

Pseudo-grading of tumor subclones using phenotype algebra

Namrata Bhattacharya^{1,5,8}, Anja Rockstroh^{1,8}, Sanket Suhas Deshpande⁴, Sam Koshy Thomas¹⁰, Smriti Chawla², Pierre Solomon⁷, Cynthia Fourgeux⁷, Gaurav Ahuja^{4,6}, Brett G. Hollier^{1,8}, Himanshu Kumar³, Antoine Roquilly⁷, Jeremie Poschmann⁷, Melanie Lehman^{1,9}, Colleen C. Nelson^{1,8*}, Debarka Sengupta^{4,5,6*}

1. Australian Prostate Cancer Research Centre-Queensland, Faculty of Health, School of Biomedical Sciences, Centre for Genomics and Personalised Health, Queensland University of Technology, Brisbane, Queensland-4000, Australia
2. Center for Computational Biomedicine, Harvard Medical School, Boston, MA-02115, USA
3. Laboratory of Immunology and Infectious Disease Biology, Department of Biological Sciences, Indian Institute of Science Education and Research (IISER), Bhopal, India
4. Department of Computational Biology, Indraprastha Institute of Information Technology-Delhi (IIIT-Delhi), Okhla, Phase III, New Delhi-110020, India
5. Department of Computer Science and Engineering, Indraprastha Institute of Information Technology-Delhi (IIIT-Delhi), Okhla, Phase III, New Delhi-110020, India
6. Centre for Artificial Intelligence, Indraprastha Institute of Information Technology-Delhi (IIIT-Delhi), Okhla, Phase III, New Delhi-110020, India.
7. Nantes Université, CHU Nantes, INSERM, Center for Research in Transplantation and Translational Immunology, UMR, 1064, Nantes, France
8. Translational Research Institute, Princess Alexandra Hospital, Woolloongabba, Queensland-4102, Australia
9. Vancouver Prostate Centre, Department of Urologic Sciences, University of British Columbia, Vancouver, Canada
10. School of Mathematical Sciences, The University of Adelaide, North Terrace, Adelaide, SA-5005, Australia

*Corresponding authors: {debarka@iiitd.ac.in, colleen.nelson@qut.edu.au}.

Abstract

Robust characterization of cellular phenotypes from single-cell gene expression data is of paramount importance in studying complex biological systems and diseases. Single-cell RNA-sequencing (scRNA-seq), coupled with robust computational analysis, facilitates characterization of phenotypic heterogeneity in tumors. Current scRNA-seq analysis pipelines are capable of accurately identifying a myriad of malignant and non-malignant cell subtypes from single-cell profiling of tumor microenvironments. Unfortunately, given the extent of phenotypic heterogeneity, it is not straightforward to assess the risk associated with individual malignant cell subpopulations in a tumor, primarily due to the complexity of the cancer phenotype space and the lack of clinical annotations associated with tumor scRNA-seq studies, involving prospectively collected tissue samples. Effective risk-stratification of individual malignant subclones holds promise for formulating tailored therapeutic interventions. To this end, we present SCellBOW, a computational approach that facilitates risk-stratification by leveraging scRNA-seq profiles and language modeling techniques. We compared SCellBOW with existing best practice methods for its ability to precisely represent phenotypically divergent cell types across multiple scRNA-seq datasets, including our in-house generated human splenocyte and matched peripheral blood mononuclear cell (PBMC) dataset. SCellBOW offers a remarkable feature for executing algebraic operations such as '+' and '-' on single-cells in the latent space while preserving the biological meanings. This feature catalyzes the simulation of the residual phenotype of tumors, following positive and negative selection of specific malignant cell subtypes in a tumor. As a proof of concept, we tested and validated *phenotype algebra* across three independent cancer types – glioblastoma multiforme, breast cancer and metastatic prostate cancer. In particular, we demonstrate how the negative selection of specific clones may lead to variable prognosis. From the metastatic prostate cancer scRNA-seq data, SCellBOW identifies a hitherto unknown and pervasive AR⁻/NE_{low} (androgen receptor negative, neuroendocrine-low) malignant cell subpopulation with a conspicuously high predictive risk score. We could trace this back in a large-scale spatial omics atlas of 141 well-characterized metastatic prostate cancer samples at the spot resolution.

Keywords - Intra-tumor heterogeneity, single-cell RNA-seq, risk stratification, transfer learning, clustering, phenotype algebra

Introduction

Intra- and inter-tumoral heterogeneity are pervasive in cancer and manifest as a constellation of molecular alterations in tumor tissues. The late-stage clonal proliferation, partial selective sweeps, and spatial segregation within the tumor mass collectively orchestrate lineage plasticity and metastasis¹. In collaboration with non-malignant cell types in the tumor stroma, malignant cells with distinct genetic and phenotypic properties create complex and dynamic ecosystems, rendering the tumors recalcitrant to therapies². Thus, the phenotypic characterization of malignant cell subpopulations is critical for understanding the underlying mechanisms of resistive behavior. The widespread adoption of single-cell RNA-sequencing (scRNA-seq) has enabled the profiling of individual cells, thereby obtaining a high-resolution snapshot of their unique molecular landscapes^{3,4}. A precise understanding of cell-to-cell functional variability captured by scRNA-seq profiles is crucial in this context. These molecular profiles assist robust deconvolution of the oncogenic processes instigated by various selection pressures exerted by anticancer agents and cross-talks between malignant and non-malignant cell types within the tumor microenvironment⁵. To effectively analyze tumor scRNA-seq data, various specialized techniques have been developed to assist in proactive investigation of complex and elusive cell populations^{6,7}, regulatory gene interactions⁸, neoplastic cell lineage trajectories⁹, and expression-based inference of copy number variations¹⁰. While these computational techniques have been successful in gaining novel biological insights; however, their adoption in clinical setups is still elusive.

Over the past years, leading consortia such as The Cancer Genome Atlas (TCGA)¹¹ and other large-scale independent studies have established reproducible molecular subtypes of cancers with divergent prognoses. For instance, metastatic prostate cancer has typically been categorized based on the androgen receptor (AR) or neuroendocrine (NE) signatures: the less aggressive AR+/NE- (AR prostate cancer, ARPC) and the highly aggressive AR-/NE+ (NE prostate cancer, NEPC)¹². More recent studies have identified additional phenotypes, such as the AR-/NE- (double-negative prostate cancer, DNPC) and AR+/NE+ (amphicrine prostate cancer, AMPC)^{13,14}. These findings underscore the importance of considering tumor heterogeneity while dictating the differential treatment regimes to improve patient outcomes¹⁵. These studies are predominantly based on bulk omics assays, which precludes the detectability of fine-grained molecular subtypes of clinical relevance. To address this, there is an urgent need for the development of novel analytical approaches that are capable of exploiting single-cell omics profiles for risk attribution to malignant cell subtypes.

Language modeling approaches are gaining popularity for their ability to deduce powerful latent representations of words by unsupervised ingestion of voluminous textual data. A growing number of natural language processing (NLP) based methods have recently been applied to scRNA-seq data analyses and achieved superior performance. For example, Exceiver¹⁶ and scBERT¹⁷ use attention-based transformers for clustering and cell type annotation, respectively. However, transformer-based gene encoding models are often susceptible to overfitting, particularly if the training datasets are small¹⁸. Additionally, transformers are designed for ordered sequences, whereas genes do not have any inherent ordering in a cell. Although the genes have regulatory dependencies, analytical methods are often not prejudiced on the same. Recent studies have shown that NLP-based tools can be used to predict patients' survival based on their health records^{19,20}. However, this is limited to the processing of electronic health

records. We demonstrate how bag-of-words-based language models can be utilized to decipher molecular heterogeneity of cancers and infer survival risks associated with individual malignant subclones.

In this work, we propose SCellBOW, an unsupervised transfer learning-based embedding technique for single-cell clustering, visualization, survival prediction, and risk stratification. SCellBOW uses a bag-of-words model under the hood, which is independent of any strict ordering of genes²¹. Moreover, the shallow neural network used in SCellBOW is less prone to overfitting on small datasets and enables faster convergence during training²². SCellBOW learned neuronal weights are transferable. This is particularly useful when a limited amount of labeled data is available for the second task. The model can use the knowledge acquired from the first dataset to warm start the learning process on the second dataset. In line with *word algebra* supported by language models, which allows exploration of word analogies (i.e., algebraic operations on word vectors, e.g., word vector associated with *royal* when added with *man*, brings it closer to *king*), cellular embeddings show biologically intuitive outcomes under algebraic operations. We exploited this feature to introduce *phenotype algebra* and applied the same in attributing prognostic value to cancer cell subpopulations identified from scRNA-seq studies. SCellBOW representations, *aka* embeddings, allow algebraic operations such as '+' and '-'. To this end, we applied SCellBOW for pseudo-grading of cancer subclones. Here '*pseudo*' indicates that the risk assessment is computational (i.e. *in silico*). Further, the term '*grading*' should not be misunderstood as the pathological grading of tumors. Rather, it's the relative risk stratification of malignant cell types within a tumor. With the help of *phenotype algebra*, it is possible to simulate the exclusion of the phenotype associated with a specific malignant cell subpopulation and relate the same to the disease outcome. Traditional survival estimation based on bulk omics readouts cannot provide such insights. This information might ultimately therapeutic decision-making and improve patient outcomes. SCellBOW-enabled clustering can help identify rare cell populations with important biological or clinical relevance. We demonstrate the utility of bag-of-words-like modeling in clustering of single-cell expression data by evaluating the method on multiple scRNA-seq datasets ranging from normal prostate to pancreas and PBMCs. As an extended validation, we apply SCellBOW to an in-house scRNA-seq dataset comprising human PBMC and matched splenocytes. We use the *phenotype algebra* component of SCellBOW for pseudo-grading malignant cell subtypes in glioblastoma, breast cancer, and metastatic prostate cancer. With the help of this, we identify a hitherto unknown prostate cancer cell type that is nearly ubiquitous across numerous metastatic prostate cancer samples.

Results

Overview of SCellBOW

Correlating genomic readouts of tumors with clinical parameters has helped us associate molecular signatures with disease prognosis. Given the prominence of phenotypic heterogeneity in tumors, it is important to understand the connection between molecular signatures of clonal populations and disease aggressiveness. Existing methods map tumor samples to a handful of well-characterized molecular subtypes, with known survival patterns. However, bulk RNA-seq measures the average level of gene expression distributed across millions of cells in a tissue sample, thereby obscuring the intra-tumoral heterogeneity²³. Tumor scRNA-seq studies have successfully revealed the extent of gene expression

variance across individual malignant cells, contributing to an in-depth understanding of the mechanisms driving cancer progression²⁴. To date, most studies use marker genes identified from past bulk RNA-seq studies to annotate malignant cell clusters identified from scRNA-seq data. This is inadequate since every tumor is unique when looked through the lenses of single-cell expression profiles²⁵. The proposed approach, SCellBOW, is an adoption of the popular natural language model, Doc2vec²¹, in single-cell representation learning. SCellBOW enables the detection and risk-assessment (*aka.* pseudo-grading) of malignant cell subtypes in tumors (**Fig 1**). Below we highlight key constructions and benefits of this new approach.

SCellBOW generates low-dimensional numeric embeddings associated with individual cells based on their gene expression patterns. SCellBOW represents cells as documents and gene names as words. It learns embeddings from documents, where each document is considered a bag-of-words. The gene expression levels are encoded by word frequencies. This encoding is intuitive and produces meaningful and robust embeddings of cells. SCellBOW uses a neural net language model to preserve the cellular similarities observed in the gene expression space. Although Doc2vec is not meant for transfer learning, our experiments involving the repurposing of a pre-trained, modified Doc2vec model establish promising transfer learning use cases (**Supplementary Note 1**). SCellBOW involves learning in two phases, i.e., the pre-training and the fine-tuning. The pre-training phase involves training a shallow source network that extracts the gene expression patterns from a relatively large scRNA-seq dataset (herein referred to as the source data). The neuronal weights estimated at training are repurposed during the fine-tuning phase based on the relatively small scRNA-seq dataset (herein referred to as the target data), which might represent a rather specialized task. Our results allude to a visible improvement in scRNA-seq analysis outcomes, even with small sample sizes and multiple covariates.

The vocabulary of SCellBOW is composed of words (genes), which are processed to preserve linear semantic relationships. We empirically show the retention of such relationships in the context of cellular phenotypes of cancer cells. SCellBOW uses algebraic operations to compare and analyze cellular phenotypes in a way that is not possible with traditional methods. For example, the '-' operation can be used to predict the likelihood of survival by eliminating the impact of a specific cellular aggressive phenotype from the whole tumor. This could potentially identify the contribution of that phenotype to patient survival. *Phenotype algebra* can be performed either on the pre-defined cancer subtypes or SCellBOW clusters. To explore *phenotype algebra*, three independent datasets are required. Out of the three datasets, two are scRNA-seq datasets used for transfer learning (source and target data), and one bulk RNA-seq dataset with patient-survival information, which is used to train the survival prediction model. Notably, embeddings of transcriptomes from scRNA-seq target data and bulk RNA-seq survival data have been determined concurrently by recalibrating the pre-trained model. The survival model has been trained on the embeddings of the survival data and then tested on the embeddings of the target data to make predictions about the degree of aggressiveness exhibited by the tumor variants.

SCellBOW robustly dissects tissue heterogeneity

Malignant cells are far more heterogeneous as compared to associated normal cells. Clustering is often the first step toward recognizing subclonal lineages in a tumor sample. Clustering single-cells based on their molecular profiles can potentially identify rare cell populations with distinct phenotypes and clinical outcomes. To evaluate the strength of the clustering ability of SCellBOW, we benchmarked our method

against five existing scRNA-seq clustering tools on four non-cancerous datasets²⁶⁻³⁰ (**Supplementary Table 3**). Among these packages, Seurat²⁶ and Scanpy²⁷ are the most popular, and both employ graph-based clustering techniques. DESC²⁸ is a deep neural network-based scRNA-seq clustering package, whereas ItClust²⁹ and scETM³⁰ are transfer learning methods. All the packages are resolution dependent except for ItClust. ItClust automatically selects the resolution with the highest silhouette score. For the objective evaluation of performance, we used an adjusted Rand index (ARI) and normalized mutual information (NMI) for comparing clusters with known cell annotations. For all methods, except for ItClust, we computed overall ARI and NMI for different values of resolution ranging between 0.2 to 2.0. To measure the signal-to-noise ratio of low-dimensional single-cell embeddings, we computed the cell type silhouette index (SI) based on known annotations. We used a number of scRNA-seq datasets to evaluate the tools cited above.

We constructed three use cases leveraging publicly available scRNA-seq datasets. Each instance constitutes a pair of single-cell expression datasets, of which the source data is used for self-supervised model training and the target data for model fine-tuning and analysis of the clustering outcomes. The target data, in all cases, has associated cell type annotations derived from fluorescence-activated cell sorting (FACS) enriched pure cell subpopulations. The first use case consists of non-cancerous cells from prostate cancer patients (120,300 cells)³¹ as the source data and cells from healthy prostate tissues (28,606 cells)³² as the target data. This use case was designed to assess the resiliency of SCellBOW to the presence of disease covariates in a large scRNA-seq dataset. The second use case comprises a large PBMC dataset³³ (68,579 cells) as the source data, whereas the target data was sourced from a relatively small FACS-annotated PBMC dataset (2,700 cells) from the same study. The third use case, as source data, comprises pancreatic cells from three independent studies processed with different single-cell profiling technologies (inDrop³⁴, CEL-Seq2³⁵, SMARTer³⁶). The target data used in this case was from an independent study processed with a different technology Smart-Seq2³⁷ (**Fig. 2, Supplementary Fig. 2**). For the normal prostate, PBMC, and pancreas datasets, SCellBOW produced the highest ARI scores across most notches of the resolution spectrum (0.2 to 2.0) (**Fig. 2g-i**). We observed a similar trend in the case of NMI (**Fig. 2j**). SCellBOW exhibited superior NMI compared to other tools for the normal prostate, PBMC, and pancreas datasets. For further deterministic evaluation of the different methods across the datasets, we set 1.0 as the default resolution for calculating cell type SI (**Fig 2k**). In the PBMC and pancreas datasets, SCellBOW yielded the highest cell type SI for both datasets. SCellBOW and Seurat were comparable in performance for the pancreas dataset, outperforming other tools. We observed poor performance by DESC and ItClust across all the datasets in terms of celltype SI. In terms of cell type SI, Seurat and scETM showed improved results in the normal prostate dataset. However, SCellBOW outperformed both Seurat and scETM in terms of overall ARI and NMI for the same dataset, indicating higher clustering accuracy. To further evaluate the cluster quality in the normal prostate dataset, we compared the known cell types to the predicted clusters. Known cell types such as basal epithelial, luminal epithelial, and smooth muscle cells were grouped into homogeneous clusters by SCellBOW (**Supplementary Fig. 3**). We observed that the majority of fibroblasts and endothelial cells were mapped by SCellBOW to single clusters, unlike Seurat and scETM. SCellBOW retained hillock cells in close proximity to both basal epithelial and club cells, in contrast to Seurat, which only includes basal epithelial cells.

Our analyses on the public datasets confirm the robustness of SCellBOW as compared to the prominent single-cell analysis tools, including the prominent transfer learning methods. To this end, we applied SCellBOW to investigate a more challenging task of analyzing an in-house scRNA-seq data comprising

splenocytes and matched PBMCs from two healthy and two brain dead donors. Given multiple covariates, such as the origin of the cells, and the physiological states of the donors, analysis of this scRNA-seq data presents a challenging use case for single-cell analysis. We used the established high-throughput scRNA-seq platform CITE-seq³⁸ to pool the eight samples into a single experiment. After post-sequencing quality control, we were left with a total of 4,819 cells. We annotated the cells using Azimuth³⁹, with occasional manual interventions (**Supplementary Note 3**). SCellBOW yielded clusters largely coherent with the independently done cell type annotation using Azimuth (**Fig 3**). Further, most clusters harbor cells from all the donors indicating that the sample pooling strategy was effective in reducing batch effects. B cells, T cells, and NK cells map to SCellBOW clusters where the respective cell types are the majority. The majority of CD4+ T cells map to CL0 and CL9 (here CL is used as an abbreviation for cluster) (**Fig 3e**). CL0 is shared between CD8+ T cells, CD4+ T cells, and Treg cells, which originate from the same lineage. CL4 majorly consists of NK cells with a small fraction of CD8+ T cells, which is not unduly deviant from the PBMC lineage tree. As a control, we performed similar analyses of the Scanpy clusters (**Fig 3f**). While Scanpy performed reasonably well, a few misalignments could be spotted. For example, CD4+ and CD8+ T cells were split across many clusters with mixed cell type mappings. SCellBOW maps CD14 monocytes to a single cluster, whereas Scanpy distributes CD14 monocytes across two clusters (CL1 and CL8), wherein CL8 is equally shared with conventional Dendritic Cells (cDC). 'Eryth' annotated cells are indignantly mapped to different clusters by both SCellBOW and Scanpy. This could be due to the Azimuth's reliance on high mitochondrial gene expression levels for annotating erythroid cells³⁹ (**Supplementary Note 3**). However, SCellBOW CL6 is dominated by Erythroid cells with marginal interference from other cell types. In the case of Scanpy, no such cluster could be detected. Further, we quantitatively evaluated SCellBOW and the rest of the benchmarking methods by measuring ARI, NMI, and cell type SI (**Fig. 3g**). While most methods did reasonably well, SCellBOW offered an edge. We observed best results in SCellBOW in terms of ARI, NMI, and cell type SI as compared to other tools. Discerning phenotypic heterogeneity from the expression profiles of seemingly similar cells is a challenging task. The performance of SCellBOW, with near-ground-truth cell type annotation, in such a scenario confirmed its ability to adequately decipher the underlying cellular heterogeneity and provide robust cell type clustering.

SCellBOW facilitates subclonal survival-risk attribution of tumors

Every cancer features unique genotypic as well as phenotypic diversity, impeding the personalized management of the disease. Patient survival is arguably the most important clinical parameter to assess the utility of clinical interventions. The widespread adoption of genomics in cancer care has allowed correlating molecular portraits of tumors with patient survival in all major cancers. This, in turn, has broadened our understanding of the aggressiveness and impact of diverse molecular subtypes associated with a particular cancer type. Single-cell expression profiling of cancers allows the detection and characterization of tumor-specific malignant cell subtypes. This presents an opportunity to gauge the aggressiveness of each cancer cell subtype in a tumor. A few studies suggested that there is an association between tumorigenicity of stromal cells in tumor microenvironments and patient survival⁴⁰. This has sparked debate about the utility of gene expression in the profiling of single tumor cells. Currently, there is no method to estimate the relative aggressiveness of different malignant cell subtypes using survival information as a surrogate. SCellBOW helps achieve this by enabling algebraic operations (subtraction or addition between two expression vectors or phenotypes) in the embedding space. We refer to this as *phenotype algebra*. The ability to perform algebraic operations in cellular embedding vector space is imparted by the use of word

embedding methods as the backbone of SCellBOW. Below is an example from materials science literature⁴¹.

$$\text{antiferromagnetic} - \text{IrMn} + \text{NiFe} \approx \text{ferromagnetic},$$

Here, by subtracting the antiferromagnetic contribution of *IrMn* and adding the ferromagnetic contribution of *NiFe*, the resulting magnetic coupling can approximate a ferromagnetic behavior. We utilized this ability to find an association between the embedding vectors, representing *total tumor – a specific malignant cell cluster* with tumor aggressiveness. This immediately opens up a way to infer the level of aggressiveness of a specific cluster of cancer cells obtained through single-cell clustering. Subtracting an aggressive phenotype from the whole tumor (average of SCellBOW embeddings across all malignant cells in a tumor) would better the odds of survival relative to dropping a subtype under negative selection pressure.

As proof of concept, we first validated our approach on glioblastoma multiforme (GBM), which has been studied widely employing single-cell technologies. GBM has three well-characterized malignant subtypes: proneural (PRO), classical (CLA), and mesenchymal (MES)^{42,43}. We obtained known markers of PRO, CLA, and MES to annotate 4,508 malignant GBM cells obtained from a single patient reported by Couturier and colleagues⁴⁴ and used it as our target data (**Fig. 4b, Supplementary Methods 2.1**). As the source data, we used the GBM scRNA-seq data from Neftel *et al.*⁴⁵. The survival data consisted of 613 bulk GBM samples with paired survival information from the TCGA consortium. We constructed the following query vectors for the survival prediction task: *total tumor* i.e., average of embeddings of all malignant cells, *total tumor – (MES+CLA)*, *total tumor – MES/CLA/PRO* (individually). We conjectured that for the most aggressive malignant cell subtype, *total tumor – subtype specific pseudo-bulk*, would yield the biggest drop in survival risk relative to the *total tumor*. Survival risk predictions associated with *total tumor – MES/CLA/PRO* thus obtained reaffirmed the clinically known aggressiveness order, i.e., CLA > MES > PRO, where CLA succeeds the rest of the subtypes in aggressiveness⁴⁶ (**Fig. 4c, d**). More complex queries can be formulated, such as *total tumor – (MES+CLA)*, which would indicate that the tumor does not comprise the two most aggressive phenotypes, CLA and MES, and instead consists only of PRO cells. This, hypothetically, represents the most favorable scenario, as testified by the *phenotype algebra*^{43,47}.

We performed a similar benchmarking on well-established PAM50-based⁴⁸ breast cancer (BRCA) subtypes: luminal A (LUMA), luminal B (LUMB), HER2-enriched (HER2), basal-like (BASAL), and normal-like (NORMAL) (**Fig. 4a**). We used 24,271 cancer cells from Wu *et al.*⁴⁹ as the source data, 545 single-cell samples from Zhou *et al.*^{50,51} as the target data. We used 1,079 bulk BRCA samples with paired survival information from TCGA as the survival data. SCellBOW-predicted survival risks for the different subtypes were generally in agreement with the clinical grading of the PAM50 subtypes^{52,53}. Exclusion of LUMA from *total tumor* yielded highest risk score indicating that LUMA has the best prognosis, followed by HER2, and LUMB; whereas BASAL and NORMAL were assigned worse prognosis^{54–56} (**Fig. 4e, f**). We observed an interesting misalignment from the general perception about the relative aggressiveness of NORMAL subtype- removal of this subtype from *total tumor* indicated the highest improvement in prognosis. The NORMAL subtype is a poorly characterized and rather heterogeneous category. Recent evidence suggests that these tumors potentially represent an aggressive molecular subtype and are often associated with highly aggressive claudin-low tumors^{52,57}. These benchmarking studies on the well-characterized cancer subtypes of GBM and BRCA affirm SCellBOW's capability to preserve the desirable

characteristics in the resulting phenotypes obtained as outputs of algebraic expressions involving other independent phenotypes and operators such as $+/−$.

Phenotype algebra enables subclonal pseudo-grading of metastatic prostate cancer cells

In prostate cancer, the processes of transdifferentiation and dedifferentiation are vital in metastasis and treatment resistance⁵⁸ (**Fig. 5**). Prostate cancer originates from secretory prostate epithelial cells, where AR, a transcription factor regulated by androgen, plays a key role in driving the differentiation^{59,60}. Androgen-targeted therapies (ATs) constitute the primary treatment options for metastatic prostate cancer, and they are most effective in well-differentiated prostate cancer cells with high AR activity⁶¹. After prolonged treatment with ATs, the cancers eventually progress towards metastatic castration-resistant prostate cancer (mCRPC), which is highly recalcitrant to therapy⁶². In response to more potent ATs, the prostate cancer cells adapt to escape reliance on AR with low AR activity⁶³. The loss of differentiation pressure results in altered states of lineage plasticity in prostate tumors^{64,65}. The most well-defined form of treatment-induced plasticity is neuroendocrine transformation. NEPC is highly aggressive that often manifests with visceral metastases, and currently lacks effective therapeutic options^{64,66}. Recent studies have pointed toward the existence of additional prostate cancer phenotypes that emerge through lineage plasticity and metastasis of malignant cells¹⁴. This includes malignant phenotypes such as low AR signaling and DNPC, which is lacking AR activity and NE features. These additional phenotypes, resulting from the mechanisms of resistance to AR inhibition, can likewise be characterized by distinct gene expression patterns. Presumably, these phenotypes represent an intermediate or transitory state of the progression trajectory from high AR activity to neuroendocrine transdifferentiation⁶⁷.

Here, we performed a pooled analysis of scRNA-seq target data consisting of 836 malignant cells derived from tumors collected from 11 tumors (mCRPC)⁶⁸. We initially classified cells into three categories - AR activity high (ARAH), AR activity low (ARAL), and neuroendocrine (NE). The classification was performed based on the known molecular signatures associated with ARPC and NEPC genes⁶⁷ (**Fig. 5b, Supplementary Methods 2.2**). We used the pre-trained model built on Karthaus *et al.* dataset for transfer learning³¹. We used 81 advanced metastatic prostate cancer patient samples with paired survival information from Abida *et al.*⁶⁹ as the survival data. We subsequently assessed the relative aggressiveness of these high-level categories using *phenotype algebra* (**Fig. 5e, f**). We observed that subtracting the latent signature associated with NE subtype from the tumor led to the largest drop in the predicted risk score, aligning with the anticipated order of survival⁷⁰. In contrast, removing the ARAH subtype from the *total tumor* had a minimal impact on the predicted risk score. Furthermore, subtracting the ARAL subtype resulted in a predicted risk score between ARAH and NE, which supports the hypothesis that ARAL represents an intermediate state between ARAH and NE⁶⁷. Upon further examination of the tumor prognosis by removing the subtypes in a combined state (ARAH+ARAL+NE), the *total tumor* had the highest positive improvement. The survival probability graph also followed the same order.

Grouping prostate cancer cells into three high-level categories is an oversimplified view of the actual heterogeneity of advanced prostate cancer biology. Herein, SCellBOW clusters could be utilized for the discovery of novel subpopulations based on the single-cell expression profiles, that results as a consequence of therapy-induced lineage plasticity. Subsequently, *phenotype algebra* can be used to assign a relative rank

to these clusters under a negative selection pressure based on their aggressiveness. We utilized this concept to cluster the malignant cells from He *et al.* using SCellBOW, resulting into eight clusters and then predicted the relative risk for each cluster (**Fig. 6a, b**). This approach enables a novel and more refined understanding of lineage plasticity states and characteristics that determine aggressiveness during prostate cancer progression. This goes beyond the conventional categorization into ARAH, ARAL, and NE.

To further elucidate these altered cellular programs, we performed a gene set variation analysis (GSVA)⁷¹ based on the AR- and NE- activity (**Fig. 6e, Supplementary Table 2**). Our result showed that CL4 is characterized by the highest expression of NE-associated genes and absence of AR-regulated genes, indicating the conventional NEPC subtype. Despite CL4 having the strongest NEPC signatures, the elimination of CL2 from the tumor conferred an even higher aggressiveness level. Overall, we observed that unlike other clusters, CL2 is composed of cells from the majority of drug-treated patients and multiple metastatic sites (**Fig. 6d**). This highlights that the clustering is not confounded by the individuals or the tissue origin, as often observed during integrative analysis of tumor scRNA-seq data. CL2 rather features a unique gene expression profile common to these cells. Moreover, CL2 has a mixed signature entailing ARAH, ARAL, and NE, indicating the emergence of a more transdifferentiated subtype as a consequence of therapy-induced lineage plasticity (**Fig. 6c**). Even though CL2 shows NE signature, it is distinguished by the gene expression signature induced by the inactivation of the androgen signaling pathway due to ATTs. As a consequence, cells manifesting this novel signature are grouped into a single cluster. Among all clusters, CL1 and CL3 resemble the traditional ARAH subtype. According to our *phenotype algebra* model, the exclusion of CL6 and CL7 from the *total tumor* yielded the highest risk score. Brady *et al.*¹³ have broadly partitioned metastatic prostate cancer into six phenotypic categories using digital spatial profiling (DSP) transcript and protein abundance data in spatially defined metastasis regions. To categorize the SCellBOW clusters into these broad phenotypes, we performed Pearson's correlation test between averaged DSP expression measurements of the six phenotypes and averaged scRNA-seq expression of the SCellBOW clusters (**Fig. 6f**). The results revealed that CL2 has the highest correlation with the phenotype defined by lack of expression of AR signature genes and low or heterogeneous expression of NE-associated genes (AR-/NE_{low}). Similarly, CL4 showed the highest correlation with the NEPC phenotype defined by positive expression NE-associated genes without AR activity (AR-/NE+) as expected. While CL6 exhibits a closer resemblance to a DNPC phenotype (AR-/NE-).

To gain a deeper understanding of these modified cellular processes in CL2 as compared to other clusters, we conducted cluster-wise functional GSVA based on the hallmarks of cancer (**Fig 6g**). We observed that CL2 exhibited the least prostate cancer signature, indicating that this cluster has deviated from prototypical prostate cancer behavior. It has rather dedifferentiated into a more aggressive phenotype as a consequence of therapy-induced lineage plasticity⁷². CL2 exhibited the highest enrichment of genes related to cancer stemness compared to other clusters. CL2 showed pronounced repression of epithelial genes that are downregulated during the epithelial-to-mesenchymal transition (EMT). Furthermore, there is a lack of expression of genes that are upregulated during the reversion of mesenchymal to epithelial phenotype (MET). Thus, the cells in CL2 are undergoing a process of dedifferentiation from being epithelial cells and activating the mesenchymal gene networks. Existing studies have reported that adaptive resistance is positively correlated with the acquisition of mesenchymal traits in cancer^{73,74}. CL2 was enriched with signatures associated with metastasis and drug resistance. Specifically, CL2 showed downregulation of the genes involved in drug metabolism and cellular response while upregulation of the genes associated with drug resistance. In cancer cells, the acquisition of stemness-like as well as cell dedifferentiation

(mesenchymal) traits can facilitate the formation of metastases and lead to the development of drug resistance⁷⁵. Further analysis of the metastatic potential of CL2 indicated that the cluster is enriched with genes associated with vasculature and angiogenesis. Tumors induce angiogenesis in the veins and capillaries of the host tissue to become vascularized, which is crucial for their growth and metastasis⁷⁶. Thus, based on our findings, we anticipate that CL2 corresponds to a highly aggressive and dedifferentiated subclone of mCRPC within the lineage plasticity continuum, correlating to poor patient survival and positive metastatic status. Our cluster-wise *phenotype algebra* results imply that the traits of the androgen and neuroendocrine signaling axes are not the exclusive defining features of the predicted risk ranking and that other yet under-explored biological programs play additional important roles.

Discussion

In this work, we introduce SCellBOW, a distributed bag-of-words-based transfer learning model for single-cell data analysis and *phenotype algebra* to estimate the relative aggressiveness of the cancer cell subtypes. SCellBOW treats a cell as a document, a specific mRNA as a word, and its expression as the word frequency. Traditionally, language models are trained with supervised learning and require a large amount of pre-annotated data to achieve good performance. However, using source data annotations as a reference limits the identification of new cell types in the target data. It is important to realize that exploratory studies involving single-cell expression profiling require unsupervised analysis of the data. SCellBOW allows self-supervised pre-training on gene expression datasets by learning the general syntax and semantic patterns in the unlabeled scRNA-seq dataset. With the SCellBOW transfer learning module, we observe considerable improvement in the quality of clusters obtained from scRNA-seq target datasets^{77,78}. The entire SCellBOW framework is independent of any supervision, i.e., the availability of cell type annotations for any of the source and target datasets is not required. Beyond robust identification of cell type clusters, SCellBOW is the first of its kind scRNA-seq based subclonal tumor aggressiveness assessment tool. SCellBOW uses algebraic operations to analyze the cellular phenotypes that could potentially identify their contribution to patient survival outcomes. We have leveraged the power of word algebra to perform pseudo-grading of cancer subtypes based on their aggressiveness. This involves comparing the likelihood of a subtype being eliminated from the entire tumor under negative selection pressure. In addition to overall survival probability, SCellBOW assigns a risk score to discern the differences between equally aggressive subclones that may be hard to decipher from their survivability profiles. The *phenotype algebra* module exhibits resilience in describing the risk associated with malignant cell subpopulations arising from various types of cancer, which otherwise cannot be accomplished using gene expression data (**Supplementary Fig. 6h-j**).

We compared the clustering capability of SCellBOW to popular single-cell analysis techniques as well as some recently proposed transfer learning approaches. An array of scRNA-seq datasets was used for the same, representing different sizes, cell types, batches, and diseases. As evident from the comparative analyses, SCellBOW exhibits robustness and consistency across all data and metrics. We also reported tangible benefits of transfer learning from apt source data. For example, transfer learning with a large number of tumor-adjacent normal cells from the prostate as the source and healthy prostate cells as the target adequately portrayed the cell type diversity therein. SCellBOW reliably clustered pancreatic islet-specific cells that were processed using varied single-cell technologies. The prevalence of single-cell

expression profiling and the heterogeneity of PBMCs makes it one of the best-studied tissue systems in humans⁷⁹. We isolated matched PBMCs and splenocytes from two healthy donors and two brain-dead donors (the cause of death is subarachnoid hemorrhage on aneurysm rupture). The utilization of this data presents a more intricate problem, where the cells originate from diverse biological and technical replicates, conditions, individuals, and organs. SCellBOW-based analysis of PBMCs, with near-ground-truth cell type annotation, confirmed its ability to decipher the underlying cellular heterogeneity adequately. Notably, SCellBOW, unlike ItClust and scETM, does not limit model training to genes intersecting between the source and the target datasets and is independent of the source data cell type annotation. Given the shallow architecture of the neural network, SCellBOW is faster than other studied transfer learning methods (**Supplementary Table 6**).

As discussed earlier, SCellBOW-based vector representation of cellular transcriptomes preserves their phenotypic relationships in a vector space. This feature can be exploited in innovative ways. We demonstrated a potential use case in estimating subclonal survival risk in three cancer types: GBM, BRCA, and mCRPC. Several examples were shown where simple algebraic operators such as '+' and '-' could derive clinically intuitive outcomes. For example, subtracting the CLA phenotype from the whole tumor resulted in an improvement in the survival risk in a GBM patient compared to MES and PRO subtypes. Upon further probing of the tumor prognosis by removal of the subtypes in combinations from the tumor, specifically CLA and MES subtypes, which are known to be the most aggressive, we observed the highest improvement. To summarize, SCellBOW allows simulations involving multiple phenotypes. Our subsequent investigation focused on a BRCA dataset that harbors a more complex subtype structure. We observed a deviance in our results from the general perception about the relative aggressiveness of the NORMAL subtype in breast cancer. Notably, the removal of the NORMAL subtype from the *total tumor* was associated with the highest improvement in prognosis. Despite indications that these tumors often do not respond to neoadjuvant chemotherapy⁵⁶, the clinical significance of the NORMAL subtype remains uncertain due to a limited number of studies. The NORMAL breast cancer cells are poorly characterized and heterogeneous in nature. As per the recent classification of the breast cancer subtypes, the NORMAL subtype has been identified to be a potentially aggressive molecular subtype, referred to as claudin-low tumors⁵⁷. This contradicts the common assumption that a NORMAL subtype is an artifact resulting from a high proportion of normal cells in the tumor specimen⁸⁰. These evidence suggest that the NORMAL subtype of breast cancer is potentially an aggressive molecular subtype, which is consistent with the prediction made by SCellBOW.

In advanced metastatic prostate cancer, molecular subtypes are still poorly defined, and characterizing novel or more fine-grained subtypes with clinical implications on patient prognosis is an active field of research. As a proof-of-concept, we executed SCellBOW on the three mCRPC subtypes namely ARAH, ARAH, NE. Eliminating the subtypes individually and in combinations (ARAH+ARAL+NE) exhibited clinically intuitive change in prognosis. To date, mCPRC classification has largely been confined to the gradients of AR and NE activities. There is a considerable scope of embracing more fine-grained subtypes to better explain clonal selection and epithelial plasticity in drug resistance in wider use cases. Convinced by the overall performance of SCellBOW, we applied *phenotype algebra* on the mCRPC clusters obtained from He et al. dataset⁶⁸. We observed a misalignment from the general perception that the neuroendocrine subtype features the worst prognosis. Our results pointed towards the existence of a more aggressive and dedifferentiated subclone in the lineage plasticity continuum. This novel subclone shares gene expression signatures with both androgen-repressed and neuroendocrine-activated genes. GSVA analysis of this hitherto unknown phenotype offered insights into its putative functional characteristics, which can be

broadly defined by stemness, EMT, MET, and drug resistance. To summarize, SCellBOW clustering combined with *phenotype algebra* can provide novel insights into facets of prostate cancer biology that can be used to delineate features of lineage plasticity and aid in better classification of molecular phenotypes with clinical significance. Further, such a tool can empower medical oncologists and oncology researchers to develop more personalized and systemic treatment regimes for cancer patients.

Materials and Methods

Description of datasets

To evaluate the performance of SCellBOW, we used fifteen publicly available scRNA-seq datasets and an in-house scRNA-seq dataset spanning different cell types, sizes, and diseases. To benchmark the efficiency of clustering, four use cases were constructed, each involving a pair of single-cell expression datasets. In the first use case, we used ~120,300 non-cancerous prostate cells from Karthaus *et al.*³¹ and ~28,600 healthy prostate cells from Henry *et al.*³². The second use case consisted of two PBMC datasets from Zheng *et al.*³³ containing approximately 68,000 cells and 2,700 cells. For the third use case, we combined three independent batches of pancreatic islet cells from Baron *et al.*³⁴, Muraro *et al.*³⁵, and Wang *et al.*³⁶, with a total of 11,181 cells as the source data and 2,068 cells from Segerstolpe *et al.*³⁷ as the target data. In the fourth use case, we used an in-house CITE-seq dataset isolated from matched spleen and PBMC samples from multiple patients.

To validate the accuracy of *phenotype algebra*, we constructed three use cases from publicly available tumor datasets comprising malignant cells from GBM, BRCA, and mCRPC patient tumors. Each instance involved a pair of single-cell expression datasets and a bulk RNAseq dataset paired with clinical information. For ease of reference, we stick to the following nomenclature – a) the source data comprises the scRNA-seq dataset used for model pre-training, b) The target data comprises malignant cells under investigation, c) The survival data comprises tumor bulk RNA-seq samples. In the case of GBM, we used two scRNA-seq datasets comprising approximately 12,074 cells and 4,508 cells obtained from Neftel *et al.*⁴⁵ and Couturier *et al.*⁴⁴, respectively. In the case of BRCA, we retrieved triple-negative breast cancer scRNA-seq datasets from Wu *et al.*⁴⁹ and Zhou *et al.*⁵¹ with approximately 24,271 cells and 545 malignant cells, respectively. In both the use cases, the bulk RNAseq was obtained from TCGA (<https://portal.gdc.cancer.gov>). In the third use case, we obtained author-annotated malignant cells from He *et al.*⁶⁸. A total of 836 malignant cells were derived from 11 patients and three metastasis sites: bone, lymph node, and liver. We retrieved 81 mCRPC bulk RNA-seq samples with paired survival from Abida *et al.*⁶⁹. We obtained the gene expression of 141 regions of interest determined by digital spatial gene expression profiling of mCRPC from Brady *et al.*¹³. We used this data to correlate the He *et al.* scRNA-seq gene expression with the DSP gene expression of the six phenotypic categories based on the AR- and NE-activity: AR+/NE-, AR_{low}/NE-, AR-/NE-, AR-/NE_{low}, AR+/NE+, and AR-/NE+. Details of the datasets are described in **Supplementary Methods**.

Data preprocessing and quality control

SCellBOW follow the same preprocessing procedure for both the source and target data. SCellBOW first filters the cells with less than 200 genes expressed and genes that are expressed in less than 20 cells. The thresholds may vary depending on the dimension of the dataset (**Supplementary Table 4**). After eliminating low-quality cells and genes, SCellBOW log-normalizes the gene expression data matrix. In the first step, CPM normalization is performed where the expression of each gene in each cell is divided by the total gene expression of the cell, multiplied by 10,000. In the second step, the normalized expression matrix is natural-log transformed after addition of 1 as a pseudo count. Highly variable genes are selected using the `highly_variable_genes()` function from the Scanpy package. SCellBOW performs a z-score scaling on the log-normalized expression matrix with the selected highly variable genes. SCellBOW can handle data in different formats, including UMI count, FPKM, and TPM. The UMI count data follow the same preprocessing procedure as above. SCellBOW skips the normalization step for both TPM and FPKM, since they have been length normalized.

Creating a corpus from source and target data

The gene expression data matrix is analogous to the term-frequency matrix, which represents the frequency of different words in a set of documents. In the context of genomic data, a similar concept can be used to represent the expression level of a gene (words) in a set of cells (documents). Let $E \in R^{G \times C}$ be the gene expression matrix obtained from a scRNA-seq experiment, where each value $E_{g,c}$ of the matrix indicates the expression value of a gene $g \in G$ in a cell $c \in C$ obtained after Scanpy preprocessing. SCellBOW generates embeddings by taking two input datasets- a source data matrix and a target data matrix. The source data set the initial weights of the neural network model, while the target data contains the cells that require clustering or malignancy potential ranking. Before generating the embeddings, SCellBOW performs feature scaling on the data matrix, rescaling each feature to a range of $[0, 10]$ as follows.

$$E'_{g,c} = \text{int} \left(\frac{E_{g,c} - \min}{\max - \min} \right), \quad (1)$$

where scaling is applied to each entry with $\max = 10$ and $\min = 0$. This establishes the equivalence between term-frequency and gene expression matrix. Here we consider that a specific gene is a word, and the expression of the gene in a cell corresponds to the number of copies of the word in a document. To scale the data matrix, we have used the `MinMaxScaler()` function from the scikit-learn package⁸¹. SCellBOW takes only the integer part of the scaled value. This scaling helps control the number of gene names copied in the document. To create the corpus, SCellBOW duplicates the name of the expressed gene in a cell as many times as the gene is expressed. For each cell, SCellBOW shuffles the genes to ensure a uniform distribution of genes across the cell. The resulting gene names are treated as the words in the document. To build the vocabulary from a sequence of cells, a tag number is associated with each cell using the `TaggedDocument()` function in the Gensim package⁸². For each gene in the documents, a token is assigned using a module called `tokenize` with a `word_tokenize()` function in NLTK package⁸³ that splits a gene names into tokens.

The training architecture of the neural network

The scRNA-seq data matrices are high-dimensional. For more effective clustering and algebraic operations, SCellBOW produces a low-dimensional fixed-length embedding of the single-cell transcriptome using transfer learning. The data matrix is transformed from $E' \in Z^{G \times C}$ feature space to $E'' \in R^{d \times C}$ latent feature space of d dimensions, with $d \ll G$. To generate the embeddings, SCellBOW train a Doc2vec distributed memory model of paragraph vectors (PV-DM) language model. The PV-DM model is similar to bag-of-words models in Word2vec⁸⁴. The training corpus in Doc2vec contains a set of documents, each containing a sequence of words $W = \{w_1, w_2, \dots, w_T\}$ that forms a vocabulary V . The words within each document are treated as shared among all documents. The training objective of the model is to maximize the probability of predicting the target word w_t , given the context words that occur within a fixed-size window of size n around w_t in the whole corpus as follows.

$$\text{maximize}(\frac{1}{T} \sum_{t=1}^T \log \Pr(w_t | w_{t-n}, \dots, w_{t+n})). \quad (2)$$

The probability can be modeled using the hierarchical softmax function as follows.

$$\Pr(w_t | w_{t-n}, \dots, w_{t+n}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}, \quad (3)$$

where each of y_i computes the log probability of the word w_t normalized by the sum of the log probabilities of all words in V . This is achieved by adjusting the weights of the hidden layer of a neural network. To build the Doc2vec language model, we used the doc2vec() function in the Gensim library. The initial learning rate was set to 0.025, and the window size to 5. We chose the PV-DM training algorithm and set the embedding vector size to 300 as the default parameter. The choice of embedding vector size may be adjusted according to the size and dimensionality of the dataset.

Fine-tuning using transfer learning with SCellBOW

At first, SCellBOW is pre-trained with a source data matrix $E'_{src} \in Z^{G_{src} \times C_{src}}$ with G_{src} genes and C_{src} cells. The source data matrix sets the initial weights of the neural network model. During transfer learning, SCellBOW fine-tunes the weights learned by the pre-trained model from E'_{src} using a target dataset $E'_{trg} \in Z^{G_{trg} \times C_{trg}}$ with G_{trg} genes and C_{trg} cells. This facilitates faster convergence of the neural network as compared to starting from randomly initialized weights. The output layer of the network is a fixed-length low-dimension embedding (a vector representation) for each cell $c \in C_{trg}$. To infer the latent structure of E'_{trg} single-cell corpus, we used the infer_vector() function in the Gensim package to produce the dimension-reduced vectors for each cell in C_{trg} . This step ensures that the network can map the target data into a low dimensional embedding space R^d ; i.e., $E' \rightarrow E''$, where $E'' \in R^{d \times C}$. The resulting embeddings can be used for various downstream analyses of the single-cell data, such as clustering, visualization of cell types, and *phenotype algebra*.

Data preparation for phenotype algebra

To perform *phenotype algebra*, SCellBOW preprocesses the source data matrix using the standard preprocessing steps. SCellBOW uses an additional bulk RNA-seq gene expression matrix (referred to as survival data). The samples are paired with survival information (e.g., vital status, days to follow-up, days to death). In the survival data, samples without follow-up time or survival status and samples with clinical information but no corresponding RNA-seq data were excluded. SCellBOW accounts for unequal cell distribution across different classes in the target data matrix. SCellBOW upsamples the imbalanced target dataset by generating synthetic samples from the minority class. The value of the synthetic minority class sample is determined by interpolating between its neighboring cells from the same class. We used SMOTE⁸⁵ from the imblearn python library⁸⁶. This confirms that the output of *phenotype algebra* is not confounded by the cell type proportions.

Generating vectors for phenotype algebra

SCellBOW constructs different phenotype vectors from the target dataset. The vectors are constructed either based on the user-defined subtypes or SCellBOW clusters. These vectors are created by calculating the mean gene expression of cells belonging to each phenotype. SCellBOW first constructs a reference vector by taking the average of all malignant cells (referred to as *total tumor*). This accounts for the tumor in its entirety and acts as the reference for the relative ranking of the different groups. Similarly, for each group in the dataset, SCellBOW then constructs a query vector by taking the average of the gene expression of the cells in that individual group. SCellBOW can also perform addition or subtraction operations on the query vectors of two or more cellular phenotypes to evaluate their risk in a combined state such as $T - (a + b)$, where T is the embedding for *total tumor*, a and b are the embeddings for query vectors of two distinct cellular phenotypes. Following this, SCellBOW concatenates the bulk RNA-seq data matrix and the reference and query vectors based on common genes. We used the concatenate() function in the AnnData python package⁸⁷. The resulting combined dataset is then passed to the pre-trained model, which maps it to a lower-dimensional embedding space.

Phenotype algebra in subclonal survival risk attribution

SCellBOW infers the survival probability and predicted risk score for the user-defined tumor subtypes and SCellBOW clusters. The survival risk of the different groups is predicted by fitting a random survival forest (RSF)⁸⁸ machine learning model with SCellBOW embeddings. At first, the RSF is trained with the survival information combined with bulk RNA-seq embeddings obtained from transfer learning. SCellBOW subtracts each query vector from the reference vector and fits the difference of these vectors as input to the RSF model to infer the survival probability $S(t)$. The survival probability computes the probability of occurrence of an event beyond a given time point t as follows.

$$Pr[T < t] = \int_{-\infty}^t f(x)dx, \quad (4)$$

$$S(t) = 1 - Pr[T < t] = Pr[T > t], \quad (5)$$

where T denotes the waiting time until the event occurs and $f(x)$ is the probability density function for the occurrence of an event. SCellBOW computes the survival probability using the `predict_survival_function()` from the scikit-survival RandomSurvivalForest python package⁸⁹ with `n_estimators = 1000`.

In addition to survival probability, SCellBOW can estimate the relative aggressiveness of different malignant phenotypes by assigning a risk score for distinct groups. To infer the predicted risk score, SCellBOW first trains 50 bootstrapped RSF models using 80% of the training set for each iteration. The training data is sub-sampled using different seeds for every iteration. We used the `predict()` function from the scikit-survival package to compute the risk score of each of the input vectors. SCellBOW derives the median of the predicted risk score for each group from the 50 bootstrapped models. The *phenotype algebra* model assigns groups with shorter survival times a lower rank by considering all possible pairs of groups in the data. The groups with a lower predicted risk score after removal from the reference pseudo-bulk are considered more aggressive, as they are associated with shorter survival times.

Data visualization of SCellBOW clusters

SCellBOW maps the cells to low-dimensional vectors in such a way that two cells with similar gene expression patterns will have the least cosine distance between their inferred vectors.

$$\text{similarity}(a, b) = \frac{a \cdot b}{\|a\| \times \|b\|}. \quad (6)$$

After generating low-dimensional embeddings for the cells, SCellBOW identifies the groups of cells with similar gene expression patterns. To identify the clusters, SCellBOW uses the Leiden algorithm⁹⁰ on the embedding matrix of the target dataset. We used the `leiden()` function in the Scanpy package with a default resolution of 1.0. The resolution might vary in the reported results depending on the dimension of the target dataset. To visualize the clusters within a two-dimensional space, SCellBOW uses the `umap()` function from Scanpy library.

Sample collection

In this study, we isolated PBMC and splenocyte from four patients (two healthy donors, HD1 18 years old male and HD2 61 years old female; two brain-dead donors, BD1 57 years old female and BD2 58 years old female). PBMCs were isolated from blood after Ficoll gradient selection. Spleen tissue was mechanically dissociated, then digested with Collagenase and DNase. Splenocytes were finally isolated after Ficoll gradient selection. Splenocytes and PBMCs were stored in DMSO with 10% FBS in liquid nitrogen at the Biological Resource Centre for Biobanking (CHU Nantes, Hotel Dieu, Centre de Ressources Biologiques). This biocollection was authorized in May 2013 by the French Agence de la biomédecine (PFS13-009).

Library preparation and sequencing

For the in-house dataset, we performed CITE-seq using Hashtag Oligos (HTO) to pool samples into a single 10X Genomics channel for scRNA-seq⁹¹. HTO binding was carried out according to the manufacturer's

guidance (BioLegend, Total-seq B). After thawing and cell washing, 1M cells were centrifuged and resuspended in 100uL PSE buffer (PBS/FBS/EDTA). Cells were incubated with 10uL human FcR blocking reagent for 10 min at 4°C. 1 uL of a Hashtag oligonucleotide (HTO) antibody (Biolegend) was then added to each sample and incubated at 4°C for 30 min. Cells were then washed in 1mL PSE, centrifuged at 500g for 5 min, and resuspended in 200uL PSE. Cells were stained with 2uL DAPI, 60uM filtered, and viable cells were sorted (ARIA). Cells were checked for counting and viability, then pooled and counted again. Cell viability was set to < 95%. Cells of HD1 spleen were lost during the washing step and thus not used during further downstream processing. Cells were then loaded on one channel of Chromium Next Controller with a 3' single-cell Next v3 kit (10X Genomics Inc. Pleasanton, CA). We then followed the supplier's protocol CG000185, Rev C, until library generation. For the HTO library, we followed the protocol of the 10X genomics 3' feature barcode kit (PN-1000079) to generate HTO libraries. 310pM of pooled libraries were sequenced on a NOVAseq6000 instrument with an S1 (v1) flow-cell (Illumina Inc. San Diego, CA) at the genomic platform core facility on the basis of a placement agreement (Nantes, IRS-Un, France). The program was run as follows: Read1 29 cycles / 8 cycles (i7) / 0 (i5)/Read2 93 cycles (Standard module, paired-end, two lanes). The FASTQ files were demultiplexed with CellRanger v3.0.1 (10X Genomics) and aligned on the GRCh38 human reference genome. We recovered a total of 6,296 cells with CellRanger, and their gene expression matrices were loaded on R.

Post-sequencing quality control and data preprocessing

We performed the downstream analysis of the in-house CITE-seq dataset in R using Seurat 4.0. For RNA and HTO quantification, we selected cell barcodes detected by both RNA and HTO. We demultiplexed the cells based on their HTO enrichment using the *HTODemux* function in the Seurat R package with default parameters^{26,92}. We subsequently eliminated doublet HTOs (maximum HTO count > 1) and negative HTOs. Singlets were used for further analysis leaving 4,819 cells and 33,538 genes. We annotated each of the cell barcodes using HTO classification as PBMC and Spleen based on the origin of cells and HD1, HD2, BD1, and BD2 based on the patients. In the preprocessing step, we eliminated cells with gene counts < 200 and genes if the number of cells expressing this gene is <3 (**Supplementary Note 2**). Post gene and cell filtering, the gene expression levels for each cell were normalized and log transformed by Scanpy preprocessing methods. The top 3,000 HVGs were used directly for the pipeline of SCellBOW.

We performed automatic cell annotation based on RNA data, using the Seurat-based Azimuth using human PBMC as the reference atlas³⁹. We defined six major cell populations: B cells, CD4 T cells, CD8 T cells, natural killer (NK) cells, monocytes, and dendritic cells. We then manually verified the annotation based on the RNA expression of known marker genes (**Supplementary Fig. 5d and Supplementary Table 7**).

Data availability

The in-house CITE-seq scRNA-seq expression data generated in this study is available in the GEO with accession number GSE221007. Details of the public datasets analyzed in this paper are described in **Supplementary Table 4**.

All source codes are available at GitHub <https://github.com/cellsemantics/SCellBOW>.

Acknowledgments:

DS acknowledges the support of the ihub-Anubhuti-iiitd Foundation set up under the NM-ICPS scheme of the DST. JP, AR, PS, and SF thank the biological resource centre for biobanking (CHU Nantes, Hôtel Dieu, Centre de Ressources Biologiques (CRB), Nantes, F-44093, France (BRIF: BB-0033-00040)) and the Genomics Core Facility GenoA, member of Biogenouest and France Genomique, and to the Bioinformatics Core Facility BiRD, member of Biogenouest and Institut Français de Bioinformatique (IFB) (ANR-11-INBS-0013) for the use of their resources and their technical support. Parts of the schematics were created with BioRender.com.

Author contributions

DS conceived the study with CCN. JP, Antoine Roquilly, PS, and CF contributed to the experimental design of the PBMC-spleen CITE-seq study. NB developed the SCellBOW python package and conducted data analyses with assistance from Anja Rockstroh and SSD. SSD and SKT assisted NB in model training. DS supervised the algorithm development. CCN supervised cancer-related analyses and interpretation. Anja Rockstroh supervised and performed specific data analysis and data interpretation tasks. Anja Rockstroh, HK, JP, GA, ML, and BGH assisted in the biological interpretation of results. DS, CCN, and ML supervised the entire study. All authors substantially contributed to manuscript writing and reviewing.

Conflict of interest

DS and GA are stockholders at CareOnco BioTech. Pvt. Ltd. DS is a stockholder at GenterpretR Inc. The remaining authors declare no competing interests.

Figures

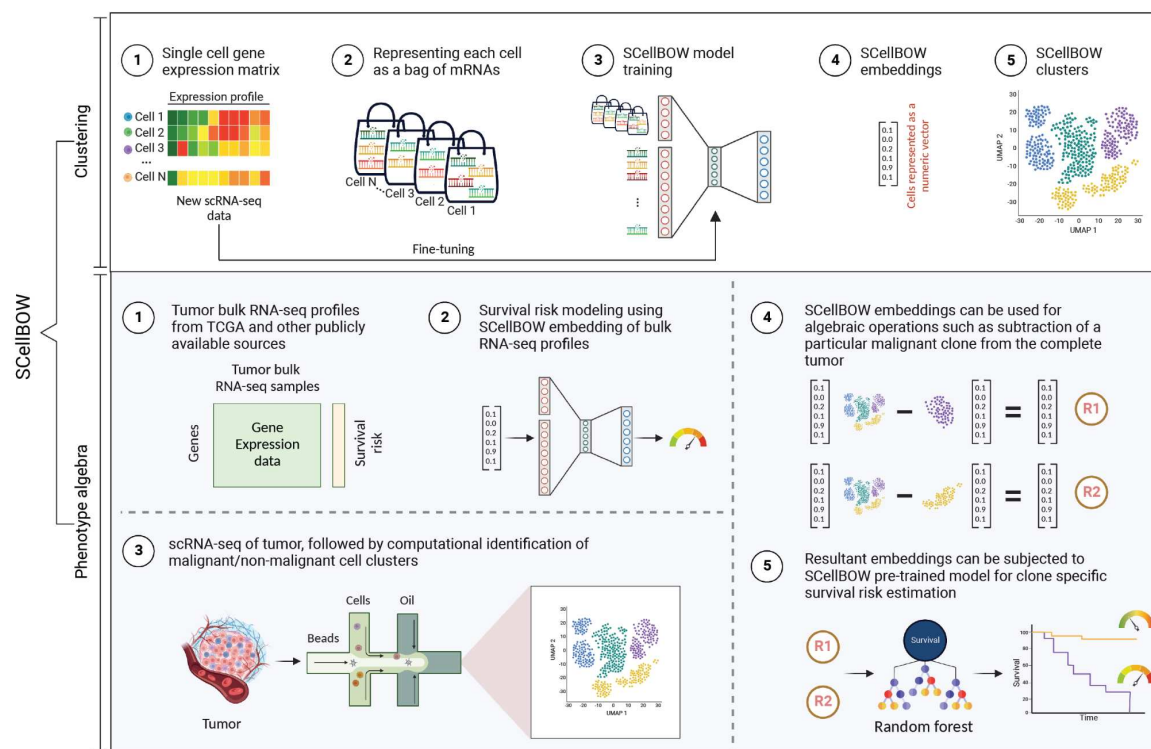


Fig. 1

Fig. 1. SCellBOW workflow. **a**, Schematic overview of SCellBOW workflow for identifying subclones and assessing subclonal tumor aggressiveness. For SCellBOW clustering, firstly a corpus was created from the gene expression matrix, where cells were analogous to documents and genes to words. Next, the pre-trained model was retrained with the vocabulary of the target dataset. Then, clustering was performed on embeddings generated from the neural network. For SCellBOW *phenotype algebra*, vectors were created for reference (*total tumor*) and queries. Then the query vector was subtracted from the reference vector to calculate the predicted risk score using bootstrapped random survival forest. Finally, survival probability was evaluated, and phenotypes were stratified by the median predicted risk score.

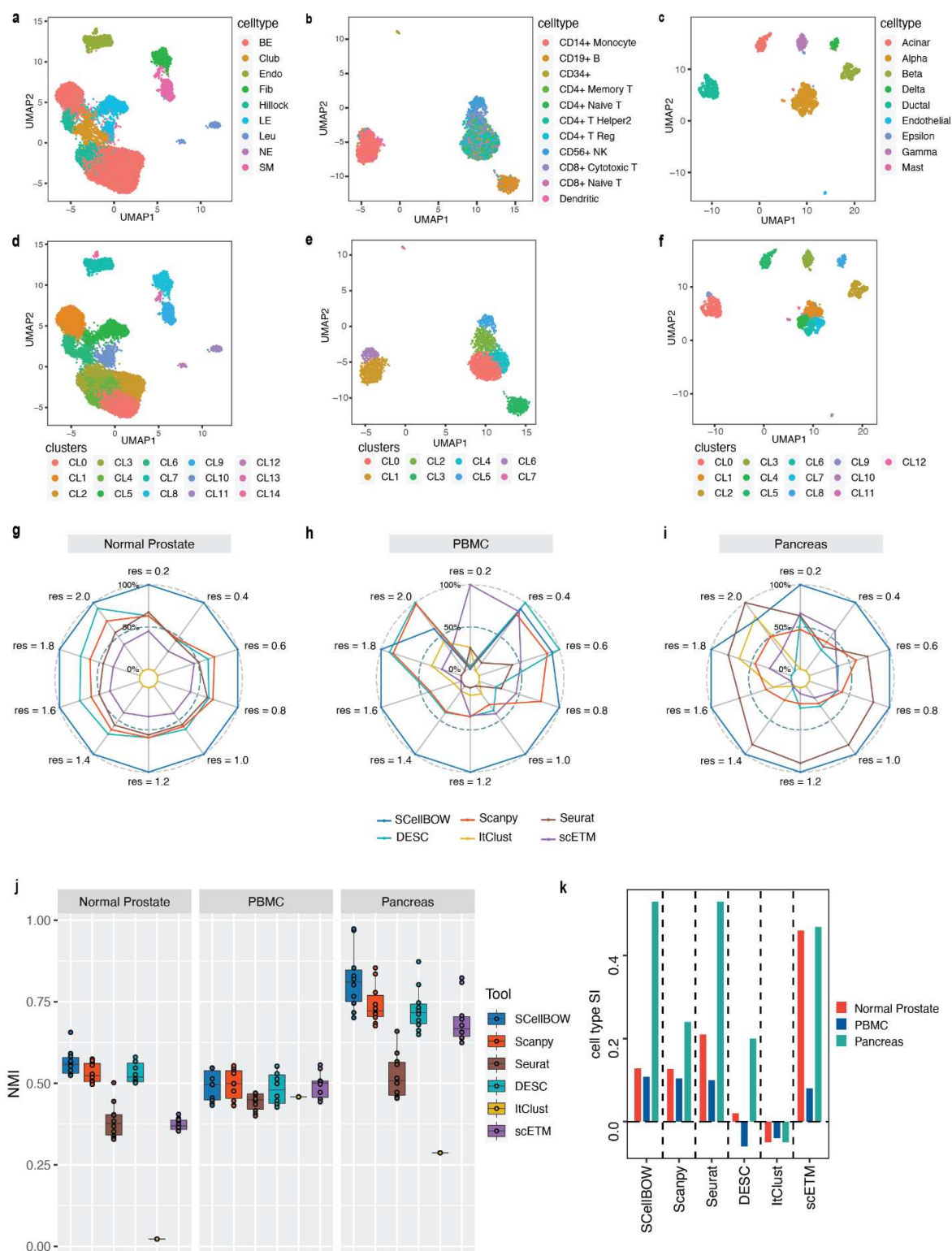


Fig. 2

Fig. 2. Evaluation of single-cell representations using SCellBOW.

a-c, UMAP plots for the normal prostate (a), PBMC (b), and pancreas (c) datasets. The coordinates are colored by cell types.

d-f, UMAP plots for normal prostate (d), PBMC (e), and pancreas (f) datasets, where the coordinates are colored by SCellBOW clusters. CL is used as an abbreviation for cluster.

g-i, Radial plot for the percentage of contribution of different methods towards ARI for various resolutions ranging from 0.2 to 2.0. ItClust is a resolution-independent method; thus the ARI is kept constant across all the resolutions.

j, Box plot for the NMI of different methods across different resolutions ranging from 0.2 to 2.0 in steps of 0.2.

k, Bar plot for the cell type silhouette index (SI) for different methods. The default resolution was set to 1.0.

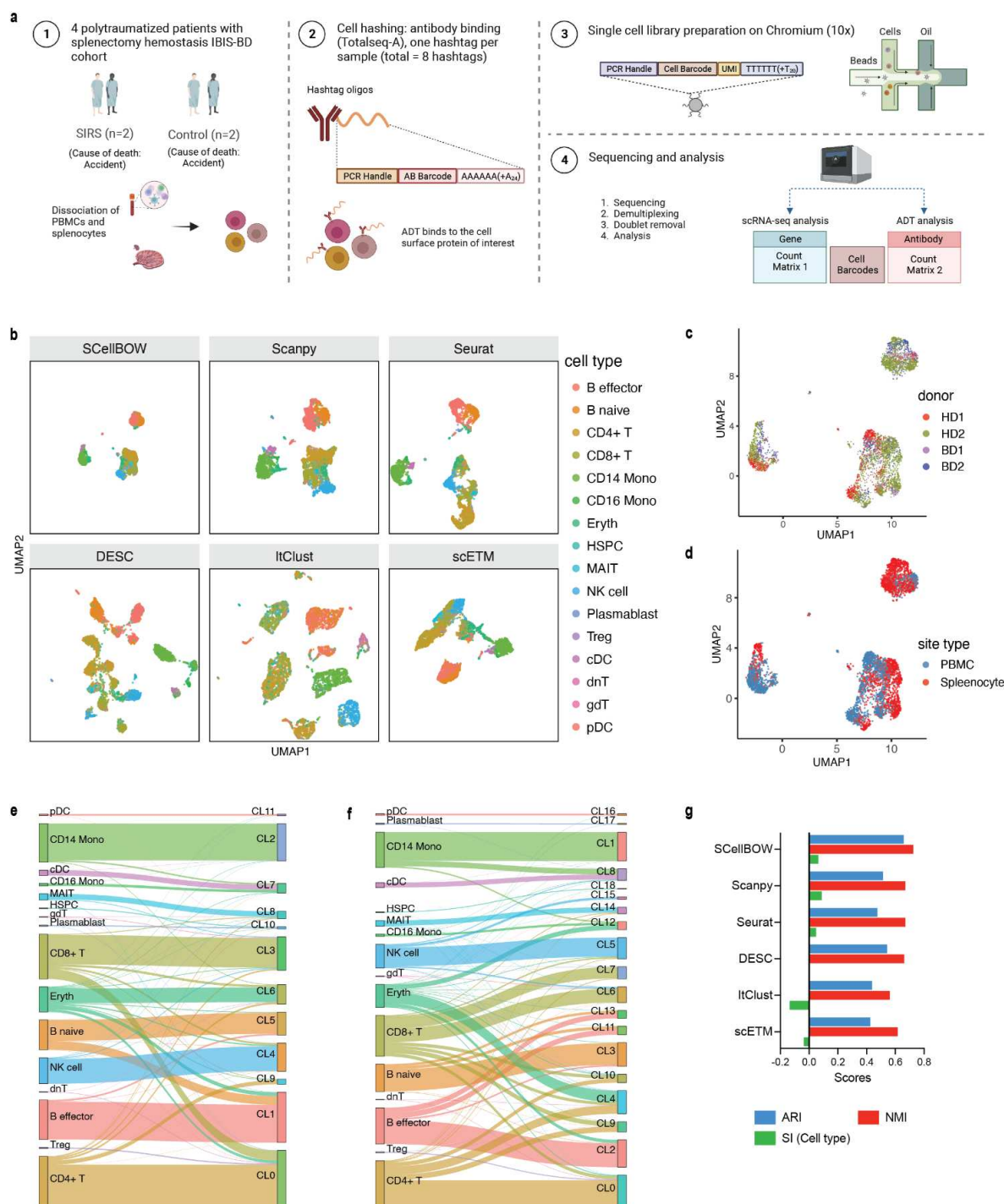


Fig. 3

Figure 3. Evaluation of in-house splenocytes and matched PBMCs.

- a,** An experiment schematic diagram highlighting the sites of the organs for tissue collection and sample processing. In this CITE-seq experiment, PBMCs and splenocytes were collected, followed by high-throughput sequencing and downstream analyses.
- b,** The UMAP plots for the embedding of SCellBOW compared to different benchmarking methods. The coordinates of all the plots are colored by cell type annotation results using Azimuth.
- c-d,** UMAP plots for SCellBOW embedding colored by donors (c) and cell types (d).
- e-f,** Alluvial plots for Azimuth cell types mapped to SCellBOW clusters (e) and Scanpy clusters (f). The resolution of SCellBOW was set to 1.0. CL is used as an abbreviation for cluster.
- g,** Bar plot for ARI, NMI, celltype SI at resolution 1.0.

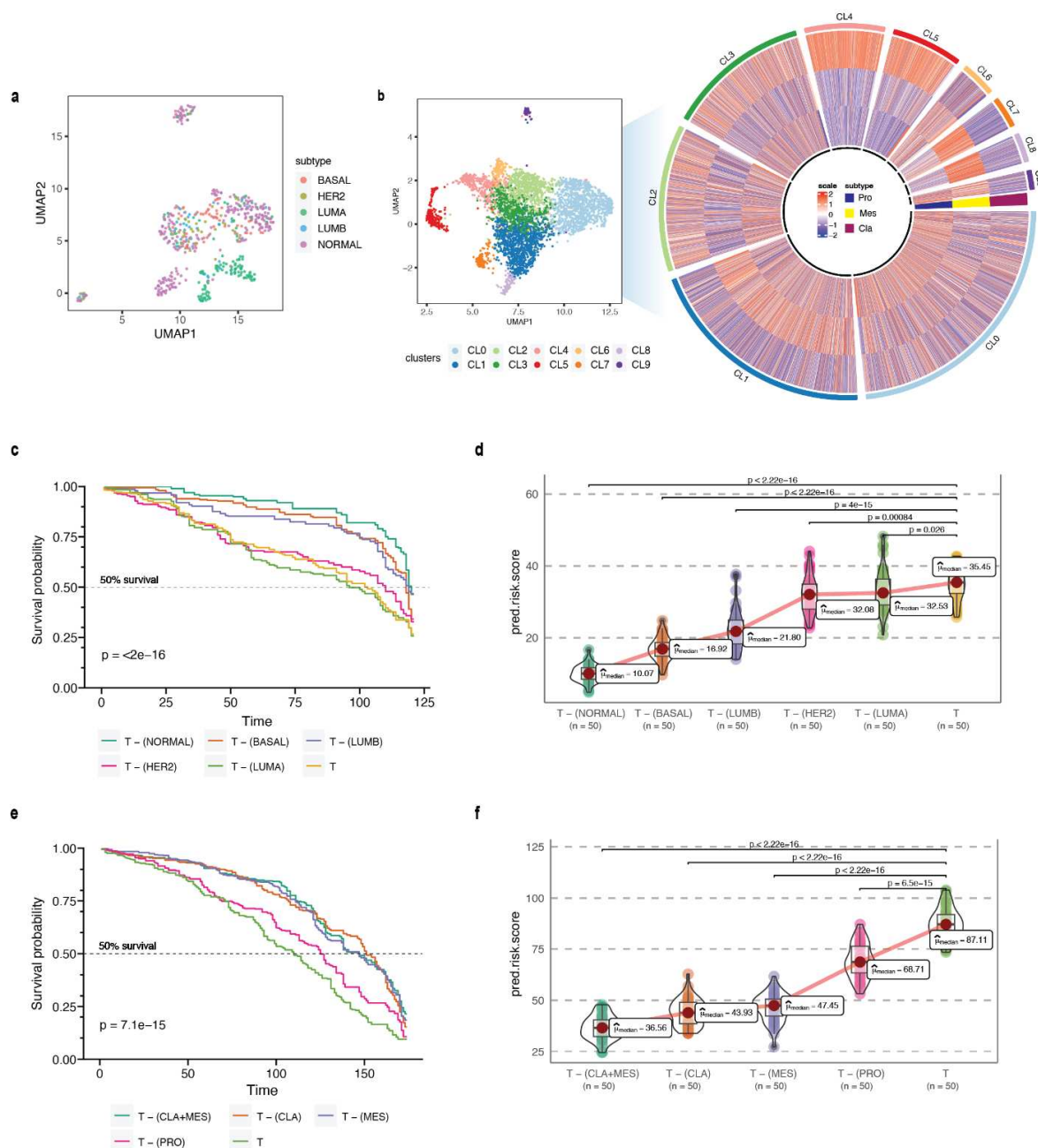


Fig. 4

Fig. 4. Subclonal survival risk inference.

- a**, UMAP plot for the embedding of BRCA target dataset colored by PAM50 molecular subtype.
- b**, Heatmap for GSVA score for three molecular subtypes of GBM: CLA, MES, and PRO, grouped by SCellBOW clusters at resolution 1.0.
- c**, Survival plot for BRCA molecular subtypes based on *phenotype algebra*. The *total tumor* is denoted by T .
- d**, Violin plot for predicted risk scores for BRCA molecular subtypes.
- e**, Survival plot for GBM molecular subtypes based on *phenotype algebra*.
- f**, Violin plot for predicted risk scores for GBM molecular subtypes.

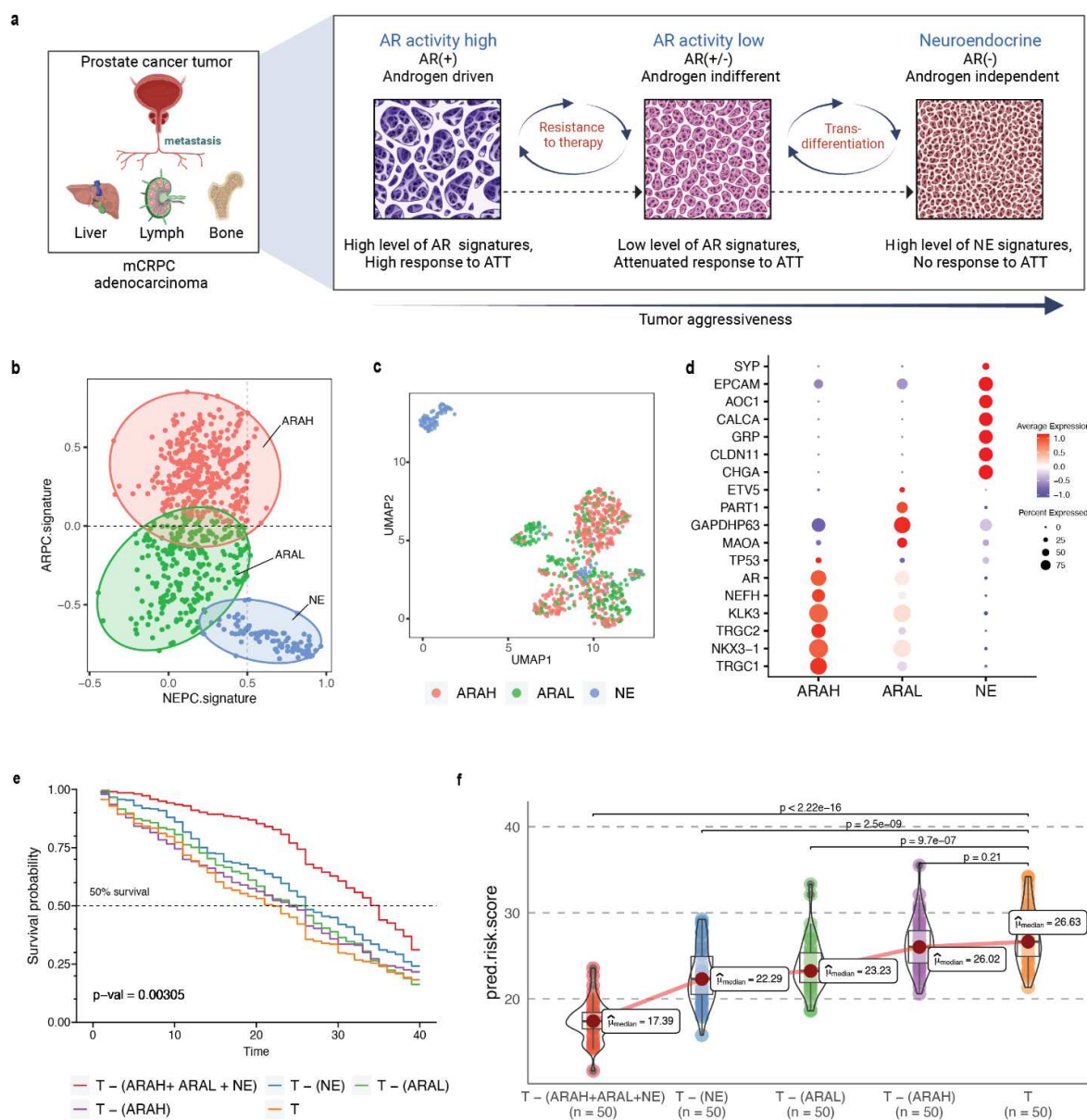


Fig. 5

Fig. 5. Phenotype algebra of metastatic prostate cancer data based on AR- and NE-activity.

- a**, Schematic of the transdifferentiation states underlying lineage plasticity that occurs during mCRPC progression from an ARPC to NEPC.
- b**, Scatter plot of GSVA scores of ARPC and NEPC gene sets, K-means clustering was used to allocate cells into the three high-level ARAH, ARAL, and NEPC categories.
- c**, UMAP plot for projection of SCellBOW embedding colored by ARAH, ARAL, and NEPC.
- d**, Heatmap showing the top differentially expressed genes (y-axis) between each high-level category (x-axis) and all other cells, tested with a Wilcoxon rank-sum test.
- e**, Survival plot for mCRPC phenotypes based on *phenotype algebra*. The *total tumor* is denoted by *T*.
- f**, Violin plot for predicted risk scores for mCRPC phenotypes - ARAH, ARAL, and NEPC.

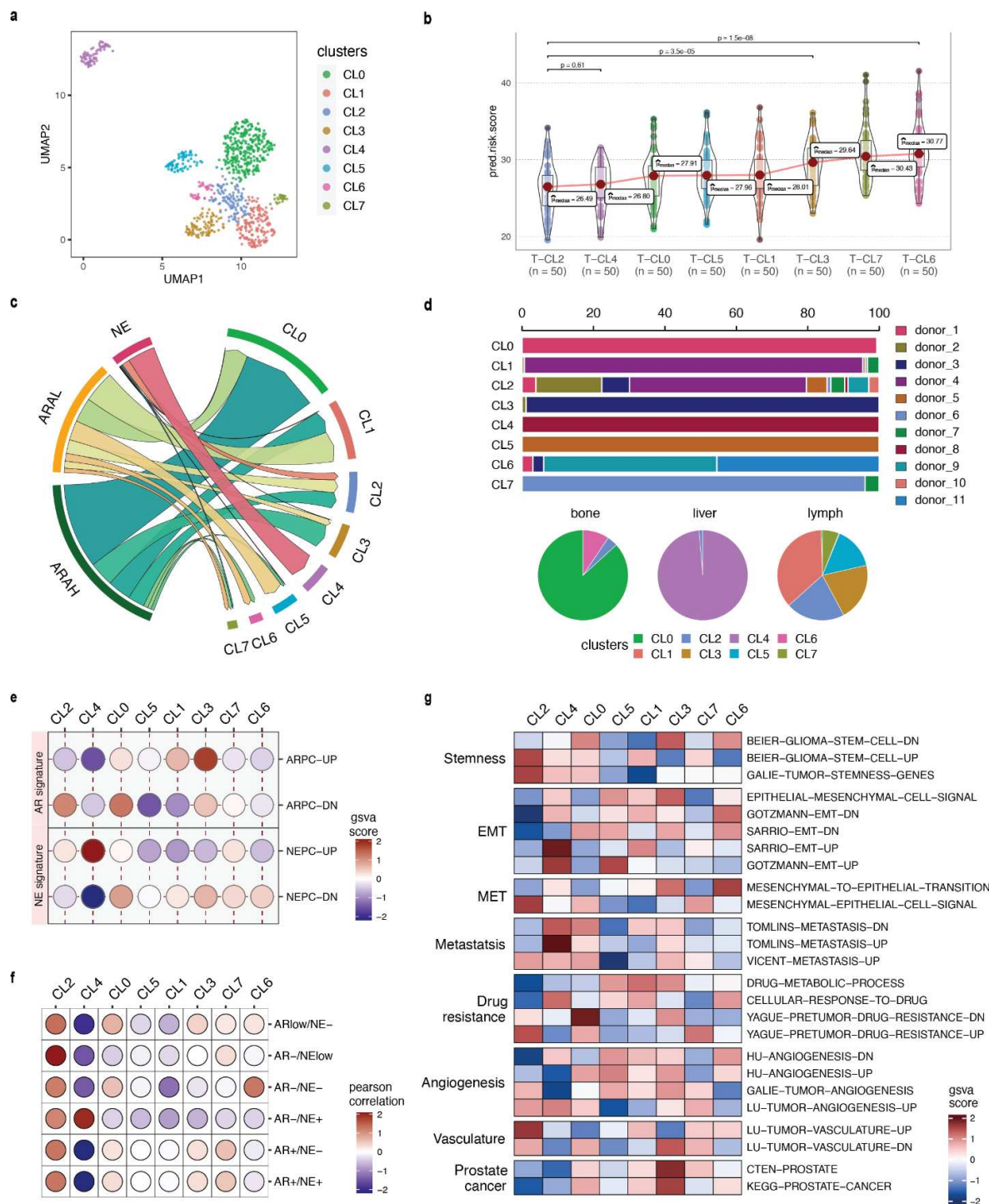


Fig. 6

Fig. 5. Phenotype algebra of metastatic prostate cancer data based on SCellBOW clusters.

- a**, UMAP plot for projection of embeddings with coloring based on the SCellBOW clusters at resolution 0.8. CL is used as an abbreviation for cluster.
- b**, Violin plot of *phenotype algebra*-based cluster-wise risk scores for SCellBOW clusters based on *phenotype algebra*-based predictions.
- c**, Illustration of distribution of cells from the three high-level groups- ARAH, ARAL, and NEPC across the SCellBOW clusters.
- d**, Patient and organ site distribution across the SCellBOW clusters.
- e**, Bubble plot of row-scaled GSVA scores for custom curated gene sets containing activated and repressed AR- and NE- signatures
- f**, Correlation plot of six phenotypic categories based on DSP gene expression correlated with the SCellBOW clusters based on scRNA-seq gene expression. The six phenotypic categories are defined by defined by Brady *et al.* based on the activity of AR and NE programs.
- g**, Top gene sets correlated with SCellBOW clusters. Signatures collected from the C2 “curated,” C5 “Gene Ontology,” and H “hallmark” gene sets from mSigDB⁹³. Ranking by row-scaled GSVA scores of one cluster against all others.

References

1. D'Entropio, S. C. *et al.* Characterizing genetic intra-tumor heterogeneity across 2,658 human cancer genomes. *Cell* **184**, 2239–2254.e39 (2021).
2. Lawson, D. A., Kessenbrock, K., Davis, R. T., Pervolarakis, N. & Werb, Z. Tumour heterogeneity and metastasis at single-cell resolution. *Nat. Cell Biol.* **20**, 1349–1360 (2018).
3. Bhattacharya, N., Nelson, C. C., Ahuja, G. & Sengupta, D. Big data analytics in single-cell transcriptomics: Five grand opportunities. *Wiley Interdiscip. Rev. Data Min. Knowl. Discov.* **11**, (2021).
4. Kanev, K. *et al.* Tailoring the resolution of single-cell RNA sequencing for primary cytotoxic T cells. *Nat. Commun.* **12**, 569 (2021).
5. Dagogo-Jack, I. & Shaw, A. T. Tumour heterogeneity and resistance to cancer therapies. *Nat. Rev. Clin. Oncol.* **15**, 81–94 (2018).
6. Pang, L. *et al.* Discovering Rare Genes Contributing to Cancer Stemness and Invasive Potential by GBM Single-Cell Transcriptional Analysis. *Cancers* **11**, (2019).
7. Poonia, S. *et al.* Marker-free characterization of full-length transcriptomes of single live circulating tumor cells. *Genome Res.* **33**, 80–95 (2023).
8. Chapman, A. R. *et al.* Correlated gene modules uncovered by high-precision single-cell transcriptomics. *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2206938119 (2022).
9. Simeonov, K. P. *et al.* Single-cell lineage tracing of metastatic cancer reveals selection of hybrid EMT states. *Cancer Cell* **39**, 1150–1162.e9 (2021).
10. Tickle, T., Tirosh, I., Georgescu, C., Brown, M. & Haas, B. inferCNV of the Trinity CTAT Project. Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, MA, USA. Preprint at (2019).
11. Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).

12. Beltran, H. *et al.* Divergent clonal evolution of castration-resistant neuroendocrine prostate cancer. *Nat. Med.* **22**, 298–305 (2016).
13. Brady, L. *et al.* Inter- and intra-tumor heterogeneity of metastatic prostate cancer determined by digital spatial gene expression profiling. *Nat. Commun.* **12**, 1426 (2021).
14. Han, H. *et al.* Mesenchymal and stem-like prostate cancer linked to therapy-induced lineage plasticity and metastasis. *Cell Rep.* **39**, 110595 (2022).
15. Chawla, S. *et al.* Gene expression based inference of cancer drug sensitivity. *Nat. Commun.* **13**, 5680 (2022).
16. Connell, W., Khan, U. & Keiser, M. J. A single-cell gene expression language model. *arXiv [q-bio.QM]* (2022).
17. Yang, F. *et al.* scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data. *Nature Machine Intelligence* **4**, 852–866 (2022).
18. Wu, X., Yang, F., Zhou, T. & Lin, X. Rethinking the Impacts of Overfitting and Feature Quality on Small-scale Video Classification. in *Proceedings of the 29th ACM International Conference on Multimedia* 4760–4764 (Association for Computing Machinery, 2021).
19. Nunez, J.-J., Leung, B., Ho, C., Bates, A. T. & Ng, R. T. Predicting the Survival of Patients With Cancer From Their Initial Oncology Consultation Document Using Natural Language Processing. *JAMA Netw Open* **6**, e230813 (2023).
20. Kim, S. *et al.* Deep-Learning-Based Natural Language Processing of Serial Free-Text Radiological Reports for Predicting Rectal Cancer Patient Survival. *Front. Oncol.* **11**, 747250 (2021).
21. Le, Q. V. & Mikolov, T. Distributed Representations of Sentences and Documents. *arXiv [cs.CL]* (2014).
22. Bejani, M. M. & Ghatee, M. A systematic review on overfitting control in shallow and deep neural networks. *Artif. Intell. Rev.* **54**, 6391–6438 (2021).
23. Zhu, S., Qing, T., Zheng, Y., Jin, L. & Shi, L. Advances in single-cell RNA sequencing and its applications in cancer research. *Oncotarget* **8**, 53763–53779 (2017).

24. Baslan, T. & Hicks, J. Unravelling biology and shifting paradigms in cancer with single-cell sequencing. *Nat. Rev. Cancer* **17**, 557–569 (2017).
25. Zhang, Y. *et al.* Precision treatment exploration of breast cancer based on heterogeneity analysis of lncRNAs at the single-cell level. *BMC Cancer* **21**, 918 (2021).
26. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* **36**, 411–420 (2018).
27. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* **19**, 15 (2018).
28. Li, X. *et al.* Deep learning enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nat. Commun.* **11**, 2338 (2020).
29. Hu, J. *et al.* Iterative transfer learning with neural network for clustering and cell type classification in single-cell RNA-seq analysis. *Nat Mach Intell* **2**, 607–618 (2020).
30. Stein-O’Brien, G. L. *et al.* Decomposing Cell Identity for Transfer Learning across Cellular Measurements, Platforms, Tissues, and Species. *Cell Syst* **12**, 203 (2021).
31. Karthaus, W. R. *et al.* Regenerative potential of prostate luminal cells revealed by single-cell analysis. *Science* **368**, 497–505 (2020).
32. Henry, G. H. *et al.* A Cellular Anatomy of the Normal Adult Human Prostate and Prostatic Urethra. *Cell Rep.* **25**, 3530–3542.e5 (2018).
33. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, 1–12 (2017).
34. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell systems* **3**, (2016).
35. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst* **3**, 385–394.e3 (2016).
36. Wang, Y. J. *et al.* Single-Cell Transcriptomics of the Human Endocrine Pancreas. *Diabetes* **65**, 3028–3038 (2016).

37. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2 Diabetes. *Cell Metab.* **24**, (2016).
38. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**, 865–868 (2017).
39. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587.e29 (2021).
40. Dwivedi, B. & Bhasin, M. Survival Genie: A Web Portal for Single-Cell Data, Gene-Ratio, and Cell Composition-Based Survival Analyses. *Blood* **138**, 276 (2021).
41. Tshitoyan, V. *et al.* Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature* **571**, 95–98 (2019).
42. Tang, Y. *et al.* Identification of five important genes to predict glioblastoma subtypes. *Neurooncol Adv* **3**, vdab144 (2021).
43. Behnan, J., Finocchiaro, G. & Hanna, G. The landscape of the mesenchymal signature in brain tumours. *Brain* **142**, 847–866 (2019).
44. Couturier, C. P. *et al.* Single-cell RNA-seq reveals that glioblastoma recapitulates a normal neurodevelopmental hierarchy. *Nat. Commun.* **11**, 3406 (2020).
45. Neftel, C. *et al.* An Integrative Model of Cellular States, Plasticity, and Genetics for Glioblastoma. *Cell* **178**, 835–849.e21 (2019).
46. Lin, N. *et al.* Prevalence and clinicopathologic characteristics of the molecular subtypes in malignant glioma: a multi-institutional analysis of 941 cases. *PLoS One* **9**, e94871 (2014).
47. Patel, A. P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396–1401 (2014).
48. Parker, J. S. *et al.* Supervised risk predictor of breast cancer based on intrinsic subtypes. *J. Clin. Oncol.* **27**, 1160–1167 (2009).
49. Wu, S. Z. *et al.* Stromal cell diversity associated with immune evasion in human triple-negative breast cancer. *The EMBO Journal* vol. 39 Preprint at <https://doi.org/10.15252/embj.2019104063> (2020).

50. Karaayvaz, M. *et al.* Unravelling subclonal heterogeneity and aggressive disease states in TNBC through single-cell RNA-seq. *Nat. Commun.* **9**, 3588 (2018).
51. Zhou, S. *et al.* Single-cell RNA-seq dissects the intratumoral heterogeneity of triple-negative breast cancer based on gene regulatory networks. *Mol. Ther. Nucleic Acids* **23**, 682–690 (2021).
52. Mathews, J. C. *et al.* Robust and interpretable PAM50 reclassification exhibits survival advantage for myoepithelial and immune phenotypes. *NPJ Breast Cancer* **5**, 30 (2019).
53. Weigelt, B. *et al.* Breast cancer molecular profiling with single sample predictors: a retrospective analysis. *Lancet Oncol.* **11**, 339–349 (2010).
54. Hennigs, A. *et al.* Prognosis of breast cancer molecular subtypes in routine clinical care: A large prospective cohort study. *BMC Cancer* **16**, 1–9 (2016).
55. Ahn, H. J., Jung, S. J., Kim, T. H., Oh, M. K. & Yoon, H.-K. Differences in Clinical Outcomes between Luminal A and B Type Breast Cancers according to the St. Gallen Consensus 2013. *Journal of Breast Cancer* vol. 18 149 Preprint at <https://doi.org/10.4048/jbc.2015.18.2.149> (2015).
56. Yersal, O. Biological subtypes of breast cancer: Prognostic and therapeutic implications. *World Journal of Clinical Oncology* vol. 5 412 Preprint at <https://doi.org/10.5306/wjco.v5.i3.412> (2014).
57. Liu, Z., Zhang, X.-S. & Zhang, S. Breast tumor subgroups reveal diverse clinical prognostic power. *Sci. Rep.* **4**, 4002 (2014).
58. Gupta, P. B., Pastushenko, I., Skibinski, A., Blanpain, C. & Kuperwasser, C. Phenotypic Plasticity: Driver of Cancer Initiation, Progression, and Therapy Resistance. *Cell Stem Cell* **24**, 65–78 (2019).
59. Formaggio, N., Rubin, M. A. & Theurillat, J.-P. Loss and revival of androgen receptor signaling in advanced prostate cancer. *Oncogene* **40**, 1205–1216 (2021).
60. Stelloo, S. *et al.* Androgen receptor profiling predicts prostate cancer outcome. *EMBO Mol. Med.* **7**, 1450–1464 (2015).
61. Lonergan, P. E. & Tindall, D. J. Androgen receptor signaling in prostate cancer development and progression. *J. Carcinog.* **10**, 20 (2011).
62. Einstein, D. J. *et al.* Metastatic Castration-Resistant Prostate Cancer Remains Dependent on

- Oncogenic Drivers Found in Primary Tumors. *JCO Precis Oncol* **5**, (2021).
63. Augello, M. A., Den, R. B. & Knudsen, K. E. AR function in promoting metastatic prostate cancer. *Cancer and Metastasis Reviews* vol. 33 399–411 Preprint at <https://doi.org/10.1007/s10555-013-9471-3> (2014).
 64. Antonarakis, E. S. Targeting lineage plasticity in prostate cancer. *The Lancet Oncology* vol. 20 1338–1340 Preprint at [https://doi.org/10.1016/s1470-2045\(19\)30497-8](https://doi.org/10.1016/s1470-2045(19)30497-8) (2019).
 65. Beltran, H. *et al.* The Role of Lineage Plasticity in Prostate Cancer Therapy Resistance. *Clin. Cancer Res.* **25**, 6916–6924 (2019).
 66. Yamada, Y. & Beltran, H. Clinical and Biological Features of Neuroendocrine Prostate Cancer. *Current Oncology Reports* vol. 23 Preprint at <https://doi.org/10.1007/s11912-020-01003-9> (2021).
 67. Labrecque, M. P. *et al.* Molecular profiling stratifies diverse phenotypes of treatment-refractory metastatic castration-resistant prostate cancer. *J. Clin. Invest.* **129**, 4492–4505 (2019).
 68. He, M. X. *et al.* Transcriptional mediators of treatment resistance in lethal prostate cancer. *Nat. Med.* **27**, 426–433 (2021).
 69. Abida, W. *et al.* Genomic correlates of clinical outcome in advanced prostate cancer. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 11428–11436 (2019).
 70. Wang, H. T. *et al.* Neuroendocrine Prostate Cancer (NEPC) Progressing From Conventional Prostatic Adenocarcinoma: Factors Associated With Time to Development of NEPC and Survival From NEPC Diagnosis—A Systematic Review and Pooled Analysis. *Journal of Clinical Oncology* vol. 32 3383–3390 Preprint at <https://doi.org/10.1200/jco.2013.54.3553> (2014).
 71. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* **14**, 7 (2013).
 72. Merkens, L. *et al.* Aggressive variants of prostate cancer: underlying mechanisms of neuroendocrine transdifferentiation. *J. Exp. Clin. Cancer Res.* **41**, 46 (2022).
 73. Hangauer, M. J. *et al.* Drug-tolerant persister cancer cells are vulnerable to GPX4 inhibition. *Nature* **551**, 247–250 (2017).

74. Catapano, J. *et al.* Acquired drug resistance interferes with the susceptibility of prostate cancer cells to metabolic stress. *Cell. Mol. Biol. Lett.* **27**, 100 (2022).
75. Castellón, E. A., Indo, S. & Contreras, H. R. Cancer Stemness/Epithelial–Mesenchymal Transition Axis Influences Metastasis and Castration Resistance in Prostate Cancer: Potential Therapeutic Target. *Int. J. Mol. Sci.* **23**, 14917 (2022).
76. Lugano, R., Ramachandran, M. & Dimberg, A. Tumor angiogenesis: causes, consequences, challenges and opportunities. *Cell. Mol. Life Sci.* **77**, 1745–1770 (2020).
77. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning. *Nat. Biotechnol.* (2021) doi:10.1038/s41587-021-01001-7.
78. Mieth, B. *et al.* Using transfer learning from prior reference knowledge to improve the clustering of single-cell RNA-Seq data. *Sci. Rep.* **9**, 20353 (2019).
79. Phongpreecha, T. *et al.* Single-cell peripheral immunoprofiling of Alzheimer’s and Parkinson’s diseases. *Sci Adv* **6**, (2020).
80. Strehl, J. D., Wachter, D. L., Fasching, P. A., Beckmann, M. W. & Hartmann, A. Invasive Breast Cancer: Recognition of Molecular Subtypes. *Breast Care* **6**, 258–264 (2011).
81. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* **12**, 2825–2830 (2011).
82. Rehurek, R. & Sojka, P. Gensim--python framework for vector space modelling. *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic* **3**, (2011).
83. Bird, S., Klein, E. & Loper, E. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. (‘O’Reilly Media, Inc.’, 2009).
84. Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient Estimation of Word Representations in Vector Space. *arXiv [cs.CL]* (2013).
85. Huang, P. J. *Classification of Imbalanced Data Using Synthetic Over-sampling Techniques*. (2015).
86. Lemaitre, G., Nogueira, F. & Aridas, C. K. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *arXiv [cs.LG]* (2016).

87. Virshup, I., Rybakov, S., Theis, F. J., Angerer, P. & Alexander Wolf, F. anndata: Annotated data. *bioRxiv* 2021.12.16.473007 (2021) doi:10.1101/2021.12.16.473007.
88. Ishwaran, H., Kogalur, U. B., Blackstone, E. H. & Lauer, M. S. Random survival forests. *The Annals of Applied Statistics* vol. 2 Preprint at <https://doi.org/10.1214/08-aos169> (2008).
89. Pölsterl, S. scikit-survival: A Library for Time-to-Event Analysis Built on Top of scikit-learn. *J. Mach. Learn. Res.* **21**, 1–6 (2020).
90. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-connected communities. *Sci. Rep.* **9**, 5233 (2019).
91. Abidi, A. *et al.* Characterization of Rat ILCs Reveals ILC2 as the Dominant Intestinal Subset. *Front. Immunol.* **11**, 255 (2020).
92. Stoeckius, M. *et al.* Cell Hashing with barcoded antibodies enables multiplexing and doublet detection for single cell genomics. *Genome Biol.* **19**, 224 (2018).
93. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**, 417–425 (2015).

