

SUPPORTING INFORMATION

PROTEOMICS

Research article

PROTEOMICS-BASED SCORING OF CELLULAR RESPONSE TO STIMULI FOR IMPROVED CHARACTERIZATION OF SIGNALING PATHWAY ACTIVITY

Elizaveta M. Kazakova^{1,2}, Elizaveta M. Solovyeva², Lev I. Levitsky², Julia A. Bubis², Daria D. Emekeeva^{1,2}, Anastasia A. Antonets³, Alexey A. Nazarov³, Mikhail V. Gorshkov², Irina A. Tarasova²

¹ Moscow Institute of Physics and Technology (National Research University), Dolgoprudny, 141701, Russia

² V.L. Talrose Institute for Energy Problems of Chemical Physics, Federal Research Center of Chemical Physics, Russian Academy of Sciences, Moscow, 119334, Russia

³ M. V. Lomonosov Moscow State University, Department of Chemistry, Leninskie Gory 1/3, 119991 Moscow, Russia

Correspondence: 38 Leninsky pr-t, bld.2, 119334 Moscow Russia, e-mail: iatarasova@yandex.ru

Table of content

Tables S1-S7 are provided as separate .xlsx files at <https://github.com/kazakova/Metrics>.

Table S1. Summary on missing value percentages across datasets considered in this study: label-free quantification at protein level with NSAF [DOI: 10.1021/pr060161n].

Table S2a. Results of statistical analysis of IFN-induced proteome changes: imputation with minimal NSAF [DOI: 10.1021/pr060161n].

Table S2b. Results of statistical analysis of IFN-induced proteome changes: imputation with k-nearest neighbors machine learning.

Table S3. Summary of virus assay-based and proteomics scores calculated for ranking the functionality of IFN-dependent antiviral mechanisms in cancer and normal cells. Table S3 contains descriptions of proteomic datasets used in this study.

Table S4. Results of GO analyses for differentially regulated proteins selected using different workflows, sorted by GO_score: a) imputation with minimal NSAF and satisfying $fdr < 0.05$ and $|\log_{2}FC| > 0.585$; b) imputation with k-nearest neighbors machine learning and satisfying $fdr < 0.05$ and $|\log_{2}FC| > 0.585$; c) imputation by minimal detected NSAFs with quartile-based dynamic selection; d) imputation with k-nearest neighbors machine learning and quartile-based dynamic selection.

Table S5. Results of Shapiro-Wilk testing for normal data distribution using two imputation strategies: a) the minimal detected NSAFs, b) the kNN-assisted imputation.

Table S6. IPAS [doi:10.1101/2020.10.19.345629v2] run against REACTOME has identified the IFN-gamma signaling on the top of the enriched processes in the IFN γ treated MRC5 cells.

Table S7. Correlation of the enriched pathways with topotecan concentration.

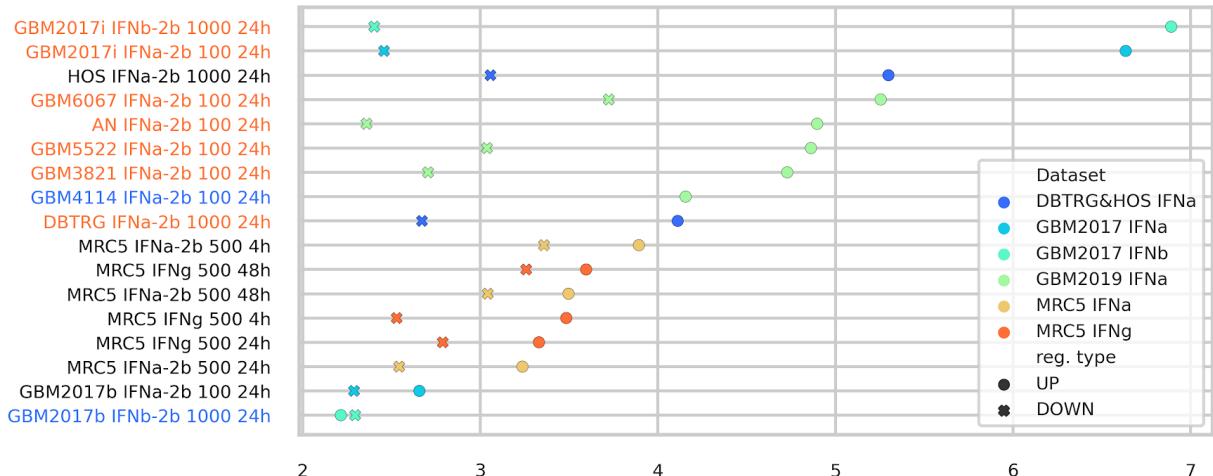
Figure S1. Summary on performance of different score equations, missing value imputation strategies (kNN or minimal NSAF) and DRF selection (static, semi dynamic or dynamic).

Figure S2. Top enriched biological processes ranked by GO score.

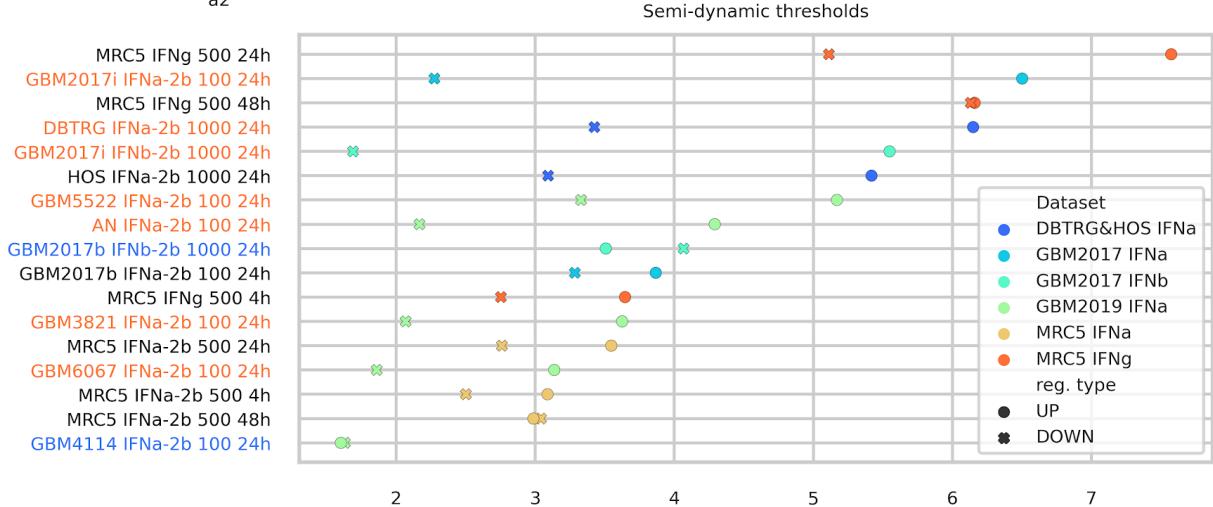
Figure S3. Heatmap for the proteins involved in biological processes correlated.

Euclidean distance

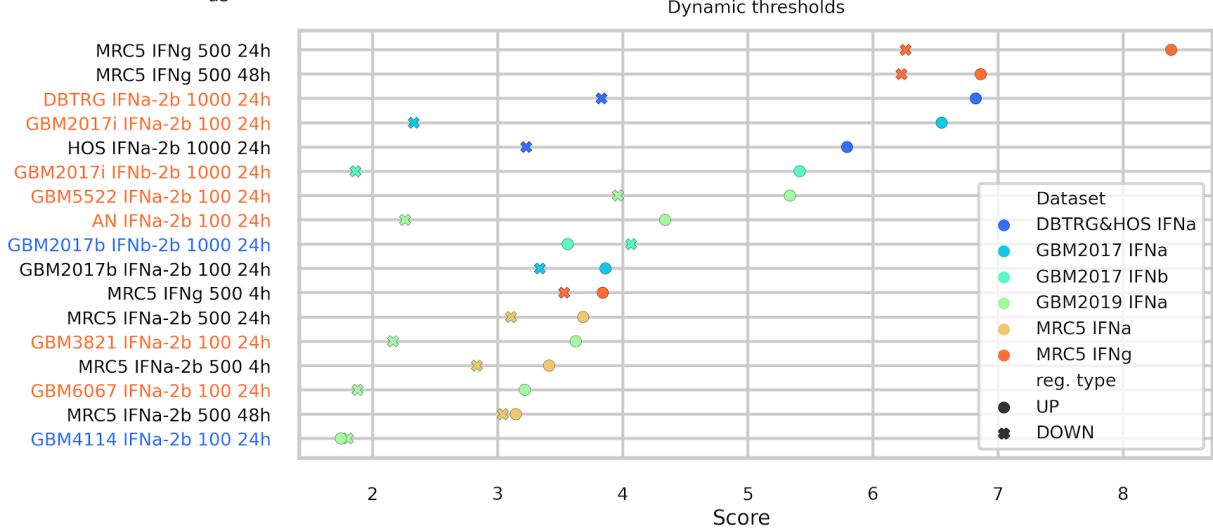
a1



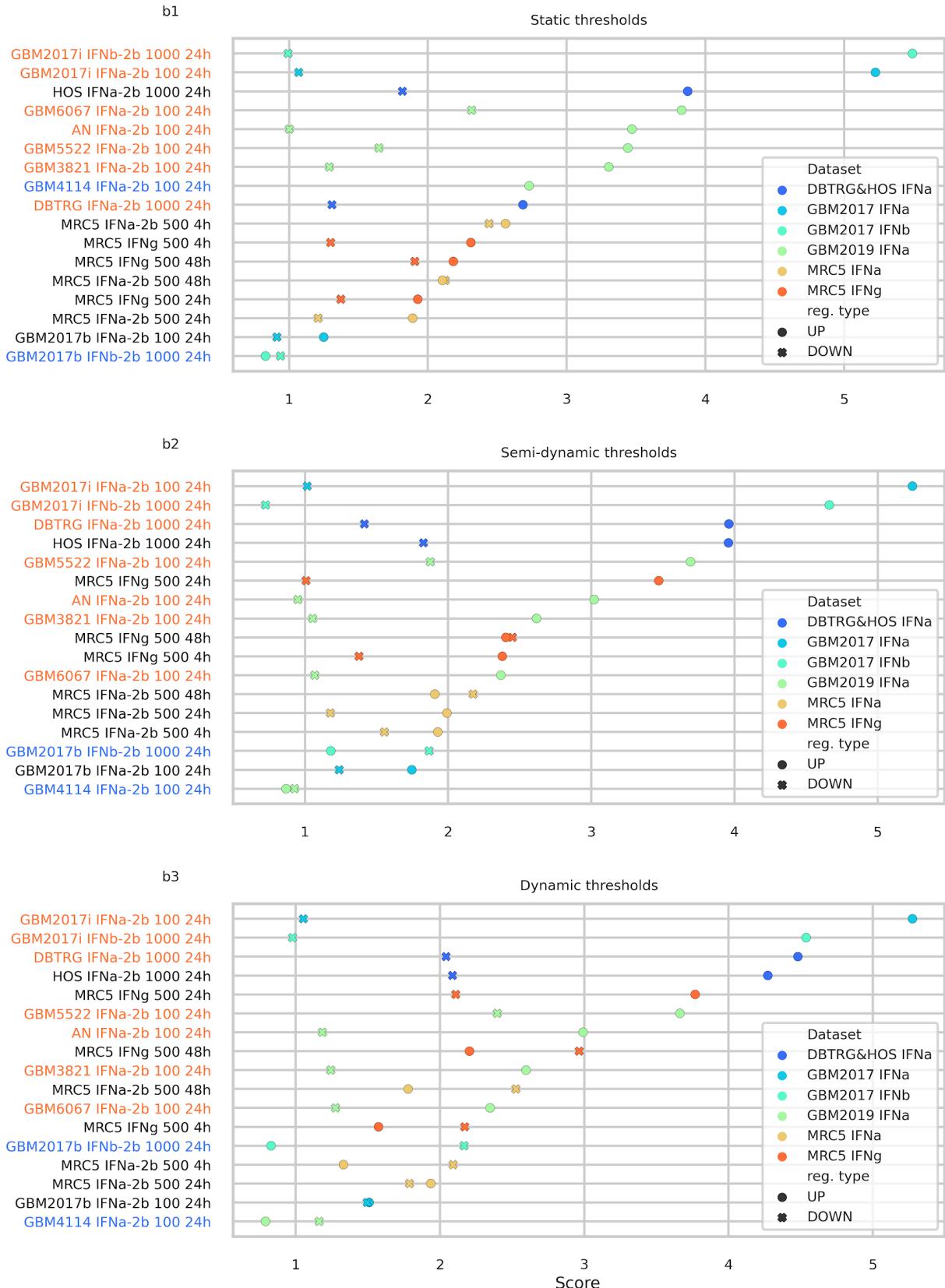
a2



a3



Modified euclidean distance

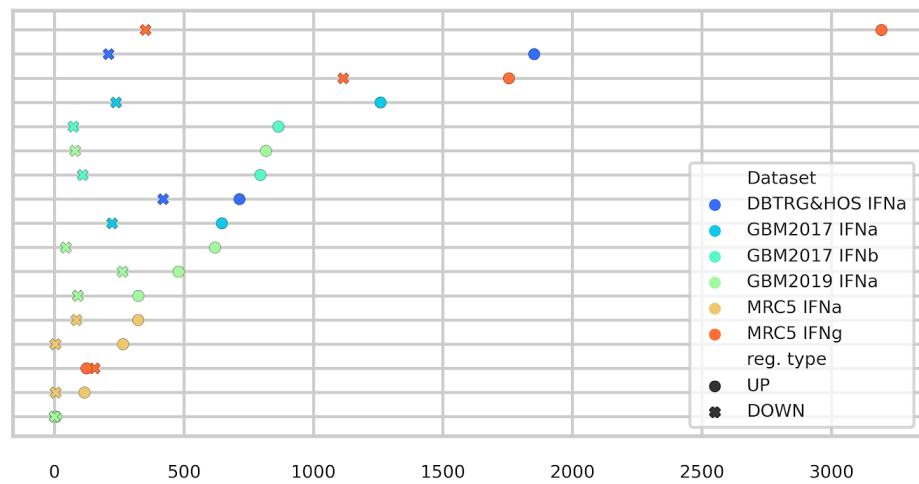


$$\pi_1 = \sum_{i=1}^n |\log_2 FC_i \cdot \log_{10} FDR_i|$$

c1

MRC5 IFNg 500 48h
 DBTRG IFNa-2b 1000 24h
 MRC5 IFNg 500 24h
 GBM2017i IFNa-2b 100 24h
 GBM2017i IFNb-2b 1000 24h
 AN IFNa-2b 100 24h
 GBM2017b IFNb-2b 1000 24h
 HOS IFNa-2b 1000 24h
 GBM2017b IFNa-2b 100 24h
 GBM6067 IFNa-2b 100 24h
 GBM5522 IFNa-2b 100 24h
 GBM3821 IFNa-2b 100 24h
 MRC5 IFNa-2b 500 24h
 MRC5 IFNa-2b 500 48h
 MRC5 IFNg 500 4h
 MRC5 IFNa-2b 500 4h
 GBM4114 IFNa-2b 100 24h

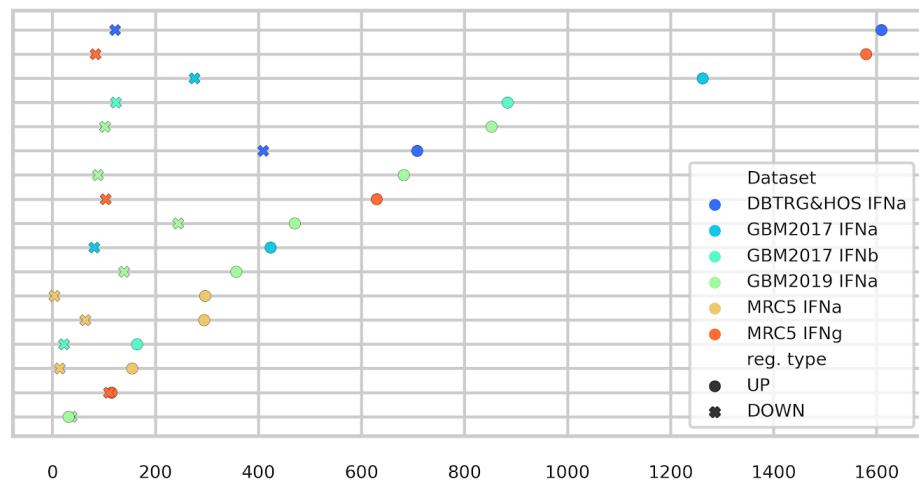
Static thresholds



c2

DBTRG IFNa-2b 1000 24h
 MRC5 IFNg 500 48h
 GBM2017i IFNa-2b 100 24h
 GBM2017i IFNb-2b 1000 24h
 AN IFNa-2b 100 24h
 HOS IFNa-2b 1000 24h
 GBM6067 IFNa-2b 100 24h
 MRC5 IFNg 500 24h
 GBM5522 IFNa-2b 100 24h
 GBM2017b IFNa-2b 100 24h
 GBM3821 IFNa-2b 100 24h
 MRC5 IFNa-2b 500 48h
 MRC5 IFNa-2b 500 24h
 GBM2017b IFNb-2b 1000 24h
 MRC5 IFNa-2b 500 4h
 MRC5 IFNg 500 4h
 GBM4114 IFNa-2b 100 24h

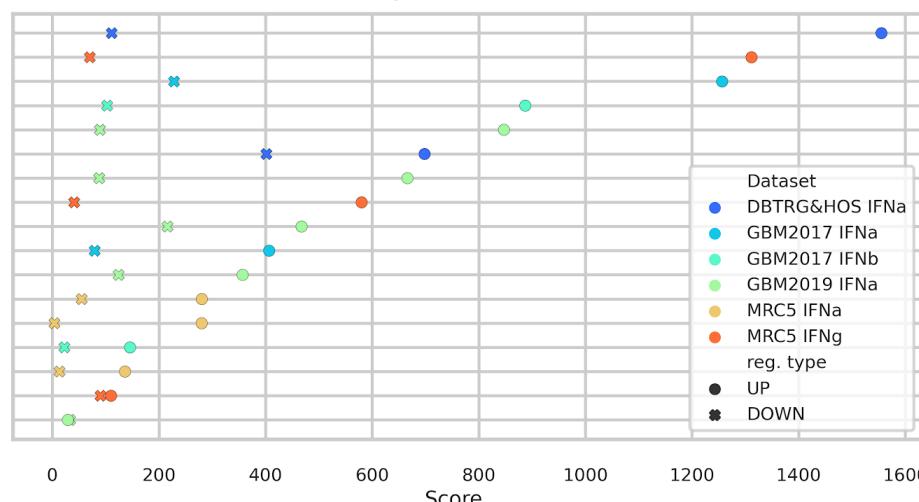
Semi-dynamic thresholds

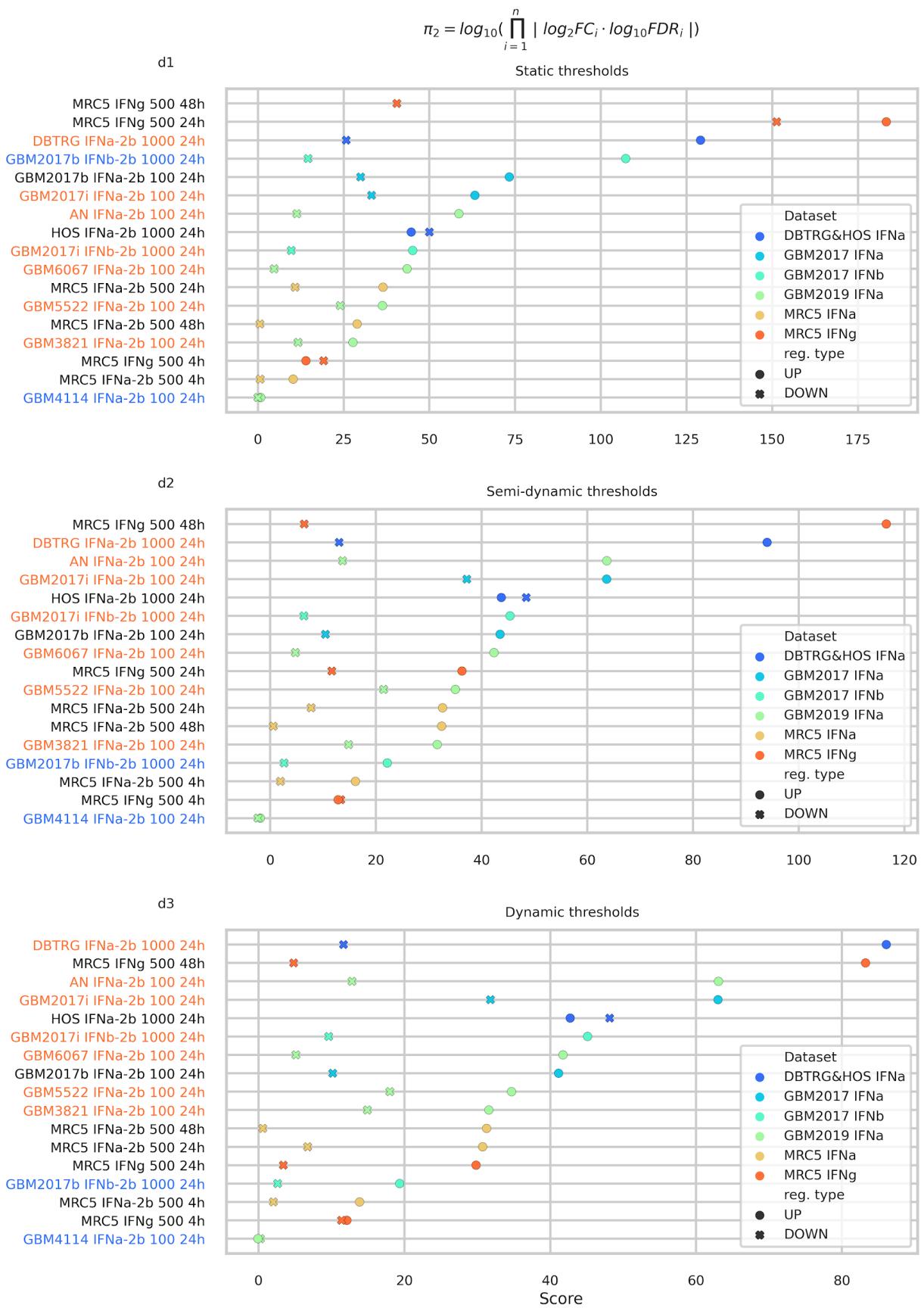


c3

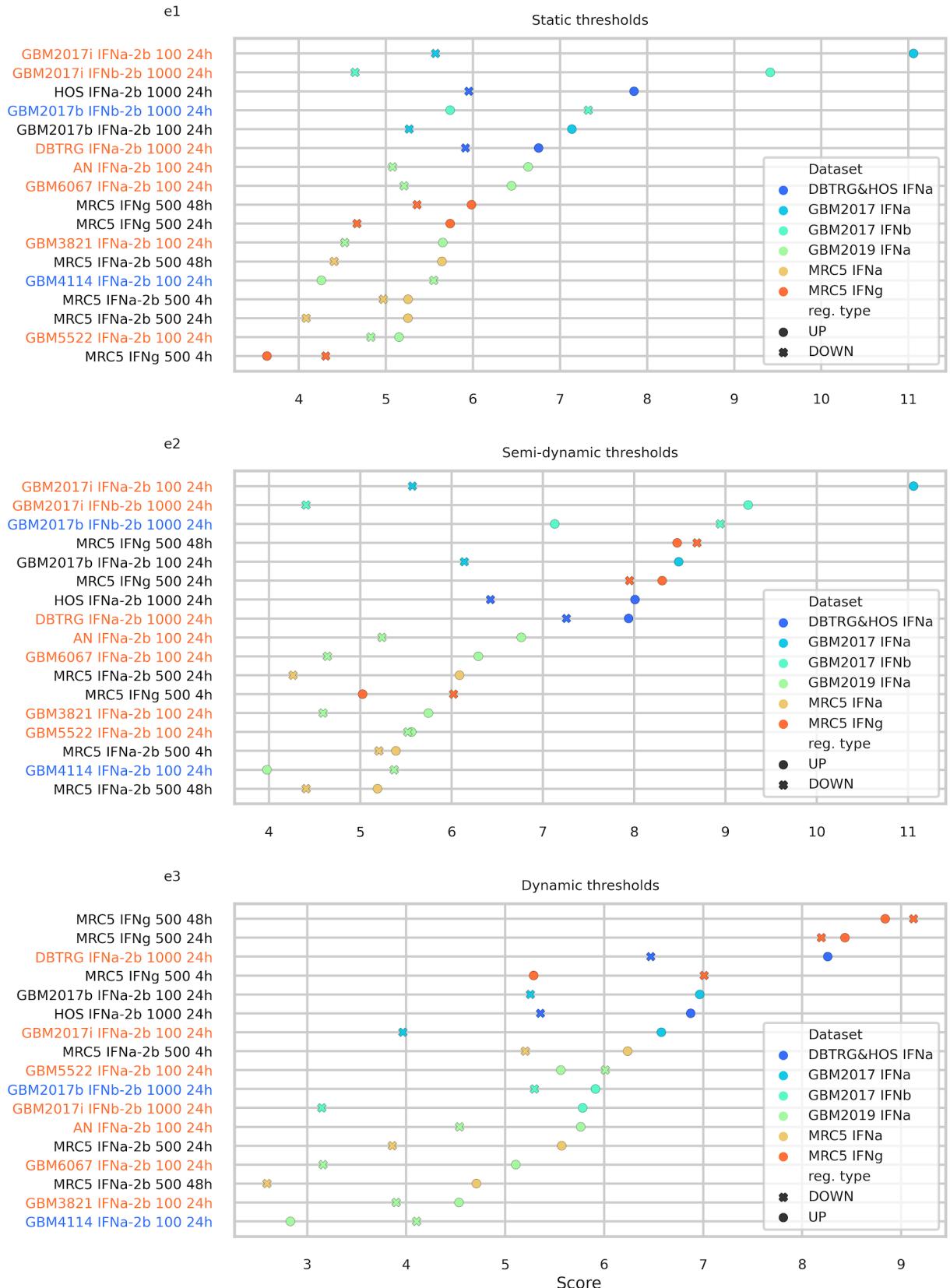
DBTRG IFNa-2b 1000 24h
 MRC5 IFNg 500 48h
 GBM2017i IFNa-2b 100 24h
 GBM2017i IFNb-2b 1000 24h
 AN IFNa-2b 100 24h
 HOS IFNa-2b 1000 24h
 GBM6067 IFNa-2b 100 24h
 MRC5 IFNg 500 24h
 GBM5522 IFNa-2b 100 24h
 GBM2017b IFNa-2b 100 24h
 GBM3821 IFNa-2b 100 24h
 MRC5 IFNa-2b 500 24h
 MRC5 IFNa-2b 500 48h
 GBM2017b IFNb-2b 1000 24h
 MRC5 IFNa-2b 500 4h
 MRC5 IFNg 500 4h
 GBM4114 IFNa-2b 100 24h

Dynamic thresholds

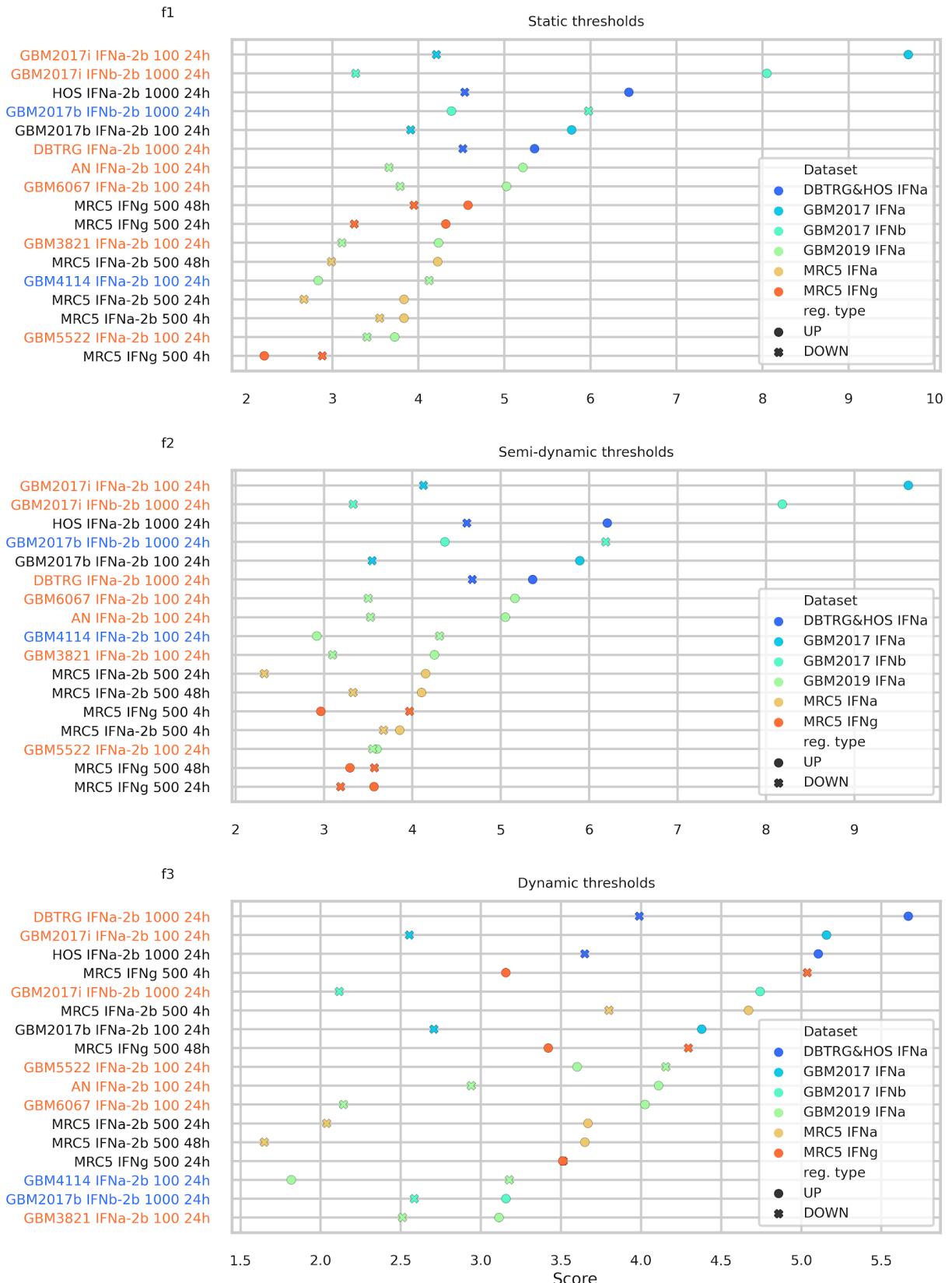




Euclidean distance

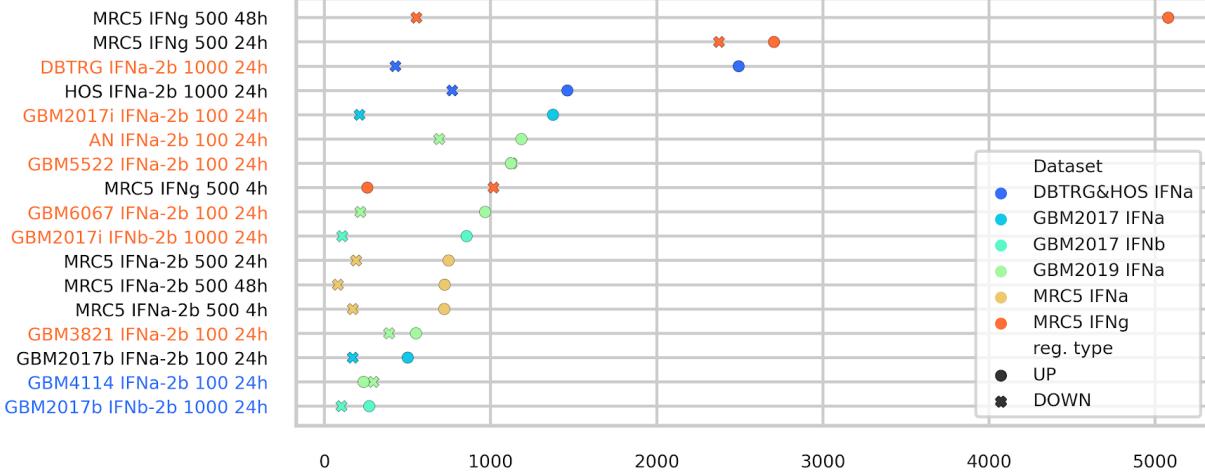


Modified euclidean distance

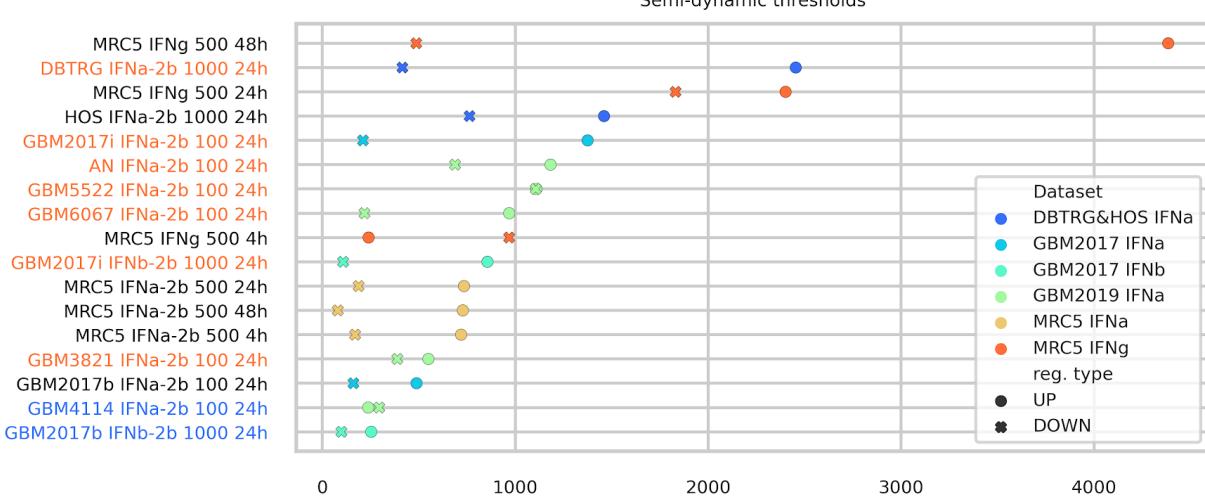


$$\pi_1 = \sum_{i=1}^n | \log_2 FC_i \cdot \log_{10} FDR_i |$$

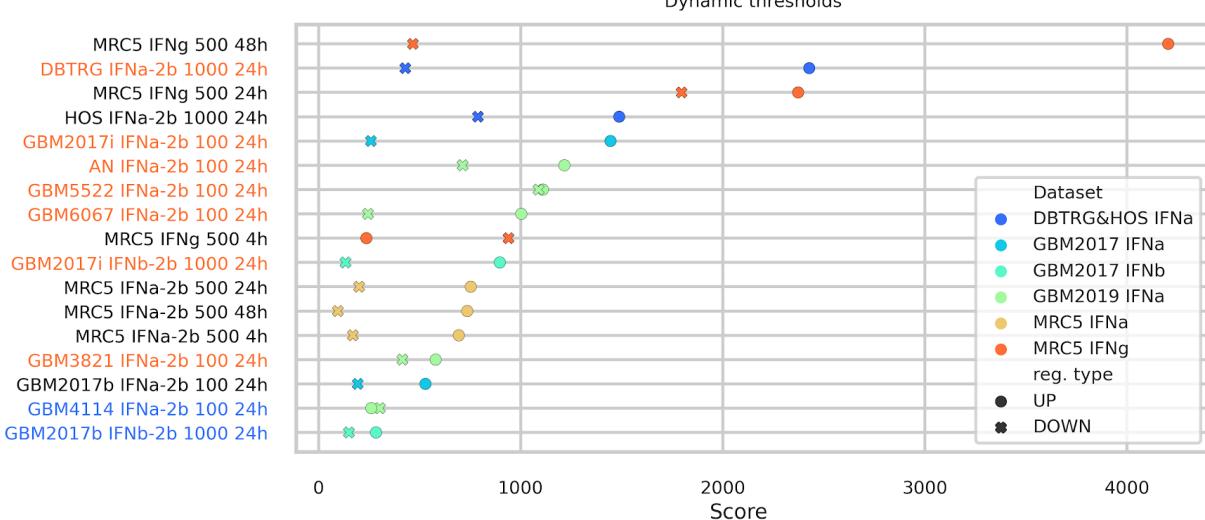
g1



g2

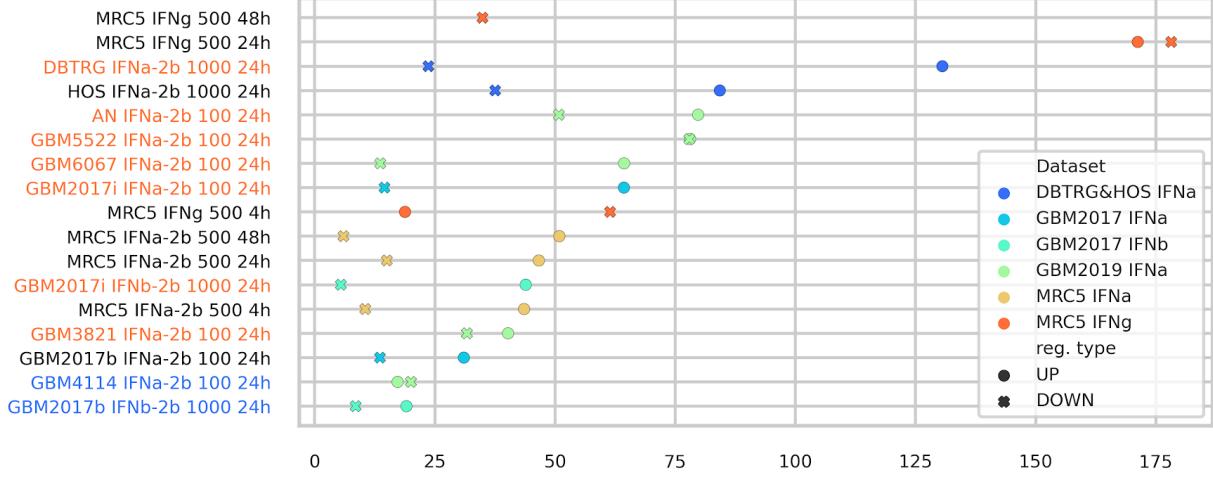


g3

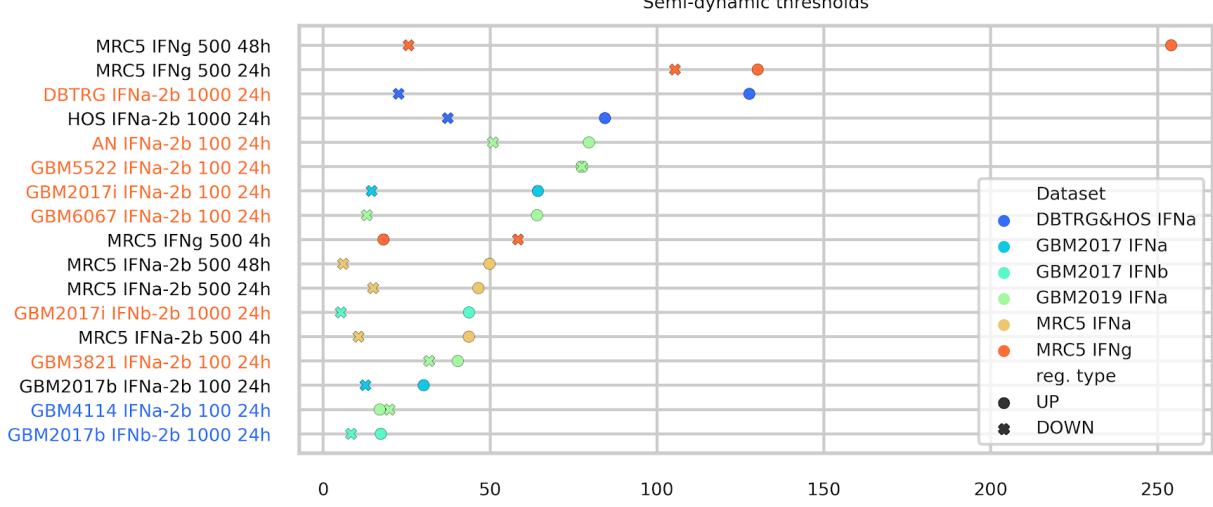


$$\pi_2 = \log_{10} \left(\prod_{i=1}^n |\log_2 FC_i \cdot \log_{10} FDR_i| \right)$$

h1



h2



h3

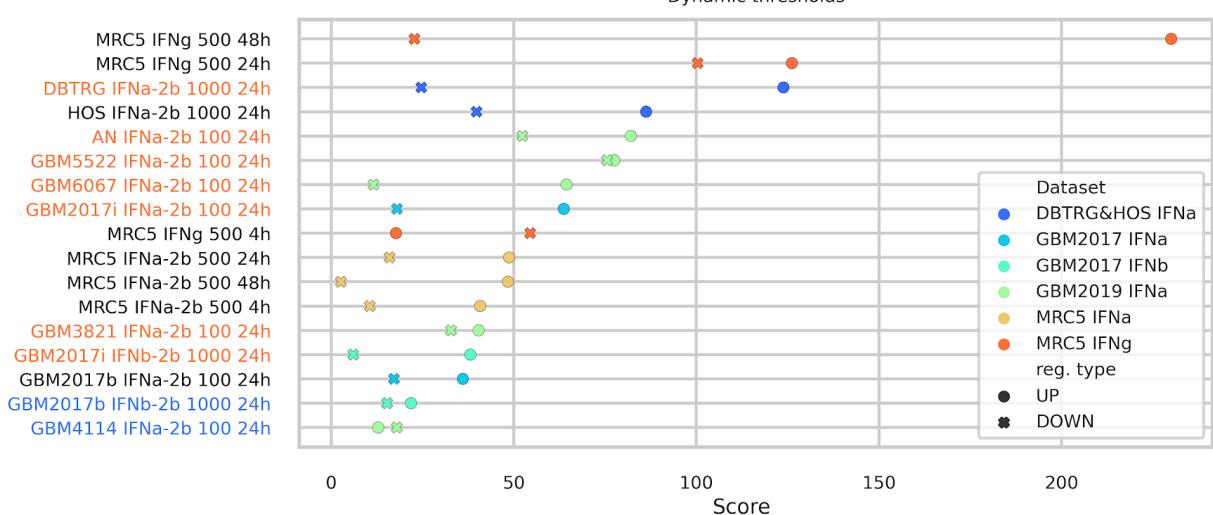
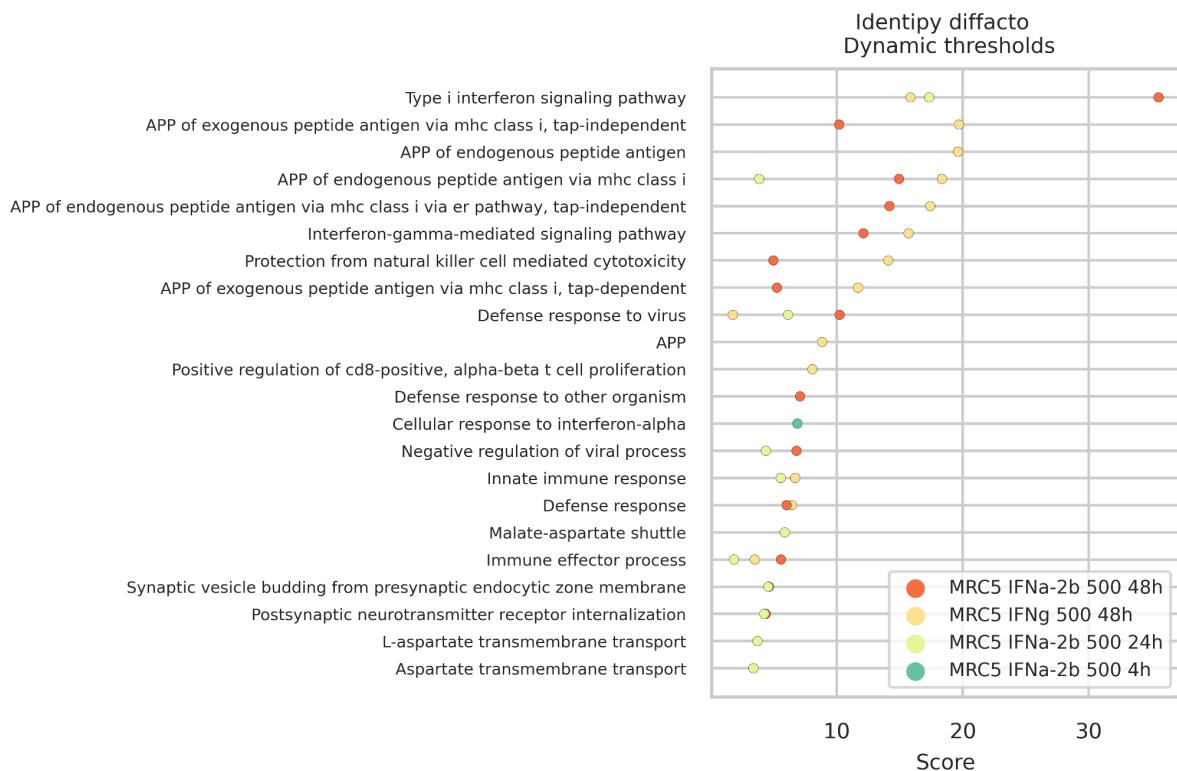


Figure S1. Ranking cellular response to type I and type II interferon treatment using different score equations, missing value imputation strategies (imputation with k-nearest neighbors machine learning or imputation with minimal NSAF detected in LC-MS/MS replicate) and thresholds for DRF selection of (static, semi dynamic or dynamic). Negative controls (blue text on Y axis) clusterized below the positive ones (red text on Y axis) differentiate the best selectivity for response scoring, imputation and DRF selection strategies. DRF selection: static thresholds - DRF selection with $fdr < 0.05$, $|\log_2 FC| > 0.585$; semi dynamic - DRF selection with $|\log_2 FC| > 0.585$ and $X_1 = Q_1 - 1.5(Q_3 - Q_1)$; dynamic thresholds stand for fdr and FC determined as $X_1 = Q_1 - 1.5(Q_3 - Q_1)$, $X_2 = Q_1 + 1.5(Q_3 - Q_1)$ based on $-\log_{10} fdr$ and $\log_2 FC$ density plots. Designation of subfigures in **Fig. S1** corresponds to **Table 2** in the main manuscript, summarizing the clusterization of positive (Pos) and negative (Neg) controls for different data processing strategies. “True” stands if all positive controls are ranked with higher scores than the negative ones.

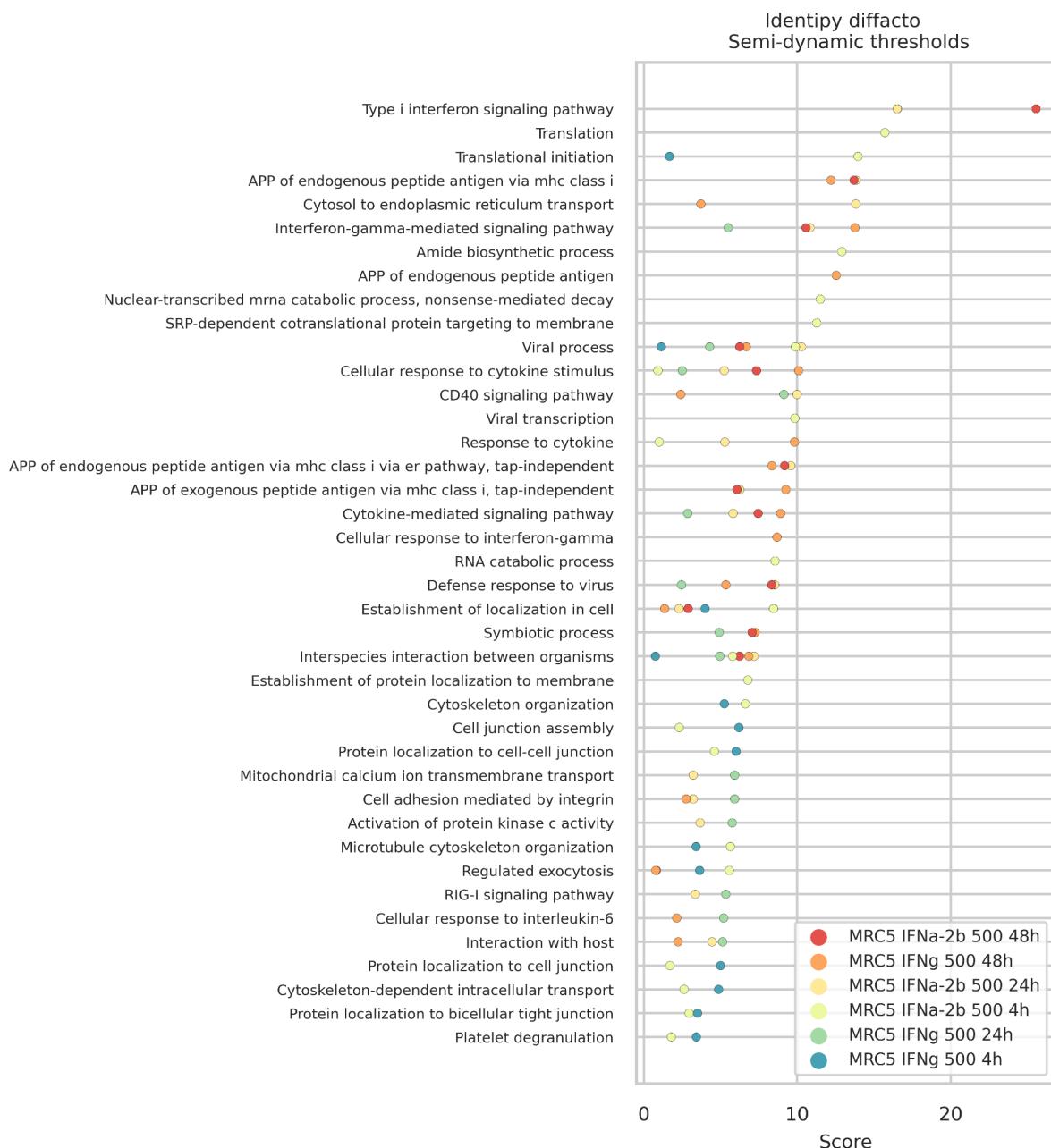
Ref. Fig.S1	Imputation	Equation	DRF selection	Pos above Neg?
a1	Minimal NSAF	E	static	False
a2			semi dyn	False
a3			dyn	False
b1		E_m	static	False
b2			semi dyn	True
b3			dyn	True
c1		π_1	static	False
c2			semi dyn	True
c3			dyn	True
d1		π_2	static	False
d2			semi dyn	True
d3			dyn	True
e1	kNN	E	static	False
e2			semi dyn	False
e3			dyn	False
f1		E_m	static	False
f2			semi dyn	False
f3			dyn	False
g1		π_1	static	True
g2			semi dyn	True
g3			dyn	True
h1		π_2	static	True
h2			semi dyn	True
h3			dyn	True

Figure S2. Top enriched biological processes ranked by GO score. IdentiPy [DOI: [10.1021/acs.jproteome.7b00640](https://doi.org/10.1021/acs.jproteome.7b00640)] and Scavager [DOI: [10.1002/pmic.201800280](https://doi.org/10.1002/pmic.201800280)] were used for peptide spectrum matching and postsearch validation. LFQ was performed using Diffacto [DOI: [10.1074/mcp.O117.067728](https://doi.org/10.1074/mcp.O117.067728)]. Differentially regulated proteins were selected using (a) dynamic, (b) semi dynamic or (c) static, $fdr < 0.05$ and $|\log_2 FC| > 0.585$, thresholds. STRING [doi: [10.1093/nar/gky1131](https://doi.org/10.1093/nar/gky1131)] was used for GO analysis. $GO_{score} = E \cdot \text{abs}(\log_{10} fdr)$, where E and fdr are the enrichment and Benjamini-Hochberg corrected p -values, respectively, for the enriched GO term reported by STRING. “APP” stands for antigen processing and presentation.

(a)



(b)



(c)



Figure S3. Heatmap for the proteins involved in biological processes correlated with topotecan concentrations, Pearson coeff. = 0.9 (GO terms shown in Fig.5 of the manuscript).

