In Search of the Perfect Model: How Cancer Cell Lines Relate to Native Cancers

Rahel Paloots^{1,2}, Ziying Yang^{1,2}, and Michael Baudis^{1,2}

¹Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich, 8057, Switzerland ²Swiss Institute of Bioinformatics (SIB), Winterthurerstrasse 190, Zurich, 8057, Switzerland

Cancer cell lines are frequently used in biological and translational research to study cellular mechanisms and explore treatment options. However, cancer cell lines may display mutational profiles divergent from native cancers or may be misidentified or contaminated. We explored how similar cancer cell lines are to native cancers to find the most suitable representations for the corresponding diseases by utilising large collections of copy number variation (CNV) profiles and applied machine learning (ML) algorithms to predict cell line classifications.

Our results confirm that cancer cell lines indeed accumulate more mutations compared to native cancers but retain similar CNV profiles. We demonstrate that many relevant oncogenes and tumor suppressor genes are altered by CNV events in both cancers and their corresponding cell lines. Based on the similarities between the two groups and the predictions of the ML model, we provide some recommendations about cell lines with good potential to represent selected cancer types in *in vitro* studies.

copy number variants | cancer | cancer cell lines Correspondence: *info@progenetix.org*

Introduction

Derived from cellular samples of primary or recurring tumors, hematopoietic neoplasias, or cancer metastases, cancer cell lines serve as tractable *in vitro* representations of human malignancies, enabling researchers to dissect cancer biology and evaluate potential therapeutic interventions. They are popular models due to established handling procedures and relatively low cost. To establish a cancer cell line, cells are isolated from a native cancer - *e.g.* a solid tumor or bone marrow in case of some lymphoid neoplasms. The ability of many cancer cells to divide perpetually is exploited for propagating them indefinitely for recurrent use in studies, potentially over the course of decades. For example, the first cancer cell line ever established, HeLa (1951) (1), is still one of the most commonly used *in vitro* model systems.

Since the establishment of HeLa and besides individual cell lines established in research projects and maintained at individual institutions, thousands of cell lines have been made available through commercial services, promising precise matching of cell lines to disease types. Many fixed sets of cell lines have been created to model different cancer types in a consistent manner under experimental conditions, e.g. for comparing drug activity or effects of targeted genetic modifications. A widely used collection of cancer cell lines is the NCI-60 panel that represents 9 distinct cancer types (2).

Cancer tissues are complex and genetically heterogeneous,

frequently comprising various clonal populations (3, 4). Genomic aberrations in cancers enable disease progression as well as metastatic growth. A major class of these genomic alterations are copy number variations (CNVs) where large regions in the genomes have been modified by amplification or deletion of a section. These CNV profiles are also cancer type specific, *e.g.* duplication of chromosome 13 in colorectal carcinomas (5–7). Moreover, genomic alterations in cancers depend on the grade as well as the stage of the disease (3), suggesting different subtypes can be detected within disease classification. For example, medulloblastomas have four well known and characterized subtypes. Two of these well-known subsets are based on the expression of Wnt and Shh genes, the other two groups are "group 3" and "group 4" since the biology behind these types is not clear (8).

Due to their restricted origin, cancer cell lines inherently capture only a limited subpopulation of the original tumor's genetic diversity. This inherent limitation is further compounded by the selective pressure exerted by in vitro culture conditions. *In vitro* systems also lack the crucial interactions with other cell types that help shape neoplastic growths. Therefore, under these disparate conditions, cancer cell lines obtain novel features, such as bear a larger amount of copy number alterations compared to primary cancers (9, 10).

In this study, we evaluated cancer cell lines of different types and compared their profiles to their respective native cancer types. We have used statistical and machine learning models to detect CN patterns in both groups. We give an overview of the genomic differences between cancers and their cell lines and unravel the key differences in the genomic features of the two groups. Based on these results, we suggest the best available cell lines per diagnostic group.

Methods

Input Data. Cancer cell line CNV data used in this study originates from cancercelllines.org- a knowledge resource for cancer cell line variants. This database currently includes 5,600 individual CNV samples (11). Progenetix, our source for native cancer samples includes over 140,000 cancer CNV profiles (12). Following an initial analysis of cell line samples, a subset of 32 cancer types was chosen for further investigation. This selection was guided by the distribution of representative samples within the cancer cell lines set.

To reduce the impact of samples with limited quality of their CNV profiling data, all cell line samples were assessed visually and flagged. Native cancer samples without any detected segments were excluded from the dataset automatically. Both native cancer and cell line samples labeled as Unspecified Tissue (NCIT:C132256) were also excluded from the dataset. A set of NCIT cancer types with a sufficient number of cancer cell line samples was used in the similarity assessment and machine learning models (Supplementary Table A1).

Binning of CNV Call Values. To bring cell line and tumor data to a uniform format, CNV data was transformed into a binned matrix. All bins are of equal specified length and represent an area in the genome. Inside each bin, CNV coverages are calculated. Duplication and deletion coverages of the bins are calculated separately. All CNV gains or losses in a bin are counted and the fraction of CNVs covering the bin are calculated. All bins with CNV gain fractions are then saved into an array, followed by the bins with deletion coverage fractions. An open-source python package (bycon) was used for bin calculations. Genomic intervals were then calculated for different bin sizes: 1-10Mb. 5 Mb was selected as bin size for further evaluations. The bin size of 5 Mb was chosen to avoid biases to small and very large segments in the genome.

Like genomic bins, CNV frequency maps can also be created with the bycon package. The frequency maps show the occurrence of a CNV in all samples (%). For example, if all lung carcinoma samples have a duplication in 8q, the frequency in this region would be 100%. These frequencies are calculated for all genome bins (length 1 Mb). Calculated frequencies were then used for similarity assessments and visualizations.

CNV Coverage Calculation. CNV coverage of cell line and tumor samples were retrieved from progenetix and cancercelllines databases. CNV coverage fraction is the amount of the genome that is affected by structural variants (gains and losses). To assess the levels of structural variants between cell lines and tumors of the same cancer types, the average CNV coverage fractions and subsequently fold changes between the two groups were calculated. Pre-calculated CNV coverages for both groups are included in cancercelllines and progenetix databases.

Cosine Similarity. Cosine similarity is a suitable measure for bins with values between 0 and 1 due to its scale invariance, which ensures consistent comparison regardless of vector magnitude. Its angle-based approach captures directional information, making it effective for assessing relative proportions or relationships between values, particularly in sparse datasets. Additionally, its normalized output provides easy interpretation and comparison across different datasets or dimensions. Pairwise cosine similarities between samples were calculated to detect outliers and assess the similarity between instances of the same cell line. Additionally, similarities between different cell lines as well as cell lines and native cancers were determined. Cosine similarity was calculated by using open source python package scikit-learn (Version: 1.2.2). To assess the overall similarity of cell line and native cancer profiles, cosine similarities of CNV frequencies were calculated.

Cosine Similarity
$$(\mathbf{A}, \mathbf{B}) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \cdot \|\mathbf{B}\|}$$

Supervised Machine Learning Algorithms. We trained support vector machine (SVM) and random forest (RF) algorithms on native cancer samples to be able to predict cancer cell line diagnostic classifications based on the profiles of neoplasms. Python package scikit-learn (Version: 1.2.2) SVM classifier with rbf kernel and RF classifiers (n estimators=200) were used for training and prediction. Data matrices with bin sizes from 1-10 Mb were tested and 5 Mb bin size selected for further use. To find the accuracy of differentiating between cancer cell line and native cancer, equal numbers of neoplasia instances to cancer cell lines were picked at random, as the number of neoplasia samples is greater. To predict cell line diagnostic labels based on tumor profiles, machine learning models were trained on selected cancer types. Cancer types were chosen based on the number of cell line samples available per NCIT classification term and child terms of the types were included (Supplementary Table A1). Testing size for both algorithms was 0.25 with random state=10.

Feature Importances. We utilized scikit-learn's PermutationImportance to determine the contributions of each feature to predictions. The script fits a PermutationImportance object to the dataset and computes feature importances using permutation tests. Class feature importances ("Cell Line" and "Tumor") were then calculated for both SVM and RF models. Feature importance results included training with "Cell Line"- "Tumor" of equal numbers to identify general differences between all neoplasias and cell lines. We identified features important for distinct cancer types by training the algorithms on all cell line and neoplasia samples of the same diagnosis. The feature importances of RF model were used for this analysis because of the inherent feature importance measures of the model and the ability to capture complex relationships. Identified features were then mapped to the COSMIC set of signature genes https://cancer.sanger.ac. uk/signatures/downloads/ (last accessed 2023-08-28) to identify relevant genes in these bins.

Matching Cell Lines with Cancer Subtypes. We partitioned cancer data by using the K-means clustering algorithm. Each cancer type (primary cancer samples only) was partitioned into clusters using 2-20 number of clusters. Median sample of each cluster in the partitioned data was calculated and the median sample of each cell line as well to represent each group. Then, cosine similarity between the medians was calculated. Only similarities equal to and above 0.7 were analyzed further.

Plotting and Visualization. All frequency maps and CNV profile plots were created with progenetix/cancercelllines online software tools and bycon package software. Other graphs were created with python packages plotly, seaborn and matplotlib.

Results

Cancer cell lines display limited heterogeneity. A cancer cell line derived from a single source would exhibit minimal genetic variation within its population and within the samples of the same cell line. Even though cancer cell lines are not necessarily monoclonal, they still only represent a small subset of a tumor cell population. To determine the uniformity of the samples of the same cell line, we performed pairwise similarity calculations for all cell lines with at least 3 samples. Then, we compared computed indices to CNV sample plots to ascertain the efficiency of the measure by visual assessment. Figure 1 depicts 2 analyzed cell lines: prostate small cell carcinoma cell line NCI-H660 (CVCL_1576) and lung adenosquamous carcinoma cell line NCI-H596 (CVCL_1571). NCI-H660 is an example of a largely homogenous cell line with all similarity scores above 0.7 and there is a visual congruence among the samples on the CNV plot as well. NCI-H596 on the other hand portrays 2 distinct subsets that also brought about lower scores. Both subsets exhibit high similarity within the group but diverge from each other significantly. These results suggest that our measure can accurately determine the similarity of individual samples.

To disclose the homogeneity of the cell lines of a distinct cancer type, we calculated pairwise similarities of all samples within the same cell line (at least 3 samples) and cancer type (Supplementary Fig. A1). We demonstrate a quantifiable level of homogeneity in cell line samples for the majority of cancer types. The average level of homogeneity across these samples surpasses 0.6. Lower similarity was detected for cervical carcinoma and fibrosarcoma cell lines. These results suggest that samples of the same cell line are overall relatively homogeneous. To detect any outliers, we set the similarity threshold within the same cell line to 0.6. This threshold was set based on visual assessment of single sample plots in combination with calculated similarity indices. Therefore cell lines with at least one paired score below 0.6, would have an outlier. 53% of the cell lines had at least one outlier and only 2% diverged greatly with no similarities above threshold. Overall, the data indicated limited variability between the instances of the same cell line.

Cancer cell lines have higher CNV coverage compared to tumors of the same cancer type. Studies in breast and ovarian carcinoma cell lines have established that compared to tumors, cell lines accumulate more mutations (9, 10). For a systematic assessment of how cancer cell line genomes relate to neoplasias, we compared the CNV coverage fractions of both genomes. As expected, the CNV coverage in most cancer types was higher in cancer cell lines than in neoplasias (Fig.2) with the average fold change of 2.6. The exception was fibrosarcoma where fold change was below 1 and CNV coverage in tumors was greater. In chronic myelogenous leukemia (CML) with BCR-ABL a 15-fold CNV fraction change was detected. This could be due to most cell lines being in the "blast phase" of CML when they acquire a higher mutation load due to increased genomic instability (13). The

link between increased mutation load and blast phase becomes evident when considering the source of CML cell lines used in research, *e.g.* a study by Drexler et al. (2000) where most analyzed CML cell lines originated from patients in blast crisis (14). Interestingly, all three highest fold change values are in different types of leukemia, highlighting the importance of genetic fusion events in many hematopoietic leukemias. Another important fusion in leukemias is PML-RARA in acute myeloid leukemia (AML) but a variety of genes have been reported to be involved in gene fusions in leukemias (15).

Native cancers are highly heterogeneous but exhibit similar CNV profiles to cancer cell lines. It has been reported that despite harboring a higher burden of CNVs compared to their tissues of origin cancer cell lines remarkably retain a similar CNV signature to native cancers (9, 10). To systematically assess the concordance between neoplasia and cell line samples, we calculated pairwise similarity scores based on CNV coverage. Our analysis of the comparisons between cell lines and neoplasias, along with the internal consistency of the neoplasia data set, revealed scores that were lower than what we had anticipated.(Fig.3). Heterogeneity was higher within neoplasia samples than in instances of different cell lines of the same cancer type (Fig. 3, Supplementary Fig. A3). Standard deviations of similarities were higher for neoplasias in all cancer types except acute myeloid leukemia.

Pairwise comparisons of samples illustrate the heterogeneity but fail to reveal recurring patterns of similarity. Therefore, we compared the CNV frequency maps of neoplasias and corresponding cell lines. Frequency maps represent the occurrence of CNVs in the dataset in percentages - indicating the presence of CNVs in the proportion of samples. We calculated a similarity index based on CNV frequency data of cell lines and neoplasias and found high similarity indices for most cancer types (Supplementary Table A2). Clearly detectable patterns emerge between cell lines and its native cancer upon visual inspection of frequency plots as well (Supplementary Fig. A2). Our results confirm that both breast carcinomas and ovarian carcinomas exhibit similar CNV patterns to their cell lines (9, 10). Even though melanoma has bee reported as a highly heterogeneous disease (16, 17), it yielded a high similarity score (Supplementary Table A2, Supplementary Fig. A2) indicating an overall good representation of genomic patterns in aggregated cell line data.

Emerging cancer CNV patterns can be used to determine the origin of some cancer cell lines. Given the high occurrence of characteristic CNV patterns in a cancer type, our goal was to employ a method to forecast the diagnostic classification of a cell line. For that we trained a support vector machine (SVM) model on the 32 cancer types (Supplementary Table A1). In the evaluation of the model's performance we found that for some diagnostic groups the CNV based diagnostic prediction was overall successful with the best results observed for breast adenocarcinoma, glioblastoma and colorectal carcinoma (Supplementary Table A3).



Fig. 1. CNV sample plots and similarity heatmaps for cell lines NCI-H660 (top) and NCI-H598 (bottom), indicating the regional coverage by copy number gains (yellow) and losses (blue) for chromosomes 1-22. For each of the 2 cell lines 5 individual instances are shown to visualize similarities and differences in CNV events. Inter-sample CNV cosine similarities are displayed on the right.



Fig. 2. Comparison of CNV coverage fractions between cell lines and patient derived samples from the same diagnostic groups. Fold change calculated per cancer type: average CNV coverage of cell line samples/average CNV coverage of native cancer samples. Error bars represent standard errors.

We then applied the trained model to predict the diagnostic classifications of our collection of cell line CNV profiles. The

prediction accuracies of the cell lines were largely similar to testing - highest percentage of correctly predicted samples

Acute Myeloid Leukemia	0.26	0.24	0.25	- 0.5
B Acute Lymphoblastic Leukemia	0.21	0.13	0.12	
Bladder Carcinoma	0.24	0.17	0.16	
Breast Adenocarcinoma	0.33	0.18	0.13	
Burkitt Lymphoma	0.35	0.20	0.16	
Cervical Carcinoma	0.36	0.23	0.19	
Chronic Myelogenous Leukemia, BCR-ABL1+	0.30	0.21	0.18	- 0.4
Colorectal Carcinoma	0.33	0.31	0.32	
Diffuse Large B-Cell Lymphoma	0.37	0.22	0.15	
Endometrial Carcinoma	0.26	0.21	0.21	
Esophageal Carcinoma	0.39	0.27	0.20	
Ewing Sarcoma	0.28	0.18	0.12	
Fibrosarcoma	0.43	0.25	0.16	- 0.3
Gastric Carcinoma	0.32	0.27	0.29	
Glioblastoma	0.23	0.18	0.21	
Head and Neck Squamous Cell Carcinoma	0.36	0.23	0.17	
Kidney Carcinoma	0.31	0.23	0.20	
Lung Adenocarcinoma	0.34	0.27	0.27	
Lung Large Cell Carcinoma	0.33	0.25	0.22	- 0.2
Lung Small Cell Carcinoma	0.30	0.20	0.15	
Lung Squamous Cell Carcinoma	0.27	0.18	0.17	
Melanoma	0.19	0.13	0.15	
Neuroblastoma	0.13	0.13	0.19	
Osteosarcoma	0.29	0.17	0.15	
Ovarian Carcinoma	0.23	0.06	0.04	- 0.1
Pancreatobiliary Carcinoma	0.28	0.06	0.05	0.1
Plasma Cell Neoplasm	0.23	0.08	0.06	
Pleural Malignant Mesothelioma	0.13	0.07	0.19	
Prostate Carcinoma	0.14	0.10	0.09	
Soft Tissue Sarcoma	0.23	0.12	0.09	
T Acute Lymphoblastic Leukemia	0.37	0.22	0.19	
	А	В	С	- 0.0

Fig. 3. Similarity heatmap of comparisons between cell lines and neoplasias. Each column shows average pairwise similarity for the cancer type: A - samples of different cell lines, B - cell line samples vs neoplasia samples, C - neoplasia samples.

belonging to breast adenocarcinoma, glioblastoma and colorectal carcinoma. Interestingly, the percentage of correctly predicted cell lines is higher for breast adenocarcinoma cell lines than tumor-tumor predictions, indicating that these cell lines are good representations of the majority of the disease. On the contrary, while the testing accuracy was the highest for acute myeloid leukemia (AML) (90.51%), the prediction accuracy of AML cell lines was below 20%. In fact, the results for all leukemia types in our data sets revealed poor performance. The fold change of CNV coverage was also the highest among leukemias (Fig. 2), suggesting poor representation of the majority of native leukemia samples. **Co-occurrence of relevant cancer genes in features important for class determination.** To interpret the performance of our machine learning model, we employed a feature importance analysis upon its implementation. This analysis identified the most significant regions within the genome that contribute to the model's classification capabilities.

Our initial objective was to examine whether any genomic features could be identified as influential in classifying "neoplasia" or "cell line". Therefore, we collected all cell line samples in our dataset and randomly picked the same number of neoplasia samples. This step ensured the same number of samples to avoid introducing bias to our model. Distinguishing "neoplasia" from "cell line" with 89% accuracy, this model also highlighted the most relevant features in the pro-



Fig. 4. Frequency maps of all neoplasias and cancer cell lines. Blue - deletions, yellow - amplifications. Top 30 detected genes are shown.

cess (Fig. 4).

Next, we analyzed the important features in our selected 32 cancer types, by training our models with both cell line and neoplasia samples and finding the features relevant for each cancer type. Then, we matched identified features (genomic bins) with COSMIC cancer gene set, to identify any underlying oncogenes or tumor suppressor genes in the region. Of all the identified features, 56% are duplications and 44% are deletions. We then sorted the features by including only features with highest average values and that exist in at least 4 cancer types to have an overview of shared relevant features in the cancer types (Fig. 5). 42 out of 50 top shared features are duplications, indicating that duplications carry more weight in class identification. Many of these duplicated features include important genes such as BRCA1, EGFR, MAPK1. Wildtype-BRCA1 is a tumor suppressor gene and mutated BRCA1 increases risk for breast and ovarian carcinoma (18). BRCA1 is detected as an important feature for breast adenocarcinoma but not ovarian carcinoma (Fig. 5). EGFR- epidermal growth factor receptor, is an oncogene that promotes tumor progression. Notably, EGFR over-expression has been detected in lung adenocarcinomas but not in small cell lung carcinomas (SCLC) (19). Duplications in EGFR in lung adenocarcinomas and SCLCs are not marked as important features in our dataset. However, EGFR duplications are important features for determining breast adenocarcinoma sample types and breast cancers are also known to express EGFR (Fig. 5) (19). Another potent oncogene in cancers, particularly bladder carcinomas is FGFR3 (20). Interestingly, FGFR3 deletions are highlighted as important features for class detection in 5 cancer types, including bladder carcinomas (Fig. 5).

Our results suggest that several relevant cancer genes are an integral part for the determination of the class in our model. These genes and features are present in both neoplasias and cancer cell lines but at significantly higher levels in cell lines. These high level CN changes in cancer cell lines make them excellent models for studying the effects of these genes and testing for possible pharmaceuticals.

Modeling Tumor Heterogeneity: The Utility of Cell Lines. Cancers are heterogeneous diseases that consist of multiple subsets. For example, medulloblastoma is a form of brain cancer that includes several distinct subtypes with characteristic mutational profiles (8). Another well-known heterogeneous cancer type is melanoma (17) where a large population of subclones have been detected (21). Based on known intratumoral heterogeneity, we hypothesized that certain cell lines would more accurately represent the unique molecular characteristics of specific cancer subsets. Therefore, we partitioned our neoplasia data using K-means clustering and picking the median sample of each cluster. We then matched these cluster medians to median samples of cell lines. Selecting the median sample allowed us to establish an "average" representation for each group. Indeed, we were able to match some subsets to some cell lines with high similarity, including melanomas and lung small cell carcinomas. For instance, Figure 6 displays a subset of lung small cell carcinoma tumor and cell line samples. We were able to identify three distinct cell lines that would be the best representations of this tumor subpopulation. Additionally, the subsets of melanoma, colorectal carcinoma, glioblastoma and kidney carcinoma matched with high similarity to several cell lines.

Selection of representative cell lines. Our analysis, utilizing a SVM model trained on 32 distinct cancer types and complemented by visual assessment of CNV profiles, has yielded a shortlist of cell lines that serve as strong candidates for faithfully representing primary neoplasia (Table 1). Notably, these cell lines were consistently predicted across all 32 cancer types within the SVM model, suggesting a broader applicability rather than a specific subtype match. Out of all 32 cancer types analysed, we were able to determine candid models for 15 types. The highest number of accurate models



Bladder Carcinoma Breast Adenocarcinoma Burkitt Lymphoma Cervical Carcinoma Colorectal Carcinoma Diffuse Large B-Cell Lymphoma Esophageal Carcinoma Glioblastoma Kidney Carcinoma Lung Adenocarcinoma Lung Small Cell Carcinoma Melanoma Neuroblastoma Ovarian Carcinoma Pancreatobiliary Carcinoma Prostate Carcinoma

Fig. 5. Known oncogenes and tumor suppressor genes important for class determination. Deleted genes are marked with blue and amplified genes with yellow. 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

pgxbs-kftvjjmk, progenetix												
pgxbs-kftvjjmv, progenetix	 1. L		U., 1		1.0.1							
pgxbs-kftvjjmn, progenetix								ш. т				
pgxbs-kftvjjn0, progenetix												
pgxbs-kftvgjf1, progenetix												
pgxbs-kftvkmw0, cancercelllines						.						
pgxbs-kftvkmw1, cancercelllines	 											
pgxbs-kftvkwo9, progenetix												
pgxbs-kftvl3vs, progenetix												
pgxbs-kftvla5k, progenetix						. 📕 🗕						
pgxbs-kftvkmue, cancercelllines												1
pgxbs-kftvkmug, cancercelllines	 											
pgxbs-kftvgix2, progenetix				_			ШШ.,					
pgxbs-kftvkwah, progenetix									• .			 ┇┛──┐╜╜╎╎
pgxbs-kftvkwdm, progenetix									_ 11			
pgxbs-kftvgiwe, progenetix												
pgxbs-kftvgnqz, progenetix												ii'
pgxbs-kftvgjcx, progenetix						 		_				— Ц
pgxbs-kftvgjf0, progenetix												
pgxbs-kftvkwc2, progenetix	 											
pgxbs-kftvkwlq, progenetix												
pgxbs-kftvgjdv, progenetix												
pgxbs-kftvl1d6, progenetix												h
pgxbs-kftvgjr6, progenetix	 											
pgxbs-kftvksxi, cancercelllines												
pgxbs-kftvj3t4, cancercelllines		Jan and and a second					ЦЦ					
pgxbs-kftvksec, cancercelllines												
pgxbs-kftvgjsg, progenetix												
pgxbs-kftvl1aa, progenetix											╜╷╴╸	
pgxbs-kftvgjd5, progenetix												
pgxbs-kftvl59w, progenetix												U
pgxbs-kftvkwtt, progenetix	1							_			_ L	
paxbs-kftvl4m7, progenetix												

Fig. 6. Clustered sample plots of similar cell line and neoplasia subset samples. Cell line samples are indicated in pink and neoplasia samples in black.

were identified for breast adenocarcinomas where prediction accuracy was also the highest. We were also able to pinpoint 3 models for AML despite overall low prediction accuracy for this cancer type.

Discussion

We analyzed the molecular variability of over 500 cell lines (around 3,400 samples) using copy number variation profiles data. Overall, different instances of the same cell line displayed mostly limited variability. Due to the shared origin of these samples, this consistency in prediction is unsurprising. Higher similarity of different cancer cell lines of the same type to each other may be due to cell lines originating from a smaller clonal population and therefore being more stable. This trade-off presents a key challenge: While simplified models may offer stability and ease of analysis, they inherently struggle to capture the full spectrum of disease heterogeneity.

Studies with cancer cell lines have demonstrated that genetic drift occurs in cancer cell lines and significantly affects their pharmacogenetic properties (22, 23). Our results suggest that the heterogeneity of cancer cell lines is limited, especially when compared to native cancers. However, further analyses are required to determine the full scope of the genetic drift in the cell lines. Additionally, as apparent genetic drift in cancer cell lines can be the result of positive selection (22), some level of heterogeneity is inevitable.

The results of this study confirm that cancer cell lines in-

Table 1. Cell lines with good correspondence to native cancer types

Cancer Types	Cell Lines
Acute Myeloid Leukemia	KO52 (CVCL_1321), KG-1 (CVCL_0374), MV4-11 (CVCL_0064)
Breast Adenocarcinoma	EFM-19 (CVCL_0253), ZR-75-1 (CVCL_0588), MDA-MB-134-VI (CVCL_0617), MDA-
	MB-361 (CVCL_0620), CAL-148 (CVCL_1106), CAL-85-1 (CVCL_1114), CAMA-1
	(CVCL_1115), COLO 824 (CVCL_1136), DU4475 (CVCL_1183), Evsa-T (CVCL_1207),
	HCC1008 (CVCL_1244), HCC1143 (CVCL_1245), HCC1428 (CVCL_1252), HCC1500
	(CVCL_1254), HCC1599 (CVCL_1256), HCC2218 (CVCL_1263), (CVCL_1267), HCC38
	(CVCL_1270), MDA-MB-175-VII (CVCL_1400), MFM-223 (CVCL_1408), ZR-75-
	30 (CVCL_1661), UACC-812 (CVCL_1781), UACC-893 (CVCL_1782), EFM-192A
	(CVCL_1812), HDQ-P1 (CVCL_2067), JIMT-1 (CVCL_2077), KPL-1 (CVCL_2094), BT-
	483 (CVCL_2319), VP229 (CVCL_2754), VP267 (CVCL_2755), OCUB-F (CVCL_3352),
	HCC712 (CVCL_3378), SUM190PT (CVCL_3423), SUM44PE (CVCL_3424), MDA-MB-
	468GFP (CVCL_DH83)
Colorectal Carcinoma	COLO 320DM (CVCL_0219), SW403 (CVCL_0545), SW948 (CVCL_0632), LS1034
	(CVCL_1382), NCI-H508 (CVCL_1564), SW1463 (CVCL_1718), SW626 (CVCL_1725),
	COLO 201 (CVCL_1987), COLO 206F (CVCL_1988), WiDr (CVCL_2760)
Diffuse Large B-Cell Lym-	HT (CVCL_1290), OCI-Ly10 (CVCL_8795), OCI-Ly19 (CVCL_1878), Ri-1 (CVCL_1885),
phoma	SU-DHL-4 (CVCL_0539)), SU-DHL-5 (CVCL_1735)
Esophageal Carcinoma	KYSE-140 (CVCL_1347), KYSE-150 (CVCL_1348), KYSE-270 (CVCL_1350), KYSE-30
	(CVCL_1351), KYSE-510 (CVCL_1354), OE21 (CVCL_2661), OE33 (CVCL_04/1)
Glioblastoma	DK-MG (CVCL_1173), GaMG (CVCL_1226), SF295 (CVCL_1690), SNB-75 (CVCL_1706),
	YKG-1 (CVCL_1/96)
Head and Neck Squamous Cell	SCC-25 (CVCL_1682), BICR 22 (CVCL_2310), CAL-27 (CVCL_1107)
Kidney Cansinance	7(0 D (CVCL 1050) A 409 (CVCL 1056) A 704 (CVCL 10(5) A CUN (CVCL 10(7)
Kidney Carcinoma	709-P (CVCL_1050), A-498 (CVCL_1050), A-704 (CVCL_1005), ACHN (CVCL_1007),
	$(CVCL_2741)$ KMPC 1 (CVCL_2082) KMPC 2 (CVCL_2084) SLD26 (CVCL_V612)
Lung Adenocarcinoma	$(CVCL_2741)$, KWRC-1 ($CVCL_2963$), KWRC-2 ($CVCL_2964$), SER20 ($CVCL_2012$)
Lung Adenocaremonia	H1603 (CVCL 1488) NCL H1003 (CVCL 1512) NCL H2030 (CVCL 1517) NCL
	H1095 (CVCL_1488), NCLH295 (CVCL_1512), NCLH2050 (CVCL_1517), NCLH201 H2087 (CVCL_1524) NCLH291 (CVCL_1546) NCLH650 (CVCL_1575) NCLH820
	(CVCL 1592) NCL-H920 (CVCL 1599) RERELC-KL (CVCL 1654)
Lung Small Cell Carcinoma	COR-L 279 (CVCL, 1140)
Lung Squamous Cell Carci-	VMRC-LCP (CVCL, 1788) LC-1/sq (CVCL, 3008)
noma	
Melanoma	C32 (CVCL 1097), C32TG (CVCL 2324), COLO 741 (CVCL 1133), COLO 829
	(CVCL_1137), G-361 (CVCL_1220), Hs 936.T (CVCL_1033), HT-144 (CVCL_0318),
	IGR-37 (CVCL_2075), Ma-Mel-36 (CVCL_A171), Malme-3M (CVCL_1438), SK-MEL-
	1 (CVCL_0068), SK-MEL-103 (CVCL_6069), SK-MEL-19 (CVCL_6025), SK-MEL-
	199 (CVCL_6104), SK-MEL-29 (CVCL_6031), SK-MEL-30 (CVCL_0039), UACC-257
	(CVCL_1779), WM1193C (CVCL_C265), WM164 (CVCL_7928), WM1799 (CVCL_A341),
	WM266-4 (CVCL_2765), WM3060 (CVCL_6796), WM3066 (CVCL_C270), WM3208V
	(CVCL_L028), WM3248 (CVCL_6798), WM51 (CVCL_6995), WM852 (CVCL_6804),
	WM902B (CVCL_6807)
Ovarian Carcinoma	59M (CVCL_2291), COV318 (CVCL_2419), FU-OV-1 (CVCL_2047), OVCAR-4
	(CVCL_1627), PEO1 (CVCL_2686), PEO4 (CVCL_2690), PEO6 (CVCL_2691)
Pancreatobiliary Carcinoma	SU.86.86 (CVCL_3881), Capan-1 (CVCL_0237)
Prostate Carcinoma	NCI-H660 (CVCL_1576), VCaP (CVCL_2235)

deed harbor a higher amount of CNVs compared to native cancers but display CNVs shared with native cancers of the same diagnosis (Supplementary Table A2, Supplementary Fig. A2) (9, 10). The highest differences in CNV coverage fold changes were detected for leukemias, CML in particular. Most cell lines currently in use for CML originate from blast phase CML (13), explaining the large differences in the genomes. Furthermore, this demonstrates that not all stages of cancers are well represented in *in vitro* models.

We examined CNV samples of 32 different cancer types and showed that neoplasias exhibit high molecular variability (Fig. 3, Supplementary Fig. A3). A useful strategy in taking advantage of cell lines would be to partition primary cancer into subsets and match these subsets to cell lines by similarity. For example, we successfully employed this strategy to match lung small cell carcinoma subset to cell lines (Fig. 6). The current model may not fully account for the significant variations within different cancer types. By partitioning of cancer types we can improve the identification of cell lines that closely resemble specific tumor subtypes.

We demonstrated that the prediction of cell line's disease classification based on CNV patterns of native cancers was highly accurate for breast adenocarcinomas, colorectal carcinomas and glioblastomas. Using ML algorithms to classify cancers and cell lines also informs us about genomic features important for classifications. We showed that several wellknown oncogenes and tumor suppressor genes might have influenced the decision-making process (Fig. 5). We also demonstrate that duplications are more important for class determination than deletions.

In summary, despite the undeniable value of cancer cell lines in elucidating tumor biology and propelling advancements in precision medicine, the inherent genomic heterogeneity observed in cancer samples and across individuals needs to be accounted for. We provide a careful selection of the models corresponding best to the target disease to best capture the genomic intricacies of cancers. Data to create neoplasia subsets and match them to appropriate cell lines are available at Progenetix and cancercelllines.org resources. Future steps in utilizing these resources could involve the creation of software tools to enable dynamic comparisons of cancer cell lines to native cancers.

Bibliography

- WF Scherer, JT Syverton, and GO Gey. Studies on the propagation in vitro of poliomyelitis viruses. iv. viral multiplication in a stable strain of human malignant epithelial cells (strain hela) derived from an epidermoid carcinoma of the cervix. J Exp Med, 97(5):695–710, 1953.
- Robert H Shoemaker. The nci60 human tumour cell line anticancer drug screen. Nature Reviews Cancer, 6(10):813–823, 2006.
- Timothy N Clinton, Ziyu Chen, Hannah Wise, Andrew T Lenis, Shweta Chavan, Mark TA Donoghue, Nima Almassi, Carissa E Chu, Shawn Dason, Pavitra Rao, et al. Genomic heterogeneity as a barrier to precision oncology in urothelial cancer. *Cell reports*, 41(12), 2022.
- Roland F Schwarz, Charlotte KY Ng, Susanna L Cooke, Scott Newman, Jillian Temple, Anna M Piskorz, Davina Gale, Karen Sayal, Muhammed Murtaza, Peter J Baldwin, et al. Spatial and temporal heterogeneity in high-grade serous ovarian cancer: a phylogenetic analysis. *PLoS medicine*, 12(2):e1001789, 2015.
- Kaja CG Berg, Peter W Eide, Ina A Eilertsen, Bjarne Johannessen, Jarle Bruun, Stine A Danielsen, Merete Bjørnslett, Leonardo A Meza-Zepeda, Mette Eknæs, Guro E Lind, et al. Multi-omics of 34 colorectal cancer cell lines-a resource for biomedical studies. *Molecular cancer*, 16:1–16, 2017.
- 6. Jordi Camps, Marian Grade, Quang Tri Nguyen, Patrick Hörmann, Sandra Becker,

Amanda B Hummon, Virginia Rodriguez, Settara Chandrasekharappa, Yidong Chen, Michael J Difilippantonio, et al. Chromosomal breakpoints in primary colon cancer cluster at sites of structural variants in the genome. *Cancer research*. 68(5):1284–1295, 2008.

- Eleanor J Douglas, Heike Fiegler, Andrew Rowan, Sarah Halford, David C Bicknell, Walter Bodmer, Ian PM Tomlinson, and Nigel P Carter. Array comparative genomic hybridization analysis of colorectal cancer cell lines and primary carcinomas. *Cancer research*, 64(14): 4817–4825, 2004.
- Michael D Taylor, Paul A Northcott, Andrey Korshunov, Marc Remke, Yoon-Jae Cho, Steven C Clifford, Charles G Eberhart, D Williams Parsons, Stefan Rutkowski, Amar Gaijar, et al. Molecular subgroups of medulloblastoma: the current consensus. *Acta neuropathologica*, 123:465–472, 2012.
- Richard M Neve, Koei Chin, Jane Fridlyand, Jennifer Yeh, Frederick L Baehner, Tea Fevr, Laura Clark, Nora Bayani, Jean-Philippe Coppe, Frances Tong, et al. A collection of breast cancer cell lines for the study of functionally distinct cancer subtypes. *Cancer cell*, 10(6): 515–527, 2006.
- Silvia Domcke, Rileen Sinha, Douglas A Levine, Chris Sander, and Nikolaus Schultz. Evaluating cell lines as tumour models by comparison of genomic profiles. *Nature communications*, 4(1):2126, 2013.
- Rahel Paloots and Michael Baudis. cancercelllines. org—a novel resource for genomic variants in cancer cell lines. *Database*, 2024:baae030, 2024.
- Qingyao Huang, Paula Carrio-Cordo, Bo Gao, Rahel Paloots, and Michael Baudis. The progenetix oncogenomic resource in 2021. *Database*, 2021, 2021.
- Jayastu Senapati, Elias Jabbour, Hagop Kantarjian, and Nicholas J Short. Pathogenesis and management of accelerated and blast phases of chronic myeloid leukemia. *Leukemia*, 37(1):5–17, 2023.
- Hans G Drexler, Yoshinobu Matsuo, and Roderick AF MacLeod. Continuous hematopoietic cell lines as model systems for leukemia–lymphoma research. *Leukemia research*, 24(11): 881–911, 2000.
- Abdul Mannan, Ibrahim N Muhsen, Eva Barragán, Miguel A Sanz, Mohamad Mohty, Shahrukh K Hashmi, and Mahmoud Aljurf. Genotypic and phenotypic characteristics of acute promyelocytic leukemia translocation variants. *Hematology/Oncology and Stem Cell Therapy*, 13(4):189–201, 2020.
- Noemi Andor, Trevor A Graham, Marnix Jansen, Li C Xia, C Athena Aktipis, Claudia Petritsch, Hanlee P Ji, and Carlo C Maley. Pan-cancer analysis of the extent and consequences of intratumor heterogeneity. *Nature medicine*, 22(1):105–113, 2016.
- Tomasz M Grzywa, Wiktor Paskal, and Paweł K Włodarski. Intratumor and intertumor heterogeneity in melanoma. *Translational oncology*, 10(6):956–975, 2017.
- Amy X Zhong, Yumay Chen, and Phang-Lang Chen. Brca1 the versatile defender: Molecular to environmental perspectives. *International Journal of Molecular Sciences*, 24(18): 14276, 2023.
- Elena Levantini, Giorgia Maroni, Marzia Del Re, and Daniel G Tenen. Egfr signaling pathway as therapeutic target in human cancers. In *Seminars in Cancer Biology*, volume 85, pages 253–275. Elsevier, 2022.
- Claudia Maria Ascione, Fabiana Napolitano, Daniela Esposito, Alberto Servetto, Stefania Belli, Antonio Santaniello, Sarah Scagliarini, Felice Crocetto, Roberto Bianco, and Luigi Formisano. Role of fgfr3 in bladder cancer: Treatment landscape and future challenges. *Cancer Treatment Reviews*, page 102530, 2023.
- Li Ding, Minjung Kim, Krishna L Kanchi, Nathan D Dees, Charles Lu, Malachi Griffith, David Fenstermacher, Hyeran Sung, Christopher A Miller, Brian Goetz, et al. Clonal architectures and driver mutations in metastatic melanomas. *PloS one*, 9(11):e111153, 2014.
- Uri Ben-David, Benjamin Siranosian, Gavin Ha, Helen Tang, Yaara Oren, Kunihiko Hinohara, Craig A Strathdee, Joshua Dempster, Nicholas J Lyons, Robert Burns, et al. Genetic and transcriptional evolution alters cancer cell line drug response. *Nature*, 560(7718):325– 330, 2018.
- Rene Quevedo, Petr Smirnov, Denis Tkachuk, Chantal Ho, Nehme El-Hachem, Zhaleh Safikhani, Trevor J Pugh, and Benjamin Haibe-Kains. Assessment of genetic drift in large pharmacogenomic studies. *Cell systems*, 11(4):393–401, 2020.