

1 **Compendium of Synovial Signatures Identifies Pathologic Characteristics for**
2 **Predicting Treatment Response in Rheumatoid Arthritis**

3

4 Ki-Jo Kim ^{1,2,3}, Minseung Kim ^{2,3} & Ilias Tagkopoulos ^{2,3}

5

6 ¹ Division of Rheumatology, St. Vincent Hospital, The Catholic University of Korea, Seoul,
7 Republic of Korea

8 ² Department of Computer Science, University of California, Davis, California, US

9 ³ Genome Center, University of California, Davis, California, US

10

11 Corresponding author:

12 Ilias Tagkopoulos, PhD

13 Genome Center, University of California, Davis, California, US

14 Address: 5313 Genome & Biomedical Science Facility, 451 Health Sciences Drive, University of
15 California, Davis, CA 95616

16 E-mail: itagkopoulos@ucdavis.edu

17

18

1 **ABSTRACT**

2 Treatment of patients with rheumatoid arthritis (RA) is challenging due to clinical heterogeneity
3 and variability. Integration of RA synovial genome-scale transcriptomic profiling of different
4 patient cohorts can provide insights on the causal basis of drug responses. A normalized
5 compendium was built that consists of 256 RA synovial samples that cover an intersection of
6 11,769 genes from 11 datasets. Differentially expression genes (DEGs) that were identified in
7 three independent methods were fed into functional network analysis, with subsequent grouping
8 of the samples based on a non-negative matrix factorization method. Finally, we built a predictive
9 model for treatment response by using RA-relevant pathway activation scores and four machine
10 learning classification techniques. We identified 876 up-regulated DEGs including 24 known
11 genetic risk factors and 8 drug targets. DEG-based subgrouping revealed 3 distinct RA patient
12 clusters with distinct activity signatures for RA-relevant pathways. In the case of infliximab, we
13 constructed a classifier of drug response that was highly accurate with an AUC/AUPR of 0.92/0.86.
14 The most informative pathways in achieving this performance were the NF κ B-, Fc ϵ RI- TCR-, and
15 TNF signaling pathways. Similarly, the expression of the HMMR, PRPF4B, EVI2A, RAB27A,
16 MALT1, SNX6, and IFIH1 genes contributed in predicting the patient outcome. Construction and
17 analysis of normalized synovial transcriptomic compendia can provide useful insights for
18 understanding RA-related pathway involvement and drug responses for individual patients. The
19 efficacy of a predictive model for personalized drug response has been demonstrated and can be
20 generalized to several drugs, co-morbidities, and other relevant features.

21

22

23

1 **Introduction**

2 Rheumatoid arthritis (RA) is a complex autoimmune disease involving a multitude of
3 environmental and genetic factors that exhibit nonlinear dynamic interactions (1). The disease is
4 characterized by chronic inflammation of the synovium, which results in irreversible damage to
5 the joints over time, leading to pain and functional impairment. Severity and clinical course of the
6 disease is highly variable across the different patients and hence difficult to predict (1). Despite an
7 introduction of tumor necrosis factor (TNF) inhibitors, over 30% of patients do not respond fully
8 to therapy (2). Moreover, a considerable subset of the patients who showed initial good response
9 experience disease flare or efficacy reduction even on drugs (2). A similarly discouraging picture
10 is painted in the case of other biologics, which creates a current and present need for better
11 understanding of the disease. An early, aggressive and personalized treatment that provides the
12 best possible drug combination for a patient is likely to improve our ability to treat RA. Despite
13 the introduction of novel drugs and the fact that RA is an active research topic, however, we still
14 have substantial gaps in our knowledge regarding the mechanistic basis of RA progression, which
15 is needed to administer personalized and precise care.

16 In RA, gene expression profiling has been used to gain insights regarding pathogenesis and drug
17 response (3). Since studies so far have been in unrelated cohorts and study groups, small sample
18 size, heterogeneity in study population (sex, age, and ethnicity), differences in technical protocols,
19 microarray platform, and data analysis methods has hindered a comprehensive analysis across all
20 available datasets. In addition, most studies have collected samples from whole blood or peripheral
21 blood mononuclear cells, which are easier to acquire but have a limited capacity to adequately
22 reflect the joint inflammation (4-6). However, integrated analysis of the compendium by the
23 accrued genome-wide datasets provided opportunities to capture the missing features, bridge the

1 gap between prior knowledge, and better understand human diseases in the field of cancer and
2 infectious disease (7,8).

3 In this study, our aim is to elucidate the various transcriptional and signaling signatures of RA by
4 performing a comprehensive meta-analysis of the publicly available datasets that have been
5 published so far. We focus on synovial tissue samples to avoid the high false discovery rates
6 coming from blood samples. We have applied several preprocessing and normalization steps to
7 create a cohesive, homogenized compendium of genome-wide gene expression signatures for
8 downstream analysis. We used this compendium to separate expression-driven subgroup,
9 understand the key cellular components in each group and then use genes and pathways with high
10 information value that we have identified to create predictive models for drug responsiveness.

11

1 **Methods**

2 **Systematic search and data collection**

3 We used the keywords “Rheumatoid Arthritis”, “Synovium or synovial tissue”, “Transcriptomics
4 or microarray”, “Dataset” in Google Scholar and PubMed to find relevant publications to the topic
5 of synovial gene signatures of patients with rheumatoid arthritis (**Fig. 1**). We retrieved all
6 publications that were accompanied by high-throughput datasets (20 studies in total). From the
7 resulting set, we removed entries that had been duplicated and selected datasets measuring over
8 10,000 genes to secure the largest size of genes and samples. Since there was a trade-off between
9 the number of studies to include and the number of genes that are within the intersection from all
10 datasets, we optimized the product of the two by selecting the point where these two trends cross
11 (**Fig. S1**). The final RA sample count was 256, the osteoarthritis (OA) count 41, and 36 normal
12 (NC) samples were included as controls. Ultimately, the final RA compendium was constructed
13 out of 11 studies with a total of 333 samples, one per patient, covering 22,721 genes total (common
14 core of 11,769 genes).

15

16 **Data normalization and removal of batch effects**

17 For one-channel arrays, the image data was first imported and then the Robust Multi-array Average
18 (RMA) method was applied for a set of replicates for background correction, normalization, probe-
19 set summarization. For dual-channel arrays, the image data were imported and background
20 correction was performed using normexp as it was shown to outperform other methods. Red and
21 green channels were separated and quantile-normalized for each set of replicates. The vectors for
22 the matrices were normalized using the quantile normalization method. Residual technical batch
23 effects arising due to heterogeneous data integration were corrected using the ComBat function

1 within the empirical Bayes package. Quality assurance and distribution bias was evaluated by
2 Principal Component Analysis (**Fig. S2**).

3

4 **The RA compendium**

5 After preprocessing, the gene expression profiles have a significant reduction of systematic,
6 dataset-specific bias in comparison with the same dataset before normalization and batch
7 correction (**Fig. S2**). The resulting compendium has a gene size of 11,769 in 333 samples,
8 including 256 RA patients, 41 OA patients, and 36 normal controls. In 105 of the RA samples,
9 synovial tissue sampling was conducted before the start of certain drug: 11 for adalimumab
10 (ADLM), 62 for infliximab (IFXM), 8 for methotrexate (MTX), 12 for rituximab (RTXM), and 12
11 for tocilizumab (TOCM). For these patients, assessment of disease activity and response was
12 performed per the EULAR response criteria (9) 12-16 weeks after initiation of therapy: 32 were
13 good, 47 were moderate, and 26 were poor responders. Information on demographics and clinical
14 characteristics including age, sex, disease duration, and RF positive were not fully annotated for
15 each RA sample (**Table S1**).

16

17 **Filtering of differentially expressed genes**

18 In order to identify the differentially expressed genes (DEGs), we employed three widely-used
19 methods: (a) an empirical Bayesian method using the Benjamini and Hochberg procedure with a
20 significance threshold at an adjusted p-value < 0.05 ; (b) the Significance Analysis of Microarray
21 (SAM) method, with a significance threshold of false discovery rate $FDR < 0.05$; (c) the Rank
22 Products (RP) method with a significance threshold set at percentage of false prediction $pfp < 0.05$.

1 The resulting list of DEGs is the intersection of the three individual DEGs sets for each method to
2 minimize the false discovery rate statistic.

3

4 **Functional enrichment analysis**

5 We performed functional enrichment analysis focusing on the up-regulated DEGs using the
6 Database for Annotation, Visualization, and Integrated Discovery (DAVID) software (10). Terms
7 were regarded significant if the p-value (EASE score) is lower than 0.05, the enrichment score
8 higher than 1.3, and the fold enrichment was larger than 1.5.

9

10 **Gene set enrichment analysis**

11 Gene set enrichment analysis (GSEA) analysis was carried out using the GSEA software from the
12 Broad Institute to assess the overrepresentation of RA-related gene sets (11). The enrichment
13 results were visualized with the Enrichment Map format, where nodes represent gene-sets and
14 weighted links between the nodes represent an overlap score depending on the number of genes
15 two gene-sets share (Jaccard coefficient) (12). To intuitively identify redundancies between gene
16 sets, the nodes were connected if their contents overlap by more than 25%. Clusters map to one or
17 more functionally enriched groups, which were manually circled and assigned a label.

18

19 **Construction of protein-protein interaction network**

20 To assess the interconnectivity of DEGs in the RA synovium samples, we constructed a protein-
21 protein network based on the interaction data obtained from public databases including BIOGRID
22 (13), HPRD (14), IntAct (15), Reactome (16), and STRING (17). In the network, nodes and edges
23 represent genes and functional or physical relationships between them, respectively. Graph theory

1 concepts such as degree, closeness, and betweenness were employed to assess the topology of this
2 network. Hub molecules were defined as the shared genes in top 10% with the highest rank in each
3 arm of the three centrality parameters (18).

4

5 **Drug target prioritization strategy**

6 To obtain a genome-wide drug target prioritization, we applied the Heat Kernel Diffusion Ranking
7 approach (19). This method prioritizes the candidate genes by diffusing the differential expression
8 values of the candidate genes through the network based on the confidence scores of the
9 associations or interactions and is a powerful network-based machine learning approach to identify
10 putative drug targets. The diffusion process is formulated by using a Laplacian exponential
11 diffusion kernel and a score is computed by multiplying the differential expression values with the
12 heat kernel. Drug targets are assumed to get a high score since these genes tend to be the central
13 nodes in subnetworks showing significant transcriptional changes following treatment. The
14 predictive power of the method is dependent on the network characteristics, the quality of the
15 expression data, the type of drug and the target (19).

16

17 **Non-negative matrix factorization and determination of the optimal number of clusters**

18 To classify the RA patients into subgroups based on their molecular signatures, we used the non-
19 negative matrix factorization (NMF) method. NMF clustering is a powerful unsupervised approach
20 to identify the disease subtype or patient subgroup and discover biologically meaningful molecular
21 pattern (8,20). We applied the consensus NMF clustering method and initialized 100 times for
22 each rank k (range from 2 to 6), where k was a presumed number of subtypes in the dataset. For
23 each k , 100 matrix factorizations were used to classify each sample 100 times. The consensus

1 matrix was used to assess how consistently sample-pairs cluster together. We then computed the
2 cophenetic coefficients and silhouette scores for each k , to quantitatively assess global clustering
3 robustness across the consensus matrix. The maximum peak of the cophenetic coefficient and
4 silhouette score plots determined the optimal number of clusters (20). To confirm unsupervised
5 clustering results, we used t -distributed stochastic neighborhood embedding (t -SNE) (21), a
6 powerful dimensionality reduction method. The t -SNE method captures the variance in the data
7 by attempting to preserve the distances between data points from high to low dimensions without
8 any prior assumptions about the data distribution.

9

10 **Scoring of pathway activation**

11 To quantify certain biological pathway activity, we calculated the gene expression z-scores (8,22).
12 Briefly, a Z-score is defined as the difference between the error-weighted mean of the expression
13 values of the genes in each pathway and the error-weighted mean of all genes in a sample after
14 normalization. BCR-, chemokine-, Jack-STAT-, MAPK-, NF κ B-, p53-, PI3K-AKT-, RIG-I-like
15 receptor-, Fc ϵ RI-, TCR-, TGF β -, TLR-, TNF-, VEGF-, and Wnt signaling pathways and their
16 gene sets were imported from Kyoto Encyclopedia of Genes and Genomes (KEGG) database (23)
17 and IFN type I- and type II signaling pathways and their gene sets referred to Reactome database
18 (16). Z-scores were computed using each pathway in the signature collection for each of the
19 samples, resulting in a matrix of pathway activation scores.

20

21 **Supervised learning analyses for the prediction of drug responsiveness**

22 We used Naïve Bayes (NB), Decision Trees (DT), k -Nearest-Neighbors (KNN), and Support
23 Vector Machines (SVM) to create drug responsiveness predictors (24,25). Each binary SVM was

1 built using Gaussian Radial Basis Function (RBF) kernel and the Sigma hyperparameter was
2 determined from the estimation based upon the 0.1 and 0.9 quantiles of the samples. For soft
3 margins, the C parameter that achieved the best performance was in the range of 2^{-4} to 2^7 . For
4 KNN, the k parameter was tuned in the range 2 to 20. All tuning hyperparameters were separately
5 determined for each bootstrapped training dataset.

6 To determine the optimal feature set that enables distinguishing ‘good’ from ‘not good’ responders
7 with the highest accuracy, we employed the wrapper feature selection method (25). The wrapper
8 method uses the classifier as a black box to rank different subsets of the features according to their
9 predictive power. In the wrapper method, a feature set is fed to the classifier and its performance
10 is scored and the feature set with the highest rank is selected as the optimal feature set. The
11 predictive power of each predictor was assessed through Receiver-Operator Characteristics (ROC)
12 and Precision-Recall (PR) curve (26). Data was separated into independent training and test sets
13 in a three-to-one sample-size ratio in a way of stratified random sampling. To make up for small
14 sample size and minimize the error, we constructed the pool of resampled dataset by applying
15 bootstrapping with 1000 iterations and subsequently applying a stratified 10-fold cross-validation
16 (CV) for each bootstrapped dataset (24,25). Tenfold CV measures the prediction performance in a
17 self-consistent way by systematically leaving out part of the dataset during the training process
18 and testing against those left-out subset of samples. Compared to the test on independent dataset,
19 CV has less bias and better predictive and generalization power. The predictive ability of the
20 models generated from all the approaches was tested by performing the CV test at all the ten
21 locations under study. Given the unequal numbers of trials in each class, balanced accuracy
22 formula was employed to calculate the accuracy (27). The baseline is estimated by random

1 expectation based on the pre-determined ratio of each condition. In case of IFXM, a probability of
2 0.29 (18/62) for a “good” and 0.71 (44/62) for a “not good” responder was applied.

3

4 **Statistical analysis**

5 For continuous distributed data, between-group comparisons were performed using the one-way
6 ANOVA, unpaired *t*-test or Mann-Whitney *U* test. Categorical or dichotomous variables were
7 compared using the chi-squared test or Fisher’s exact test. To investigate the difference of pathway
8 activation score across the subgroups, we fitted the one-way ANOVA model using logistic
9 regression. All analyses were conducted in *R* (The R Project for Statistical Computing, [www.r-](http://www.r-project.org)
10 [project.org](http://www.r-project.org)) and R packages used in the analysis and their references were summarized in the **Table**

11 **S2.**

12

13

1 **Results**

2 **The RA transcriptomics compendium**

3 To get a list of RA-related DEGs, gene expression profiles of RA patients were compared with
4 samples from the OA and NC groups. We identified 2762 DEGs for RA versus OA, and 3087
5 DEGs for RA versus NC (**Fig. 1**). Distribution of DEGs was assessed after the DEGs were divided
6 into up- and down-regulated groups (**Fig. 2A**). The number of up-regulated DEGs was 1486 for
7 RA versus OA and 1774 for RA versus NC. The intersection between two up-regulated DEG sets
8 was 876, which we considered as RA-unique (**Fig. 2A** and **Supplementary File S1**).

9

10 **Enriched biological processes and protein-to-protein interaction network**

11 Through GSEA, we performed a functional enrichment analysis where 206 gene ontology
12 processes were identified (**Fig. 2B** and **Fig. S3**). As expected, immune-related biological processes
13 including adaptive and innate immune response, T cell- and B cell activation and response, and
14 cytokine-related responses, were enriched. These occupied the main positions in the network and
15 closely connected to each other. Among cytokine-related processes, interferon- β (IFN- β),
16 interferon- γ (IFN- γ), interleukin (IL)-4, IL-10, IL-12, IL-17, toll-like receptor (TLR), and tumor
17 necrosis factor (TNF)-related processes stood out as being substantially more enriched.

18 Interestingly, several biological processes associated with viral invasion and defense response
19 against viruses were newly identified (See **Fig. S4**). Metabolic processes such as calcium ion
20 regulation and protein synthesis/transportation were enriched (all $P < 0.01$), suggestive of active
21 intracellular signaling and enhanced protein production and enzyme activity.

22 Identification of central attractors in the gene and protein network can provide targets for further
23 experimentation and/or drug discovery. For this reason, we constructed the protein-to-protein

1 interaction network of RA (**Fig. 2C**). We identified 3563 interactions among the 876 DEGs.
2 Thirty-one of DEGs were overlapped with RA genetic susceptibility loci previously discovered
3 (28) (**Fig. S5**) and a total of 56 genes were ranked as hub molecules based on the centrality analysis.
4 The *CD2*, *PTPRC* (protein tyrosine phosphatase, receptor type C, also known as *CD45*), and
5 *PRKCC* (protein kinase C theta) were RA-susceptible genes having hub position in the network
6 and products of these genes are involved in signal transduction of T cells. Eight genes including
7 primary targets (*JAK2*, *SYK*, *CTLA4*, *MS4A1*) and counterpart receptor molecules (*TNFRSF14*,
8 *TNFRSF17*, *TNFRSF18*, and *IL21R*) of cytokines targeted by the drugs currently in use or under
9 clinical trial or development are also differentially expressed (29,30). Interestingly, the targets of
10 small molecule therapeutics, *JAK2* and *SYK* are central hub nodes, in contrast to the targets of
11 biologic agents, such as *CTLA4*, *MS4A1* (also known as *CD20*), *TNFRSF14*, *TNFRSF17*, and
12 *TNFRSF18*. We found 219 RA-associated genes from the DisGeNet database (31), which are
13 genes and variants having a responsible role in disease. Forty-six of them were overlapped with
14 the RA synovial DEG. To assess topological proximity between RA-associated genes and drug
15 targets in PPI network of synovial DEGs, the shortest distance between nodes was calculated (**Fig.**
16 **S6**). Mean distance of *JAK2* and *SYK* was 2.11 ± 0.69 S.D. and 2.09 ± 0.68 , respectively, and
17 significantly shorter than those of other target molecules (range, 2.65 ~ 3.39) (in all cases p-value
18 < 0.05).

19

20 **Identification and characterization of molecular subgroups**

21 Next, we assessed whether RA patients can be categories in subgroups based on their expression
22 profiles through consensus NMF clustering (20). To identify the optimal number of clusters and
23 to assess robustness of the clustering result, we computed the cophenetic coefficient and silhouette

1 score for different numbers of clusters from 2 to 6, where we found that 3 clusters are the optimal
2 representation of the data (**Fig. 3A, Fig. S7, and Methods**). Segregation of RA subgroups was also
3 reproduced by *t*-SNE analysis and principal component analysis (**Fig. 3B**). To identify
4 characteristic molecular signaling pathways enriched in each cluster, we performed functional
5 enrichment analyses for the predicted genes of each cluster (**Fig. 3C**). Nine enriched pathways
6 were identified across the 3 clusters. RA cluster 1 (541 DEGs, 112 samples) exhibited activity for
7 chemokine-, p53-, TNF- and TLR signaling pathways. Cluster 2 (130 DEGs, 64 samples) was also
8 enriched in chemokine-, and p53 signaling pathways, while cluster 3 (205 DEGs, 80 samples) had
9 a high enrichment of chemokine-, Fc ϵ RI-, Jak-STAT-, and T cell receptor (TCR) signaling
10 pathways.

11 Some limitations lie with functional enrichment analysis using DAVID: dependency on the gene
12 list not their expression levels, missing gene sets, and bias towards well-studied genes, and
13 dropouts of some genes during combining the datasets. Thus, to better understand the differences
14 among the three clusters, we analyzed the activation of individual pathways with adding six more
15 pathways that are known to be associated with RA (30,32,33): transforming growth factor (TGF)-
16 β -, vascular endothelial growth factor (VEGF)-, Wnt-, B cell receptor (BCR)-, NF κ B-, PI3K-AKT
17 (**Fig. 4**). As shown in the chord diagram, these pathways are strongly connected, with only TGF β -,
18 P53-, and Wnt signaling pathways more isolated than others (less shared DEGs). Especially TGF β -
19 and Wnt, have an opposite trend in their DEG expression (higher in cluster 1, mid in cluster 2 and
20 low in cluster 3), which is the opposite of the trend we observe in most of the other pathways (**Fig.**
21 **4 and Fig. S8**). P53 signaling pathways shared fewer genes with other pathways but strongly
22 correlated with BCR-, chemokine-, Fc ϵ RI-, TCR-, TLR-, and TNF signaling pathways.

1 While the activation scores of all pathways exhibited significant difference across the various
2 clusters, all clusters exhibited one of the two trends in a statistically significant manner ($P < 0.05$ in
3 all cases) and in accordance with the observation through DEG-driven enrichment (all cases except
4 TNF). It was found that IFN-related processes were enriched in GSEA and this was reproduced in
5 the precedent studies (5,34-37). Since KEGG database did not provide IFN pathways data, we
6 imported data on gene set of IFN pathways from the Reactome database. Levels of activation score
7 of Type I and II IFN pathways were significantly different across the clusters (all $P < 0.001$) and
8 showed a tendency to increase from cluster 1 to 3 (**Fig. S9**). Compared with RA cluster 2 and 3,
9 RA cluster 1 had moderate activation scores for most of the proinflammatory signaling pathways
10 but high for PI3K-AKT-, TGF β - and Wnt signaling pathways, which are principally involved in
11 synovial proliferation and tissue remodeling (38). RA cluster 2 and 3 showed comparable activities
12 for most of the proinflammatory pathways. More active in RA cluster 2 were the P53- and PI3K-
13 AKT signaling pathways, which were reported to play a role in regulating apoptosis of
14 synoviocytes or macrophages (39,40). In RA cluster 3, TCR-, Jak-STAT-, and NF κ B signaling
15 pathways were remarkable and it is noteworthy that IFN signaling pathways were most scored.
16 Cellular processes affected by these pathways are in agreement with the DEG-driven enriched GO
17 terms in each cluster (**Fig. S10**). This result indicates that there exist RA subgroups representing a
18 distinct mode of inflammation deflected toward a certain combination of signaling pathways
19 (**Table S3**). To prioritize the drug candidate targets, we ranked the DEGs using the Heat Kernel
20 Diffusion Ranking approach (19). The identified drug target candidates are the central nodes in the
21 subnetworks, with the highest disruption to the network under perturbation.
22 Next, we examined the relationship between identified 3 subgroups and the pertinent clinical
23 features based on the provided information. There was no difference in gender ratio, age

1 distribution, and tissue sampling method across the subgroups ($P > 0.10$ in all cases, see **Fig. S11**).
2 Because data on the disease duration, activity, and RF positive were not provided individually for
3 each sample, we compared two distinctively opposing datasets from compendium: the first
4 (GSE45867) includes naïve, untreated RA patients with disease duration of <1 year, moderate
5 disease activity and with arthroscopic needle biopsy performed before MTX or TCZM therapy
6 (41). The second (GSE21537) is a cohort of the long-standing RA patients with high disease
7 activity who had failed at least two DMARDs (including MTX) and did arthroscopic needle biopsy
8 before IFXM therapy (42). Disease duration and activity were significantly longer and higher in
9 the latter dataset (all $P < 0.001$) while there was no difference in age, gender, and RF positive
10 between two datasets (all $P > 0.10$). Distribution of 3 subgroups did not differ between two datasets
11 ($P = 0.8664$) (**Fig. S11**), indicating gene expression pattern by 3 subgroups have little direct
12 relevance to disease duration and activity.

13

14 **Towards a predictor of drug response**

15 For 105 RA samples that we had drug effectiveness data, we tested the hypothesis that there is an
16 association between drug responsiveness and cluster membership. Out of the 5 drugs that we had
17 data on (ADLM, IFXM, MTX, RTX, and TOCM) we were not able to identify any such
18 association (**Fig. S12**). Cluster 1 patients had an encouraging response to TOCM but at a low
19 statistical significance level (p-value equal to 0.082). In addition to the intricacy of the pertinent
20 pathways, the small size of samples treated by the specific drug, and their potential heterogeneity
21 make the association between drug responsiveness and RA clusters difficult.
22 Since the differential expression of genes and pathways is at a higher resolution than general
23 clustering signatures, we tested whether drug response can be predicted by using such features.

1 We focused on the patients that were treated with IFXM due to the larger sample size (n=62). To
2 test this hypothesis, we applied outcome to a binary classification (labels of “good” and “not good”
3 responder) and tried two approaches: pathway-driven and DEG-driven models. Note that PCA
4 analysis does not reveal separating distributions between the “good” and “not good” responders
5 both for pathway activation score and DEG values (**Fig. S13**).

6 As features, we used the 17 pathways that are represented by continuous variables through their
7 activation scores (refer to the pathway activation score for each pathway in the **Supplementary**
8 **File S2**). To reduce the number of dimensions we performed feature selection through recursive
9 elimination (**Table S4**). Based on those results made a predictive model using 4 supervised
10 machine learning methods (NB, DT, KNN, and SVM) for selected key pathway scores and
11 calculated the performance. All models outperformed the baseline (all $P < 0.001$) (**Fig. 5A**, left
12 plot) and SVM, the best performing model, had an average performance AUC/AUPR of 0.87/0.78
13 (all $P < 0.001$) (**Fig. 5A**, middle and right plots). The selected key predictors for SVM model were
14 NFκB-, FcεRI-, TCR-, and TNF signaling pathways. Next, models based on expression values of
15 DEG were fit in order to sort out the informative genes and compare their performance with
16 pathway-driven models. DEG-driven models showed superior performance as compared with
17 pathway-driven models (**Fig. 5B**, left plot). The overall AUC of the ROC curves exceeds 0.85 (**Fig.**
18 **5B**, middle and right plots). SVM showed the best performance AUC/AUPR of 0.92/0.86 and with
19 the *HMMR*, *PRPF4B*, *EVI2A*, *RAB27A*, *MALT1*, *SNX6*, and *IFIH1* genes as features. The
20 expression of these genes provide a distinct signature between two different outcomes ($P < 0.05$
21 in all cases, see **Fig. S14**).

1 **Discussion**

2 Here, we built the largest RA compendium made by synovial transcriptomes. DEGs extracted from
3 this compendium encompassed the susceptible genes and target molecules. Their topology in the
4 network has opened new possibilities to elucidate biological roles and offer a cue for existing
5 clinical questions. Unbiased cluster analysis of RA compendium resulted in meaningful categories
6 of RA patients with distinct activity for relevant pathways. The pathway-based analysis allowed
7 refinement in our understanding of RA subgroups and it was also feasible to construct pathway-
8 or DEG-driven predictive model for intended treatment by machine learning methods.

9 Synovial tissues are considerably more difficult samples to obtain, as they are obtained during
10 joint replacement surgery, synovectomy or by arthroscopy at 4-8 sites of the affected joint.
11 However, they are more suitable to understand the mechanism and response to RA, since blood-
12 derived samples are a distant and hence more noisy proxy to the disease, with known quality issues
13 (4-6). Moreover, to refine the RA-unique genes, we compared RA samples with two control sets
14 (OA and NC groups) and adopted the DEGs shared by three independent methods. We found that
15 24 of the DEGs are the known RA-associated genetic loci and take a central position in the synovial
16 network. Since functional implications of risk allele were often obscure, it would be helpful to
17 elucidate the biological mechanisms in which risk alleles operate. *STAT1*, a transcription factor
18 downstream of IFN signaling pathway, highlighted as a key molecule in the previous reports
19 (35,43), was found to be one of the hub genes. Other hub genes, such as *JAK2*, *SYK*, and *BTK* are
20 small molecules that have increasingly drawn attention as novel therapeutic targets following the
21 cytokine-targeting biologics (30). In contrast, molecules such as TNF receptor molecules, *CTLA4*,
22 *IL6R*, and *MS4A1* were located at the functional periphery of the network although drugs against
23 these molecules are widely used in clinical practice. Moreover, these molecules were placed farther

1 from RA-associated genes than *JAK2* and *SYK* in the network, inferring part of their less potent
2 efficacy in active RA. This was in good harmony with a recent clinical trial that baricitinib, an
3 inhibitor of the Janus kinases *JAK1* and *JAK2*, showed a stronger therapeutic effect as compared
4 with ADLM, a TNF inhibitor (44).

5 Biological processes and pathways identified from RA compendium show what is happening in
6 the inflamed synovium of RA and are in good line with the previous studies (5,35,36). It is worthy
7 of note that processes concerning viral cycle and anti-viral response were found to be enriched.
8 This could be the internal process analogous to or the vestige of viral infection such as
9 *Chikungunya* virus (45,46). A series of studies pointed out activation of IFN-related gene
10 signatures in a subset of RA patients and its substantial similarity to viral infection (5,34-37,46)
11 and one reported that the type I IFN signature negatively predicts the clinical response to rituximab
12 treatment in patients with RA (34). Here, our results suggest that such a probable link between the
13 IFN signature and the anti-viral response may exist (46).

14 Interestingly, we were able to identify three distinct subgroups through NMF analysis of the RA
15 compendium and they differed in activation level of RA-relevant signaling pathways (8,20).
16 Various combinations of molecular perturbations might converge to dysregulation of common
17 pathways and lead to the similar phenotype (47). Since combinations of genomic perturbations are
18 variable across the patients, pathway- or module-based approaches are desirable for a better
19 understanding of complex inflammatory disease like RA. We looked at the enriched pathways
20 derived from DEGs, which were commensurate with the pathway activation scores calculated from
21 the whole gene list in the compendium. The RA cluster 1 was weighted toward signals regarding
22 synovial proliferation and tissue remodeling (PI3K-AKT-, TGF β - and Wnt signaling pathways)
23 (38). RA clusters 2 and 3 showed a strong disposition for proinflammatory signaling pathways

1 (Chemokine-, TNF-, TLR- and VEGF signaling pathways). Apoptosis-related pathways (P53- and
2 PI3K-AKT signaling pathway) were much prominent in RA cluster 2 (39,40), while BCR-, Jak-
3 STAT-, NF κ B-, and TCR signaling pathways were stronger in RA cluster 3. It is known that
4 synoviocytes are the main culprit of invasive synovium and quantitative and qualitative activities
5 of synovial macrophage reflect therapeutic efficacy (48,49). They add to the cellular resistance to
6 apoptosis and increase of the potential for proliferation, hence they contribute to the progression
7 and perpetuation of destructive joint inflammation. Therefore, we speculate that an aggressive
8 suppression of pro-inflammatory signals would be better pertinent to RA cluster 3, while
9 therapeutic strategies to control propagation and survival of synoviocytes and macrophages
10 together with anti-inflammatory treatment should be considered in RA cluster 1 and 2 (50) (**Table**
11 **S3**). This insight, together with the candidate gene targets for drug development that we have
12 identified in each cluster, may provide good starting points for delivering precision and
13 personalized treatment.

14 Machine learning has become ubiquitous and indispensable for solving complex problems in most
15 sciences (51). Since the problem of unresolved heterogeneity is prevalent to medicine, the same
16 methods are expected to open up vast new possibilities in medicine and actively employed in a
17 variety of clinical research (51). We tried to make a predictive model for 62 samples that were
18 obtained from the synovial tissue of RA patients before administration of IFXM. Because key
19 features are informative for predicting the outcome rather than being directly implicated in the
20 major pathways or usual suspects related to the RA synovium, they could be different depending
21 on drugs and models. The fact that we achieved high performance scores in RA response prediction
22 from mining the RA compendium, despite this was not attainable through individual statistical
23 techniques in the past (42), argues that similar techniques can guide us to narrow choices for more

1 effective drugs. Interestingly, DEG-driven models outperformed models that were relying on
2 pathways as features. Among 7 featured genes in SVM model, HMMR (Hyaluronan-mediated
3 motility receptor, also known as RHAMM) exacerbated collagen-induced arthritis by supporting
4 cell migration and up-regulating genes involved with inflammation (52) and MATL1 (Mucosa
5 associated lymphoid tissue lymphoma translocation gene 1) was recently identified to play a
6 crucial role in the pathogenesis of RA as MATL1-deficient mice were completely resistant to
7 collage-induced arthritis (53). Direct connection to RA was not revealed for the rest of the
8 identified informative genes so far and it remains to be investigated how and why these features
9 are indicative of drug response.

10 There are some limitations to be addressed in this study. First, removal of batch effects is not ideal
11 which adds to the noise in the compendium. Second, we did not fully address the association of
12 RA subgroup with clinical factors including age, sex, disease duration, and RF positive due to lack
13 of complete annotation for each RA sample. Third, a limited number of samples were treated with
14 other drugs except for IFXM precluded us from making a predictive model. In general, more meta-
15 data would be desired, although this is to be expected as these studies were performed in different
16 clinical environments, with different procedures and goals, which did not include their aggregation
17 to a single compendium and application of advanced machine learning techniques. In the future,
18 we anticipate that the construction of datasets with sufficient metadata for machine learning
19 analysis would enable critical insights and may lead to novel drug targets for RA treatment.

20

References

- 1
- 2
- 3 1. Lee DM, Weinblatt ME. Rheumatoid arthritis. *Lancet* 2001;358:903-11.
- 4 2. Smolen JS, Aletaha D. Rheumatoid arthritis therapy reappraisal: strategies, opportunities
- 5 and challenges. *Nat Rev Rheumatol* 2015;11:276-89.
- 6 3. Burska AN, Roget K, Blits M, Soto Gomez L, van de Loo F, Hazelwood LD, et al. Gene
- 7 expression analysis in RA: towards personalized medicine. *Pharmacogenomics J*
- 8 2014;14:93-106.
- 9 4. You S, Cho CS, Lee I, Hood L, Hwang D, Kim WU. A systems approach to rheumatoid
- 10 arthritis. *PLoS One* 2012;7:e51508.
- 11 5. van Baarsen LG, Wijbrandts CA, Timmer TC, van der Pouw Kraan TC, Tak PP, Verweij
- 12 CL. Synovial tissue heterogeneity in rheumatoid arthritis in relation to disease activity
- 13 and biomarkers in peripheral blood. *Arthritis Rheum* 2010;62:1602-7.
- 14 6. Haupl T, Stuhlmuller B, Grutzkau A, Radbruch A, Burmester GR. Does gene expression
- 15 analysis inform us in rheumatoid arthritis? *Ann Rheum Dis* 2010;69 Suppl 1:i37-42.
- 16 7. Sweeney TE, Braviak L, Tato CM, Khatri P. Genome-wide expression for diagnosis of
- 17 pulmonary tuberculosis: a multicohort analysis. *Lancet Respir Med* 2016;4:213-24.
- 18 8. You S, Knudsen BS, Erho N, Alshalalfa M, Takhar M, Al-Deen Ashab H, et al.
- 19 Integrated Classification of Prostate Cancer Reveals a Novel Luminal Subtype with Poor
- 20 Outcome. *Cancer Res* 2016;76:4948-58.
- 21 9. van Gestel AM, Prevoo ML, van 't Hof MA, van Rijswijk MH, van de Putte LB, van Riel
- 22 PL. Development and validation of the European League Against Rheumatism response
- 23 criteria for rheumatoid arthritis. Comparison with the preliminary American College of

- 1 Rheumatology and the World Health Organization/International League Against
2 Rheumatism Criteria. *Arthritis Rheum* 1996;39:34-40.
- 3 10. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large
4 gene lists using DAVID bioinformatics resources. *Nat Protoc* 2009;4:44-57.
- 5 11. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al.
6 Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide
7 expression profiles. *Proc Natl Acad Sci U S A* 2005;102:15545-50.
- 8 12. Merico D, Isserlin R, Stueker O, Emili A, Bader GD. Enrichment map: a network-based
9 method for gene-set enrichment visualization and interpretation. *PLoS One*
10 2010;5:e13984.
- 11 13. Stark C, Breitkreutz BJ, Reguly T, Boucher L, Breitkreutz A, Tyers M. BioGRID: a
12 general repository for interaction datasets. *Nucleic Acids Res* 2006;34:D535-9.
- 13 14. Peri S, Navarro JD, Kristiansen TZ, Amanchy R, Surendranath V, Muthusamy B, et al.
14 Human protein reference database as a discovery resource for proteomics. *Nucleic Acids*
15 *Res* 2004;32:D497-501.
- 16 15. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, et al. The
17 MIntAct project--IntAct as a common curation platform for 11 molecular interaction
18 databases. *Nucleic Acids Res* 2014;42:D358-63.
- 19 16. Fabregat A, Sidiropoulos K, Garapati P, Gillespie M, Hausmann K, Haw R, et al. The
20 Reactome pathway Knowledgebase. *Nucleic Acids Res* 2016;44:D481-7.
- 21 17. Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, et al. The STRING
22 database in 2017: quality-controlled protein-protein association networks, made broadly
23 accessible. *Nucleic Acids Res* 2017;45:D362-d8.

- 1 18. Koschutski D, Schreiber F. Centrality analysis methods for biological networks and their
2 application to gene regulatory networks. *Gene Regul Syst Bio* 2008;2:193-201.
- 3 19. Laenen G, Thorrez L, Bornigen D, Moreau Y. Finding the targets of a drug by integration
4 of gene expression data with a protein interaction network. *Mol Biosyst* 2013;9:1676-85.
- 5 20. Brunet JP, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern discovery
6 using matrix factorization. *Proc Natl Acad Sci U S A* 2004;101:4164-9.
- 7 21. Maaten LVD, Hinton GE. Visualizing Data using t-SNE. *J Machine Learning Res*
8 2008;9:2579-605.
- 9 22. Levine DM, Haynor DR, Castle JC, Stepaniants SB, Pellegrini M, Mao M, et al. Pathway
10 and gene-set activation measurement from mRNA expression data: the tissue distribution
11 of human pathways. *Genome Biol* 2006;7:R93.
- 12 23. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new perspectives
13 on genomes, pathways, diseases and drugs. *Nucleic Acids Res* 2017;45:D353-d61.
- 14 24. James G, Witten D, Hastie T, Tibshirani R. An introduction to statistical learning:
15 Springer; 2013.
- 16 25. Kuhn M, Johnson K. Applied predictive modeling: Springer; 2013.
- 17 26. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot
18 when evaluating binary classifiers on imbalanced datasets. *PLoS One* 2015;10:e0118432.
- 19 27. Brodersen KH, Ong CS, Stephan KE, Buhmann JM. The Balanced Accuracy and Its
20 Posterior Distribution. 2010 20th International Conference on Pattern Recognition 2010.
21 p. 3121-4.
- 22 28. Yamamoto K, Okada Y, Suzuki A, Kochi Y. Genetics of rheumatoid arthritis in Asia--
23 present and future. *Nat Rev Rheumatol* 2015;11:375-9.

- 1 29. Koenders MI, van den Berg WB. Novel therapeutic targets in rheumatoid arthritis.
2 Trends Pharmacol Sci 2015;36:189-95.
- 3 30. Kelly V, Genovese M. Novel small molecule therapeutics in rheumatoid arthritis.
4 Rheumatology (Oxford) 2013;52:1155-62.
- 5 31. Pinero J, Bravo A, Queralt-Rosinach N, Gutierrez-Sacristan A, Deu-Pons J, Centeno E, et
6 al. DisGeNET: a comprehensive platform integrating information on human disease-
7 associated genes and variants. Nucleic Acids Res 2017;45:D833-d9.
- 8 32. Choy E. Understanding the dynamics: pathways involved in the pathogenesis of
9 rheumatoid arthritis. Rheumatology (Oxford) 2012;51 Suppl 5:v3-11.
- 10 33. Rabelo Fde S, da Mota LM, Lima RA, Lima FA, Barra GB, de Carvalho JF, et al. The
11 Wnt signaling pathway and rheumatoid arthritis. Autoimmun Rev 2010;9:207-10.
- 12 34. Thurlings RM, Boumans M, Tekstra J, van Roon JA, Vos K, van Westing DM, et al.
13 Relationship between the type I interferon signature and the response to rituximab in
14 rheumatoid arthritis patients. Arthritis Rheum 2010;62:3607-14.
- 15 35. van der Pouw Kraan TC, van Gaalen FA, Kasperkovitz PV, Verbeet NL, Smeets TJ,
16 Kraan MC, et al. Rheumatoid arthritis is a heterogeneous disease: evidence for
17 differences in the activation of the STAT-1 pathway between rheumatoid tissues.
18 Arthritis Rheum 2003;48:2132-45.
- 19 36. Woetzel D, Huber R, Kupfer P, Pohlers D, Pfaff M, Driesch D, et al. Identification of
20 rheumatoid arthritis and osteoarthritis patients by transcriptome-based rule set generation.
21 Arthritis Res Ther 2014;16:R84.
- 22 37. van der Pouw Kraan TC, Wijbrandts CA, van Baarsen LG, Voskuyl AE, Rustenburg F,
23 Baggen JM, et al. Rheumatoid arthritis subtypes identified by genomic profiling of

- 1 peripheral blood cells: assignment of a type I interferon signature in a subpopulation of
2 patients. *Ann Rheum Dis* 2007;66:1008-14.
- 3 38. Miao CG, Yang YY, He X, Li XF, Huang C, Huang Y, et al. Wnt signaling pathway in
4 rheumatoid arthritis, with special emphasis on the different roles in synovial
5 inflammation and bone remodeling. *Cell Signal* 2013;25:2069-78.
- 6 39. Pope RM. Apoptosis as a therapeutic tool in rheumatoid arthritis. *Nat Rev Immunol*
7 2002;2:527-35.
- 8 40. Smith MD, Walker JG. Apoptosis a relevant therapeutic target in rheumatoid arthritis?
9 *Rheumatology (Oxford)* 2004;43:405-7.
- 10 41. Ducreux J, Durez P, Galant C, Nzeusseu Toukap A, Van den Eynde B, Houssiau FA, et
11 al. Global molecular effects of tocilizumab therapy in rheumatoid arthritis synovium.
12 *Arthritis Rheumatol* 2014;66:15-23.
- 13 42. Lindberg J, Wijbrandts CA, van Baarsen LG, Nader G, Klareskog L, Catrina A, et al. The
14 gene expression profile in the synovium as a predictor of the clinical response to
15 infliximab treatment in rheumatoid arthritis. *PLoS One* 2010;5:e11310.
- 16 43. Yoshida S, Arakawa F, Higuchi F, Ishibashi Y, Goto M, Sugita Y, et al. Gene expression
17 analysis of rheumatoid arthritis synovial lining regions by cDNA microarray combined
18 with laser microdissection: up-regulation of inflammation-associated STAT1, IRF1,
19 CXCL9, CXCL10, and CCL5. *Scand J Rheumatol* 2012;41:170-9.
- 20 44. Taylor PC, Keystone EC, van der Heijde D, Weinblatt ME, Del Carmen Morales L,
21 Reyes Gonzaga J, et al. Baricitinib versus Placebo or Adalimumab in Rheumatoid
22 Arthritis. *N Engl J Med* 2017;376:652-62.

- 1 45. Miner JJ, Aw Yeang HX, Fox JM, Taffner S, Malkova ON, Oh ST, et al. Chikungunya
2 viral arthritis in the United States: a mimic of seronegative rheumatoid arthritis. *Arthritis*
3 *Rheumatol* 2015;67:1214-20.
- 4 46. Nakaya HI, Gardner J, Poo YS, Major L, Pulendran B, Suhrbier A. Gene profiling of
5 Chikungunya virus arthritis in a mouse model reveals significant overlap with rheumatoid
6 arthritis. *Arthritis Rheum* 2012;64:3553-63.
- 7 47. Kim YA, Wuchty S, Przytycka TM. Identifying causal genes and dysregulated pathways
8 in complex diseases. *PLoS Comput Biol* 2011;7:e1001095.
- 9 48. Bottini N, Firestein GS. Duality of fibroblast-like synoviocytes in RA: passive responders
10 and imprinted aggressors. *Nat Rev Rheumatol* 2013;9:24-33.
- 11 49. Haringman JJ, Gerlag DM, Zwinderman AH, Smeets TJ, Kraan MC, Baeten D, et al.
12 Synovial tissue macrophages: a sensitive biomarker for response to treatment in patients
13 with rheumatoid arthritis. *Ann Rheum Dis* 2005;64:834-8.
- 14 50. Martinez-Lostao L, Garcia-Alvarez F, Basanez G, Alegre-Aguaron E, Desportes P,
15 Larrad L, et al. Liposome-bound APO2L/TRAIL is an effective treatment in a rabbit
16 model of rheumatoid arthritis. *Arthritis Rheum* 2010;62:2272-82.
- 17 51. Obermeyer Z, Emanuel EJ. Predicting the Future - Big Data, Machine Learning, and
18 Clinical Medicine. *N Engl J Med* 2016;375:1216-9.
- 19 52. Nedvetzki S, Gonen E, Assayag N, Reich R, Williams RO, Thurmond RL, et al.
20 RHAMM, a receptor for hyaluronan-mediated motility, compensates for CD44 in
21 inflamed CD44-knockout mice: a different interpretation of redundancy. *Proc Natl Acad*
22 *Sci U S A* 2004;101:18081-6.

- 1 53. Gilis E, Staalj J, Beyaert R, Elewaut D. The Paracaspase MALT1 Plays a Central Role in
- 2 the Pathogenesis of Rheumatoid Arthritis [Abstract]. *Arthritis Rheumatology* 2017;69.
- 3
- 4

1 **Author Contributions**

2 K-J Kim, I Tagkopoulos designed the study and K-J Kim acquired the data. K-J Kim, M Kim, I
3 Tagkopoulos contributed to the analysis and interpretation of data. K-J Kim conducted all data
4 analysis and drafted the initial manuscript. All authors were involved in drafting the article or
5 revising it critically for important intellectual content, and all authors approved the final version
6 to be published.

7

8 **Additional Information**

9 **Competing interests**

10 The authors declare that they have no competing interests.

11

12

1 **Figure Legends**

2

3 **Figure 1. Overview of the data processing steps.** (A) Twenty studies maximally covering 20,511
4 genes were retrieved from the literature. (B) Selected were 11 datasets adequate to integrated
5 analysis, which included 256 RA, 41 OA, and 36 NC samples covering 11,769 gene. (C) The
6 merged dataset was normalized using quantile method and its batch effect was corrected. (D) DEG
7 of RA compared to OA or NC were obtained using three methods, eBayes, SAM, and RP.
8 Intersection of three DEG sets was chosen as significant DEG. The number of DEG was 2762 in
9 RA versus OA and 3087 in RA versus NC. (E) A list of strategies for integrated analysis.
10 (Abbreviation: RA, rheumatoid arthritis; OA, osteoarthritis; NC, normal controls; DEG,
11 differentially expressed genes; eBayes, empirical Bayes; SAM, significance analysis of microarray;
12 RP, rank products).

13

14 **Figure 2. Differentially expressed genes and their functional network.** (A) Venn diagram
15 showing the overlap of up- and down-regulated DEG between RA versus OA and RA versus NC.
16 (B) Gene-Set enrichment map for up-regulated DEG. Nodes represent GO-termed gene-sets. Their
17 color intensity and size is proportional to the enrichment significance and the gene size,
18 respectively. Edge thickness represents the degree of overlap between gene sets and only edges
19 with a Jaccard coefficient larger than 0.25 were visualized. Clusters of functionally related gene-
20 sets were manually curated based on the GO parent-child hierarchy and assigned a label. (C)
21 Protein-Protein interaction network of up-regulated DEG. Red and blue nodes indicate the known
22 RA-susceptible genes and drug target molecules, respectively. Drug targets were defined subject
23 to the targets of drugs currently in use or under clinical trial and development. Yellow nodes
24 correspond to the hub molecules, which are determined as the shared genes in top 10% with the

1 highest rank in each arm of three centrality parameters; degree, closeness, and betweenness.
2 Orange, green, and purple colored-nodes are the overlapped between red and yellow, yellow and
3 blue, and red and blue ones, respectively. Right-side inset box is the schematic diagram of the
4 interesting genes.

5
6 **Figure 3. Identification of novel RA subgroups according to synovial signatures. (A)**
7 Reordered consensus matrices on RA compendium. The samples were clustered using average
8 linkage and 1-correlation distances. Deep-red color indicates perfect agreement of the solution,
9 whilst blue color indicates no agreement (Right-side color bar). Basis and consensus represent
10 clusters based on the basis and consensus matrices, respectively. The silhouette score is a similarity
11 measure within its own cluster compared to other clusters. **(B)** t-SNE (upper plot) and PCA (lower
12 plot) reduces the dimensions of a multivariate dataset. Each data point is assigned a location in a
13 two-dimensional map to illustrate potential clusters of neighboring samples, which contain similar
14 gene expression patterns. **(C)** KEGG pathways enriched by the clustered DEG. Heatmap with red
15 gradient represents the level of fold enrichment for each KEGG pathway, which was determined
16 by DAVID software.

17
18 **Figure 4. Pathway activation scores according to RA subgroups.** Chord diagram shows
19 interrelationship among pathways and link thickness is proportional to the overlap between two
20 pathways, calculated using the Jaccard coefficients. Turkey boxplots reveals pathway activation
21 scores across the RA subgroups and ANOVA test was used to analyze the differences among
22 groups. *, $P < 0.05$; **, $P < 0.01$; ***, $P < 0.001$.

23

1 **Figure 5. Predictive models and their performance. (A) Pathway-driven models. (B) DEG-**
2 **driven models.** (Left plot) The training and testing balanced accuracy for each classifier as
3 compared with the baseline. All models outperformed the baseline (all $P < 0.001$) and the
4 performance of the trained models was significantly compromised in testing sets (all $P < 0.001$).
5 (Middle and right plots) Averaged ROC and PR curves showing the performance of each classifier.
6
7
8

Fig. 1. Overview of data processing steps

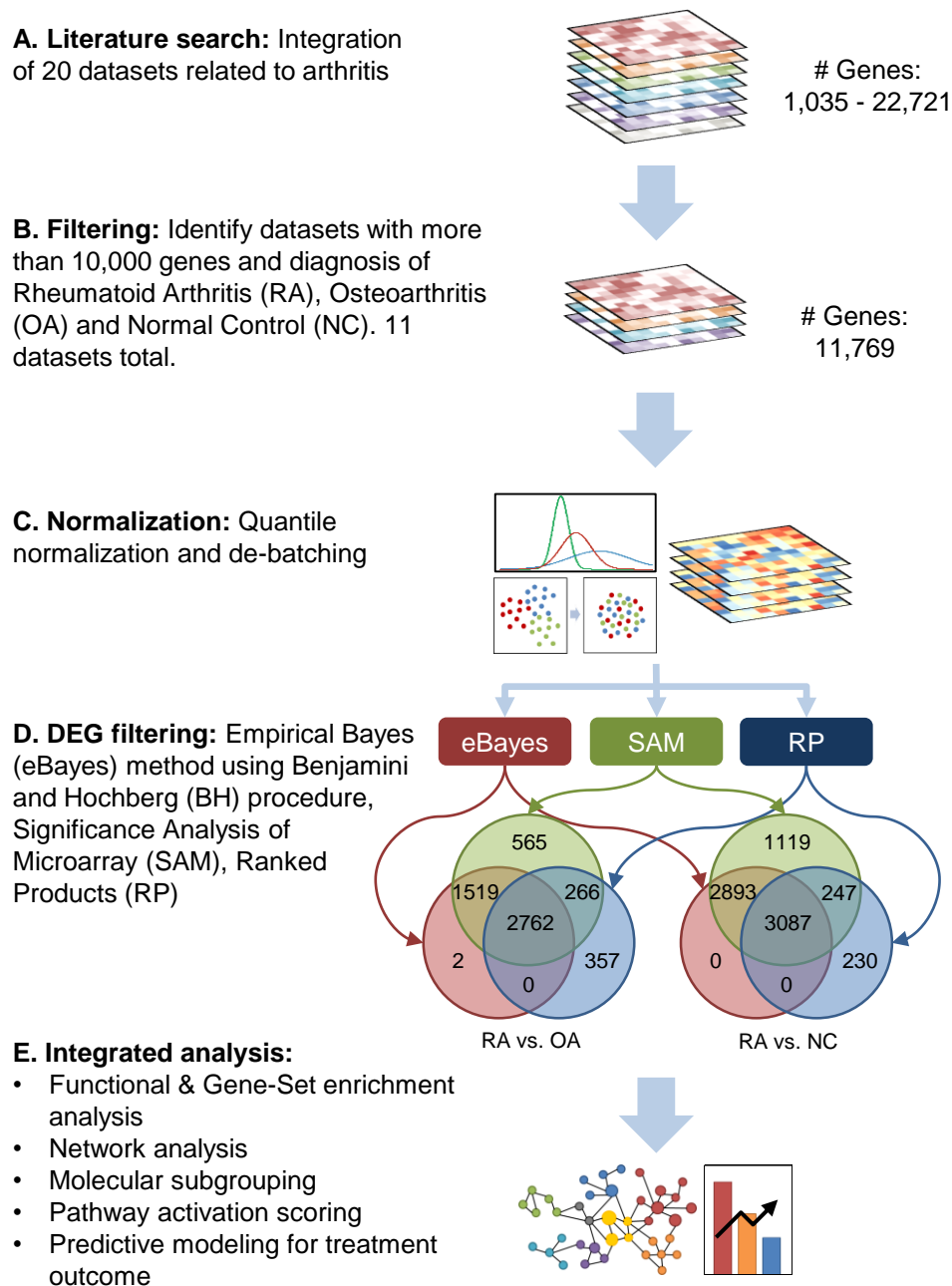
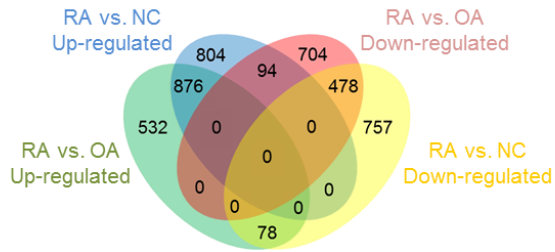
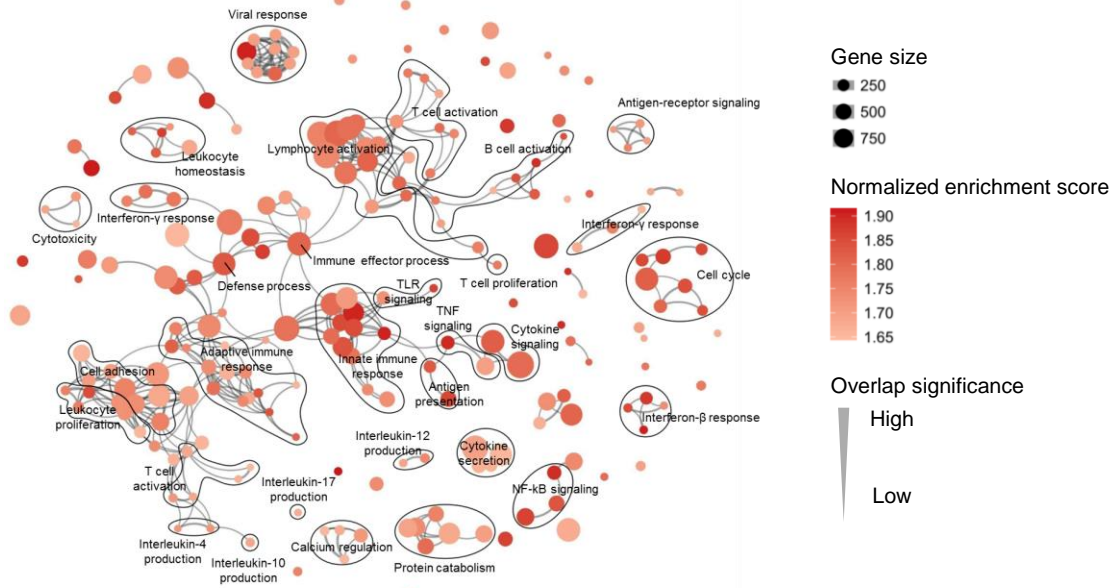


Fig. 2. Differentially expressed genes and their functional network

A



B



C

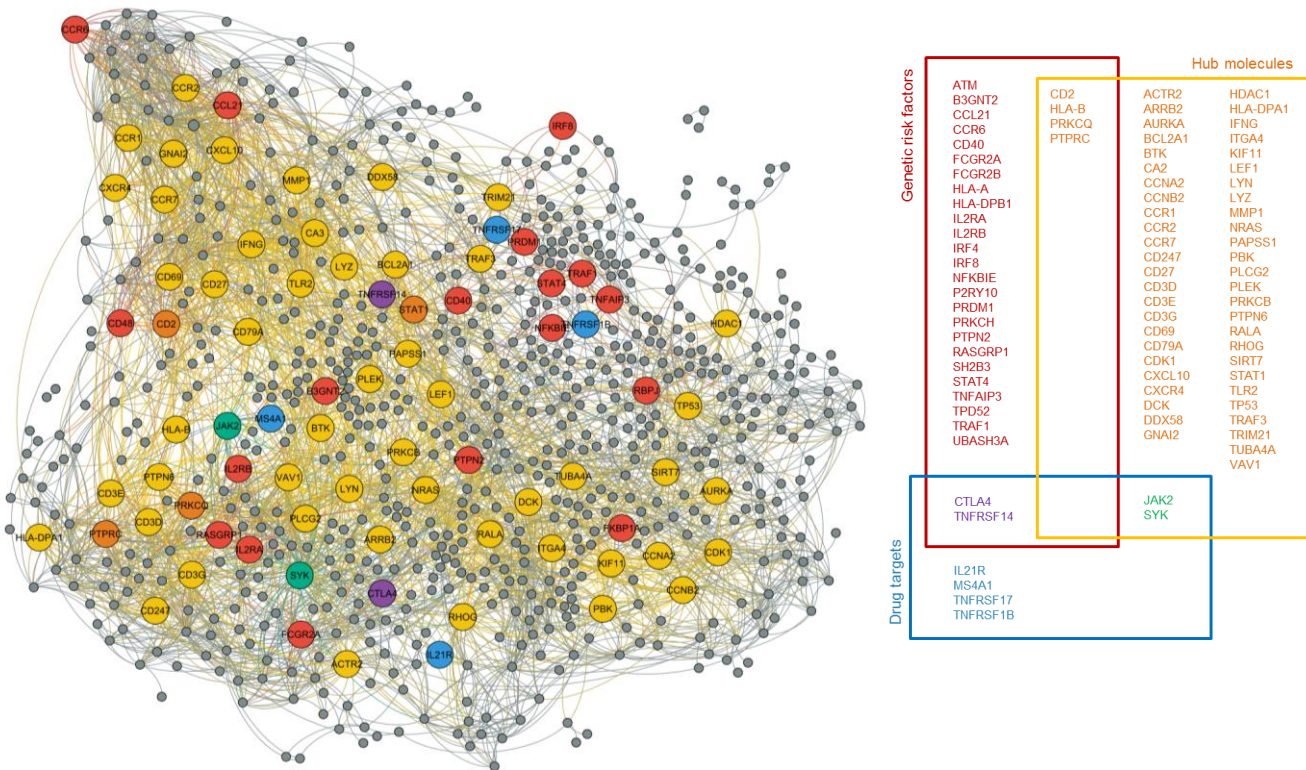
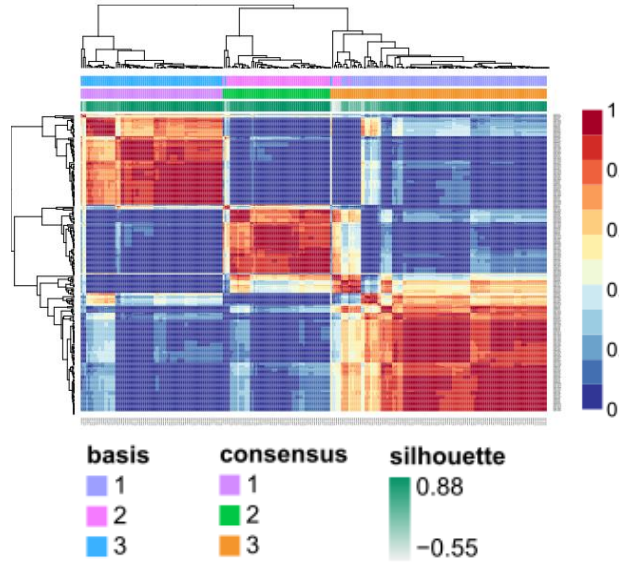
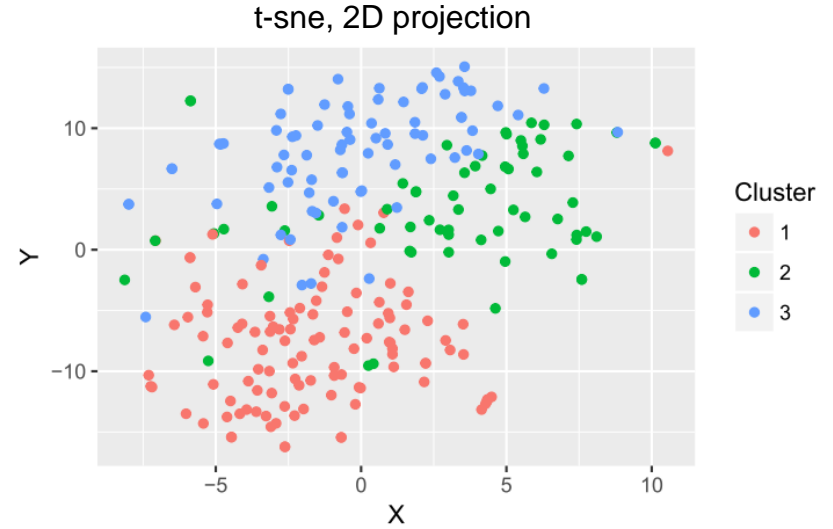


Fig. 3. Identification of novel RA subgroups according to synovial gene signatures.

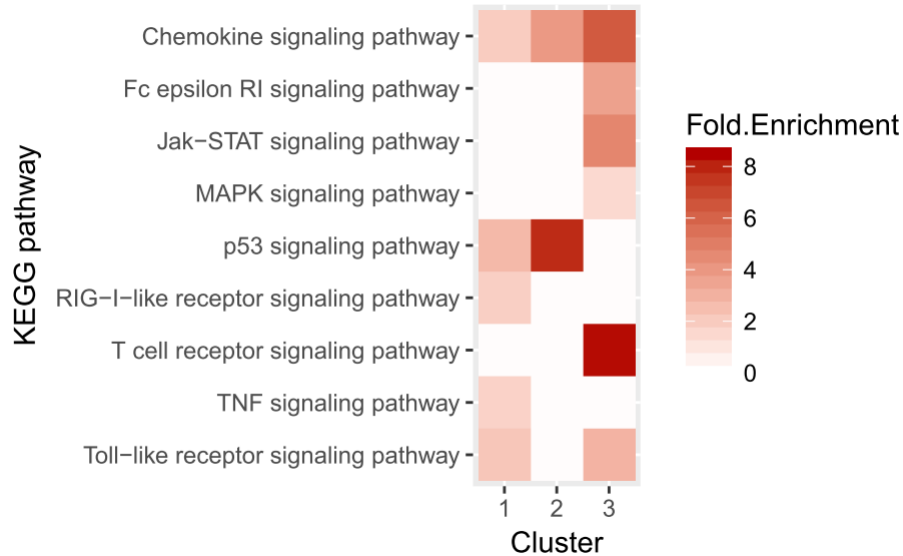
A



B



C



PCA

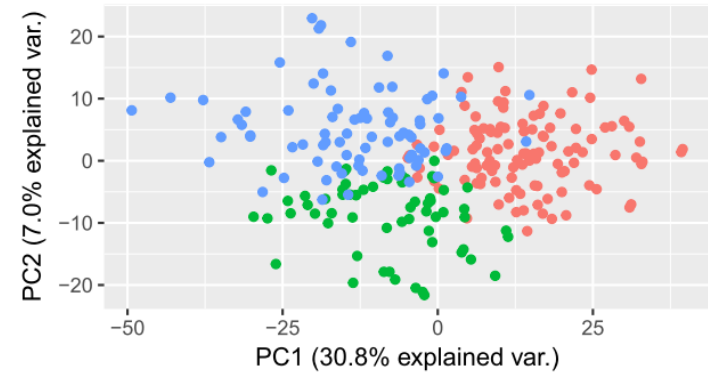


Fig. 4. Pathway activation scores according to RA subgroups.

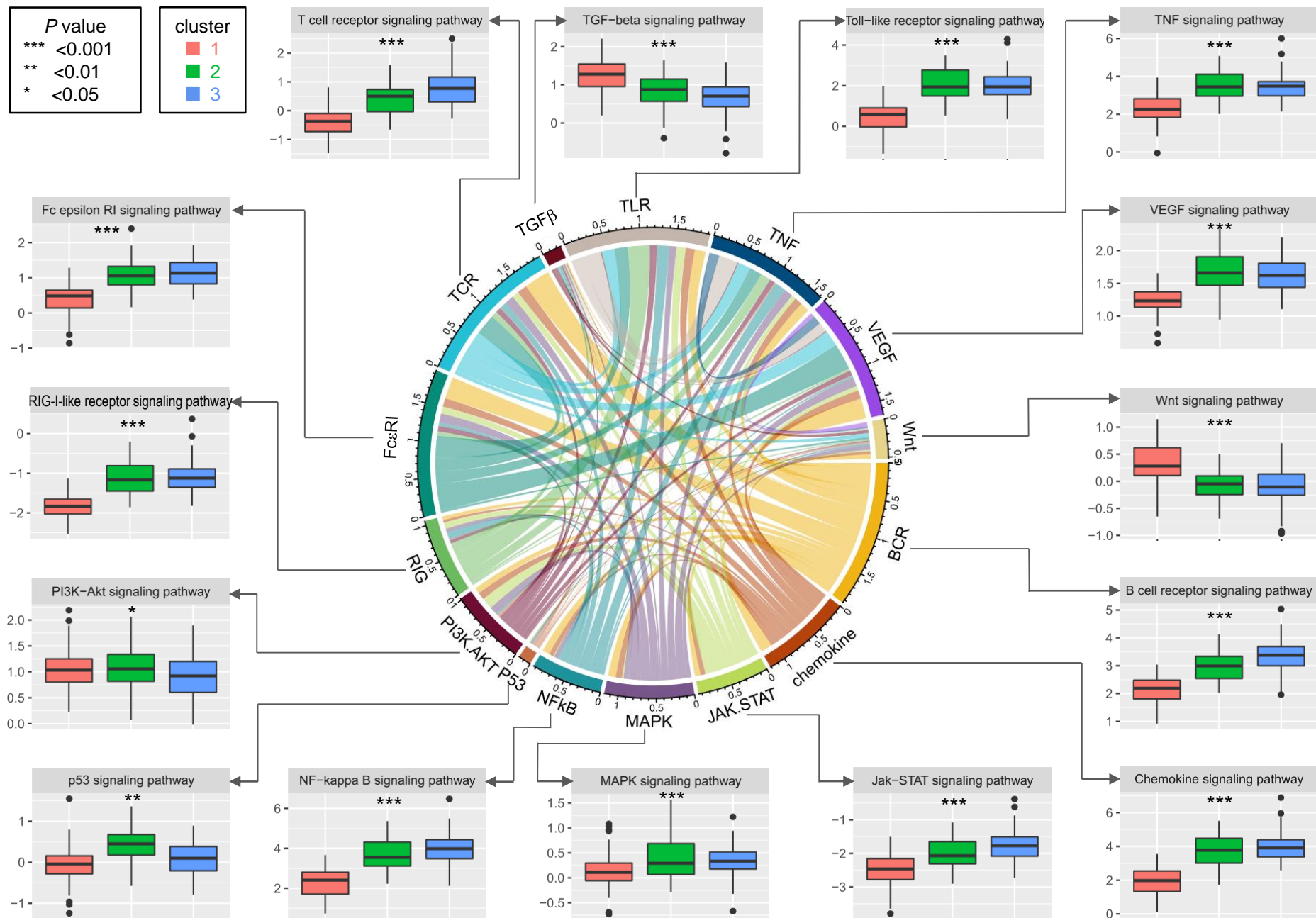
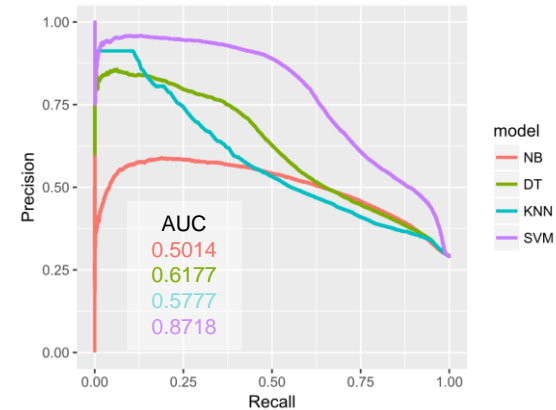
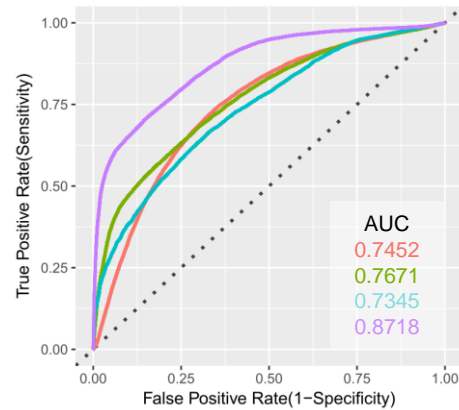
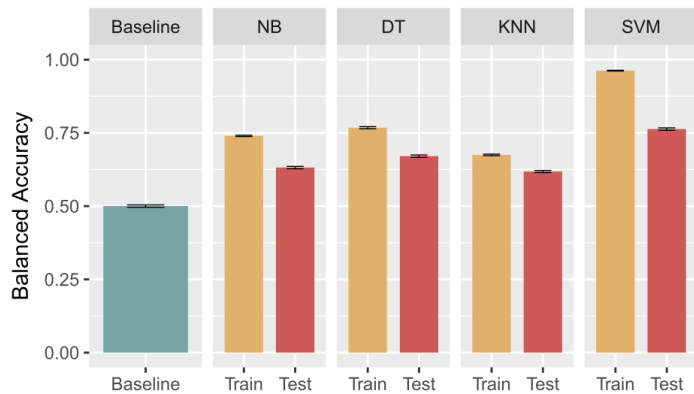


Fig. 5. Predictive models and their performance

A



B

