

Cell Painting predicts impact of lung cancer variants

Juan C. Caicedo, John Arevalo, Federica Piccioni, Mark-Anthony Bray, Cathy L. Hartland, Xiaoyun Wu, Angela N. Brooks, Alice H. Berger, Jesse S. Boehm, Anne E. Carpenter*, and Shantanu Singh*

Broad Institute of Harvard and MIT, Cambridge, MA 02142

ABSTRACT Most variants in most genes across most organisms have an unknown impact on the function of the corresponding gene. This gap in knowledge is especially acute in cancer, where clinical sequencing of tumors now routinely reveals patient-specific variants whose functional impact on the corresponding genes is unknown, impeding clinical utility. Transcriptional profiling was able to systematically distinguish these variants of unknown significance as impactful vs. neutral in an approach called expression-based variant-impact phenotyping. We profiled a set of lung adenocarcinoma-associated somatic variants using Cell Painting, a morphological profiling assay that captures features of cells based on microscopy using six stains of cell and organelle components. Using deep-learning-extracted features from each cell's image, we found that cell morphological profiling (cmVIP) can predict variants' functional impact and, particularly at the single-cell level, reveals biological insights into variants that can be explored at our public online portal. Given its low cost, convenient implementation, and single-cell resolution, cmVIP profiling therefore seems promising as an avenue for using non-gene specific assays to systematically assess the impact of variants, including disease-associated alleles, on gene function.

Monitoring Editor

Jennifer Lippincott-Schwartz
Howard Hughes Medical
Institute

Received: Nov 2, 2021

Revised: Jan 26, 2022

Accepted: Mar 22, 2022

INTRODUCTION

Lung cancer is the leading cause of cancer-related mortality and presents high mutation rates (Lawrence *et al.*, 2013; Bray *et al.*, 2018). New variants are found every year in clinical studies, most of them variants of unknown significance (VUS). Although custom-tailored assays might be created to assess the function of each gene in the presence or absence of each variant, this is exceptionally time-consuming. It is only practical for a small number of known onco-

genes and tumor suppressors and is impossible for genes whose function is unknown. This limits the expansion of precision medicine, where cancer patients are tested to identify their specific mutations and ultimately receive targeted treatments.

High-dimensional profiling assays have been proposed as an accelerant for determining the significance of VUS: by measuring many phenotypic properties of cells exposed to each variant in each gene of interest, the strategy is to capture many genes' functions in a single assay and therefore assess many variants' impact. This strategy was successfully demonstrated using high-throughput transcriptional profiling in an approach called expression-based variant impact phenotyping (eVIP; Berger *et al.*, 2017; Thornton *et al.*, 2021), where the transcriptional profiles of overexpressed reference genes (wild-type) are systematically compared with those of their variants (mutants) to assess impact. In this case, a bead-based, high-throughput transcriptional profiling method called L1000 was used (Peck *et al.*, 2006; Subramanian *et al.*, 2017).

We hypothesized that another profiling readout, image-based profiling, could also be used for variant impact phenotyping. Image-based profiling has proven powerful in more than a dozen applications in biological research and drug discovery (Chandrasekaran *et al.*, 2020). We sought to develop cell morphology-based variant impact phenotyping (cmVIP) as a way to assess

This article was published online ahead of print in MBoC in Press (<http://www.molbiolcell.org/cgi/doi/10.1091/mbc.E21-11-0538>) on March 30, 2022.

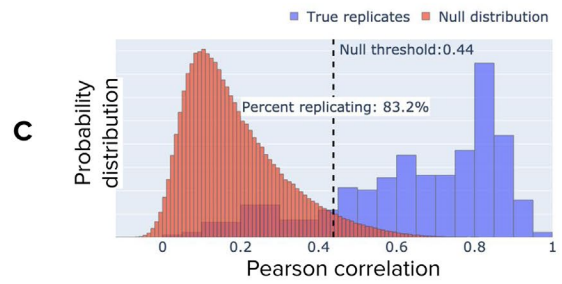
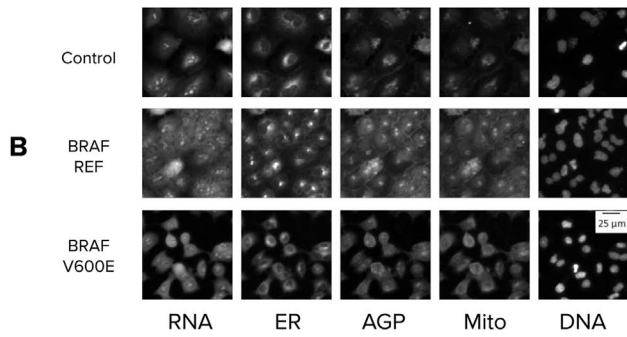
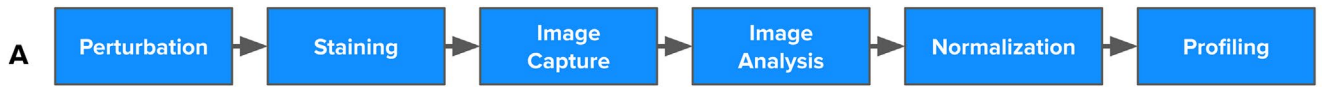
Competing financial interests: The authors declare no competing financial interests.

*Address correspondence to: Shantanu Singh (shantanu@broadinstitute.org); Anne E. Carpenter (anne@broadinstitute.org).

Abbreviations used: CCA, canonical correlation analysis; cmVIP, cell morphology-based variant impact phenotyping; COF, change of function; CTL, empty control; GOF, gain of function; L1000, high-throughput transcriptional profiling platform; LOF, loss of function; LUAD, lung adenocarcinoma; NT, neutral; ORF, open reading frame; REF, reference gene; UMAP, uniform manifold approximation and projection; VAR, variant allele; VIP, expression-based variant impact phenotyping; VUS, variants of unknown significance.

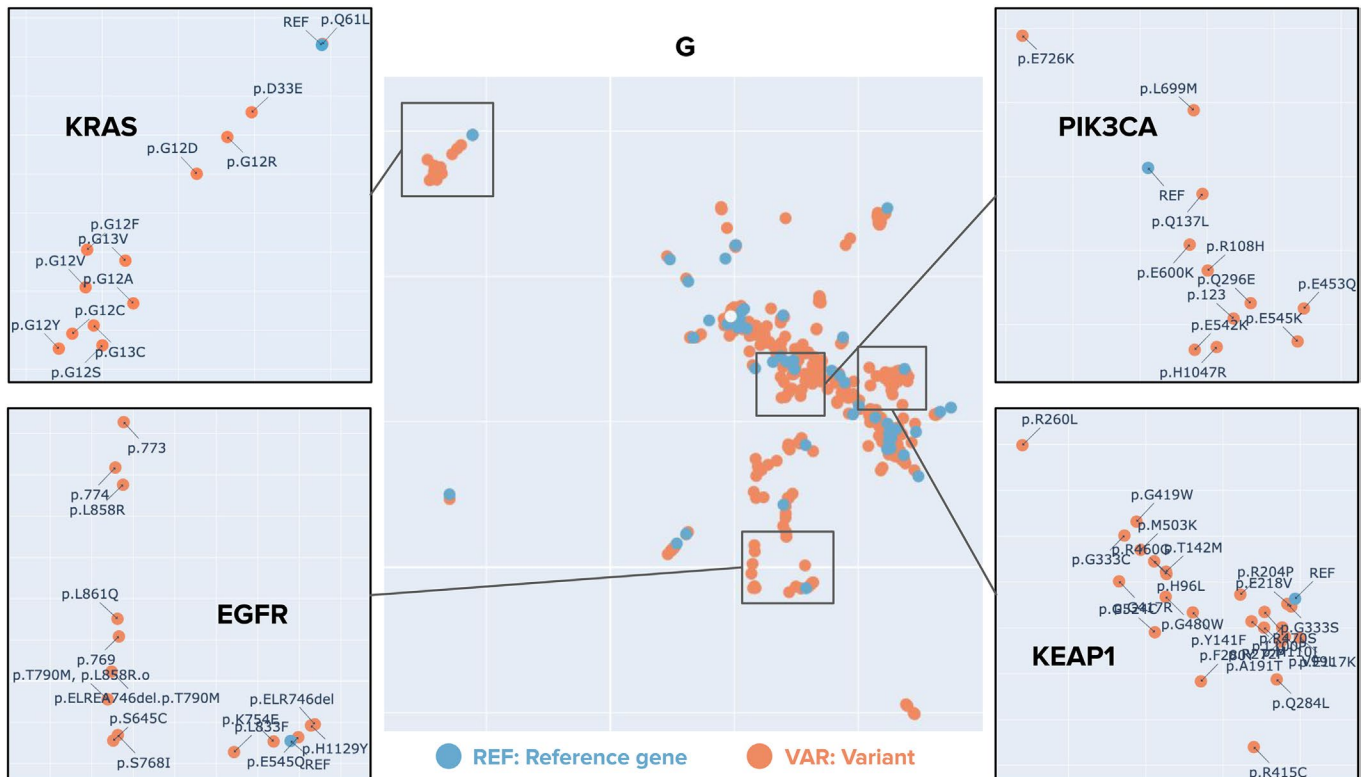
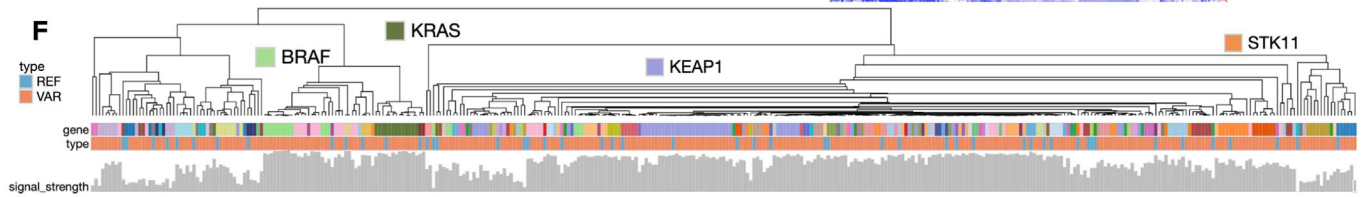
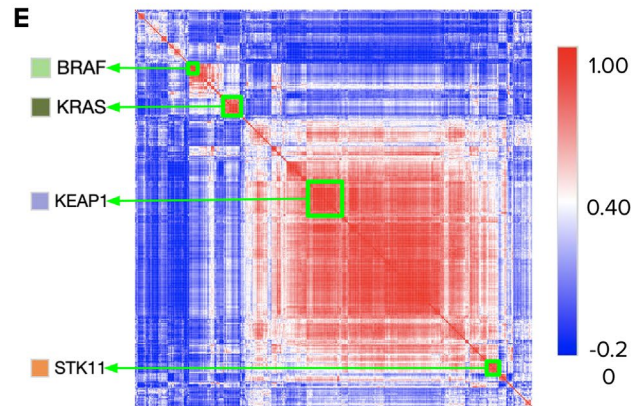
© 2022 Caicedo *et al.* This article is distributed by The American Society for Cell Biology under license from the author(s). Two months after publication it is available to the public under an Attribution-Noncommercial-Share Alike 4.0 International Creative Commons License (<http://creativecommons.org/licenses/by-nc-sa/4.0>).

"ASCB®," "The American Society for Cell Biology®," and "Molecular Biology of the Cell®" are registered trademarks of The American Society for Cell Biology.



D

ABCB9	EGFR	MAPK7	RB1
AKR1B1	ERBB2	MAX	RBM10
AKT1	FAM173B	MDM2	RBM45
ARAF	FBXW7	MET	RIT1
ATF2	FCGR3B	MYC	SERPINB5
BRAF	FGFR3	NFE2L2	STC2
CASP8	GPR137B	NRAS	STK11
CCND1	HAVCR2	PIK3CA	TPK1
CGREF1	HRAS	PIK3R3	U2AF1
CTNNB1	IDH2	PPP2R1A	ZBTB24
DARC	KEAP1	PTEN	ZNF597
DCAF8	KRAS	RAF1	
DOK1	MAP2K1	RARA	



the functional impact of coding variants inexpensively for many genes using the same systematic assay. If scaled up, a catalogue might be created of all possible variants in a given oncogene or tumor suppressor to help guide clinicians.

Here, we present a systematic study of the ability of image-based profiling to characterize lung cancer variants. We conducted a high-throughput Cell Painting (Bray *et al.*, 2016) experiment using gene overexpression in A549 cells to investigate the extent to which cell morphology can reveal sufficient phenotypic differences between reference genes and variants. We developed deep learning-based computational methods to transform images of cells into high-dimensional phenotypic profiles and used them to quantify the impact of variants. In addition, we compare the performance of image-based profiling with respect to gene expression profiling to capture phenotypic changes induced by variants and to predict their functional impact.

RESULTS

Cell Painting captures a diversity of allele phenotypes

We tested 375 overexpression perturbations (50 reference genes and 325 variants) in A549 cells using the Cell Painting assay in 384 well plates with eight replicates each (*Materials and Methods*). The overexpression construct set was previously created to test the expression-based variant impact phenotyping (eVIP) method (Berger *et al.*, 2017) and contains variants previously identified by exome sequencing primary lung adenocarcinomas (Cancer Genome Atlas Research Network, 2012), as well as their reference genes. They include many known impactful variants as well as many variants of unknown significance (VUS). As negative controls, we used wells with untreated cells that we call EMPTY controls.

We found that the Cell Painting assay can detect phenotypic signals for the majority of alleles (83.2%); this is an important first step in determining the impact of variants. We evaluated this as follows: after acquiring Cell Painting images for each sample (Figure 1B), we transformed them into replicate-level allele profiles using a deep learning-based workflow (Pawlowski *et al.*, 2016; Michael Ando, McLean and Berndt, 2017; Figure 1A; see also *Materials and Methods*). We evaluated the quality of profiles using the percent replicating score (Way *et al.*, 2021), measured as the percentage of perturbation reagents whose replicates consistently have higher similarity (reproducible signal) than random sets of perturbations; in this case 83.2% (Figure 1C).

Variant phenotypes cluster consistently with the corresponding reference gene's phenotype

Having determined that most reference genes and their variants' overexpression produced a replicable profile, we next sought to assess the structure of the relationships (similarities) among those profiles. Cell Painting has been shown to recapitulate genetic pathway relationships between reference genes in overexpression perturbation experiments (Rohban *et al.*, 2017); here, we tested a large number of variants per gene, together with the reference alleles. After aggregating replicate-level profiles into perturbation-level profiles to obtain a high-dimensional representation of each allele in our experiment, we clustered them.

The correlation matrix (Figure 1E) displays a large set of alleles that have closely similar phenotypic characteristics, which indicates that within this dataset most cancer variants share the same major phenotype. Cell Painting profiles are still able to capture subtle and meaningful variations between alleles as reflected in the continuous groups of reference genes and their corresponding variants in the hierarchical clustering (Figure 1F, color bar marked "gene") and in the UMAP data visualization (Becht *et al.*, 2018) (Figure 1G).

Because the profiles of most variants tend to cluster together within each gene, as observed in the hierarchical clustering of the correlation matrix (Figure 1F), we conclude that the phenotypic changes of variants remain closely related to the reference allele and rarely result in a major phenotypic disruption that places them in a different cluster. This type of closely related variation is consistent with previous studies in morphological and transcriptional profiling (Squires *et al.*, 2020; Way *et al.*, 2021), which report that the major factor of variation detected by profiling platforms is first associated with cell lines, then with groups of perturbations that share similar mechanisms, and finally with specific effects of each perturbation.

Interestingly, for a subset of alleles with functional annotation, Cell Painting profiles cluster the data in two major parts in the correlation matrix (Supplemental Figure 1): one part is enriched with variants from known oncogenes such as BRAF, EGFR, KRAS, and CTNNB1, and the other part is enriched with variants from known tumor suppressor genes, including FBXW7, KEAP1, and STK11. This result confirms that morphology captures relevant cellular changes associated with known cancer biology.

FIGURE 1: Cell morphology captures phenotypic variation of lung cancer alleles. (A) Workflow to create image-based profiles by transforming Cell Painting images into quantitative, multivariate representations of the states of cells impacted by each allele (whether a reference gene or variant overexpressed in the cells). (B) Example Cell Painting images under three experimental conditions: empty controls, BRAF reference gene overexpression, and BRAF V600E allele overexpression. Images are random crops of 200 × 200 pixels from a field of view (1080 × 1080), and each channel has been independently rescaled to fit the visible intensity range. Fluorescent channels: RNA/nucleoli and cytoplasmic RNA (SYTO 14), ER/endoplasmic reticulum (concanavalin A), AGP/actin, Golgi and plasma membrane (phalloidin and WGA), Mito/mitochondria (MitoTracker Deep Red), and DNA/nucleus (Hoechst 33342). (C) Distribution of true replicates vs. a null distribution of randomized replicates in this experiment, resulting in 83.2% of all perturbations having high self-correlation. Note that the null threshold (above which significant correlations are detected) is 0.44 in the Pearson correlation scale of [-1,1]. (D) List of genes included in our study; some genes whose variants are grouped in the dendrogram are outlined. For each gene, we tested several variants. (E) Correlation matrix between all pairs of perturbations (reference and variant overexpression) sorted according to the hierarchical clustering of the rows and columns. (F) Dendrogram depicting groups found by the hierarchical clustering in the correlation matrix. The type bar coloring refers to whether the perturbation is a *reference* sequence or a *variant*. The *gene* bar is colored according to the color code in D. (G) UMAP plots of reference genes' and variants' perturbation-level profiles (combining data from all replicate wells). Clusters of reference genes and their variants are observed and four examples are zoomed in (full-scale figures available at <http://broad.io/cmvip/umap.html>).

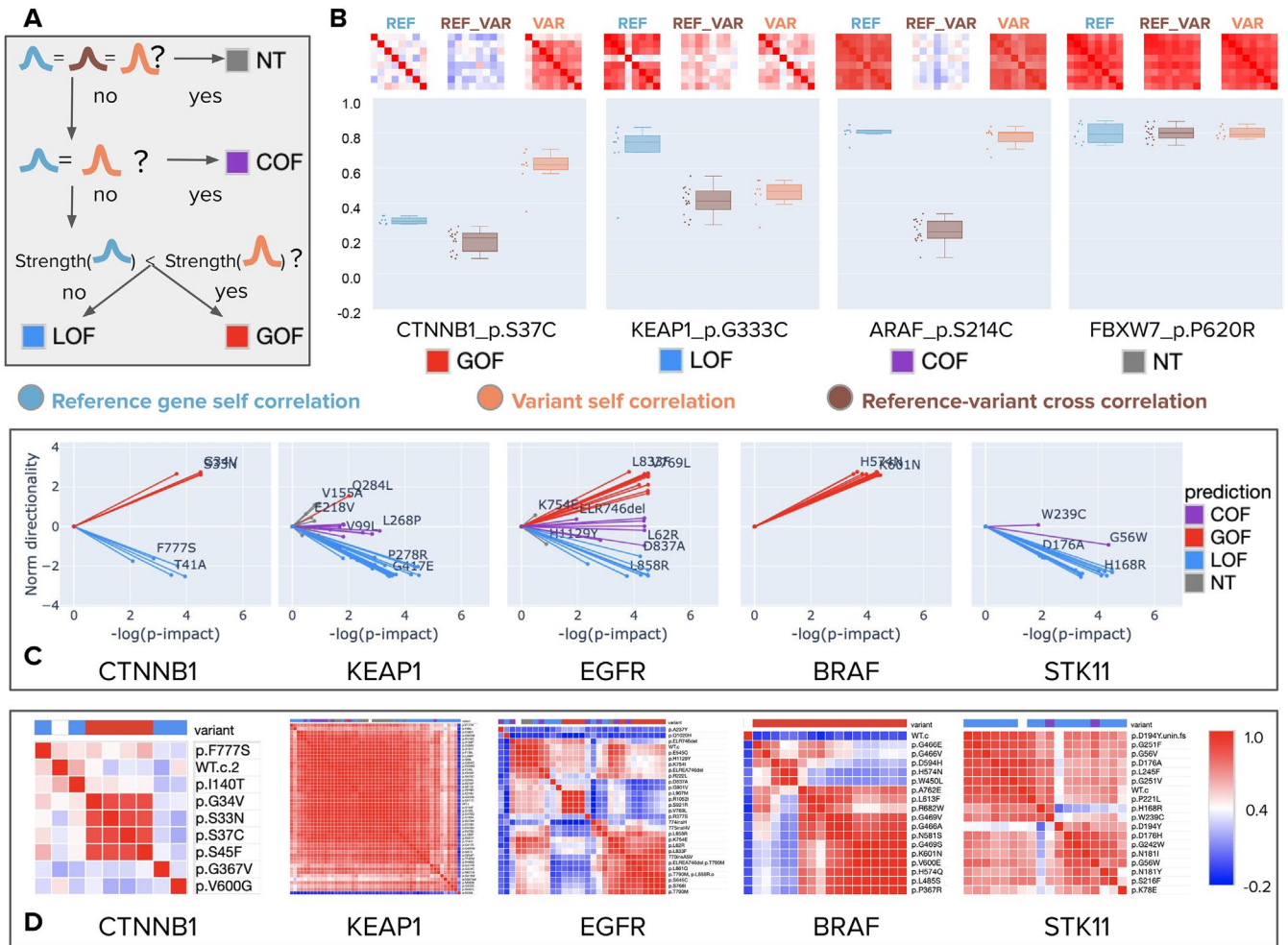


FIGURE 2: Morphology-based variant impact phenotyping (cmVIP) and resulting predictions in a diverse set of variants. (A) Decision tree of the VIP algorithm (Berger *et al.*, 2017; Thornton *et al.*, 2021), which we adopt for classifying variants by their Cell Painting profiles as gain of function (GOF), loss of function (LOF), change of function (COF), and neutral (NT) mutations. (B) Example predictions by cmVIP on four variants, one of each type. The correlation matrices at the top show how similar the replicates of each pair are (reference gene self-correlation, reference-variant cross-correlation, and variant self-correlation). The correlation matrix colors represent the correlation values in the same color scale as in D. The boxplots below the matrices show the distribution of median values of the matrices' rows (self-correlation) and columns (cross-correlation). (C) Sparkler plots display the magnitude and directionality of predictions for all variants in a gene set. The x-axis represents the negative log p-value of the impact test (the larger the more impactful), and the y-axis represents the log p-value of the directionality test polarized by the result of the strength test. All variants for these genes are displayed, but only a few are annotated to aid visualization. All the plots and annotations can be queried at full scale in the interactive website: <http://broad.io/cmvip>. (D) Correlation matrices for the groups of variants presented in C, together with the corresponding reference alleles.

Cell morphology-based variant impact phenotyping correctly classifies benchmark variants

We next tested whether the detected differences in morphology can predict each variant's impact on gene function. Using the decision tree from prior eVIP (Berger *et al.*, 2017; Thornton *et al.*, 2021), we tested for significant differences in the similarity between cell morphological profiles of reference genes and their variants (*Materials and Methods*). We call this extension of VIP *cell morphology-based VIP* (cmVIP), which interprets replicate correlations among alleles as probability distributions that can be compared using statistical tests (Figure 2A; *Materials and Methods*).

We found that cmVIP correctly classified 100% of the set of 20 well-characterized variants (Supplemental Table 1) that Berger *et al.* previously used in evaluating eVIP. This set of 20 variants has been characterized previously using functional assays. We also predicted

the directionality of the variants in this benchmark set and found that cmVIP correctly classifies 16 out of the 20 variants into one of two groups: change of function (COF) or gain of function (GOF) variant vs loss of function (LOF) variant (Supplemental Table 1).

Finally, we also estimated the false positive rate of cmVIP with mock alleles using a set of high-replicate controls. We collected 64 replicates for each of these control alleles (known to have high phenotypic activity), and then we sampled random groups of eight replicates without replacement to simulate reference genes and variant pairs. Next, we ran the cmVIP analysis to determine if this mock pair had an impact; we expected a negative answer as a result. We ran this simulation 1000 times and found that cmVIP falsely calls the mock alleles impactful 6.75% of the time on average (Supplemental Table 2), close to the false discovery rate of 5% at which the testing procedure is controlled. These results suggest that Cell

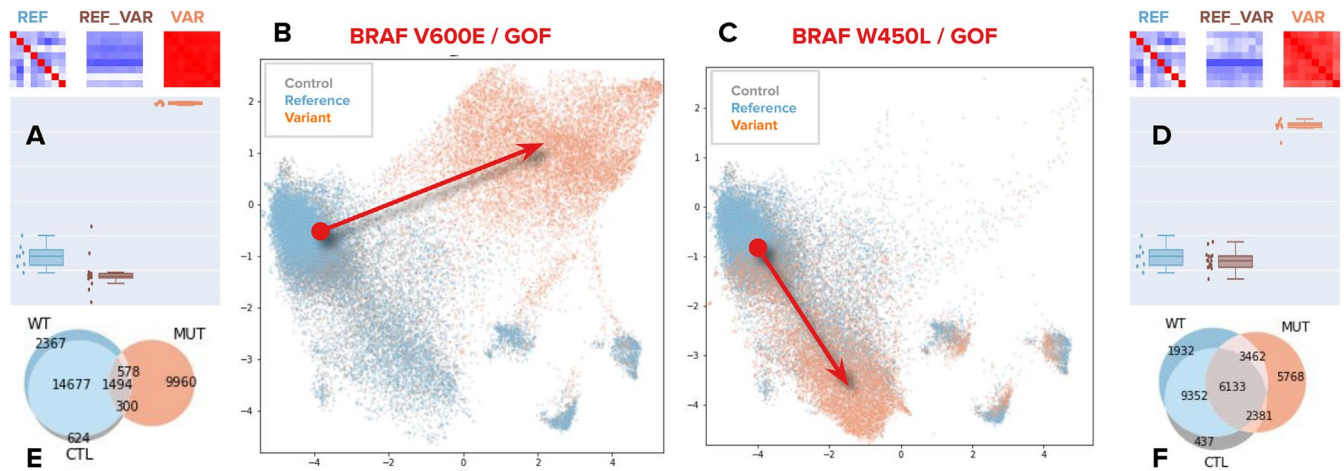


FIGURE 3: Single-cell heterogeneity of variants. Different mutations of the same gene result in different phenotypes. (A, D) Correlation matrices and box/dot plots of bulk-level profiles for the corresponding alleles, as in Figure 2. These matrices are used to obtain the impact and directionality predictions with cmVIP. (B, C) UMAP visualizations of three populations of cells, the empty control population (in blue), the reference gene population (in green), and the variant population (in orange). Each point in the plots is a single cell extracted from the Cell Painting images using segmentation. The UMAP embedding for all panels is computed using a fixed sampling of negative control wells. Arrows indicate the shift in phenotypic space from the reference gene population to the variant population. Note that variants of the same gene move in different directions. (E, F) Venn diagrams of the overlap among the reference gene, variant, and control populations of cells. These counts are obtained using graph analysis in the original feature space (*Materials and Methods*).

Painting can reliably predict the impact status of variants of unknown significance.

The impact of variants of unknown significance can be predicted at high throughput with Cell Painting

We next explored the full set of reference and variant alleles. cmVIP found 258 variants (79.3%) to be impactful; from these 158 variants (48.6%) were classified as GOF or COF variants, and 100 as LOF variants (30.7%). We show examples (Figure 2B) and provide an online resource to explore all genes and their variants (<http://broad.io/cmvip>).

Similarly to eVIP (Berger *et al.*, 2017; Thornton *et al.*, 2021), the cmVIP decision tree (Figure 2A) starts by looking at the correlation matrices of reference gene replicates (REF self-correlation) and variant replicates (VAR self-correlation) as probability distributions. Given that the image-based profiling workflow involves control-based normalization (*Materials and Methods*), we expect self-correlation matrices (correlation values between true replicates) to have high signal when the underlying phenotype is different from negative controls. This interpretation applies to reference gene and variant self-correlation matrices (REF_REF and VAR_VAR in Figure 2B). Finally, the reference gene versus variant cross-correlation matrix (REF_VAR) reveals how similar is the variant in question to its corresponding reference gene.

cmVIP interprets statistically significant changes in these three distributions of similarities among replicates in a biologically meaningful way. For instance, CTNNB1 has a relatively low signal in its reference form (REF_REF median signal strength = 0.30, Figure 2B), meaning that overexpressing it in cells changes their morphology only marginally. A gain-of-function (GOF) variant in this gene (e.g., CTNNB1 S37C in Figure 2B), by contrast, yields a relatively stronger signal (VAR_VAR median signal strength = 0.50, Figure 2B) and is different from the reference (REF_VAR median signal strength = 0.25). Loss-of-function (LOF) variants, on the other hand, are usually characterized by variants with a weak phe-

notype relative to a reference allele that has a strong phenotype (e.g., KEAP1 G333C, Figure 2B). Change-of-function variants show strong phenotypes for the reference gene and variants, and they differ from each other (e.g., KRAS G12V, Figure 2B). Finally, neutral mutations show high similarity between the reference gene and variant, indicating no detectable phenotypic change (e.g., FBXW7 P620R, Figure 2B).

The statistical tests of cmVIP provide *p*-values for such differences, which can be visualized to compare the impact and directionality of a group of variants using sparkler plots (Figure 2C). These show, for example, that the KEAP1 and STK11 variants tested in our study mainly present a LOF or COF variant pattern; BRAF variants have GOF behavior, while CTNNB1 and EGFR variants present a diverse range of GOF, COF, and LOF variants.

Cell Painting reveals allele heterogeneity at single-cell resolution

Image-based profiling inherently offers single-cell resolution while being the lowest cost even among bulk profiling methods. We investigated whether single-cell morphological profiling might provide insights into the heterogeneity of allele subpopulations or other phenotypic mechanisms that cannot be observed using bulk-level data (Ursu *et al.*, 2022). We extract single cells from Cell Painting images using the seeded-watershed segmentation algorithm and then compute deep-learning feature embeddings for each one individually (*Materials and Methods*). The feature representation of single cells has been transformed using a sphering transformation with respect to a set of 1.5 million negative control cells to minimize the impact of technical variation across batches.

We found that single-cell data visualizations for each allele allowed qualitatively observing cell heterogeneity and the relationship among cells overexpressing a particular variant allele relative to its reference allele counterpart. For example, the two BRAF variants in Figure 3 (V600E and W450L) were classified as impactful GOF variants using the cmVIP algorithm: both showed a strong phenotype in

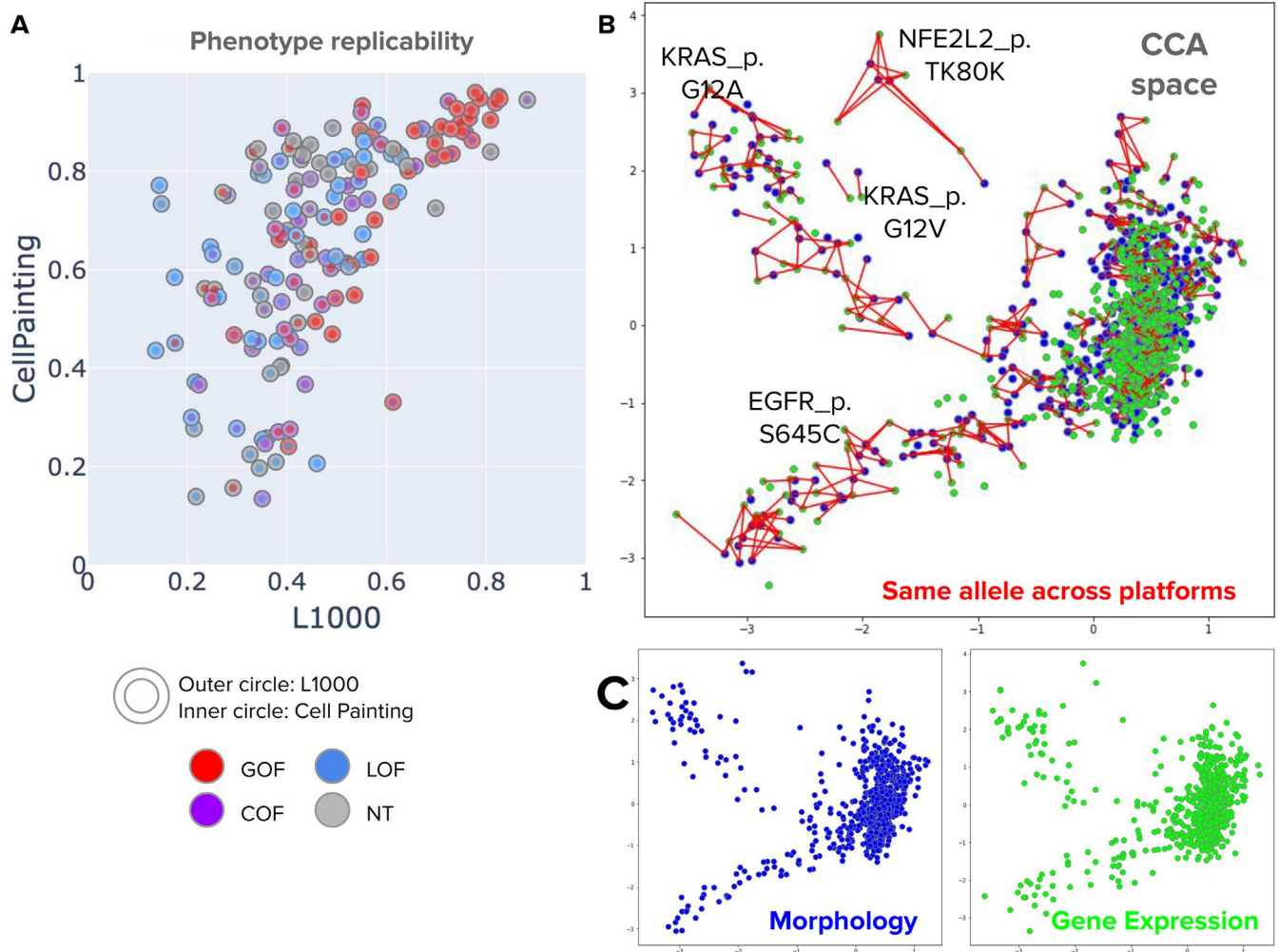


FIGURE 4: Correlation between Cell Painting profiles and L1000 profiles for a common subset of 160 alleles. (A) Signal replicability, defined as the median pairwise correlation between replicates of the same allele, was calculated for each variant in the common subset in both profiling platforms. The x-axis corresponds to the signal strength in L1000 and the y-axis represents the signal strength in Cell Painting. The Spearman correlation coefficient is 0.69. (B) Canonical correlation analysis (CCA) in the multidimensional feature space for both profiling platforms at the perturbation level. CCA obtains a common latent space by finding the directions of maximal correlation between two multivariate datasets, allowing us to project data points from Cell Painting and L1000 in the same subspace. The axes in this plot are the first and second CCA directions. Points in blue are morphology profiles and points in green are gene expression profiles. The red lines connect two points of different modalities that represent the same allele. (C) Same representation of Cell Painting profiles (morphology) and L1000 profiles (gene expression) in CCA space as in B, but using an independent plot for each platform.

the variant replicate correlation matrices compared with the reference gene, whose replicate correlation was weak (Figure 3, A and D). When looking at single cells in reduced-dimensional space (Figure 3, B and C), we observe that each variant's phenotypes move to different regions of the phenotypic space compared with the reference allele. These two different regions are not exclusive of these two variants; they are also occupied by other BRAF variants (W450L is similar to H574N and D594H, while V600E is similar to L485S, K601N, and H574Q; interactive website http://broad.io/cmvip/variants/BRAF_p.W450L/). This suggests different mechanisms between the two groups of variants; in fact, it is well known that V600E and other constitutively activating variants have different behavior than W450L and other variants of the same gene (Yao *et al.*, 2015, 2017; Dankner *et al.*, 2018).

We quantify and summarize these variations in single-cell states using graph analysis and nearest neighbors (*Materials and Methods*),

which can be observed in the Venn diagrams (Figure 4, E and F) that summarize single-cell counts with shared phenotypes (*Materials and Methods*). UMAP plots that allow single cell visualizations, as well as the corresponding Venn diagrams, are available for all the variants in our study at <http://broad.io/cmvip>.

Cell Painting phenotypic variations are highly correlated with gene expression variations

A subset of 160 variants that we profiled for this study were previously profiled using transcriptional profiling with the L1000 platform. Given the pairs of profiles for the same perturbations, we investigated the extent to which phenotypic variation captured with Cell Painting profiles corresponds with L1000 variation. Although they are not identical, we found high correlation between the two platforms in this subset of alleles by conducting two different correlation analyses (Figure 4).

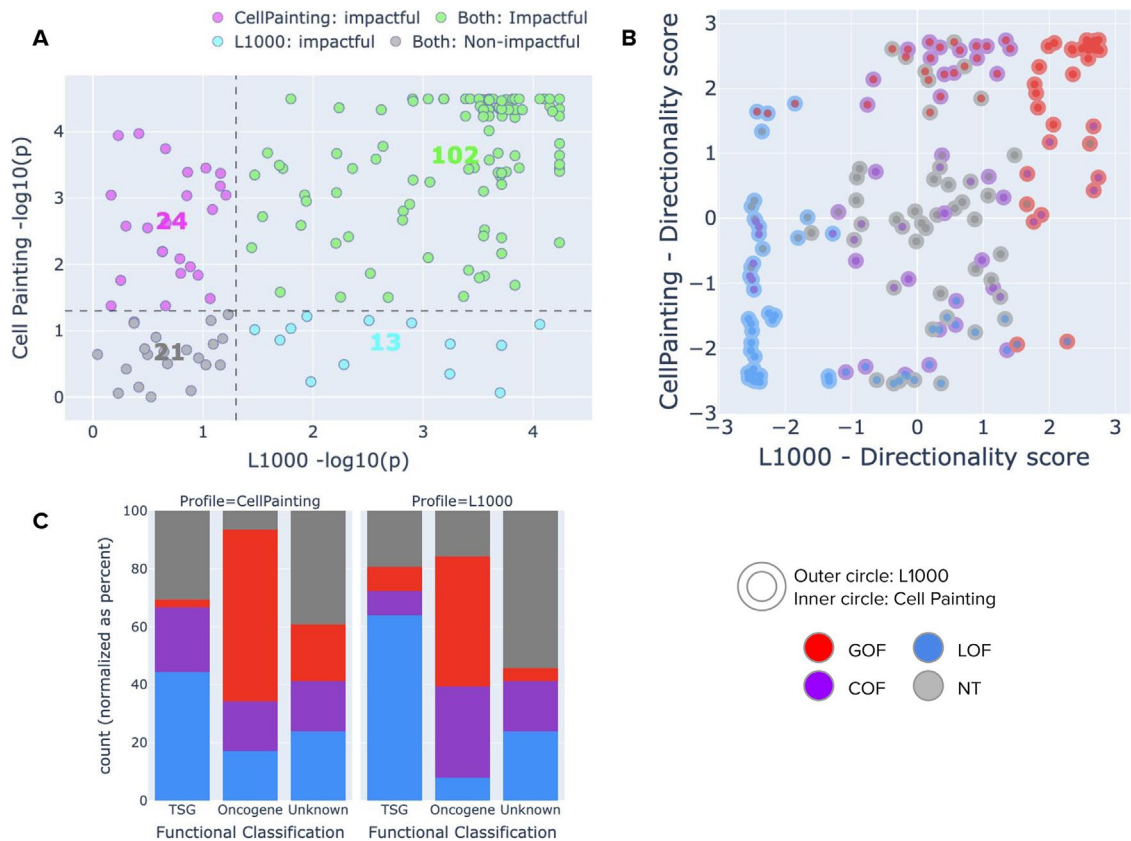


FIGURE 5: Comparison of VIP predictions using Cell Painting (morphological profiling) and L1000 (transcriptional profiling). Both platforms use the same underlying statistical tests of the VIP algorithm. (A) Impact test results. The x-axis presents the negative log p -value obtained by eVIP (L1000), and the y-axis represents the negative log p -value obtained by cmVIP. The dotted lines represent the significance threshold adopted in this study (0.05). Each point is one variant and its color indicates the prediction agreement between the two platforms: green is impactful by both platforms, gray is neutral by both platforms, pink is impactful by Cell Painting only, and blue is impactful by L1000 only. (B) Directionality test results. The x-axis indicates the polarized log p -value obtained by L1000, and the same for Cell Painting on the y-axis. Each point is one variant with the inner circle colored according to the predictions obtained by each platform. (C) Distribution of cmVIP and eVIP predictions in known oncogenes, known tumor suppressor genes (TSGs), or genes of unknown function. The distribution of oncogenes is enriched with GOF/COF calls in both platforms, and similarly, the distribution of tumor suppressor genes is enriched with LOF calls.

First, when measuring the phenotype replicability of alleles, we found a high correlation between the signal of Cell Painting profiles and the signal of L1000 profiles (Figure 4A). Phenotype replicability is defined as the median replicate correlation among true replicates of the same allele; high correlation values indicate that the underlying condition is detectable by the profiling platforms and reproducible among replicates—that is, when an allele has a high signal in L1000, it is likely to be detected with high signal in Cell Painting as well.

Second, we projected perturbation-level profiles of both platforms to the same latent space using canonical correlation analysis (CCA), which finds directions of maximal correlation between two paired multidimensional datasets. We found high agreement between profiles from both platforms when projected into the first two CCA components (Figure 4, B and C). This alignment confirms that the relative similarities and differences observed between allele phenotypes in our study can be reproduced with different assays under different experimental settings, increasing the confidence that the signal captured by both platforms is reliable and biologically meaningful.

Cell Painting predictions are consistent with transcriptional profiling predictions

We next explored how well cmVIP’s predictions matched known observations about cancer genes and variants. Beyond the 20 benchmark genes tested above (Supplemental Table 1), 140 additional variants in our study were previously characterized using transcriptional profiling via expression-based variant impact phenotyping (eVIP; Berger *et al.*, 2017).

We found that eVIP and cmVIP platforms agree on the predicted impact for 123 of the 160 variants (76.8%; Figure 5A). From those concordant predictions, 102 alleles were found to be impactful while 21 were found to be neutral. This level of agreement increases the confidence that both phenotypic profiling platforms are consistently quantifying relevant cancer biology in the underlying experiment, and also confirm that the VIP strategy generalizes well to diverse phenotypic readouts.

Next, we evaluated the agreement between the two platforms on the predicted directionality of impactful variants, and we found consistency in 21 LOF variants, 29 GOF variants, and 7 COF variants (Figure 4B). A common disagreement appears with variants that are

called GOF by one platform and COF by the other (23 variants). Other disagreements are observed between LOF and NT (17 variants) and between COF and NT (10 variants), which happen when one platform has higher phenotype strength for those variants than the other, that is, one platform detects the phenotype and the other does not. A few unexpected disagreements also appeared in five cases with LOF versus GOF directionality classifications: IDH2 K130del (CP:GOF/L1000:LOF), IDH2 S249G (CP:GOF/L1000:LOF), PIK3CA E600K (CP:GOF/L1000:LOF), RIT1 R122L (CP:LOF/L1000:GOF), and CTNNB1 V600G (CP:LOF/L1000:GOF). These may represent either occasional technical errors, or cases where the function of the reference or variant allele is undetectable by one platform versus the other.

Finally, we looked at the functional classification of genes for a few variants in the common set (Figure 5C). Our set of 160 alleles in common between the two platforms has not been completely characterized as to their GOF, LOF, COF, or NT status, but many of their genes are classified as tumor suppressors or oncogenes. One would expect that variants found in tumor suppressors are more likely to be LOF than GOF/COF, whereas variants found in oncogenes are more likely to be GOF/COF than LOF. We found that both cmVIP and eVIP make predictions consistent with these expected trends (Figure 5C).

DISCUSSION

Here we demonstrate that images of cells overexpressing given cancer-associated variants can be used to predict their impact on a diverse array of genes' functions at high throughput using the cmVIP strategy. The signal obtained from image-based profiling was sensitive to morphological variations of lung cancer variants in this experiment and was useful to characterize and make predictions for 325 variants. The accuracy appears comparable to that of transcriptional profiling, and the two platforms' predictions are generally concordant. Resolving the impact of variants at high throughput has the potential to accelerate precision oncology (Suzuki *et al.*, 2019; Vichas *et al.*, 2021).

Unbiased cell morphological profiling based on the Cell Painting assay has been shown to be a powerful approach to drug discovery and functional genomics (Caicedo, Singh and Carpenter, 2016; Chandrasekaran *et al.*, 2020). Our work expands the application of image-based profiling with Cell Painting to cancer variant phenotyping, indicating that it might be scaled up to much larger collections of variants efficiently and cost-effectively. The approach may be extended from somatic variations found in cancer to investigate the impact of germ-line variations of unknown significance in humans. Exploring a variety of cell lines and examining their concordance for variant impact prediction would be particularly interesting.

Extending this even further, it would be informative to test whether the allele-associated morphologies we observed in a cultured cell line would be identical to morphologies in cultured organoids or even tumors with the given allele. We actually suspect that this will often *not* be the case; the simplified system described here may cause certain phenotypes to manifest very differently than they would in the complex context of an organoid or tumor with all its three-dimensional cell interactions and mixtures of cell types. An interesting aspect of profiling technologies is their ability to detect similar *patterns* of morphological changes rather than precisely matching individual morphological features, and we would expect those patterns of similarity among alleles (or dissimilarity from reference allele to a given variant allele) to be more consistent with tumor samples than particular shape changes

themselves. Nevertheless, some such cases are likely to exist, where the specific morphological impact of a given allele would recapitulate in a more complex cell system such as organoids or tumors.

Image-based profiling provides single-cell resolution to investigate cellular heterogeneity across perturbations. We observed single-cell phenotypic differences between variants of the same gene, which could provide insights into functional differences of alleles. The richness of single-cell variation and the ease of implementation suggests that phenotypic studies could be performed using image-based profiling with fewer technical replicates while maintaining the ability to detect meaningful morphological variations. We leave it to future research to further investigate particular cases where single cells reveal interesting heterogeneity patterns to uncover novel cancer biology, as well as potential confounders therein.

In this work, we also used novel computational methods based on deep learning models to transform images of cells into quantitative phenotypic profiles, an approach just starting to be used in the field (Pratapa, Doron and Caicedo, 2021). The sensitivity of image-based profiling can be further increased with the advent of more powerful machine learning algorithms that extract precise patterns from images using computer vision. Our methods are open source and can be adopted for similar applications in the future, and we also expect contributions from the imaging community to develop new techniques that harness the morphology of cells for studying cellular biology.

Future studies might aim to integrate imaging and mRNA data types (if both are available) to explore whether their predictive power increases when they are combined. Our results indicate that morphology and gene expression, as captured by the Cell Painting and L1000 assays, measure highly correlated phenotypic variation, which mutually confirms their ability to detect meaningful biological events. This suggests the possibility of modeling their correspondences using computational approaches to translate one data type from the other or to understand their causal relationships. Our dataset has been used simultaneously in a study to identify which gene expression variations correspond to which morphology variations, and vice versa (Haghighi *et al.*, 2021). While this has been explored at the bulk level, our results and previous work based on scRNAseq (Ursu *et al.*, 2022) indicate that this type of analysis could be extended to understand multiomics connections at the single-cell level.

We publicly provide all data used and created in this study, including the raw images and the computed profiles (*Materials and Methods*). Further, we provide a public portal where researchers can explore alleles of interest to see the distribution of signal strength, impact and directionality predictions, VIP calls, and UMAP plots of alleles (<http://broad.io/cmvip/>).

MATERIALS AND METHODS

Profiling cancer variants with Cell Painting

Cells were grown, stained, fixed, and imaged as described in our protocol (Bray *et al.*, 2016). Briefly, A549 cells are grown in a 384-well format and infected with lentiviral open reading-frame (ORF) constructs that induce overexpression of various ORFs and alleles therein. After 96 h, MitoTracker stain was added to live cells to label the mitochondria. Cells were then fixed with formaldehyde, permeabilized with Triton X-100, and stained with the remaining dyes to identify the nucleus (Hoechst), nucleoli and cytoplasmic RNA (SYTO 14), endoplasmic reticulum (concanavalin A), Golgi and plasma membrane (wheat germ agglutinin), and actin cytoskeleton (phalloidin). Plates were imaged using an ImageXpress Micro XLS

automated microscope (Molecular Devices). We captured images from nine fields of view (sites) per well in five fluorescent channels, each using a 20x lens. Separate, grayscale image files for each channel were then stored in 16-bit TIFF format. All raw image data are publicly available at the Cell Painting Image Collection (<https://registry.opendata.aws/cell-painting-image-collection/>).

The alleles in the ORF library represent a subset of those identified in an analysis of 412 primary lung adenocarcinomas that were previously sequenced (Cancer Genome Atlas Research Network, 2012; Imielinski *et al.*, 2012), which detected 518 unique missense and in-frame insertions or deletions in the 50 genes prioritized in this study (Berger *et al.*, 2017). In all, ORF constructs for 325 variants (and reference versions) of these 50 genes were successfully generated and assayed. An additional 88 constructs are included in the dataset, representing TP53 variants that inadvertently had double mutations. A comprehensive description of the process for selecting the constructs that were analyzed is presented in Supplemental Figure 2. The additional alleles have been included in the dataset for completeness. Eight replicates were assayed for two of the plates of constructs; a third plate—comprising multiple replicate wells of a small number of “control” alleles—was assayed in two replicates.

Cell line

A549 cells (adenocarcinomic human alveolar basal epithelial cells), RRID:CVCL_0023, were obtained from ATCC; they were not additionally authenticated before this experiment. The cell line tested negative for mycoplasma before this experiment.

Mutated cDNA Library

The cDNA library is identical to that described in Berger *et al.* (2017): wild-type ORF constructs were obtained from the human ORFeome library version 5.1 (<http://horfdb.dfci.harvard.edu>) and used as templates for site-directed mutagenesis to generate mutated cDNAs in the pDONR223 Gateway entry vector. All constructs used in downstream analyses were validated by Sanger sequencing to include the intended mutation and no other identified sequence differences relative to the wild-type construct. After sequence verification, mutated ORFs were shuttled into the pLX317 lentiviral expression vector by LR recombination.

Image analysis

Illumination correction. TIFF images were corrected for non-homogeneous illumination variation across the image field using a retrospective approach (Singh *et al.*, 2014). Briefly, the method computes illumination correction functions by averaging all images of the same channel in a multiwell plate, followed by a median filter. Images in the plate are corrected by dividing their intensity values by the corresponding illumination correction function. For visualization purposes (e.g., example images reported in Figure 1), we rescale intensity values to fit the range of 255 grayscale values separately for each channel.

Segmentation. Single-cell identification was performed using CellProfiler (McQuin *et al.*, 2018; Stirling *et al.*, 2021) with the Identify Primary (nuclei) and Secondary (cell bodies) objects functionality. This approach runs thresholding and seeded watershed to identify the structures of interest. The single-cell analysis presented in this work was conducted by recording the center of the nucleus of each cell and then cropping a fixed-size region around these coordinates (see *Feature Extraction*). Cell masks were not used to isolate cells from the background.

Feature extraction. Feature extraction computes a numerical representation of the image content. Standard approaches use hand-crafted descriptors such as texture or shape features (Stirling *et al.*, 2021). Although widely used to quantify cellular morphology, they still require careful hyperparameter tuning to get high-quality representations, and, due to the high variability in the acquisition process, different datasets require custom adjustment. In contrast, representation and deep learning methods aim to find transformations automatically that yield a compact and meaningful representation based solely on image pixels. Previous empirical exploration showed promising results using deep learning models trained in the natural images and then using them to extract features from cellular images (Pawlowski *et al.*, 2016). Motivated by this and the success of transfer learning in computer vision applications, we use a pretrained EfficientNet neural network (Tan and Le, 2021) to obtain embeddings for Cell Painting images.

First, we compute a feature vector that represents the content for each segmented cell. Bounding boxes are centered on the centers of segmented cells, cropped to 128 × 128 pixels, and rescaled to 224 × 224 pixels to match the expected input of EfficientNet B0. We process each of the five Cell Painting channels independently as if they were separate RGB images by replicating their grayscale values in three channels and then running them through the EfficientNet. We keep the feature vectors of the second-to-last layer, which produces a 1280-dimensional representation for one image, and then concatenate the five vectors (one per channel), generating 6400 features to represent a single-cell profile. This process was executed using the DeepProfiler open source tool (<https://github.com/cytomining/DeepProfiler>).

Image-based profiling

In general, we followed the image-based profiling best practices defined by the community for transforming images into quantitative readouts (Caicedo *et al.*, 2017). More specifically, in order to get perturbation-level (or bulk-level) profiles, we first aggregate single-cell profiles into replicate-level (or well-level) profiles by computing their means, and then aggregate replicate-level profiles by computing their medians. In our study, we conducted a multilevel analysis of image-based profiles including perturbation-level profiles to verify associations among alleles and with gene expression data; replicate-level profiles to make impact and directionality predictions using the cmVIP algorithm; and single cell-level profiles to explore phenotype heterogeneity.

Data normalization and batch correction

As is the case in many biological experiments, imaging assays may also be prone to nuisance variation due to technical artifacts. We used *negative control sphering*, which has been shown to be effective in other studies (Michael Ando, McLean, and Berndl, 2017; Moshkov *et al.*, 2021; Way *et al.*, 2021), to correct for batch effect biases. The sphering transform used in this work makes the assumption that negative controls sampled from different batches ought to be similar to each other in the biological sense, and any deviations from this normal-looking phenotype are rather technical. Therefore, by finding a new embedding space where controls have roughly the same amount of variation in every dimension, the patterns of interest naturally emerge while batch effects are minimized. This is the same principle used in the Typical Variation Normalization (TVN) transform (Michael Ando, McLean, and Berndl, 2017).

Sphering is achieved by computing a singular value decomposition of the covariance matrix of control profiles and then scaling all the directions of the orthogonal basis by the inverse of the

corresponding eigenvalues (Kessy, Lewin, and Strimmer, 2015). The rescaled dimensions define a new representation space where large variations (usually associated with nuisance variations) are reduced, and rare variations (usually phenotypic variations) are amplified. We calculated the transformation matrix using control samples at the replicate level and used it to project all other perturbation profiles in our experiment into the corrected feature space. The sphering transform has a regularization parameter for safely inverting the eigenvalues of the covariance matrix, which was set to 0.01 in our analysis.

Cell morphology-based Variant Impact Phenotyping

Our procedure closely follows the eVIP algorithm (Berger *et al.*, 2017; Thornton *et al.*, 2021). For any given variant and its corresponding reference gene, cmVIP estimates the impact and directionality of the variant based on three correlation sets: 1) variant self-correlation: median correlation values in the rows of the replicate correlation matrix of the variant, 2) reference gene self-correlation: median correlation values in the rows of the replicate correlation matrix of the reference gene, and 3) reference-variant cross-correlation: median correlation values in the rows and columns of the correlation matrix between variant and reference gene replicates.

cmVIP follows the rule-based decision tree depicted in Figure 2A. The first stage determines if there is a statistically significant difference between any of the three correlation sets using the Kruskal-Wallis test, which is a nonparametric test. If the test rejects the null hypothesis, that is, there is a difference, then the variant is considered to be impactful; otherwise, the variant is considered to be neutral.

For impactful variants, cmVIP determines their functional directionality by running a Wilcoxon statistical test on variant self-correlations vs reference gene self-correlations. If the test rejects the null hypothesis, that is, there is a difference between variant and reference gene, then their medians are directly compared. If the median of the variant is higher than the reference one, we predict that it is a gain-of-function variant; otherwise, we call it a loss-of-function variant. In case the Wilcoxon test fails to reject the null hypothesis, that is, there is no difference between variant and reference, we predict that it is a change-of-function variant.

The Benjamini-Hochberg multiple-hypothesis correction procedure is used to control the false discovery rate of each step to be less than 5%.

Single-cell analysis

We used single-cell profiles to explore phenotypic differences between variants of the same reference gene. The first step before using single-cell profiles for quantitative analysis was to sphere the control distribution at the single-cell level (see *Data Normalization and Batch Correction* for more details). To accomplish this, we used ~1.5 million single-cell profiles taken from all 320 control wells in our experiment to compute the sphering transform. Then we projected all other single cells coming from overexpression perturbations in the corrected space. The regularization parameter used for sphering single cells was set to 0.01 (the same as in the aggregated profiles case).

Corrected single-cell profiles were then used to compute visualizations using the UMAP projection one gene at a time, including the reference gene and all its available variants. We observed that, when single cells in this UMAP visualization were colored with plate identifiers, the different replicates were well mixed and integrated (random coloring patterns; see <http://broad.io/cmvip> for examples). By computing visualizations for all alleles of the same gene at the same time, we can also qualitatively assess the relative differ-

ences among their phenotypes. We used the UMAP algorithm default parameters in their Python implementation in all cases to reveal the structure of the feature space in the most unbiased way possible.

Beyond qualitative single-cell analysis using UMAP visualizations, we used graph analysis based on nearest neighbors to objectively quantify the overlap between populations of cells in the original feature space. In this analysis, we first created a five-nearest neighbor graph using a sample of 15,000 single cells coming from three populations (5000 from each): reference gene, variant, and negative controls. The sample from each population comes from a mix of all replicates. In this graph, we proceed to classify the phenotype of single cells in one of seven categories: 1) pure reference gene phenotype, 2) pure variant phenotype, or 3) pure control phenotype, if all five nearest neighbors are from one of these three populations; 4) shared reference-variant phenotype, 5) shared reference-control phenotype, or 6) shared variant-control phenotype, if the five nearest neighbors are a mix of these two populations; finally, 7) combined phenotype, if the five nearest neighbors are a mix of the three populations. The classification of single cells into these seven categories is used to create the Venn diagrams of single-cell phenotypic overlap presented in Figure 3 and at the interactive website <http://broad.io/cmvip>.

Data and code availability

We make the data used in this project publicly available. The raw images can be downloaded from the AWS Open Data-Cell Painting Image Collection (<https://registry.opendata.aws/cell-painting-image-collection/>) in the following path: `cytodata/datasets/LUAD-BBBC043-Caicedo/`. CellProfiler was used to prepare and segment cells. The code used to process raw images and obtain deep learning features, which is based on TensorFlow (Abadi *et al.*, 2016), is available at <https://github.com/cytomining/DeepProfiler/>.

After image-based profiles were obtained, all our analysis was developed using the data science Python ecosystem, including NumPy (Harris *et al.*, 2020), SciPy (Virtanen *et al.*, 2020), Pandas, and JupyterLab, among others. All our scripts and notebooks are available at <https://github.com/broadinstitute/luad-cell-painting>. Finally, an interactive website with the aggregated data, predictions for all variants, and full-resolution figures presented in this manuscript is available at <http://broad.io/cmvip>.

ACKNOWLEDGMENTS

We thank Mukta Bagul, J.T. Neal, and Oana Ursu for helpful discussions. Funding for the project was provided by the National Institutes of Health (NIH R35 GM122547 to AEC), the Broad Institute V2F Initiative (to AEC), the Broad Institute Schmidt Fellowship program (to JCC), and the Slim Initiative for Genomic Medicine, a project funded by the Carlos Slim Foundation in Mexico.

REFERENCES

- Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G, Isard M, *et al.* (2016). TensorFlow: A system for large-scale machine learning. in OSDI. usenix.org, pp. 265–283.
- Becht E, McInnes L, Healy J, Dutertre C-A, Kwok IWH, Ng LG, Ginhoux F, Newell EW (2018). Dimensionality reduction for visualizing single-cell data using UMAP. *Nature Biotechnol* [Preprint]. doi:10.1038/nbt.4314.
- Berger AH, Brooks AN, Wu X, Shrestha Y, Chouinard C, Piccioni F, Bagul M, Kamburov A, Imielinski M, Hogstrom L, *et al.* (2017). High-throughput phenotyping of lung cancer somatic mutations. *Cancer Cell* 32, 884.
- Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A (2018). Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians* 68, 394–424.

- Bray M-A, Singh S, Han H, Davis CT, Borgeson B, Hartland C, Kost-Alimova M, Gustafsdottir SM, Gibson CC, Carpenter AE, et al. (2016). Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nature Protocols* 11, 1757–1774.
- Caicedo JC, Cooper S, Heigwer F, Warchal S, Qiu P, Molnar C, Vasilevich AS, Barry JD, Bansal HS, Kraus O, et al. (2017). Data-analysis strategies for image-based cell profiling. *Nature Methods* 14, 849–863.
- Caicedo JC, Singh S, Carpenter AE (2016). Applications in image-based profiling of perturbations. *Curr Opin Biotechnol* 39, 134–142.
- Cancer Genome Atlas Research Network (2012). Comprehensive genomic characterization of squamous cell lung cancers. *Nature* 489, 519–525.
- Chandrasekaran SN, Ceulemans H, Boyd JD, Carpenter AE (2020). Image-based profiling for drug discovery: due for a machine-learning upgrade?. *Nature Rev Drug Discovery* 1–15.
- Dankner M, Rose AAN, Rajkumar S, Siegel PM, Watson IR (2018). Classifying BRAF alterations in cancer: new rational therapeutic strategies for actionable mutations. *Oncogene* 37, 3183–3199.
- Haghighi M, Singh S, Caicedo JC, Carpenter AE (2021). High-dimensional gene expression and morphology profiles of cells across 28,000 genetic and chemical perturbations. *bioRxiv* doi:10.1101/2021.09.08.459417.
- Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith N, et al. (2020). Array programming with NumPy. *Nature* 585, 357–362.
- Imielinski M, Berger A, Hammerman PS, Hernandez B, Pugh TJ, Hodis E, Cho J, Suh J, Capelletti M, Sivachenko A, et al. (2012). Mapping the hallmarks of lung adenocarcinoma with massively parallel sequencing. *Cell* 150, 1107–1120.
- Kessy A, Lewin A, Strimmer K (2015). Optimal whitening and decorrelation. *arXiv [stat.ME]* Available at: <http://arxiv.org/abs/1512.00809>.
- Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, Carter SL, Stewart C, Mermel CH, Roberts SA, et al. (2013). Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* 499, 214–218.
- McQuin C, Goodman A, Chernyshev V, Kametsky L, Cimini BA, Karhohs KW, Doan M, Ding L, Rafelski SM, Thirstrup D, et al. (2018). CellProfiler 3.0: next-generation image processing for biology. *PLoS Biol* 16, e2005970.
- Michael Ando D, McLean C, Berndt M (2017). Improving phenotypic measurements in high-content imaging screens. *bioRxiv* doi:10.1101/161422.
- Moshkov N, Becker T, Yang K, Horvath P, Dancik V, Wagner BK, Clemons PA, Singh S, Carpenter AE, Caicedo JC (2021). Predicting compound activity from phenotypic profiles and chemical structures. *bioRxiv* doi:10.1101/2020.12.15.422887.
- Pawlowski N, Caicedo JC, Singh S, Carpenter AE, Storkey A (2016). Automating morphological profiling with generic deep convolutional networks. *bioRxiv* doi:10.1101/085118.
- Peck D, Crawford ED, Ross KN, Stegmaier K, Golub TR, Lamb J (2006). A method for high-throughput gene expression signature analysis. *Genome Biol* 7, R61.
- Pratapa A, Doron M, Caicedo JC (2021). Image-based cell phenotyping with deep learning. *Curr Opin Chem Biol* 65, 9–17.
- Rohban MH, Singh S, Wu X, Berthet JB, Bray M-A, Shrestha Y, Varelas X, Boehm JS, Carpenter AE (2017). Systematic morphological profiling of human gene and allele function via Cell Painting. *eLife* 6. doi:10.7554/eLife.24060.
- Singh S, Bray M-A, Jones TR, Carpenter AE (2014). Pipeline for illumination correction of images for high-throughput microscopy. *J Microsc* 256, 231–236.
- Squires C, Shen D, Agarwal A, Shah D, Uhler C (2020). Causal imputation via synthetic interventions. *arXiv [stat.ME]* Available at: <http://arxiv.org/abs/2011.03127>.
- Stirling DR, Swain-Bowden MJ, Lucas AM, Carpenter AE, Cimini BA, Goodman A (2021). CellProfiler 4: improvements in speed, utility and usability. *BMC Bioinform* 22, 433.
- Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, Gould J, Davis JF, Tubelli AA, Asiedu JK, et al. (2017). A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell* 171, 1437–1452.e17.
- Suzuki A, Onodera K, Matsui K, Seki M, Esumi H, Soga T, Sugano S, Kohno T, Suzuki Y, Tsuchihara K, et al. (2019). Characterization of cancer omics and drug perturbations in panels of lung cancer cells. *Sci Rep* 9, 19529.
- Tan M, Le QV (2021). EfficientNetV2: smaller models and faster training. *arXiv [cs.CV]* Available at: <http://arxiv.org/abs/2104.00298>.
- Thornton AM, Fang L, Lo A, McSharry M, Haan D, O'Brien C, Berger AH, Giannakis M, Brooks AN (2021). eVIP2: Expression-based variant impact phenotyping to predict the function of gene variants. *PLoS Comput Biol* 17, e1009132.
- Ursu O, Neal JT, Shea E, Thakore PI, Jerby-Arnon L, Nguyen L, Dionne D, Diaz C, Bauman J, Mosaad MM, et al. (2022). Massively parallel phenotyping of coding variants in cancer with Perturb-seq. *Nat Biotechnol*, 1–10.
- Vichas A, Riley AK, Nkinsi NT, Kamalpurkar S, Parrish PCR, Lo A, Duke F, Chen J, Fung I, Watson J, et al. (2021). Integrative oncogene-dependency mapping identifies RIT1 vulnerabilities and synergies in lung cancer. *Nature Commun* 12, 4789.
- Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson P, Weckesser W, Bright J, et al. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods* 17, 261–272.
- Way GP, Natoli T, Adebeye A, Litchevskiy L, Yang A, Lu X, Caicedo JC, Cimini BA, Karhohs K, Logan DJ, et al. (2021). Morphology and gene expression profiling provide complementary information for mapping cell state. *bioRxiv* doi:10.1101/2021.10.21.465335.
- Yao Z, Torres NM, Tao A, Gao Y, Luo L, Li Q, Stanchina E, Abdel-Wahab O, Solit DB, Poulikakos P, et al. (2015). BRAF mutants evade ERK-dependent feedback by different mechanisms that determine their sensitivity to pharmacologic inhibition. *Cancer Cell* 28, 370–383.
- Yao Z, Yaeger R, Rodrik-Outmezguine VS, Tao A, Torres NM, Chang MT, Drost M, Zhao H, Cecchi F, Hembrough T, et al. (2017). Tumours with class 3 BRAF mutants are sensitive to the inhibition of activated RAS. *Nature* 548, 234–238.