

Deep representation learning for temporal inference in cancer omics: a systematic literature review

Guillermo Prol-Castelo ^{1,2}, Davide Cirillo ^{1,*}, Alfonso Valencia ^{1,3}

¹Barcelona Supercomputing Center (BSC), C/ Jordi Girona 29, 08034 Barcelona, Spain

²Department of Medicine and Life Sciences, Universitat Pompeu Fabra, C/Dr Aiguader 88, Barcelona 08003, Spain

³Institució Catalana de Recerca i Estudis Avançats (ICREA), Pg. Lluís Companys 23 08010, Barcelona, Spain

*Corresponding author. Barcelona Supercomputing Center (BSC), C/ Jordi Girona 29, 08034 Barcelona, Spain. E-mail: davide.cirillo@bsc.es

Abstract

Deep learning methods, including deep representation learning (DRL) approaches such as variational autoencoders (VAEs), have been widely applied to cancer omics data to address the high dimensionality of these datasets. Despite remarkable advances, cancer is a complex and dynamic disease, making it challenging to study, and the temporal resolution of cancer progression captured by omics-based studies remains limited. In this systematic literature review, we explore the use of DRL, particularly the VAE, in cancer omics studies for modeling time-related processes, such as tumor progression and evolutionary dynamics. Our work reveals that these methods most commonly support subtyping, diagnosis, and prognosis in this context, but rarely emphasize temporal information. We observed that the scarcity of longitudinal omics data currently limits deeper temporal analyses that could enhance these applications. We propose that applying the VAE as a generative model to study cancer in time, particularly focusing on cancer staging, could lead to meaningful advancements in our understanding of the disease.

Keywords longitudinal omics data, cancer progression, variational autoencoder, generative artificial intelligence

Introduction

Deep learning (DL) is capable of capturing complex patterns inherent in health data, including molecular, imaging, and clinical data [1]. Given the high-dimensionality of omics data, deep representation learning (DRL) models have emerged as a way to model complex distributions in lower-dimensional representations [2]. Moreover, specific DRL methodologies add a generative capability to the model [3]. Such is the case of the variational autoencoder (VAE), a DRL method that uses Bayesian principles through variational inference to approximate complex relationships between observations and latent variables [4]. The VAE, introduced over a decade ago, consists of an encoder and a decoder; the former infers a low-dimensional representation called the latent space, while the latter learns to reconstruct the latent space into the original data. This architecture makes it particularly useful to both model a real data distribution and generate synthetic data [5].

Given the VAE is able to learn nonlinear relationships in heterogeneous, high-dimensional data, its use is appropriate for the study of complex diseases. Such is the case of cancer, the second most common cause of death worldwide, with an increasing incidence [6, 7]. Well-established techniques for data collection, including omics approaches, have produced large amounts of information, also referred to as big data, about cancer [8]. High-throughput

molecular and imaging omics data are commonly high-dimensional, noisy, sparse, and heterogeneous. These are also highly informative, including data at different levels of granularity: patient- [9] and/or cell-level [10] resolution; information from different biological sources: genes [11], proteins [12], or metabolites [13]; and even spatial information [14]. Leveraging DRL methods, such as the VAE, and cancer omics, could have a translational impact by discovering clinically relevant biomarkers [15] and patient subgroups [16].

Navigating such a vast amount of information becomes even more challenging due to the fact that cancer is a dynamic process [17], adding a temporal component. However, two main challenges arise in modeling time-dependent processes in cancer. First, cancer progression is not uniform across individuals [18, 19], meaning that samples collected at the same nominal time points may reflect different progression states across patients. Second, sequencing-based assays are typically destructive, preventing repeated measurements of the same biological sample. As a result, temporally annotated omics datasets generally contain unaligned samples, derived from either multiple individuals or separate samples from the same individual.

Here we present a systematic literature review (SLR) on the use of DRL, paying special attention to the VAE, in cancer studies with omics data for modeling time-related processes, such as tumor progression and evolutionary dynamics. To do so, we

Received: June 5, 2025. **Revised:** January 6, 2026. **Accepted:** March 5, 2026

© The Author(s) 2026. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

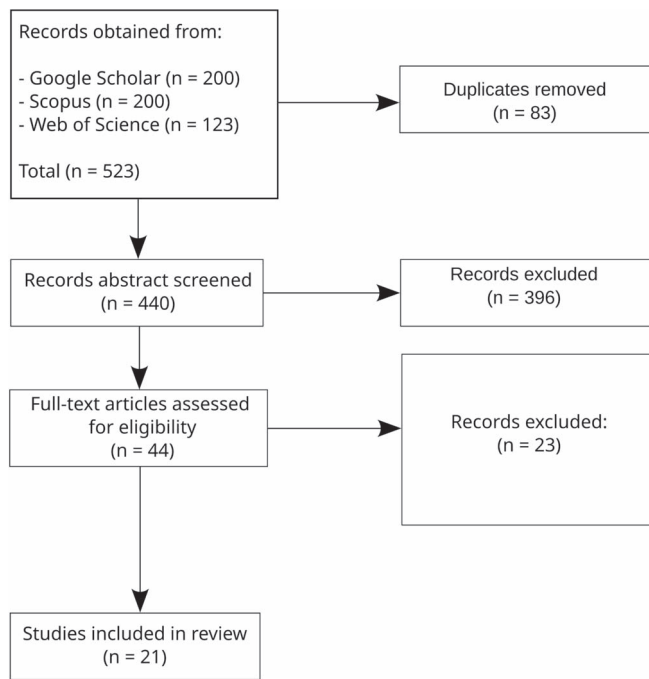


Figure 1 Flow diagram. Summary of the SLR process, following PRISMA guidelines.

consider the following research questions (answered in section Results):

1. How has representation learning, particularly the VAE, been used to study cancer with omics data?
2. How do cancer studies take longitudinal data into account?
3. How is cancer's temporal dimension studied with DL?
4. Which type of omics are most commonly used in time-related cancer studies?
5. What are the key limitations and challenges in the use of DRL for cancer studies with omics data in a temporal context?

The process of the SLR is summarized in Fig. 1 (see Methods for details), and our main insights are summarized in Fig. 2. We found 440 papers related to our search terms and reviewed their titles and abstracts. These reveal that the most common use of DRL and VAEs in cancer studies are related to subtyping, diagnosis, and prognosis. From the abstracts screening, we selected 44 papers for full-text review, of which 21 were included in our final analysis (summarized in Table 1). We learned that single-cell omics, with a recent surge of spatially resolved data, is the most common type of data used in time-related cancer studies, albeit in pseudotime, applying the methods included in our search. There is a lack of longitudinal omics records in the study of cancer in time, and even the temporally annotated datasets available contain unaligned samples. Stages have been considered as a way to provide a time unit for the study of cancer's advancement, but these usually consider different patients. Overall, we observed a lack of systematic approaches for leveraging temporal omics data to study cancer progression. To address this gap, in the Discussion we outline approaches to address these challenges, such as leveraging the VAE to synthesize temporally aligned instances across cancer stages.

Methods

We undertook this SLR in accordance with the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) recommendations [20].

Search process

We performed an SLR of DRL methods, with emphasis on VAEs, applied to cancer omics data for modeling time-related processes. Given the complex nature of the topic and the diversity of possible terms related to our research, we considered the following keywords:

- Variational autoencoder: We considered its abbreviation, VAE, as well. Since the name may not have been explicitly used in the title or abstract, especially during the years after its first publication in 2013 [4], we also included the term “representation learning” and “deep learning.”
- omics: We included the main omics data types, such as genomics, transcriptomics, proteomics, and metabolomics, as well as their combination by including multiomics.
- cancer: May also be referred to as tumor.
- time: To include possible references to studies that consider cancer progression through time, such as longitudinal studies, the queries included the keywords progression, evolution, trajectory, time series, and longitudinal.

For this review, we used Google Scholar, Scopus, and Web of Science to search the literature, including results from 2014 through 2024.

Given the keywords, we defined the following queries for each search engine:

- Google Scholar: (“deep learning” OR “representation learning” OR “Variational Autoencoder” OR VAE) AND (omics OR genomics OR transcriptomics OR proteomics OR metabolomics OR “multi-omics” OR multiomics) AND (“cancer” OR “tumor” OR “tumour”) AND (progression OR evolution OR trajectory OR “time series” OR longitudinal) since:2014 before:2025
- Web of Science: TS=(“deep learning” OR “representation learning” OR “Variational Autoencoder” OR VAE) AND TS=(omics OR genomics OR transcriptomics OR proteomics OR metabolomics OR “multi-omics” OR multiomics) AND TS=(“cancer” OR “tumor” OR “tumour”) AND TS=(progression OR evolution OR trajectory OR “time series” OR longitudinal) AND PY=(2014–2024)
- Scopus: TITLE-ABS-KEY(“deep learning” OR “representation learning” OR “Variational Autoencoder” OR VAE) AND (omics OR genomics OR transcriptomics OR proteomics OR metabolomics OR “multi-omics” OR multiomics) AND (“cancer” OR “tumor” OR “tumour”) AND (progression OR evolution OR trajectory OR “time series” OR longitudinal) AND PUBYEAR > 2013 AND PUBYEAR < 2025

Selection criteria

To keep our search within scope and guarantee good quality, we enforced the following inclusion criteria: (i) peer-review journal articles and conference proceedings; (ii) omics data studies; and (iii) papers that study cancer evolution or progression using DL and/or representational learning. For the same reasons, we also apply these

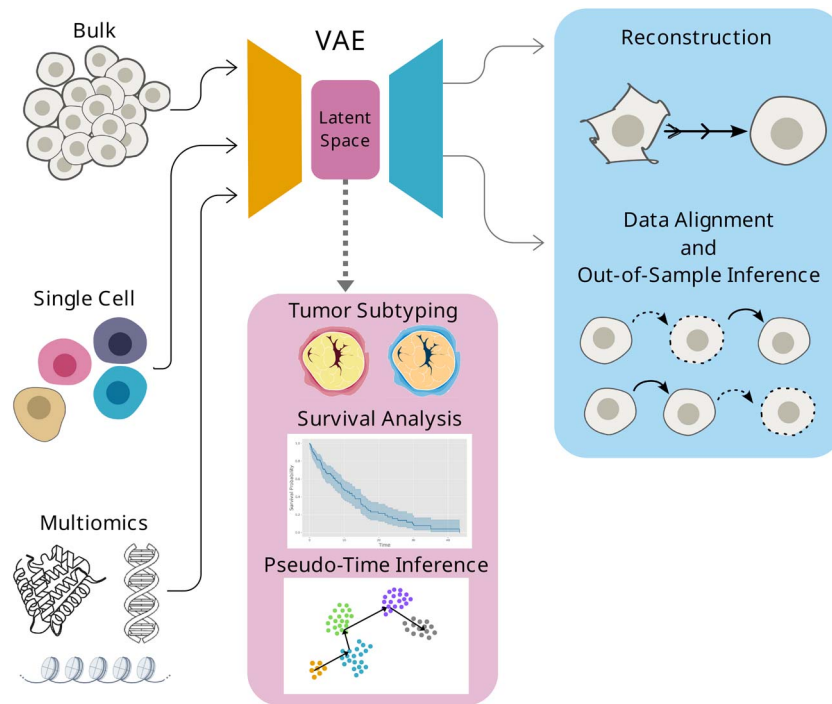


Figure 2 Summary figure of the most common methodologies involving representation learning and cancer studies. On the left side of the figure are shown the most common types of data used to train VAEs particularly, and DRL methods more generally, in cancer studies. After being passed to the encoder, data are embedded into a lower-dimensional representation, the latent space. The latent space has been commonly used for subtyping, survival analyses (commonly as part of prognostic analyses), and, when it comes to single-cell data, also for pseudo-time trajectories inference. Possible further applications, shown on the right-hand side of the figure, include leveraging the decoder's generative capabilities to reconstruct noisy data, and align and infer out-of-sample data, which may be used to study cancer progression in time. However, such applications of the decoder remain underexplored in biomedical research.

exclusion criteria: (i) opinions, letters to the editor, reviews, systematic reviews, and gray literature; (ii) non-English language works; (iii) papers focusing solely on survival and prognostic clinical analyses without providing biological and/or mechanistic insights; (iv) papers lacking temporal characterization. The rationale for inclusion and exclusion criteria was motivated by the scope and objectives of this review, which focuses on peer-reviewed research papers the omics-based DL and representation learning for studying temporal aspects related to cancer.

When in doubt about borderline cases, we followed a set of decision rules: (i) we included papers that considered the temporal dimension either explicitly, or implicitly through a proxy (such as cancer stages, pseudo-time inference, or spatially resolved progression patterns); (ii) papers focusing solely in prognosis prediction not using omics data were excluded; (iii) studies that included only cross-sectional data were excluded unless they incorporated a temporal proxy to study progression (such as cancer stages, pseudo-time inference, or spatially resolved progression patterns); (iv) non-cancer temporal studies were included if all other criteria were met, given the potential translational implications for cancer research. Moreover, Placido *et al.* [21] was included because it represents a singular showcase of a DRL architecture with generative capabilities (the Transformer) applied to available longitudinal data.

Data extraction

The search on Scopus and Web of Science was performed through the respective websites. Publish or Perish [22] was used to extract

the results from Google Scholar. The searches were limited to 200 results for each query. For each query, a table with information about the publication was downloaded from the respective platform. We then identified the duplicated titles and removed them. Next, we read the titles and abstracts, narrowing down our selection. On the remaining entries a full-text screening was performed. The process is summarized in Fig. 1. One reviewer (primary researcher) conducted the systematic review process, with results independently verified by the senior author.

Risk of bias

While carrying out this SLR, we identified several possible sources of bias, while taking steps to reduce their impact. The screening process was conducted by one primary reviewer (primary researcher), with a second reviewer (senior author) independently assessing the results to mitigate any subjective bias. While formal inter-rater reliability metrics were not calculated, consensus was reached on all included studies. The data extraction process was designed to mitigate a database coverage bias: the use of three different databases that are not limited to specific fields or literature types (such as preprint services). We limited the data extraction from Google Scholar to the first 200 results, whereas Scopus and Web of Science returned <200 results each from our queries (see Selection criteria). This is because of the practical feasibility of screening thousands of papers and the diminishing returns of relevance. Even though some relevant papers might have been excluded, the use of three databases was intended to keep the most relevant results. Temporal scope was deliberately restricted to publications from 2014 (corresponding to the introduction of the VAE,

Table 1 Summary table. This table contains the information of the final list of 21 papers obtained from the SLR. Papers have been grouped by the used Temporal Proxy, as reflected in the Results section: [23–30] provide longitudinal data, [31–36] rely on a temporal proxy such as stages, and [37–42] use single-cell data and pseudo-time

Reference	Summary	Model	Data type	Data source	Disease(s)	Temporal proxy
[23]	Classify colorectal cancer biopsies into no recurrence, recurrence, and stable disease.	Neural network	extracellular vesicle RNA (evRNA)	Self-obtained	Colorectal cancer	Recurrence status
[21]	Encode aligned trajectories in a latent space, and predict an out-of-sample, time-dependent risk of occurrence.	Gated recurrence units (GRUs) and Transformers	Electronic Health Records (EHRs)	Danish National Patient Registry (DNPR) and US Veterans Affairs (US-VA)	General Cancer	Actual time series
[24]	Learn a representation of time-longitudinal and multimodal data representation, used for classification.	GRUs	Whole-Slide Images and Cognitive assessment	ADNI	Alzheimer's Disease (AD)	Actual time series
[25]	Create embeddings to cluster temporal growth patterns in time-dependent data.	LSTM-VAE	Multi-omics: proteomics and metabolomics	Self-obtained and literature	Cardio vascular disease (CVD)	Actual time series
[26]	Combine multiomics to discover biomarkers of progress toward myocardial infarction.	Neural network	Microarray gene expression	Literature	Myocardial infarction	Actual time series
[27]	Segment histology images to monitor of kidney function deterioration.	Convolutional neural network	Whole-slide images	Self-obtained and, for validation, from literature	Kidney function deterioration	Pseudo-time
[28]	Prediction of early and late effects of radiation in exposed mice lungs.	DL (from the precision-medicine-toolbox)	Radiomics	Self-obtained	Radiation-induced lung injuries	Early and late status
[29]	First application of DL on tumor evolution. Infer most likely tumor phylogeny.	Reinforcement Learning (RL) and LSTM	Single-cell genotype matrix from single-cell sequencing	Literature	BRCA, Acute lymphoblastic leukemia, and colorectal cancer	Phylogeny
[30]	Learn graph embeddings of relationships between subclones.	Encoder–Decoder Neural network	Mutation trees from scDNA-Seq	Literature	Acute myeloid leukemia (AML), clear cell Renal Cell Carcinoma (ccRCC), mesothelioma, breast cancer, non-small cell lung cancer, and uveal melanoma	Phylogeny
[31]	Compare normal tissue against the different stages by transforming a source sample into the feature space of a target sample.	Adversarial network	bulk RNA-Seq gene expression	TCGA and GTEx	Thyroid cancer (THCA)	Stages
[32]	Classify early and late stages from RNA-Seq.	Neural network	bulk RNA-Seq gene expression	TCGA	Kidney renal papillary cell carcinoma (KIRP)	Stages
[33]	Inferred pseudo-time trajectories in the latent space of a Convolutional VAE between non-tumor, carcinoma <i>in situ</i> , and invasive stages.	Convolutional VAE	Chromatin images	Self-obtained	Ductal carcinoma in situ (DCIS)	Pseudo-time and invasive stage
[34]	Distinguish between normal tissue, non-invasive abnormally grown tissue, and invasive cancer.	Neural network	Gene expression microarrays	GEO and ArrayExpress	Cervical cancer (CxCa)	Invasive stage
[35]	Compares several binary and multiclass classifications on normalized data and on representations.	Stacked denoising autoencoder (SDAE) and VAE	bulk RNA-Seq Gene Expression	GTEx, TCGA, and SRA	BRCA, COAD, KIRC, LIHC, LUAD, LUSC, SKCM, STAD, THCA, UCEC	Stages
[36]	Data generation at different stages of BRCA.	Conditional VAE (CVAE)	bulk RNA-Seq gene expression and miRNA expression	TCGA	BRCA	Stages

(continued)

Table 1 Continued.

Reference	Summary	Model	Data type	Data source	Disease(s)	Temporal proxy
[37]	Distinguish cell subpopulations and infer cell differentiation trajectories.	VAE	Gene expression from single-cell RNA-Seq and copy number variation (CNV) from scDNA-Seq	Literature	Oligodendroglioma, Glioblastoma, Melanoma, Astrocytoma	Cell differentiation trajectory
[38]	Embed ST data for pseudo-time and trajectory inference, organizing the different tissue layers in different tissues.	Graph Autoencoder (AE)	Spatial transcriptomics (ST) and histology	Literature	Dorsolateral prefrontal cortex from human, mouse brain, BRCA, PDAC, Phalaenopsis (orchid tissue), and Mouse olfactory bulb	Pseudo-time
[39]	Embed ST of mouse visual cortex for trajectory and pseudo-time inference.	Variational graph autoencoder (VGAE)	ST and histological images	Mouse Brain from Literature and BRCA from 10× Genomics	Mouse Brain (no disease) and BRCA (invasive ductal carcinoma, IDC)	Pseudo-time
[40]	Embed a joint multiomics representation to cluster metabolite state.	VAE	ST, including metabolomic and proteomic	Self-obtained	Human lung cancer, tonsil tissue, and endometrium tissue	Metabolite state
[41]	Combine ST and histological data into embeddings that identify cancer-related cell states in the brain.	Graph Attention Autoencoders (GATE)	ST, histological images, and segmentation	10× Genomics	Human DLPFC, ovarian, and breast cancers	Pseudo-time
[42]	Embed spatial CNV and histology data from six cancers to identify spatial subclones and reveal evolutionary trees.	AE with a Vision Transformer (ViT) backbone	ST and Histology	Self-obtained	BRCA, CRC, PDAC, ccRCC, UCEC, and CHOL	Pseudo-time

first appeared in arXiv at the end of 2013 [4] and later published at ICLR 2014) through the end of 2024, given our literature search was conducted in January 2025. Publication status bias may also be present: preprints were not excluded, however, only one appeared in our final selection, which may introduce a bias against negative results.

During the screening process, we observed that the majority of journal articles returned by our queries addressed cancer diagnosis, subtyping, and prognosis (see Fig. 3), highlighting the considerable research efforts dedicated to these areas. The screening process also revealed a potential application bias within the field, with breast and colorectal cancers being overrepresented (four studies each) relative to other cancer types, and a strong emphasis on analyses using the TCGA dataset.

Results

The five research questions stated on the Introduction are addressed in the following subsections.

VAEs have been used to study cancer diagnosis, prognosis, and subtyping

The most common applications of DRL or DL in cancer are related to diagnosis, prognosis, or subtyping (Fig. 3). Studies of cancer diagnosis aim to distinguish between cancerous and non-cancerous samples [43], or between cancer types [44]. Prognosis studies aim to predict the outcome of a disease, which usually involves a survival analysis [45], recurrence prediction [46], or treatment response [47]. Subtyping

studies aim to identify distinct subtypes, i.e. groups into which a given cancer type can be classified [48, 49]. All three applications are crucial for developing more personalized treatments that would improve the outcome of the disease.

All of these and other cancer studies rely on the use of high-throughput omics data, which yield information at varying levels of resolution. Bulk omics data have been analyzed using VAEs to learn biologically meaningful low-dimensional representations, enabling the classification of different cancer types in large-scale resources such as The Cancer Genome Atlas (TCGA) [15]. Classically, the most accessible data have been bulk data, where the sampled unit is a whole tissue or group of cells. However, the use of single-cell data has been increasing, as it provides more detailed resolution, sampling information from individual cells. For instance, VAEs have been applied to single-cell data to detect CNVs with the objective of probing intratumor heterogeneity [50]. Moreover, spatially resolved omics data have recently surged in popularity [51], additionally revealing the spatial distribution of cells in a tissue. Precisely Chang *et al.* [51] developed a VAE-based methodology capable of distinguishing healthy tissue and different breast cancer types. Omics data also vary in the type of information they provide: transcriptomics data reveal gene expression levels, proteomics data reveal protein levels, and metabolomics data reveal metabolite levels. Besides, omics may be studied in combination, known as multi-omics data [45, 52]. Namely, Chaudhary *et al.* [16] used an AE to link multi-omic features to the varying survival outcomes in hepatocellular carcinoma. Recent studies have focused on multiomics and single-cell data, with a recent

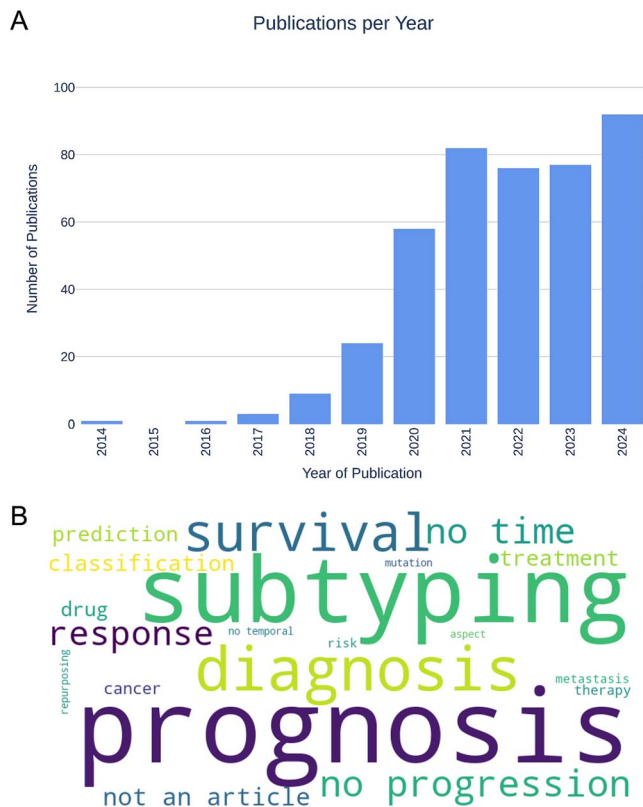


Figure 3 Results from queries. (A) Number of publications found, per year. There is a clear increase in relevant publications from 2020. (B) Word cloud of exclusion reasons. Most publications use deep and representation learning applied to cancer for subtyping, diagnosis, and studying prognosis and survival.

emphasis on spatially resolved data, given they provide very detailed information and, as such, are expected to reveal new insights into cancer biology. Wang *et al.* [53] used multi-omics, single-cell data to discern rare cell populations in melanoma. The data included in these studies may go beyond molecular data, including imaging data such as radiomics [54] and histological [55] data. Histological images and DRL have been used to cluster the different areas of lung cancer whole-slide images in order to improve diagnosis [56]. Altogether, the surveyed literature shows that DRL methods, especially VAEs, are widely applied to cancer diagnosis, prognosis, and subtyping using omics data, but less frequently focus on temporal information.

Longitudinal omics data remain limited in cancer research

The data referenced in the previous section tend to be captured at a single time point. Here, we distinguish between primary and metastatic cancer, which exhibit different growth and behavioral patterns, and focus our analysis on primary tumors. Given the dynamic nature of primary cancer, longitudinal omics data are crucial for understanding disease progression. However, collecting such data from real patients is challenging due to factors including retention difficulties, measurement timing, resource limitations, and ethical considerations. In addition, sequencing-based assays are typically destructive, precluding repeated measurements of the same biological sample. As a result, direct observation of cancer progression in real time is not feasible, either at the patient or cellular level. Even when temporally annotated data are available, inter-patient heterogeneity

leads to unaligned observations across individuals. Consistent with these limitations, we observed that the use of longitudinal omics data in primary cancer studies remains scarce.

We found a singular case [23], where extracellular vesicle RNA (EVRNA) data from liquid biopsies of colorectal cancer were collected from three cohorts: no recurrence, recurrence, and stable disease. The authors used a DL model to classify temporal subtypes from Gene Set Enrichment Analysis (GSEA) pathways. They were able to find patients that changed subtype during follow-up from the liquid biopsy data, ahead of imaging-based diagnosis. This finding confirms an extra layer of complexity in the study of cancer progression: cancer subtypes may be time-dependent and may be observed through molecular data.

A study following up on cancer patients is [21], where the authors used EHR data of patients with pancreatic cancer. They used a model based on GRU and Transformers to embed event features, encode trajectories in a latent space, and predict a time-dependent risk of occurrence. The study shows how representation learning enables out-of-sample predictions when longitudinal, aligned data are available.

The use of longitudinal data and its study with DL is more common in other diseases. During our review, we found examples of longitudinal studies on CVD, AD, kidney function deterioration, and the effects of radiation. Lee *et al.* [24] trained GRUs on radiomics images of the brain and clinical assessments to learn a representation for binary classification of cognitive decline in AD. Chung *et al.* [25] used a Long Short-Term Memory (LSTM)-VAE trained on multi-omics (proteomics and metabolomics) to embed time-dependent data obtained during cardiac remodeling. The embeddings were then used to cluster temporal growth patterns, determined to be biologically relevant through Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway analysis. Time-dependent microarray data in [26] was collected at different time points pre- and postmyocardial infarction to predict the risk of heart failure. They used a DL model to combine interactomic, molecular, and clinical data, used alongside GSEA to discover biomarkers of progress toward myocardial infarction. Another study, Hölischer *et al.* [27] assessed progression of kidney function deterioration with histomorphometry data, applying DL for semantic segmentation. The authors used trajectory and pseudo-time unsupervised analysis, adapted from single-cell methodologies, to elucidate patterns in disease progression towards kidney failure. Pseudo-time trajectories recapitulated the patterns of damage observed in histology. Brown *et al.* [28] collected longitudinal data of blood samples and radiomics from mice lungs exposed to radiation, using DL for survival analyses from the radiomics data. Altogether, these studies highlight the importance of DL and DRL in the study of disease progression, and the potential of longitudinal data to provide insights into the evolution of diseases.

The study of cancer evolution with DL and DRL is common in phylogeny, albeit mostly to discern subclones. A 2020 study [29] claimed to be the first to apply DL to study tumor evolution. Specifically, RL was used to infer the most likely tumor phylogeny, and identify tumor subclones. Baciu-Drăgan and Beerenwinkel [30] used unsupervised learning to learn graph embeddings of relationships between subclones. Hence, phylogeny studies of cancer subtypes offer a way to sort subtypes in terms of evolutionary terms. However, while phylogenetic data and studies reflect the evolutionary history of tumors and the subclone ordering, they do not provide explicit temporal resolution in the same way as clinically annotated stages or time-resolved sampling.

Stages as a way to study cancer in a longitudinal manner

Thus far, we have considered actual time as the unit for the study of cancer progression. However, collecting large-scale longitudinal data at regular intervals across multiple patients is often impractical. In addition, tumors evolve at different rates in different patients, so progression is not uniform. In this context, clinically defined stages could serve as a surrogate time unit to study cancer progression. Stages refer to the extent of cancer [57, 58]. The most common staging system is the Tumor, Nodes, and Metastasis (TNM), which takes into account the size of the tumor, its spread to nearby lymph nodes, and whether it has metastasized. Thus, there are five stages: stage 0 describes a non-invasive tumor, or carcinoma *in situ*; stages I through stage III denote an invasive tumor, the more invasive, the larger the number; while stage IV refers to a metastatic tumor that has spread to distant body parts. Hence, from a conceptual standpoint, cancer stages such as TNM classifications could act as a proxy for temporal progression in the absence of longitudinal omics data, allowing cross-sectional samples to be aligned for trajectory inference, as shown in works such as [59] for GRN reconstruction.

The results from our search highlight a common trend: the difficulty of differentiating cancer stages. As demonstrated in Nelligan *et al.* [31], this is addressed by comparing normal tissue to samples from successive cancer stages. Authors used adversarial networks to transform a source sample, normal tissue from TCGA and the Genotype-Tissue Expression (GTEx), into the feature space of a target sample, the different stages, finding several up- and down-regulated genes across stages. Another way to circumvent this issue is to group stages together [32–35]. In [35], data from TCGA was used to classify either stages I and II or stages II and III across different cancers. The classification was compared on real data across different DRL techniques, namely Principal Component Analysis (PCA), Stacked Denoising AE, and VAE. The authors found that representation learning techniques did not improve classification performances. Lee *et al.* [32] grouped stages I and II as an “early” stage, and stages III and IV as a “late” stage from papillary renal cancer data. Stages may be simplified even further to distinguish between normal tissue, non-invasive abnormally-grown tissue, and invasive cancer. Long *et al.* [34] used DL to distinguish between normal, neoplasm, and cancer samples from cervical cancer microarray data. Besides molecular data, histological images may also be used to classify stages. As shown by Zhang *et al.* [33], a Convolutional VAE was used to learn embeddings of cell states from chromatin images of breast cancer. The authors considered three stages: non-tumor, DCIS, and IDC. The latent embeddings were used to infer a pseudo-time ordering of the stages with Partition-based graph abstraction (PAGA), a single-cell trajectory inference method, revealing an ordering from non-tumor to tumor progression. Hence, there is a growing interest in the correct characterization of stages, but its hard separability has led to simplifications in the stages considered.

The lack of separability between stages may be due to the typically unbalanced distribution of samples across stages. In this context, generative modeling may be useful, creating synthetic instances in order to balance out the data. Dogan *et al.* [36] argued that the use of the conditional variational autoencoder (CVAE) may be used to generate the stages of breast cancer in mRNA and miRNA data from TCGA. They trained the CVAE on the molecular data and used labels of the different stages: I to IV and took solid tissue normal as pseudo-stage 0. The CVAE is then able to generate new, synthetic samples at the trained labels (here, stages). These synthetic stages then showed high separability in a hierarchical clustering analysis. While the study

did not report the model loss, it provides an opportunity to reflect on the importance of ensuring that generative models accurately capture the distribution of the original data, and whether synthetic samples are representative and properly validated.

Single-cell omics data are commonly used to infer pseudo-time trajectories of cancer cells

Single-cell data provide higher-resolution measurements than bulk data, as they profile individual cells rather than aggregate signals from tissue samples. In terms of disease progression, single-cell data can more accurately capture the individual differences between cells within the same tissue. Hence, single-cell data enable analyses that order cells according to transcriptional similarity, allowing the inference of differentiation trajectories. These approaches organize cells along a pseudo-time axis, and a variety of methods have been developed for this purpose, such as Monocle [60] and PAGA [61]. During our search, we found that single-cell data were more commonly used for trajectory inference than bulk or imaging data. These data were often processed with DRL techniques to generate lower-dimensional embeddings, which were then used as representations to infer pseudo-time trajectories. For example, the VAE was used by Rashid *et al.* [37] to embed data from Glioblastoma and Oligodendroglioma, and these embeddings were used to infer lineage and differentiation trajectories.

We also found that most of the results from our search that used single-cell data included ST data, providing information on the spatial localization of cells within the tissue. Besides pseudo-time inference, spatial data also allow to link the different tissue layers. Chen *et al.* [38] used a contrast learning (CL) method to learn latent embeddings of cells and their spots from ST data of the mouse brain. The embeddings were then used for pseudo-time and trajectory inference with PAGA, organizing the different tissue layers. CL was also used by Zong *et al.* [39] on the mouse visual cortex, combining ST with histological images to learn embeddings of similar spots across tissue regions. PAGA and Monocle were used for pseudo-time trajectory inference, indicating the progression in the developmental process from the white matter layer towards outer layers. Thomas *et al.* [40] used a VAE to embed a joint representation of spatial proteomics and metabolomics data from B cells in the tonsils, and then used these embeddings for metabolite state clustering. The trajectory inference was able to reconstruct the cell differentiation process of B cells in a spatial domain. Zuo *et al.* [41] used a graph attention AE and multiview collaborative learning to embed ST and histological data from human brain slices. The embeddings were used to identify cancer-related cell states: stemness, migration, and metastasis. Mo *et al.* [42] combined an AE with a visual transformer to obtain embeddings from spatial CNV and histology data from six cancer types in order to identify spatial subclones and reveal their evolutionary tree. Altogether, the literature indicates the relevance of combining ST data and DRL to relate pseudo-time trajectories with the spatial conformation of healthy and cancerous tissues.

Limitations and challenges

From the studies included in our review, we observed some common trends that limit the study of cancer progression with DRL and DL. First, applications of DL and DRL in cancer omics studies are mostly centered on subtyping, diagnosis, and prognosis, which typically do not explicitly model the longitudinal aspect of the disease. Studies that do take into account a temporal aspect tend to rely on single-cell data. These approaches typically focus on cellular-level

dynamics rather than patient-level progression. Compared with bulk data, single-cell data add more persample noise, as it includes a level of cellular heterogeneity that is not present in bulk data. Moreover, single-cell studies frequently capture relative cellular progression based on transcriptional similarity, rather than the true temporal dynamics of the disease. Finally, clinically defined stages, even when used as proxy time units for cancer progression, are difficult to classify accurately, making robust inference of stage-wise progression from existing omics datasets challenging.

These limitations imply some challenges in the study of cancer progression inference from omics data. First, the lack of longitudinal data is a major challenge, primarily because collecting repeated samples from patients over time is difficult due to logistical, ethical, and clinical constraints, including the impact of ongoing treatments on disease progression. Even when data are collected longitudinally, which is usually more frequent in experimental animal data than human data [62, 63], the rate of cancer progression varies across individuals, so measurements taken at the same nominal time point may not correspond to the same disease stage. A second challenge is the limited availability of methodologies to infer real-time progression. The use of generative models to augment existing data may help address this issue, potentially enabling inference of out-of-sample behavior and effectively introducing a temporal dimension to the dataset. However, this approach leads to a third challenge: the scarcity of data available for model validation. These challenges must be addressed in order to provide a more accurate understanding of cancer progression, and to develop more personalized treatments.

Discussion

DL is a commonly used methodology to learn complex patterns in health data. DRL comprises a specific type of DL methods that have the ability to learn a lower-dimensional representation of complex data distributions. This ability has been combined with generative capabilities in the VAE, a Bayesian-statistics-based representation learning method. Given the complexity of high-throughput data obtained in cancer research, these methods have been applied to study the disease. However, the dynamic nature of cancer implies multiple temporal aspects that remain unexplored.

In this study, we conducted an SLR on the use of DRL in cancer studies with omics data for modeling time-related processes. To encompass the complexity of the topic, we created a query that was adapted to three different search engines (see Methods). Our search returned 440 papers, of which 21 were relevant to our investigation (see Tables 1 and 2 for a summary). We summarized these papers and categorized them into four main groups. (i) Discarded papers were analyzed to identify the reasons for their exclusion, revealing that the most common applications of DRL and VAEs in cancer are related to subtyping, diagnostics, and prognosis studies. (ii) These techniques have been applied in some diseases, taking into consideration a temporal aspect, but their application in longitudinal cancer studies remains limited. (iii) Stages may serve as a proxy to align cancer data in a temporal manner. (iv) Single-cell omics data are the most frequently used type in cancer studies, with a recent surge in spatially resolved data; however, these studies arrange cells in pseudo-time. Collectively, these findings reveal that the study of cancer progression in time has not yet been adequately explored in omics-based studies.

The temporal capabilities and the interpretability of the models used in the 21 relevant papers have been summarized and compared

in Table 3. Overall, this comparison highlights a trade-off between explicit temporal modeling (e.g. recurrent and attention-based architectures) and latent-space interpretability and generative flexibility, with VAEs uniquely balancing both aspects in omics-based cancer studies. VAEs offer distinct advantages for temporal omics analysis through their latent space representations. Unlike standard neural networks, VAEs learn structured low-dimensional embeddings that preserve data geometry while enabling generative capabilities [36]: synthetic data augmentation, missing value imputation, and trajectory interpolation along temporal dimensions.

The probabilistic latent space enhances interpretability, both by preserving the original data distribution and revealing feature contributions to latent dimensions [64]. The VAE may also be further specialized when used in combination with models specific for sequential learning, such as LSTMs [25]. These properties make VAEs particularly well suited, compared with non-generative approaches, for reconstructing cancer dynamics from cross-sectional data. However, the VAE still faces some limitations. First, standard VAEs are typically trained in an unsupervised manner with respect to temporal information; conditional variants, such as the CVAE [36], partially address this limitation by incorporating temporal or clinical metadata during training. Second, the highly nonlinear and heterogeneous nature of cancer progression poses challenges for modeling complex temporal distributions within the assumptions of conventional VAEs. During our review, we observed that, to overcome these limitations, the VAE has been used either in combination with other models, such as LSTMs [25], or designed with more complex architectures, such as Stacked AEs [35]. Nonetheless, these approaches often come at the expense of reduced interpretability. Finally, VAEs may struggle with out-of-sample generalization, as latent representations learned from limited or homogeneous cohorts may not transfer robustly across patient populations, increasing the risk of generating spurious or biologically implausible trajectories. This issue is analogous to the phenomenon of hallucinations commonly observed in generative AI models [65].

We have attempted to make our search as comprehensive as possible. For this reason, we used multiple search engines and created comprehensive queries, combining DRL terms, omics data types, cancer terminology, and temporal concepts. These queries were designed to cover relevant keywords related to our search questions, while also considering synonyms and alternative terms for the keywords. However, we acknowledge these limitations may have missed relevant VAE-specific temporal studies using alternative terminology or methodological descriptions and did not return papers that could be relevant in the domain. Namely, we identified a number of articles that are not captured by our specific queries but we decided to discuss anyway here due to their relevance. For example, Mora *et al.* [66] used a VAE to study the differences across stages of clear renal cell carcinoma using different layers of multiomics data (DNA methylation, gene and protein expression). They revealed genes that were significantly differentially expressed between early (stages I and II) and late stages (stages III and IV). Although not focusing on cancer, another interesting related approach can be found in [64], where they used a VAE to embed bulk RNA-seq data during mouse Central Nervous System (CNS) development. The VAE identified genes with qualitatively different functional profiles and multivariate trends. The VAE was compared with other embedding methods (PCA, tSNE, UMAP, and PHATE), and was determined to be the most trustworthy at distinguishing genes with known but different anterior-posterior axis association and groups that displayed similar biological enrichment.

Table 2 Data and code summary table. This table contains the data and code information of the final list of 21 papers obtained from the SLR. The columns represent: Reference, Dataset size, External validation, Data availability, and Code availability.

Reference	Dataset size	External validation	Data availability	Code availability
[23]	155 plasma samples from 71 patients + 70 healthy controls	Consensus Molecular Subtypes (CMS) classifier applied on the TCGA colorectal cancer ($n = 647$).	GSE255775	Not available
[21]	Denmark: 6.2 M patients; US-VA: 3 M patients	Cross-application (Denmark to US-VA).	Under request to Danish Health Data Authority and VA Boston Healthcare System IRB	Available
[24]	Four modalities of data: cognitive performance ($n = 802$), CSF ($n = 601$), MRI ($n = 865$), and demographics ($n = 865$)	Simulation data of 1000 samples, 500 of which have all data modalities.	Under request to ADNI	Not available
[25]	72 treatment mice and 84 controls	Pathway enrichment to identify which method enriches more significant pathways.	Not available	Not available
[26]	9 premedicated pigs + 3 controls	Pathway enrichment with GSEA.	GSE34569	Not available
[27]	1743 WSIs from 1043 cases	External cohorts (KPMP and HuBMAP).	Raw WSIs available under request to authors; external data from respective consortia	Available
[28]	36 irradiated mice (18 C57BL6, 18 C3H/HeN) + 36 controls	Histology and pathway enrichment.	Not available	Not available
[29]	1000 simulation samples and 3 real datasets (BRCA 16, Acute Lymphoblastic Leukemia 146 cells, and Colorectal cancer 72)	External datasets with known phylogenetic patterns.	Simulation data available in code repository	Available
[30]	16 groups of synthetic trees, 43 tumor mutation trees of 6 cancers, and 123 Acute Myeloid Leukemia trees	Validation against known cancer evolution biology (linear vs branching patterns).	Available in code repository	Available
[31]	432 THCA tumors + 369 normal tissue samples	None reported.	TCGA-THCA	Not available
[32]	259 KIRP samples + 32 normal tissue samples	None reported.	TCGA-KIRP	Available as supplementary material
[33]	560 tissue microarray (TMA) samples from 122 female patients across 3 disease stages (non-tumor, DCIS, and IDC)	Pathologist blind-test assignment of tumor grade.	PSI Public Data Repository	Available
[34]	202 cancer, 115 cervical intraepithelial neoplasia, and 105 normal samples from 8 datasets	Two datasets used for external validation.	Multiple datasets from GEO (see Table 1 in the publication)	Not available
[35]	Almost 40 000 samples	None reported.	recount2	Available
[36]	TCGA-BRCA: 1072 tumor + 104 normal samples. starBase : 863 066 miRNA-mRNA interactions	Pathway enrichment.	TCGA-BRCA and starBase	Available
[37]	6523 cells across 4 tumors	Separation of malignant and non-malignant cells; differentially expressed genes coincide with literature.	GSE70630 , GSE57872 , GSE72056 , GSE89567	Not available
[38]	Human Dorsolateral prefrontal cortex (12 slices), human ovarian cancer (1 sample, 3493 spots), human breast cancer (1 sample, 4727 spots), mouse hypothalamus MERFISH (1,365 cells)	Cross-species and marker genes.	Human DLPFC , Mouse brain , Phalaenopsis , Human breast cancer , Human PDAC (GSE111672) , Mouse olfactory bulb , Mouse hypothalamus	Available
[39]	Human Dorsolateral prefrontal cortex (12 slices), rest unspecified: mouse hypothalamus MERFISH, mouse visual cortex seqFISH, human breast cancer (10× Visium), mouse olfactory bulb Stereo-seq	Embeddings validated with Monocle3 and PAGA.	Multiple public repositories: spatialLIBD , 10× Genomics, and literature	Available
[40]	Human lung cancer (7 patients, 21 regions, and 19 507 cells), human tonsil (2 samples, 11 regions, and 31 156 cells), human endometrium (5 regions, 8215 cells)	Multi-tissue application.	IMC and 3D-SMF image data , proteomics , and metabolomics	Available

(continued)

Table 2 Continued.

Reference	Dataset size	External validation	Data availability	Code availability
[41]	Human Dorsolateral prefrontal cortex (12 slices), human ovarian cancer (for training, 1 sample, 3493 spots; for validation, 20 samples with 24,489 cells), human breast cancer (for training, 1 sample, 4,727 spots; for validation, 1 sample, 4,081 cells), mouse primary visual cortex (1 sample, 1,365 cells)	Multi-tissue application.	10 × Genomics (ovarian & breast cancer), spatialLIBD (human DLPFC), ClusterMap (mouse visual cortex), GSE176078 (breast cancer scRNA-seq), E-MTAB-8859 (ovarian cancer scRNA-seq), TCGA via Xena	Available
[42]	131 Visium ST sections from 78 cases across 6 cancer types (54 BRCA, 30 CRC, 23 PDAC, 12 RCC, 5 UCEC, and 7 CHOL)	Co-detection by indexing (CODEX) samples matching.	HTAN dbGaP phs002371.v3.p1, HTAN DCC Portal (WUSTL Atlas)	Available

Table 3 Model architectures. Comparison of key characteristics of the models assessed during the SLR process

Model	Temporal modeling	Interpretability	Reference
VAE (including convolutional, conditional, graph, and attention-based)	Able to embed several disease stages, and to generate data at learned and intermediate stages.	Interpretable latent space.	[33, 35–41]
LSTM and GRUs	Learns long-term relationships. Flexible time intervals.	Requires specific algorithms [67]	[21–25, 29]
Stacked AEs	Learns very complex relationships between variables in data by concatenating several embeddings. Good for feature extraction.	Limited because of the complex relationships between embeddings.	[35]
Transformers	Learns long-range relationships.	Requires attention algorithms.	[21, 42]
Neural networks	Mostly used for classification and image segmentation.	Low, but lots of algorithms available.	[23, 26–28, 30–32, 34]
RL	For sequential decision-making.	Needs specific algorithms to analyze.	[29]

The lack of full matches from our queries may reflect substantial challenges and limitations in the field. Research primarily focused on addressing the high dimensionality of omics data through embedding inference, which is subsequently used for downstream tasks. These tasks rarely capture temporal cancer dynamics, except in single-cell studies that leverage embeddings to infer pseudo-time trajectories of cells. However, pseudo-temporal ordering does not provide additional longitudinal information, such as clinically diagnosed stages or the actual elapsed time between cellular states. Such longitudinal omics data are difficult to obtain, and when available, it is often collected during treatment and involve destructive sampling. As a result, omics studies of cancer progression frequently rely on data, i.e. not temporally aligned across patients or samples.

The application of the VAE specifically, and DRL more generally, to model temporal progression in cancer omics studies remains under-explored. The field faces significant challenges. On the one hand, some challenges are not specific to time-dependent datasets. Data collected from cancer patients are often heterogeneous, whether it is because of extrinsic reasons, such as different sequencing technologies available

and data-processing techniques used, and/or intrinsic reasons, such as the inherent inter-personal heterogeneity of human biology. On the other hand, some challenges are specific to the temporal aspect. These include unaligned data, either because of missing patients follow-ups, or because of sampling at unique time-points or stages of the disease. Furthermore, additional challenges arise when considering multiomics data, due to missing profiles across different omics layers and due to the integration of clinical data [68].

A critical, specific challenge is the need for validation strategies, which would require proper datasets that include the longitudinal dimension. Potential resources include TCGA, which provides clinical stage annotations, TRACERx, which tracks patients over time with multi-region sampling [69], Patient-Derived Xenograft (PDX) models [70], and pseudo-temporal information derived from large-scale single-cell atlases [71]. Validation of synthetic data necessitate ensuring that generated samples meet three key criteria: fidelity, meaning that the original data distributions and biological structure are preserved; utility, ensuring that synthetic data remain informative for downstream tasks such as subtyping and prognosis; and privacy,

preventing the re-identification of patients from generated data [72]. When detailed, time-annotated datasets are available, they enable benchmarking trained models on their ability to produce temporally aligned and out-of-sample trajectories, as well as to simulate potential interventions, such as treatment responses. This approach has been implemented successfully in single-cell studies for pseudotemporal reconstruction, where generative models are evaluated against known developmental hierarchies [73, 74]. Extending similar frameworks to bulk tumor longitudinal data remains an important goal. Addressing these challenges will require open, curated temporal datasets, standardized validation frameworks for synthetic omics data as well as multidisciplinary collaboration. Developing such resources and practices will be critical to advancing VAEs and DRL from exploratory tools toward validated methods for studying temporal cancer progression.

An important goal of the application of AI to cancer research is to have an impact in the clinical oncological setting. Generative AI is no exception. It has shown significant potential at improving cancer diagnostic [75], enabling multimodal data integration [37], and compensating for insufficient data [21]. However, there are significant challenges when it comes to the translation of generative AI to the clinic. A key limitation is the lack of interpretability of model outputs [76]. Although explainability techniques such as SHAP [77] and LIME [78] exist, our systematic review indicates that deep generative models are frequently applied without accompanying explainable analyses. Besides, generative AI relies on the characteristics of the data that it learns from, meaning it could replicate or even amplify potential existing biases [79, 80]. These challenges raise ethical and regulatory uncertainties [81] that would require rigorous validation to ensure clinical safety and efficacy [82, 83] in order to prevent synthetic biases and hallucinations.

There has been a recent increase in interest in generative models, which can be leveraged to align data, including in the context of temporal analyses. European projects, such as [EVENFLOW](#), are

delving into the application of such methodologies to overcome current limitations in cancer studies. Moreover, synthetic data generation methods based on simulation processes, such as agent-based modeling of cellular dynamics [84, 85], are making strides in integrating physics into generative frameworks to produce more realistic and reliable synthetic data. Other initiatives aim to unite experts across diverse disciplines to tackle the challenges of AI in biomedical applications, including generative AI and cancer, as exemplified by the [AHEAD](#) project, which brings together specialists in biomedicine, AI, ethics, law, psychology, and related fields. These approaches, together with investments in generating suitable datasets for validation, could pave the way for the development of personalized treatments that account for individual differences in cancer progression.

List of abbreviations

AD, Alzheimer's Disease; AE, Autoencoder; AI, Artificial Intelligence; CL, Contrast Learning; CNS, Central Nervous System; CNV, Copy Number Variation; CVAE, Conditional Variational Autoencoder; CVD, Cardiovascular Disease; DCIS, Ductal Carcinoma In Situ; DL, Deep Learning; EHR, Electronic Health Record; evRNA, extracellular vesicle RNA; GSEA, Gene Set Enrichment Analysis; GRU, Gated Recurrence Unit; GRN, Gene Regulatory Network; GTEx, Genotype-Tissue Expression; IDC, Invasive Ductal Carcinoma; KEGG, Kyoto Encyclopedia of Genes and Genomes; LIME, Local Interpretable Model Agnostic Explanations; LSTM, Long Short-Term Memory; ML, Machine Learning; NN, Neural Network; DRL, Deep Representation Learning; RNN, Recurrent Neural Network; PAGA, Partition-based Graph Abstraction; PCA, Principal Component Analysis; PRISMA, Preferred Reporting Items for Systematic Reviews and Meta-Analyses; PDX, Patient-Derived Xenograft; RL, Reinforcement Learning; SC, Single-Cell; scRNA-seq, Single-Cell RNA Sequencing; SHAP, SHapley Additive exPlanations; SLR, Systematic Literature Review; ST, Spatial Transcriptomics; TCGA, The Cancer Genome Atlas; TNM, Tumor, Nodes, and Metastasis; TRACERx,

Dictionary of terms

Term	Definition
Autoencoder (AE)	DRL, unsupervised technique that encodes unlabeled data into a lower-dimensional representation, the latent space. It consists of an encoder, which embeds the input data, and a decoder, recreating the input data from the learned embeddings.
Conditional Variational Autoencoder (CVAE)	A VAE used for semi- or fully-supervised learning. In this case, the encoder and decoder are passed the labels of the input data and learn to encode and decode the input data depending on, i.e., conditioned to, the input data.
Convolutional Variational Autoencoder	A VAE that entails the use of convolutional layers in the encoder and decoder, making it suitable to learn latent-space embeddings of images.
Contrast Learning	A supervised DRL method that learns embeddings of data by comparing, or contrasting, between similar and dissimilar instances within the input data.
Deep Learning (DL)	Discipline within Machine Learning that uses artificial, multilayer neural-networks to learn patterns present in a given input dataset.
Deep Representation Learning (DRL)	Also known as Deep Feature Learning, consists of techniques to automatically estimate a representation of the features included on a given dataset. We use the adjective Deep to remark the fact that we consider Deep-Learnig-based techniques.
Long Short-Term Memory (LSTM)	Specialized DL model designed to learn when to retain or forget relevant information present in sequential data.
Reinforcement Learning (RL)	Area of Machine Learning where an automated agent learns a process of decision-making while interacting with an environment and receiving feedback (a reward or a penalty) from an external interpreter.
Variational Autoencoder (VAE)	DRL technique that follows the AE ability to learn embeddings from unlabeled data, albeit considering a variational Bayesian inference to learn a probabilistic latent space.

TRacking Cancer Evolution through therapy (Rx); VAE, Variational Autoencoder

Key Points

- There is a growing interest on the application of deep learning methods, such as deep representation learning (DRL), to cancer studies.
- Cancer is a complex and dynamic disease, whose temporal dynamics are not yet fully captured in omics-based studies.
- Among DRL methods, the Variational Autoencoder (VAE) using omics-based data has been widely used in cancer studies, particularly for subtyping, diagnosis, and prognosis.
- The temporal aspects of cancer progression are often insufficiently captured in omics-based studies, primarily due to the scarcity of longitudinal data.
- Applying the VAE as a generative model to study cancer in time, such as focusing on cancer staging, could lead to significant advancements in our understanding of cancer.

Author contributions

G.P.C. developed the SLR, and wrote and revised earlier and consolidated versions of the manuscript. A.V. and D.C. conceived the study and oversaw its execution. Study selection, data extraction, and risk of bias assessment were performed by G.P.C., with results independently verified by D.C. All authors reviewed and contributed to the final manuscript.

Conflicts of interest

None declared.

Funding

This work has been supported by the EU project EVENFLOW under Horizon Europe agreement no. 101070430.

Data availability

The results from the queries, as well as the process of literature selection, including independent verification and conflict resolution, are publicly available at the GitHub repository [gprocastelo/SLR-VAE](https://github.com/gprocastelo/SLR-VAE).

References

- Miotto R, Wang F, Wang S *et al*. Deep learning for healthcare: review, opportunities and challenges. *Brief Bioinform* 2018;**19**:1236–46. <https://doi.org/10.1093/bib/bbx044>
- Bengio Y, Courville A, Vincent P. Representation learning: a review and new perspectives. *IEEE Trans Pattern Anal Mach Intell* 2013;**35**:1798–828. <https://doi.org/10.1109/TPAMI.2013.50>
- Bercea CI, Wiestler B, Rueckert D *et al*. Evaluating normative representation learning in generative AI for robust anomaly detection in brain imaging. *Nat Commun* 2025;**16**:1624. <https://doi.org/10.1038/s41467-025-56321-y>
- Kingma DP, Welling M. *Auto-Encoding Variational Bayes* 2013. arXiv Version Number: 11. <https://arxiv.org/abs/1312.6114>
- Wang Y, Miller AC, Blei DM. Comment: variational autoencoders as empirical Bayes. *Stat Sci* 2019;**34**:229–233. <https://doi.org/10.1214/19-STS710>
- Roth GA, Abate D, Abate KH *et al*. Global, regional, and national age-sex-specific mortality for 282 causes of death in 195 countries and territories, 1980–2017: a systematic analysis for the global burden of disease study 2017. *The Lancet* 2018;**392**:1736–88. [https://doi.org/10.1016/S0140-6736\(18\)32203-7](https://doi.org/10.1016/S0140-6736(18)32203-7)
- Hulvat MC. Cancer incidence and trends. *Surg Clin North Am* 2020;**100**:469–81. <https://doi.org/10.1016/j.suc.2020.01.002>
- Cirillo D, Valencia A. Big data analytics for personalized medicine. *Curr Opin Biotechnol* 2019;**58**:161–7. <https://doi.org/10.1016/j.copbio.2019.03.004>
- International Human Genome Sequencing Consortium. Initial sequencing and analysis of the human genome. *Nature* 2001;**409**:860–921. <https://doi.org/10.1038/35057062>
- Tang F, Barbacioru C, Wang Y *et al*. mRNA-Seq whole-transcriptome analysis of a single cell. *Nat Methods* 2009;**6**:377–82. <https://doi.org/10.1038/nmeth.1315>
- Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet* 2009;**10**:57–63. <https://doi.org/10.1038/nrg2484>
- Patterson SD, Aebersold RH. Proteomics: the first decade and beyond. *Nat Genet* 2003;**33**:311–23. <https://doi.org/10.1038/ng1106>
- Idle JR, Gonzalez FJ. Metabolomics. *Cell Metab* 2007;**6**:348–51. <https://doi.org/10.1016/j.cmet.2007.10.005>
- Ståhl PL, Salmén F, Vickovic S *et al*. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science* 2016;**353**:78–82. <https://doi.org/10.1126/science.aaf2403>
- Withnell E, Zhang X, Sun K *et al*. XOMiVAE: an interpretable deep learning model for cancer classification using high-dimensional omics data. *Brief Bioinform* 2021;**22**:bbab315.
- Chaudhary K, Poirion OB, Lu L *et al*. Deep learning-based multi-omics integration robustly predicts survival in liver cancer. *Clin Cancer Res* 2018;**24**:1248–59. <https://doi.org/10.1158/1078-0432.CCR-17-0853>
- Dagogo-Jack I, Shaw AT. Tumour heterogeneity and resistance to cancer therapies. *Nat Rev Clin Oncol* 2018;**15**:81–94. <https://doi.org/10.1038/nrclinonc.2017.166>
- Cao K, Schwartz R. Computationally reconstructing the evolution of cancer progression risk. *Comput Adv Bio Med Sci* 2024:212–223.
- Frank SA. Somatic evolutionary genomics: Mutations during development cause highly variable genetic mosaicism with risk of cancer and neurodegeneration. *Proc Natl Acad Sci USA* 2010;**107**:1725–30. <https://doi.org/10.1073/pnas.0909343106>
- Page MJ, McKenzie JE, Bossuyt PM *et al*. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *BMJ* 2021;**372**:71. <https://doi.org/10.1136/bmj.n71>
- Placido D, Yuan B, Hjaltelin JX *et al*. A deep learning algorithm to predict risk of pancreatic cancer from disease trajectories. *Nat Med* 2023;**29**:1113–22. <https://doi.org/10.1038/s41591-023-02332-5>
- Harzing AW. *Publish or Perish* 2007. <https://harzing.com/resources/publish-or-perish>.
- Bahrambeigi V, Lee JJ, Branchi V *et al*. Transcriptomic profiling of plasma extracellular vesicles enables reliable annotation of the cancer-specific transcriptome and molecular subtype. *Cancer Res* 2024;**84**:1719–32. <https://doi.org/10.1158/0008-5472.CAN-23-4070>
- Lee G, Kang B, Nho K *et al*. MildInt: deep learning-based multimodal longitudinal data integration framework. *Front Genet* 2019;**10**:617. <https://doi.org/10.3389/fgene.2019.00617>

25. Chung NC, Mirza B, Choi H *et al.* Unsupervised classification of multi-omics data during cardiac remodeling using deep learning. *Methods* 2019;**166**:66–73.
26. Iborra-Egea O, Gálvez-Montón C, Prat-Vidal C *et al.* Deep learning analyses to delineate the molecular Remodeling process after myocardial infarction. *Cells* 2021;**10**:3268. <https://doi.org/10.3390/cells10123268>
27. Hölscher DL, Bouteldja N, Joodaki M *et al.* Next-generation morphometry for pathomics-data mining in histopathology. *Nat Commun* 2023;**14**:470. <https://doi.org/10.1038/s41467-023-36173-0>
28. Brown KH, Ghita-Pettigrew M, Kerr BN *et al.* Characterisation of quantitative imaging biomarkers for inflammatory and fibrotic radiation-induced lung injuries using preclinical radiomics. *Radiother Oncol* 2024;**192**:110106. <https://doi.org/10.1016/j.radonc.2024.110106>
29. Azer ES, Ebrahimabadi MH, Malikić S *et al.* Tumor phylogeny topology inference via deep learning. *iScience* 2020;**23**:101655.
30. Baciú-Drăgan M-A, Beerenwinkel N. Oncotree2vec — a method for embedding and clustering of tumor mutation trees. *Bioinformatics* 2024;**40**:i180–8. <https://doi.org/10.1093/bioinformatics/btae214>
31. Nelligan NM, Bender MR, Feltus FA. Simulating the restoration of normal gene expression from different thyroid cancer stages using deep learning. *BMC Cancer* 2022;**22**:612. <https://doi.org/10.1186/s12885-022-09704-z>
32. Lee S, Jung J, Park I *et al.* A deep learning and similarity-based hierarchical clustering approach for pathological stage prediction of papillary renal cell carcinoma. *Comput Struct Biotechnol J* 2020;**18**:2639–46. <https://doi.org/10.1016/j.csbj.2020.09.029>
33. Zhang X, Venkatachalapathy S, Paysan D *et al.* Unsupervised representation learning of chromatin images identifies changes in cell state and tissue organization in DCIS. *Nat Commun* 2024;**15**:6112. <https://doi.org/10.1038/s41467-024-50285-1>
34. Long NP, Jung KH, Yoon SJ *et al.* Systematic assessment of cervical cancer initiation and progression uncovers genetic panels for deep learning-based early diagnosis and proposes novel diagnostic and prognostic biomarkers. *Oncotarget* 2017;**8**:109436–56.
35. Smith AM, Walsh JR, Long J *et al.* Standard machine learning approaches outperform deep representation learning on phenotype prediction from transcriptomics data. *BMC Bioinformatics* 2020;**21**:119. <https://doi.org/10.1186/s12859-020-3427-8>
36. Dogan H, Hakguder Z, Madadjim R *et al.* Elucidation of dynamic microRNA regulations in cancer progression using integrative machine learning. *Brief Bioinform* 2021;**22**:bbab270. <https://doi.org/10.1093/bib/bbab270>
37. Rashid S, Shah S, Bar-Joseph Z *et al.* Dhaka: variational autoencoder for unmasking tumor heterogeneity from single cell genomic data. *Bioinformatics* 2021;**37**:1535–43. <https://doi.org/10.1093/bioinformatics/btz095>
38. Chen N, Yu X, Li W *et al.* A signal-diffusion-based unsupervised contrastive representation learning for spatial transcriptomics analysis. *Bioinformatics* 2024;**40**:btae663. <https://doi.org/10.1093/bioinformatics/btae663>
39. Zong Y, Yu T, Wang X *et al.* conST: an interpretable multi-modal contrastive learning framework for spatial transcriptomics. *bioRxiv* 2022:2022.01.14.476408. <https://doi.org/10.1101/2022.01.14.476408>
40. Thomas H, Allam M, Cai S *et al.* Single-cell spatial metabolomics with cell-type specific protein profiling for tissue systems biology. *Nat Commun* 2023;**14**:8260. <https://doi.org/10.1038/s41467-023-43917-5>
41. Zuo C, Zhang Y, Cao C *et al.* Elucidating tumor heterogeneity from spatially resolved transcriptomics data by multi-view graph collaborative learning. *Nat Commun* 2022;**13**:5962. <https://doi.org/10.1038/s41467-022-33619-9>
42. Mo C-K, Liu J, Chen S *et al.* Tumour evolution and microenvironment interactions in 2D and 3D space. *Nature* 2024;**634**:1178–86. <https://doi.org/10.1038/s41586-024-08087-4>
43. Sado IT, Fitime LF, Pelap GF *et al.* Early multi-cancer detection through deep learning: an anomaly detection approach using Variational autoencoder. *J Biomed Inform* 2024;**160**:104751.
44. Khorshed T, Moustafa MN, Rafea A. Deep learning for multi-tissue cancer classification of gene expressions (GeneXNet). *IEEE Access* 2020;**8**:90615–29. <https://doi.org/10.1109/ACCESS.2020.2992907>
45. Tan K, Huang W, Hu J *et al.* A multi-omics supervised autoencoder for pan-cancer clinical outcome endpoints prediction. *BMC Med Inform Decis Mak* 2020;**20**:129. <https://doi.org/10.1186/s12911-020-1114-3>
46. Fernandez-Mateos J, Cresswell GD, Trahearn N *et al.* Tumor evolution metrics predict recurrence beyond 10 years in locally advanced prostate cancer. *Nat Cancer* 2024;**5**:1334–51. <https://doi.org/10.1038/s43018-024-00787-0>
47. Zhang L, Cui Y, Zhou G *et al.* Leveraging mitochondrial-programmed cell death dynamics to enhance prognostic accuracy and immunotherapy efficacy in lung adenocarcinoma. *J Immunother Cancer* 2024;**12**:e010008. <https://doi.org/10.1136/jitc-2024-010008>
48. Li J, Li L, You P *et al.* Towards artificial intelligence to multi-omics characterization of tumor heterogeneity in esophageal cancer. *Semin Cancer Biol* 2023;**91**:35–49. <https://doi.org/10.1016/j.semcancer.2023.02.009>
49. Hassan AM, Naeem SM, Eldosoky MAA *et al.* Multi-omics-based machine learning for the subtype classification of breast cancer. *Arab J Sci Eng* 2025;**50**:1339–52. <https://doi.org/10.1007/s13369-024-09341-7>
50. Kurt S, Chen M, Toosi H *et al.* CopyVAE: a variational autoencoder-based approach for copy number variation inference using single-cell transcriptomics. *Bioinformatics* 2024;**40**:btae284. <https://doi.org/10.1093/bioinformatics/btae284>
51. Chang X, Jin X, Wei S *et al.* DeepST: identifying spatial domains in spatial transcriptomics by deep learning. *Nucleic Acids Res* 2022;**50**:e131–1.
52. Yang B, Xin TT, Pang SM *et al.* Deep subspace mutual learning for cancer subtypes prediction. *Bioinformatics* 2021;**37**:3715–22. <https://doi.org/10.1093/bioinformatics/btab625>
53. Wang X, Duan M, Li J *et al.* MarsGT: multi-omics analysis for rare population inference using single-cell graph transformer. *Nat Commun* 2024;**15**:338. <https://doi.org/10.1038/s41467-023-44570-8>
54. Wang S, Zhu C, Jin Y *et al.* A multi-model based on radiogenomics and deep learning techniques associated with histological grade and survival in clear cell renal cell carcinoma. *Insights Imaging* 2023;**14**:207. <https://doi.org/10.1186/s13244-023-01557-9>
55. Saillard C, Delecourt F, Schmauch B *et al.* Pacpaint: a histology-based deep learning model uncovers the extensive intratumor molecular heterogeneity of pancreatic adenocarcinoma. *Nat Commun* 2023;**14**:3459. <https://doi.org/10.1038/s41467-023-39026-y>
56. Quiros AC, Coudray N, Yeaton A *et al.* Mapping the landscape of histomorphological cancer phenotypes using self-supervised learning on unannotated pathology slides. *Nat Commun* 2024;**15**:4596. <https://doi.org/10.1038/s41467-024-48666-7>
57. James R. Egner. AJCC cancer staging manual. *JAMA* 2010;**304**:1726. <https://doi.org/10.1001/jama.2010.1525>
58. Feng Y, Spezia M, Huang S *et al.* Breast cancer development and progression: risk factors, cancer stem cells, signaling pathways,

- genomics, and molecular pathogenesis. *Genes Dis* 2018;**5**:77–106. <https://doi.org/10.1016/j.gendis.2018.05.001>
59. Sun X, Zhang J, Nie Q. Inferring latent temporal progression and regulatory networks from cross-sectional transcriptomic data of cancer samples. *PLoS Comput Biol* 2021;**17**:e1008379. <https://doi.org/10.1371/journal.pcbi.1008379>
 60. Trapnell C, Cacchiarelli D, Grimsby J *et al*. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* 2014;**32**:381–6. <https://doi.org/10.1038/nbt.2859>
 61. Alexander Wolf F, Hamey FK, Plass M *et al*. PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome Biol* 2019;**20**:59.
 62. Song X, Garcia-Saldivar P, Kindred N *et al*. Strengths and challenges of longitudinal non-human primate neuroimaging. *NeuroImage* 2021;**236**:118009. <https://doi.org/10.1016/j.neuroimage.2021.118009>
 63. Sherwani MK, Ruuskanen MO, Feldner-Busztin D *et al*. Multi-omics time-series analysis in microbiome research: a systematic review. *Brief Bioinform* 2025;**26**:bbaf502.
 64. Mora A, Rakar J, Cobeta IM *et al*. Variational autoencoding of gene landscapes during mouse CNS development uncovers layered roles of Polycomb repressor complex 2. *Nucleic Acids Res* 2022;**50**:1280–96. <https://doi.org/10.1093/nar/gkac006>
 65. Maleki N, Padmanabhan B, Dutta K. AI hallucinations: a misnomer worth clarifying. In: *2024 IEEE Conference on Artificial Intelligence (CAI)*, pp. 133–8. Singapore: IEEE, 2024. <https://doi.org/10.1109/CAI59869.2024.00033>
 66. Mora A, Schmidt C, Balderson B *et al*. SiRCle (signature regulatory clustering) model integration reveals mechanisms of phenotype regulation in renal cancer. *Genome Med* 2024;**16**:144. <https://doi.org/10.1186/s13073-024-01415-3>
 67. Arras L, Arjona-Medina J, Widrich M *et al*. Explaining and Interpreting LSTMs. Wojciechand Montavon S, Vedaldi G, Hansen A *et al*. (eds). In: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning* Cham: Springer International Publishing, 2019. 211–38. https://doi.org/10.1007/978-3-030-28954-6_11.
 68. Acharya D, Mukhopadhyay A. A comprehensive review of machine learning techniques for multi-omics data integration: challenges and applications in precision oncology. *Brief Funct Genomics* 2024;**23**:549–60. <https://doi.org/10.1093/bfgp/ela013>
 69. Bailey C, Black JRM, Reading JL *et al*. Tracking cancer evolution through the disease course. *Cancer Discov* 2021;**11**:916–32. <https://doi.org/10.1158/2159-8290.CD-20-1559>
 70. Hynds RE, Huebner A, Pearce DR *et al*. Representation of genomic intratumor heterogeneity in multi-region non-small cell lung cancer patient-derived xenograft models. *Nat Commun* 2024;**15**:4653. <https://doi.org/10.1038/s41467-024-47547-3>
 71. Hrovatin K, Sikkema L, Shitov VA *et al*. Considerations for building and using integrated single-cell atlases. *Nat Methods* 2025;**22**:41–57. <https://doi.org/10.1038/s41592-024-02532-y>
 72. Hernandez M, Epelde G, Alberdi A *et al*. Synthetic data generation for tabular health records: a systematic review. *Neurocomputing* 2022;**493**:28–45. <https://doi.org/10.1016/j.neucom.2022.04.053>
 73. Rohbeck M, De Brouwer E, Bunne C *et al*. Modeling complex system dynamics with flow matching across time and conditions. 2024. <https://openreview.net/forum?id=hwnObmOTrV>
 74. Bunne C, Stark SG, Gut G *et al*. Learning single-cell perturbation responses using neural optimal transport. *Nat Methods* 2023;**20**:1759–68. <https://doi.org/10.1038/s41592-023-01969-x>
 75. Moulaei K, Yadegari A, Baharestani M *et al*. Generative artificial intelligence in healthcare: a scoping review on benefits, challenges and applications. *Int J Med Inform* 2024;**188**:105474. <https://doi.org/10.1016/j.ijmedinf.2024.105474>
 76. Kolla L, Parikh RB. Uses and limitations of artificial intelligence for oncology. *Cancer* 2024;**130**:2101–7. <https://doi.org/10.1002/cncr.35307>
 77. Lundberg SM, Lee S-I. A unified approach to interpreting model predictions. Guyon I, Von Luxburg U, Bengio S *et al*. (eds). In: *Advances in Neural Information Processing Systems*, Vol. **30**. Curran Associates, Inc., 2017. https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf
 78. Ribeiro MT, Singh S, Guestrin C. "Why should I trust you?": explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–44. San Francisco, CA: ACM, 2016. <https://doi.org/10.1145/2939672.2939778>
 79. Wagner AD, Oertelt-Prigione S, Adjei A *et al*. Gender medicine and oncology: report and consensus of an ESMO workshop. *Ann Oncol* 2019;**30**:1914–24. <https://doi.org/10.1093/annonc/mdz414>
 80. White K, Lawrence JA, Tchangalova N *et al*. Socially-assigned race and health: a scoping review with global implications for population health equity. *Int J Equity Health* 2020;**19**:25. <https://doi.org/10.1186/s12939-020-1137-5>
 81. Chua IS, Gaziel-Yablowitz M, Korach ZT *et al*. Artificial intelligence in oncology: path to implementation. *Cancer Med* 2021;**10**:4138–49. <https://doi.org/10.1002/cam4.3935>
 82. Maria Iannantuono G, Bracken-Clarke D, Floudas CS. Applications of large language models in cancer care: current evidence and future perspectives. *Front Oncol* 2023;**13**:1268915. <https://doi.org/10.3389/fonc.2023.1268915>
 83. El Kababji S, Mitsakakis N, Jonker E *et al*. Augmenting insufficiently accruing oncology clinical trials using generative models: validation study. *J Med Internet Res* 2025;**27**:e66821. <https://doi.org/10.2196/66821>
 84. Ghaffarizadeh A, Heiland R, Friedman SH *et al*. PhysiCell: an open source physics-based cell simulator for 3-D multicellular systems. *PLoS Comput Biol* 2018;**14**:e1005991. <https://doi.org/10.1371/journal.pcbi.1005991>
 85. Leon M P-d, Montagud A, Noël V *et al*. PhysiBoSS 2.0: a sustainable integration of stochastic Boolean and agent-based modelling frameworks. *npi Syst Biol Appl* 2023;**9**, 54.