, and late fusion, as illustrated in 582 Figure 9. In the following, we present and compare data processing steps for these four methods 583 using two data modalities as an example. The same analysis can be extended to multiple modalities. 584

4.2.1 585 Early Fusion

586

587

588

589

590

591

592

594

595

596

597

598

601

Early fusion is the simplest way to combine data from multiple modalities. The data from different modalities are concatenated to a single input before being fed to the Transformer model, which processes the input as a single entity. Mathematically, the concatenation operation is represented as $x_{cat} = [x_1, x_2]$, where x_1 and x_2 are the inputs from two data modalities, and x_{cat} is the concatenated input to the model. Early fusion is simple and efficient. However, it assumes that all modalities are equally important and relevant for the task at hand (Kalfaoglu et al., 2020), which may not always be practically true (Zhong et al., 2023).

Cross-Attention Fusion 4.2.2 593

Cross-attention is a relatively more flexible approach to modeling the interactions between data modalities and learning their joint representations. The self-attention layers attend to different modalities at different stages of data processing. Cross-attention allows the model to selectively attend to different modalities based on their relevance to the task (Li et al., 2021a) and capture complex interactions between the modalities (Rombach et al., 2022).

4.2.3 Hierarchical Fusion 599

600 Hierarchical fusion is a complex approach to combining multiple modalities. For instance, the Depth-supervised Fusion Transformer for Salient Object Detection (DFTR) employs hierarchical feature extraction to improve salient object detection performance by fusing low-level spatial 602 features and high-level semantic features from different scales (Zhu et al., 2022). Yang et al. 603 (2020) introduced a hierarchical approach to fine-grained classification using a fusion Transformer.

Furthermore, the Hierarchical Multimodal Transformer (HMT) for video summarization can capture global dependencies and multi-hop relationships among video frames (Zhao et al., 2022).

607 4.2.4 Late Fusion

608

609

610

611

612 613 In late fusion, each data modality is processed independently by its own Transformer model, the branch outputs are concatenated and passed through the final classifier. Late fusion allows the model to capture the unique features of each modality while still learning their joint representation. Sun et al. (2021) proposed a multi-modal adaptive late fusion Transformer network for estimating the levels of depression. Their model extracts long-term temporal information from audio and visual data independently and then fuses weights at the end to learn a joint representation of data.

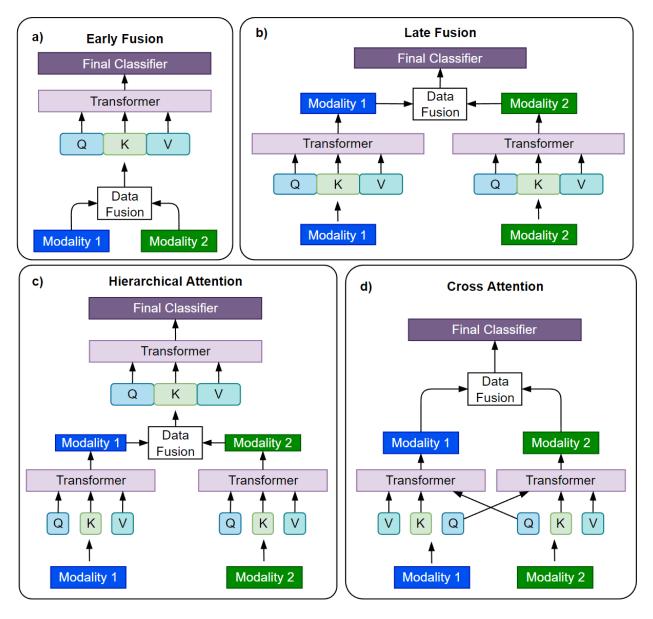


Figure 9. Four different strategies of fusing information from various data modalities in multimodal Transformers are presented.

14 4.3 Transformers for Processing Oncology Datasets

615 Transformers have been successfully applied to various tasks in oncology, including cancer screening, diagnosis, prognosis, treatment selection, and prediction of clinical variables (Boehm 616 617 et al., 2021; Shao et al., 2021; Liang et al., 2022a; Lian et al., 2022; Chen et al., 2021). For instance, a Transformer-based model was used to predict the presence and grade of breast cancer using a 618 combination of imaging and genomics data (Boehm et al., 2021). TransMIL (Shao et al., 2021), a 619 620 Transformer model, was proposed to process histopathology images using self-attention to learn 621 and classify histopathology slides by overcoming the challenges faced by multi-instance learning (MIL). Recently, a Transformer and convolution parallel network, TransConv (Liang et al., 2022a), 622 623 was proposed for automatic brain tumor segmentation using MRI data. Transformers and GNNs 624 have also been combined in MML for early-stage NSCLC prognostic prediction using the patient's 625 clinical and pathological features and by modeling the patient's physiological network (Lian et al... 626 2022). Similarly, a multimodal co-attention Transformer was proposed for survival prediction using 627 WSIs and genomic sequences (Chen et al., 2021). The authors used a co-attention mechanism to learn the interactions between the two data modalities. 628

Reinforcement learning with human feedback (RLHF) has emerged as a promising technique 629 to infuse large language models with domain knowledge and human preferences for healthcare 630 applications. Sun et al. (2023) proposed an approach to continuously improve a conversational agent 631 for behavioral interventions by integrating few-shot generation, prompt engineering, and RLHF to 632 leverage human feedback from therapists and clients. Giuffrè et al. (2024) discussed strategies to 633 optimize large language models for digestive disease by using RLHF to infuse domain knowledge 634 through supervised fine-tuning. Basit et al. (2024) introduced MedAide, an on-premise healthcare 635 chatbot that employs RLHF during training to enhance its medical diagnostic capabilities on edge 636 637 devices. Dai et al. (2023) presented Safe RLHF, a novel algorithm that decouples human preferences for helpfulness and harmlessness during RLHF to improve the safety and value alignment of large 638 language models in sensitive healthcare domains. 639

5 MML - CHALLENGES AND OPPORTUNITIES

Learning from multimodal oncology data is a complex and rapidly growing field that presents both challenges and opportunities. While MML has shown significant promise, there are many challenges owing to the inductive biases of the ML models (Ektefaie et al., 2023a). In this context, we present major challenges of MML in oncology settings that, if addressed, could unlock the full potential of this emerging field.

5.1 Large Amounts of High-quality Data

645

646 647

648

649

650

651

DL models are traditionally trained on large datasets with enough samples for training, validation, and testing, such as JFT-300M (Sun et al., 2017) and YFCC100M (Thomee et al., 2016), which are not available in the cancer domain. For example, the largest genomics data repository, the Gene Expression Omnibus (GEO) database, has approximately 1.1 million samples with the keyword 'cancer' compared to 3 billion images in JFT-300M (Jiang et al., 2022b). Annotating medical and oncology data is a time-consuming and manual process that requires significant expertise in many

7D 1 1 4	D C	D' 1	•	C
I ahia /i	References	Luccucced	1n	Section /
Table 7.	References	Discusseu	111	SCCHOII 4.

Sections	References Discussed
Multimodal Transformers	Vaswani et al. (2017), Otter et al. (2021), Xu et al. (2023), Dosovitskiy et al. (2020), Zhang (2020), Ahmed et al. (2022b), Ma et al. (2021a), Boehm et al. (2021), Devlin et al. (2019), Radford et al. (2018), Zhuang et al. (2021), Radford et al. (2021), Raffel et al. (2020), Lewis et al. (2019), Scao et al. (2022), Jia et al. (2021), Yu et al. (2022), Hu and Singh (2021), Tang et al. (2022), Radford et al. (2021), Alayrac et al. (2022), Yu et al. (2022), Jaegle et al. (2021), Achiam et al. (2023), Otter et al. (2021), Xu et al. (2023), Han et al. (2023), Galassi et al. (2021)
MML Data Fusion Stages	Kalfaoglu et al. (2020), Zhong et al. (2023), Li et al. (2021a), Rombach et al. (2022), Zhu et al. (2022), Yang et al. (2020), Zhao et al. (2022), Sun et al. (2021)
Transformers for Oncology Datasets	Boehm et al. (2021), Shao et al. (2021), Liang et al. (2022a), Lian et al. (2022), Chen et al. (2021), Sun et al. (2023), Giuffrè et al. (2024), Basit et al. (2024), Dai et al. (2023)

different areas of medical sciences. Factors like heterogeneity of the disease, noise in data recording, background, and training of medical professionals leading to inter- and intra-operator variability cause lack of reproducibility and inconsistent clinical outcomes (Lipkova et al., 2022).

5.2 Data Registration and Alignment

Data alignment and registration refer to the process of combining and aligning data from different modalities in a useful manner (Zhao et al., 2023). In multimodal oncology data, this process involves aligning data from multiple modalities such as CT, MRI, PET, and WSIs, as well as genomics, transcriptomics, and clinical records. Data registration involves aligning the data modalities to a common reference frame and may involve identifying common landmarks or fiducial markers. If the data is not registered or aligned correctly, it may be difficult to fuse the information from different modalities (Liang et al., 2022b).

5.3 Pan-Cancer Generalization and Transference

Transference in MML aims to transfer knowledge between modalities and their representations to improve the performance of a model trained on a primary modality (Liang et al., 2022b). Because of the unique characteristics of each cancer type and site, it is challenging to develop models that can generalize across different cancer sites. Furthermore, models trained on a specific modality, such as radiology images, will not perform well with other imaging modalities, such as histopathology slides. Fine-tuning the model on a secondary modality, multimodal co-learning, and model induction are techniques to achieve transference and generalization (Wei et al., 2020). To overcome this challenge, mechanisms for improved universality of ML models need to be devised.

672 5.4 Missing Data Samples and Modalities

- The unavailability of one or more modalities or the absence of samples in a modality affects the
- 674 model learning, as most of the existing DL models cannot process the "missing information". This
- 675 requirement, in turn, constrains the already insufficient size of datasets in oncology. Almost all
- 676 publicly available oncology datasets have missing data for a large number of samples (Jiang et al.,
- 677 2022b). Various approaches for handling missing data samples and modalities are gradually gaining
- 678 traction. However, this is still an open challenge (Mirza et al., 2019).

679 5.5 Imbalanced Data

- Class imbalance refers to the phenomenon when one class (e.g., cancer negative/positive) is
- 681 represented significantly more in the data than another class. Class imbalance is common in oncology
- 682 data (Mirza et al., 2019). DL models struggle to classify underrepresented classes accurately.
- 683 Techniques such as data augmentation, ensemble, continual learning, and transfer learning are used
- 684 to counter the class imbalance challenge (Mirza et al., 2019).

685 5.6 Explainability and Trustworthiness

- The explainability in DL, e.g., how GNNs and Transformers make a specific decision, is still an
- area of active research (Li et al., 2022b; Nielsen et al., 2022). GNNExplainer (Ying et al., 2019),
- 688 PGM-Explainer (Vu and Thai, 2020), and SubgraphX (Yuan et al., 2021) are some attempts to
- 689 explain the decision-making process of GNNs. The explainability methods for Transformers have
- 690 been analyzed in (Remmer, 2022). Existing efforts and a roadmap to improve the trustworthiness of
- 691 GNNs have been presented in the latest survey (Zhang et al., 2022a). However, the explainability
- and trustworthiness of multimodal GNNs and Transformers is an open challenge.

693 5.7 Over-smoothing in GNNs

- One particular challenge in using GNNs is over-smoothing, which occurs when the GNN is trained
- 695 for too long, causing the node representations to become almost similar (Wu et al., 2020). This
- 696 leads to a loss of information, a decrease in the model's performance, and a lack of generalization
- 697 (Valsesia et al., 2021). Regularization techniques such as dropout, weight decay, skip-connection,
- 698 and incorporating higher-order structures, such as motifs and graphlets, have been proposed.
- 699 However, building deep architectures that can scale and adapt to varying structural patterns of
- 700 graphs is still an open challenge.

701 5.8 Modality Collapse

- Modality collapse is a phenomenon that occurs in MML, where a model trained on multiple
- 703 modalities may become over-reliant on a single modality, to the point where it ignores or neglects
- 704 the other modalities (Javaloy et al., 2022). Recent work explored the reasons and theoretical
- 705 understanding of modality collapse (Huang et al., 2022). However, the counter-actions needed to
- 706 balance model dependence on data modalities require active investigation by the ML community.

5.9 **Dynamic and Temporal Data** 707

Dynamic and temporal data refers to the data that changes over time (Wu et al., 2020). Tumor 708 surveillance is a well-known technique to study longitudinal cancer growth over multiple data 709 modalities (Waqas et al., 2021). Spatio-temporal methods such as multiple instance learning, GNNs, 710 and hybrid of multiple models can capture complex change in the data relationships over time; 711 however, learning from multimodal dynamic data is very challenging and an active area of research 712 (Fritz et al., 2022). 713

5.10 **Data Privacy** 714

715

717

719

721

722

723

724

725

726

727 728

729

730

731

732

733

734

735

736

737

738

739

740

741

Given the sensitive nature of medical data, privacy and security are critical considerations in the development and deployment of MML models for oncology applications. With the increased 716 adoption of MML in healthcare settings, it is essential to adapt these techniques to enable local data processing and protect patient privacy while fostering collaborative research and analysis across 718 different sites and institutions. Federated learning (FL) has emerged as a promising approach to 720 train large multimodal models across various sites without the need for direct data sharing (Pati et al., 2022). In an FL setup, each participating site trains a local model on its own data and shares only the model updates with a central server, which aggregates the updates and sends the updated global model back to the sites. This allows for collaborative model development while keeping the raw data securely within each site's premises.

To further enhance privacy protection in FL and other distributed learning scenarios, differential privacy (DP) can be integrated into the model training process. DP is a rigorous mathematical framework that involves adding carefully calibrated noise to data or model updates before sharing, in order to protect individual privacy while preserving the utility of the data for analysis (Nampalle et al., 2023; Islam et al., 2022; Akter et al., 2022). Secure multi-party computation (SMPC) is another powerful technique for enabling joint analysis and model training on private datasets held by different healthcare providers or research institutions, without revealing the raw data to each other (Şahinbaş and Catak, 2021; Alghamdi et al., 2023; Yogi and Mundru, 2024). SMPC protocols leverage advanced cryptographic techniques to allow multiple parties to compute a function over their combined data inputs securely, such that each party learns only the output of the computation and nothing about the other parties' inputs. In addition to these solutions, implementing appropriate access control and authentication mechanisms is crucial for restricting access to sensitive healthcare data to only authorized individuals and entities (Orii et al., 2024). This involves defining and enforcing strict policies and procedures for granting, managing, and revoking access privileges based on the principle of least privilege and the need-to-know basis. Regular security risk assessments should also be conducted to identify and mitigate potential vulnerabilities proactively, ensuring the ongoing protection of patient data.

Other Challenges 5.11 742

MML requires extensive computational resources to train models on a variety of datasets and tasks. 743 Robustness and failure detection (Ahmed et al., 2022a) are critical aspects of MML, particularly 744 in applications such as oncology. Uncertainty quantification techniques, such as Bayesian neural 745 networks (Dera et al., 2021), are still under-explored avenues in the MML. By addressing these

751

752

754

755

756 757

758

759

761

762

763

764

challenges, it is possible to develop MML models that are able to surpass the performance offered by single-modality models. 748

Potential Future Directions 5.12 749

The future of MML in oncology holds immense potential. A critical direction is the integration of large amounts of high-quality data from diverse modalities, such as imaging, genomic, and clinical data, to enhance the accuracy and comprehensiveness of cancer diagnostics and treatment predictions 753 in an end-to-end manner. Overcoming challenges in data registration and alignment is crucial to ensure seamless integration and accurate interpretation of multimodal data. Developing robust models capable of pan-cancer generalization and transference can enable more universal applications across different cancer types. Addressing issues of missing data samples and modalities, and tackling imbalanced datasets, will be essential to improve model robustness and fairness. Enhancing explainability and trustworthiness in these models is vital for clinical adoption, necessitating transparent and interpretable AI systems. Preventing modality collapse is important for maintaining the distinct contributions of each data modality. Moreover, leveraging dynamic and temporal data 760 can offer deeper insights into cancer progression and treatment responses. Ensuring data privacy and ethical considerations will be paramount as the field advances, balancing innovation with the protection of patient information. Lastly, implementing MML applications in clinical settings is crucial to fully realize the benefits of MML in cancer research.

765 5.13 Limitations of the Study

MML is a broad research field that has recently gained traction. In this review, we have focused on 766 the application of MML on oncology data. However, MML is widely being adopted in applications 767 such as autonomous vehicles, education, earth science, climate change, and space exploration (Xiao 768 et al., 2020; Li et al., 2024; Hadid et al., 2024; Sanders et al., 2023). Moreover, beyond GNNs and 769 Transformers, MML has been explored using encoder-decoder methods, constraint-based methods, 770 canonical correlations, Restricted Boltzmann Machines (RBMs), and many more (Zhao et al., 2024; 771 Qi et al., 2020). Each of these topics require an extensive review of the literature in the form of 772 separate articles. 773

6 MULTIMODAL ONCOLOGY DATA SOURCES

Unifying the various collections of oncology data into central archives necessitates a focused effort. We have assembled a list of datasets from data portals maintained by the National 775 Institute of Health and other organizations, although this list is not exhaustive. The goal of this 776 777 compilation is to offer machine learning researchers in oncology a consolidated data resource. The collection, which is updated regularly, can be accessed at https://lab-rasool.github. 778 io/pan-cancer-dataset-sources/ (Tripathi et al., 2024a). The compilation of pan-779 780 cancer datasets from sources such as The Cancer Imaging Archive (TCIA), Genomic Data Commons (GDC), and Proteomic Data Commons (PDC) serves as a valuable resource for cancer research. 781 By providing a unified view of multimodal data that includes imaging, genomics, proteomics, and 782 clinical records, this compilation facilitates the development of adaptable and scalable datasets 783 specifically designed for machine learning applications in oncology (Tripathi et al., 2024a). The 784

Table 5. References Discussed in Section 5.

Sections	References Discussed
Large Amounts of High-quality Data	Ektefaie et al. (2023a), Sun et al. (2017), Thomee et al. (2016), Jiang et al. (2022b), Lipkova et al. (2022)
Data Registration and Alignment	Zhao et al. (2023), Liang et al. (2022b)
Pan-Cancer Generalization and Transference	Liang et al. (2022b), Wei et al. (2020)
Missing Data Samples and Modalities	Jiang et al. (2022b), Mirza et al. (2019)
Imbalanced Data	Mirza et al. (2019)
Explainability and Trustworthiness	Li et al. (2022b), Nielsen et al. (2022), Ying et al. (2019), Vu and Thai (2020), Yuan et al. (2021), Remmer (2022), Zhang et al. (2022a)
Over-smoothing in GNNs	Wu et al. (2020), Valsesia et al. (2021)
Modality Collapse	Javaloy et al. (2022), Huang et al. (2022)
Dynamic and Temporal Data	Wu et al. (2020), Waqas et al. (2021), Fritz et al. (2022)
Data Privacy	Pati et al. (2022), Nampalle et al. (2023), Islam et al. (2022), Akter et al. (2022), Şahinbaş and Catak (2021), Alghamdi et al. (2023), Yogi and Mundru (2024), Orii et al. (2024)
Other Challenges	Ahmed et al. (2022a), Dera et al. (2021)
Limitations of the Study	Xiao et al. (2020), Li et al. (2024), Hadid et al. (2024), Sanders et al. (2023), Zhao et al. (2024), Qi et al. (2020)

compiled datasets encompass a broad spectrum of data modalities, such as radiology images (CT, MRI, PET), pathology slides, genomic data (DNA, RNA), proteomics, and clinical records. This multimodal nature enables the integration of different data types to capture the intricacies of cancer. Moreover, the compilation covers 32 cancer types, ranging from prevalent cancers like breast, lung, and colorectal to less common forms such as mesothelioma and uveal melanoma. The inclusion of hundreds to thousands of cases for each cancer type provides a substantial resource for training machine learning models, especially deep learning algorithms.

Standardizing the diverse data formats, annotations, and metadata across different sources is essential for creating datasets that are suitable for machine learning. The HoneyBee framework, a modular system designed to streamline the creation of machine learning-ready multimodal oncology datasets from diverse sources, can help address this challenge (Tripathi et al., 2024b). HoneyBee supports data ingestion from various sources, handles different data formats and modalities, and ensures consistent data representation. It also facilitates the integration of multimodal data, enabling the creation of datasets that combine imaging, genomics, proteomics, and clinical data for a holistic view of each patient case. Furthermore, HoneyBee incorporates pre-trained foundational embedding models for different data modalities, such as image encoders, genomic sequence embedders, and

clinical text encoders. These embeddings can serve as input features for downstream machine learning models, leveraging transfer learning and reducing the need for extensive labeled data. The framework's scalable and modular architecture allows for efficient processing of large-scale datasets and easy integration of new data sources, preprocessing techniques, and embedding models. By utilizing the HoneyBee framework, researchers can create high-quality, multimodal oncology datasets tailored to their specific research objectives, promoting collaboration and advancing

7 CONCLUSION

Existing research into the integration of data across various modalities has already yielded promising 808 outcomes, highlighting the potential for significant advancements in cancer research. However, the 809 lack of a comprehensive framework capable of encompassing the full spectrum of cancer dataset 810 811 modalities presents a notable challenge. The synergy between diverse methodologies and data across different scales could unlock deeper insights into cancer, potentially leading to more accurate 812 prognostic and predictive models than what is possible through single data modalities alone. In our 813 814 survey, we have explored the landscape of multimodal learning applied to oncology datasets and the specific tasks they can address. Looking ahead, the key to advancing this field lies in the development 815 of robust, deployment-ready MML frameworks. These frameworks must not only scale efficiently 816 across all modalities of cancer data but also incorporate capabilities for uncertainty quantification, 817 interpretability, and generalizability. Such advancements will be critical for effectively integrating 818 oncology data across multiple scales, modalities, and resolutions. The journey towards achieving 819 820 these goals is complex, yet essential for the next leaps in cancer research. By focusing on these 821 areas, future research has the potential to significantly enhance our understanding of cancer, leading 822 to improved outcomes for patients through more informed and personalized treatment strategies.

CONFLICT OF INTEREST STATEMENT

machine learning applications in cancer research.

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

A.W, A.A., R.R, P.S., and G.R. conceived the manuscript. A.W., A.A. reviewed the literature and wrote the initial draft. All authors reviewed, improved, and contributed to the manuscript.

FUNDING

- 827 This work has been supported in part by the National Science Foundation awards 1903466, 2008690,
- 828 2234836, and 2234468, and in part by the Biostatistics and Bioinformatics Shared Resource at
- 829 the H. Lee Moffitt Cancer Center & Research Institute, an NCI-designated Comprehensive Cancer
- 830 Center (P30-CA076292).

DATA AVAILABILITY STATEMENT

- 831 The datasets analyzed for this study can be found in the pan-cancer-dataset-sources repository
- 832 https://lab-rasool.github.io/pan-cancer-dataset-sources/.

REFERENCES

- 833 Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., et al. (2023). Gpt-4
- technical report. arXiv preprint arXiv:2303.08774
- 835 Ahmed, S., Dera, D., Hassan, S. U., Bouaynaya, N., and Rasool, G. (2022a). Failure detection in
- deep neural networks for medical imaging. Frontiers in Medical Technology 4
- 837 Ahmed, S., Nielsen, I. E., Tripathi, A., Siddiqui, S., Rasool, G., and Ramachandran, R. P. (2022b).
- 838 Transformers in time-series analysis: A tutorial. arXiv preprint arXiv:2205.01138
- 839 Ahmedt-Aristizabal, D., Armin, M. A., Denman, S., Fookes, C., and Petersson, L. (2022). A survey
- on graph-based deep learning for computational histopathology. *Computerized Medical Imaging*
- and Graphics 95, 102027. doi:https://doi.org/10.1016/j.compmedimag.2021.102027
- 842 Akter, M., Moustafa, N., and Lynar, T. (2022). Edge intelligence-based privacy protection framework
- for iot-based smart healthcare systems. In IEEE INFOCOM 2022-IEEE Conference on Computer
- 844 Communications Workshops (INFOCOM WKSHPS) (IEEE), 1–8
- 845 Al-jabery, K. K., Obafemi-Ajayi, T., Olbricht, G. R., and Wunsch II, D. C. (2020). Data
- 846 preprocessing. In Computational Learning Approaches to Data Analytics in Biomedical
- 847 Applications, eds. K. K. Al-jabery, T. Obafemi-Ajayi, G. R. Olbricht, and D. C. Wunsch II
- 848 (Academic Press). 7–27. doi:https://doi.org/10.1016/B978-0-12-814482-4.00002-4
- 849 Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., et al. (2022). Flamingo: a
- visual language model for few-shot learning. Advances in neural information processing systems
- 851 35, 23716–23736
- 852 Alghamdi, W., Salama, R., Sirija, M., Abbas, A. R., and Dilnoza, K. (2023). Secure multi-party
- computation for collaborative data analysis. In E3S Web of Conferences (EDP Sciences), vol.
- 854 399, 04034
- 855 Almasan, P., Suárez-Varela, J., Rusek, K., Barlet-Ros, P., and Cabellos-Aparicio, A. (2022). Deep
- reinforcement learning meets graph neural networks: Exploring a routing optimization use case.
- 857 Computer Communications 196, 184–194
- 858 Anand, D., Gadiya, S., and Sethi, A. (2020). Histographs: graphs in histopathology. In Medical
- 859 *Imaging 2020: Digital Pathology* (SPIE), vol. 11320, 150–155
- 860 Angermueller, C., Lee, H. J., Reik, W., and Stegle, O. (2017). DeepCpG: accurate prediction of
- single-cell DNA methylation states using deep learning. *Genome biology* 18, 1–13
- 862 Asan, O., Nattinger, A. B., Gurses, A. P., Tyszka, J. T., and Yen, T. W. (2018). Oncologists' Views
- Regarding the Role of Electronic Health Records in Care Coordination. *JCO Clinical Cancer*
- 864 Informatics, 1–12doi:10.1200/CCI.17.00118. PMID: 30652555
- 865 Bai, S., Zhang, F., and Torr, P. H. (2021). Hypergraph convolution and hypergraph attention. *Pattern*
- 866 *Recognition* 110, 107637
- 867 Baltrušaitis, T., Ahuja, C., and Morency, L.-P. (2018). Multimodal machine learning: A survey and
- taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 423–443

- Barhoumi, Y., Bouaynaya, N. C., and Rasool, G. (2023). Efficient scopeformer: Towards scalable and rich feature extraction for intracranial hemorrhage detection. *IEEE Access*
- Basit, A., Hussain, K., Hanif, M. A., and Shafique, M. (2024). Medaide: Leveraging large language models for on-premise medical assistance on edge devices. *arXiv preprint arXiv:2403.00830*
- Boehm, K. M., Khosravi, P., Vanguri, R., Gao, J., and Shah, S. P. (2021). Harnessing multimodal data integration to advance precision oncology. *Nature Reviews Cancer*, 1–13
- 875 [Dataset] Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., et al. (2022).
 876 On the opportunities and risks of foundation models
- Borisov, V., Leemann, T., Seßler, K., Haug, J., Pawelczyk, M., and Kasneci, G. (2022). Deep Neural
- Networks and Tabular Data: A Survey. *IEEE Transactions on Neural Networks and Learning*
- 879 Systems, 1–21doi:10.1109/TNNLS.2022.3229161
- Çalışkan, M. and Tazaki, K. (2023). Ai/ml advances in non-small cell lung cancer biomarker
 discovery. *Frontiers in Oncology* 13
- Cao, Z.-J. and Gao, G. (2022). Multi-omics single-cell data integration and regulatory inference with graph-linked embedding. *Nature Biotechnology* 40, 1458–1466
- Chan, H.-P., Hadjiiski, L. M., and Samala, R. K. (2020). Computer-aided diagnosis in the era of deep learning. *Medical Physics* 47, e218–e227. doi:https://doi.org/10.1002/mp.13764
- 886 Chao, C.-H., Zhu, Z., Guo, D., Yan, K., Ho, T.-Y., Cai, J., et al. (2020). Lymph node gross
- tumor volume detection in oncology imaging via relationship learning using graph neural
- network. In Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd
- International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part VII 23 (Springer),
 772–782
- Chatzianastasis, M., Vazirgiannis, M., and Zhang, Z. (2023). Explainable multilayer graph neural network for cancer gene prediction. *arXiv* preprint arXiv:2301.08831
- 893 Chen, B., Jin, J., Liu, H., Yang, Z., Zhu, H., Wang, Y., et al. (2023). Trends and hotspots in research
- on medical images with deep learning: a bibliometric analysis from 2013 to 2023. *Frontiers in Artificial Intelligence* 6, 1289669
- 896 Chen, R. J., Lu, M. Y., Wang, J., Williamson, D. F., Rodig, S. J., Lindeman, N. I., et al. (2020).
- Pathomic Fusion: an integrated framework for fusing histopathology and genomic features for cancer diagnosis and prognosis. *IEEE Transactions on Medical Imaging* 41, 757–770
- 899 Chen, R. J., Lu, M. Y., Weng, W.-H., Chen, T. Y., Williamson, D. F., Manz, T., et al. (2021).
- Multimodal co-attention transformer for survival prediction in gigapixel whole slide images. In
- 901 2021 IEEE/CVF International Conference on Computer Vision (ICCV). 3995–4005. doi:10.1109/
- 902 ICCV48922.2021.00398
- 903 Choi, E., Bahadori, M. T., Song, L., Stewart, W. F., and Sun, J. (2017). GRAM: graph-based attention model for healthcare representation learning. In *Proceedings of the 23rd ACM SIGKDD*
- international conference on knowledge discovery and data mining. 787–795
- 906 Cui, H., Xuan, P., Jin, Q., Ding, M., Li, B., Zou, B., et al. (2021). Co-graph attention reasoning based
- 907 imaging and clinical features integration for lymph node metastasis prediction. In *Medical Image*
- 908 Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference,
- 909 Strasbourg, France, September 27–October 1, 2021, Proceedings, Part V 24 (Springer), 657–666
- 910 Dai, J., Pan, X., Sun, R., Ji, J., Xu, X., Liu, M., et al. (2023). Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*

- 912 Dara, S. and Tumma, P. (2018). Feature extraction by using deep learning: A survey. In 2018
- 913 Second International Conference on Electronics, Communication and Aerospace Technology
- 914 (*ICECA*). doi:10.1109/ICECA.2018.8474912
- 915 Defferrard, M., Bresson, X., and Vandergheynst, P. (2016). Convolutional neural networks on
- graphs with fast localized spectral filtering. Advances in neural information processing systems
- 917 29
- 918 Dera, D., Bouaynaya, N. C., Rasool, G., Shterenberg, R., and Fathallah-Shaykh, H. M. (2021).
- 919 PremiUm-CNN: Propagating Uncertainty Towards Robust Convolutional Neural Networks. *IEEE*
- 920 Transactions on Signal Processing 69, 4669–4684
- 921 Dera, D., Rasool, G., and Bouaynaya, N. (2019). Extended variational inference for propagating
- 922 uncertainty in convolutional neural networks. In 2019 IEEE 29th International Workshop on
- 923 *Machine Learning for Signal Processing (MLSP)* (IEEE), 1–6
- 924 Derrow-Pinion, A., She, J., Wong, D., Lange, O., Hester, T., Perez, L., et al. (2021). ETA prediction
- 925 with graph neural networks in google maps. In Proceedings of the 30th ACM International
- 926 Conference on Information & Knowledge Management. 3767–3776
- 927 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). Bert: Pre-training of deep bidirectional
- 928 transformers for language understanding. In North American Chapter of the Association for
- 929 Computational Linguistics
- 930 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2020).
- An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale doi:10.48550/
- 932 ARXIV.2010.11929
- 933 Du, H., Feng, J., and Feng, M. (2019). Zoom in to where it matters: a hierarchical graph based
- model for mammogram analysis. arXiv preprint arXiv:1912.07517
- 935 Ektefaie, Y., Dasoulas, G., Noori, A., Farhat, M., and Zitnik, M. (2023a). Geometric multimodal
- 936 representation learning. arXiv preprint arXiv:2209.03299
- 937 Ektefaie, Y., Dasoulas, G., Noori, A., Farhat, M., and Zitnik, M. (2023b). Multimodal learning with
- 938 graphs. *Nature Machine Intelligence*, 1–11
- 939 Farooq, H., Chen, Y., Georgiou, T. T., Tannenbaum, A., and Lenglet, C. (2019). Network curvature
- as a hallmark of brain structural connectivity. *Nature communications* 10, 1–11
- 941 Fathinezhad, F., Adibi, P., Shoushtarian, B., and Chanussot, J. (2023). Graph neural networks and
- 942 reinforcement learning: A survey
- 943 Feng, Y.-H., Zhang, S.-W., and Shi, J.-Y. (2020). DPDDI: a deep predictor for drug-drug interactions.
- 944 *BMC bioinformatics* 21, 1–15
- 945 Fout, A., Byrd, J., Shariat, B., and Ben-Hur, A. (2017). Protein interface prediction using graph
- onvolutional networks. Advances in neural information processing systems 30
- 947 Fritz, C., Dorigatti, E., and Rügamer, D. (2022). Combining graph neural networks and spatio-
- temporal disease models to improve the prediction of weekly covid-19 cases in germany. *Scientific*
- 949 Reports 12, 3930
- 950 Galassi, A., Lippi, M., and Torroni, P. (2021). Attention in Natural Language Processing. *IEEE*
- 951 Transactions on Neural Networks and Learning Systems 32, 4291–4308. doi:10.1109/TNNLS.
- 952 2020.3019893

- 953 Ghaffari Laleh, N., Ligero, M., Perez-Lopez, R., and Kather, J. N. (2023). Facts and hopes on the
- use of artificial intelligence for predictive immunotherapy biomarkers in cancer. *Clinical Cancer*
- 955 Research 29, 316–323
- 956 Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message
- passing for quantum chemistry. In International conference on machine learning (PMLR),
- 958 1263–1272
- 959 Giuffrè, M., Kresevic, S., Pugliese, N., You, K., and Shung, D. L. (2024). Optimizing large
- language models in digestive disease: strategies and challenges to improve clinical outcomes.
- 961 Liver International
- 962 Gonzalez Zelaya, C. V. (2019). Towards explaining the effects of data preprocessing on machine
- learning. In IEEE 35th International Conference on Data Engineering (ICDE). 2086–2090.
- 964 doi:10.1109/ICDE.2019.00245
- 965 Grossman, R. L., Heath, A. P., Ferretti, V., Varmus, H. E., Lowy, D. R., Kibbe, W. A., et al.
- 966 (2016). Toward a shared vision for cancer genomic data. New England Journal of Medicine 375,
- 967 1109–1112. doi:10.1056/nejmp1607591
- 968 Hadid, A., Chakraborty, T., and Busby, D. (2024). When geoscience meets generative ai and large
- language models: Foundations, trends, and future challenges. Expert Systems, e13654
- 970 Hamilton, W., Ying, Z., and Leskovec, J. (2017). Inductive representation learning on large graphs.
- 971 Advances in neural information processing systems 30
- 972 Hamilton, W. L. (2020). Graph representation learning. Synthesis Lectures on Artifical Intelligence
- 973 and Machine Learning 14
- 974 Han, K., Wang, Y., Chen, H., Chen, X., Guo, J., Liu, Z., et al. (2023). A Survey on Vision
- 975 Transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 87–110.
- 976 doi:10.1109/TPAMI.2022.3152247
- 977 Hartsock, I. and Rasool, G. (2024). Vision-language models for medical report generation and
- 978 visual question answering: A review. arXiv preprint arXiv:2403.02469
- 979 Hook, D. W., Porter, S. J., and Herzog, C. (2018). Dimensions: Building Context for Search and
- Evaluation. Frontiers in Research Metrics and Analytics 3, 23. doi:10.3389/frma.2018.00023
- 981 Hu, R. and Singh, A. (2021). Unit: Multimodal multitask learning with a unified transformer. In
- 982 Proceedings of the IEEE/CVF International Conference on Computer Vision. 1439–1449
- 983 Huang, Y. and Chung, A. C. (2020). Edge-variational graph convolutional networks for uncertainty-
- aware disease prediction. In Medical Image Computing and Computer Assisted Intervention—
- 985 MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part
- 986 *VII 23* (Springer), 562–572
- 987 Huang, Y., Du, C., Xue, Z., Chen, X., Zhao, H., and Huang, L. (2021). What makes multi-modal
- learning better than single (provably). In Advances in Neural Information Processing Systems,
- eds. A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan
- 990 Huang, Y., Lin, J., Zhou, C., Yang, H., and Huang, L. (2022). Modality competition: What
- makes joint training of multi-modal network fail in deep learning?(provably). In *International*
- 992 *Conference on Machine Learning* (PMLR), 9226–9259
- 993 Ibrahim, A., Mohamed, H. K., Maher, A., and Zhang, B. (2022). A survey on human cancer
- oategorization based on deep learning. Frontiers in Artificial Intelligence 5. doi:10.3389/frai.
- 995 2022.884749

- 996 Iqbal, M. S., Ahmad, W., Alizadehsani, R., Hussain, S., and Rehman, R. (2022). Breast cancer dataset, classification and detection using deep learning. In *Healthcare* (MDPI), vol. 10, 2395
- 998 Iqbal, M. S., Luo, B., Mehmood, R., Alrige, M. A., and Alharbey, R. (2019). Mitochondrial
- organelle movement classification (fission and fusion) via convolutional neural network approach.
- 1000 *IEEE Access* 7, 86570–86577
- 1001 Islam, T. U., Ghasemi, R., and Mohammed, N. (2022). Privacy-preserving federated learning model
- 1002 for healthcare data. In 2022 IEEE 12th Annual Computing and Communication Workshop and
- 1003 *Conference (CCWC)* (IEEE), 0281–0287
- 1004 Jaegle, A., Borgeaud, S., Alayrac, J.-B., Doersch, C., Ionescu, C., Ding, D., et al. (2021). Perceiver
- 1005 IO: A General Architecture for Structured Inputs &; Outputs doi:10.48550/ARXIV.2107.14795
- 1006 Jansen, C., Ramirez, R. N., El-Ali, N. C., Gomez-Cabrero, D., Tegner, J., Merkenschlager, M., et al.
- 1007 (2019). Building gene regulatory networks from scATAC-seq and scRNA-seq using linked self
- organizing maps. *PLoS computational biology* 15, e1006555
- 1009 Javaloy, A., Meghdadi, M., and Valera, I. (2022). Mitigating Modality Collapse in Multimodal
- 1010 VAEs via Impartial Optimization. In International Conference on Machine Learning (PMLR),
- 1011 9938–9964
- 1012 Ji, Y., Zhou, Z., Liu, H., and Davuluri, R. V. (2020). DNABERT: pre-trained Bidirectional
- 1013 Encoder Representations from Transformers model for DNA-language in genome. bioRxiv
- 1014 doi:10.1101/2020.09.17.301879
- 1015 Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., et al. (2021). Scaling up visual and
- 1016 vision-language representation learning with noisy text supervision. In International Conference
- 1017 *on Machine Learning* (PMLR), 4904–4916
- 1018 Jiang, J., Dun, C., Huang, T., and Lu, Z. (2018). Graph convolutional reinforcement learning. arXiv
- 1019 *preprint arXiv:1810.09202*
- 1020 Jiang, P., Sinha, S., Aldape, K., Hannenhalli, S., Sahinalp, C., and Ruppin, E. (2022a). Big Data
- in basic and translational cancer research. *Nature Reviews Cancer* 22, 625–639. doi:10.1038/
- 1022 s41568-022-00502-0
- 1023 Jiang, P., Sinha, S., Aldape, K., Hannenhalli, S., Sahinalp, C., and Ruppin, E. (2022b). Big data in
- basic and translational cancer research. *Nature Reviews Cancer* 22, 625–639
- 1025 Jiao, L., Chen, J., Liu, F., Yang, S., You, C., Liu, X., et al. (2022). Graph representation learning
- meets computer vision: A survey. *IEEE Transactions on Artificial Intelligence*
- 1027 Jin, D., Huo, C., Dang, J., Zhu, P., Zhang, W., Pedrycz, W., et al. (2022). Heterogeneous
- graph neural networks using self-supervised reciprocally contrastive learning. arXiv preprint
- 1029 *arXiv:2205.00256*
- 1030 Joo, S., Ko, E., Kwon, S., Jeon, E., Jung, H., Kim, J., et al. (2021). Multimodal deep learning
- models for the prediction of pathologic response to neoadjuvant chemotherapy in breast cancer.
- 1032 *Sci Rep* 11, 18800. doi:10.1038/s41598-021-98408-8
- 1033 Kaczmarek, E., Jamzad, A., Imtiaz, T., Nanayakkara, J., Renwick, N., and Mousavi, P. (2021). Multi-
- omic graph transformers for cancer classification and interpretation. In PACIFIC SYMPOSIUM
- 1035 ON BIOCOMPUTING 2022 (World Scientific), 373–384
- 1036 Kalfaoglu, M. E., Kalkan, S., and Alatan, A. A. (2020). Late Temporal Modeling in 3D CNN
- 1037 Architectures with BERT for Action Recognition. In Computer Vision–ECCV 2020 Workshops:
- 1038 Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16 (Springer), 731–747

- 1039 Khan, M., Ashraf, I., Alhaisoni, M., Damaševičius, R., Scherer, R., Rehman, A., et al. (2020).
- Multimodal brain tumor classification using deep learning and robust feature selection: A
- machine learning application for radiologists. Diagnostics (Basel) 10, 565. doi:10.3390/
- 1042 diagnostics10080565
- 1043 Khan, S., Ali, H., and Shah, Z. (2023). Identifying the role of vision transformer for skin cancer—a
- scoping review. Frontiers in Artificial Intelligence 6. doi:10.3389/frai.2023.1202990
- 1045 Kipf, T. N. and Welling, M. (2016). Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*
- 1047 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. nature 521, 436–444
- 1048 Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., et al. (2019). BioBERT: a pre-trained
- biomedical language representation model for biomedical text mining. *Bioinformatics* 36,
- 1050 1234–1240. doi:10.1093/bioinformatics/btz682
- 1051 Leng, D., Zheng, L., Wen, Y., Zhang, Y., Wu, L., Wang, J., et al. (2022). A benchmark study of
- deep learning-based multi-omics data fusion methods for cancer. Genome Biology 23, 1–32
- 1053 Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., et al. (2019). BART:
- Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and
- 1055 Comprehension. In Annual Meeting of the Association for Computational Linguistics
- 1056 Li, M. M., Huang, K., and Zitnik, M. (2022a). Graph representation learning in biomedicine and
- healthcare. *Nature Biomedical Engineering*, 1–17
- 1058 Li, P., Gu, J., Kuen, J., Morariu, V. I., Zhao, H., Jain, R., et al. (2021a). Selfdoc: Self-supervised
- document representation learning. In Proceedings of the IEEE/CVF Conference on Computer
- 1060 Vision and Pattern Recognition. 5652–5660
- 1061 Li, P., Wang, J., Qiao, Y., Chen, H., Yu, Y., Yao, X., et al. (2021b). An effective self-supervised
- framework for learning expressive molecular global representations to drug discovery. *Briefings*
- in Bioinformatics 22, bbab109
- 1064 Li, P., Yang, Y., Pagnucco, M., and Song, Y. (2022b). Explainability in graph neural networks: An
- experimental survey. arXiv preprint arXiv:2203.09258
- 1066 Li, Z., Pardos, Z. A., and Ren, C. (2024). Aligning open educational resources to new taxonomies:
- How ai technologies can help and in which scenarios. Computers & Education 216, 105027
- 1068 Lian, J., Deng, J., Hui, E. S., Koohi-Moghadam, M., She, Y., Chen, C., et al. (2022). Early
- stage NSCLS patients' prognostic prediction with multi-information using transformer and graph
- neural network model. *eLife* 11, e80547. doi:10.7554/eLife.80547
- 1071 Liang, J., Yang, C., Zeng, M., and Wang, X. (2022a). TransConver: transformer and convolution
- parallel network for developing automatic brain tumor segmentation in MRI images. *Quantitative*
- 1073 Imaging in Medicine and Surgery 12
- 1074 Liang, P. P., Zadeh, A., and Morency, L.-P. (2022b). Foundations and recent trends in multimodal
- machine learning: Principles, challenges, and open questions. arXiv preprint arXiv:2209.03430
- 1076 Lipkova, J., Chen, R. J., Chen, B., Lu, M. Y., Barbieri, M., Shao, D., et al. (2022). Artificial
- intelligence for multimodal data integration in oncology. Cancer Cell 40, 1095–1110. doi:https:
- 1078 //doi.org/10.1016/j.ccell.2022.09.012
- 1079 Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., et al. (2017). A
- survey on deep learning in medical image analysis. *Medical image analysis* 42, 60–88

- Liu, J., Pandya, P., and Afshar, S. (2021). Therapeutic advances in oncology. *International Journal* of *Molecular Sciences* 22. doi:10.3390/ijms22042008
- Liu, T., Huang, J., Liao, T., Pu, R., Liu, S., and Peng, Y. (2022). A hybrid deep learning model for predicting molecular subtypes of human breast cancer using multimodal data. *Irbm* 43, 62–74
- 1085 Ma, J., Liu, J., Lin, Q., Wu, B., Wang, Y., and You, Y. (2021a). Multitask learning for visual
- $1086 \qquad \text{question answering. } \textit{IEEE Transactions on Neural Networks and Learning Systems} \text{ , } 1-15 \\ \text{doi:} 10.$
- 1087 1109/TNNLS.2021.3105284
- 1088 Ma, L., Yang, Z., Miao, Y., Xue, J., Wu, M., Zhou, L., et al. (2019). NeuGraph: Parallel Deep Neural
- Network Computation on Large Graphs. In USENIX Annual Technical Conference. 443–458
- 1090 Ma, X. and Jia, F. (2020). Brain tumor classification with multimodal MR and pathology images. In
- 1091 Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 5th International
- 1092 Workshop, BrainLes 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October
- 1093 *17, 2019, Revised Selected Papers, Part II 5* (Springer), 343–352
- 1094 Ma, Y., Liu, X., Zhao, T., Liu, Y., Tang, J., and Shah, N. (2021b). A unified view on graph neural
- networks as graph signal denoising. In *Proceedings of the 30th ACM International Conference*
- on Information & Knowledge Management. 1202–1211
- 1097 Ma, Y. and Tang, J. (2021). *Deep learning on graphs* (Cambridge University Press)
- 1098 Miotto, R., Li, L., Kidd, B. A., and Dudley, J. T. (2016). Deep patient: an unsupervised
- representation to predict the future of patients from the electronic health records. *Scientific*
- 1100 reports 6, 1–10
- 1101 Mirza, B., Wang, W., Wang, J., Choi, H., Chung, N. C., and Ping, P. (2019). Machine learning and
- integrative analysis of biomedical big data. Genes 10, 87
- 1103 Mo, S., Cai, M., Lin, L., Tong, R., Chen, Q., Wang, F., et al. (2020). Multimodal priors
- guided segmentation of liver lesions in MRI using mutual information based graph co-attention
- networks. In Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd
- 1106 International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part IV 23 (Springer),
- 1107 429-438
- 1108 Muhammad, L. J. and Bria, A. (2023). Editorial: Ai applications for diagnosis of breast cancer.
- 1109 Frontiers in Artificial Intelligence 6. doi:10.3389/frai.2023.1247261
- 1110 Muhammad, W., Ahmed, S.-b.-S., Naeem, S., Khan, A. A. M. H., Qureshi, B. M., Hussain, A., et al.
- 1111 (2024). Artificial neural network-assisted prediction of radiobiological indices in head and neck
- cancer. Frontiers in Artificial Intelligence 7, 1329737
- 1113 Nampalle, K. B., Singh, P., Narayan, U. V., and Raman, B. (2023). Vision through the veil:
- Differential privacy in federated learning for medical image classification. arXiv preprint
- 1115 *arXiv:2306.17794*
- 1116 Nguyen, T., Le, H., Quinn, T. P., Nguyen, T., Le, T. D., and Venkatesh, S. (2021). GraphDTA:
- predicting drug-target binding affinity with graph neural networks. *Bioinformatics* 37, 1140–1147
- 1118 Nie, M., Chen, D., and Wang, D. (2023). Reinforcement learning on graphs: A survey. IEEE
- 1119 Transactions on Emerging Topics in Computational Intelligence
- 1120 Nielsen, I. E., Dera, D., Rasool, G., Ramachandran, R. P., and Bouaynaya, N. C. (2022). Robust
- explainability: A tutorial on gradient-based attribution methods for deep neural networks. *IEEE*
- 1122 Signal Processing Magazine 39, 73–84

- 1123 Orii, L., Feldacker, C., Tweya, H., and Anderson, R. (2024). ehealth data security and privacy:
- Perspectives from diverse stakeholders in malawi. Proceedings of the ACM on Human-Computer
- 1125 *Interaction* 8, 1–26
- 1126 Otter, D. W., Medina, J. R., and Kalita, J. K. (2021). A survey of the usages of deep learning for
- 1127 natural language processing. *IEEE Transactions on Neural Networks and Learning Systems* 32,
- 1128 604–624. doi:10.1109/TNNLS.2020.2979670
- 1129 Park, J., Cho, J., Chang, H. J., and Choi, J. Y. (2021). Unsupervised hyperbolic representation
- learning via message passing auto-encoders. In *Proceedings of the IEEE/CVF Conference on*
- 1131 Computer Vision and Pattern Recognition. 5516–5526
- 1132 Park, M.-K., Lim, J.-M., Jeong, J., Jang, Y., Lee, J.-W., Lee, J.-C., et al. (2022). Deep-Learning
- 1133 Algorithm and Concomitant Biomarker Identification for NSCLC Prediction Using Multi-Omics
- Data Integration. *Biomolecules* 12, 1839
- 1135 Pati, P., Jaume, G., Fernandes, L. A., Foncubierta-Rodríguez, A., Feroce, F., Anniciello, A. M., et al.
- 1136 (2020). HACT-Net: A Hierarchical Cell-to-Tissue Graph Neural Network for Histopathological
- 1137 Image Classification. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging*,
- and Graphs in Biomedical Image Analysis: Second International Workshop, UNSURE 2020, and
- 1139 Third International Workshop, GRAIL 2020, Held in Conjunction with MICCAI 2020, Lima, Peru,
- 1140 October 8, 2020, Proceedings 2 (Springer), 208–219
- 1141 Pati, S., Baid, U., Edwards, B., Sheller, M., Wang, S.-H., Reina, G. A., et al. (2022). Federated
- learning enables big data for rare cancer boundary detection. *Nature communications* 13, 7346
- 1143 Qi, G., Sun, Y., Li, M., and Hou, X. (2020). Development and application of matrix variate restricted
- boltzmann machine. *IEEE Access* 8, 137856–137866
- 1145 Quinn, M., Forman, J., Harrod, M., Winter, S., Fowler, K. E., Krein, S. L., et al. (2019). Electronic
- health records, communication, and data sharing: challenges and opportunities for improving the
- diagnostic process. *Diagnosis* 6, 241–248. doi:doi:10.1515/dx-2018-0036
- 1148 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., et al. (2021). Learning
- transferable visual models from natural language supervision. In *International conference on*
- 1150 *machine learning* (PMLR), 8748–8763
- 1151 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. (2018). Improving language
- understanding by generative pre-training
- 1153 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., et al. (2020). Exploring the
- Limits of Transfer Learning with a Unified Text-to-Text Transformer. The Journal of Machine
- 1155 *Learning Research* 21
- 1156 Rajadhyaksha, N. and Chitkara, A. (2023). Graph contrastive learning for multi-omics data. arXiv
- 1157 *preprint arXiv:2301.02242*
- 1158 Rao, J., Zhou, X., Lu, Y., Zhao, H., and Yang, Y. (2021). Imputing single-cell RNA-seq data by
- 1159 combining graph convolution and autoencoder neural networks. *Iscience* 24, 102393
- 1160 Rasmy, L., Xiang, Y., Xie, Z., Tao, C., and Zhi, D. (2021). Med-BERT: Pretrained contextualized
- embeddings on large-scale structured electronic health records for disease prediction. *npj Digital*
- 1162 *Medicine* 4. doi:10.1038/s41746-021-00455-y
- 1163 Remmer, E. (2022). Explainability Methods for Transformer-based Artificial Neural Networks: a
- 1164 Comparative Analysis. Ph.D. thesis

- 1165 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2022). High-resolution image
- synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*
- vision and pattern recognition. 10684–10695
- 1168 Rowe, S. P. and Pomper, M. G. (2022). Molecular imaging in oncology: Current impact and future
- directions. CA: A Cancer Journal for Clinicians 72, 333–352
- 1170 Rozemberczki, B., Gogleva, A., Nilsson, S., Edwards, G., Nikolov, A., and Papa, E. (2022).
- MOOMIN: Deep Molecular Omics Network for Anti-Cancer Drug Combination Therapy.
- 1172 In Proceedings of the 31st ACM International Conference on Information & Knowledge
- 1173 *Management*. 3472–3483
- 1174 Şahinbaş, K. and Catak, F. O. (2021). Secure multi-party computation based privacy preserving
- data analysis in healthcare iot systems. arxiv e-prints, article. arXiv preprint arXiv:2109.14334
- 1176 Sanders, L. M., Scott, R. T., Yang, J. H., Qutub, A. A., Garcia Martin, H., Berrios, D. C., et al.
- 1177 (2023). Biological research and self-driving labs in deep space supported by artificial intelligence.
- 1178 *Nature Machine Intelligence* 5, 208–219
- 1179 Sankar, A., Wu, Y., Gou, L., Zhang, W., and Yang, H. (2018). Dynamic graph representation
- learning via self-attention networks. arXiv preprint arXiv:1812.09430
- 1181 Saueressig, C., Berkley, A., Kang, E., Munbodh, R., and Singh, R. (2021). Exploring graph-based
- neural networks for automatic brain tumor segmentation. In From Data to Models and Back: 9th
- 1183 International Symposium, DataMod 2020, Virtual Event, October 20, 2020, Revised Selected
- 1184 *Papers* 9 (Springer), 18–37
- 1185 Scao, T. L., Fan, A., Akiki, C., Pavlick, E., Ilić, S., Hesslow, D., et al. (2022). BLOOM: A
- 176B-Parameter Open-Access Multilingual Language Model. *arXiv preprint arXiv:2211.05100*
- 1187 Schulz, S., Woerl, A.-C., Jungmann, F., Glasner, C., Stenzel, P., Strobl, S., et al. (2021). Multimodal
- deep learning for prognosis prediction in renal cancer. Frontiers in Oncology 11. doi:10.3389/
- 1189 fonc.2021.788740
- 1190 Shang, J., Ma, T., Xiao, C., and Sun, J. (2019). Pre-training of graph augmented transformers for
- medication recommendation. arXiv preprint arXiv:1906.00346
- 1192 Shao, Z., Bian, H., Chen, Y., Wang, Y., Zhang, J., Ji, X., et al. (2021). TransMIL: Transformer
- based Correlated Multiple Instance Learning for Whole Slide Image Classification. Advances in
- Neural Information Processing Systems 34
- 1195 Shi, J., Wang, R., Zheng, Y., Jiang, Z., and Yu, L. (2019). Graph convolutional networks for cervical
- 1196 cell classification. In MICCAI 2019 Computational Pathology Workshop COMPAY
- 1197 Siam, A., Alsaify, A. R., Mohammad, B., Biswas, M. R., Ali, H., and Shah, Z. (2023). Multimodal
- deep learning for liver cancer applications: a scoping review. Frontiers in Artificial Intelligence 6
- 1199 Siegel, R. L., Miller, K. D., Wagle, N. S., and Jemal, A. (2023). Cancer Statistics, 2023. CA: A
- 1200 Cancer Journal for Clinicians 73, 17–48. doi:https://doi.org/10.3322/caac.21763
- 1201 Singh, A., Hu, R., Goswami, V., Couairon, G., Galuba, W., Rohrbach, M., et al. (2022). FLAVA: A
- foundational language and vision alignment model. In CVPR
- 1203 Singh, G., Manjila, S., Sakla, N., True, A., Wardeh, A. H., Beig, N., et al. (2021). Radiomics and
- radiogenomics in gliomas: a contemporary update. *British Journal of Cancer* 125, 641–657
- 1205 Sleeman IV, W. C., Kapoor, R., and Ghosh, P. (2022). Multimodal classification: Current landscape,
- taxonomy and future directions. ACM Computing Surveys 55, 1–31

- 1207 Song, J., Zheng, Y., Zakir Ullah, M., Wang, J., Jiang, Y., Xu, C., et al. (2021a). Multiview
- multimodal network for breast cancer diagnosis in contrast-enhanced spectral mammography
- images. International Journal of Computer Assisted Radiology and Surgery 16, 979–988
- 1210 Song, Q., Su, J., and Zhang, W. (2021b). scGCN is a graph convolutional networks algorithm for
- knowledge transfer in single cell omics. *Nature communications* 12, 3826
- 1212 Stark, S. G., Ficek, J., Locatello, F., Bonilla, X., Chevrier, S., Singer, F., et al. (2020). SCIM:
- universal single-cell matching with unpaired feature sets. *Bioinformatics* 36, i919–i927
- 1214 Sun, C., Shrivastava, A., Singh, S., and Gupta, A. (2017). Revisiting unreasonable effectiveness
- of data in deep learning era. In *Proceedings of the IEEE international conference on computer*
- 1216 *vision*. 843–852
- 1217 Sun, H., Liu, J., Chai, S., Qiu, Z., Lin, L., Huang, X., et al. (2021). Multi-Modal Adaptive Fusion
- 1218 Transformer Network for the estimation of depression level. Sensors 21, 4764. doi:10.3390/
- 1219 s21144764
- 1220 Sun, X., Bosch, J. A., De Wit, J., and Krahmer, E. (2023). Human-in-the-loop interaction for
- 1221 continuously improving generative model in conversational agent for behavioral intervention.
- 1222 In Companion Proceedings of the 28th International Conference on Intelligent User Interfaces.
- 1223 99–101
- 1224 Syed, K., Sleeman IV, W. C., Hagan, M., Palta, J., Kapoor, R., and Ghosh, P. (2021). Multi-view
- data integration methods for radiotherapy structure name standardization. Cancers 13, 1796
- 1226 Talebi, R., Celis-Morales, C. A., Akbari, A., Talebi, A., Borumandnia, N., and Pourhoseingholi,
- M. A. (2024). Machine learning-based classifiers to predict metastasis in colorectal cancer
- patients. Frontiers in Artificial Intelligence 7
- 1229 Tang, J., Li, K., Hou, M., Jin, X., Kong, W., Ding, Y., et al. (2022). MMT: Multi-way Multi-modal
- 1230 Transformer for Multimodal Learning. In Proceedings of the Thirty-First International Joint
- 1231 Conference on Artificial Intelligence, IJCAI-22 (International Joint Conferences on Artificial
- 1232 Intelligence Organization), 3458–3465. doi:10.24963/ijcai.2022/480
- 1233 Thangudu, R. R., Rudnick, P. A., Holck, M., Singhal, D., MacCoss, M. J., Edwards, N. J., et al.
- 1234 (2020). Abstract lb-242: Proteomic data commons: A resource for proteogenomic analysis.
- 1235 Cancer Research 80, LB-242
- 1236 Thomee, B., Shamma, D. A., Friedland, G., Elizalde, B., Ni, K., Poland, D., et al. (2016).
- 1237 YFCC100M: The new data in multimedia research. *Communications of the ACM* 59, 64–73
- 1238 Tian, Z., Li, X., Zheng, Y., Chen, Z., Shi, Z., Liu, L., et al. (2020). Graph-convolutional-network-
- based interactive prostate segmentation in MR images. *Medical physics* 47, 4164–4176
- 1240 Tortora, M., Cordelli, E., Sicilia, R., Nibid, L., Ippolito, E., Perrone, G., et al. (2023).
- Radiopathomics: Multimodal learning in non-small cell lung cancer for adaptive radiotherapy.
- 1242 *IEEE Access*
- 1243 Tripathi, A., Waqas, A., Venkatesan, K., Yilmaz, Y., and Rasool, G. (2024a). Building flexible,
- scalable, and machine learning-ready multimodal oncology datasets. Sensors 24, 1634
- 1245 Tripathi, A., Wagas, A., Yilmaz, Y., and Rasool, G. (2024b). Honeybee: A scalable modular
- framework for creating multimodal oncology datasets with foundational embedding models.
- 1247 *arXiv preprint arXiv:2405.07460*

- 1248 Tripathi, A., Waqas, A., Yilmaz, Y., and Rasool, G. (2024c). Multimodal transformer model
- improves survival prediction in lung cancer compared to unimodal approaches. Cancer Research
- 1250 84, 4905–4905
- 1251 Tripathi, S., Moyer, E. J., Augustin, A. I., Zavalny, A., Dheer, S., Sukumaran, R., et al. (2022).
- RadGenNets: Deep learning-based radiogenomics model for gene mutation prediction in lung
- 1253 cancer. Informatics in Medicine Unlocked 33, 101062. doi:10.1016/j.imu.2022.101062
- 1254 Valsesia, D., Fracastoro, G., and Magli, E. (2021). RAN-GNNs: Breaking the Capacity Limits of
- 1255 Graph Neural Networks. IEEE Transactions on Neural Networks and Learning Systems
- 1256 Varlamova, E. V., Butakova, M. A., Semyonova, V. V., Soldatov, S. A., Poltavskiy, A. V., Kit, O. I.,
- et al. (2024). Machine learning meets cancer. Cancers 16, 1100
- 1258 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). Attention
- is all you need. Advances in neural information processing systems 30
- 1260 Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. (2017). Graph
- attention networks. arXiv preprint arXiv:1710.10903
- 1262 Vu, M. and Thai, M. T. (2020). PGM-Explainer: Probabilistic Graphical Model Explanations for
- 1263 Graph Neural Networks. *Advances in neural information processing systems* 33, 12225–12235
- 1264 Waikhom, L. and Patgiri, R. (2022). A survey of graph neural networks in various learning
- paradigms: methods, applications, and challenges. Artificial Intelligence Review, 1–70
- 1266 Wang, D., Su, J., and Yu, H. (2020a). Feature extraction and analysis of natural language processing
- for deep learning english language. *IEEE Access* 8, 46335–46345. doi:10.1109/ACCESS.2020.
- 1268 2974101
- 1269 Wang, J., Chen, R. J., Lu, M. Y., Baras, A., and Mahmood, F. (2020b). Weakly supervised
- prostate TMA classification via graph convolutional networks. In 2020 IEEE 17th International
- 1271 Symposium on Biomedical Imaging (ISBI) (IEEE), 239–243
- 1272 Wang, J., Ma, A., Chang, Y., Gong, J., Jiang, Y., Qi, R., et al. (2021a). scGNN is a novel graph
- neural network framework for single-cell RNA-Seq analyses. *Nature communications* 12, 1882
- 1274 Wang, S., Sun, S., Li, Z., Zhang, R., and Xu, J. (2017). Accurate de novo prediction of protein
- 1275 contact map by ultra-deep learning model. *PLoS computational biology* 13, e1005324
- 1276 Wang, T., Shao, W., Huang, Z., Tang, H., Zhang, J., Ding, Z., et al. (2021b). MOGONET
- integrates multi-omics data using graph convolutional networks allowing patient classification
- and biomarker identification. *Nature communications* 12, 3445
- Wang, Y., Wang, Y. G., Hu, C., Li, M., Fan, Y., Otter, N., et al. (2022). Cell graph neural networks
- enable the precise prediction of patient survival in gastric cancer. NPJ precision oncology 6, 45
- 1281 Waqas, A., Bui, M. M., Glassy, E. F., El Naga, I., Borkowski, P., Borkowski, A. A., et al.
- 1282 (2023). Revolutionizing digital pathology with the power of generative artificial intelligence and
- foundation models. *Laboratory Investigation*, 100255
- 1284 Waqas, A., Dera, D., Rasool, G., Bouaynaya, N. C., and Fathallah-Shaykh, H. M. (2021). Brain
- tumor segmentation and surveillance with deep artificial neural networks. Deep Learning for
- 1286 Biomedical Data Analysis, 311–350
- 1287 Waqas, A., Farooq, H., Bouaynaya, N. C., and Rasool, G. (2022). Exploring robust architectures for
- deep artificial neural networks. Communications Engineering 1, 46

- 1289 Waqas, A., Tripathi, A., Ahmed, S., Mukund, A., Stewart, P., Naeini, M., et al. (2024a). SeNMo:
- A self-normalizing deep learning model for enhanced multi-omics data analysis in oncology.
- 1291 *Cancer Research* 84, 908–908
- 1292 Waqas, A., Tripathi, A., Stewart, P., Naeini, M., and Rasool, G. (2024b). Embedding-based
- multimodal learning on pan-squamous cell carcinomas for improved survival outcomes. arXiv
- 1294 *preprint arXiv:2406.08521*
- 1295 Wei, H., Feng, L., Chen, X., and An, B. (2020). Combating noisy labels by agreement: A joint
- training method with co-regularization. In *Proceedings of the IEEE/CVF conference on computer*
- vision and pattern recognition. 13726–13735
- 1298 Wei, Y., Wang, X., Nie, L., He, X., Hong, R., and Chua, T.-S. (2019). MMGCN: Multi-modal graph
- 1299 convolution network for personalized recommendation of micro-video. In *Proceedings of the*
- 1300 *27th ACM international conference on multimedia.* 1437–1445
- 1301 Wen, H., Ding, J., Jin, W., Wang, Y., Xie, Y., and Tang, J. (2022). Graph neural networks for
- multimodal single-cell data integration. In *Proceedings of the 28th ACM SIGKDD Conference on*
- 1303 Knowledge Discovery and Data Mining. 4153–4163
- 1304 Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., and Philip, S. Y. (2020). A Comprehensive Survey
- on Graph Neural Networks. IEEE Transactions on Neural Networks and Learning Systems 32,
- 1306 4–24
- 1307 Xiao, Y., Codevilla, F., Gurram, A., Urfalioglu, O., and López, A. M. (2020). Multimodal end-to-end
- autonomous driving. IEEE Transactions on Intelligent Transportation Systems 23, 537–547
- 1309 Xie, Y., Zhang, J., Shen, C., and Xia, Y. (2021). Cotr: Efficiently bridging cnn and transformer
- for 3d medical image segmentation. In Medical Image Computing and Computer Assisted
- 1311 Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September
- 1312 27–October 1, 2021, Proceedings, Part III 24 (Springer), 171–180
- 1313 Xu, P., Zhu, X., and Clifton, D. A. (2023). Multimodal learning with transformers: A survey.
- 1314 IEEE Transactions on Pattern Analysis &; Machine Intelligence 45, 12113–12132. doi:10.1109/
- 1315 TPAMI.2023.3275156
- 1316 Xu, Y., Das, P., and McCord, R. P. (2022). SMILE: mutual information learning for integration of
- single-cell omics data. *Bioinformatics* 38, 476–486
- 1318 Yang, K. D., Belyaeva, A., Venkatachalapathy, S., Damodaran, K., Katcoff, A., Radhakrishnan, A.,
- et al. (2021a). Multi-domain translation between single-cell imaging and sequencing data using
- autoencoders. *Nature communications* 12, 31
- 1321 Yang, L., Ng, T. L. J., Smyth, B., and Dong, R. (2020). HTML: Hierarchical Transformer-
- Based Multi-Task Learning for Volatility Prediction. In *Proceedings of The Web Conference*
- 1323 2020 (New York, NY, USA: Association for Computing Machinery), WWW '20, 441–451.
- doi:10.1145/3366423.3380128
- 1325 Yang, T., Hu, L., Shi, C., Ji, H., Li, X., and Nie, L. (2021b). HGAT: Heterogeneous graph attention
- networks for semi-supervised short text classification. ACM Transactions on Information Systems
- 1327 *(TOIS)* 39, 1–29
- 1328 Yap, J., Yolland, W., and Tschandl, P. (2018). Multimodal skin lesion classification using deep
- learning. Experimental dermatology 27, 1261–1267
- 1330 Yi, H.-C., You, Z.-H., Huang, D.-S., and Kwoh, C. K. (2022). Graph representation learning in
- bioinformatics: trends, methods and applications. *Briefings in Bioinformatics* 23, bbab340

- 1332 Ying, Z., Bourgeois, D., You, J., Zitnik, M., and Leskovec, J. (2019). GNNExplainer: Generating
- Explanations for Graph Neural Networks. Advances in neural information processing systems 32
- 1334 Yogi, M. K. and Mundru, Y. (2024). Genomic data analysis with variant of secure multi-party
- 1335 computation technique. *Journal of Trends in Computer Science and Smart Technology* 5, 450–470
- 1336 Yousefi, S., Amrollahi, F., Amgad, M., Dong, C., Lewis, J. E., Song, C., et al. (2017). Predicting
- clinical outcomes from large scale cancer genomic profiles with deep survival models. *Scientific* reports 7
- 1339 Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., and Wu, Y. (2022). Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*
- 1341 Yuan, H., Yu, H., Wang, J., Li, K., and Ji, S. (2021). On explainability of graph neural networks via
- subgraph explorations. In *International Conference on Machine Learning* (PMLR), 12241–12252
- 1343 Zeng, Y., Zhou, X., Rao, J., Lu, Y., and Yang, Y. (2020). Accurately clustering single-cell RNA-seq
- data by capturing structural relations between cells through graph convolutional network. In 2020
- 1345 IEEE International Conference on Bioinformatics and Biomedicine (BIBM) (IEEE)
- Thang, H., Wu, B., Yuan, X., Pan, S., Tong, H., and Pei, J. (2022a). Trustworthy graph neural networks: Aspects, methods and trends. *arXiv* preprint arXiv:2205.07424
- 1348 Zhang, N. (2020). Learning adversarial transformer for symbolic music generation. IEEE
- 1349 Transactions on Neural Networks and Learning Systems, 1–10doi:10.1109/TNNLS.2020. 1350 2990746
- 1351 Zhang, Y.-D., Satapathy, S. C., Guttery, D. S., Górriz, J. M., and Wang, S.-H. (2021). Improved
- breast cancer classification through combining graph convolutional network and convolutional
- neural network. Information Processing & Management 58, 102439
- 1354 Zhang, Z., Yang, C., and Zhang, X. (2022b). scDART: integrating unmatched scRNA-seq and
- scATAC-seq data and learning cross-modality relationship simultaneously. *Genome Biology* 23,
- 1356 139
- 1357 Zhao, B., Gong, M., and Li, X. (2022). Hierarchical multimodal transformer to summarize videos.
- 1358 Neurocomputing 468, 360–369. doi:https://doi.org/10.1016/j.neucom.2021.10.039
- 1359 Zhao, F., Zhang, C., and Geng, B. (2024). Deep multimodal data fusion. ACM Computing Surveys
- 1360 Zhao, M., Huang, X., Jiang, J., Mou, L., Yan, D.-M., and Ma, L. (2023). Accurate Registration of
- 1361 Cross-Modality Geometry via Consistent Clustering. *IEEE Transactions on Visualization and*
- 1362 Computer Graphics
- 1363 Zhong, Z., Schneider, D., Voit, M., Stiefelhagen, R., and Beyerer, J. (2023). Anticipative Feature
- Fusion Transformer for Multi-Modal Action Anticipation. In 2023 IEEE/CVF Winter Conference
- on Applications of Computer Vision (WACV). 6057–6066. doi:10.1109/WACV56688.2023.00601
- 1366 Zhu, H., Sun, X., Li, Y., Ma, K., Zhou, S. K., and Zheng, Y. (2022). DFTR: Depth-supervised
- Fusion Transformer for Salient Object Detection doi:10.48550/ARXIV.2203.06429
- 1368 Zhuang, L., Wayne, L., Ya, S., and Jun, Z. (2021). A robustly optimized BERT pre-training approach
- with post-training. In *Proceedings of the 20th Chinese National Conference on Computational*
- 1370 Linguistics, eds. S. Li, M. Sun, Y. Liu, H. Wu, K. Liu, W. Che, S. He, and G. Rao (Huhhot, China:
- 1371 Chinese Information Processing Society of China), 1218–1227
- 1372 Zitnik, M., Agrawal, M., and Leskovec, J. (2018). Modeling polypharmacy side effects with graph
- 1373 convolutional networks. *Bioinformatics* 34, i457–i466