



T9 - Network-medicine-based drug repurposing: Nextflow pipeline for drug repurposing

Markus List

[BC]² 2025

Basel, 08.09.2025





Structure

- **~ 30 min Pipeline introduction**
 - Motivation
 - Introduction to Nextflow
 - Pipeline overview and some analysis insights
 - Hands-on session overview
- **~ 15 min Coffee break**
- **~ 30 min Hands-on session**
 - Run the pipeline on Huntington's Disease data
 - Output interpretation



Motivation



Disease modules

*Disease proteins do **not** act in **isolation***

Example **Huntington's Disease**:

- Neurodegenerative disorder caused by **mutation** in the **HTT** gene



mHtt

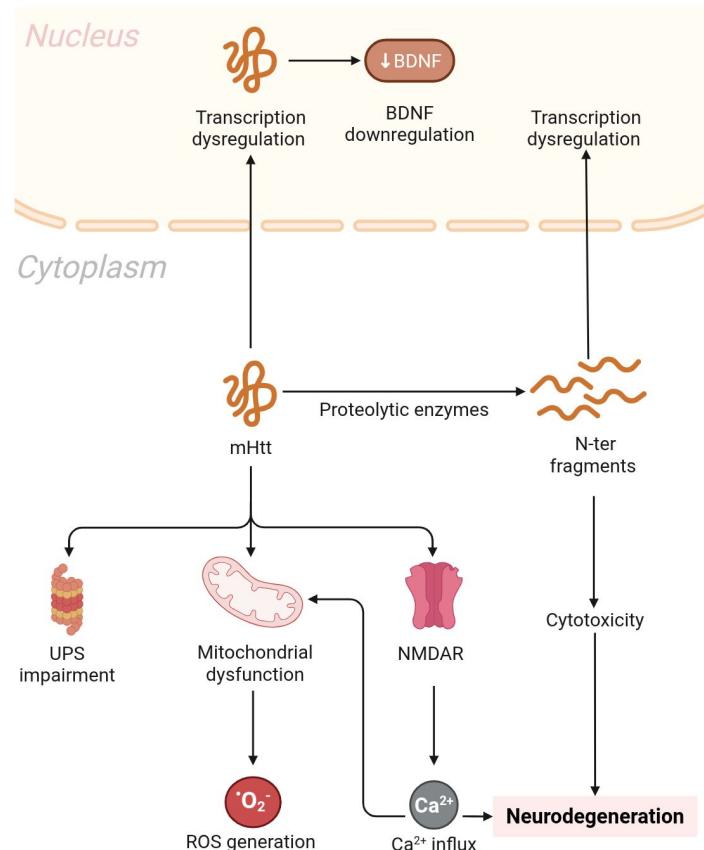


Disease modules

*Disease proteins do **not** act in isolation*

Example Huntington's Disease:

- Neurodegenerative disorder caused by **mutation** in the **HTT** gene
- **Mutant protein (mHtt) disrupts** various cellular processes



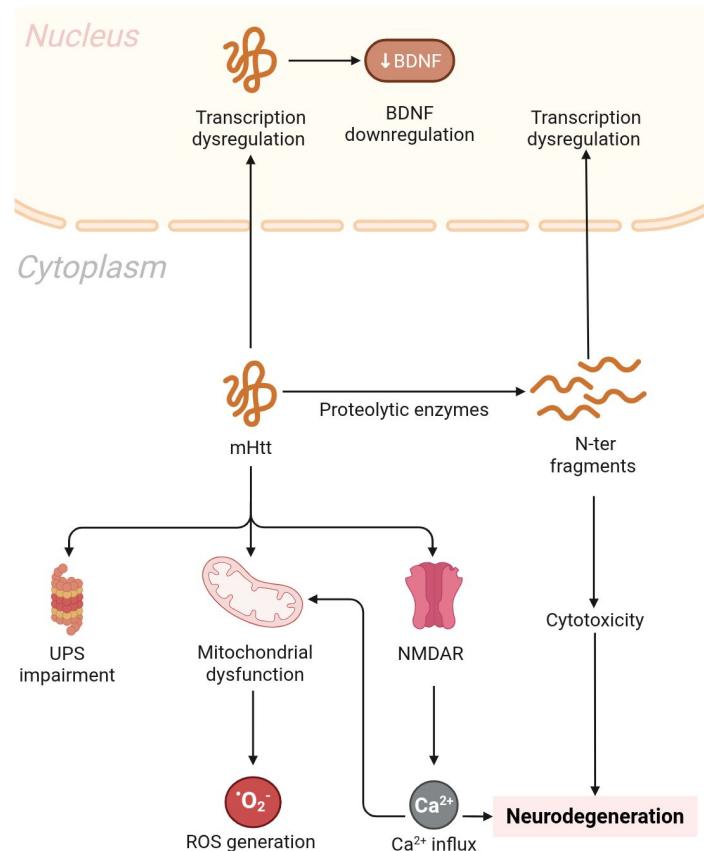


Disease modules

*Disease proteins do **not** act in isolation*

Example Huntington's Disease:

- Neurodegenerative disorder caused by **mutation** in the **HTT** gene
- **Mutant protein (mHtt) disrupts** various cellular processes
- Exact mechanism **incompletely understood**



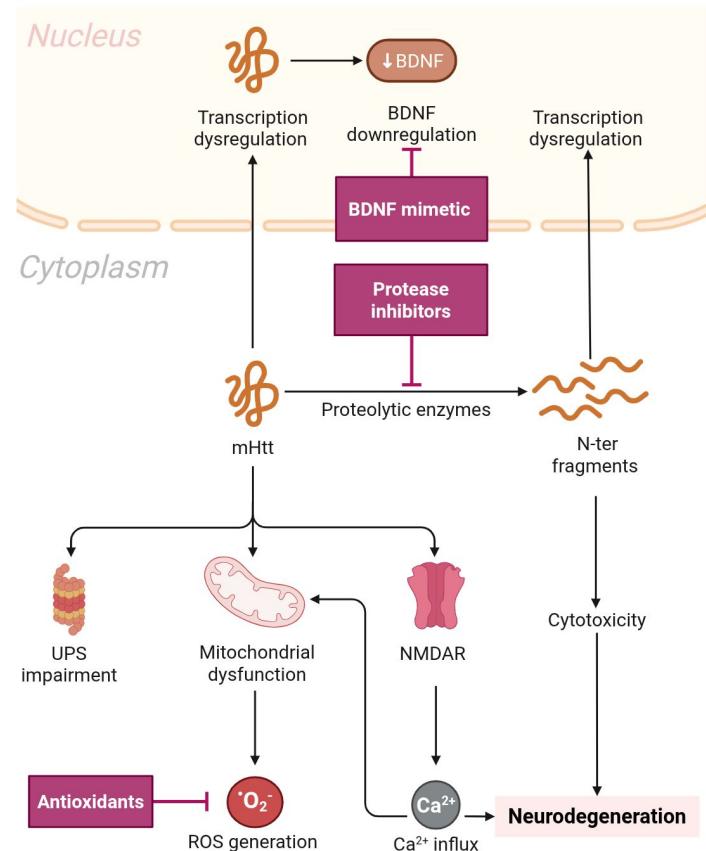


Disease modules

*Disease proteins do **not** act in isolation*

Example Huntington's Disease:

- Neurodegenerative disorder caused by **mutation** in the **HTT** gene
- **Mutant protein (mHtt) disrupts** various cellular processes
- Exact mechanism **incompletely understood**
- **Understanding the mechanism is key** to guiding treatment research





Disease module inference

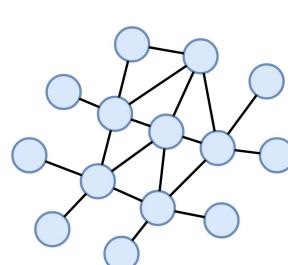
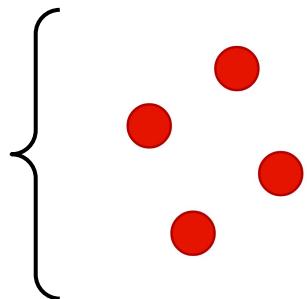


Disease module inference

Seed nodes

Disease-associated genes/proteins from

- Databases
- Literature
- Experiments



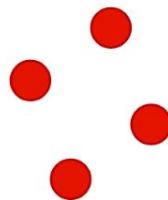
Interaction network

Usually protein-protein interactions (PPIs)

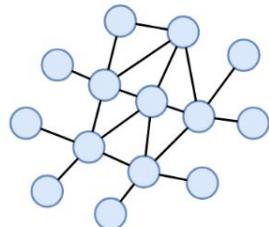


Disease module inference

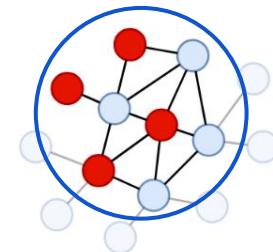
Seeds



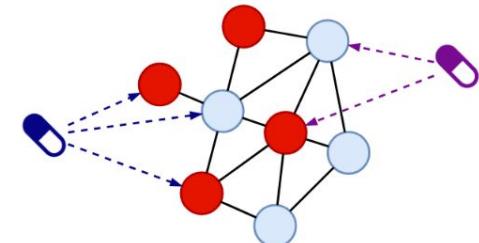
Interaction network



Disease module



Drug prioritization / repurposing



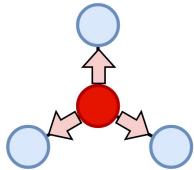


Algorithms for disease module inference



1st Neighbors:

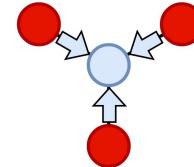
Adds all nodes that **directly interact** with a seed node



ROBUST:

Connects the seeds using **diverse** prize-collecting Steiner trees

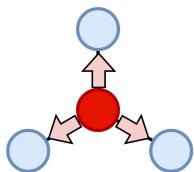
Bennett et al. *Bioinformatics* (2022)



DIAMOnD:

Iteratively expands the seeds **integrating the node** with the **most significant connectivity**

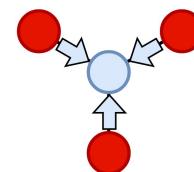
Ghiassian et al. *PLoS Computational Biology* (2015)



ROBUST (bias aware):

Works like ROBUST but **penalizes** the inclusion of **overstudied genes/proteins**

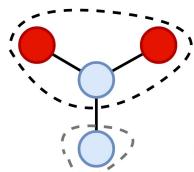
Sarkar et al. *Bioinformatics* (2023)



DOMINO:

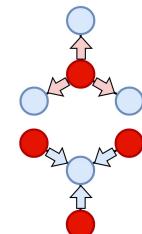
Partitions the network into **clusters** and returns those which are **enriched with seeds**

Levi et al. *Molecular Systems Biology* (2021)



Random Walk with Restart (RWR):

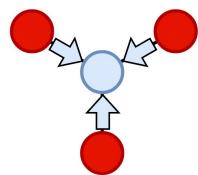
Expands the seeds using **network propagation** **until** all seeds are connected



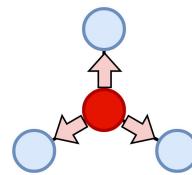


Methods use different modeling approaches

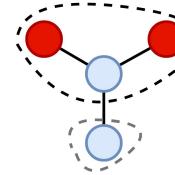
Connect



Expand



Cluster

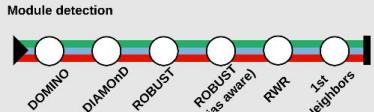


Evaluating disease modules is challenging



Module Detection

- Six module detection **algorithms** included
- Subworkflows for input and output parsing



Evaluation

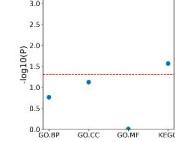
Over-Representation

g:Profiler

Functional Coherence



Empirical P-value based on Jaccard index

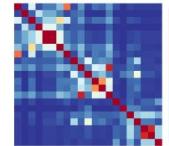


Topology Measures

General Statistics

	Nodes	Edges	Seedset	Components
enrich_seeds_Collared	203	196	10	3
enrich_seeds_Lobster	9	8	1	1
enrich_seeds_Mitochondria	10	9	1	1
enrich_seeds_Vesicle	24	23	2	2
enrich_seeds_Vesicle_transf_79	79	78	3	1
enrich_seeds_Vesicle_transf_100	100	99	3	1
enrich_seeds_Collared	22	21	1	1
enrich_seeds_Lobster	7	7	1	1
enrich_seeds_Mitochondria	34	33	1	1
enrich_seeds_Vesicle	37	36	2	1
enrich_seeds_Vesicle_transf_36	36	35	2	1

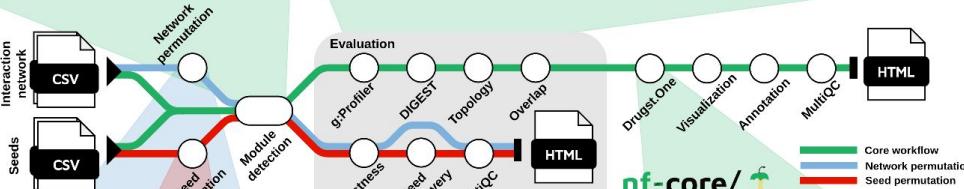
Module Overlaps



Launch

Run pipeline

```
nextflow run \
    nf-core/diseasemodulediscovery \
    -profile docker \
    --network PPI.csv \
    --seeds seed_genes.csv \
    --id_space entrez \
    --outdir results
```



nf-core/ diseasemodulediscovery

Output

Annotated Modules



- Targeting drugs
- Associated disorders
- Cellular localization

Unified Output
Annotated modules saved in BioPAX format

Pipeline Report
MultiQC report with summary statistics and figures

Network Visualizations

Interaction Network
Provide files or choose predefined

Seeds
Provide files with seed nodes

Software Deployment

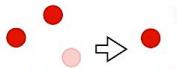


Network Permutation

Degree-preserving rewiring



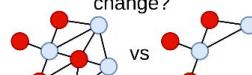
Seed Permutation
Leave-one-out approach (repeat for every seed gene)



Permutation-Based Evaluation

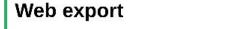
Robustness

How much do the modules change?

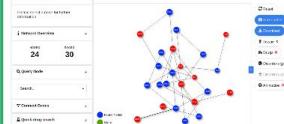


Seed Rediscovery

Can the removed seeds be recovered?



Web export



Drug prioritization

- TrustRank
- Degree Centrality
- Harmonic Centrality



Nextflow



Nextflow is a *workflow management language* used to automate scientific analysis pipelines

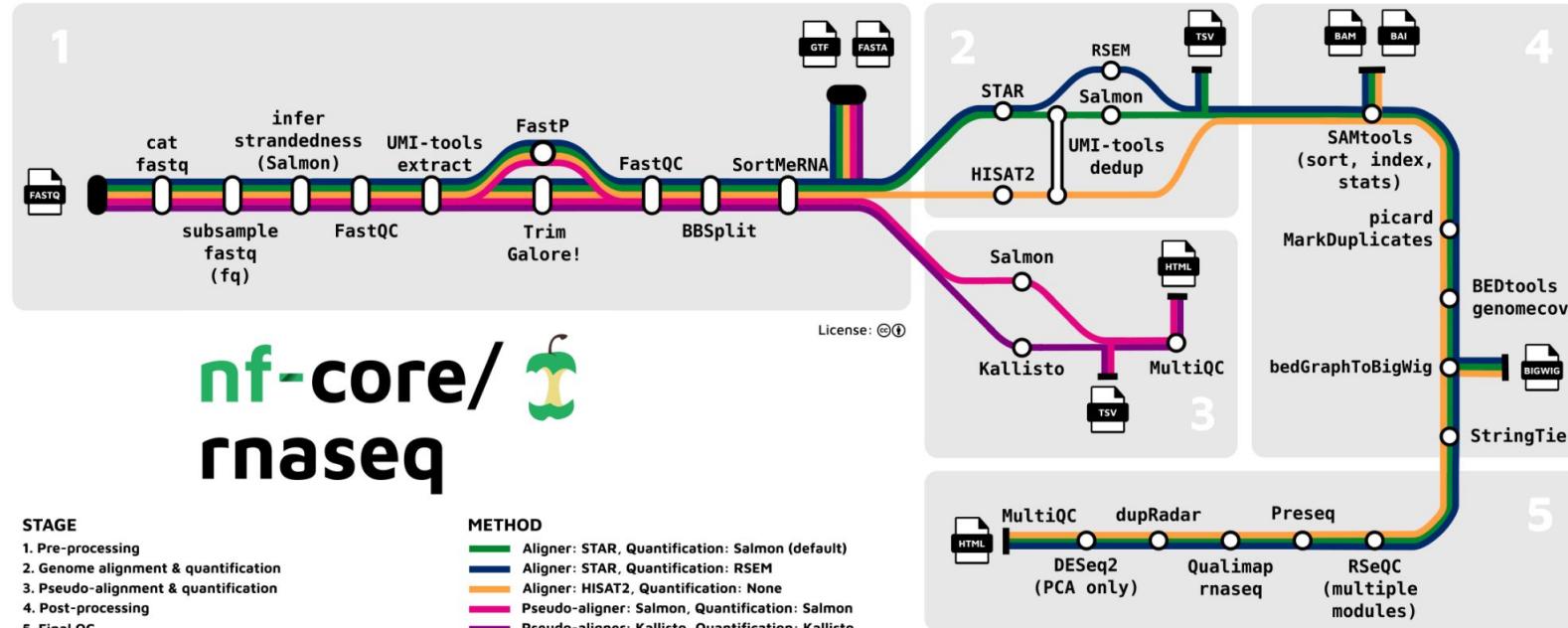


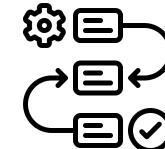
Image source: https://raw.githubusercontent.com/nf-core/rnaseq/3.14.0/docs/images/nf-core-rnaseq_metro_map_grey.png

Di Tommaso, Paolo, et al. *Nature biotechnology* 35.4 (2017): 316-319.

Nextflow is a *workflow management language* used to automate scientific analysis pipelines

Orchestration

- Passes input and output files between different processes
- Runs processes in parallel whenever possible
- Tracks intermediate results → allows resuming a cached workflow in case of failure



Reproducibility

- Nextflow automatically deploys all required software
- Various container solutions supported



Portability

- Supports execution on local machines, cloud platforms, and various workload schedulers
- Examples:   



nf-core provides a collection of **community-curated, best-practice** Nextflow **pipelines**

- Pipelines for most **standard workflows** (RNA-seq, variant calling, ChIP-seq, ...)
- **Cooperation** instead of duplication
 - Avoids redundant pipelines
 - Extend and maintaining one pipeline per analysis type
- Input, output, and usage **documentation**
- Automated continuous integration **tests** for all pipelines



Pipelines

Browse the 137 pipelines that are currently available as part of nf-core.

Search

Released 88 Under development 44 Archived 10 ⚡ Stars 28

rnaseq	1085	New release!
RNA sequencing analysis pipeline using STAR, RSEM, HISAT2 or Salmon with gene/transcript counts and extensive quality control.		
rna	rna-seq	3.20.0 released about 8 hours ago

sarek	467	
Analysis pipeline to detect germline or somatic variants (pre-processing, variant calling and annotation) from WGS / targeted sequencing		
somatic	targeted	variant-calling
whole-genome	sequencing	
3.5.1	released 8 months ago	

scrnaseq	287	
Single-cell RNA-Seq pipeline for barcode-based protocols such as 10x, DropSeq or SmartSeq, offering a variety of aligners and empty-droplet detection		
10x	genomics	metagenomics
rna	single-cell	assembly
3.6.1	released 5 months ago	

mag	244	
Assembly and binning of metagenomes		
assembly	binning	metagenomics
metagenomics	metaproteomics	metapreprocessing
4.0.0	released 3 months ago	

chipseq	218	
ChIP-seq peak-calling, QC and differential analysis pipeline.		
chip	chip-seq	chromatin-immunoprecipitation
3.0.2	released 11 months ago	

ampliseq	213	
Amplicon sequencing analysis workflow using DADA2 and QIIME2		
amplicon	amplicon-seq	eda
3.1.0	illumina	qiime
2.14.0	released 3 months ago	

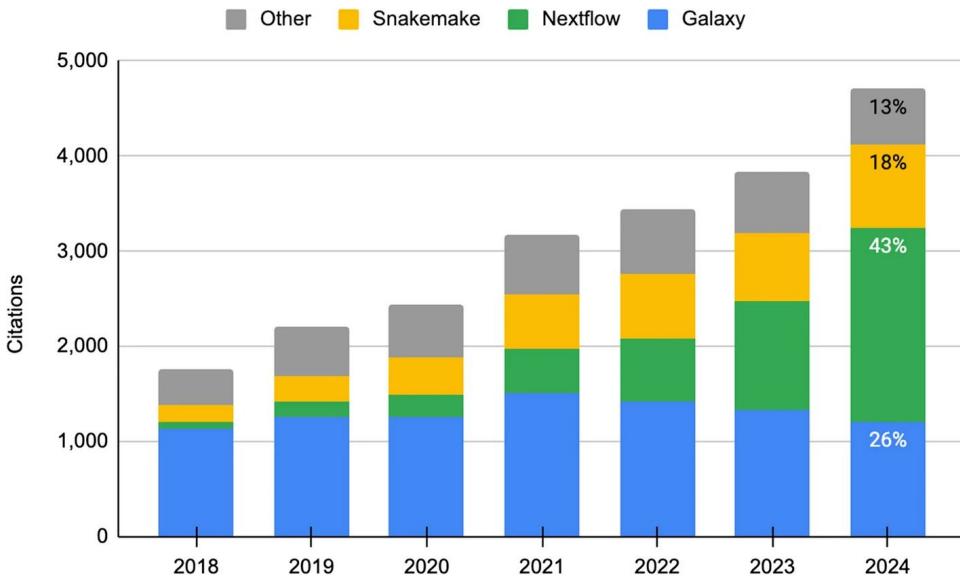


Nextflow and nf-core

- **Nextflow:** Currently most used bioinformatics pipeline framework
- Popularity in large parts thanks to nf-core
- **nf-core:** Community developing high-quality pipelines for different use cases, mostly in bioinformatics

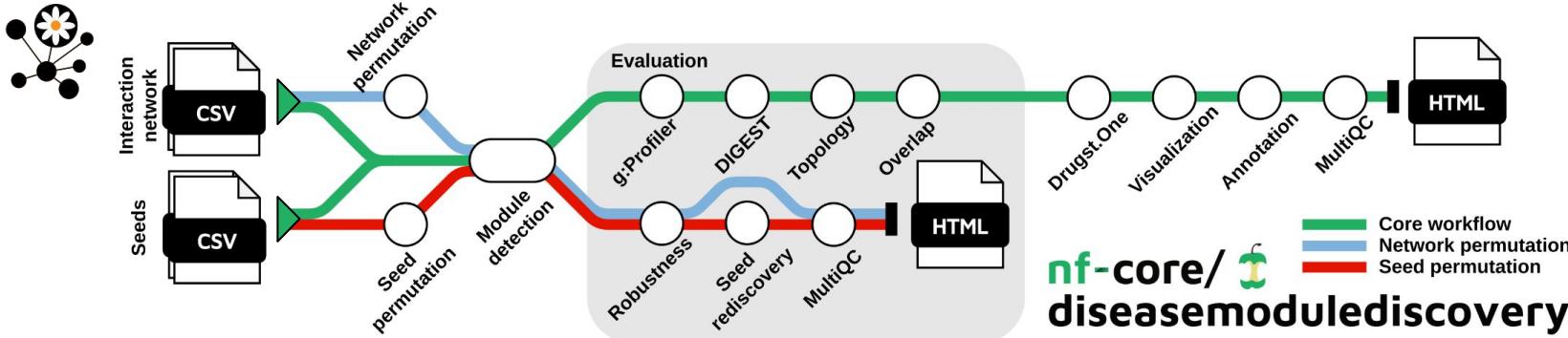
 nextflow

nf-core 





Pipeline overview



Run pipeline

```
nextflow run \
nf-core\diseasemodulediscovery \
--profile docker \
--network PPI.csv \
--seeds seed_genes.csv \
--id_space entrez \
--outdir results
```

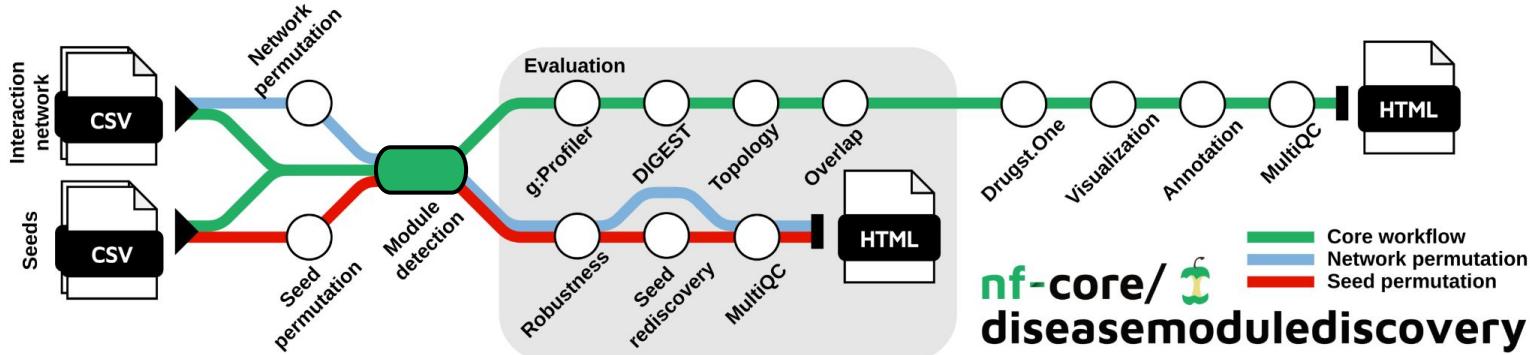
Software Deployment



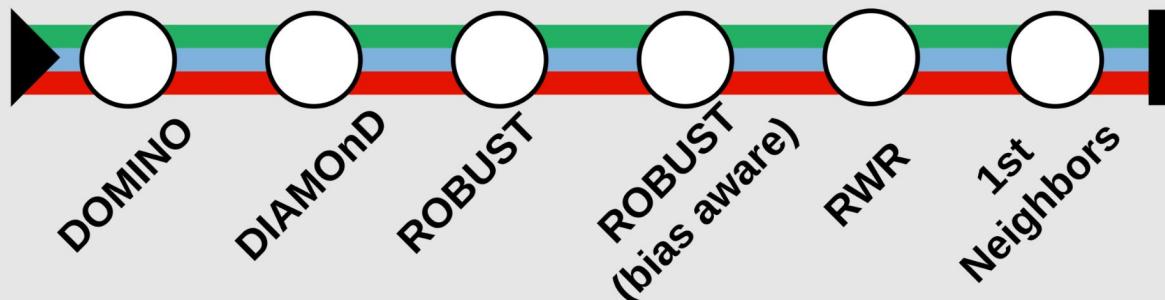
Interaction Network
Provide files *or*
choose predefined

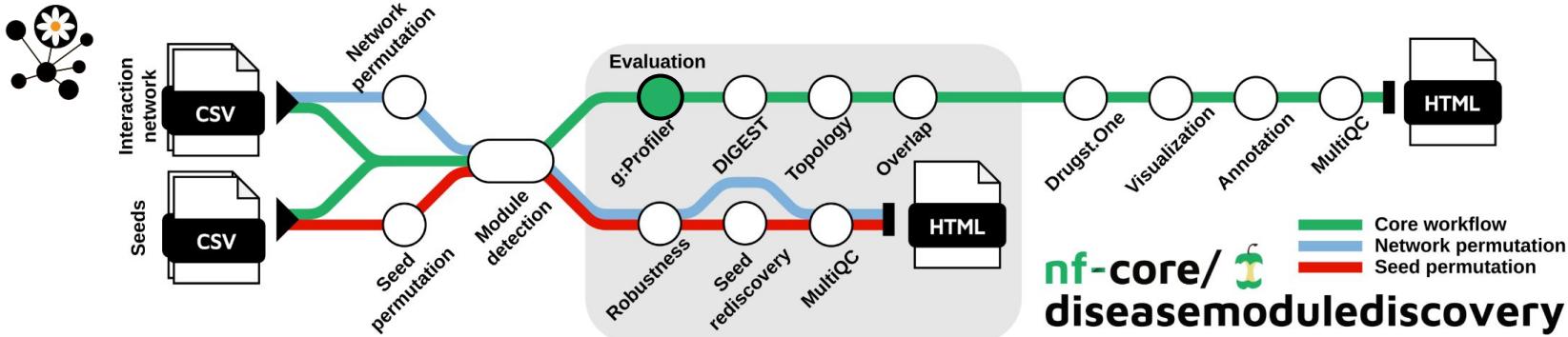


Seeds
Provide files with
seed nodes



Module detection

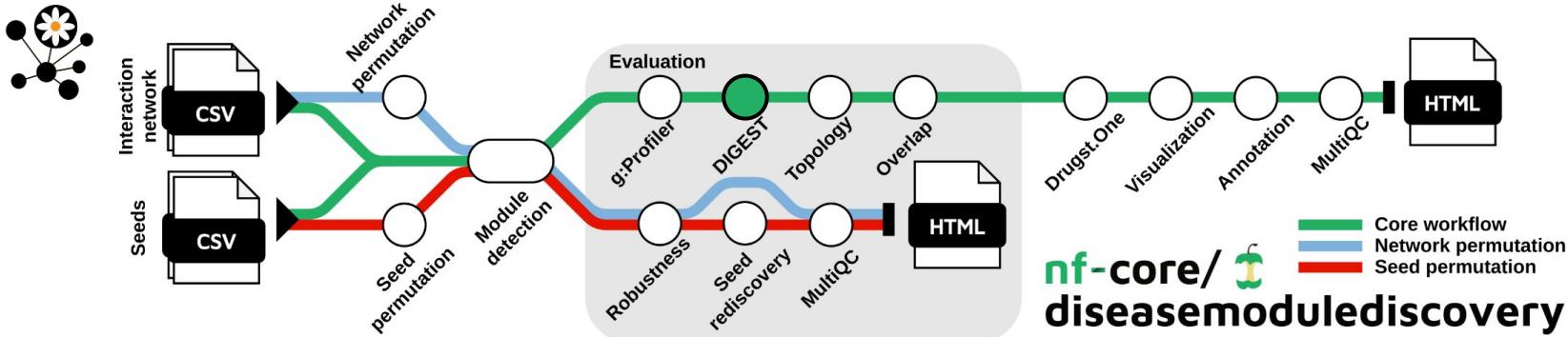




Over-Representation

g:Profiler





Over-Representation

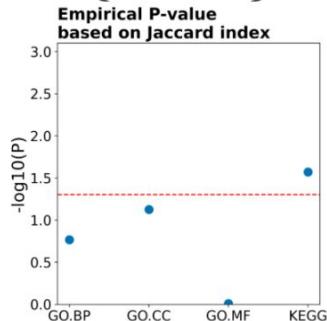
g:Profiler



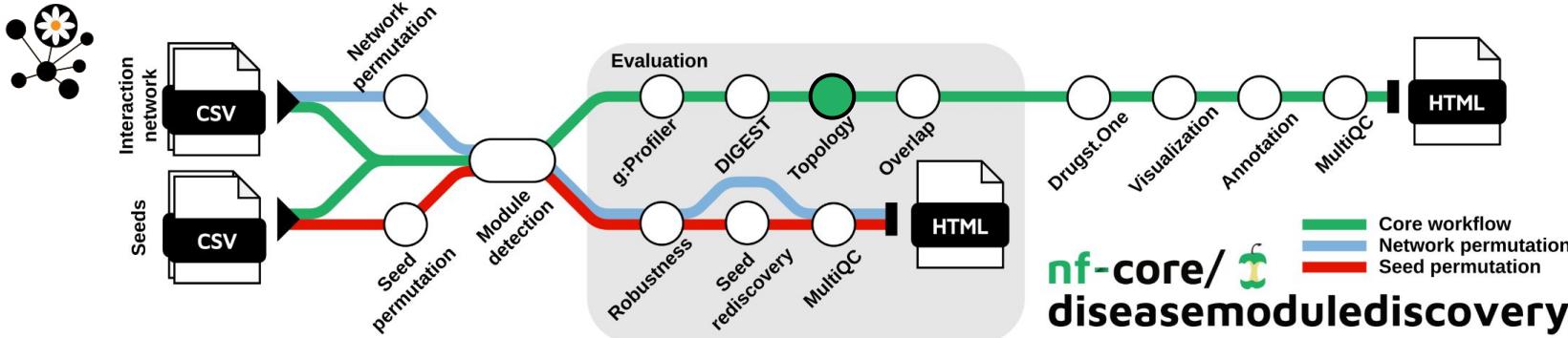
Kolberg et al. *F1000Research* (2020)

Functional Coherence

DIGEST^p

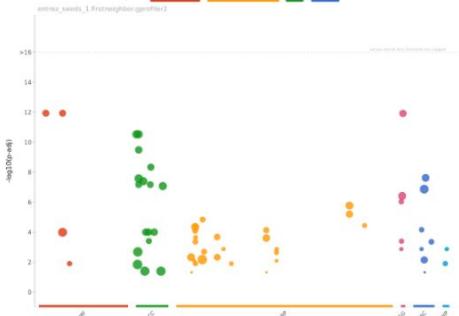


Adamowicz et al.
Briefings in Bioinformatics (2022)



Over-Representation

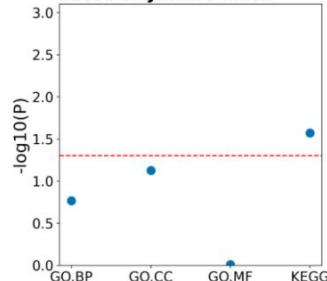
g:Profiler



Functional Coherence

 **DIGEST**^p

Empirical P-value
based on Jaccard index



Kolberg et al. *F1000Research* (2020)

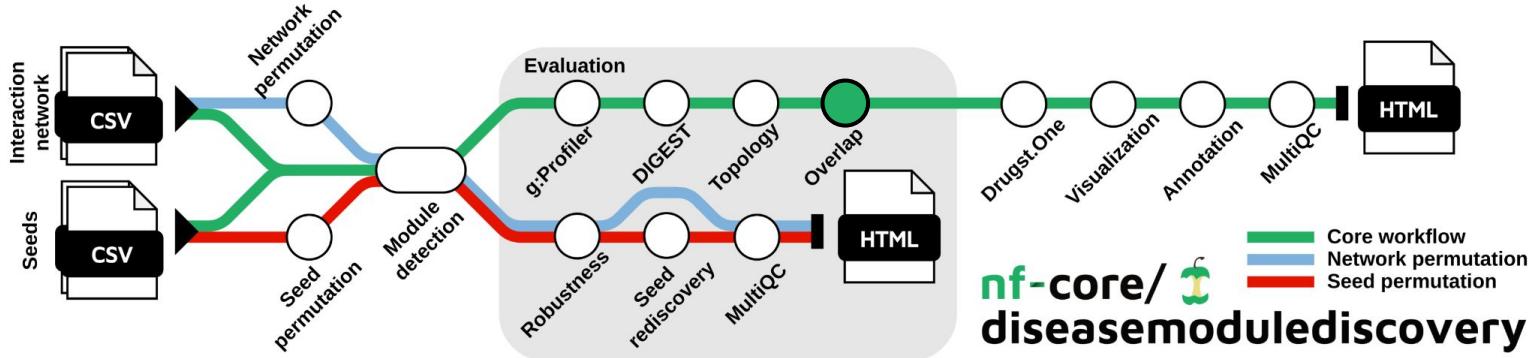
Adamowicz et al.
Briefings in Bioinformatics (2022)

Topology Measures

General Statistics

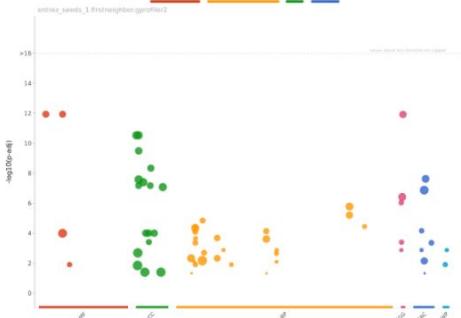
10¹⁰ rows and 5¹⁰ columns.

Sample Name	Nodes	Edges	Seeds	Diameter	Components
entrez_seeds_1_diamond	210	696	10	8	1
entrez_seeds_1_domino	5	6	4	3	1
entrez_seeds_1_firstneighbor	83	264	10	9	1
entrez_seeds_1_robust	24	29	10	8	1
entrez_seeds_1_robust_bias_aware	24	31	10	8	1
entrez_seeds_2_diamond	207	672	7	8	2
entrez_seeds_2_domino	7	15	4	2	1
entrez_seeds_2_firstneighbor	54	207	7	6	3
entrez_seeds_2_robust	31	43	6	6	1
entrez_seeds_2_robust_bias_aware	31	43	6	6	1



Over-Representation

g:Profiler

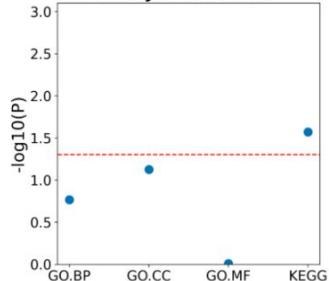


Kolberg et al. *F1000Research* (2020)

Functional Coherence

DIGEST^p

Empirical P-value
based on Jaccard index



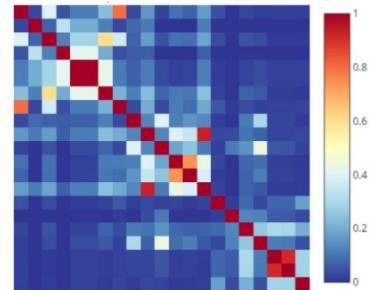
Adamowicz et al.
Briefings in Bioinformatics (2022)

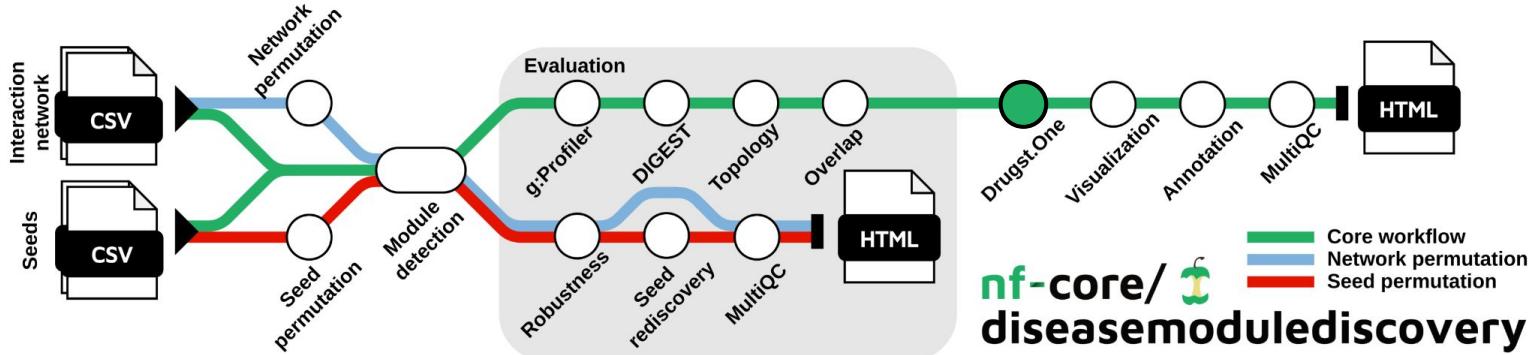
Topology Measures

General Statistics

Sample Name	Nodes	Edges	Seeds	Diameter	Components
entrez_seeds_1_diamond	210	696	10	8	1
entrez_seeds_1_domino	5	6	4	3	1
entrez_seeds_1_firstrneighbo	83	264	10	9	1
entrez_seeds_1_robust	24	29	10	8	1
entrez_seeds_1_robust_bias_aware	24	31	10	8	1
entrez_seeds_2_diamond	207	672	7	8	2
entrez_seeds_2_domino	7	15	4	2	1
entrez_seeds_2_firstrneighbo	54	207	7	6	3
entrez_seeds_2_robust	31	43	6	6	1
entrez_seeds_2_robust_bias_aware	31	43	6	6	1

Module Overlaps

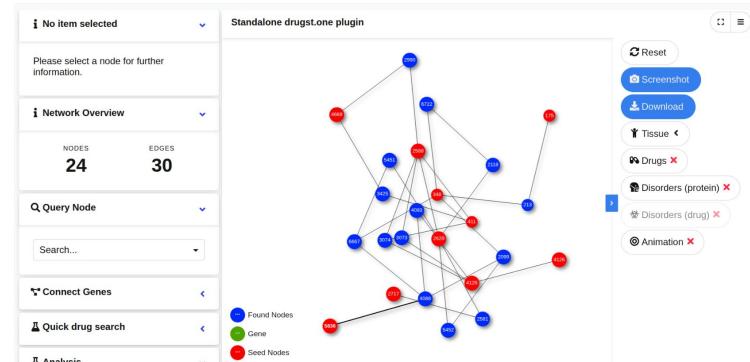




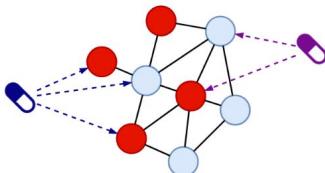
nf-core/  **diseasemodulediscovery**



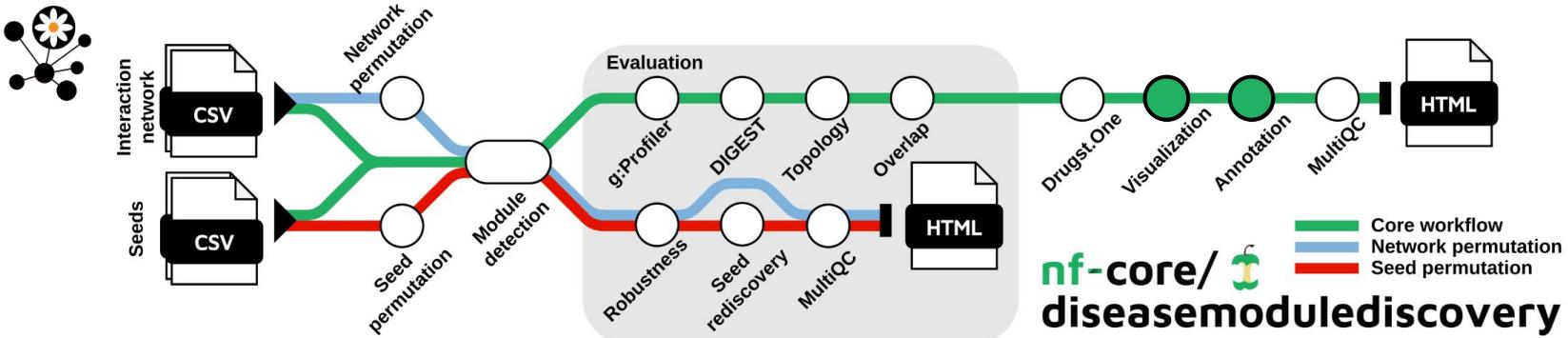
Web Export



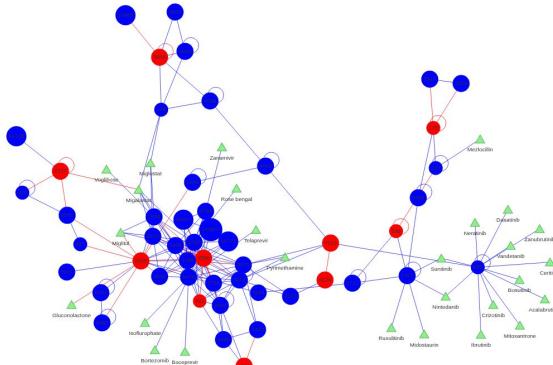
Drug Prioritization



- TrustRank
- Degree Centrality
- Harmonic Centrality



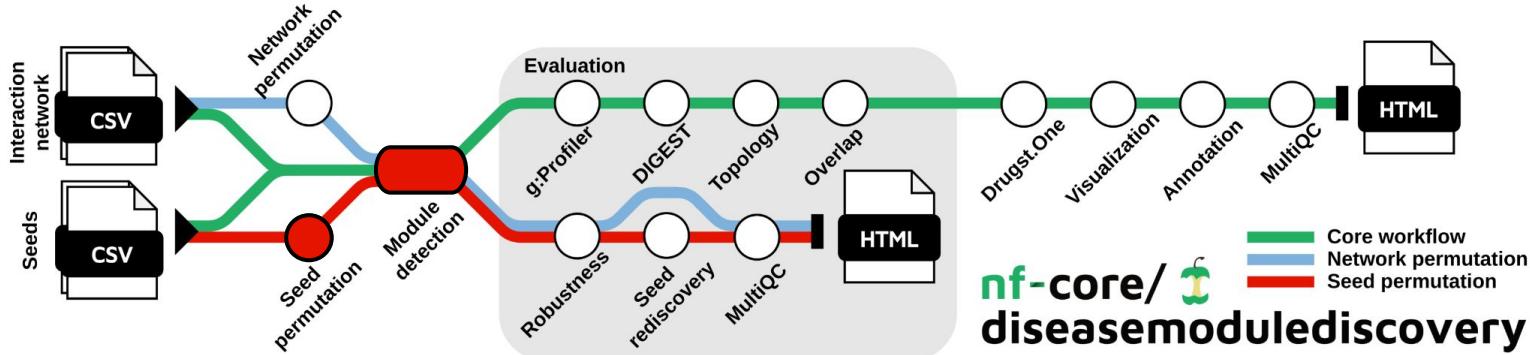
Network Visualization



Module Annotation

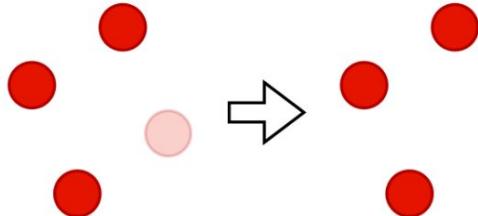


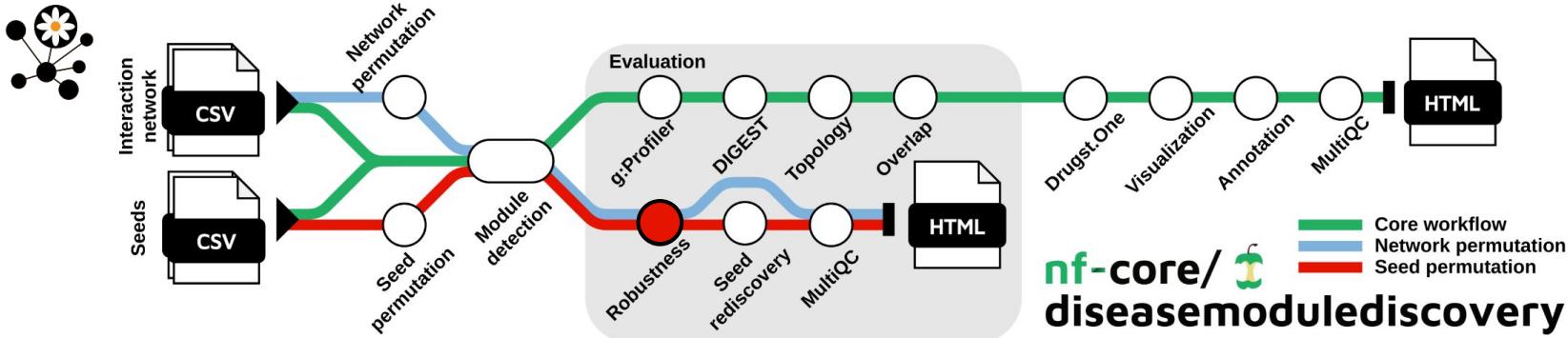
- Targeting drugs
 - Side effects
 - Disorders
 - Cellular component



Seed Permutation

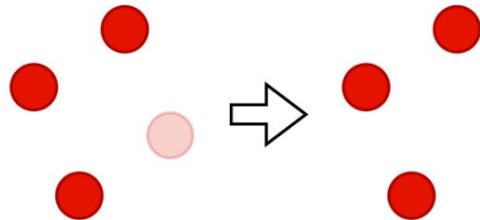
Leave-one-out approach
(repeat for every seed gene)





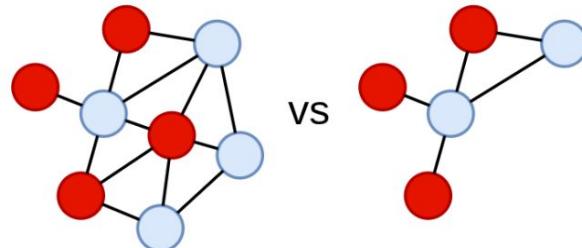
Seed Permutation

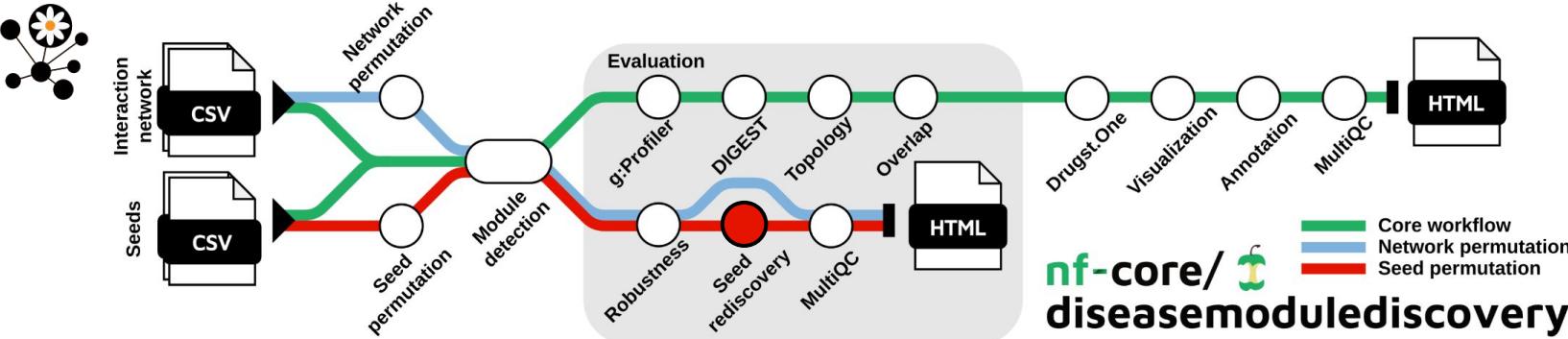
Leave-one-out approach
(repeat for every seed gene)



Robustness

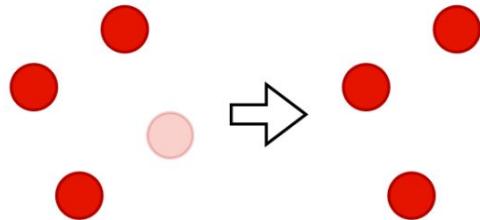
How much do the modules change?





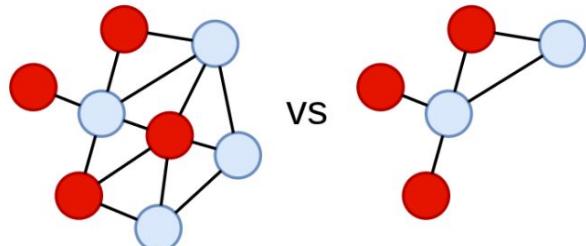
Seed Permutation

Leave-one-out approach
(repeat for every seed gene)



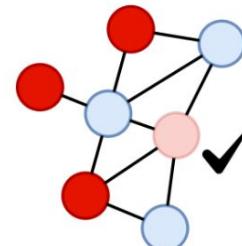
Robustness

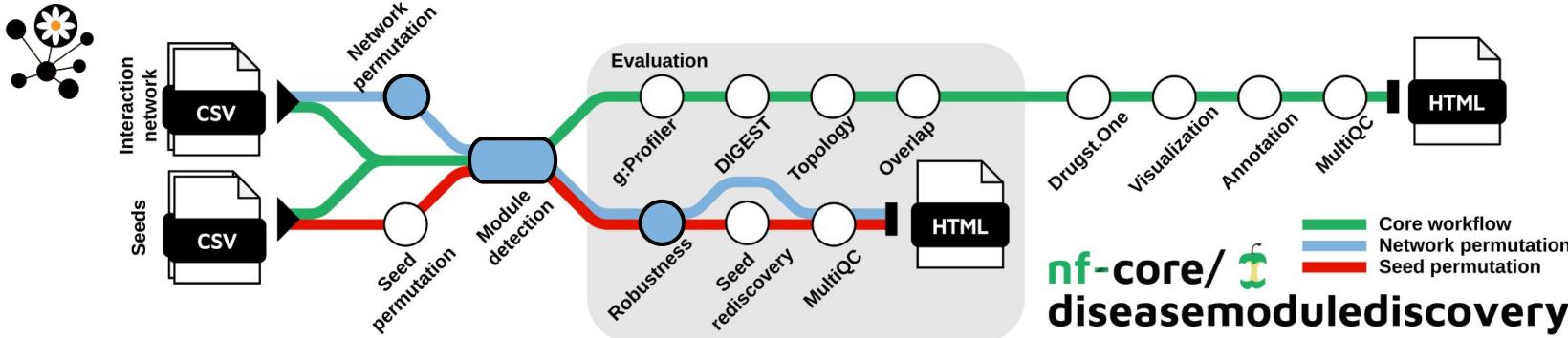
How much do the modules change?



Seed Rediscovery

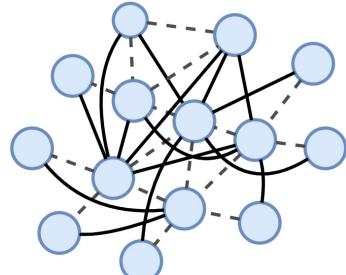
Can the removed seed genes be recovered?





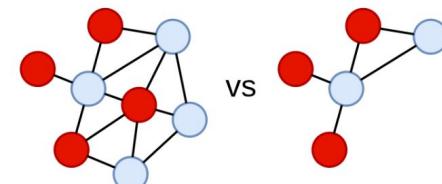
Network Permutation

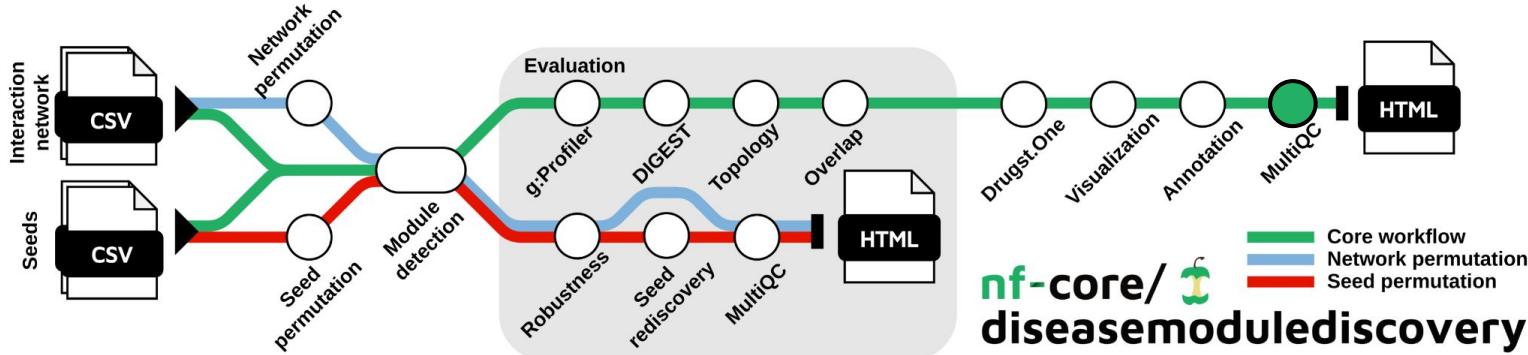
Degree-preserving rewiring



Module detection methods are prone to learning from node degrees, instead of individual interactions.

- If module remains largely **the same**, it is **based on node degrees**.





TUM

nf-core/  **diseasemodulediscovery**

 multiqc

- HTML pipeline report
- Summarizes modules and evaluation results
- Execution details and software versions





Available PPI networks

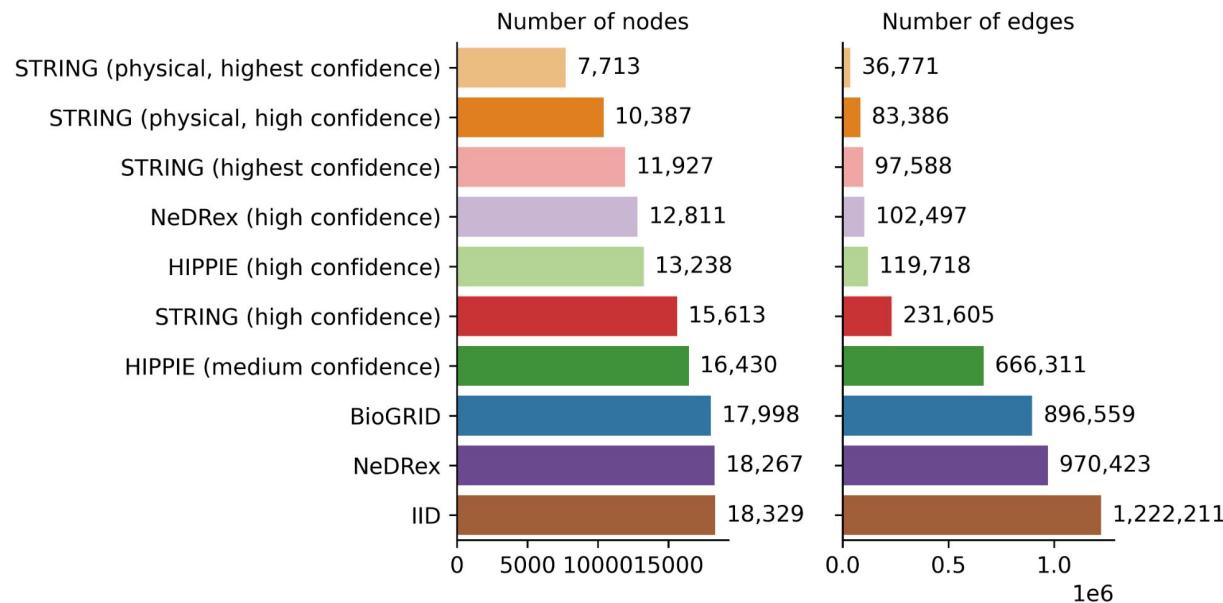
Pipeline provides selection of PPI networks in different ID spaces

PPI networks:

- BioGRID
- IID
- HIPPIE (2 variants)
- NeDRex (2 variants)
- STRING (4 variants)

ID spaces:

- Uniprot proteins
- Ensembl genes
- Entrez genes
- Gene symbols





Output

Most **output files are named** according to the seed file, network file, and method used, for example:
`<seed_file_name>.<network_file_name>.<method_name>.tsv`

A typical run directory has the following **structure**:

└── {outdir}	
└── input	Pre-processed input files
└── modules	Disease modules in different formats
└── modules_visualized	Network visualizations of the modules
└── evaluation	Module evaluation results
└── digest	Functional coherence analysis
└── gprofiler	Over-representation analysis
└── network_permutation	Network permutation analysis
└── seed_permutation	Seed set permutation analysis
└── drug_prioritization	Drug rankings
└── modules_visualized_with_drugs	Network visualizations including drugs
└── multiqc	Pipeline summary report
└── pipeline_info	Run statistics (memory usage, etc...)
└── reports	Additional reports
└── work	Cached intermediate files



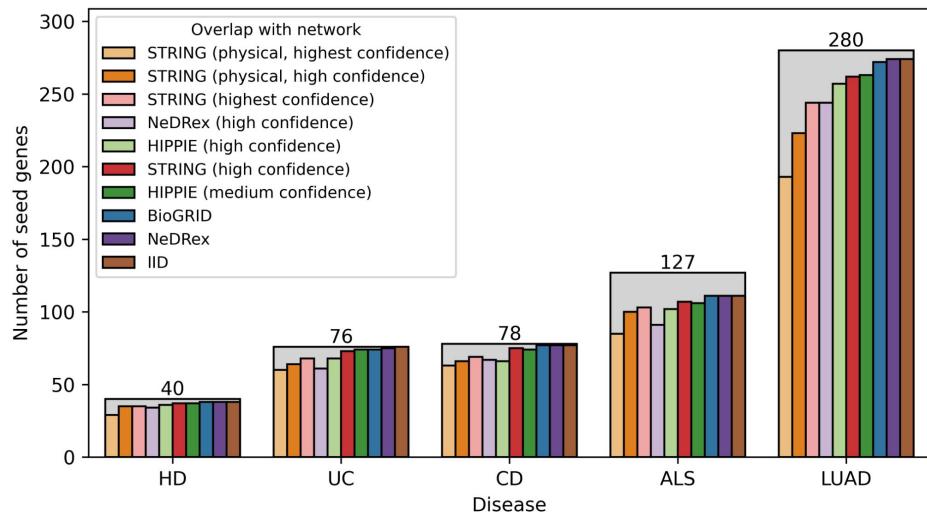
Some results



Pipeline demonstration

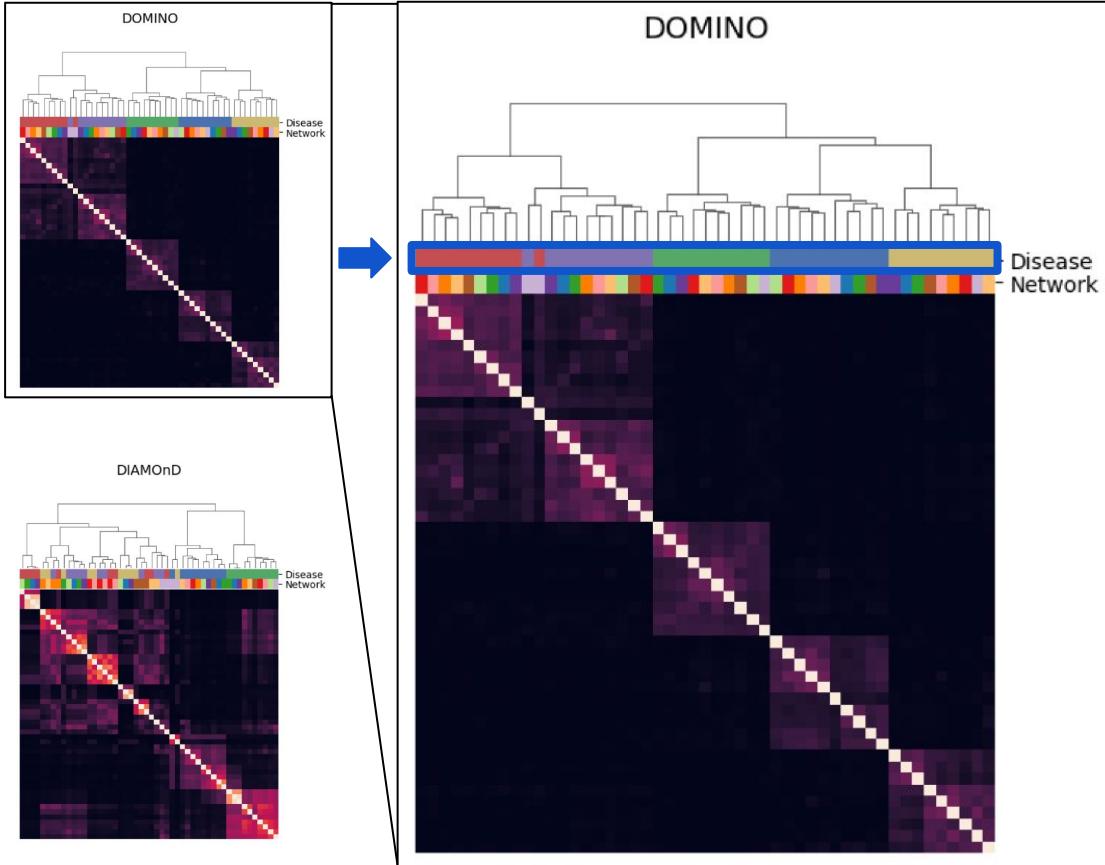
10 networks x 5 diseases = 50 input combinations

- Huntington's disease (**HD**)
- Ulcerative colitis (**UC**)
- Crohn's disease (**CD**)
- Amyotrophic lateral sclerosis (**ALS**)
- Lung adenocarcinoma (**LUAD**)

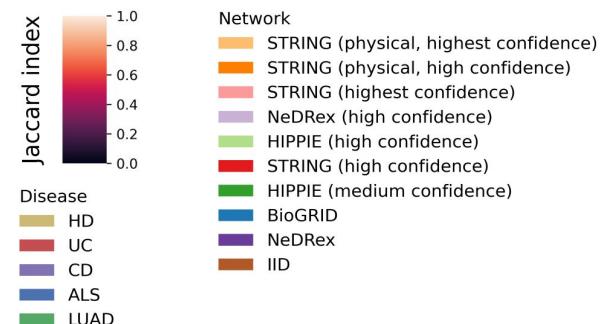




Module overlaps

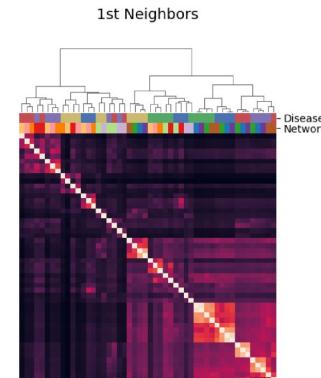
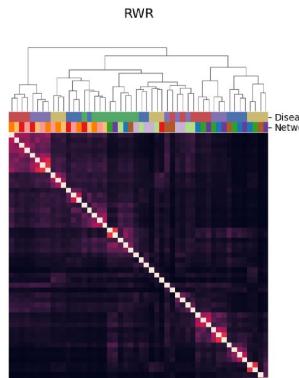
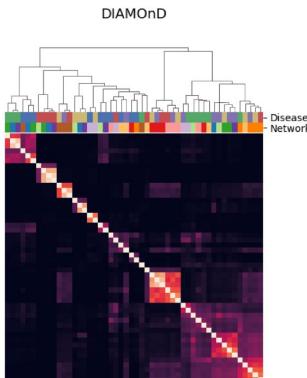
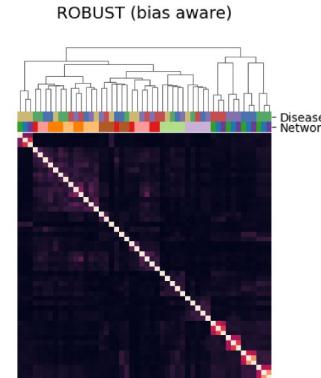
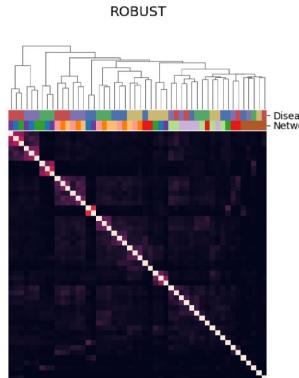
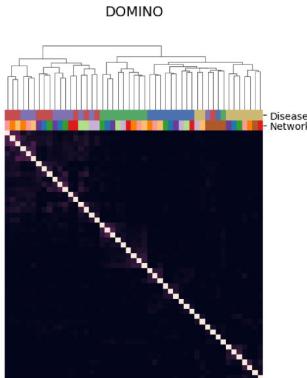


***DOMINO, ROBUST, and
ROBUST (bias aware)
yield disease-specific modules***

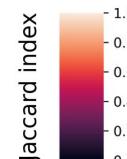




Module overlaps – excluding the seeds...



However, the disease specificity primarily arises from the seed nodes...



Network

- STRING (physical, highest confidence)
- STRING (physical, high confidence)
- STRING (highest confidence)
- NeDReX (high confidence)
- HIPPIE (high confidence)
- STRING (high confidence)
- HIPPIE (medium confidence)
- BioGRID
- NeDReX
- IID

Disease

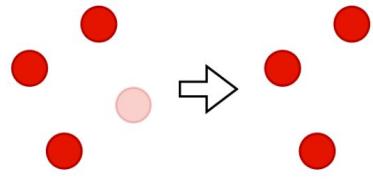
- HD
- UC
- CD
- ALS
- LUAD



Seed rediscovery

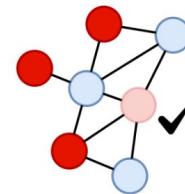
Seed Permutation

Leave-one-out approach
(repeat for every seed gene)



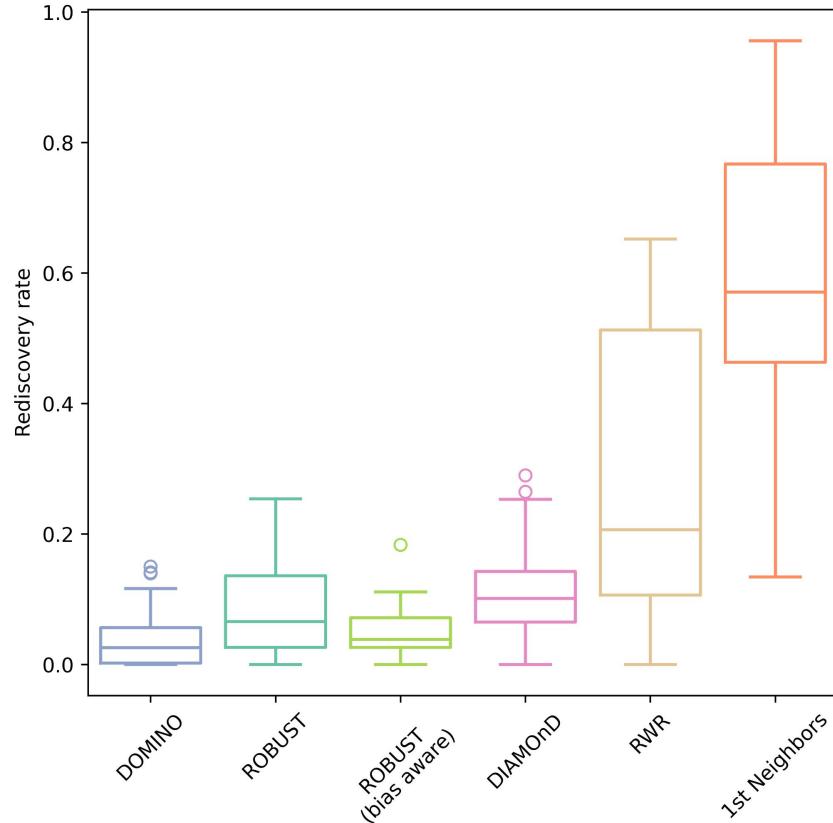
Seed Rediscovery

Can the removed seed genes
be recovered?



Including all direct interactors is the most reliable way to recover individual seeds

However, there can be a lot of direct interactors...



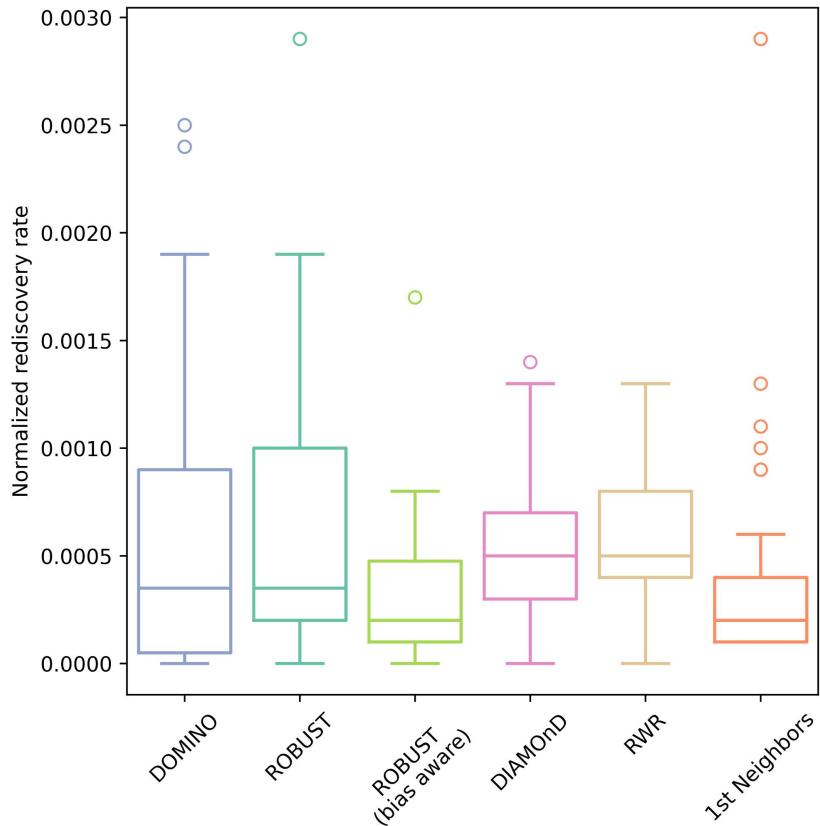


Seed rediscovery – normalized...

Normalized rediscovery rate is divided by the **number of nodes** in the module

Differences between tools are less prominent

DIAMOnD and **RWR** yield highest median recovery





Seed rediscovery – gene level

*Can individual **highly important genes** be rediscovered?*

HTT is among the most frequently rediscovered seeds in Huntington's Disease

*But: Most methods **cannot rediscover it reliably***

Removed seeds	1st Neighbors	RWR	ROBUST	ROBUST (bias aware)	DIAMOnD	DOMINO
IL6R	10	8	5	1	2	1
HTT	10	7	5	0	1	0
MAOB	10	0	1	3	0	0
PPARGC1A	9	5	4	2	1	1
IL6	10	7	1	0	3	0
NPY2R	9	5	3	1	0	0
NPY	8	5	3	2	0	0
GRIN2B	9	6	3	0	0	0
CNTF	8	6	2	0	1	0
SIRT1	9	5	2	0	0	0
MAOA	10	4	0	1	0	0
HAP1	10	3	2	0	0	0
GRIN2A	9	5	0	0	0	0
NRF1	9	4	0	0	0	0
BDNF	6	4	1	0	1	0
CNR1	5	2	2	2	0	0
IGF1	3	3	0	0	3	0
MAP3K5	6	1	0	0	0	0
PRKAA1	6	1	0	0	0	0
EIF2AK2	5	1	0	0	0	0
AIFM1	5	1	0	0	0	0
PRNP	4	1	1	0	0	0
GDNF	3	3	0	0	0	0
DIABLO	5	0	0	0	0	0
FAAH	2	1	0	2	0	0
TFAM	3	2	0	0	0	0
GRIK2	4	0	0	0	1	0
HSF1	5	0	0	0	0	0
ATF5	1	0	0	0	0	0
ABAT	1	0	0	0	0	0
SLC2A3	1	0	0	0	0	0
RCAN1	1	0	0	0	0	0
GLUL	1	0	0	0	0	0
NOG	1	0	0	0	0	0
JPH3	0	0	0	0	0	0
IP6K2	0	0	0	0	0	0
OGG1	0	0	0	0	0	0
SLC29A1	0	0	0	0	0	0



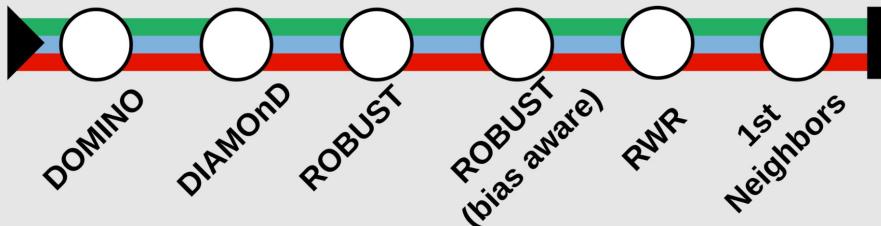
Summary



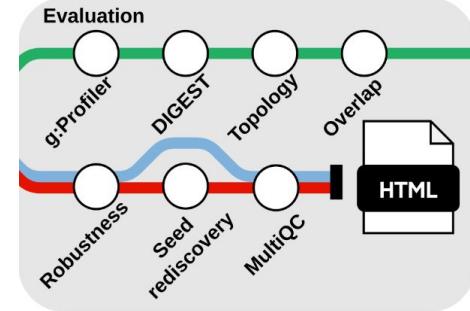
Pipeline offers easy access to...

Disease module discovery methods

Module detection



Evaluation framework



Systematically
compare methods

Modular
expandable design

Track **progress**
in the field



Acknowledgments



Developers and co-authors

Johannes Kersting

Lisa Spindler

Joaquim Aguirre-Plans

Chloé Bucheron

Quirin Manz

Fernando Delgado-Chaves

Tanja Pock

Mo Tan

Cristian Nogales

Andreas Maier

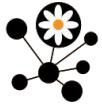
Jan Baumbach

Jörg Menche

Emre Guney



Funded by the
European Union



Hands-on



Hands-on

Focuses on:

- Running the pipeline on the Huntington's Disease data obtained earlier
- Working with the pipeline output and MultiQC report

Instructions are in the **Session3_nextflow_pipeline** folder:

https://github.com/REPO4EU/network_medicine_drug_repurposing_tutorial/tree/main/Ses sion3_nextflow_pipeline

Best to work with it directly through the [GitHub web view](#):

The screenshot shows a GitHub repository interface. At the top, there are three columns: 'Name', 'Last commit message', and 'Last commit da...'. Below this, there are two commits:

Name	Last commit message	Last commit da...
...		
README.md	completed Task 4	19 minutes ago
nextflow.config	added custom memory configuration	19 hours ago

At the bottom of the repository page, the text 'Nextflow pipeline for drug repurposing: hands-on' is displayed in large white font on a dark background, followed by a horizontal line and the text 'Before starting these hands-on tasks, ensure that you have completed the [Nextflow pipeline setup](#)'.