Supplementary information: <u>Supplementary Table 1</u> and <u>Supplementary Figures 1 - 9</u>

**A statistical framework for high-content phenotypic profiling using cellular feature distributions**

Yanthe E. Pearson[1], Stephan Kremb[1], Glenn L. Butterfoss[1], Xin Xie[1], Hala Fahs[1], Kristin C. Gunsalus[1,2*]

[1]Center for Genomics and Systems Biology, New York University Abu Dhabi, P. O. Box 129188, Abu Dhabi, United Arab Emirates

[2]Department of Biology and Center for Genomics and Systems Biology, New York University, New York, NY 10003, USA

* Corresponding author: Kristin C. Gunsalus
Email: kcg1@nyu.edu

**Supplementary Table 1:** R scripts and data files underlying each manuscript figure.

| Manuscript Figure | Description | R script | Data files |
|---|---|---|---|
| Figure 2a | Heatmap of compound similarity based on Tanimoto distance | Figure2A_tanimoto.R | Tanimoto_figure2.csv |
| Figure 2c | Cell counts shown as heatmap form | figure2C_count_heatmap.R | raw_medians_fig2C.csv |
| Figure 2d | Cell counts shown as a scatterplot | figure2D_count_scatter.R | raw_medians_fig2C.csv |
| Figure 2e | Per well density curve of controls | Figure2E_dmso_curves.R | cell_data_DMSO_figure2g.csv |
| Figure 2f | Cell cycle dose response plot of mitoxantrone treated cells | Figure2F_mitoxantrone.R | cell_data_mitoxantrone_figure2f.csv |
| Figure 2g | Density plots showing quartiles for a group of compounds | figure2G_quartiles.R | Data_fig2G.csv |
|  |  |  |  |
| Figure 3a | Heatmap of control well medians | figure3A.R | nucfeatures_medians_fig3A.csv |
| Figure 3b | Summary of two-way ANOVA test | figure3B_anova.R | anova_output_fig3b.csv |
| Figure 3c | Heatmaps showing positional adjustment | figure3C_positioneffect.R | Well_medians_fig3c.csv |
| Figure 3d | Cell cycle distributions | figure3D.R | cell_data_fig3D.csv |
|  |  |  |  |
| Figure 4a | Plots for replicate feature distributions, both treatment and control | figure4A.R | cell_data_figure4A_control.csv cell_data_figure4A_treatment1.csv cell_data_figure4A_treatment2.csv |
| Figure 4b | Plots of density curves and CDF curves for two data samples | figure4B.R | No data files required |
| Figure 4c | Plot for distribution of statistical scores | figure4C-D.R | control_data_fig4C.csv treatment_data_fig4C.csv |
| Figure 4d | Plot of sorted features | figure4C-D.R | treatment_replicate_test_allpanels_fig4D.csv control_replicate_test_allpanels_fig4D.csv |
|  |  |  |  |
| Figure 5a | Plot of replicate distributions R script Includes EMD calculation | figure5A.R | cell_data_figure5a_vincristine_control.csv |
| Figure 5b,c | Global control and individual treatment groups | figure5B-C.R | treated_cells_data_fig5B-C.csv control_cells_data_fig5B-C.csv |
| Figure 5d,e Supplementary Figure 4 | Full feature EMD fingerprints of individual control samples. | figure5D-E_fullfeature.R radarchart2_new.R | raw_emd_profile174.csv |
| Figure 5f | Plot of Vincristine fingerprint | Figure5F_vincristine.R radarchart2_new.R | ScaledEMDprofile_69feats_785treatments.csv |
|  |  |  |  |
| Figure 6a | Heatmap of similarity among treatments | Figure6A_heatmap.R | fullfeature_EMDprofile_785treatments.csv |
| Figure 6b | Run the umap and plot phenotypic trajectories while projecting cell count onto UMAP. | line_segments.R Figure6B_umap.R | fullfeature_EMDprofile_785treatments.csv |
|  |  |  |  |
| Figure 7a | UMAP paths | line_segments.R figure7A_umap_paths.R | UMAP_dimensions_69feats_785treatments_new.csv |
| Figure 7b | Cell count dose response | figure7B.R | cell_counts.csv |
| Figure 7c | Cell cycle distributions | Figure7C_cellcycle_four.R | figure7c_treatments.csv figure7c_control.csv |
| Figure 7d | Radial plots (see Figure 5f) |  |  |

**Supplementary Table 1**: A summary of all R scripts and data files used for each figure in the manuscript. The files can be found on the GitHub within individual figure folders at https://github.com/GunsalusPiano/EMD.
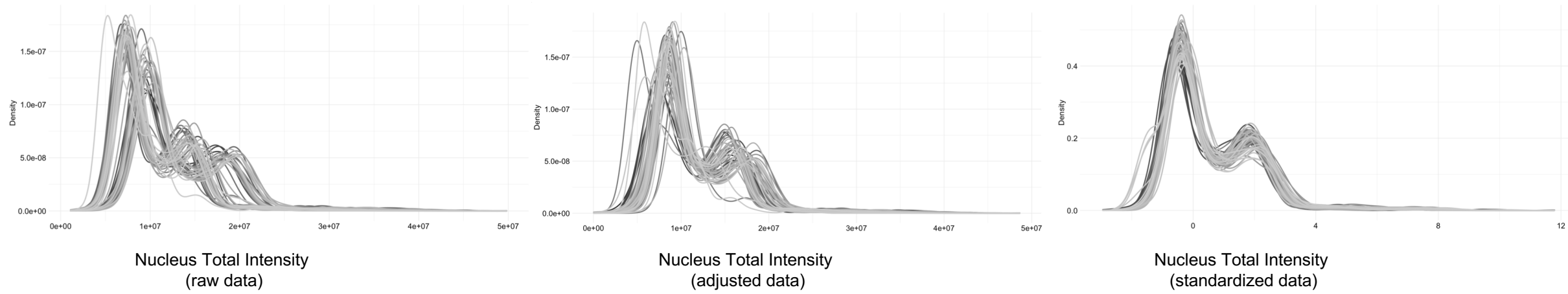
Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.
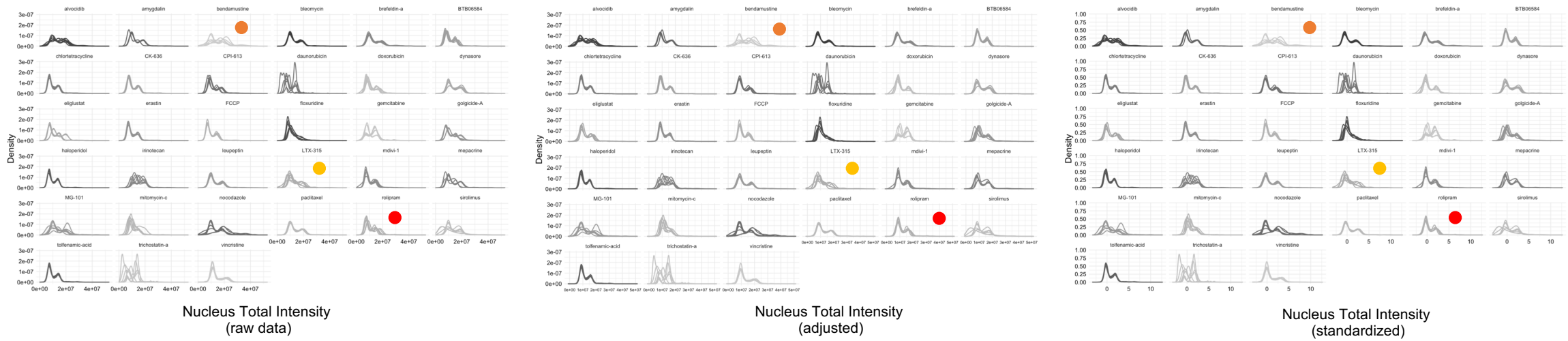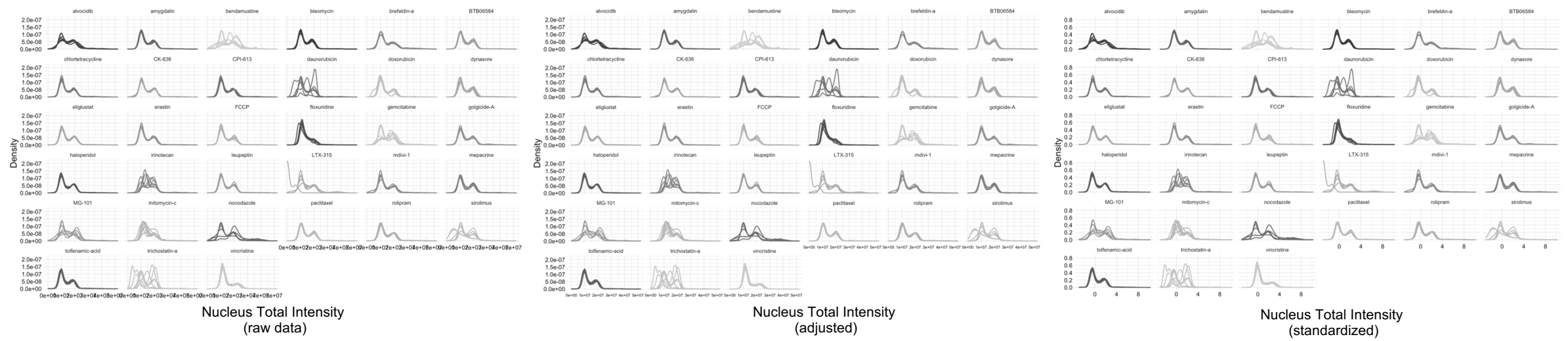
C

DMSO- control distributions



Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.
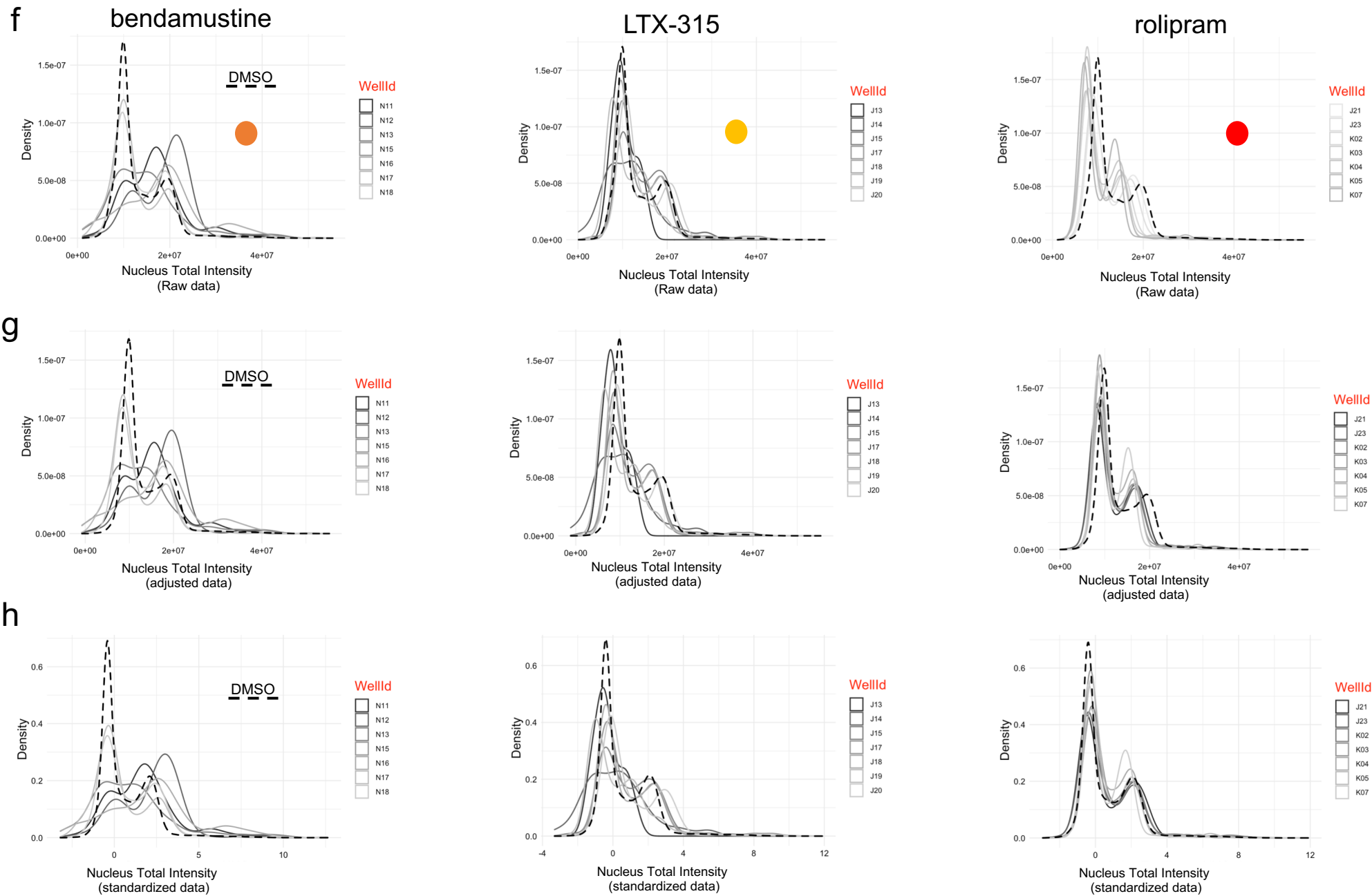
d

Plate 1 replicate 1: Cell cycle feature under different chemical perturbations

Nucleus Total Intensity
(raw data)

Nucleus Total Intensity
(adjusted)

Nucleus Total Intensity
(standardized)

e

Plate 1 replicate 3: Cell cycle feature under different chemical perturbations

Nucleus Total Intensity
(raw data)

Nucleus Total Intensity
(adjusted)

Nucleus Total Intensity
(standardized)

Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.

Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.

Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.

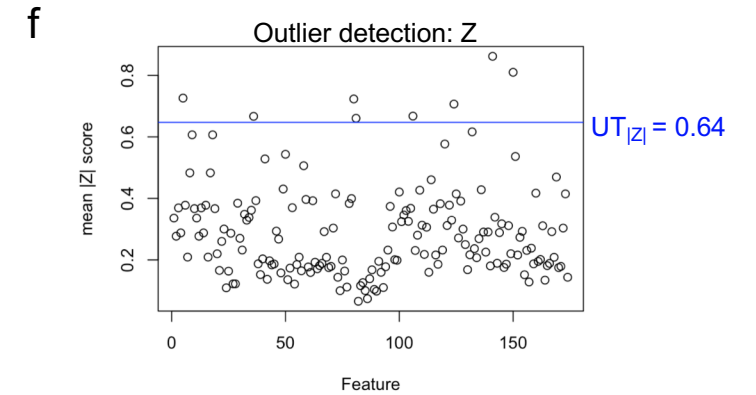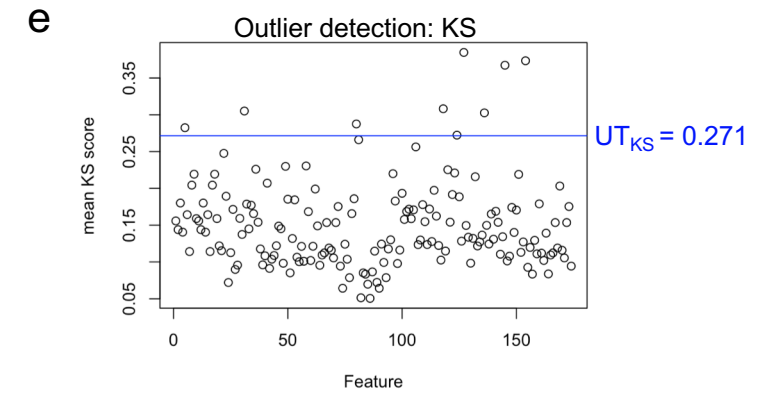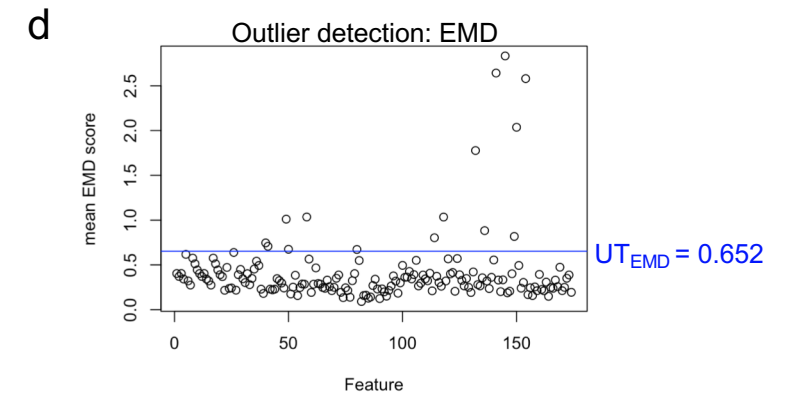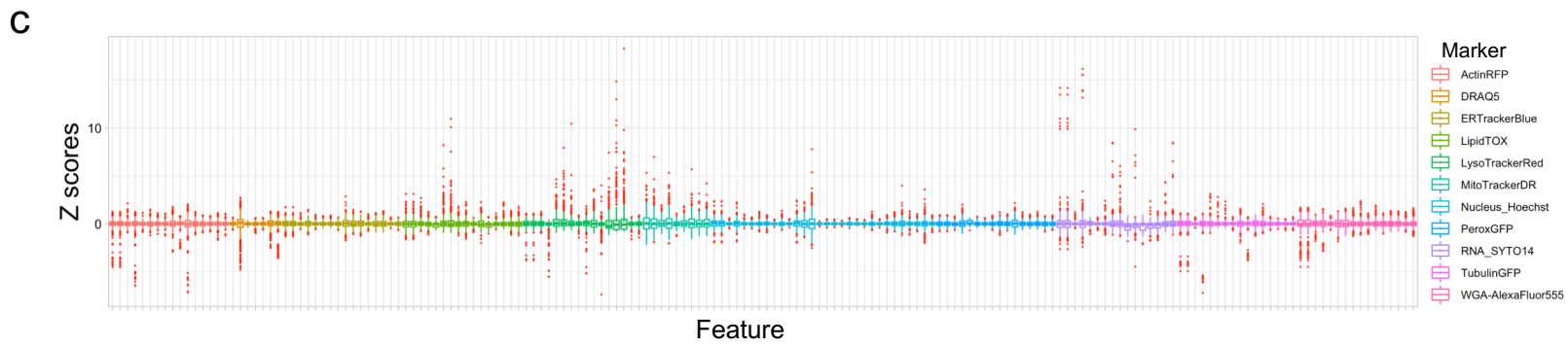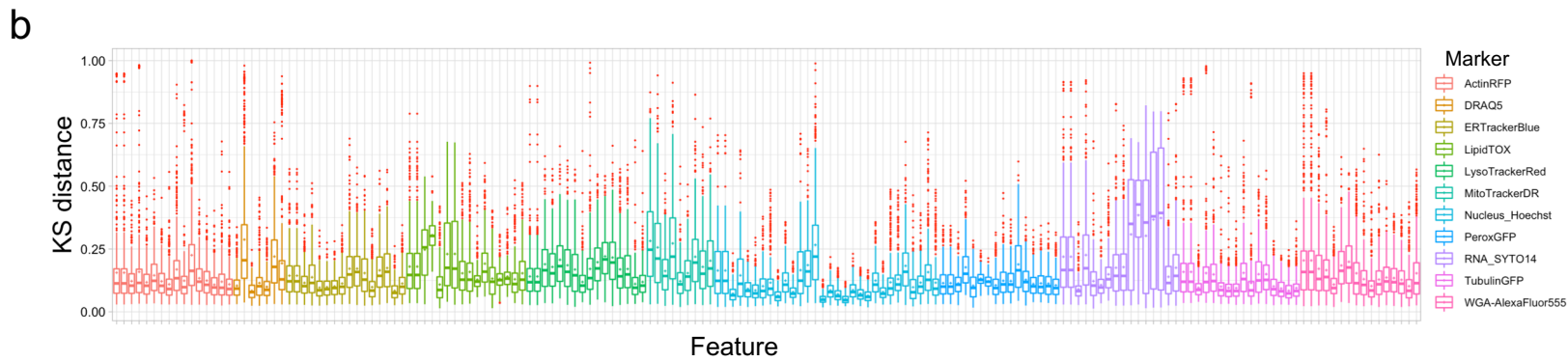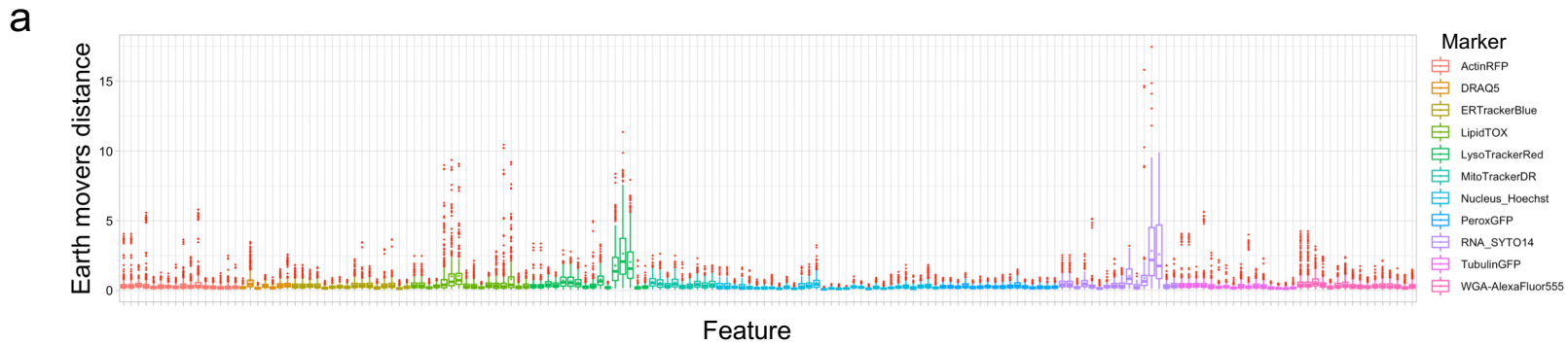**Supplementary Figure 1: Adjustment of positional effects and data standardization across different plates.**

(**a** - **b**) Positional effect detection by two-way ANOVA is applied to control well medians, for all measured features and all plates. Plots show row (**a**) and column (**b**) dependencies for all control wells, ordered by decreasing significance (p-value increasing) for each assay panel. Features related to Intensity (gray) tend to be more sensitive to positional effects than non-Intensity (cyan) features related to area, texture, and shape.

(**c**) Per well distributions of raw, adjusted, and standardized DMSO-control cell populations from plate1 1 rep 1 (see Fig. 3d).

(**d** - **e**) Two replicate plates of cell cycle feature in response to 231 unique treatment conditions including 33 compounds at seven different concentrations. Three compounds bendamustine, LTX-315, and rolipram color labeled as orange (column 1), yellow (column 2) and red (column 3) demonstrate the data correction effects on cell feature distributions (see **f** - **h**).

(**i** - **k**) Correction for row and column positional effects and standardization of plate-to-plate variation for representative features (by row): RNA, Nucleus Area, Nucleus Texture, Mitochondria, Plasma membranes and Golgi (PMG).
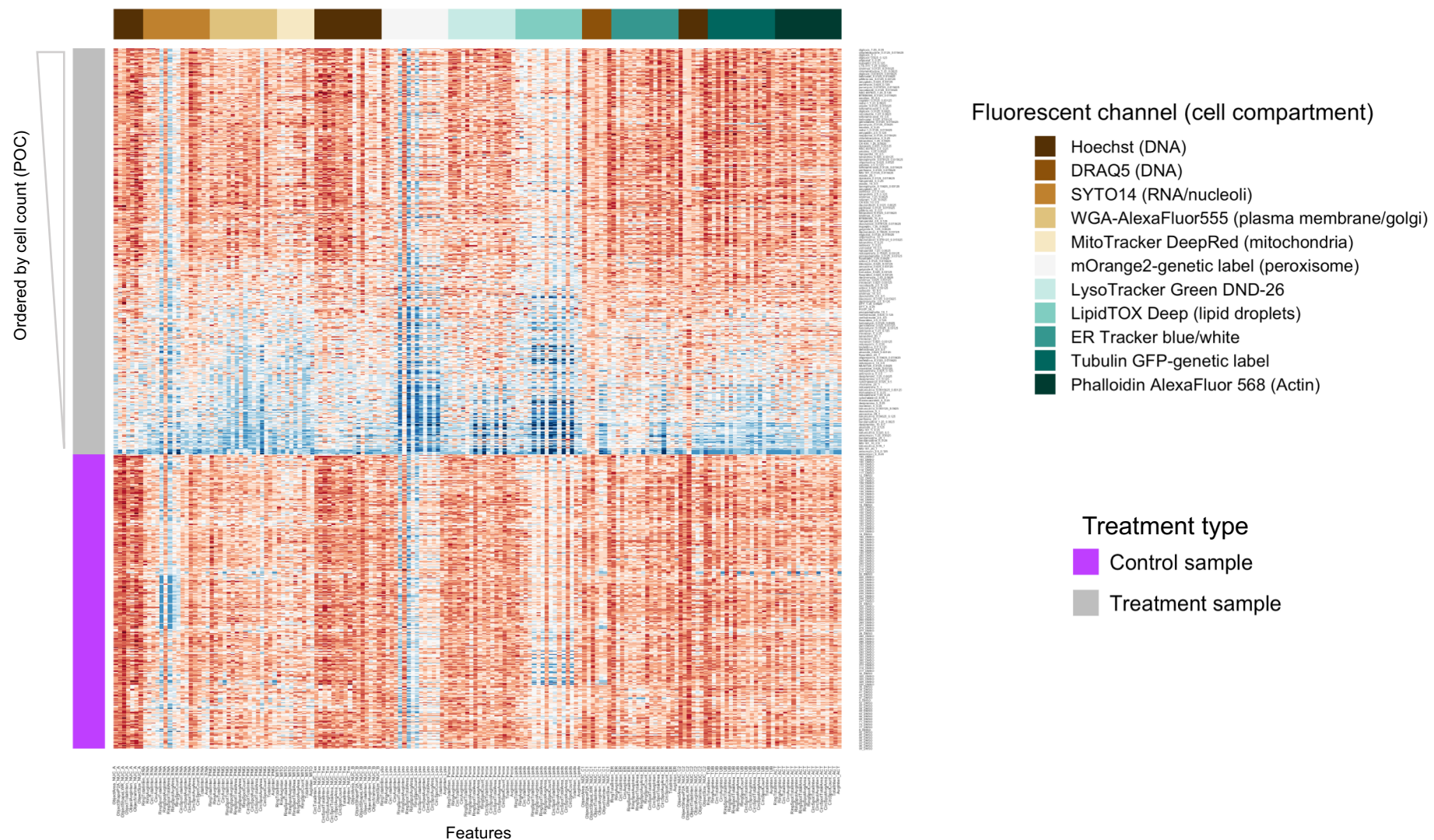
Supplementary Figure 2: Statistical distance measurement among replicates.

**Supplementary Figure 2. Statistical distance measurement among replicates.**

(**a** – **c**) Per-feature distributions for each of the three statistical metrics, color coded by cytological marker. (**a**) EMD scores are strictly positive and highly sensitive to feature distribution differences; (**b**) KS scores are bounded between 0 and 1; (**c**) robust Z-scores account for both positive and negative differences. EMD and robust Z-scores identify RNA (SYTO14), Lysosome (LysoTracker Green) channels (SYTO 14 and LysoTracker Green), and Lipid (LipidTOX) features with extreme outliers, whereas KS has difficulty discriminating noisy features.

(**d** – **f**) Outlier detection: Using per-feature averages and interquartile range (IQR) outlier detection (Upper threshold value (UT) = 1.5 × IQR + upper quartile, blue line) to identify features which fall too far from the expected range of values we find that EMD score (**d**) detects and separates outliers from the group more efficiently than KS (**e**) and robust Z-scores (**f**).

**Supplementary Figure 3: Full feature cytological profile.**
Heatmap summary of EMD profiles for all 455 treatments (grey rows); including 65 compounds at 7 concentrations each, and 330 DMSO-control samples (purple rows). The full profile is log transformed and features are min-max scaled to the range [0, 1]. Treatment profiles (grey rows) are further sorted by their cell count as percent of control (POC).

**Supplementary Figure 4: Full feature EMD fingerprints of individual control samples.**
(**a**) Radial plot of EMD scores spanning 174 measured features among individual controls (gray lines) and the median EMD score of all controls (black dashed line).
(**b**) Radial plot of residual EMD scores of individual controls (gray lines) relative to the null control line (black dashed line). Residual score is defined as the difference between the score of the individual control and the median of all controls. Residual fingerprints naturally fluctuate around median zero (black dashed line), here the values have been offset by 0.5 to expand the plot for better visualization.

Supplementary Figure 5: Summary of phenotypes for 65 compounds.

Supplementary Figure 5: Summary of phenotypes for 65 compounds.

Supplementary Figure 5: Summary of phenotypes for 65 compounds.

**j** Active and cytotoxic compounds

**Supplementary Figure 5: Summary of phenotypes for 65 compounds.**

**Supplementary Figure 5: Summary of phenotypes for 65 compounds.**

(**a** – **c**) 16 compounds elicited little bioactivity in U2OS cells. (**a**) Phenotypic trajectories of the "low stress" group. Profiles for some compounds at higher concentrations are not well separated from controls in the first 3 UMAP dimensions. The 16 compounds showed little or no effect on cell counts (**b**) and cell cycle distributions (**c**) across concentration gradients. In (**b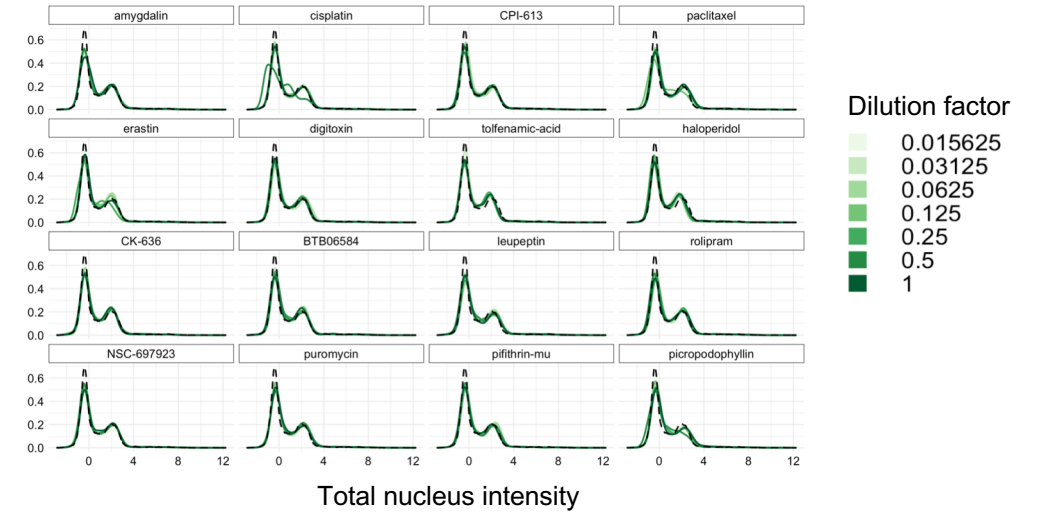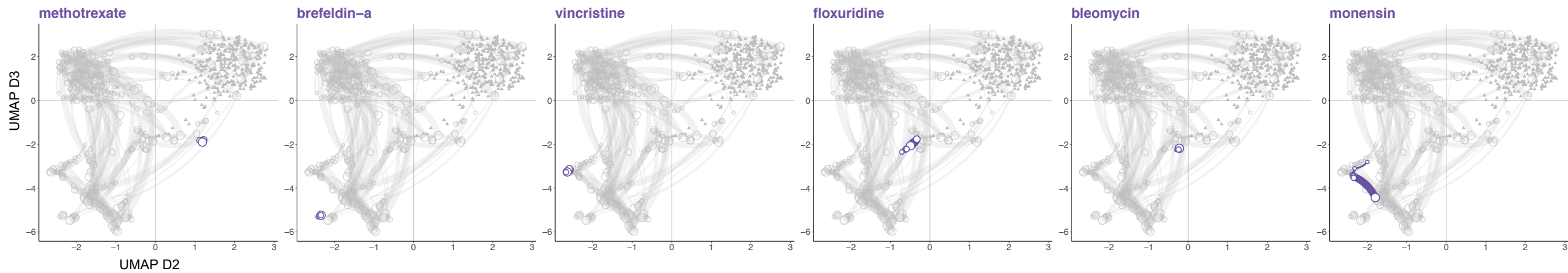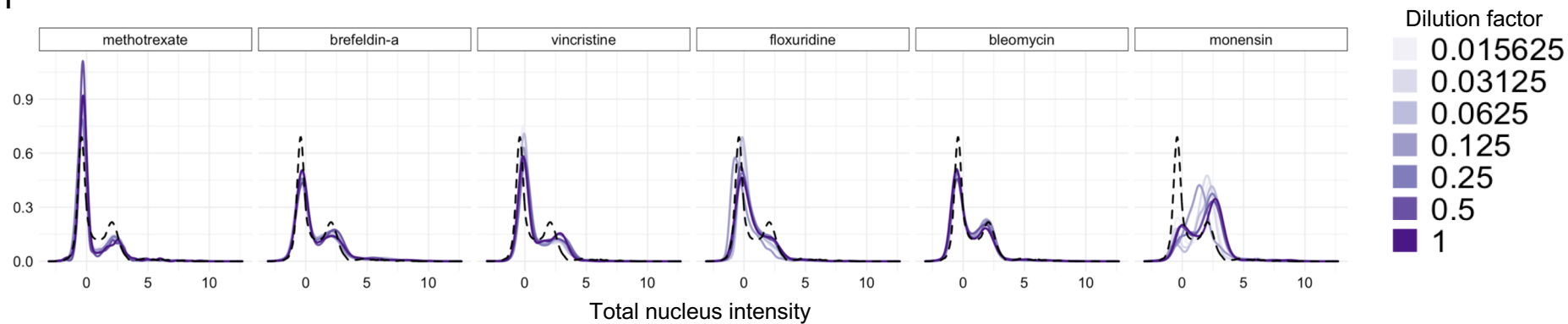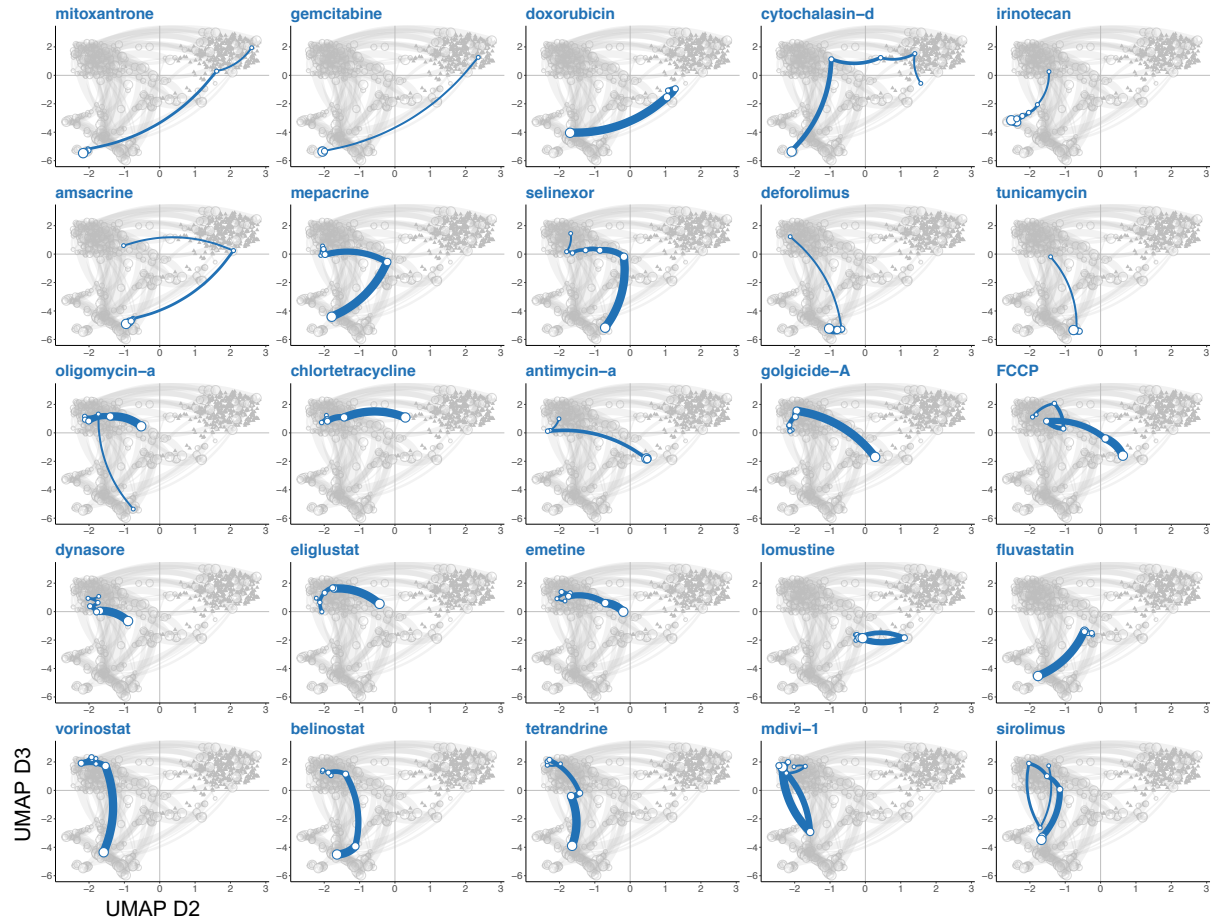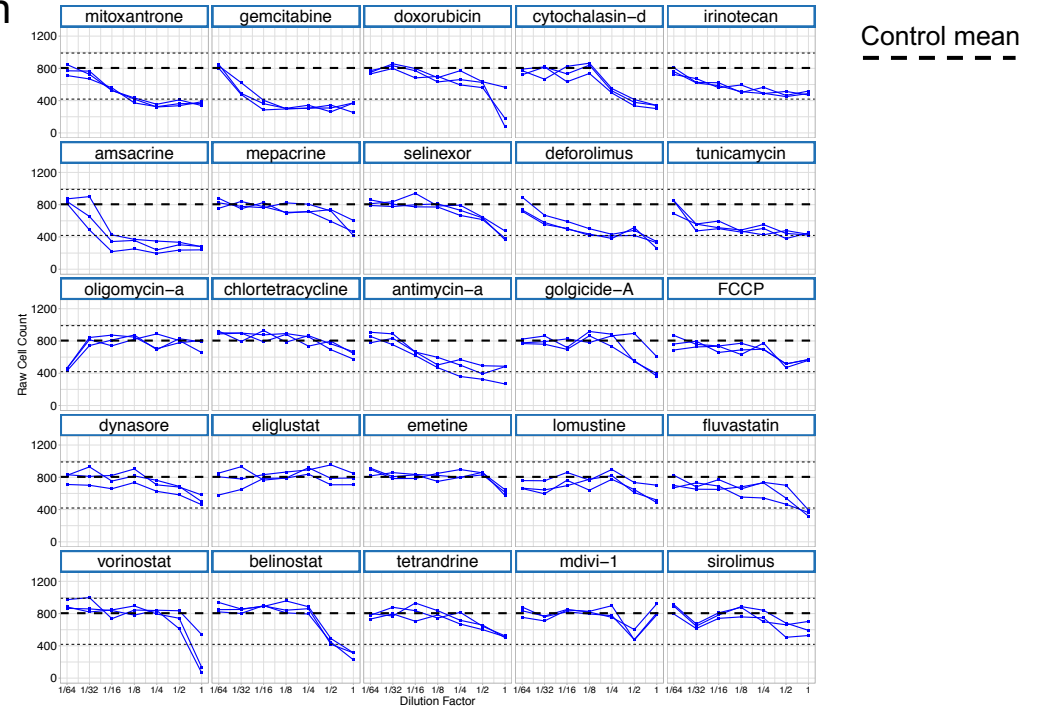**), each curve represents a replicate and three dashed lines represent max, min, and mean of control cell counts. In (**c**), cell cycle is measured by total nuclear intensity; the global control is shown as a dashed line.

(**d** – **f**) 6 compounds with diverse MOA induce phenotypic activity largely independent of dosage/concentration. (**d**) Phenotypic trajectories of the "active (dose-insensitive)" group. These compounds show a distinctive phenotype from the control and low stress group, and the phenotypic response varies little with dosage. These compounds show mildly reduced cell counts in comparison to the control (**e**) and elicit differing effects on the cell cycle (**f**). Visual cues for (**e**, **f**) are as in (**b**, **c**).

(**g** – **i**) 25 compounds elicit dose-dependent phenotypes and are characterized as the "active (dose-responsive)" group. (**g**) Phenotypic trajectories of these compounds travel from the "low stress" region to other areas in the UMAP. Some, but not all, compounds show concentration-dependent effects on cell counts (**h**) and cell cycle responses (**i**). Visual cues for (**h**, **i**) are as in (**b**, **c**).

(**j** – **l**) 18 compounds show a dose-dependent phenotypic change and become toxic at high concentrations. (**j**) Phenotypic trajectories of the "active and cytotoxic" group. These compounds show strong decreases in cell counts with increasing concentration (**k**) and elicit diverse cell cycle responses (**l**). Due to extreme cytotoxicity, cell cycle distributions are not available for some compounds at higher concentrations. Visual cues for (**k**, **l**) are as in (**b**, **c**).

Supplementary Figure 6: Divergent phenotypes.

Supplementary Figure 6: Divergent phenotypes.

**Supplementary Figure 6: Divergent phenotypes.**
Phenotypic fingerprint of both Nocodazole and Irinotecan superimposed at their (**a**) lowest and (**b**) highest concentrations. (**c**) Cell count and (**d**) cell cycle dose response reveal increased levels of stress across the concentration gradient, with diverse cell cycle effects. Cell count is measured from three replicate plates, dashed lines represent the max, min and mean of the DMSO control. (**e**) Raw biological images of cells treated with 20uM of Nocodazole (top row), 5uM of Irinotecan (middle row), and DMSO (bottom row). Columns display a selection of channels imaged in the multi-panel assay. Scale bar: 20 μm. (**f**) Cell feature distributions display diverse phenotypic responses to Nocodazole and (**g**) Irinotecan.

a

LysoTracker (green)    Peroxisome (RFP)

b

features

features

c

$$\frac{\sigma_T^2}{\sigma_C^2} < 2$$

Relative Feature Variance

Supplementary Figure 7. Feature reduction.

**Supplementary Figure 7. Feature reduction.**
(**a**) The lysosome and peroxisome reporters are highly correlated due to overlap in their fluorescent emission spectra arising from the weak lysosomal staining by LysoTracker Green. The right-hand panel shows the linear correlation between the parameters of these two channels. (**b**) Linear correlation analysis of all remaining features. For any pair of highly correlated features (correlation coefficient > 0.9), the feature showing a larger mean correlation with all other features is removed. (**c**) Relative feature activity: per-feature variance is calculated for both control and treatment samples. Features whose variance among treatments is less than twice the variance of the control are deemed "inactive".

**Supplementary Figure 8: Comparison of cytological profiles**
(**a - b**) EMD profiles without adjusting for positional effects and plate to plate variation. (**a**) Heatmap visualization of similarity by hierarchical clustering shows several treatment groups (grey) clustered with the controls (purple). Dimension reduction by UMAP of (**b**) unprocessed EMD profiles and (**c**) processed full feature EMD profiles strongly separates the brefeldin-a cluster from other treatments. (**c**) UMAP of processed full feature profile also fails to separate the low stress cluster from the control when compared to the (**d**) UMAP of the processed and optimally reduced feature profile (described in Fig. 7a).

**Supplementary Figure 9: Radial plot fingerprints of 65 compounds**
Radial fingerprints of residual EMD scores for 65 compounds at multiple concentrations, relative to the control median (black dotted line). Labels for the reduced 69 feature set are color-coded by corresponding cellular marker. Per compound legends include compound name, concentration and dilution factor (*compound_concentration_dilution factor*). Lowest and highest concentrations are colored yellow and dark purple respectively.

**Compound: CPI−613**

Legend:
- DMSO_median
- CPI−613_20_1
- CPI−613_10_0.5
- CPI−613_5_0.25
- CPI−613_2.5_0.125
- CPI−613_1.25_0.0625
- CPI−613_0.625_0.03125
- CPI−613_0.3125_0.015625

**Compound: digitoxin**

Legend:
- DMSO_median
- digitoxin_5_1
- digitoxin_2.5_0.5
- digitoxin_1.25_0.25
- digitoxin_0.625_0.125
- digitoxin_0.3125_0.0625
- digitoxin_0.15625_0.03125
- digitoxin_0.078125_0.015625

**Compound: cisplatin**

Legend:
- DMSO_median
- cisplatin_10_1
- cisplatin_5_0.5
- cisplatin_2.5_0.25
- cisplatin_1.25_0.125
- cisplatin_0.625_0.0625
- cisplatin_0.3125_0.03125
- cisplatin_0.15625_0.015625

**Compound: 9−aminoacridine**

Legend:
- DMSO_median
- 9−aminoacridine_20_1
- 9−aminoacridine_10_0.5
- 9−aminoacridine_5_0.25
- 9−aminoacridine_2.5_0.125
- 9−aminoacridine_1.25_0.0625
- 9−aminoacridine_0.625_0.03125
- 9−aminoacridine_0.3125_0.015625

**Compound: chlortetracycline**

Legend:
- DMSO_median
- chlortetracycline_20_1
- chlortetracycline_10_0.5
- chlortetracycline_5_0.25
- chlortetracycline_2.5_0.125
- chlortetracycline_1.25_0.0625
- chlortetracycline_0.625_0.03125
- chlortetracycline_0.3125_0.015625

**Compound: vorinostat**

Legend:
- DMSO_median
- vorinostat_20_1
- vorinostat_10_0.5
- vorinostat_5_0.25
- vorinostat_2.5_0.125
- vorinostat_1.25_0.0625
- vorinostat_0.625_0.03125
- vorinostat_0.3125_0.015625

Compound: BTB06584

Legend:
- DMSO_median
- BTB06584_20_1
- BTB06584_10_0.5
- BTB06584_5_0.25
- BTB06584_2.5_0.125
- BTB06584_1.25_0.0625
- BTB06584_0.625_0.03125
- BTB06584_0.3125_0.015625

**Compound: pifithrin−mu**

Legend:
- DMSO_median
- pifithrin−mu_10_1
- pifithrin−mu_5_0.5
- pifithrin−mu_2.5_0.25
- pifithrin−mu_1.25_0.125
- pifithrin−mu_0.625_0.0625
- pifithrin−mu_0.3125_0.03125
- pifithrin−mu_0.15625_0.015625

Compound: belinostat

**Compound: eliglustat**

Legend:
- DMSO_median
- eliglustat_20_1
- eliglustat_10_0.5
- eliglustat_5_0.25
- eliglustat_2.5_0.125
- eliglustat_1.25_0.0625
- eliglustat_0.625_0.03125
- eliglustat_0.3125_0.015625

# Compound: tolfenamic−acid



Legend:
- DMSO_median
- tolfenamic−acid_20_1
- tolfenamic−acid_10_0.5
- tolfenamic−acid_5_0.25
- tolfenamic−acid_2.5_0.125
- tolfenamic−acid_1.25_0.0625
- tolfenamic−acid_0.625_0.03125
- tolfenamic−acid_0.3125_0.015625

**Compound: leupeptin**

Legend:
- DMSO_median
- leupeptin_20_1
- leupeptin_10_0.5
- leupeptin_5_0.25
- leupeptin_2.5_0.125
- leupeptin_1.25_0.0625
- leupeptin_0.625_0.03125
- leupeptin_0.3125_0.015625

**Compound: amygdalin**

Legend:
- DMSO_median (dashed)
- amygdalin_20_1
- amygdalin_10_0.5
- amygdalin_5_0.25
- amygdalin_2.5_0.125
- amygdalin_1.25_0.0625
- amygdalin_0.625_0.03125
- amygdalin_0.3125_0.015625

**Compound: LTX−315**

Legend:
- DMSO_median
- LTX−315_20_1
- LTX−315_10_0.5
- LTX−315_5_0.25
- LTX−315_2.5_0.125
- LTX−315_1.25_0.0625
- LTX−315_0.625_0.03125
- LTX−315_0.3125_0.015625

**Compound: erastin**

Legend:
- DMSO_median
- erastin_20_1
- erastin_10_0.5
- erastin_5_0.25
- erastin_2.5_0.125
- erastin_1.25_0.0625
- erastin_0.625_0.03125
- erastin_0.3125_0.015625

**Compound: emetine**

Legend:
- DMSO_median
- emetine_20_1
- emetine_10_0.5
- emetine_5_0.25
- emetine_2.5_0.125
- emetine_1.25_0.0625
- emetine_0.625_0.03125
- emetine_0.3125_0.015625

**Compound: sirolimus**

Legend:
- DMSO_median
- sirolimus_20_1
- sirolimus_10_0.5
- sirolimus_5_0.25
- sirolimus_2.5_0.125
- sirolimus_1.25_0.0625
- sirolimus_0.625_0.03125
- sirolimus_0.3125_0.015625

**Compound: rolipram**

Legend:
- DMSO_median
- rolipram_20_1
- rolipram_10_0.5
- rolipram_5_0.25
- rolipram_2.5_0.125
- rolipram_1.25_0.0625
- rolipram_0.625_0.03125
- rolipram_0.3125_0.015625

**Compound: perifosine**

Legend:
- DMSO_median
- perifosine_20_1
- perifosine_10_0.5
- perifosine_5_0.25
- perifosine_2.5_0.125
- perifosine_1.25_0.0625
- perifosine_0.625_0.03125
- perifosine_0.3125_0.015625

**Compound: golgicide-A**

Legend:
- DMSO_median
- golgicide-A_20_1
- golgicide-A_10_0.5
- golgicide-A_5_0.25
- golgicide-A_2.5_0.125
- golgicide-A_1.25_0.0625
- golgicide-A_0.625_0.03125
- golgicide-A_0.3125_0.015625

**Compound: puromycin**

Legend:
- DMSO_median
- puromycin_5_1
- puromycin_2.5_0.5
- puromycin_1.25_0.25
- puromycin_0.625_0.125
- puromycin_0.3125_0.0625
- puromycin_0.15625_0.03125
- puromycin_0.078125_0.015625

**Compound: CK−636**

Legend:
- DMSO_median
- CK−636_20_1
- CK−636_10_0.5
- CK−636_5_0.25
- CK−636_2.5_0.125
- CK−636_1.25_0.0625
- CK−636_0.625_0.03125
- CK−636_0.3125_0.015625

**Compound: nocodazole**

Legend:
- DMSO_median
- nocodazole_20_1
- nocodazole_10_0.5
- nocodazole_5_0.25
- nocodazole_2.5_0.125
- nocodazole_1.25_0.0625
- nocodazole_0.625_0.03125
- nocodazole_0.3125_0.015625

# Compound: paclitaxel



Legend:
- DMSO_median
- paclitaxel_20_1
- paclitaxel_10_0.5
- paclitaxel_5_0.25
- paclitaxel_2.5_0.125
- paclitaxel_1.25_0.0625
- paclitaxel_0.625_0.03125
- paclitaxel_0.3125_0.015625

**Compound: NSC−697923**

Legend:
- DMSO_median
- NSC−697923_10_1
- NSC−697923_5_0.5
- NSC−697923_2.5_0.25
- NSC−697923_1.25_0.125
- NSC−697923_0.625_0.0625
- NSC−697923_0.3125_0.03125
- NSC−697923_0.15625_0.015625

**Compound: antimycin−a**

Legend:
- DMSO_median
- antimycin−a_10_1
- antimycin−a_5_0.5
- antimycin−a_2.5_0.25
- antimycin−a_1.25_0.125
- antimycin−a_0.625_0.0625
- antimycin−a_0.3125_0.03125
- antimycin−a_0.15625_0.015625

Compound: oligomycin−a

Legend:
- DMSO_median
- oligomycin−a_10_1
- oligomycin−a_5_0.5
- oligomycin−a_2.5_0.25
- oligomycin−a_1.25_0.125
- oligomycin−a_0.625_0.0625
- oligomycin−a_0.3125_0.03125
- oligomycin−a_0.15625_0.015625

# Compound: mdivi-1



Legend:
- DMSO_median
- mdivi-1_20_1
- mdivi-1_10_0.5
- mdivi-1_5_0.25
- mdivi-1_2.5_0.125
- mdivi-1_1.25_0.0625
- mdivi-1_0.625_0.03125
- mdivi-1_0.3125_0.015625

**Compound: amsacrine**

Legend:
- DMSO_median
- amsacrine_20_1
- amsacrine_10_0.5
- amsacrine_5_0.25
- amsacrine_2.5_0.125
- amsacrine_1.25_0.0625
- amsacrine_0.625_0.03125
- amsacrine_0.3125_0.015625

**Compound: selinexor**

Legend:
- DMSO_median
- selinexor_20_1
- selinexor_10_0.5
- selinexor_5_0.25
- selinexor_2.5_0.125
- selinexor_1.25_0.0625
- selinexor_0.625_0.03125
- selinexor_0.3125_0.015625

**Compound: doxorubicin**

Legend:
- DMSO_median
- doxorubicin_5_1
- doxorubicin_2.5_0.5
- doxorubicin_1.25_0.25
- doxorubicin_0.625_0.125
- doxorubicin_0.3125_0.0625
- doxorubicin_0.15625_0.03125
- doxorubicin_0.078125_0.015625

**Compound: gemcitabine**

Legend:
- DMSO_median
- gemcitabine_20_1
- gemcitabine_10_0.5
- gemcitabine_5_0.25
- gemcitabine_2.5_0.125
- gemcitabine_1.25_0.0625
- gemcitabine_0.625_0.03125
- gemcitabine_0.3125_0.015625

**Compound: lomustine**

Legend:
- DMSO_median
- lomustine_20_1
- lomustine_10_0.5
- lomustine_5_0.25
- lomustine_2.5_0.125
- lomustine_1.25_0.0625
- lomustine_0.625_0.03125
- lomustine_0.3125_0.015625

Compound: haloperidol

DMSO_median
haloperidol_20_1
haloperidol_10_0.5
haloperidol_5_0.25
haloperidol_2.5_0.125
haloperidol_1.25_0.0625
haloperidol_0.625_0.03125
haloperidol_0.3125_0.015625

**Compound: mepacrine**

Legend:
- DMSO_median
- mepacrine_20_1
- mepacrine_10_0.5
- mepacrine_5_0.25
- mepacrine_2.5_0.125
- mepacrine_1.25_0.0625
- mepacrine_0.625_0.03125
- mepacrine_0.3125_0.015625

**Compound: dynasore**

Legend:
- DMSO_median
- dynasore_20_1
- dynasore_10_0.5
- dynasore_5_0.25
- dynasore_2.5_0.125
- dynasore_1.25_0.0625
- dynasore_0.625_0.03125
- dynasore_0.3125_0.015625

Compound: tetrandrine

Legend:
- DMSO_median
- tetrandrine_20_1
- tetrandrine_10_0.5
- tetrandrine_5_0.25
- tetrandrine_2.5_0.125
- tetrandrine_1.25_0.0625
- tetrandrine_0.625_0.03125
- tetrandrine_0.3125_0.015625

# Compound: cytochalasin-d



DMSO_median
cytochalasin-d_0.05_1
cytochalasin-d_0.025_0.5
cytochalasin-d_0.0125_0.25
cytochalasin-d_0.00625_0.125
cytochalasin-d_0.003125_0.0625
cytochalasin-d_0.0015625_0.03125
cytochalasin-d_0.000781_0.015625

**Compound: tanespimycin**

Legend:
- DMSO_median
- tanespimycin_5_1
- tanespimycin_2.5_0.5
- tanespimycin_1.25_0.25
- tanespimycin_0.625_0.125
- tanespimycin_0.3125_0.0625
- tanespimycin_0.15625_0.03125
- tanespimycin_0.078125_0.015625

Compound: MG−101

Legend:
- DMSO_median
- MG−101_20_1
- MG−101_10_0.5
- MG−101_5_0.25
- MG−101_2.5_0.125
- MG−101_1.25_0.0625
- MG−101_0.625_0.03125
- MG−101_0.3125_0.015625

**Compound: tunicamycin**

Legend:
- DMSO_median
- tunicamycin_5_1
- tunicamycin_2.5_0.5
- tunicamycin_1.25_0.25
- tunicamycin_0.625_0.125
- tunicamycin_0.3125_0.0625
- tunicamycin_0.15625_0.03125
- tunicamycin_0.078125_0.015625

**Compound: picropodophyllin**

Legend:
- DMSO_median
- picropodophyllin_10_1
- picropodophyllin_5_0.5
- picropodophyllin_2.5_0.25
- picropodophyllin_1.25_0.125
- picropodophyllin_0.625_0.0625
- picropodophyllin_0.3125_0.03125
- picropodophyllin_0.15625_0.015625

**Compound: deforolimus**

Legend:
- DMSO_median
- deforolimus_5_1
- deforolimus_2.5_0.5
- deforolimus_1.25_0.25
- deforolimus_0.625_0.125
- deforolimus_0.3125_0.0625
- deforolimus_0.15625_0.03125
- deforolimus_0.078125_0.015625

**Compound: FCCP**

Legend:
- DMSO_median
- FCCP_20_1
- FCCP_10_0.5
- FCCP_5_0.25
- FCCP_2.5_0.125
- FCCP_1.25_0.0625
- FCCP_0.625_0.03125
- FCCP_0.3125_0.015625

**Compound: daunorubicin**

Legend:
- DMSO_median
- daunorubicin_5_1
- daunorubicin_2.5_0.5
- daunorubicin_1.25_0.25
- daunorubicin_0.625_0.125
- daunorubicin_0.3125_0.0625
- daunorubicin_0.15625_0.03125
- daunorubicin_0.078125_0.015625

Compound: mitoxantrone

DMSO_median
mitoxantrone_5_1
mitoxantrone_2.5_0.5
mitoxantrone_1.25_0.25
mitoxantrone_0.625_0.125
mitoxantrone_0.3125_0.0625
mitoxantrone_0.15625_0.03125
mitoxantrone_0.078125_0.015625

**Compound: irinotecan**

Legend:
- DMSO_median
- irinotecan_20_1
- irinotecan_10_0.5
- irinotecan_5_0.25
- irinotecan_2.5_0.125
- irinotecan_1.25_0.0625
- irinotecan_0.625_0.03125
- irinotecan_0.3125_0.015625

**Compound: bendamustine**

Legend:
- DMSO_median
- bendamustine_20_1
- bendamustine_10_0.5
- bendamustine_5_0.25
- bendamustine_2.5_0.125
- bendamustine_1.25_0.0625
- bendamustine_0.625_0.03125
- bendamustine_0.3125_0.015625

Compound: bleomycin

**Compound: dactinomycin**

Legend:
- DMSO_median
- dactinomycin_20_1
- dactinomycin_10_0.5
- dactinomycin_5_0.25
- dactinomycin_2.5_0.125
- dactinomycin_1.25_0.0625
- dactinomycin_0.625_0.03125
- dactinomycin_0.3125_0.015625

**Compound: fluvastatin**

Legend:
- DMSO_median
- fluvastatin_20_1
- fluvastatin_10_0.5
- fluvastatin_5_0.25
- fluvastatin_2.5_0.125
- fluvastatin_1.25_0.0625
- fluvastatin_0.625_0.03125
- fluvastatin_0.3125_0.015625

**Compound: trichostatin-a**

Legend:
- DMSO_median
- trichostatin-a_20_1
- trichostatin-a_10_0.5
- trichostatin-a_5_0.25
- trichostatin-a_2.5_0.125
- trichostatin-a_1.25_0.0625
- trichostatin-a_0.625_0.03125
- trichostatin-a_0.3125_0.015625

**Compound: DTT**

Legend:
- DMSO_median
- DTT_20_1
- DTT_10_0.5
- DTT_5_0.25
- DTT_2.5_0.125
- DTT_1.25_0.0625
- DTT_0.625_0.03125
- DTT_0.3125_0.015625

**Compound: sirtinol**

Legend:
- DMSO_median
- sirtinol_20_1
- sirtinol_10_0.5
- sirtinol_5_0.25
- sirtinol_2.5_0.125
- sirtinol_1.25_0.0625
- sirtinol_0.625_0.03125
- sirtinol_0.3125_0.015625

**Compound: desipramine**

Legend:
- DMSO_median
- desipramine_20_1
- desipramine_10_0.5
- desipramine_5_0.25
- desipramine_2.5_0.125
- desipramine_1.25_0.0625
- desipramine_0.625_0.03125
- desipramine_0.3125_0.015625

**Compound: floxuridine**

Legend:
- DMSO_median
- floxuridine_20_1
- floxuridine_10_0.5
- floxuridine_5_0.25
- floxuridine_2.5_0.125
- floxuridine_1.25_0.0625
- floxuridine_0.625_0.03125
- floxuridine_0.3125_0.015625

**Compound: alvocidib**

Legend:
- DMSO_median
- alvocidib_20_1
- alvocidib_10_0.5
- alvocidib_5_0.25
- alvocidib_2.5_0.125
- alvocidib_1.25_0.0625
- alvocidib_0.625_0.03125
- alvocidib_0.3125_0.015625

**Compound: anisomycin**

Legend:
- DMSO_median
- anisomycin_20_1
- anisomycin_10_0.5
- anisomycin_5_0.25
- anisomycin_2.5_0.125
- anisomycin_1.25_0.0625
- anisomycin_0.625_0.03125
- anisomycin_0.3125_0.015625

Compound: methotrexate

**Compound: monensin**

Legend:
- DMSO_median
- monensin_20_1
- monensin_10_0.5
- monensin_5_0.25
- monensin_2.5_0.125
- monensin_1.25_0.0625
- monensin_0.625_0.03125
- monensin_0.3125_0.015625

**Compound: mitomycin−c**

Legend:
- DMSO_median
- mitomycin−c_20_1
- mitomycin−c_10_0.5
- mitomycin−c_5_0.25
- mitomycin−c_2.5_0.125
- mitomycin−c_1.25_0.0625
- mitomycin−c_0.625_0.03125
- mitomycin−c_0.3125_0.015625

**Compound: MLN0128**

Legend:
- DMSO_median
- MLN0128_5_1
- MLN0128_2.5_0.5
- MLN0128_1.25_0.25
- MLN0128_0.625_0.125
- MLN0128_0.3125_0.0625
- MLN0128_0.15625_0.03125
- MLN0128_0.078125_0.015625
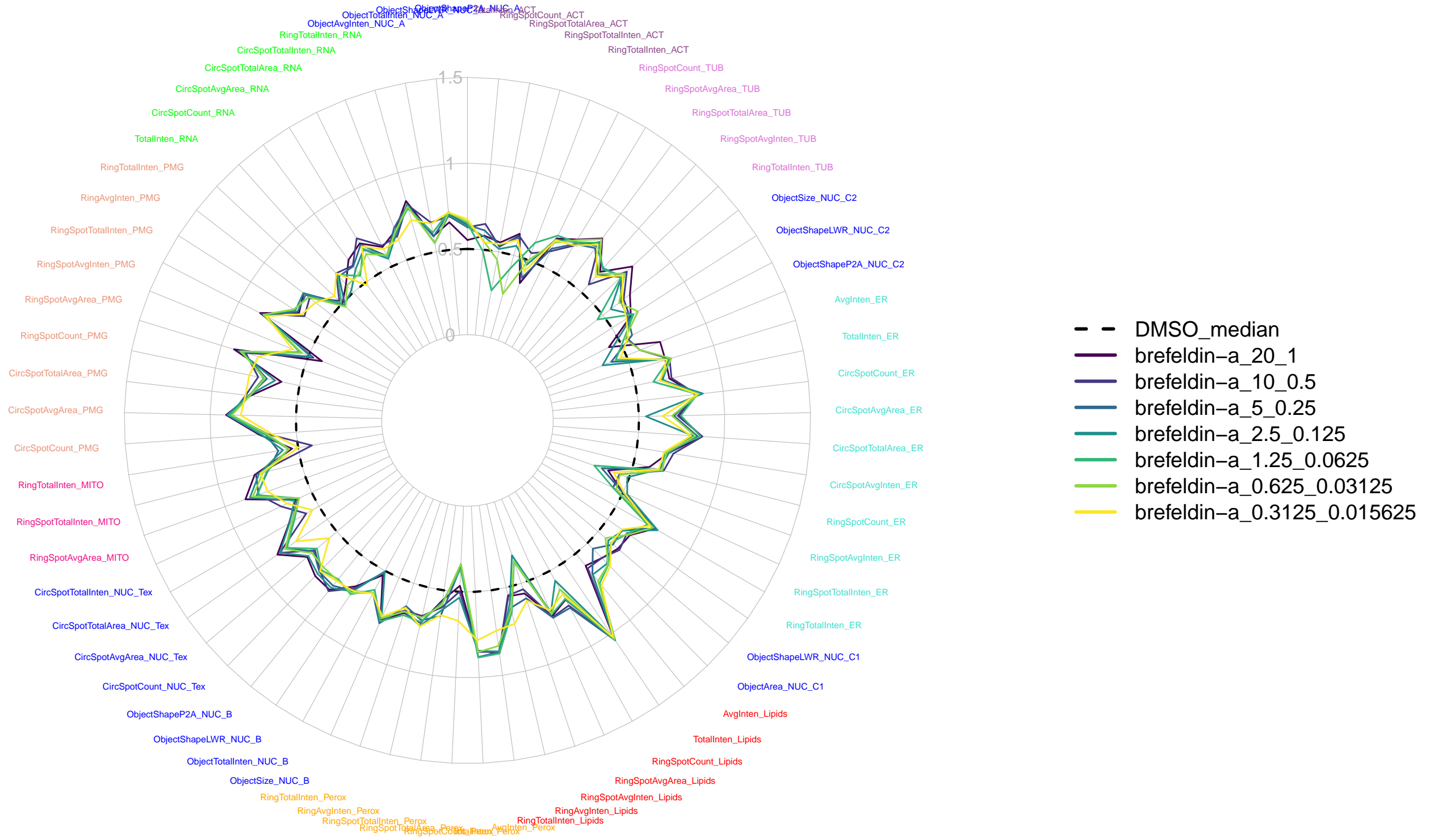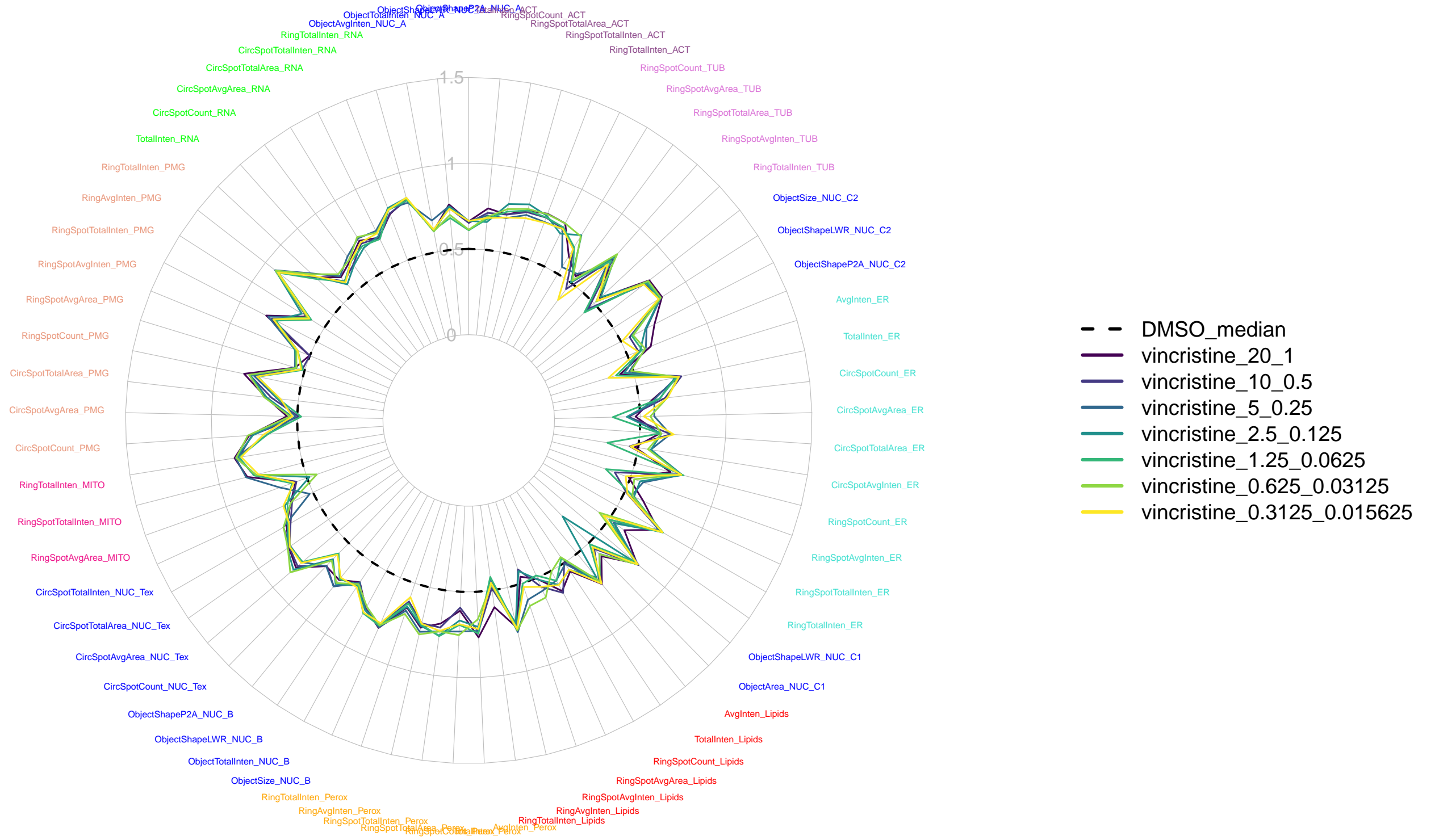
Compound: brefeldin–a

Compound: vincristine

Compound: latrunculin−a