

No Metric to Rule Them All



Toward Principled Evaluations
of Graph-Learning Datasets





The Fellowship



 **Corinna
Coupette** **A!**



 **Jeremy
Wayland** 



 **Emily
Simons** 



**UNI
FR** **Bastian
Rieck** 

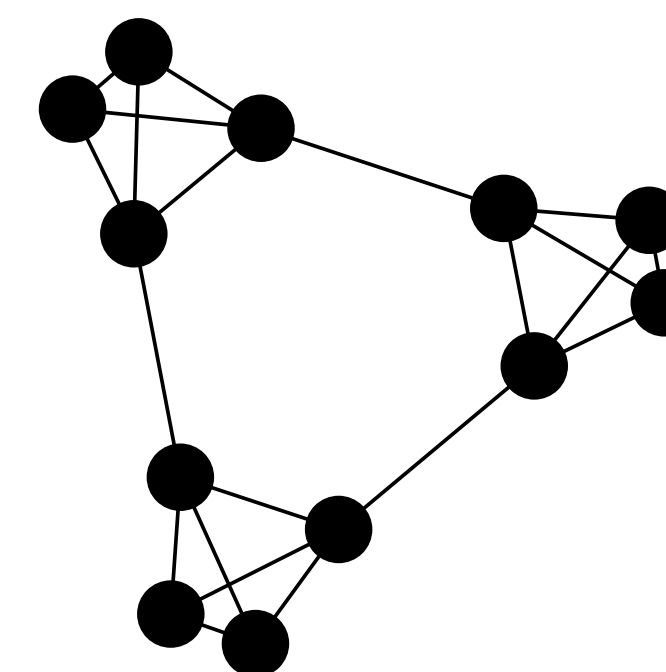


Graph Learning Benchmarks need a reality check.

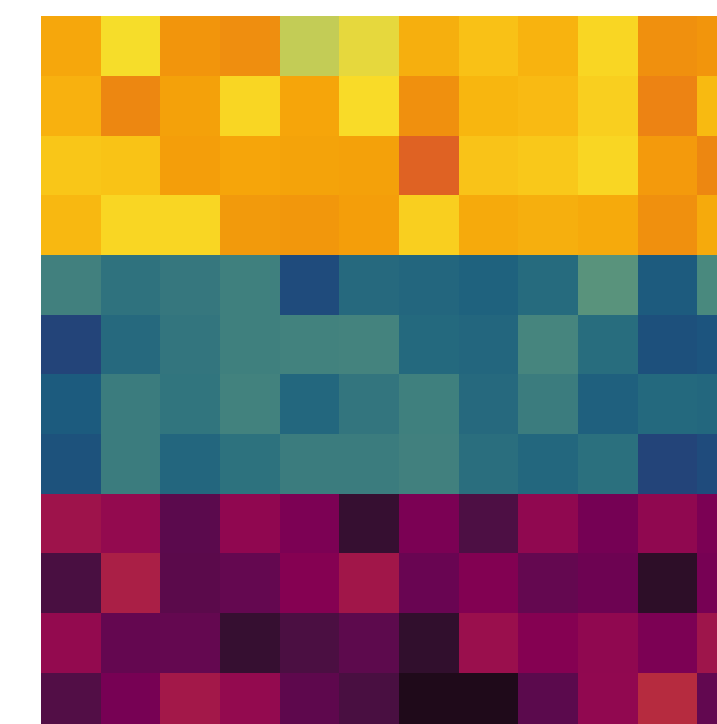


Guiding Questions

Q1: What characterizes a good graph-learning dataset?



Q2: How can we evaluate dataset quality in graph learning?

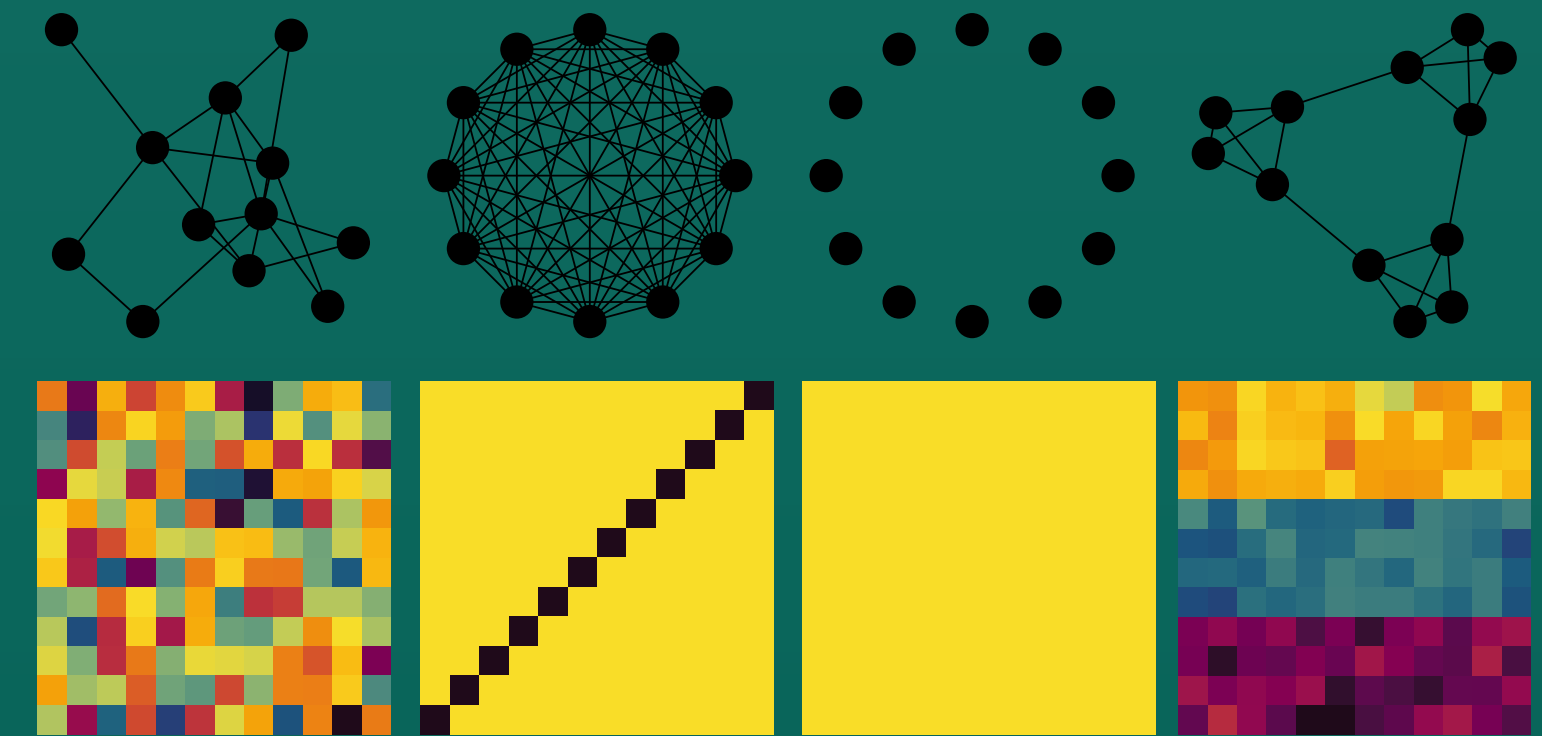




RINGSO

Relevant **I**nformation in **N**ode
features and **G**raph **S**tructure

Perturbation
Framework



Dataset
Evaluation

Benchmarks	
Keep	✓
Realign	⚠
Deprecate	✗



Guiding Principles

P1: *Task-Relevant Information*

→ Both structure and features should matter for a given task.

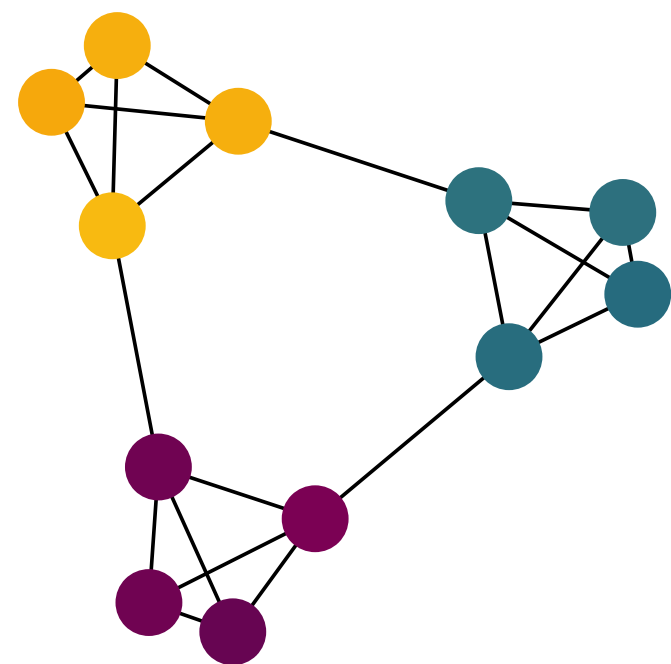
P2: *Complementary Information*

→ Structure and features should offer distinct views of the data.

Associated Metric

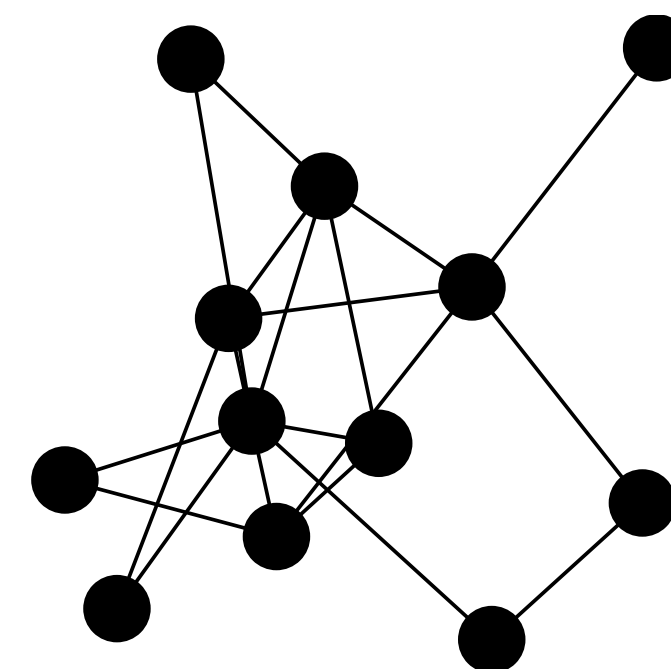
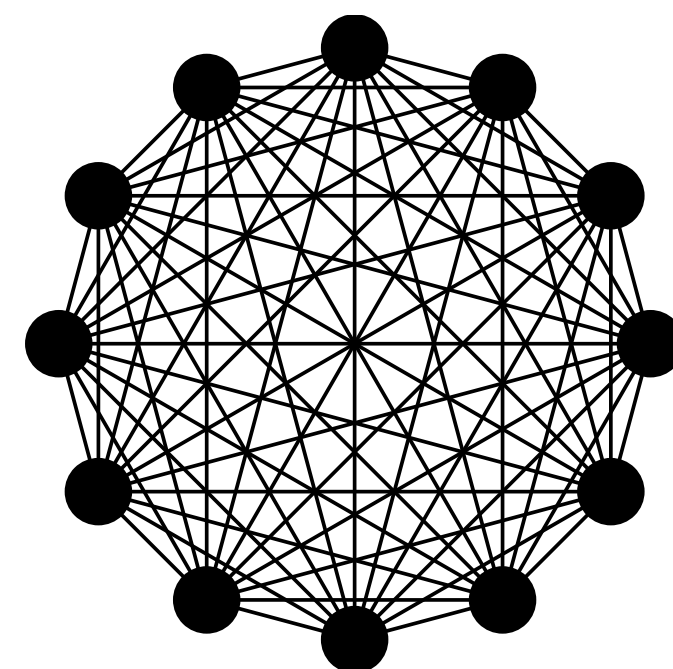
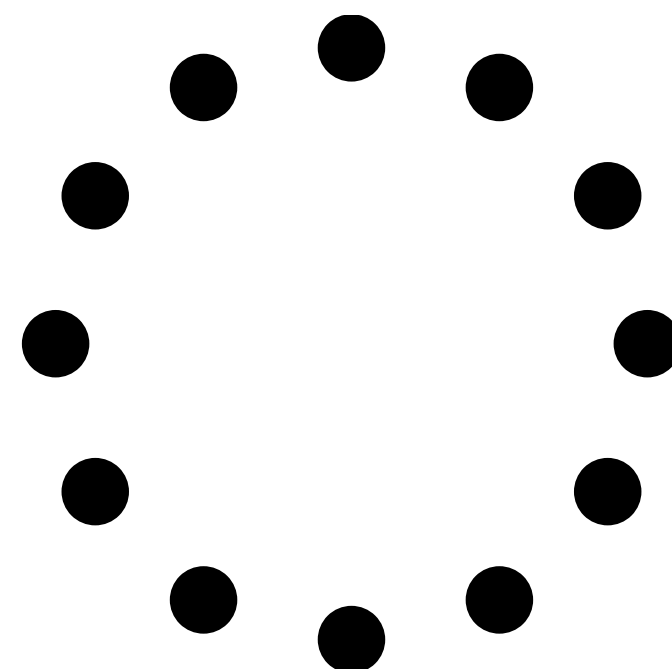
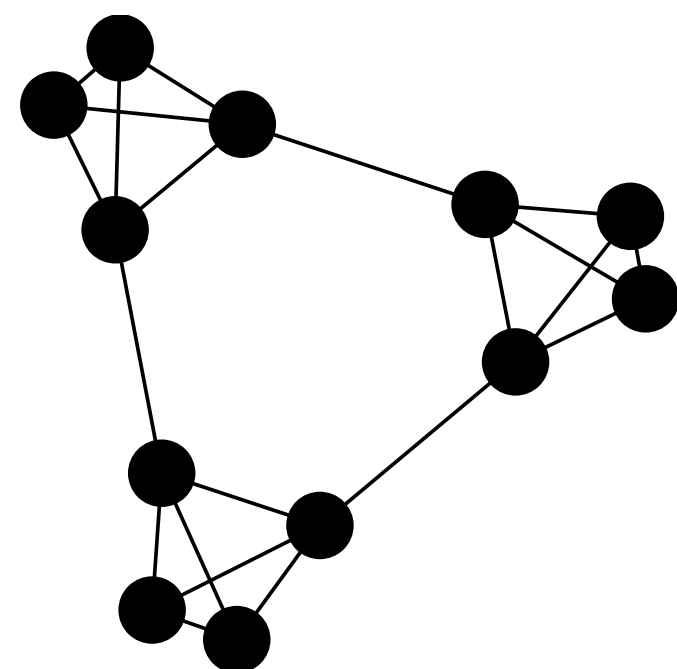
Performance
separability

Mode
complementarity

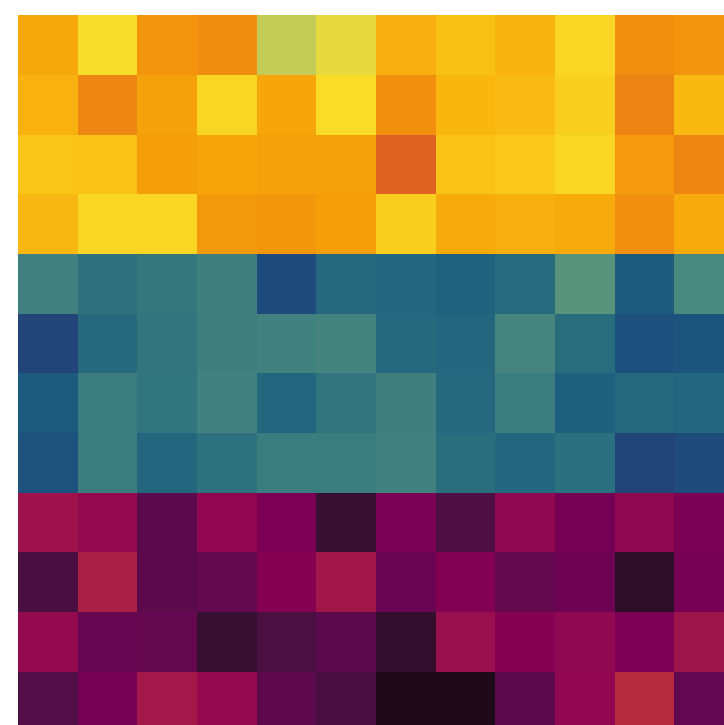


Mode Perturbations

Graph



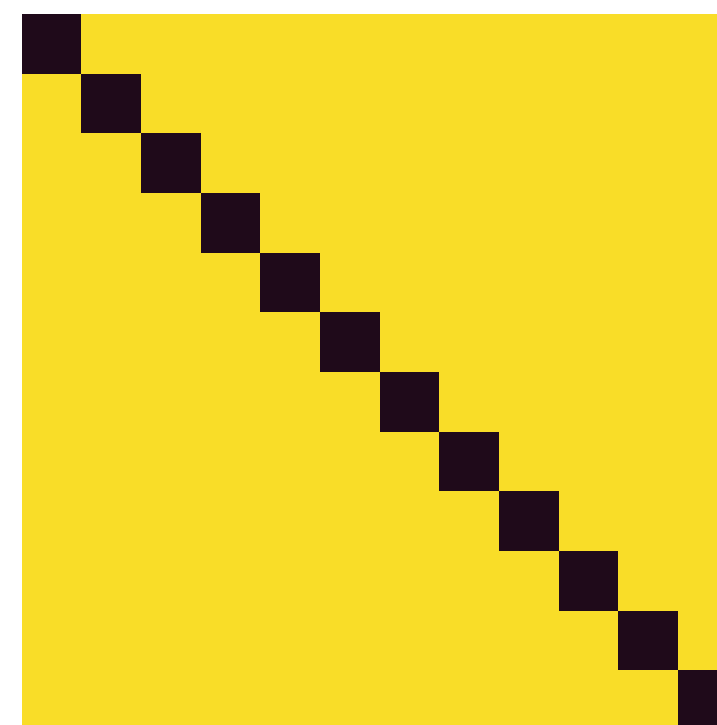
Features



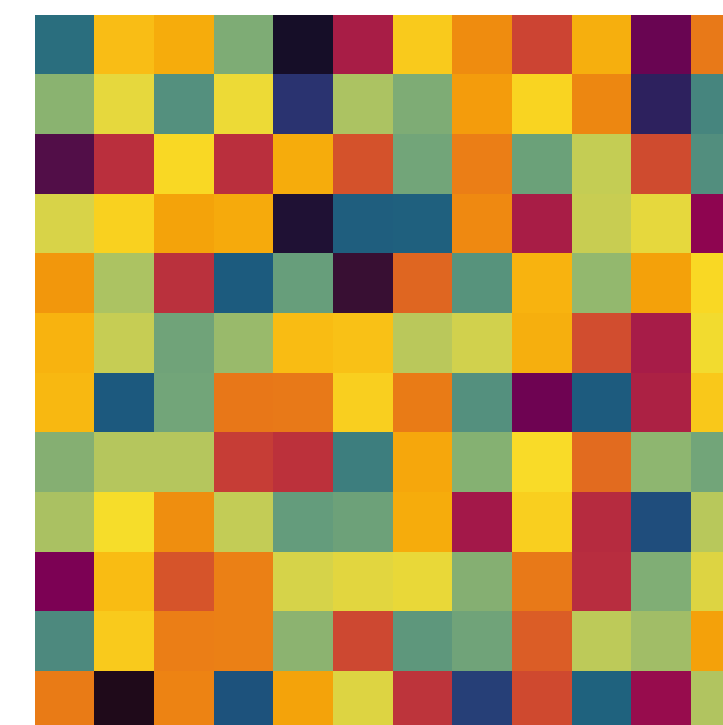
Identity
Perturbation



Empty
Perturbation



Complete
Perturbation



Random
Perturbation



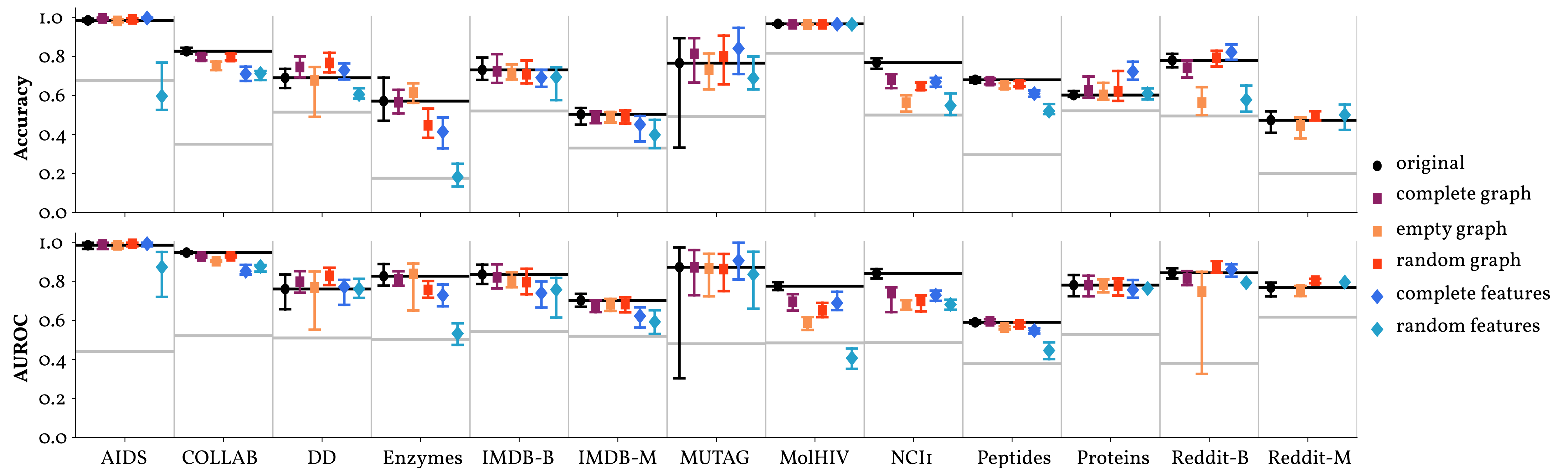
Performance Separability



A Task-Dependent Measure



Many popular datasets fail to statistically outperform their perturbed versions.



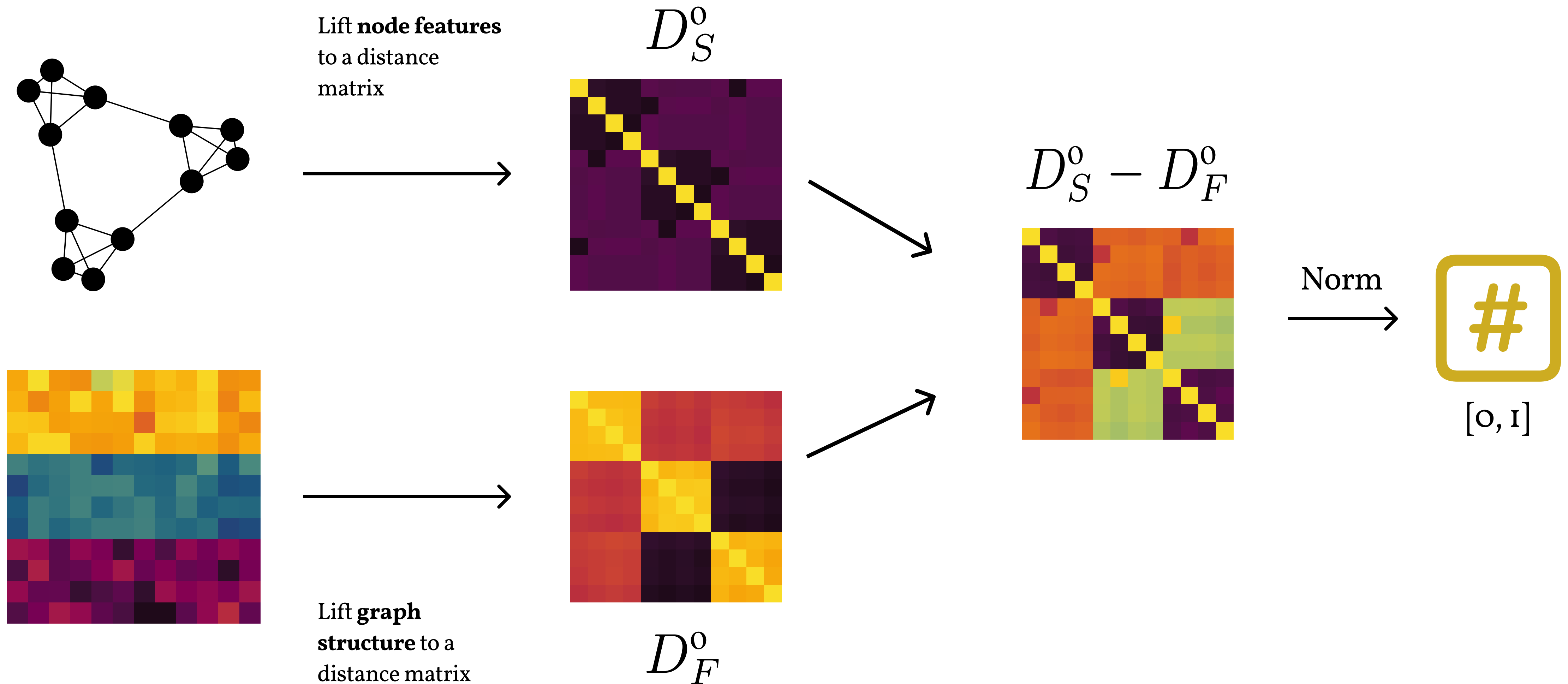


Mode Complementarity



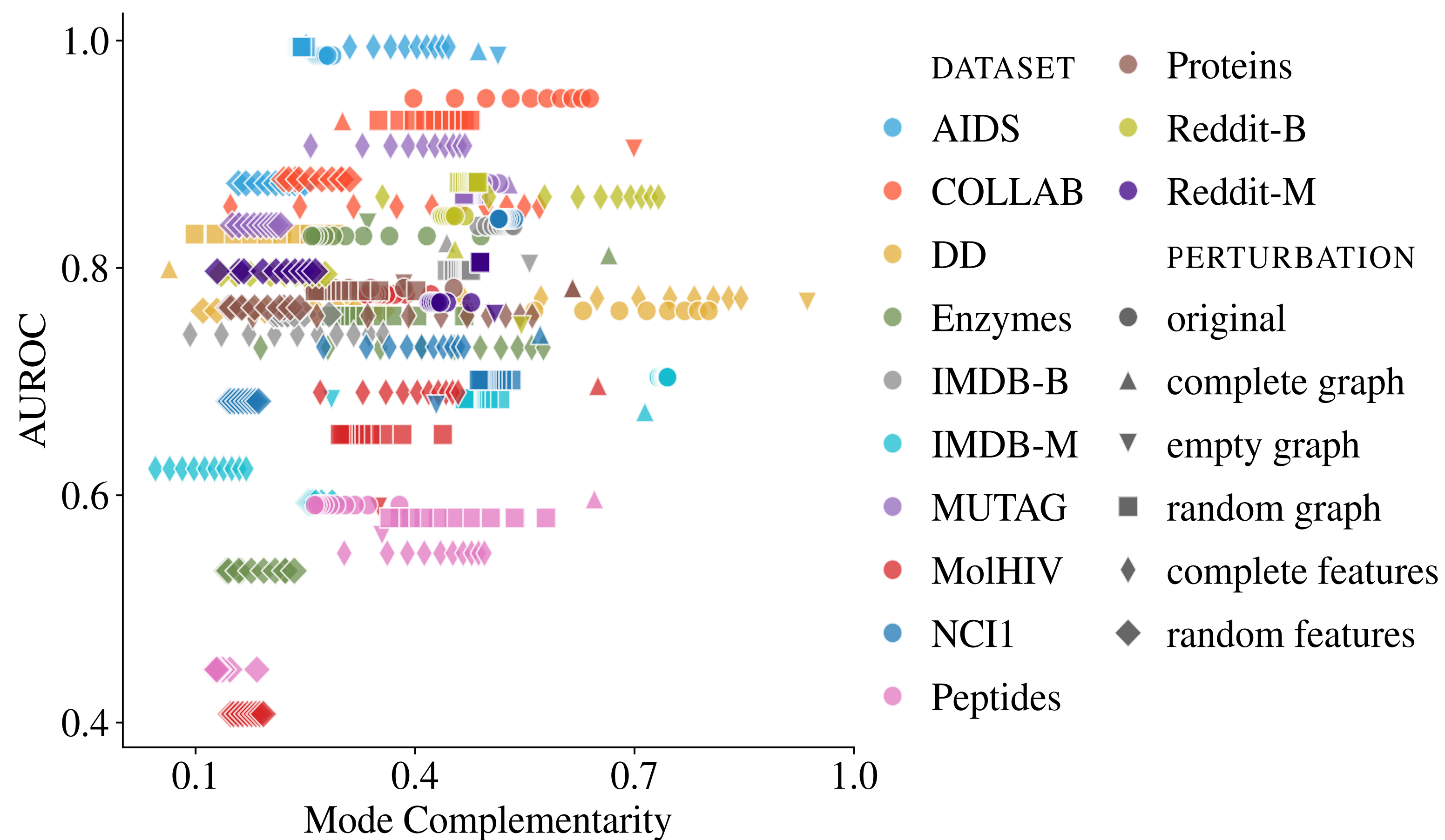
A Task-Independent Measure

Calculating Mode Complementarity





Mode complementarity correlates with performance.



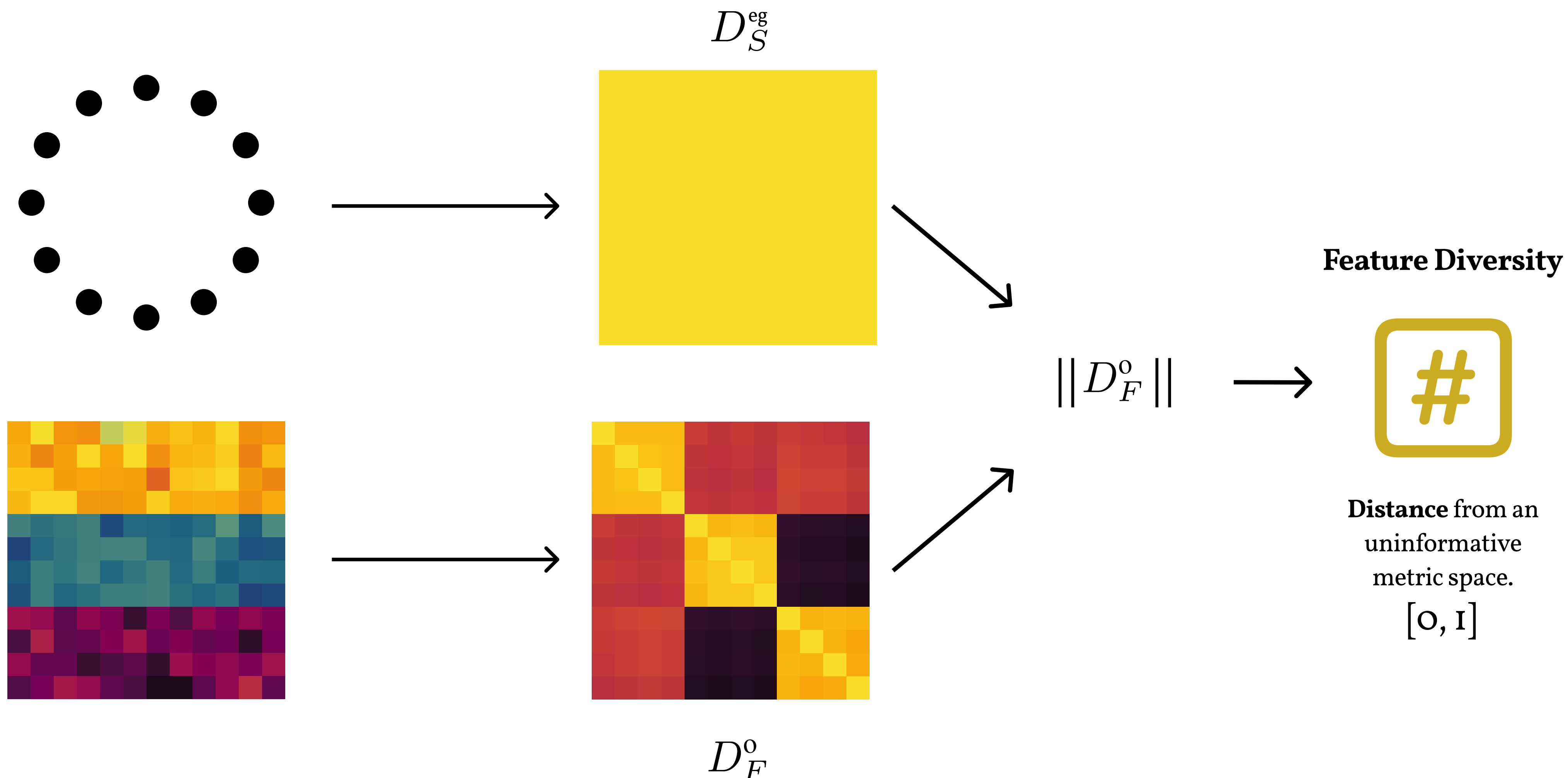
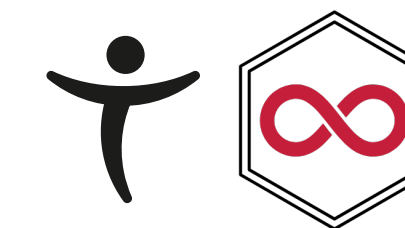


Mode Diversity

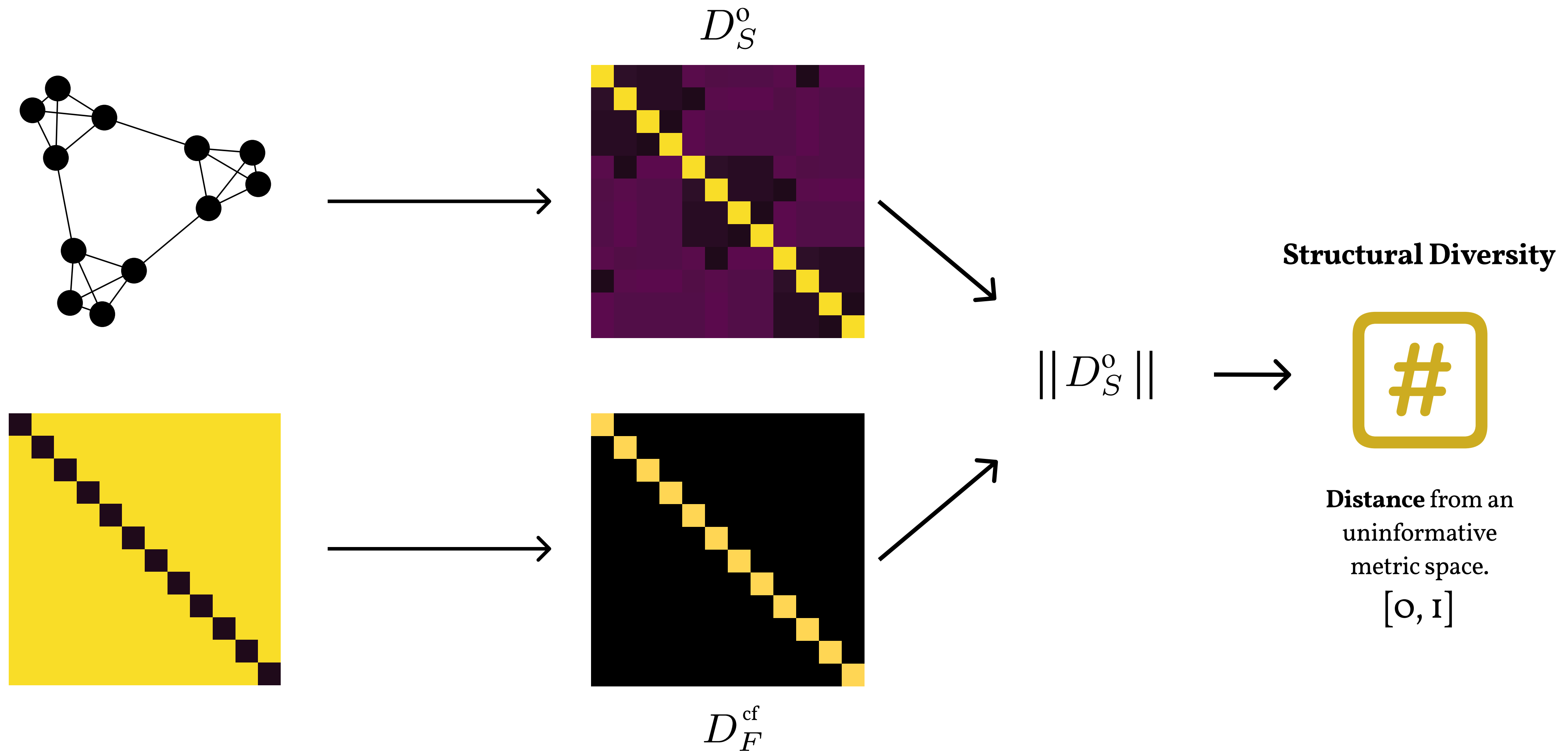


A Task-Independent Measure

Calculating Feature Mode Diversity



Calculating **Structural** Mode Diversity





Final Evaluation



A Benchmark Taxonomy

Result Overview

Dataset	Performance Separability (†)			Mode Diversity (‡)		Action
	Structure	Features	Overall	Structure ($\mu:\sigma$)	Features ($\mu:\sigma$)	
MolHIV	+	+	++	○:--	++:++	Keep († ‡)
NCI1	+	+	++	○:--	++:++	
Peptides	○	+	+	++:--	++:++	
AIDS	-	-	--	○:-	++:○	Realign († ‡)
DD	-	-	--	++:-	--:--	
MUTAG	○	-	-	○:--	++:○	
Reddit-{B,M}	-	-	--	++:--	++:○	
COLLAB	+	+	++	-:++	-:++	Deprecate (‡)
IMDB-B	○	○	○	--:○	○:++	
IMDB-M	○	+	+	--:○	-:++	
Enzymes	-	+	-	-:-	++:○	Deprecate († ‡)
Proteins	-	-	--	-:○	++:-	

RINGSO

Relevant **I**nformation in **N**ode
features and **G**raph **S**tructure



Desiderata

Both modes should matter for a
given task.



Performance Separability

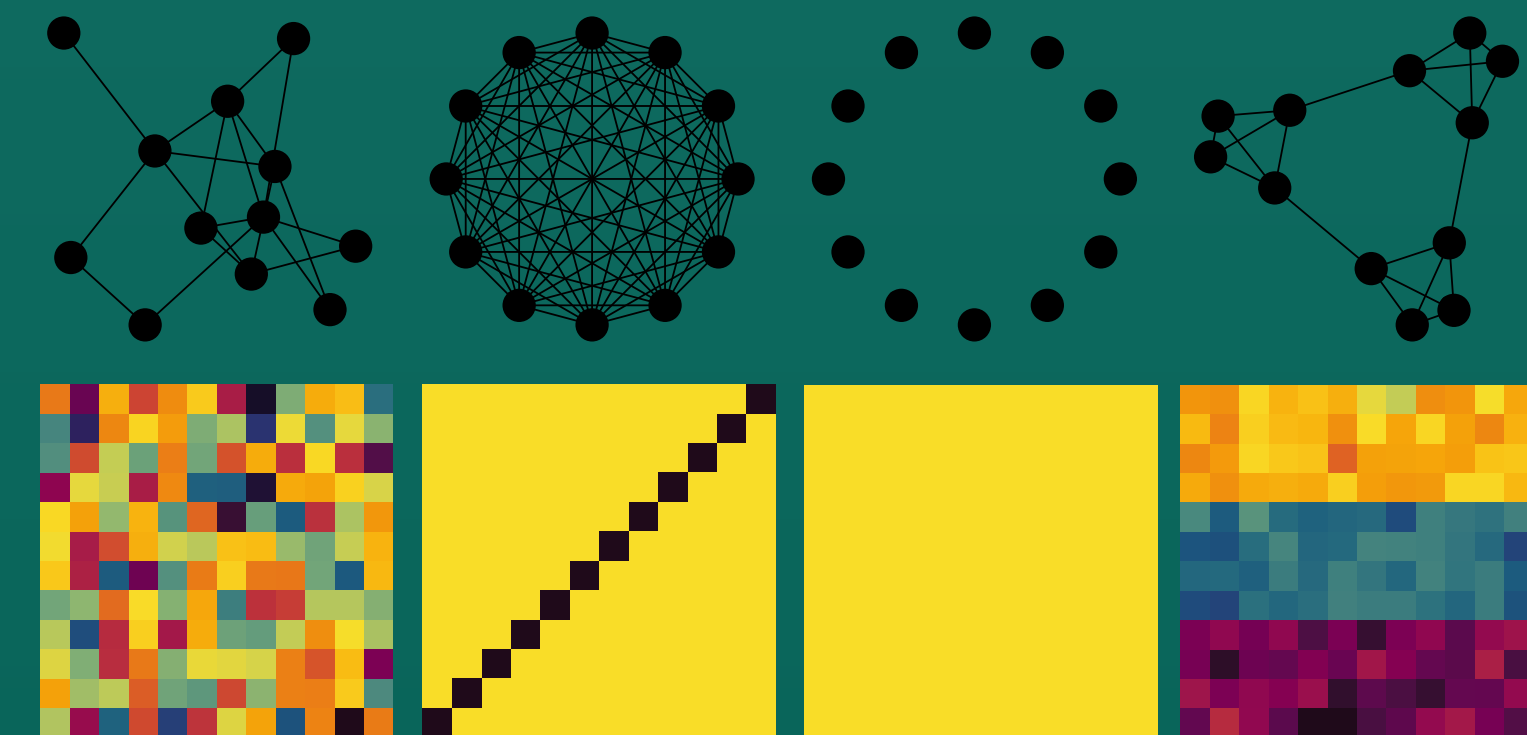
Modes should offer distinct and
interesting views of the data.



**Mode Complementarity
& Mode Diversity**

Metrics

Perturbation Framework



Dataset Evaluation

Benchmarks	
Keep	✓
Realign	⚠
Deprecate	✗

