**BMC Genomics**

# A validated generally applicable approach using the systematic assessment of disease modules by GWAS reveals a multi-omic module strongly associated with risk factors in multiple sclerosis

Tejaswi V. S. Badam[1,2†], Hendrik A. de Weerd[1,2†], David Martínez-Enguita[2], Tomas Olsson[3], Lars Alfredsson[3,4], Ingrid Kockum[3], Maja Jagodic[3], Zelmina Lubovac-Pilav[1†] and Mika Gustafsson[2*†]

## Abstract

**Background:** There exist few, if any, practical guidelines for predictive and falsifiable multi-omic data integration that systematically integrate existing knowledge. Disease modules are popular concepts for interpreting genome-wide studies in medicine but have so far not been systematically evaluated and may lead to corroborating multi-omic modules.

**Result:** We assessed eight module identification methods in 57 previously published expression and methylation studies of 19 diseases using GWAS enrichment analysis. Next, we applied the same strategy for multi-omic integration of 20 datasets of multiple sclerosis (MS), and further validated the resulting module using both GWAS and risk-factor-associated genes from several independent cohorts. Our benchmark of modules showed that in immune-associated diseases modules inferred from clique-based methods were the most enriched for GWAS genes. The multi-omic case study using MS data revealed the robust identification of a module of 220 genes. Strikingly, most genes of the module were differentially methylated upon the action of one or several environmental risk factors in MS ($n = 217$, $P = 10^{-47}$) and were also independently validated for association with five different risk factors of MS, which further stressed the high genetic and epigenetic relevance of the module for MS.

**Conclusions:** We believe our analysis provides a workflow for selecting modules and our benchmark study may help further improvement of disease module methods. Moreover, we also stress that our methodology is generally applicable for combining and assessing the performance of multi-omic approaches for complex diseases.

**Keywords:** Benchmark, Multi-omics, Network modules, Multiple sclerosis, Risk factors, Disease modules, Network analysis, Protein network analysis, Transcriptomics, Methylomics, Data integration, Genome-wide association analysis

* Correspondence: mika.gustafsson@liu.se
†Zelmina Lubovac-Pilav and Mika Gustafsson share senior authorship.
²Bioinformatics, Department of Physics, Chemistry and Biology, Linköping university, Linköping, Sweden
Full list of author information is available at the end of the article

Badam *et al. BMC Genomics*      (2021) 22:631

Page 2 of 13

## Summary
Our benchmark of multi-omic modules and validated translational systems medicine workflow for dissecting complex diseases resulted in multi-omic module of 220 genes highly enriched for risk factors associated with multiple sclerosis.

## Background
Complex diseases are the result of disruptions of many interconnected multimolecular pathways, reflected in multiple omic layers of regulation of cellular function, rather than perturbations of a single gene or protein [1]. Systems and network medicine aim to translate observed omic differences in patients using networks, to personalize medicine [2]. Importantly, genes that are associated with diseases are more likely to interact with each other rather than with non-disease associated genes, forming multi-omic network disease modules [3, 4]. Owing to the incompleteness of the underlying multi-omic interactions, the networks are often modeled as effective gene-gene interactions, using for example STRING database [5]. Thus, network modules might be ideal tools for multi-omic analysis. However, the evaluation of performance of different module inference methods remains a poorly understood topic, which creates the need for transparent evaluation of these methods based on objective benchmarks across various diseases and omics. Genomic concordance has been suggested as a multi-omic validation principle [4, 6], i.e., modules derived from one omic, such as gene expression or DNA methylation should be enriched for disease-associated single nucleotide polymorphisms (SNPs).

The variety of algorithms that have been proposed and applied for identification of disease modules can be categorized into two main groups. On the one hand, there are methods which rely purely on clustering of the genes in relevant disease networks [7]. On the other hand, there are algorithms which make use of disease-associated molecules or genetic loci to reveal disease modules that correlate with disease function, such as the disease module detection (DIAMOnD) algorithm [8], clique-based methods [9, 10] and weighted gene co-expression network analysis (WGCNA) [11]. The data-derived information can either be differentially expressed genes or differentially correlated or co-expressed genes. Methods following the former approach were recently benchmarked by a metric utilizing genomic concordance within the DREAM consortia [6]. However, so far, algorithms from the latter group have not been benchmarked.
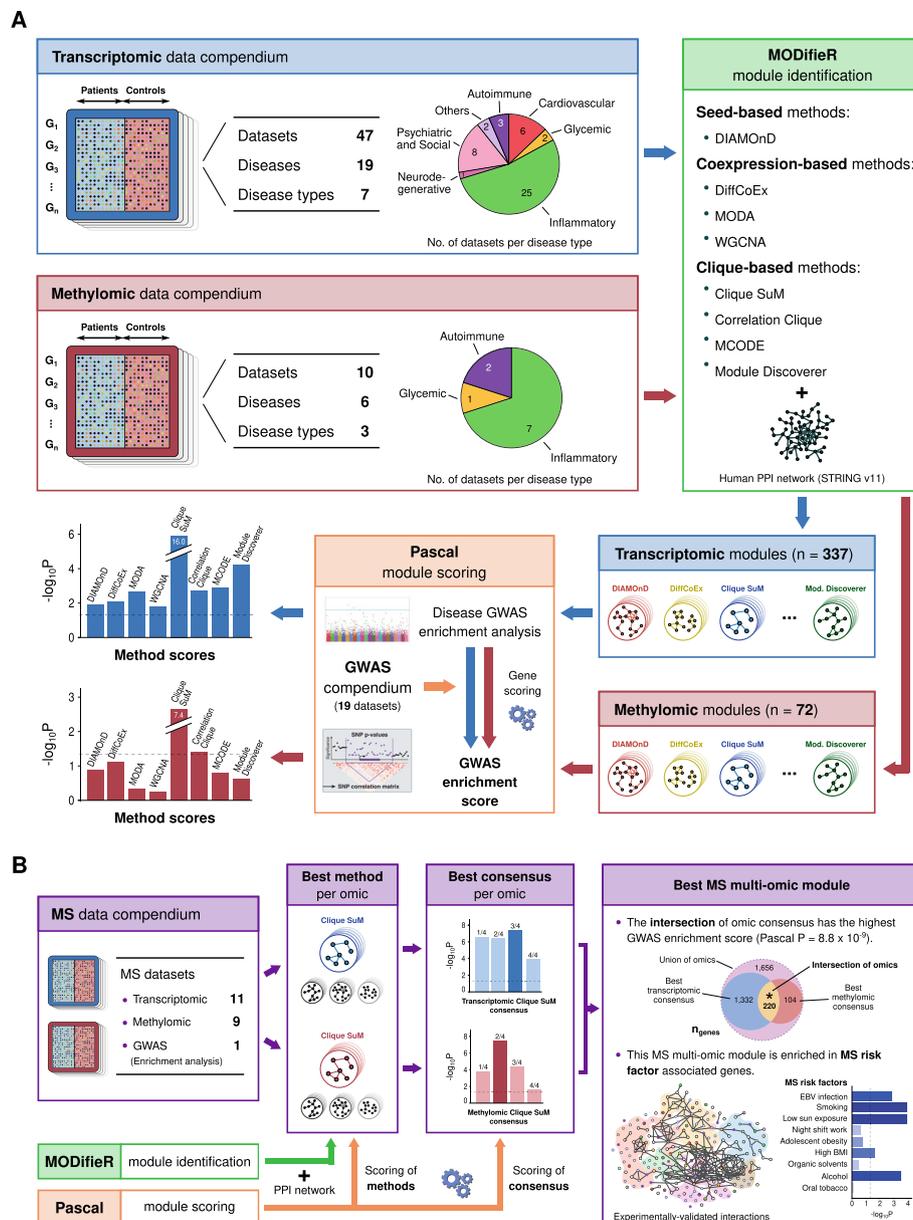
In this study we analyzed, assessed, and compared the performance of eight of the most popular methods for disease module analysis using the R package MODifieR [12] on 19 different diseases including 47 expression and ten methylation datasets. We assessed the performance of the methods using genome-wide association (GWAS) enrichment analysis from the summary statistics of all assayed SNPs similarly as in DREAM [6]. The resulting workflow provided a systematic procedure for selecting the best method for each disease and set the stage for method development in the disease module area. Moreover, it allowed the predictive assessment of combining multiple datasets across several omics using GWAS, which we tested in multiple sclerosis (MS), a heterogeneous complex disease. Briefly, we derived multi-omic modules in a stepwise optimization of GWAS enrichment from transcriptomic and methylomic analyses of MS. We further evaluated the identified multi-omic MS module of 220 genes for its enrichment across DNA methylation studies of eight known lifestyle-associated risk factors of MS. Additionally, we validated the identified significant enrichment risk factors in an independent DNA methylation MS study which indeed showed a very strong and significant MS enrichment for both module genes and risk factor associations. In summary, we provide a robust multi-omic strategy that can be used to disentangle networks of affected genes in complex diseases from both genetic and environmental levels.

## Results
### A benchmark comparing 337 transcriptionally derived disease modules from 19 different diseases
We compiled a benchmark source of disease modules and summary statistics of GWAS datasets from 19 well-powered case-control studies (Supplementary Table 1), some of which were previously used in the DREAM topological disease module challenge [6]. For these datasets we assessed modules using the same metric as in the recent DREAM study [6], based on the pathway scoring algorithm (Pascal) [13]. For each disease we compiled one to five publicly available transcriptomic datasets considering both easily assessable tissues (e.g., blood) and target tissues, thereby covering 47 transcriptomic datasets in total (Fig. 1a). Modules were created using eight different methods from MODifieR [12] and as underlying network we used 631,782 high confidence interactions from STRING database [5] (see methods; complete results are found in Additional file 1: Table S4, S5). In addition, we also tested if genes detected by several methods, hereafter called consensus module genes, had higher enrichment scores than single-method module genes. Enrichment scores for the non-empty modules ($n = 337$) from this analysis were summarized for each method and dataset (Fig. 2a). In total, we found significantly GWAS-enriched modules in 17.8% (60/337) of the single-method modules and 25.5% (12/47) of the non-empty consensus modules that combined at least

**Fig. 1** Overview of the benchmark assessment of disease modules and the integration workflow for MS. (**a**) Transcriptomic and methylomic datasets from 19 different diseases were used as inputs for eight MODifieR module identification methods. The resulting single-omic disease modules ($n = 456$) were independently assessed by GWAS enrichment analysis of the same disease using Pascal module scoring. MODifieR methods were evaluated by the combined enrichment score of their respective disease modules. (**b**) Multi-omic integrative workflow for multiple sclerosis (MS)-associated modules. Data from 20 case-control comparisons were used as input for module detection with MODifieR methods. Clique SuM modules presented the highest GWAS enrichment score and were therefore used to generate single-omic consensus modules. The intersection of the best transcriptomic and methylomic consensus modules resulted in an MS multi-omic module ($n = 220$ genes) with the highest GWAS enrichment, which was independently found to be enriched for genes associated with five known lifestyle MS risk factors using public omic data from healthy individuals

three methods as a criterion. These numbers seemed higher than expected, which might have been a consequence of the same GWAS being used to evaluate multiple transcriptomic datasets of the same disease. Hence, we aggregated scores of the same disease and method as meta-*P*-values (see Methods). Out of the 152 possible disease-method combinations, 18% of the pairs showed a significant GWAS Pascal enrichment, which is more than expected by chance ($n = 27$, $P = 1.0 \times 10^{-8}$). The most enriched method was Clique SuM, which showed significant enrichment in seven out of 19 diseases (binomial test $P = 2.3 \times 10^{-5}$). Many methods exhibited strong

**Fig. 2** Genomic concordance of MODifieR modules on transcriptomic datasets. (**a**) Heatmap of PASCAL *p*-values for eight single-method and eight consensus MODifieR modules, identified for 47 publicly available transcriptomic datasets. Module performance *P*-values are shown in a white to blue scale, where any shade of blue represents a significant module (< 0.05; the darker, the more significant), white represents a non-significant module, and grey represents a module of size zero. Datasets are classified into six disease types: cardiovascular (red), glycemic (golden), inflammatory (green), neurodegenerative (fuchsia), psychiatric and social (pink), autoimmune (dark purple), and others (light purple); and two cell types: blood (maroon), and others (light yellow). Datasets are ranked by meta-*P*-values using Fisher's method of the single-method module P-values across and within their disease types (dataset score, bottom boxplot). MODifieR methods are organized by algorithm type: seed-based (green), co-expression-based (yellow), and clique-based (red), plus the consensus modules (blue). Single-methods and consensus were scored by meta-P-values across datasets (method score, right boxplot). Consensus x/8 indicates that the module genes are found in at least x methods out of eight. (**b**) Scatter plot showing Spearman correlation between module score and betweenness centrality. Modules are represented with a different shape depending on their method and colored based on the disease type. (**c**) Scatter plot showing Spearman correlation between module score and module size. Modules are represented with a different shape depending on their method and colored based on the disease type

enrichments in coronary artery disease (CAD), type 2 diabetes, multiple sclerosis (MS), rheumatoid arthritis (RA), and the inflammatory bowel diseases (IBD) ulcerative colitis (UC) and Crohn's disease (CD), while no significant enrichments were found for asthma, hepatitis C, type 1 diabetes, narcolepsy, Parkinson's disease, or for any psychiatric and social diseases. If we instead ranked methods based on their respective module GWAS

enrichment, Clique SuM was again the most enriched method, with significant associations found for 34% (16/47) of its modules, corresponding to seven different diseases, followed by the consensus modules identified by two out of three methods. Lastly, DIAMOnD and co-expression-based methods all achieved significant results, although worse than Clique SuM. To test the sensitivity of our results to the utilized backbone network,

we recomputed modules for the methods utilizing prior networks using smaller network datasets, namely 1) 120,000 only experimentally verified interactions from STRING-db and 2) 27,719 curated interactions from Reactome database. Analyzing these results showed on average smaller modules and some difference in method performances, e.g., MCODE scoring very good on the Reactome network but poorer on the others (Additional file: Table S4). However, in general these scorings correlated well for these three datasets (Spearman rho in the range 0.18–0.26, with $7.2 \times 10^{-5} > P > 6.8 \times 10^{-3}$) and analyzing the rankings of the methods on their worst dataset showed similar rankings as above. Next, we tested the impact of network centrality and module size as potential confounding factors of the applied performance metric. We found a significant but very modest correlation for module size (Fig. 2c, Spearman rho = 0.165, $P = 2.3 \times 10^{-3}$) and a non-significant correlation for interactome centrality (Fig. 2b, rho = 0.068, $P = 0.21$). Thus, it is meaningful to compare results with differences in those module properties. In summary, we found that the Clique SuM method resulted in the highest disease enrichment for most diseases, while not producing significant modules for others, such as type 2 diabetes, where co-expression-based methods and DIAMOnD scored best. In general, we observed stronger enrichments for cardiovascular and inflammatory diseases, and weaker results for psychiatric and social diseases. Considering that the transcriptomic modules showed that Clique SuM was the best performing method, and that the cardiovascular and inflammatory diseases were the most enriched within the Clique SuM modules, we wanted to test whether this was true for methylomic data as well.
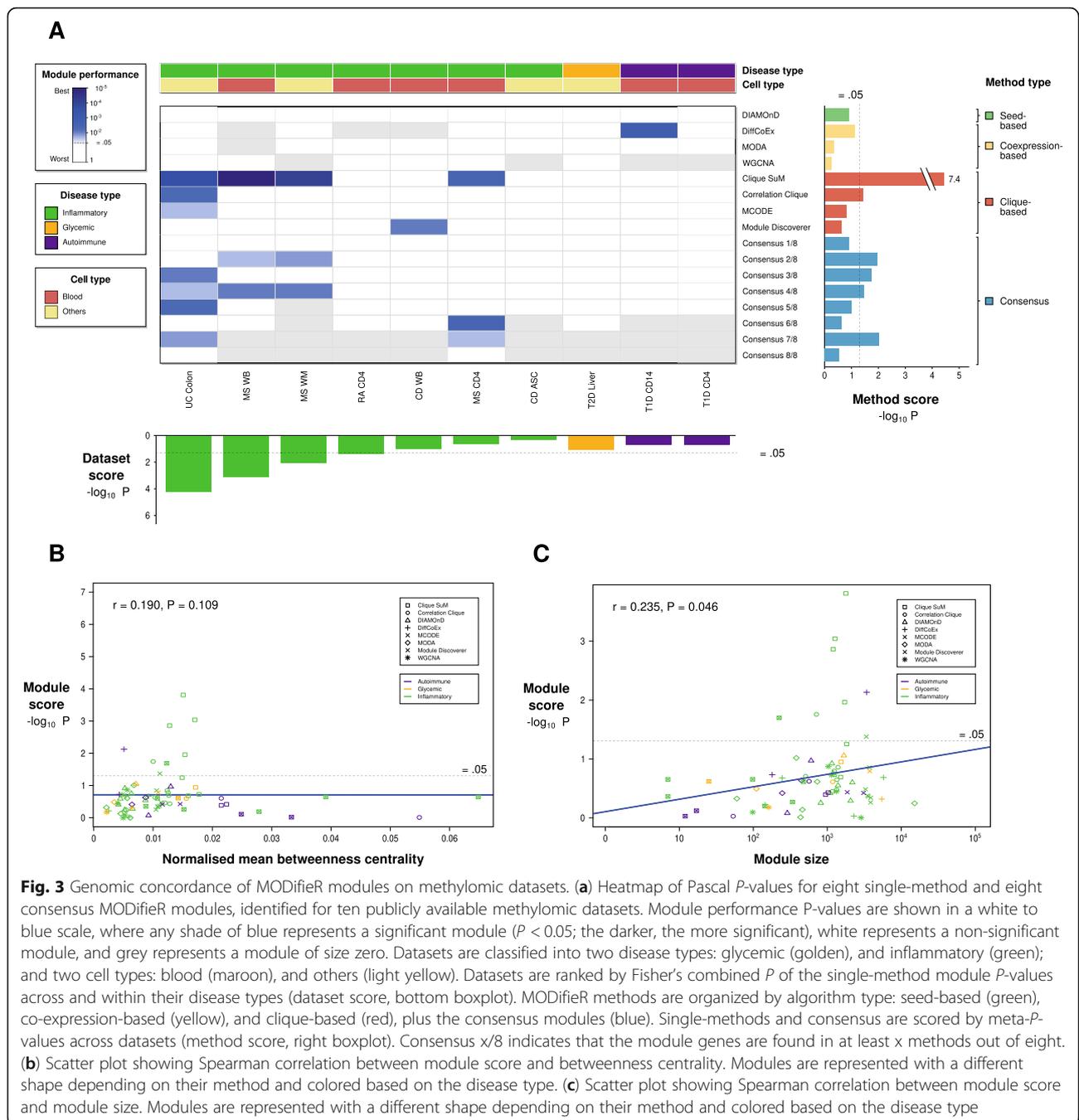
## A benchmark comparing 72 methylation-based disease modules from six different diseases using GWAS

Following the same logic of the transcriptomic benchmark, we performed a similar benchmark study for methylation modules. We collected ten datasets from three different disease categories, including six complex diseases, and ran the eight MODifieR methods on them (Fig. 1a; complete results are found in Additional file 1: Table S6, S7). In addition, we constructed consensus modules for each of the datasets. Modules were then tested for GWAS enrichment using Pascal. Inspecting the overall performance, we found nine single-method modules with a significant GWAS enrichment (9/72, 11.8%). Though this might be due to disease and cell type heterogeneity, the enrichment is more than expected by chance ($P = 9.6 \times 10^{-3}$). Interestingly, inflammatory diseases such as MS and UC showed a more significant GWAS enrichment. Since the evaluation of module performance by GWAS enrichment may be

biased due to differences in module sizes and interactome centrality, we again assessed the correlation between these values. We found a significant correlation between GWAS enrichment and module size (Fig. 3c, rho = 0.235, $P = 0.046$) and a non-significant correlation between GWAS enrichment and interactome centrality (Fig. 3b, rho = 0.190, $P = 0.109$). We found that 12.5% of the disease-method combinations yielded significant GWAS enrichment, which is more than expected from an independent random selection of modules (Fisher's exact test $P = 0.031$, $n = 6$). The highly enriched disease modules belong to MS, UC, and CD. Two out of the six diseases showed significant GWAS enrichment by using the Clique SuM modules ($P = 0.032$). In summary, the Clique SuM method resulted in a more significant GWAS enrichment for most diseases also for the methylomic benchmark.

## Multi-omic approach revealed a module enriched for MS-associated genes

Considering genomic concordance as the guidance principle for the modules that show enrichment for GWAS SNPs, differentially methylated genes and differentially expressed genes, we further wanted to evaluate multiple datasets of one specific complex disease, MS. We compiled 11 MS transcriptomic and nine methylation (see Additional file 1: Table S2) datasets from GEO which satisfy the pre-defined dataset criteria (see Methods). For each dataset we implemented the pipeline for module identification and scoring shown in Fig. 1b. We evaluated each module using MS SNP enrichment analysis and selected the most enriched modules per omic from this metric (complete results are found in Additional file 1: Table S8, S9). This analysis again showed that Clique SuM yielded the far highest average enrichment score (meta-$P = 3.2 \times 10^{-12}$) and was significantly enriched ($P < 0.05$) in 9/11 transcriptomic datasets (Fig. 4a) and 4/9 of the methylation datasets (Fig. 4b). From the significant modules generated by Clique SuM, we chose the top four modules from each of the gene transcription and methylation sets, and prioritized genes detected in modules from multiple datasets in each omic (see Additional file 1: Table S10). This analysis showed that the strongest MS SNP enrichment was found for genes in at least three out of four transcriptomic modules ($n = 1552$; $P = 6.0 \times 10^{-7}$) and two out of four methylomic modules ($n = 324$, $P = 1.5 \times 10^{-6}$). Next, we used the same principle to combine these two single-omic consensus modules and found that the intersection between them resulted in a module ($n = 220$ genes, Fig. 4) enriched for MS-associated genes (75/220, $P < 2.2 \times 10^{-16}$, OR = 7.8) and with the highest GWAS enrichment ($P = 8.8 \times 10^{-9}$), which we hereafter is referred to as the multi-omic MS module. To test if such

**Fig. 3** Genomic concordance of MODifieR modules on methylomic datasets. (**a**) Heatmap of Pascal *P*-values for eight single-method and eight consensus MODifieR modules, identified for ten publicly available methylomic datasets. Module performance P-values are shown in a white to blue scale, where any shade of blue represents a significant module (*P* < 0.05; the darker, the more significant), white represents a non-significant module, and grey represents a module of size zero. Datasets are classified into two disease types: glycemic (golden), and inflammatory (green); and two cell types: blood (maroon), and others (light yellow). Datasets are ranked by Fisher's combined *P* of the single-method module *P*-values across and within their disease types (dataset score, bottom boxplot). MODifieR methods are organized by algorithm type: seed-based (green), co-expression-based (yellow), and clique-based (red), plus the consensus modules (blue). Single-methods and consensus are scored by meta-*P*-values across datasets (method score, right boxplot). Consensus x/8 indicates that the module genes are found in at least x methods out of eight. (**b**) Scatter plot showing Spearman correlation between module score and betweenness centrality. Modules are represented with a different shape depending on their method and colored based on the disease type. (**c**) Scatter plot showing Spearman correlation between module score and module size. Modules are represented with a different shape depending on their method and colored based on the disease type
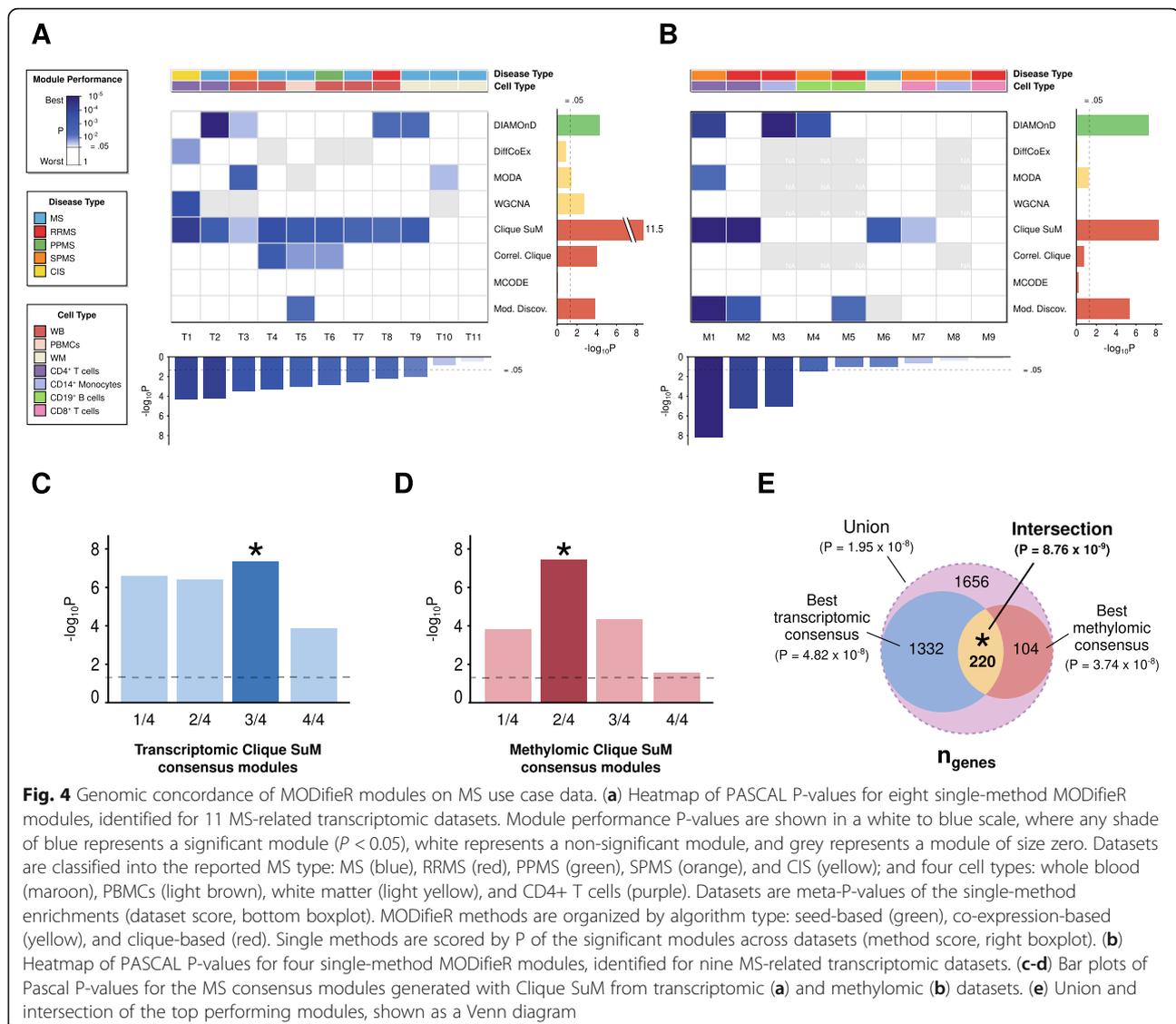
high concordance was expected even if wouldn't have used modules, we computed the overlap statistics between the differentially expressed and methylated genes of each of the MS datasets. This led to 190 pairs of studies (Additional file 1: Table S11), where the overlap between any two expression studies (*n* = 55) had average odds ratio (OR) of 1.58, while any two methylation studies (*n* = 36) had OR 2.81 and across overlap (*n* = 99) was tremendously low, merely OR = 1.03. Thus, we conclude that the overlapping multi-modules is not trivially observed in studies of genes lists, rather a result of our careful integration of omics.

## The multi-omic MS module was enriched in genes associated with major MS pathways

As we used GWAS enrichment as a selection criterion, the high GWAS enrichment of the final module was partly expected, which led us to analyze its biological functions and their potential epigenetic associations to MS. First, pathway enrichment analysis showed that the
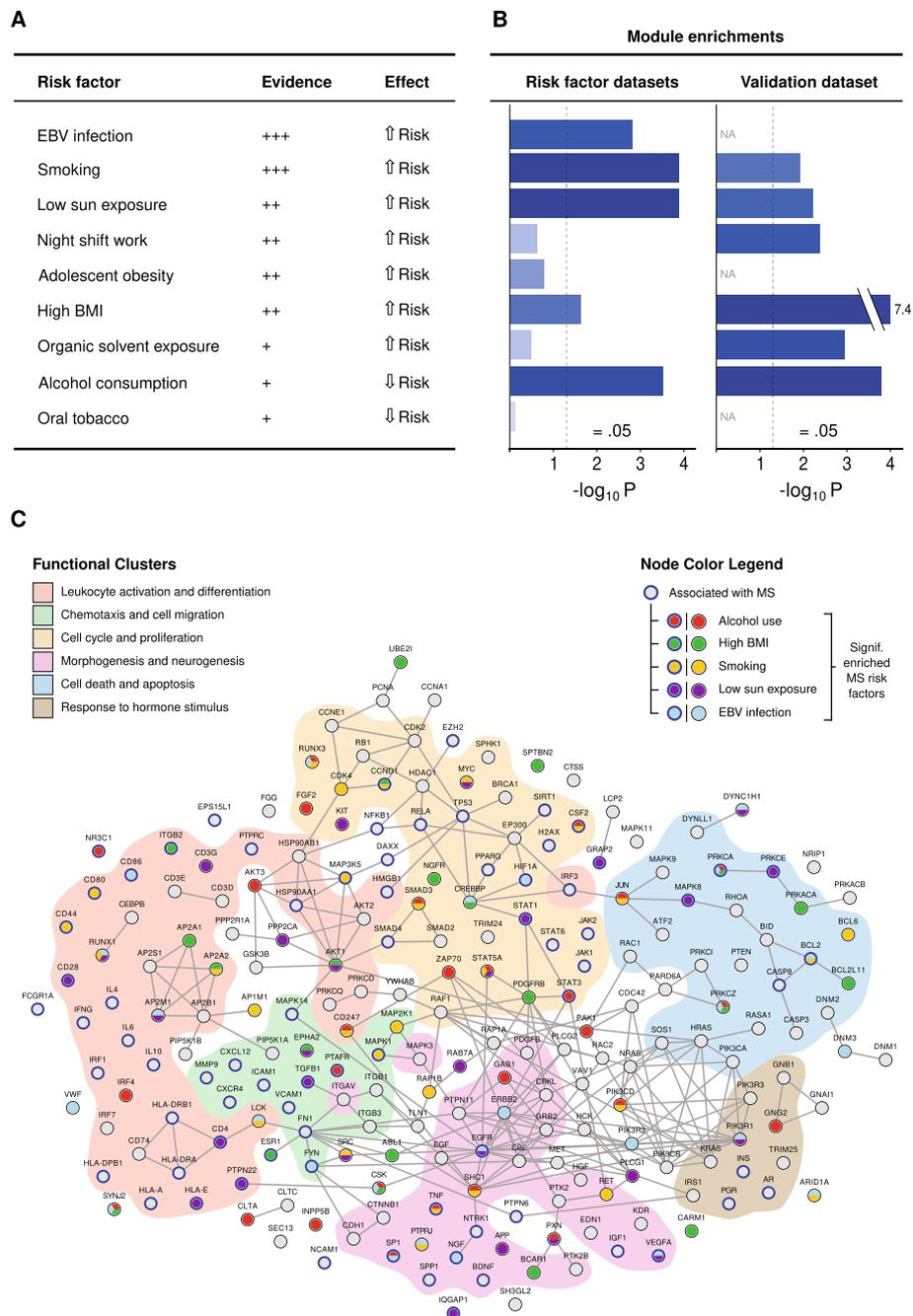
Badam *et al. BMC Genomics*      (2021) 22:631

Page 7 of 13

**Fig. 4** Genomic concordance of MODifieR modules on MS use case data. (**a**) Heatmap of PASCAL P-values for eight single-method MODifieR modules, identified for 11 MS-related transcriptomic datasets. Module performance P-values are shown in a white to blue scale, where any shade of blue represents a significant module ($P < 0.05$), white represents a non-significant module, and grey represents a module of size zero. Datasets are classified into the reported MS type: MS (blue), RRMS (red), PPMS (green), SPMS (orange), and CIS (yellow); and four cell types: whole blood (maroon), PBMCs (light brown), white matter (light yellow), and CD4+ T cells (purple). Datasets are meta-P-values of the single-method enrichments (dataset score, bottom boxplot). MODifieR methods are organized by algorithm type: seed-based (green), co-expression-based (yellow), and clique-based (red). Single methods are scored by P of the significant modules across datasets (method score, right boxplot). (**b**) Heatmap of PASCAL P-values for four single-method MODifieR modules, identified for nine MS-related transcriptomic datasets. (**c-d**) Bar plots of Pascal P-values for the MS consensus modules generated with Clique SuM from transcriptomic (**a**) and methylomic (**b**) datasets. (**e**) Union and intersection of the top performing modules, shown as a Venn diagram

multi-omic module genes are significantly involved in several inter-linked immune-related pathways, most of which have been previously associated to MS, including the T cell receptor [14] (adjusted $P = 3.6 \times 10^{-47}$), PI3K/Akt [15] ($P = 4.6 \times 10^{-35}$), ErbB [16] ($P = 7.7 \times 10^{-32}$), Fc epsilon RI [17] ($P = 8.3 \times 10^{-30}$), chemokine [18, 19] ($P = 2.6 \times 10^{-28}$), MAPK [20, 21] ($P = 2.0 \times 10^{-25}$), and B cell receptor [21] ($P = 3.9 \times 10^{-19}$) signaling pathways; Th17 ($P = 9.6 \times 10^{-29}$), and Th1 and Th2 ($P = 6.9 \times 10^{-19}$) cell differentiation [22]; natural killer cell mediated cytotoxicity ($P = 1.6 \times 10^{-27}$); and leukocyte trans-endothelial migration ($P = 3.9 \times 10^{-20}$), which indeed supports their relevance in MS. Interestingly, the module was also highly enriched in morphogenetic and neurogenetic signaling pathways, such as the neurotrophin (adjusted $P = 1.3 \times 10^{-36}$), Ras ($P = 1.4 \times 10^{-36}$), Rap1 ($P = 2.2 \times 10^{-35}$), vascular endothelial growth factor (VEGF,

$P = 1.7 \times 10^{-27}$), FoxO ($P = 3.6 \times 10^{-27}$), and mTOR ($P = 4.1 \times 10^{-14}$) signaling pathways; and in growth hormone synthesis, secretion and action ($P = 6.6 \times 10^{-31}$).

## The multi-omic MS module was enriched in genes associated with five known environmental MS risk factors validated in an independent cohort

Second, from a literature study [23, 24] we found nine environmental MS risk factors of varying evidence for which we could identify methylation studies in healthy controls. For each of these risk factors we derived the top 1000 differentially methylated genes (DMGs) and tested their enrichment with the module. Intriguingly, the module was significantly enriched for genes associated with five risk factors (Fig. 5b), which included the top associated risk factors, i.e., Epstein-Barr virus (EBV) infection (Fisher's exact test $P = 1.5 \times 10^{-3}$, OR = 2.1),

**Fig. 5** Risk factor enrichment and network visualization of the MS multi-omic module. (**a**) Evidence levels and effect on MS of the risk factor. (**b**) Enrichment overlap of multi-omic MS module genes in the top 1000 DMGs in risk factor datasets and independent risk factor methylation dataset (see Methods) shown as Fisher's exact test P-values (threshold α = 0.05). (**c**) Visualization of the module. Nodes (module genes) are arranged in functional clusters according to their overrepresented GO terms. Genes with a known association to MS are marked with a blue circle. Node colors display the associations to an MS risk factor for which the module is significantly enriched (red, alcohol use; green, high BMI; yellow, smoking; purple, low sun exposure; light blue, EBV infection; grey, no association). Edges were extracted from the STRING dB v11 human PPI network of experimentally validated interactions (confidence score > 700)

smoking ($P = 1.2 \times 10^{-4}$, OR = 2.3), low sun exposure ($P = 1.2 \times 10^{-4}$, OR = 2.3), high BMI ($P = 0.023$, OR = 1.7), and alcohol consumption ($P = 2.9 \times 10^{-4}$, OR = 2.2). Then, we asked whether these putative gene-risk factor

associations could be validated using an independent omic dataset with paired risk factor associations. For this purpose, we utilized methylation arrays of peripheral blood from 139 MS patients and 140 controls, which

Badam *et al. BMC Genomics*     (2021) 22:631

Page 9 of 13

have been described previously [25]. In this analysis we also considered risk factor associations for each individual including age, sex, BMI at age of 20, smoking, alcohol consumption, sun exposure, night shift work, contact with organic solvents. This enabled the analysis of DMGs for the MS and risk factor status as covariates in linear mixed effect analysis (see Additional file 1: Table S12). Indeed, the module genes were highly significantly enriched for MS ($n = 217$; permutation test $P = 1.2 \times 10^{-47}$), but also for all the tested risk factors (EBV was not included, Methods) and non-significantly associated to age and sex having 104–135 of the genes in each factor ($3.9 \times 10^{-8} < P < 0.013$; Fig. 5b). Combining these results, we found 90 of the 220 module genes to be associated with a risk factor from both studies, 25 genes were associated with two risk factors, and seven genes were associated with three risk factors (CSK, PRKCA, PRKCZ, RUNX1, RUNX3, STAT5A, and SYNJ2) (Fig. 5c). These suggest that the multi-omic module is capturing a key disease network with both genetically and epigenetically driven alterations, thereby providing the possibility to use it to identify potential novel biomarkers or therapeutic targets for MS.

Lastly, to check the robustness of our GWAS selection procedure we tested each of the utilized 88 expression- and 41 methylation-based MS modules for general risk factor enrichments by combining the enrichments from the seven MS risk factors, which resulted in alternative assessments of MS modules. Intriguingly, we found a highly significant correlation of this measure with our previous GWAS enrichment $P$-values (Spearman rho = 0.44, $P < 7.3 \times 10^{-7}$). Inspecting the individual methods, we again found Clique SuM to score far higher than the rest of the methods, with 45% significant modules, whereas each of the other methods had less than 25% significant modules (see Additional file 1: Table S13, S14). In summary, these results confirm the robustness and general applicability of our findings.

## Discussion
The analysis of case-control data in the context of networks has gained increased interest to detect consistent robust gene signatures of individual diseases. The application of disease modules might vary for different researchers, but here we systematically aimed at the detection of disease genes supported by genetic association. For this purpose, our study of the transcriptome and methylome profiles of 19 diseases showed significant GWAS enrichments for several inflammatory and heart diseases, while psychiatric disorders showed no enrichments and might not be suitable for GWAS validation of modules, potentially due to differences in affected tissue types and sampling points. However, the analysis of the significant results showed that methods based on

differentially expressed cliques in the protein-protein interaction network (PPI) achieved the strongest enrichments (highest scoring for Clique SuM), while those based primarily on correlations, like WGCNA, showed weaker enrichments. A potential reason for this could be that GWAS has shown to be mostly associated to the central genes of the PPI network, but our analysis demonstrated that the correlation between GWAS enrichment and centrality was non-significant in either of the omic benchmarks. We also tested whether there was an improvement using consensus approaches that counted the frequency of the result of multiple methods but found this not to increase performance. Moreover, we tested the same strategy on a set of inflammatory, glycemic, and autoimmune methylation datasets and found similar results. We would like to emphasize that, rather than scoring a single best working method, our study provides a pipeline for evaluating modules using independent high-throughput enrichments.

The work on transcription and methylation datasets suggested that MS is a disease highly enriched for GWAS, and we therefore tested if increased enrichments could be derived by their integration. We found 20 publicly available datasets and ran the assessment for both omics independently, which again showed Clique SuM to score highest. We then tested if improved results could be obtained using modules from multiple datasets of these two omics using consensus modules from Clique SuM. This resulted in a module of 220 genes highly enriched for GWAS ($P = 8.8 \times 10^{-9}$). The multiomic module was also enriched in immune-associated pathways, such as T cell and B cell receptor signaling, Th1/Th2 differentiation, or leukocyte transendothelial migration. These results conform with the current hypothesis that MS is mediated by an autoreactive response of CD4+ T cells against myelin surrounding neuronal axons, preceded by their migration across the blood-brain barrier (BBB) [26]. This auto proliferation of brain-targeting Th1 cells has been shown to be driven by memory B cells, in a process mediated by HLA-DR15 [27]. In addition, another enriched pathway was VEGF signaling. MS patients present high serum VEGF levels, which is related to pro-inflammatory functions and can alter the permeability of the BBB [28]. As GWAS was used for method prioritization, we asked if modules instead could be validated using epigenetics and lifestyle risk factor genes that we identified to associate with MS. With this aim, we compiled a set of publicly available data from methylomic studies of these risk factors in healthy individuals. This analysis demonstrated that five out of eight risk factors were enriched in our module. To validate the use of an environmental assessment using published risk factor associations, we found an independent methylome study of MS comprising

environmental data for MS patients and healthy individuals. This analysis showed a remarkable enrichment of 217 of the 220 module genes in MS-associated DMGs ($P = 1.2 \times 10^{-47}$), and a majority to be associated with the tested risk factors.

In contrast to previously known community challenges, in our study we not only used the topological properties of the network, but we also combined the methods to use an omic-based input to uncover the disease modules that might be dysregulated at each omic level, contributing to the diverse causative mechanisms behind complex diseases. Although using the PPI network as background may lead to certain knowledge bias, this kind of benchmark allowed us to look at the relevant risk factors. In our assessment of the disease modules, Clique SuM showed more robust performance on average, compared to the other methods and to the community-based consensus predictions, independently of omics using GWAS or risk factor enrichments as scorings. This robustness in result could stem from the underlying idea that fully connected sub-graphs on average represent small functional building blocks. However, for specific research problems we recommend that researchers test multiple methods and assess the outcome using independent data sources as we presented in this manuscript.

## Conclusions

In summary, our study provides a practical integrative workflow that enables system-level analysis of heterogeneous diseases, in terms of multi-omic disease modules, as well as the validation of these by using both disease-specific GWAS and risk factors enrichment. We believe that this analysis validates our integrated datasets and suggest a pipeline that could readily be tested in at least other inflammatory and cardiovascular diseases. Lastly, our study did not aim to optimize parameters for individual disease module identification methods, instead used default values from the MODifieR R package implementation of the methods, when possible [12]. However, this might be a key step in the analysis of specific diseases. Thus, the code and processed datasets are available at GitLab (https://gitlab.com/Gustafsson-lab/modifier-benchmark). In future work, this approach can be expanded to include diverse and context-specific networks to determine whether multi-omic modules are able to capture various other levels of granularity.

## Methods
### Benchmark data

A total of 47 publicly available datasets for the transcriptomic benchmark and ten publicly available datasets for the methylomic benchmark were used. To avoid bias due to subtypes of diseases and drug treatments, we searched for datasets that have only patient and control samples, and that are available for download from the Gene Expression Omnibus (GEO) database. We categorized the datasets into seven distinct disease types based on the disease-trait type associations used in Choobdar et al. [6], i.e., autoimmune, cardiovascular, glycemic, inflammatory, neurodegenerative, and psychiatric and social disorders. A total of 19 complex diseases were used in the transcriptomic benchmark analysis, while six complex diseases were used in the methylation benchmark analysis. The methylation benchmark diseases belong to inflammatory, autoimmune, and glycemic disease types (Additional file 1: Table S1).

### MS use case data

A total of 14 publicly available and one non-publicly available transcriptomic and methylomic MS-related datasets were used in the MS multi-omic integration use case. In general, every dataset in the MODifieR benchmark was also used in the MS use case, with exceptions according to certain criteria (Additional File 1: Table S2). The inclusion of transcriptomic MS datasets followed the criteria: 1) The largest dataset by sample number, per tissue, is shown in the MODifieR benchmark; 2) Replication cohorts are not included in the MS use case. Criteria for inclusion of methylomic MS datasets were the following: 1) The largest dataset by sample number, per tissue or cell type, was included in the MODifieR benchmark; 2) A single dataset for every cell-specific tissue was included in the benchmark; 3) Methylation studies that reported using whole blood as sample tissue were excluded from the MS use case, due to the high heterogeneity of this type of data.

For the additional independent validation, we utilized the methylation microarray analysis of 279 blood samples analyzing from Kular et al. [25]. For each of these MS patients ($n_{MS} = 139$) and healthy controls ($n_{HC} = 140$), we also collected their lifestyle-associated risk factors from questionnaires that were part of the Epidemiological Investigation of Multiple Sclerosis (EIMS) study. Those factors were smoking status, prior EBV infection, sunbathing, nightshift work, alcohol consumption, as well as phenotypic features (age, sex, BMI at age of 20).

### Pre-processing and quality control of risk factor methylation data

DNA methylation datasets were downloaded from GEO as raw IDAT files, when available, or matrices of beta values. Pre-processing of the data was performed using the Chip Analysis Methylation Pipeline (ChAMP) R package [29], version 2.16.2. Default parameters were used for probe and sample filtering. Probes with a detection $P$-value above 0.01, probes with a fraction of failed (bead count less than 3) samples over 0.05, non-CpG

probes, SNP-related probes, multi-hit probes, and probes located on chromosomes X and Y were removed. Samples with a proportion of failed (NA) probe *P*-values over 0.1 were also removed from the analysis. Post-filtering imputation of NA values was conducted on the beta matrices, with default parameters ("combine" method, k = 5, probe cutoff = 0.2, sample cutoff = 0.1). Filtered imputed matrices were normalized applying the Beta-Mixture Quantile dilation (BMIQ) normalization method [30], including correction of Type-I and Type-II probe effects. Data quality was assessed by producing multi-dimensional scaling (MDS) plots of the top 1000 most variable positions per sample, density plots for the distribution of beta values, and hierarchical clustering of samples, before and after normalization. Singular value decomposition (SVD) was used to detect the most significant components of variation in the data. Unwanted sources of variation in the normalized data were corrected using ComBat batch effect correction [31].

### Module identification

The MODifieR [12] R package offers nine different methods for producing disease modules for which we included all but Clique SuM exact as it is highly similar to Clique SuM. The included methods will produce modules based on the provided omic input and background network and do not include prioritization of pathway association. MODifieR methods used for module identification through this study are listed in the Additional file: Table S3. For the methods that require a network, we used the high confidence interactions (cut-off 700) from the human PPI network from STRING [5] database version 11, which included 816,352 interactions between 16,770 proteins (used in the MS use case). In the two benchmark sections when we compared expression and methylation methods, we further limited the network to 631,782 very high confidence interactions (cutoff > 900) between 12,123 proteins as the running times of some methods (e.g., Module Discoverer) took too long time to compute. Additionally, the benchmark was performed using also the 120,000 experimentally verified STRING interactions between 8000 proteins and the 27,719 Reactome [32] curated interactions of 5147 proteins.

The processed matrix for each dataset and their respective phenotypic information were downloaded from GEO. The input object is prepared using the *create_input_microarray* function from the MODifieR package which is then used for creating the modules. The input function applies linear model using limma for comparison of patients vs. controls to get the differentially methylated or expressed genes. A dynamic cutoff of 5% in the differentially methylated or expressed genes was applied for input seed genes, for the methods that require them.

### Differential methylation analysis of risk factor data

Differentially methylated probes (DMPs) were found by fitting a linear model to the data using the limma R package [33], version 3.42.2 implemented in the ChAMP function *champ. DMP*. *P*-values were adjusted for multiple testing using Benjamini-Hochberg False Discovery Rate (FDR) correction. Differentially methylated genes (DMGs) were obtained and annotated using the org. Hs.eg.db R package, version 3.10.0. DMG lists were cross-checked against the STRING database version 11 PPI network used for module identification in the MS multi-omic approach (high confidence interactions, combined score > 700). DMGs that were not present in the PPI network were removed. In case of the additional MS validation dataset, a linear mixed effect model with risk factors (age, sex, BMI at age of 20, smoking, alcohol consumption, sun exposure, night shift work, contact with organic solvents) as categorical covariates was implemented to find the differentially methylated genes after the preprocessing step, as described in the preprocessing section of the methods. Since all the patients were EBV positive, we did not include it in the linear mixed effect model.

### Validation of modules

The final modules produced from each single algorithm and the consensus were evaluated using Pascal [13] (Pathway scoring algorithm). Pascal implements a fast and rigorous gene scoring and pathway enrichment pipeline that can be run on a local machine. The SNP values are converted to gene scores by computing pairwise SNP-by-SNP correlations and obtaining Z-scores from their distribution, where SNPs are mapped to the closest gene. These obtained gene scores are fused with the pathway enrichment analysis to recompute a chi-square *P*-value for the given set of module genes. Thus, the obtained chi-square *P*-value serves as the significance of the module in its enrichment of the disease-associated pathway gene loci. A combined *P*-value was computed for each of the methods using Fisher's method [34], diseases, and datasets for ranking the performance of the modules in each criterion (see Additional file 1: Table S1 for details).

### Integration of MS single-omic modules

Clique SuM was ranked as the best performing method on average for both transcriptomic and methylomic data, according to the MS GWAS enrichment of the modules calculated by Pascal. Therefore, significant Clique SuM modules ($P < 0.05$) were selected for further analysis (nine transcriptomic and four methylomic modules). Consensus modules were generated across each omic by applying a module count-based method, where the criteria for gene inclusion in the consensus is its presence

in a certain number of single-method modules. To balance the weight of each omic in the multi-omic integration, the top four significant modules per omic were used to create each consensus (Fig. 4a and b). Single-omic Clique SuM consensus were ranked again by GWAS enrichment, and the best performing consensus per omic was selected for integration into the multi-omic module.

## Enrichment analyses of the MS multi-omic module

Disease enrichment analysis of the multi-omic module was performed by Fisher's exact test, with a significance threshold of $P < 0.05$. MS-associated genes were obtained from the gene-disease association summary provided by DisGeNET database 6.0 [35]. All genes with a known association to the disease "multiple sclerosis" (Unified Medical Language System unique identifier C0026769) were considered MS-associated genes ($n$ = 1105). Pathway enrichment analysis was carried out using the function *enrichKEGG* from the clusterProfiler R package [36], version 3.14.3. $P$-values were adjusted for multiple testing using Benjamini-Hochberg FDR correction, with a significance threshold of adj. $P < 0.05$. Enrichment of the multi-omic module in MS risk-factor-associated genes was performed by Fisher's exact test, with a significance threshold of $P < 0.05$. To provide a uniform comparison of MS risk factor-associated genes across datasets, the module was tested for enrichment in the top 1000 DMGs (with at least $P < 0.05$) obtained from the differential methylation analysis with ChAMP for each risk factor dataset.

## Representation of the MS multi-omic module

Experimentally validated interactions for the multi-omic module genes were obtained from STRING database version 11 (experimental score > 700) and imported into Cytoscape [37] version 3.7.2. To determine representative functional clusters of module genes, overrepresented Gene Ontology (GO) Biological Process (BP) terms in the module were found using BiNGO [38] version 3.0.4, with Benjamini-Hochberg FDR for multiple testing correction, and a significance threshold of adj. $P < 0.05$. Then, enriched GO terms with adj. $P < 10^{-10}$ were summarized using REVIGO [39] server tool (medium allowed similarity = 0.7) and categories of interest were selected by uniqueness (> = 80%), dispensability (> = 50%), and frequency (<=10%) criteria. Further manual assessment was performed to group similar terms with an adequate number of genes in the network.

## Abbreviations
GWAS: Genome-wide association studies; MS: Multiple sclerosis; SNP: Singe Nucleotide Polymorphism; FDR: False discovery rate; DMG: Differentially methylated gene; DMP: Differentially methylated probe; EBV: Epstein-Barr virus; CpG: Region of DNA where a cytosine nucleotide is followed by a

guanine nucleotide in the linear sequence of bases along its 5′ → 3′ direction

## Supplementary Information

## Availability of data and materials
The data used for transcriptomic benchmark and methylation benchmark are downloaded from GEO. The disease specific GWAS files are downloaded from the latest Pascal version. The processed Data for analysis is available at https://gitlab.com/Gustafsson-lab/modifier-benchmark.The risk factor (EIMS) data will be made available on request. The R-package MODifieR is available on the GitLab: https://gitlab.com/Gustafsson-lab/MODifieR; the code used for benchmark analysis and risk factor analysis is available on GitLab: https://gitlab.com/Gustafsson-lab/modifier-benchmark; the latest Pascal version: https://www2.unil.ch/cbg/index.php?title=Pascal.

## Declarations

### Ethics approval and consent to participate
Not applicable.

### Consent for publication
Not applicable.

### Competing interests
The authors declare no competing interests.

### Author details
[1]School of Bioscience, Systems Biology Research Center, University of Skövde, Skövde, Sweden. [2]Bioinformatics, Department of Physics, Chemistry and Biology, Linköping university, Linköping, Sweden. [3]Department of Clinical Neuroscience, Karolinska Institutet, Center for Molecular Medicine, Karolinska University Hospital, SE-171 76 Stockholm, Sweden. [4]Institute of Environmental Medicine, Karolinska Institutet, Center for Molecular Medicine, Karolinska University Hospital, SE-171 76 Stockholm, Sweden.

### References

1. Naylor S, Chen JY. NIH public access. Natl Institutes Heal. 2011;7:275–89.
2. Santiago JA, Bottero V, Potashkin JA. Dissecting the molecular mechanisms of neurodegenerative diseases through network biology. Front Aging Neurosci [Internet]. 2017;9:1–13. Available from:. https://doi.org/10.3389/fnagi.2017.00166/full.
3. Barabási AL, Gulbahce N, Loscalzo J. Network medicine: a network-based approach to human disease. Nat rev genet [internet]. Nat Publ Group. 2011;12(1):56–68. Available from:. https://doi.org/10.1038/nrg2918.
4. Gustafsson M, Nestor CE, Zhang H, Barabási A-L, Baranzini S, Brunak S, et al. Modules, networks and systems medicine for understanding disease and aiding diagnosis. Genome Med [Internet]. 2014;6:82. Available from:. https://doi.org/10.1186/s13073-014-0082-6.
5. Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-cepas J, et al. STRING v11 : protein – protein association networks with increased coverage , supporting functional discovery in genome-wide experimental datasets. Nucleic Acids Res. 2019;47:607–13 Oxford University Press.
6. Choobdar S, Ahsen ME, Crawford J, Tomasoni M, Fang T, Lamparter D, et al. Assessment of network module identification across complex diseases. Nat Methods. 2019;16(9):843–52. https://doi.org/10.1038/s41592-019-0509-5.
7. Schadt EE. Molecular networks as sensors and drivers of common human diseases. Nature [Internet]. 2009;461(7261):218–23. Available from:. https://doi.org/10.1038/nature08454.
8. Ghiassian SD, Menche J, Barabási AL. A DIseAse MOdule Detection (DIAMOnD) Algorithm Derived from a Systematic Analysis of Connectivity Patterns of Disease Proteins in the Human Interactome. Rzhetsky A, editor. PLoS Comput Biol [Internet]. 2015;11:e1004120. Available from:. https://doi.org/10.1371/journal.pcbi.1004120.
9. Hellberg S, Eklund D, Gawel DR, Köpsén M, Zhang H, Nestor CE, et al. Dynamic response genes in CD4+ T cells reveal a network of interactive proteins that classifies disease activity in multiple sclerosis. Cell Rep. 2016;16(11):2928–39. https://doi.org/10.1016/j.celrep.2016.08.036.
10. Wang H, Rogers G, Benson M, Jarvelin M-R, Chavali S, Ramasamy A, et al. Highly interconnected genes in disease-specific networks are enriched for disease-associated polymorphisms. Genome Biol. 2012;13:R46.
11. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics. 2008;9(1):1–13. https://doi.org/10.1186/1471-2105-9-559.
12. de Weerd HA, Badam TVS, Martínez-Enguita D, Åkesson J, Muthas D, Gustafsson M, et al. MODifieR: an ensemble R package for inference of disease modules from transcriptomics networks. Bioinformatics. 2020;1–2:3918–9.
13. Lamparter D, Marbach D, Rueedi R, Kutalik Z, Bergmann S. Fast and rigorous computation of gene and Pathway scores from SNP-based summary statistics. PLoS Comput Biol. 2016;12:1–20.
14. Carbone F, De Rosa V, Carrieri PB, Montella S, Bruzzese D, Porcellini A, et al. Regulatory T cell proliferative potential is impaired in human autoimmune disease. Nat Med. 2014;20(1):69–74. https://doi.org/10.1038/nm.3411.
15. Mammana S, Bramanti P, Mazzon E, Cavalli E, Basile MS, Fagone P, et al. Preclinical evaluation of the PI3K/Akt/mTOR pathway in animal models of multiple sclerosis. Oncotarget. 2018;9(9):8263–77. https://doi.org/10.18632/oncotarget.23862.
16. Holley JE, Gveric D, Newcombe J, Cuzner ML, Gutowski NJ. Astrocyte characterization in the multiple sclerosis glial scar. Neuropathol Appl Neurobiol. 2003;29(5):434–44. https://doi.org/10.1046/j.1365-2990.2003.00491.x.
17. Pedotti R, DeVoss JJ, Youssef S, Mitchell D, Wedemeyer J, Madanat R, et al. Multiple elements of the allergic arm of the immune response modulate autoimmune demyelination. Proc Natl Acad Sci U S A. 2003;100(4):1867–72. https://doi.org/10.1073/pnas.252777399.
18. Cui LY, Chu SF, Chen NH. The role of chemokines and chemokine receptors in multiple sclerosis. Int Immunopharmacol [internet]. 2020;83:106314. Elsevier, Available from. https://doi.org/10.1016/j.intimp.2020.106314.
19. Krumbholz M, Theil D, Cepok S, Hemmer B, Kivisäkk P, Ransohoff RM, et al. Chemokines in multiple sclerosis: CXCL12 and CXCL13 up-regulation is differentially linked to CNS immune cell recruitment. Brain. 2006;129(1):200–11. https://doi.org/10.1093/brain/awh680.
20. Krementsov DN, Thornton TM, Teuscher C, Rincon M. The emerging role of p38 mitogen-activated protein kinase in multiple sclerosis and its models. Mol Cell Biol. 2013;33(19):3728–34. https://doi.org/10.1128/MCB.00688-13.
21. Kotelnikova E, Kiani NA, Messinis D, Pertsovskaya I, Pliaka V, Bernardo-Faura M, et al. MAPK pathway and B cells overactivation in multiple sclerosis revealed by phosphoproteomics and genomic analysis. Proc Natl Acad Sci U S A. 2019;116(19):9671–6. https://doi.org/10.1073/pnas.1818347116.
22. Kunkl M, Frascolla S, Amormino C, Volpe E, Tuosto L. T helper cells: the modulators of inflammation in multiple sclerosis. Cells. 2020;9(2):482. https://doi.org/10.3390/cells9020482.
23. Waubant E, Lucas R, Mowry E, Graves J, Olsson T, Alfredsson L, et al. Environmental and genetic risk factors for MS: an integrated review. Ann Clin Transl Neurol. 2019;6(9):1905–22. https://doi.org/10.1002/acn3.50862.
24. Olsson T, Barcellos LF, Alfredsson L. Interactions between genetic, lifestyle and environmental risk factors for multiple sclerosis. Nat Rev Neurol Nature Publishing Group. 2016;13:26–36.
25. Kular L, Liu Y, Ruhrmann S, Zheleznyakova G, Marabita F, Gomez-Cabrero D, et al. DNA methylation as a mediator of HLA-DRB1 15:01 and a protective variant in multiple sclerosis. Nat Commun. 2018;9(1):2397. https://doi.org/10.1038/s41467-018-04732-5.
26. Compston A, Coles A. Multiple sclerosis. The Lancet. 2008;372(9648):1502–17. https://doi.org/10.1016/S0140-6736(08)61620-7.
27. Jelcic I, Al Nimer F, Wang J, Lentsch V, Planas R, Jelcic I, et al. Memory B Cells Activate Brain-Homing, Autoreactive CD4+ T Cells in Multiple Sclerosis. Cell. 2018;175:85–100.e23.
28. Lange C, Storkebaum E, De Almodóvar CR, Dewerchin M, Carmeliet P. Vascular endothelial growth factor: a neurovascular target in neurological diseases. Nat rev Neurol [internet]. Nat Publ Group. 2016;12(8):439–54. Available from:. https://doi.org/10.1038/nrneurol.2016.88.
29. Tian Y, Morris TJ, Webster AP, Yang Z, Beck S, Feber A, et al. Genome analysis ChAMP : updated methylation analysis pipeline for Illumina BeadChips. Bioinformatics. 2017;33:3982–4.
30. Teschendorff AE, Marabita F, Lechner M, Bartlett T, Tegner J, Gomez-cabrero D, et al. Gene expression A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data. Bioinformatics. 2013;29:189–96.
31. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. Biostatistics. 2007;8(1):118–27. https://doi.org/10.1093/biostatistics/kxj037.
32. Vastrik I, D'Eustachio P, Schmidt E, Joshi-Tope G, Gopinath G, Croft D, et al. Reactome: a knowledge base of biologic pathways and processes. Genome Biol. 2007;8(3):R39. https://doi.org/10.1186/gb-2007-8-3-r39.
33. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers differential expression analyses for RNAsequencing and microarray studies. Nucleic Acids Res. 2015;43(7):e47.
34. Mosteller F, Fisher RA. The American Statistician. 1948;2(5);30–1. https://doi.org/10.2307/2681650.
35. Piñero J, Ramírez-Anguita JM, Saüch-Pitarch J, Ronzano F, Centeno E, Sanz F, et al. The DisGeNET knowledge platform for disease genomics: 2019 update. Nucleic Acids Res. 2020;48(D1):D845–55. https://doi.org/10.1093/nar/gkz1021.
36. Yu G, Wang LG, Han Y, He QY. ClusterProfiler: an R package for comparing biological themes among gene clusters. Omi A J Integr Biol. 2012;16(5):284–7. https://doi.org/10.1089/omi.2011.0118.
37. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: A Software Environment for Integrated Models. Genome Res [Internet]. 1971;13:426 Available from: http://ci.nii.ac.jp/naid/110001910481/.
38. Maere S, Heymans K, Kuiper M. BiNGO: a Cytoscape plugin to assess overrepresentation of gene ontology categories in biological networks. Bioinformatics. 2005;21(16):3448–9.
39. Supek F, Bošnjak M, Škunca N, Šmuc T. Revigo summarizes and visualizes long lists of gene ontology terms. PLoS One. 2011;6(7):e21800. https://doi.org/10.1371/journal.pone.0021800.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.