# Supplementary information

**Supplementary 1. Curated Datasets**

**Datasets List**

The details of each dataset, including the type of perturbation, number of perturbations, cell line, and number of replicates, are below. All experiments were performed in multi-well plates (384-well). The plate design varies across the datasets. The design can be visualized for each plate by inspecting any of the CSV files for each dataset (for example, for Cell Painting, *Metadata_Plate* and *Metadata_Well* provide the coordinates, and the rest of the columns with a "*Metadata_*" prefix have details of each perturbation).

**CDRP-BBBC047-Bray-CP [35]- GE [7]:**

The U2OS cell line (ATCC catalog #HTB-96) was used in this chemical perturbation dataset. For CP and GE data, data was collected after 6 and 48 hours of compound treatment, respectively; the plate maps are available at s3://cytodata/datasets/CDRPBIO-BBBC036-Bray/metadata/CDRPBIO-BBBC036-Bray /platemap/. Cell Painting was performed using the Cell Painting v1 protocol [36]. Widefield imaging was performed on an ImageXPress microscope with a 20X objective with 2x2 binning, resulting in 696x520 images. Illumination correction, segmentation, and feature measurement were performed in CellProfiler 2.2.0, per pipelines available at https://github.com/gigascience/paper-bray2017/tree/master/pipelines. L1000 was performed using the protocol described in [37].

There are 30,430 and 21,782 unique compounds for CP and GE datasets, respectively.

For CP dataset, the median number of replicates for each compound in the set is 4 and there are 26,572 replicates of control wells.

For GE dataset, the median number of replicates for each compound in the set is 3 and there are 3,478 replicates of control wells. Gene expression data for all datasets was collected using the L1000 assay [7].

20,131 compounds are present in both datasets. 6% percent of these compounds have MoA annotations. Only 3/20,131 compounds have replicate correlation more than 90th percentile of random distribution in both modalities.

**CDRP-bio-BBBC036-Bray-CP[35]-GE [7]:**

This is a subset of the previous dataset, containing the bioactive subset of compounds. See above (CDRP-BBBC047-Bray-CP) for dataset details.

There are 2,242 and 1,917 unique compounds for CP and GE datasets, respectively.

For CP dataset, The median number of replicates for each compound in the set is 8 and there are 3,528 replicates of control wells.
For GE dataset, The median number of replicates for each compound in the set is 2 and there are 3,478 replicates of control wells.

1,916 compounds are present in both datasets. 69% percent of these compounds have MoA annotations. 131/1,916 compounds have replicate correlation more than 90th percentile of random distribution in both modalities.

## LUAD-BBBC041-Caicedo-CP [38]- GE [11]:

The A549 cell line was used in this genetic perturbation dataset. For CP and GE, data was collected after 96 hours of lentiviral overexpression; the plate maps are available at s3://cytodata/datasets/LUAD-BBBC043-Caicedo/metadata/LUAD-BBBC043-Caicedo /platemap/. Cell Painting was performed using the Cell Painting v1 protocol [36]. Widefield imaging was performed on an ImageXPress microscope with a 20X objective with 2x2 binning, resulting in 1080x1080 pixel images. Illumination correction, segmentation, and feature measurement were performed in CellProfiler 2.2.0, per pipelines available at https://github.com/broadinstitute/imaging-platform-pipelines/tree/master/cellpaint ing_a549_20x_imagexpress. L1000 was performed using the protocol described in [37]

There are 593 and 529 unique alleles for CP and GE datasets, respectively.

For CP and GE datasets, the median number of replicates for each allele in the set is 8.

525 alleles are present in both datasets. 197/525 of these alleles have replicate correlation more than 90th percentile of random distribution in both modalities.

## TA-ORF-BBBC037-Rohban-CP [3] - GE:

The U2OS cell line was used in this genetic perturbation dataset; cells were obtained from ATCC and propagated in the William Hahn lab; they were not additionally authenticated prior to this experiment.
 For CP and GE, data was collected after 72 hours of lentiviral overexpression; the platemaps are available at s3://cytodata/datasets/TA-ORF-BBBC037-Rohban/metadata/TA-ORF-BBBC037-Rohb an/platemap/. Cell Painting was performed per the Cell Painting v1 protocol [36]. Widefield imaging was performed on an ImageXPress microscope with a 20X objective with 2x2 binning, resulting in 1080x1080 pixel images. Illumination correction, segmentation, and feature measurement were performed in CellProfiler 2.1, per

pipelines available as supplementary files on [3]. L1000 was performed using the protocol described in [37]

There are 323 and 327 unique alleles for CP and GE datasets, respectively.

For CP dataset, the median number of replicates for each allele in the set is 5 and there are 268 replicates of control wells.

For GE dataset, the median number of replicates for each allele in the set is 2 and there are 56 replicates of control wells.

150 alleles are present in both datasets. 36/150 of these alleles have replicate correlation more than 90th percentile of random distribution in both modalities.

**LINCS-Pilot1-CP [39] - GE [40]:**

The A549 cell line (ATCC catalog #CCL-185) was used in this chemical perturbation dataset. For CP and GE, data was collected after 24 and 48 hours of compound treatment, respectively; the platemaps are available at https://github.com/broadinstitute/lincs-cell-painting/tree/74f7a0d132e75b86aa8908689f3ce21a58a9b0f6/metadata/platemaps/2016_04_01_a549_48hr_batch1/platemap. Cell Painting was performed using the Cell Painting v2 protocol [8]; this differs from the v1 protocol only in the concentrations and Alexa dye species of the conjugated phalloidin and wheat germ agglutinin dyes (see https://github.com/carpenterlab/2016_bray_natprot/wiki#changes-in-the-official-protocol-to-create-v2-bray-et-al-2016. Widefield imaging was performed on an Opera Phenix microscope with a 20X objective with 1x1 binning, resulting in 2160x2160 pixel images. Illumination correction, segmentation, and feature measurement were performed in CellProfiler 2.2.0, per pipelines available at https://github.com/broadinstitute/imaging-platform-pipelines/tree/master/cellpainting_a549_20x_phenix_bin1. L1000 was performed using the protocol described in [7].

There are 1,570 unique compounds across 7 doses for CP dataset. There are 1,402 unique compounds across 7 doses for GE dataset. A small fraction compounds have fewer than 7 doses.

There are 9,394 and 8,369 unique compounds-dose for CP and GE datasets, respectively.

For CP dataset, the median number of replicates for each compound in the set is 5 and there are 3,264 replicates of control wells.

For GE dataset, the median number of replicates for each compound in the set is 3 and there are 1,485 replicates of control wells.

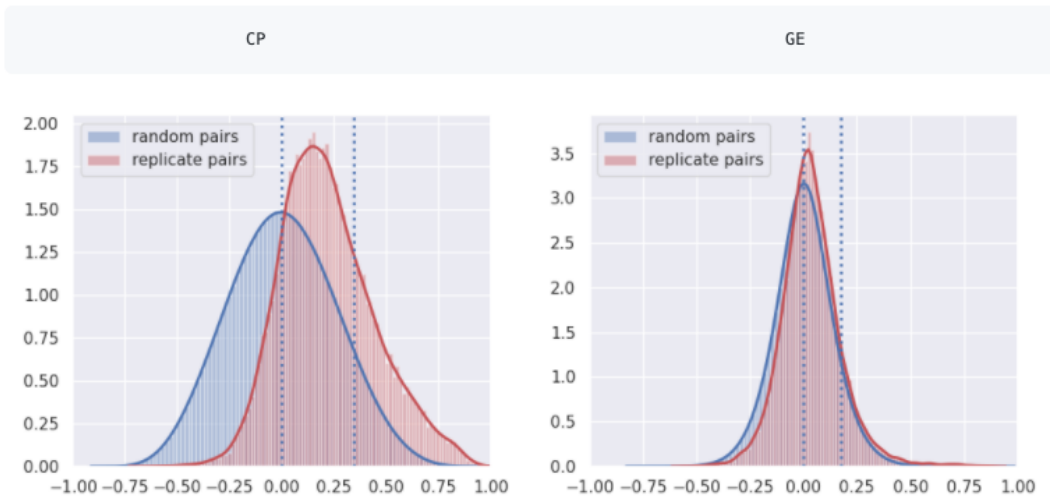6984 compound-dose pairs are present in both datasets. 100% of these compounds have MoA annotations.

Among 6984 unique compounds-dose overlapping compounds, 1140 compound-dose have replicate correlation more than 90th percentile of random distribution in both modalities.

| | | Datasets | | | |
|---|---|---|---|---|---|
| Name | Cell Line | Perturbation Type | $N_p/N_d/N_r$ CP | $N_p/N_d/N_r$ GE | $N_p$ Intersection |
| CDRP-BBBC047-Bray | U2OS | chemical | 30,430/1/4 | 21,782/1/3 | 20,131 |
| CDRP-bio-BBBC036-Bray | U2OS | chemical | 2,242/1/8 | 1,917/1/2 | 1,916 |
| LUAD-BBBC041-Caicedo | A549 | genetic | 593/-/8 | 529/-/8 | 525 |
| TA-ORF-BBBC037-Rohban | U2OS | genetic | 323/-/5 | 327/-/2 | 150 |
| LINCS-Pilot1 | A549 | chemical | 1,570/7/ 5 | 1,402/7/3 | 6,984 |

**Supplementary Table. 1.** GE and CP dataset descriptions. $N_p$ denotes the number of perturbations. $N_d$ denotes the number of doses. $N_r$ denotes the number of replicates.

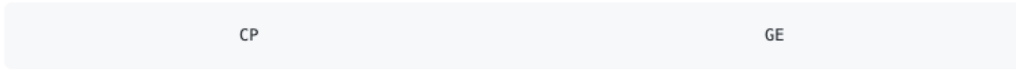**Supplementary 2. Data Quality: Replicate reproducibility**

**CDRP-BBBC047-Bray:**



**CDRPBIO-BBBC036-Bray:**
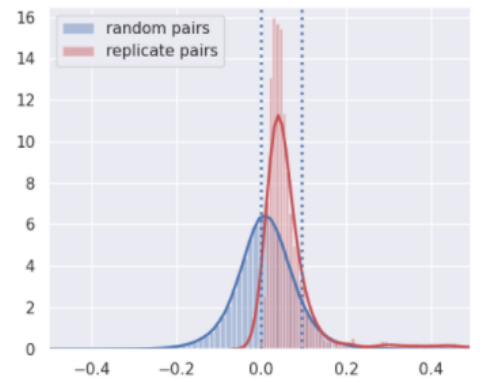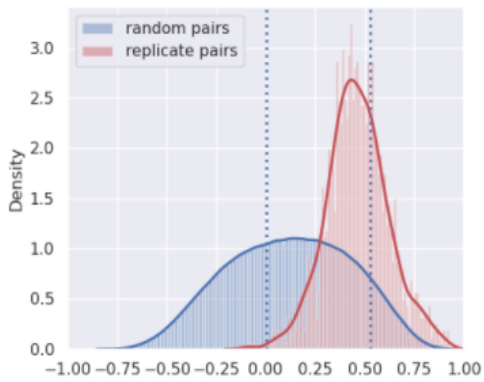


**LUAD-BBBC041-Caicedo-CP-GE :**
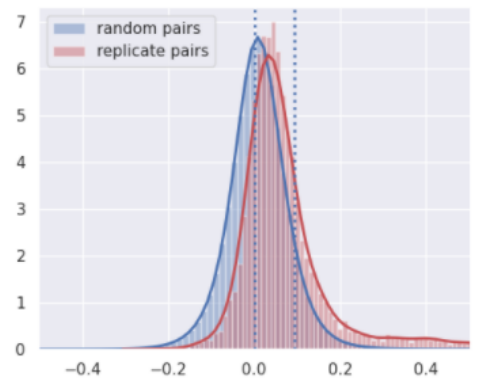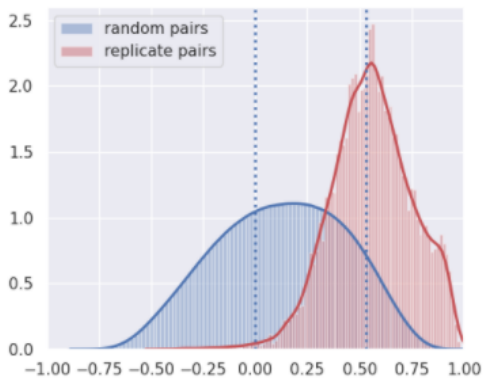
## TA-ORF-BBBC037-Rohban-CP-GE :



## LINCS-Pilot1-CP-GE :

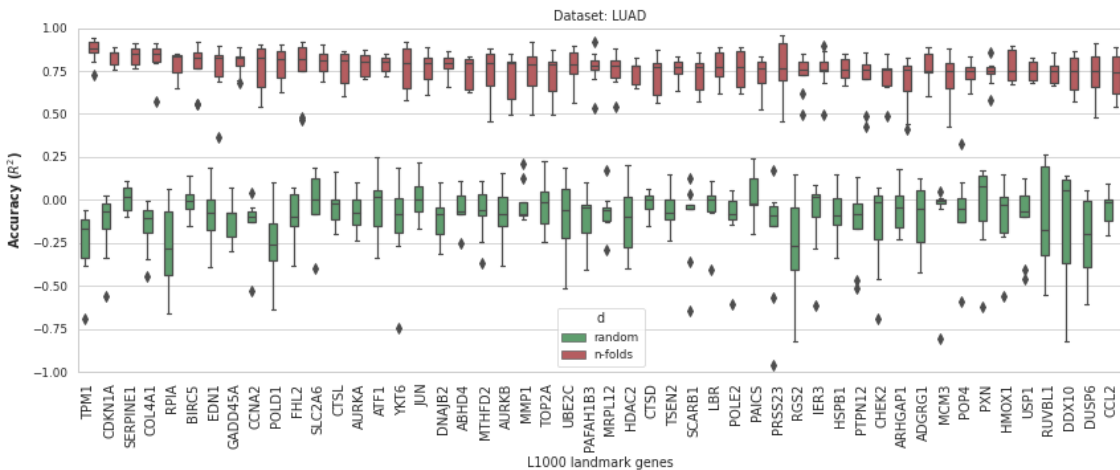- CP - All doses together (Metadata_pert_id)
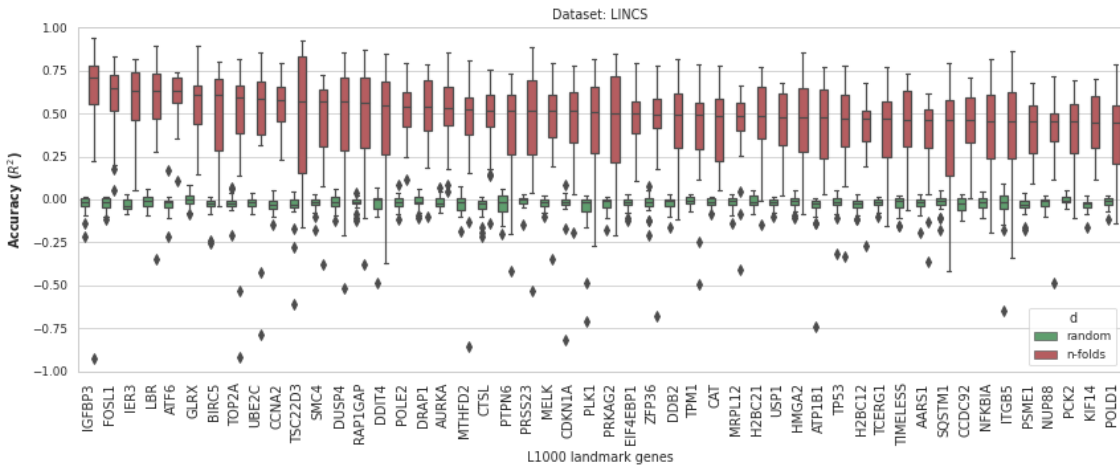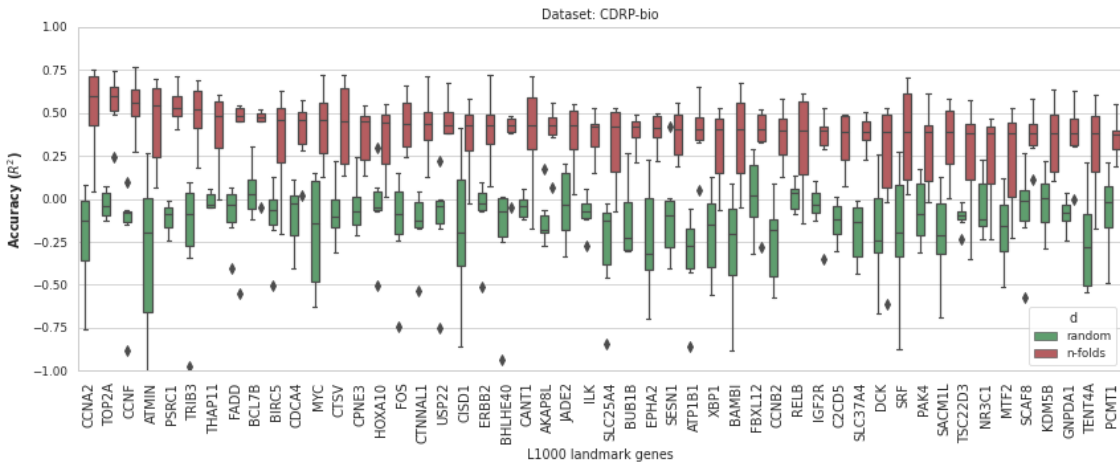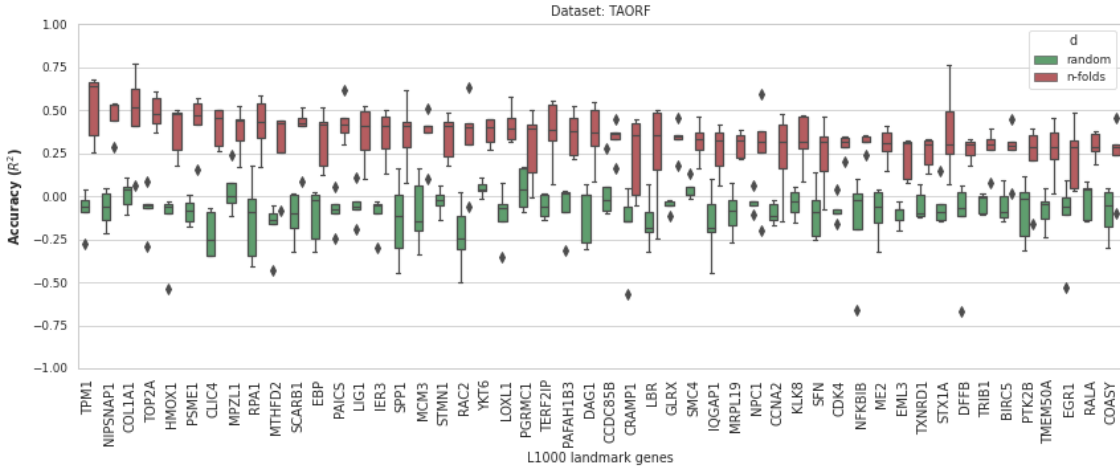- GE - All doses together (pert_id)



- CP - each sample/dose (Metadata_pert_id_dose)
- GE - each sample/dose (pert_id_dose)

**Supplementary Figure 1.** To inspect the quality of each dataset, we calculate the consistency of profiles across different replicates of the same perturbation as follows. We standardized the profiles per plate to have zero mean and unit variance. Next, we calculated the Pearson correlation coefficient between each pair of profiles for the same perturbation (red curve) and for different perturbations (blue curve). Dotted vertical lines are shown at zero and 90th percentile of the random pairs (blue) distribution.

**Supplementary 3. Top 50 highly predictable L1000 genes by Cell Painting morphological features**
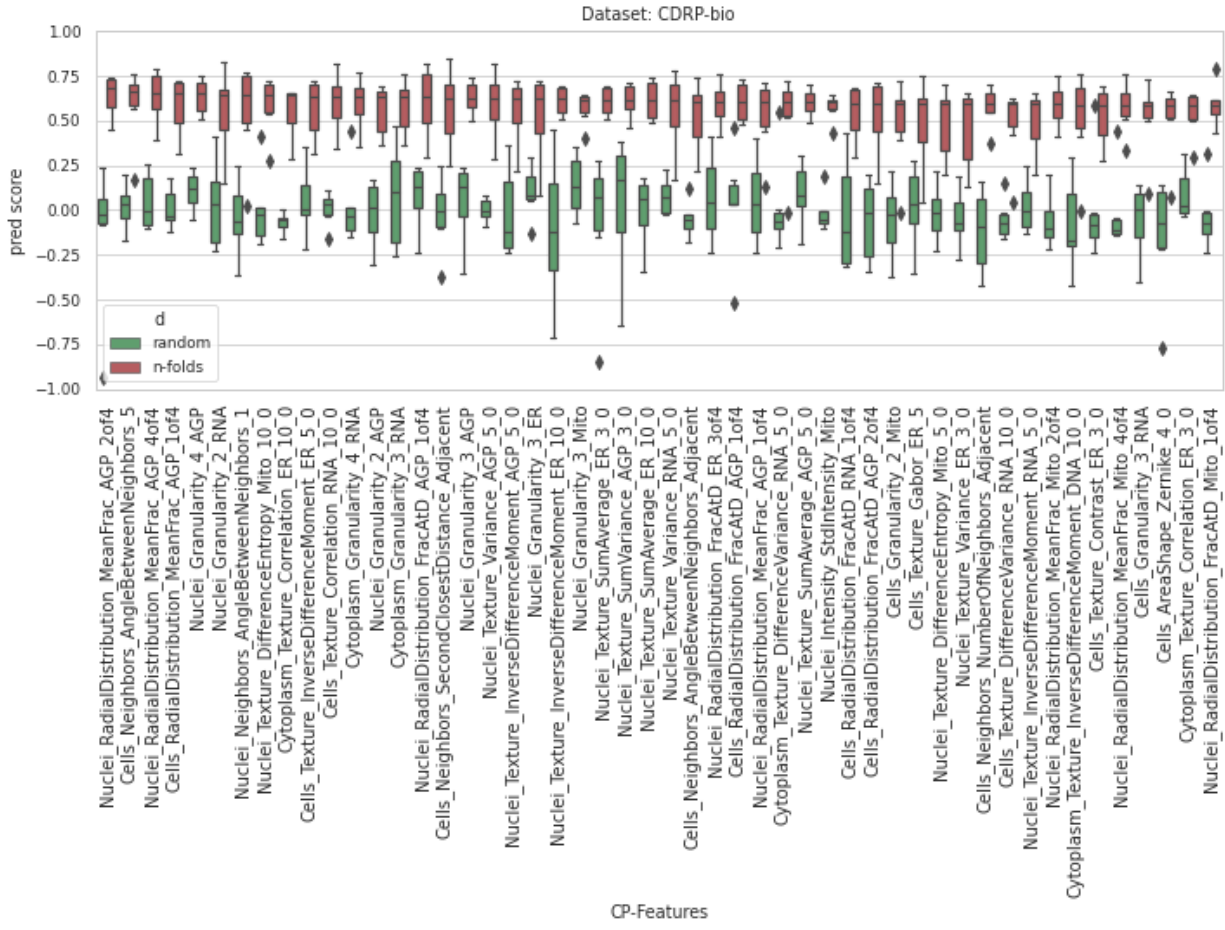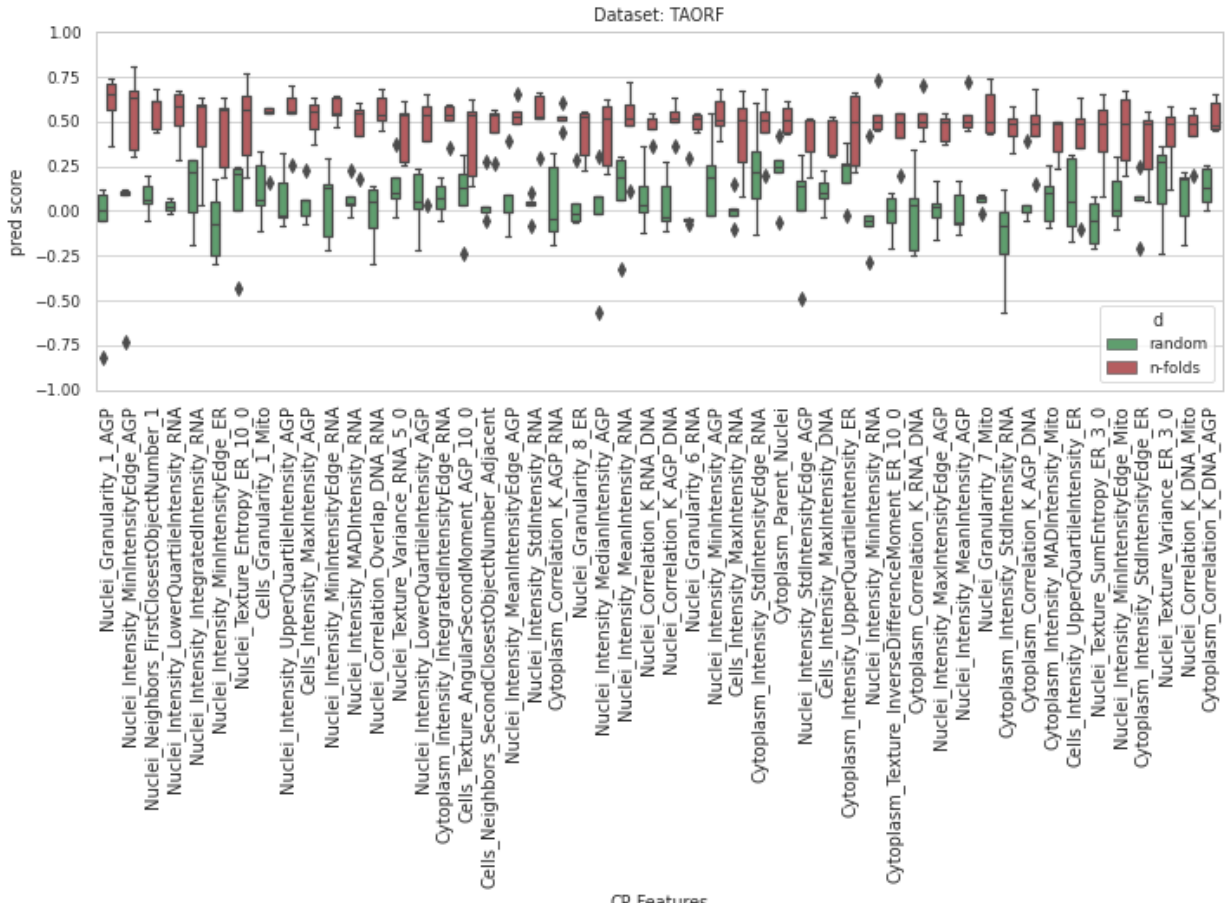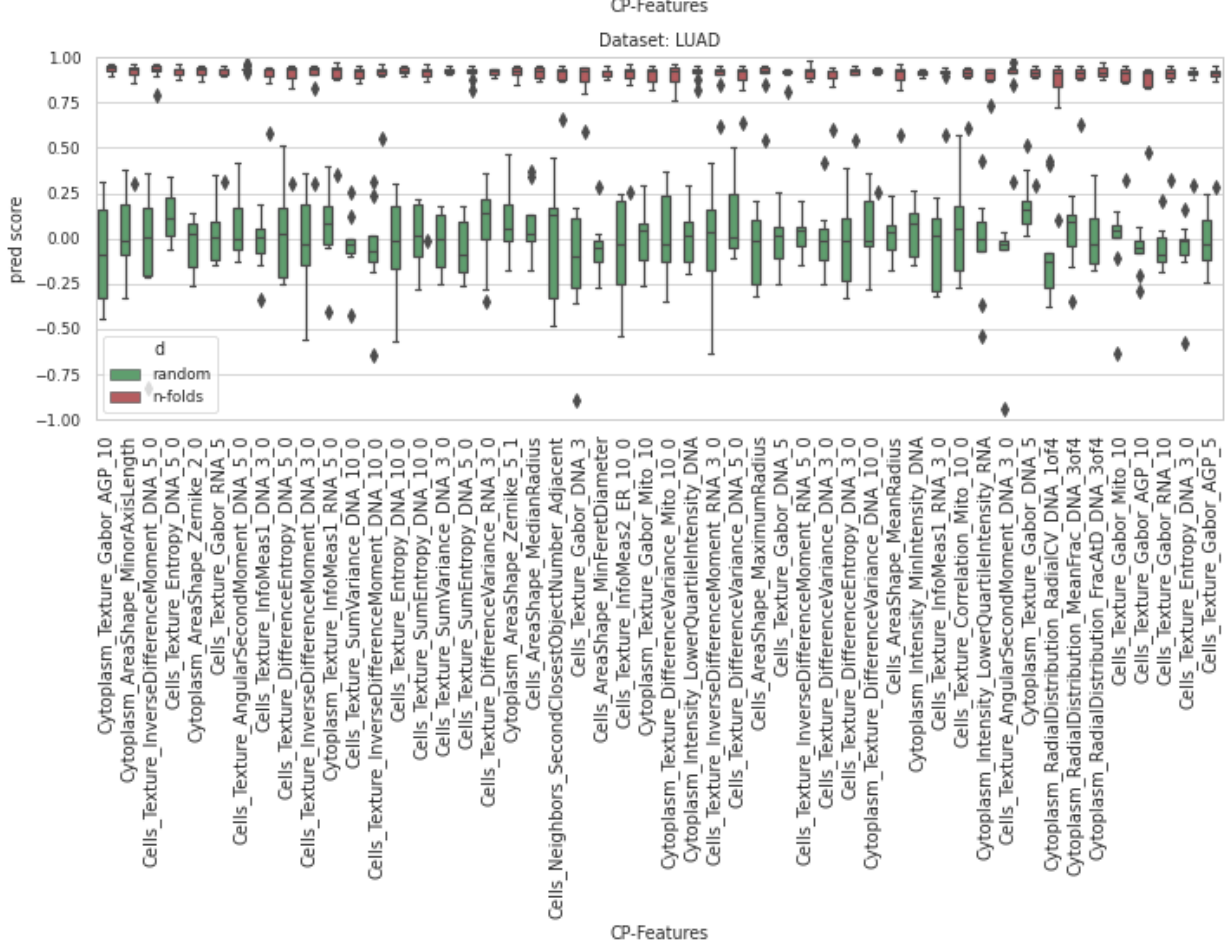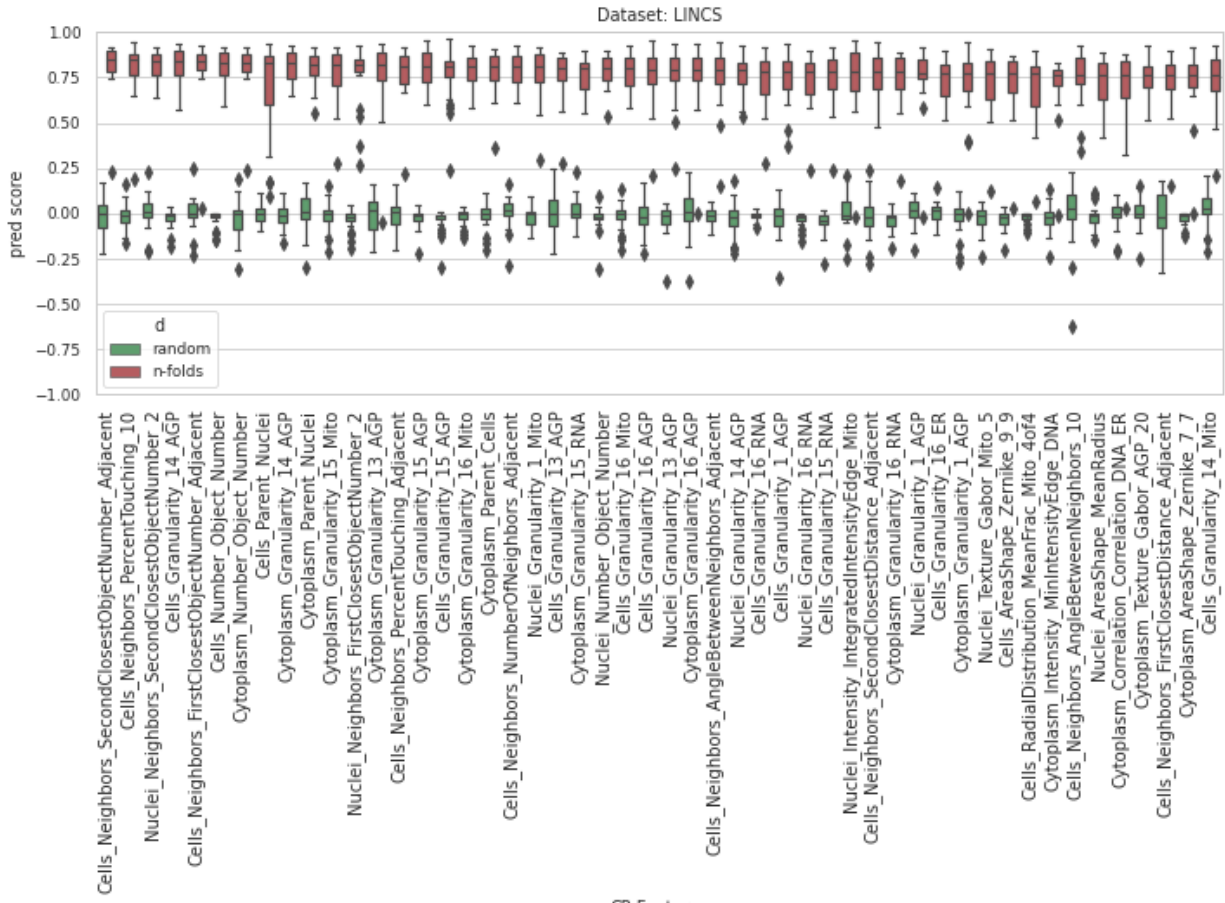
**Supplementary Figure 2.** Prediction of L1000 mRNA levels by Cell Painting features: for each dataset, the distribution of MLP baseline prediction scores for the ordered top-50 landmark genes with the highest $R^2$ median prediction scores are shown. Each (red) boxplot represents a set of k $R^2$ values that result from running k-fold cross-validation for each corresponding landmark gene. The green boxplots are the corresponding null distributions – here we apply the exact same procedure as for the blue boxplots except that we first shuffle the output. The y-axis is trimmed at -1 for clarity. Distributions are presented as boxplots, with center line being median, box limits being upper and lower quartiles and whiskers being 1.5× interquartile range; n=k=6(CDRP-bio), 5(TAORF), 9(LUAD), 25(LINCS)

## Supplementary 4. Median Prediction scores for each landmark gene across each datasets and models

Supplementary_D.csv

**Supplementary 5. Top 50 highly predictable Cell Painting morphological features by L1000 genes**

Dataset: TAORF

Dataset: CDRP-bio

Dataset: LINCS

Dataset: LUAD

**Supplementary Figure 3.** Prediction of each Cell Painting feature by L1000 mRNA levels: for each dataset, the distribution of MLP baseline prediction scores for the ordered top-50 Cell Painting features with the highest $R^2$ median prediction scores are shown. Each (red) boxplot represents a set of k $R^2$ values that result from running k-fold cross-validation for each corresponding Cell Painting feature. The green boxplots are the corresponding null distributions – here we apply the exact same procedure as for the blue boxplots except that we first shuffle the output. The y-axis is trimmed at -1 for clarity. Distributions are presented as boxplots, with center line being median, box limits being upper and lower quartiles and whiskers being 1.5× interquartile range; n=k=6(CDRP-bio), 5(TAORF), 9(LUAD), 25(LINCS).