# Supplemental information

# RummaGEO: Automatic mining of human

# and mouse gene sets from GEO

Giacomo B. Marino, Daniel J.B. Clarke, Alexander Lachmann, Eden Z. Deng, and Avi Ma'ayan

# Supplemental Information

for

## RummaGEO: Automatic Mining of Human and Mouse Gene Sets from GEO

Giacomo B. Marino[1], Daniel J. B. Clarke[1], Eden Z. Deng[1], Avi Ma'ayan[1,*]

[1]Mount Sinai Center for Bioinformatics, Department of Pharmacological Sciences, Department of Artificial Intelligence and Human Health, Icahn School of Medicine at Mount Sinai, New York 10029, NY USA

*To whom correspondence should be addressed:

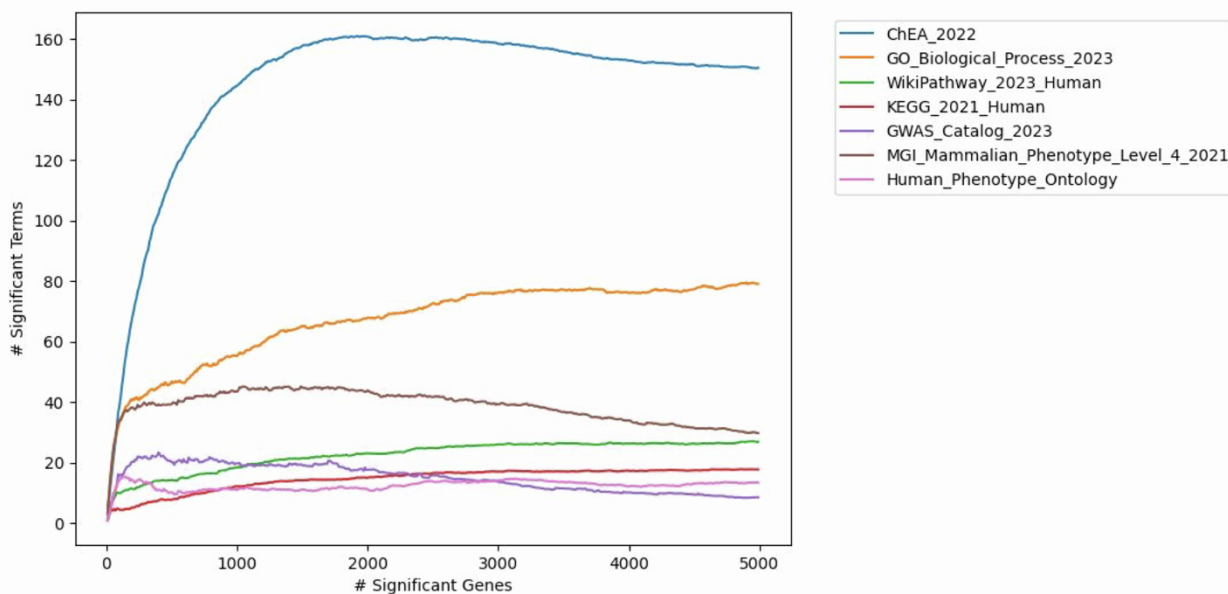Lead contact and corresponding author: avi.maayan@mssm.edu

Figure S1



**Fig. S1 Significantly enriched terms from selected Enrichr libraries varying based on the cutoff number of genes.**
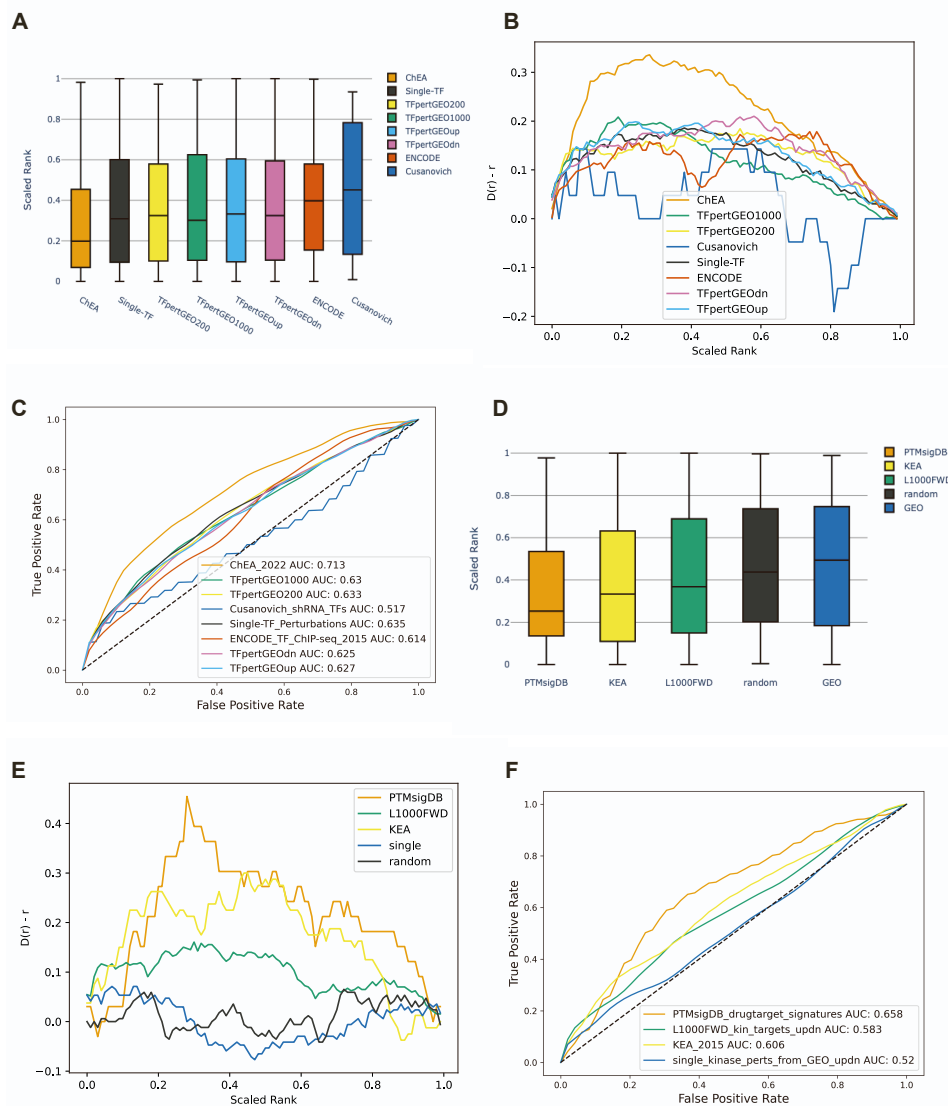
**Figure S2**

**Fig. S2 RummaGEO kinase and transcription factor libraries partial intersection benchmarking. A.** Scaled rank (0 highest rank, 1 lowest rank) for transcription factor benchmarking libraries as computed by the Fisher's exact test; **B.** Deviation of the cumulative distribution for scaled ranks of each transcription factor from uniform distribution (Kolmogorov-Smirnov test for goodness of fit compared to uniform distribution: ChEA 2022 p = 3.83E-35; TFpertGEO1000 p = 1.07E-24; TFpertGEO200 p = 2.56E-21; Cusanovich shRNA TFs p = 5.47E-01; Single-TF Perturbations p = 2.01E-25; ENCODE TF ChIP-seq 2015 p = 1.01E-10; TFpertGEOdn p = 4.33E-19; TFpertGEOup p = 6.22E-21; **C.** 5,000 bootstrapped curves with downsampled negative class were generated to compute mean receiver operating characteristic (ROC) curves and mean area under the ROC curves (AUC) for transcription factors. **D.** Scaled rank (0 highest rank, 1 lowest rank) for kinase benchmarking libraries as computed by Fisher's exact test; **E.** Deviation of the cumulative distribution for scaled ranks of each kinase from uniform distribution (Kolmogorov-Smirnov test for goodness of fit compared to uniform distribution: PTMsigDB drugtarget signatures p = 1.05E-03; L1000FWD kin targets updn p = 1.06E-08; KEA 2015 p = 1.14E-03; single kinase perts from GEO updn p = 1.50E-01; random P = 5.59E-01; **F.** 5,000 bootstrapped curves with downsampled negative class generated to compute mean receiver operating characteristic (ROC) curves and mean area under the ROC curves (AUC) for kinases.

Figure S3



**Fig. S3 Alzheimer's gene set on Rummagene (use case 1).**

Figure S4



**Fig. S4 Results of submitting the Alzheimer's gene set from Rummagene on RummaGEO (use case 1).**
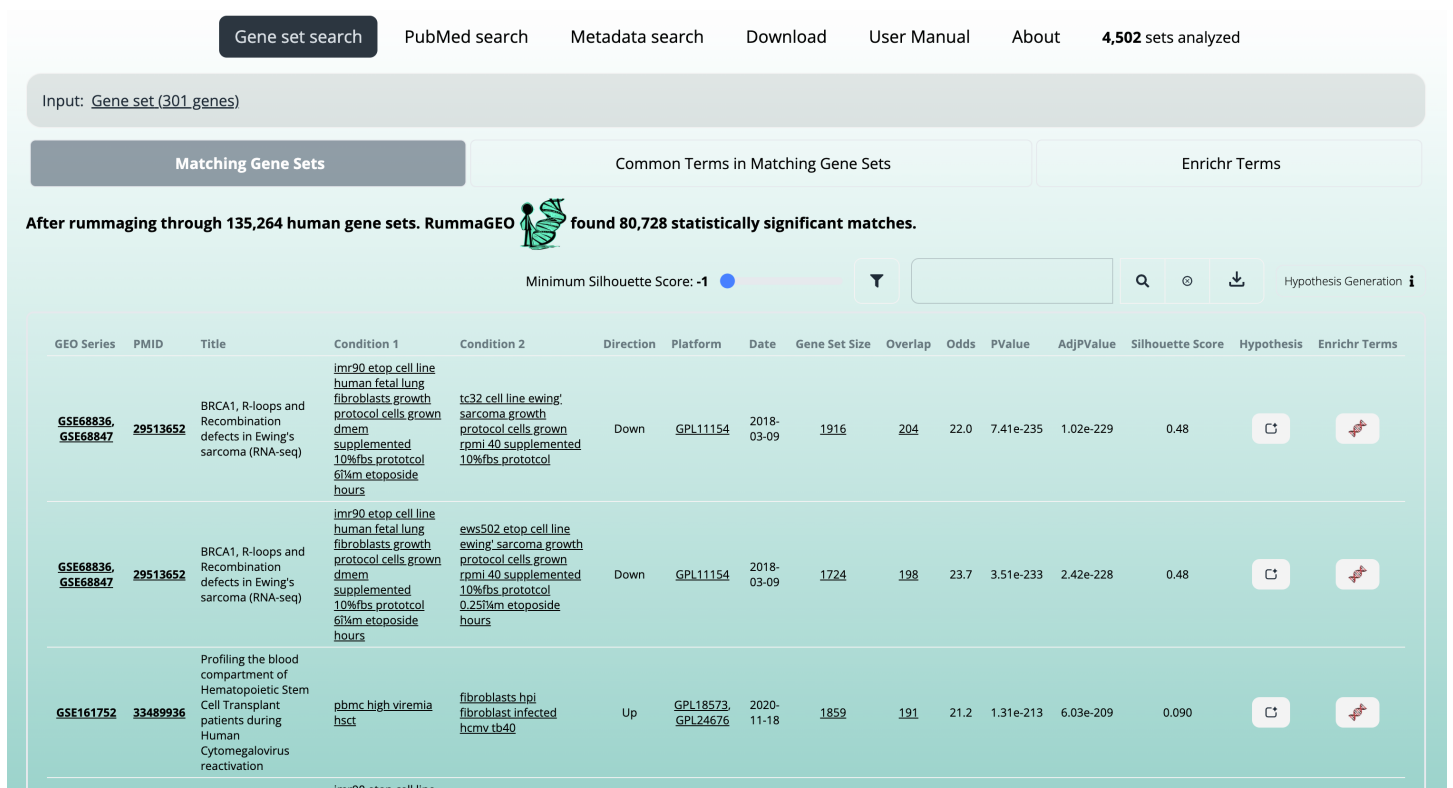
## Figure S5



**Fig. S5 Results from submitting SenoRanger gene set on RummaGEO (use case 2).**