
Machine Learning enabled Pooled Optical Screening in Human Lung Cancer Cells

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Pooled CRISPR-based gene knockout (KO) screening has emerged as a powerful
2 method to uncover gene effects on various phenotypes [1, 2]. Recently, an optical
3 pooled CRISPR screening method was developed [3] in which gene targeting
4 guide-RNA (gRNA) are determined using *in situ* sequencing coupled with mi-
5 croscopy imaging of cellular structure and spatial features [3–6]. Pooled optical
6 screening is very scalable and cost-effective. It can be coupled with different
7 imaging assays to perform large-scale high-content image-based CRISPR-based
8 KO screens. However, development of automated and general approaches for data
9 processing and analysis are required to unlock its full potential as a tool for drug
10 target discovery. Here, we introduce a machine-learning enabled computational
11 framework for *in situ* sequencing, segmentation and feature representations of cell
12 morphology from pooled optical screens and apply it to human lung cancer cells
13 (A549). We develop a convolutional neural network (CNN) method for gRNA
14 sequence calling, and show that it increases the cell yield by 10% and enables
15 automation. We suggest self-supervised single-cell embeddings as a method to
16 create informative representations of cell morphology, moderately improving upon
17 commonly used engineered features. We demonstrate that such embeddings, aggre-
18 gated for each gene KO, are more similar for gene pairs that are known to interact
19 and cluster genetic perturbations by their cellular components, biological pathways,
20 and molecular functions. We also highlight ways to use the perturbation clusters to
21 generate hypotheses about gene functions, which are consistent with results from
22 orthogonal studies. Put together, we develop a scalable and general computational
23 approach to process and analyze pooled CRISPR-based morphological screens that
24 can be applied to screen for various disease relevant phenotypes.

25 1 Introduction

26 Pooled CRISPR KO screening technologies have been widely used for conducting large scale
27 investigation of gene effects on diverse sets of phenotypes. Recently, Feldman et al. introduced a
28 methodology for performing optical pooled screens in human cells [3, 7] by obtaining high-content
29 image-based data with their corresponding perturbation identities from pooled CRISPR screens.
30 Briefly, this approach involves transfection of a pool of cells with gRNAs to enable targeted CRISPR
31 editing. Cells are then run through a phenotyping assay such as antibody staining and fixed. The
32 gRNAs within the cells are then amplified using rolling circle amplification, and *in situ* sequencing is
33 then conducted on the plates to read out the gRNA and respective CRISPR knockout within each cell.
34 This approach was applied to study the NF κ B pathway using p65 protein localization as a readout [3],
35 and in a later work to study essential genes using intensity features derived from fluorescence markers
36 [6], but not yet to a general morphology assay. Cell Painting [8] is a morphology imaging technique
37 that is known to contain rich information about cell state, allowing practitioners to cluster compounds
38 by their MOA. However, Cell Painting phenotypic screening is typically performed in arrayed format
39 which is costly, labor intensive and is subject to batch effects [9]. Combining pooled optical screens

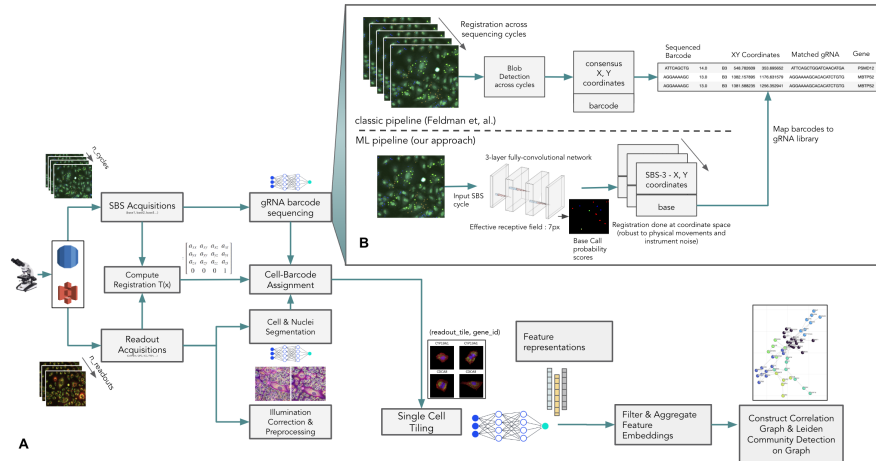


Figure 1: (A) Machine learning enabled image processing workflow for in situ sequencing, cell and nucleus segmentation and feature extraction from pooled optical screens. (B) Proposed approach for in situ sequencing using a 3-layer fully convolutional neural network followed by coordinate transformation for stitching of gRNA barcodes

40 with a general morphological profiling assay such as Cell Painting can provide an efficient and general
 41 assay for morphological screening in large genetic perturbation screens. Processing pooled optical
 42 screening data is challenging. It requires accurate gRNA sequencing, accurate segmentation of cell
 43 extents and correct association of guides to target-cells. In this work, we describe a computational
 44 framework for analyzing a screen that combines an adapted form of Cell Painting (high-content) with
 45 a pooled optical screen (high-throughput). In the following sections we present machine learning
 46 enabled methodologies for *in situ* sequencing and self-supervised feature extraction followed by
 47 the construction of a gene-gene phenotype similarity network. We demonstrate and evaluate the
 48 application of the above methodologies in learning gene similarity networks from a 300 gene (4
 49 gRNAs per gene) pooled CRISPR knock-out screen dataset containing ~ 1.5 million cells.

50 2 Machine learning improves *in situ* sequencing

51 To scale up pooled optical screening we developed a fully automated pipeline for processing (Figure
 52 1A). Hereinbelow are some of the method improvements that enabled this pipeline.

53 During the *in situ* sequencing step, each plate is processed to amplify the gRNA sequence present
 54 in each cell. These gRNAs are sequenced by synthesis (SBS) by labeling each nucleotide with a
 55 unique fluorophore, stripping, then relabeling with the next nucleotide in the sequence in a cyclic
 56 manner. This leads to a dataset in which the full plate is imaged several times, with stationary dots
 57 showing variable fluorescent signatures that need to be converted to sequencing base calls. A major
 58 step in optical screens is *in situ* sequencing of the gRNA. Feldman et. al [3], presents a computational
 59 methodology for *in situ* gRNA sequencing that requires manual alignment of field-of-view images
 60 during acquisition followed by local image registration and blob detection that requires manual fine
 61 tuning of parameters. Here, we propose an improved methodology for base-calling by training a
 62 3-layer fully convolutional neural network that takes as input a sequencing-by-synthesis fluorescence
 63 base call image with channels corresponding to fluorescent nucleotide signals (A, C, T, G) and
 64 produces a probability mask corresponding to each channel (Figure 1B). We then use the probability
 65 mask to identify base locations and the corresponding base call. The base calls are stitched based on
 66 spatial correspondence across all the SBS acquisition cycles to generate a gRNA barcode readout
 67 corresponding to each spatial location in the image (the first k ($k=10$) bases of the gRNA is
 68 referred to as a barcode in subsequent sections). Our method does not require manual alignment of
 69 field-of-view images at acquisition time, does not require manual parameter tuning, and increases the
 70 percentage of cells recovered with valid gRNA barcode from 68.6% to 78.79% in our test dataset
 71 (Table 1).

Table 1: The number of cells with a valid barcode recovered using different *in situ* sequencing methodologies

Method/Metric	Number of cells with a valid barcode	% of cells recovered with a valid barcode
SBS cycle aligned blob detector (Feldman et. al, [3])	1288234	68.60%
Blob detector + coordinate space alignment (this work)	1251669	66.65%
FCN spot detector + coordinate space alignment (this work)	1479631	78.79%

72 The gRNA barcode locations computed from the above step are projected onto the Cell Painting
73 fluorescence images using a coordinate transformation matrix constructed by image registration
74 between the acquisitions. The Cell Painting images are then preprocessed to correct for illumination
75 and intensity artifacts and single cell and nuclei contexts are segmented using CellPose [10]. Finally,
76 a single cell dataset is generated by cropping tiles centered on each nucleus and masked by its
77 corresponding cell mask. Each tile is associated with a gRNA identity based on the mapped barcode
78 locations.

79 3 Self-supervised models generate biologically informative embeddings

80 High-content image-based screens using Cell Painting have been shown to be useful in learning
81 representations and morphological profiling of perturbation effects in cells [9, 11]. Funk et. al [6]
82 demonstrated that simple fluorescence intensity and cell shape features derived from pooled optical
83 screens can be useful in defining the functional landscape of human essential genes. While explicit
84 features such as intensity and shape features can be useful, they need to be manually engineered and do
85 not capture all the kinds of variation that can occur in a perturbation dataset. Self-supervised learning
86 methods [12, 13] have been shown to improve the quality of learned representations compared to
87 supervised learning methods. Recently, SimCLR [12] and DINO-ViT [13] have achieved state-of-the-
88 art performance in learning representations from natural images. Here, we utilize these frameworks to
89 learn single-cell phenotype representations that can be used to create gene-gene phenotype similarity
90 networks. The process is as follows: 1) we extract single-cell representations 2) represent each genetic
91 perturbation by the median over all cells of that gene perturbation. 3) reduce the dimensionality
92 of these to the top 200 principal components. 4) form a correlation matrix 5) threshold to keep
93 significant correlations 6) cluster with a community detection method (Leiden).

94 To assess the performance of the methodology in learning biologically meaningful feature represen-
95 tations, we obtained evaluation metrics based on the overlap of our learned gene-gene phenotype
96 network with publicly available gene-network and ontology (STRING DB and Gene Ontology)
97 databases. For STRING DB [14] evaluation, we used the protein-protein interaction network (com-
98 bined_score > 900 as positive interaction, combined_score = 0 as no evidence of interaction) dataset
99 as the ground truth dataset. For gene ontology evaluation, we constructed a ground truth network
100 by adding an edge between a pair of genes if they belong to a common gene ontology (GO) set in
101 each of the datasets (Cellular Component, Biological Process and Molecular Function) as obtained
102 from MSigDB [15] (gene sets containing > 25 genes (out of the 300 genes in the screen) were
103 not considered for evaluation). For each of the above ground truth datasets, we computed the area
104 under the receiver operating characteristic curve (AUC) of the overlap between the latent space
105 correlation matrix and the ground truth gene network. We trained and evaluated self-supervised
106 models using the evaluation metric on the 300-gene perturbation dataset: 1. SimCLR model trained
107 using a resnet-50 backbone, 2. DINO model trained with vision transformer backbone (vit-small,
108 patch_size=16) (DINO-ViT), 3. DINO-ViT with positive pairs sampled from samples having the
109 same gene perturbation, 4. DINO-ViT with positive pairs sampled from samples having the same
110 gRNA barcode, 5. DINO-ViT model trained on ImageNet pretrained weights and fine-tuned with
111 positive pairs sampled from samples having the same gRNA barcode. We also compare the above
112 models against an ImageNet pretrained DINO-ViT model and explicitly engineered cell intensity and
113 morphology features as baselines. The results show that pre trained or fine tuned DINO-ViT features
114 outperform the commonly used engineered features (Figure 2, Table A1).

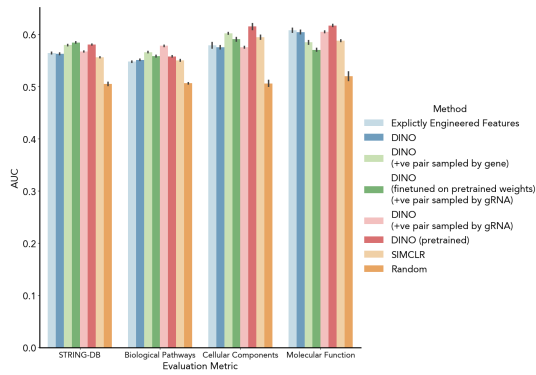


Figure 2: Comparison of feature embedding methodologies based on their ability to represent known gene-gene relationships as measured by area under the receiver operating characteristic curve (AUC)

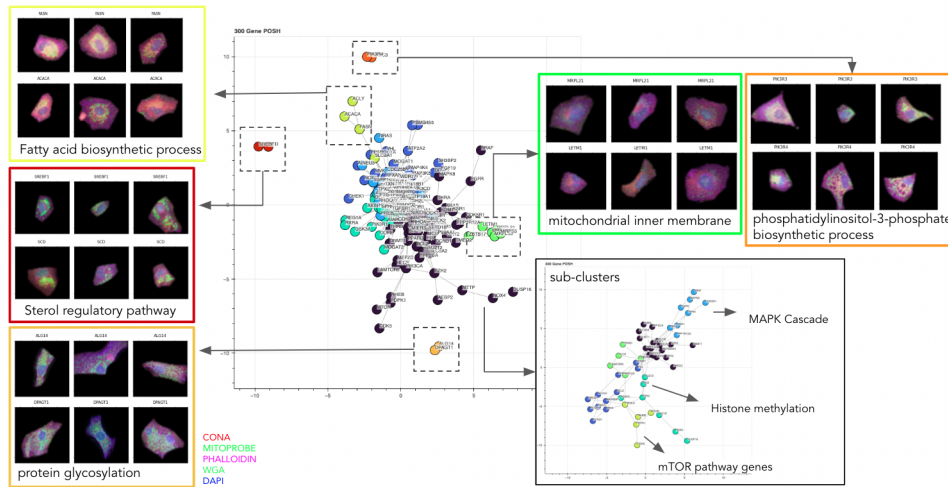


Figure 3: Community detection on pretrained DINO-ViT gene embedding correlation graph clusters genes by biological process/pathways

115 The DINO-ViT representation of the information-rich Cell Painting assay allows reconstruction of
 116 gene networks, and identification of new pathway components in a hypothesis-free way. Specifically,
 117 we were able to reconstruct the genetic modifiers of PI3K/Akt activation, protein glycosylation, fatty
 118 acid biosynthesis, sterol regulatory pathway, mitochondrial-inner membrane genes and more (Figure
 119 3). Interestingly, the network was capable of clustering key components of the lipogenesis pathways.
 120 Namely, the core fatty acid synthesis enzymes (ACLY: ATP Citrate Lyase, ACACA: Acetyl-CoA
 121 Carboxylase Alpha, and FASN: Fatty Acid Synthase), upstream AKT signaling regulators (PIK3R3:
 122 PI3K Regulatory Subunit 3, PIK3R4: PI3K Regulatory Subunit 4), and downstream palmitate and
 123 mevalonate pathway regulators (SCD: Stearoyl-CoA Desaturase, SREBF1: Sterol Regulatory Element
 124 Binding Transcription Factor 1) all contribute to lipogenesis [16]; our screening and ML approach
 125 was capable of grouping these regulators in an unsupervised manner using a generic morphological
 126 readout. While typical CRISPR screens would be capable of finding genes that increase or decrease
 127 lipogenesis, this model appears to have achieved a higher level of granularity by producing these
 128 sub-clusters. Another striking observation is the clustering of CDK5 with genes from the mTOR
 129 pathway (RHEB, mTOR and PDPK1; Figure 3, bottom left). CDK5 was recently identified to
 130 phosphorylate S6 [17], our result supports its role as an mTOR pathway regulator. In contrast to that
 131 NRAS, which is closely related to the KRAS, has a distinct morphological phenotype, different from
 132 the KRAS cluster (KRAS, BRAF and EGFR; Figure 3, top right). This observation may reflect the
 133 different function of these Ras isoforms [18].

134 In summary, streamlined data processing pipelines and *in situ* sequencing methods via automation
 135 allow for increased scale, gRNA coverage, and improved gene network reconstructions. This
 136 technique opens the possibility of whole-genome screening for numerous imaging-based phenotypes
 137 in a variety of cellular models.

References

- 138
- 139 [1] Ophir Shalem, Neville E Sanjana, Ella Hartenian, Xi Shi, David A Scott, Tarjei S Mikkelsen, Dirk Heckl,
140 Benjamin L Ebert, David E Root, John G Doench, et al. Genome-scale crispr-cas9 knockout screening in
141 human cells. *Science*, 343(6166):84–87, 2014.
- 142 [2] Tim Wang, Jenny J Wei, David M Sabatini, and Eric S Lander. Genetic screens in human cells using the
143 crispr-cas9 system. *Science*, 343(6166):80–84, 2014.
- 144 [3] David Feldman, Avtar Singh, Jonathan L Schmid-Burgk, Rebecca J Carlson, Anja Mezger, Anthony J
145 Garrity, Feng Zhang, and Paul C Blainey. Optical pooled screens in human cells. *Cell*, 179(3):787–799,
146 2019.
- 147 [4] Xiaowei Yan, Nico Stuurman, Susana A Ribeiro, Marvin E Tanenbaum, Max A Horlbeck, Christina R
148 Liem, Marco Jost, Jonathan S Weissman, and Ronald D Vale. High-content imaging-based pooled crispr
149 screens in mammalian cells. *Journal of Cell Biology*, 220(2), 2021.
- 150 [5] Michael Lawson and Johan Elf. Imaging-based screens of pool-synthesized cell libraries. *Nature Methods*,
151 18(4):358–365, 2021.
- 152 [6] Luke Funk, Kuan-Chung Su, David Feldman, Avtar Singh, Britannia Moodie, Paul C Blainey, and Iain M
153 Cheeseman. The phenotypic landscape of essential human genes. *bioRxiv*, 2021.
- 154 [7] David Feldman, Luke Funk, Anna Le, Rebecca J Carlson, Michael D Leiken, FuNien Tsai, Brian Soong,
155 Avtar Singh, and Paul C Blainey. Pooled genetic perturbation screens with image-based phenotypes. *Nature*
156 *Protocols*, 17(2):476–512, 2022.
- 157 [8] Mark-Anthony Bray, Shantanu Singh, Han Han, Chadwick T Davis, Blake Borgeson, Cathy Hartland,
158 Maria Kost-Alimova, Sigrun M Gustafsdottir, Christopher C Gibson, and Anne E Carpenter. Cell painting,
159 a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nature*
160 *protocols*, 11(9):1757–1774, 2016.
- 161 [9] Marzieh Haghghi, Shantanu Singh, Juan Caicedo, and Anne Carpenter. High-dimensional gene expression
162 and morphology profiles of cells across 28,000 genetic and chemical perturbations. *bioRxiv*, 2021.
- 163 [10] Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a generalist algorithm
164 for cellular segmentation. *Nature methods*, 18(1):100–106, 2021.
- 165 [11] Nikita Moshkov, Michael Bornholdt, Santiago Benoit, Claire McQuin, Matthew Smith, Allen Goodman,
166 Rebecca Senft, Yu Han, Mehrtash Babadi, Peter Horvath, et al. Learning representations for image-based
167 profiling of perturbations. *bioRxiv*, 2022.
- 168 [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
169 contrastive learning of visual representations. In *International conference on machine learning*, pages
170 1597–1607. PMLR, 2020.
- 171 [13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand
172 Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF*
173 *International Conference on Computer Vision*, pages 9650–9660, 2021.
- 174 [14] Damian Szklarczyk, Annika L Gable, Katerina C Nastou, David Lyon, Rebecca Kirsch, Sampo Pyysalo,
175 Nadezhda T Doncheva, Marc Legeay, Tao Fang, Peer Bork, et al. The string database in 2021: customizable
176 protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic*
177 *acids research*, 49(D1):D605–D612, 2021.
- 178 [15] Arthur Liberzon, Aravind Subramanian, Reid Pinchback, Helga Thorvaldsdóttir, Pablo Tamayo, and Jill P
179 Mesirov. Molecular signatures database (msigdb) 3.0. *Bioinformatics*, 27(12):1739–1740, 2011.
- 180 [16] T Mashima, H Seimiya, and T Tsuruo. De novo fatty-acid synthesis and related pathways as molecular
181 targets for cancer therapy. *British journal of cancer*, 100(9):1369–1372, 2009.
- 182 [17] Abul Arif, Jie Jia, Belinda Willard, Xiaoxia Li, and Paul L Fox. Multisite phosphorylation of s6k1 directs
183 a kinase phospho-code that determines substrate selection. *Molecular cell*, 73(3):446–457, 2019.
- 184 [18] Oliver Rocks, Anna Peyker, and Philippe IH Bastiaens. Spatio-temporal segregation of ras signals: one
185 ship, three anchors, many harbors. *Current opinion in cell biology*, 18(4):351–357, 2006.

Table A1: Comparison of feature embedding methods based on their ability to represent known gene-gene relationships as measured by area under the receiver operating characteristic curve of the feature correlation matrix overlap with gene relationships obtained from the respective database (STRING-DB, GO-BP = Gene Ontology Biological Process, GO-CC = Gene Ontology Cellular Components, GO-MF = Gene Ontology Molecular Function)

Embeddings/ Evaluation Metric	STRING-DB overlap (AUC)	GO-BP overlap (AUC)	GO-CC overlap (AUC)	GO-MF overlap (AUC)
Engineered Explicit Features	0.5643	0.5446	0.5882	0.6082
SIMCLR	0.5615	0.5531	0.6063	0.5812
DINO-ViT	0.5712	0.5512	0.5956	0.6137
DINO-ViT (+ve pair sampled by gene)	0.5815	0.5675	0.6050	0.5822
DINO=ViT (+ve pair sampled by gRNA)	0.5668	0.5754	0.5888	0.6062
DINO-ViT (fine-tuned on pretrained weights) (+ve pair sampled by gRNA)	0.5877	0.5617	0.5884	0.5710
DINO-ViT (pretrained)	0.5842	0.5553	0.6073	0.6182
Random	0.5084	0.4920	0.5465	0.5010