



CSE COMPUTER SCIENCE
AND ENGINEERING
UNIVERSITY OF MICHIGAN



Exploratory Analysis of Graph Data by Leveraging Domain Knowledge

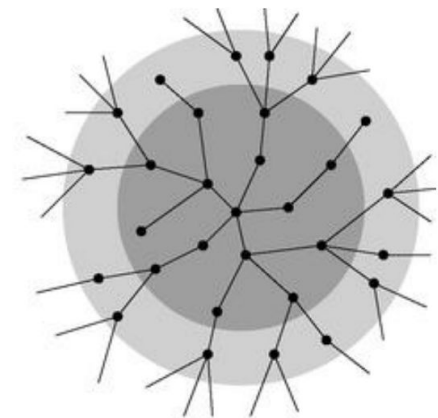
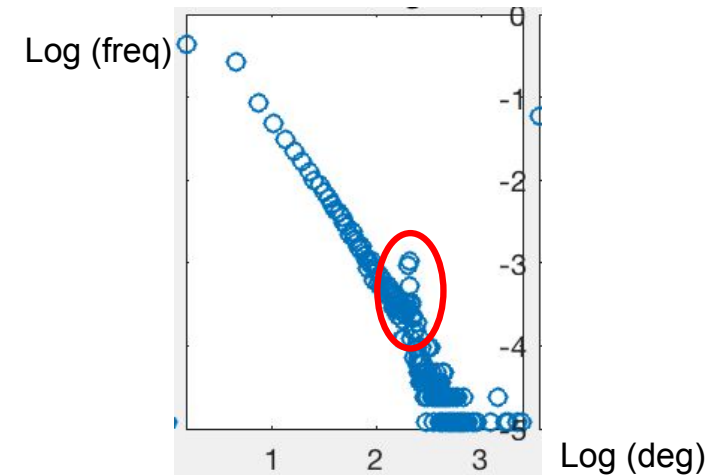
Di Jin

Danai Koutra

IEEE International Conference on Data Mining (ICDM), 2017

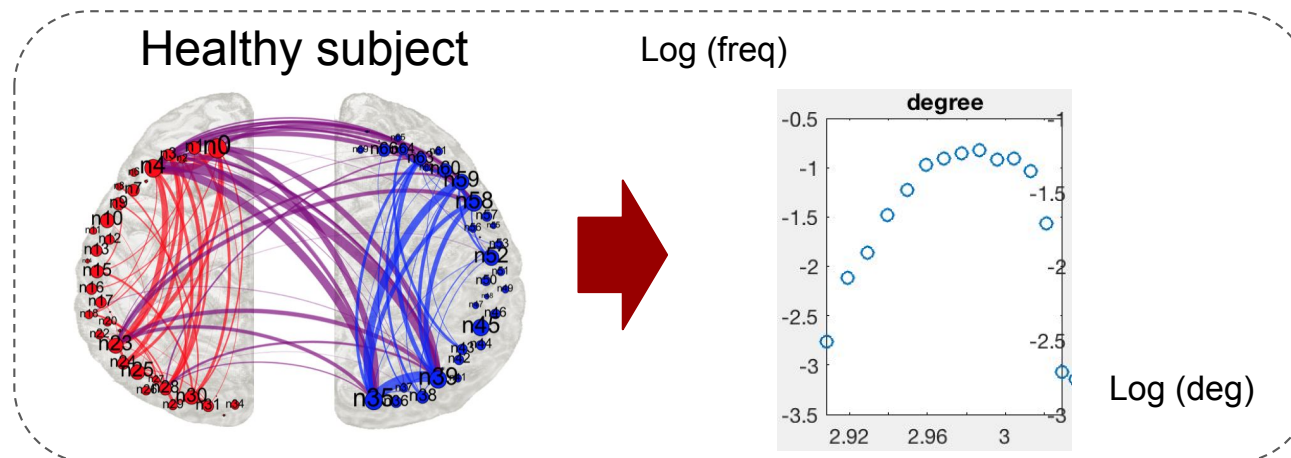
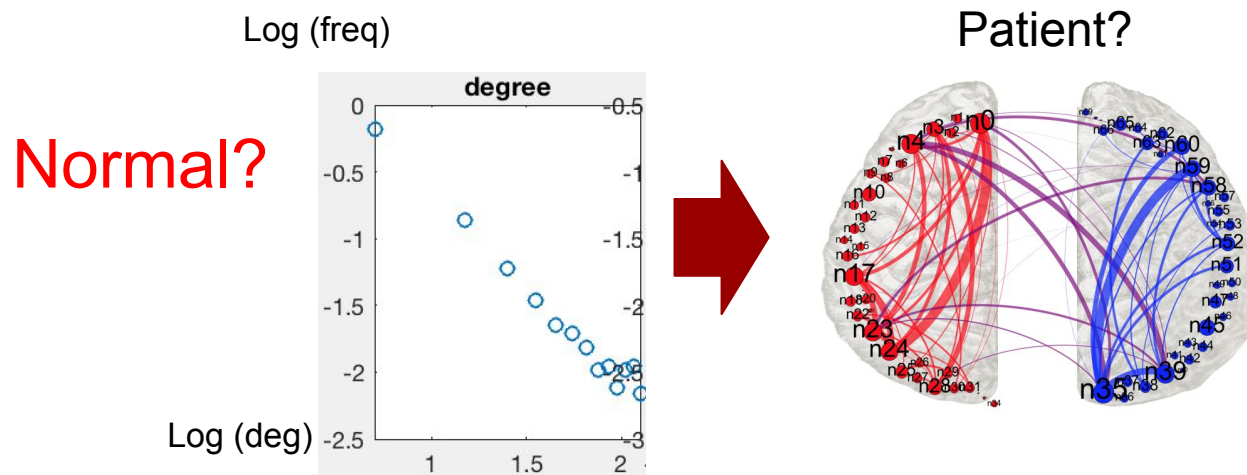
Graph invariants are prevalent

- In many tasks (e.g., anomaly detection, classification, ...)
 - Different invariants
 - Degrees
 - Betweenness
 - Average path length
 - Giant components
 -
 - Compare them with “common” laws
 - The power-like laws
 - 6 degree of separation ($\log(N) / \log(c)$)
 - 1 giant component
 -



Are graph invariants enough to understand?

- The brain connectivity correlation graph



Are graph invariants enough to understand?

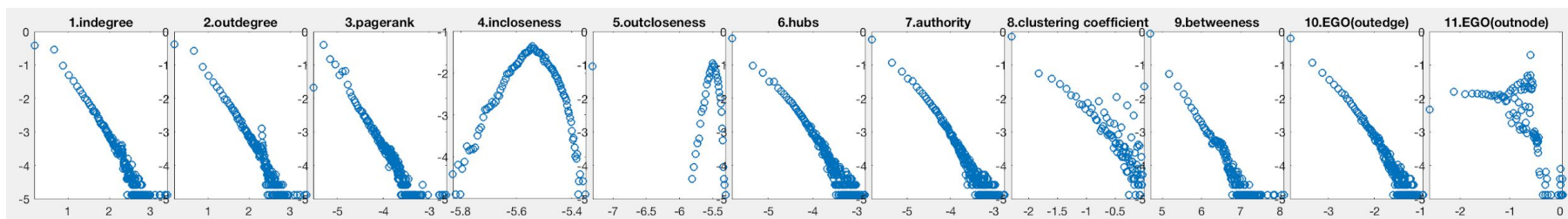
- Common laws are not golden
 - Graph invariant distributions are good.
 - But bigger picture should not be neglected.

Are graph invariants enough to understand?

- Common laws are not golden
 - Graph invariant distributions are good.
 - But bigger picture should not be neglected.
- **Prior/Domain knowledge is important**
 - Graphs are everywhere, but the domain experts are NOT.
 - “What patterns are expected?”

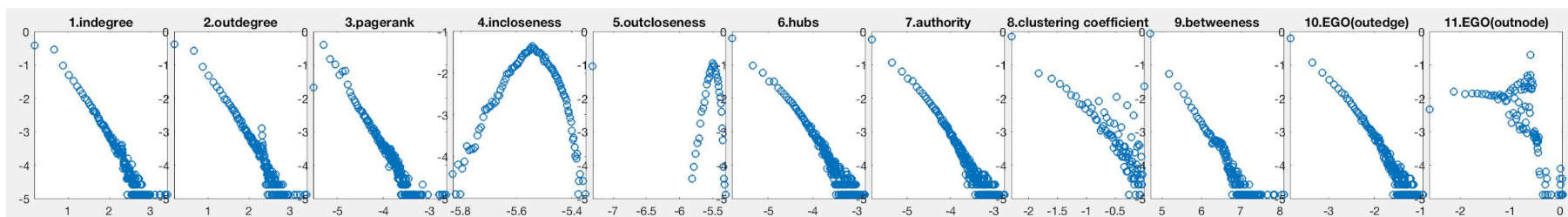
Are graph invariants enough to understand?

- Common laws are not golden
 - Graph invariant distributions are good.
 - But bigger picture should not be neglected.
- Prior/Domain knowledge is important
 - Graphs are everywhere, but the domain experts are NOT.
 - “What patterns are expected?”
- Useful graph invariant distributions (**features**) vary
 - Tons of features can be extracted, few are useful
 - “Which features to explore?”



Are graph invariants enough to understand?

- Common laws are not golden
 - Graph invariant distributions are good.
 - But bigger picture should not be neglected.
- Prior/Domain knowledge is important
 - Graphs are everywhere, but the domain experts are NOT.
 - **“What patterns are expected?”**
- Useful graph invariant distributions (**features**) vary
 - Tons of features can be extracted, few are useful
 - **“Which features to explore?”**



EAGLE: Exploratory Analysis of Graphs with domain knowLEdge

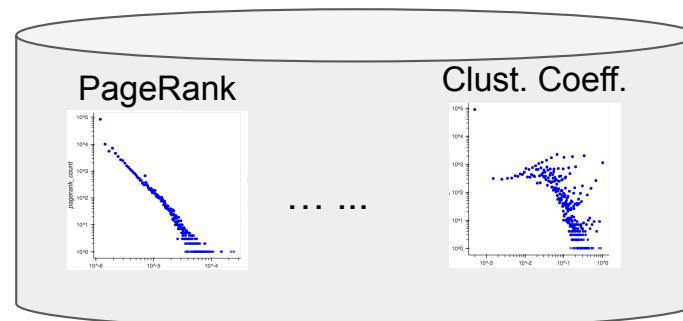
Given: an input graph & *domain knowledge*

Find: brief summary consisting of **representative** features that satisfies a set of desired properties (e.g., diversity)

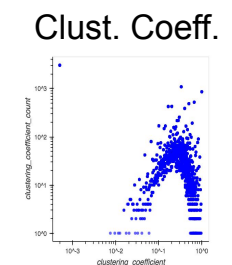
- **“What patterns are expected?”**
 - *Domain knowledge:* a collection of graphs with all features in the feature space.
- **“Which features to explore?”**
 - *Representative features:* graph invariant distributions (PDF) with desired properties.



Unknown graph

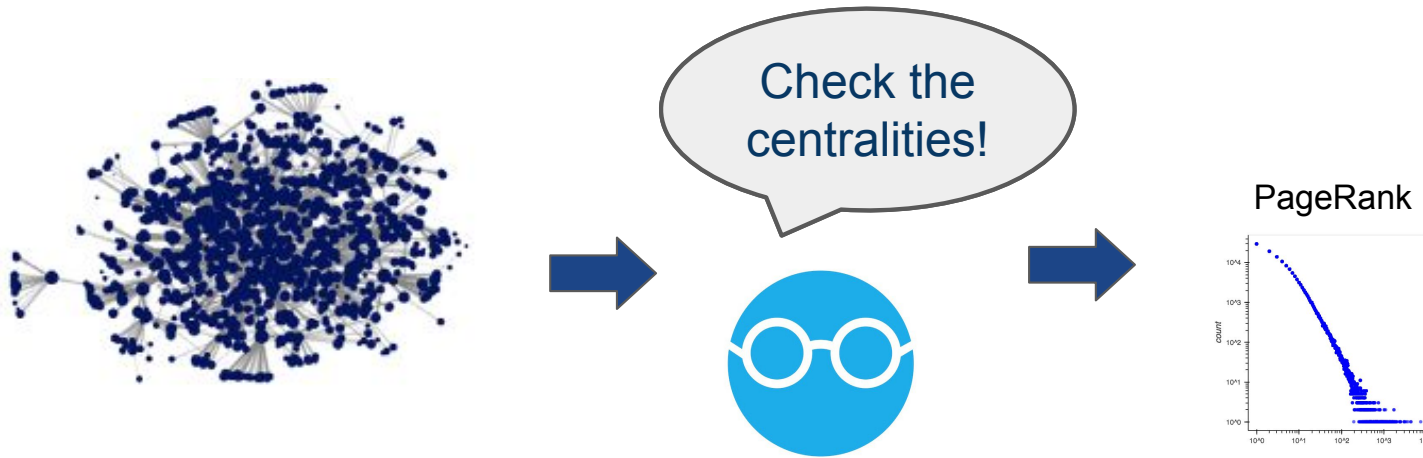


Domain knowledge

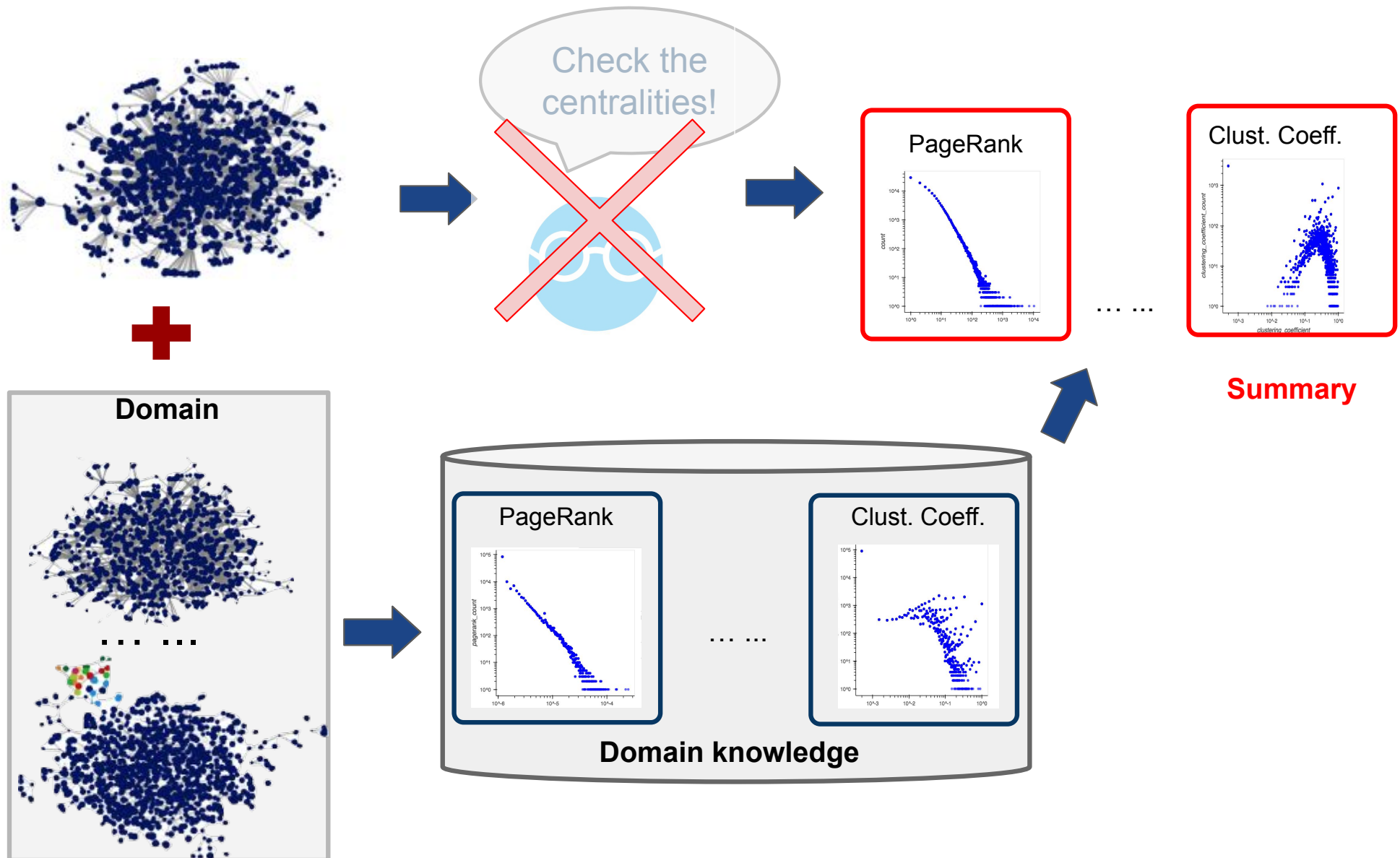


Representative
graph invariants

Many Existing Methods

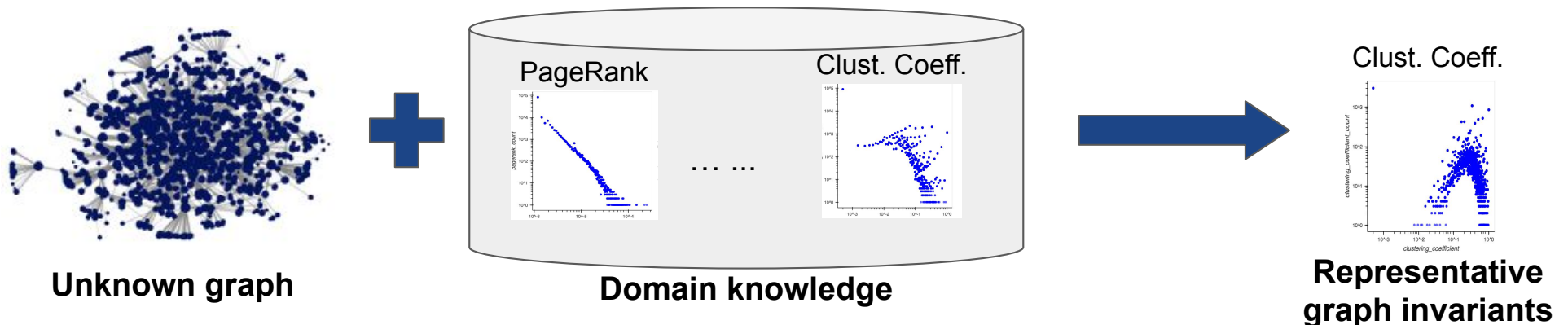


Proposed Solution: key idea



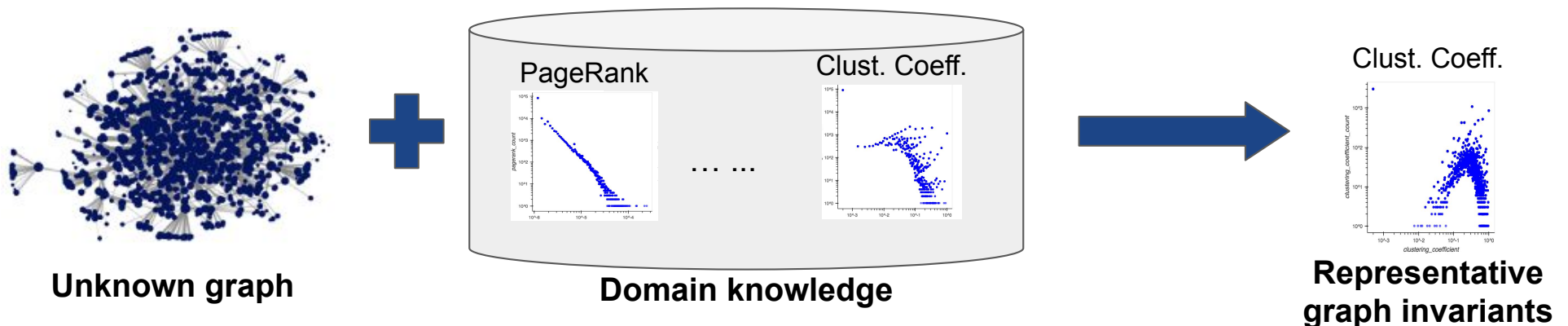
Proposed Solution: key idea

- “Summarize an unknown graph from known ones”.
 - “Known graphs”: the *domain knowledge*.
 - Summarize through *representative* graph invariants.
 - Discovers domain-specific patterns **automatically**.



Proposed Solution: key idea

- **“Summarize an unknown graph from known ones”**.
 - “Known graphs”: the *domain knowledge*.
 - Summarize through *representative* graph invariants.
 - Discovers domain-specific patterns **automatically**.
- Not a traditional graph summarization problem.
 - No compressed representation of an input graph.



EAGLE: Desired properties

The **summary** should be:

- Diverse
- Concise
- Domain-specific



The **process** should be:

- Efficient
- Interpretable: univariate graph statistics (PDF)

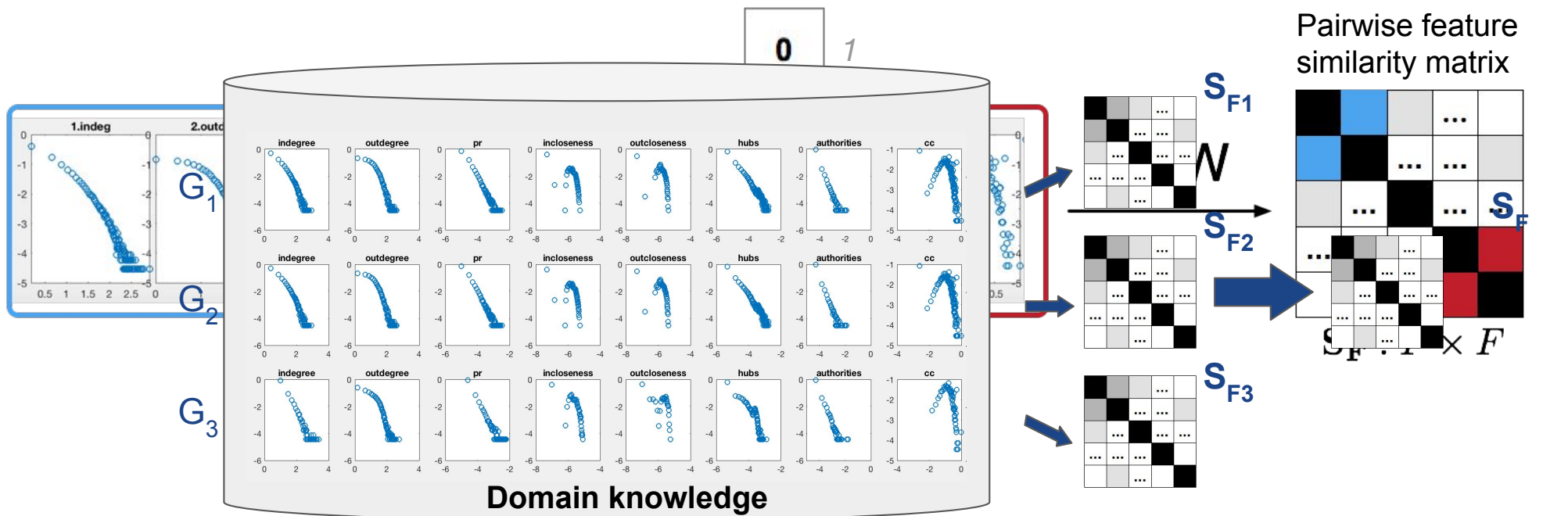
Input graph
 g



selected features (surprising)

EAGLE: Formulation

$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$



$$\mathbf{S}_F(f_j, f_l) = \sum_{i=1}^K w_i \cdot \mathbf{S}_{Fi}(f_j, f_l)$$

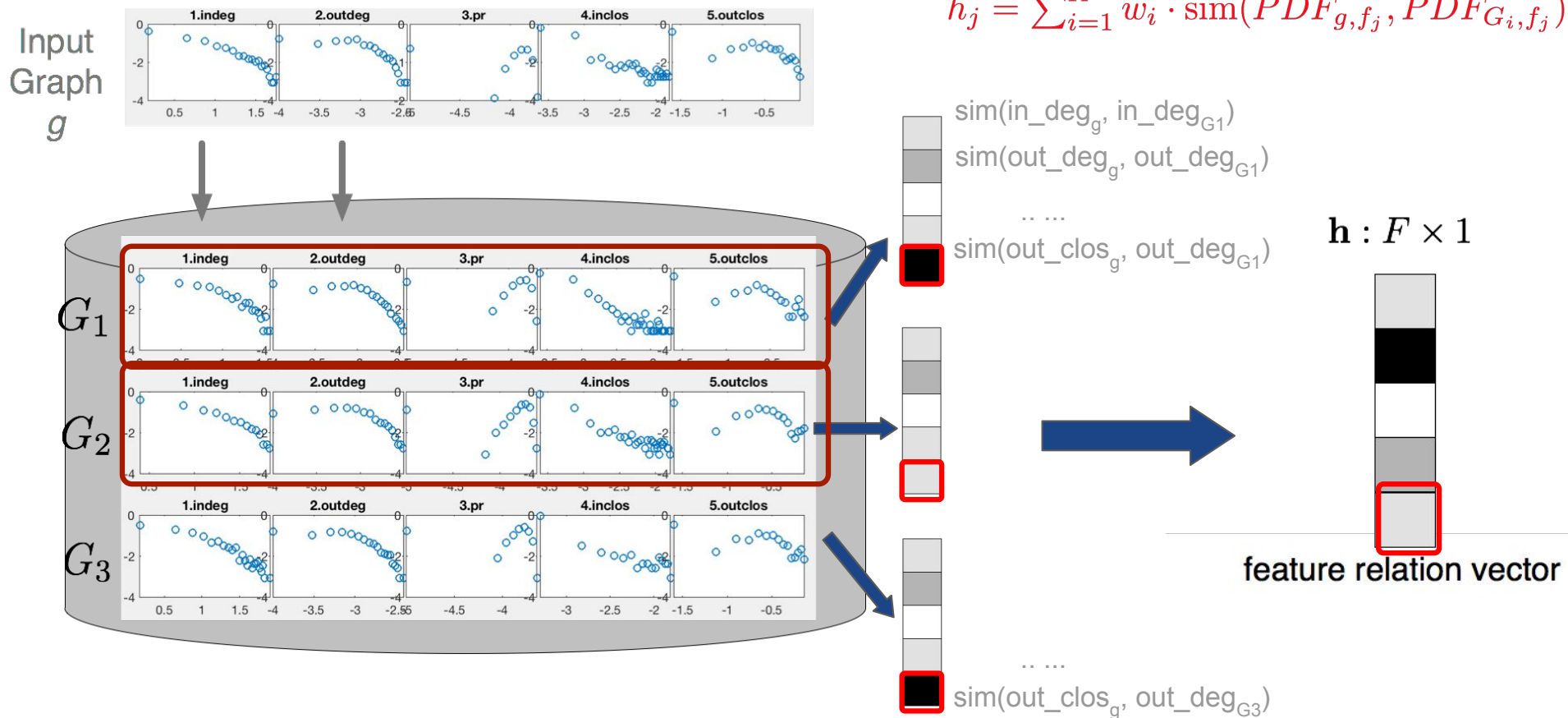
EAGLE: Formulation

$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$

$$\|\mathbf{f}\|_0 \xrightarrow{\text{Relax}} \|\mathbf{f}\|_2 = \mathbf{f}^T \mathbf{f}$$

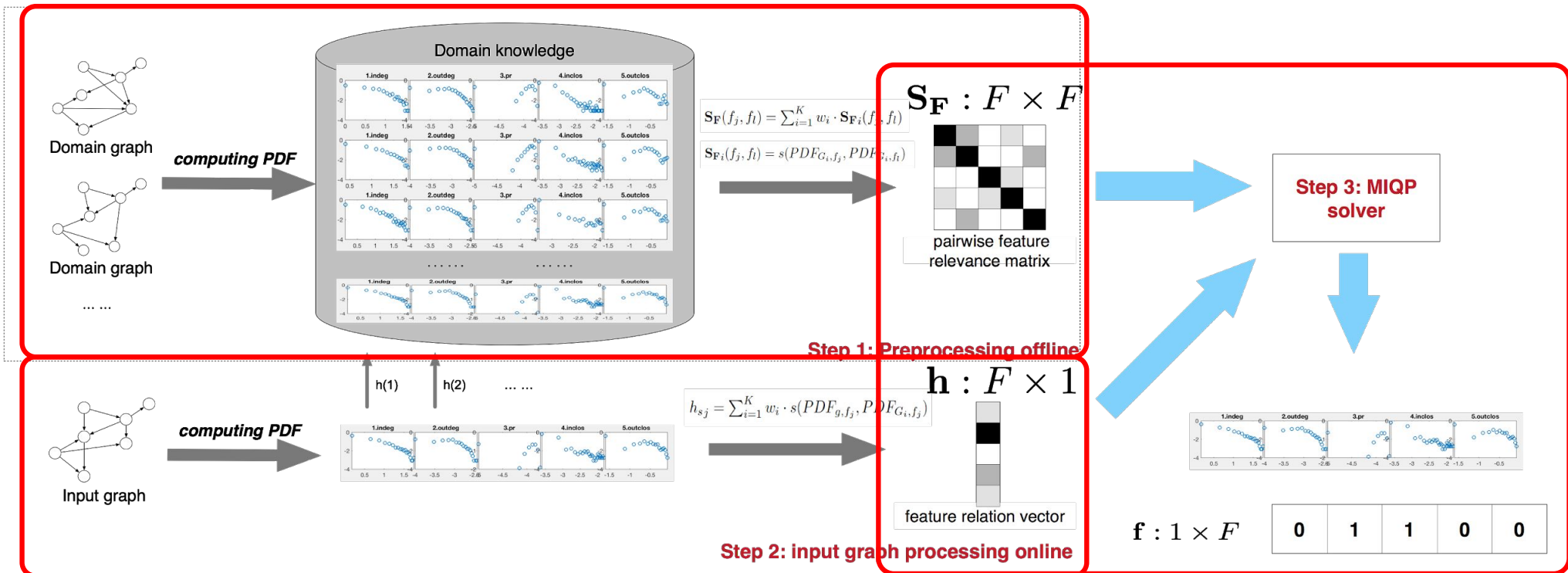
EAGLE: Formulation

$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$



Workflow: 3 steps

$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_{\mathbf{F}} \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$



How to solve? MIQP solver

- Before

$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$

- Rewriting relaxed form

$$\begin{aligned} \arg \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \lambda_1 \mathbf{f}^T \mathbf{S}_F \mathbf{f} + \lambda_2 \mathbf{f}^T \mathbf{I}_F \mathbf{f} + \lambda_3 \mathbf{f}^T \mathbf{h} \\ = \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \mathbf{f}^T \underbrace{(\lambda_1 \mathbf{S}_F + \lambda_2 \mathbf{I}_F)}_{\mathbf{Q}} \mathbf{f} + \mathbf{f}^T \underbrace{\lambda_3 \mathbf{h}}_{\mathbf{r}} \end{aligned}$$

How to solve? MIQP solver

- Before
$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$

- Rewriting relaxed form
$$\begin{aligned} \arg \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \lambda_1 \mathbf{f}^T \mathbf{S}_F \mathbf{f} + \lambda_2 \mathbf{f}^T \mathbf{I}_F \mathbf{f} + \lambda_3 \mathbf{f}^T \mathbf{h} \\ = \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \mathbf{f}^T \underbrace{(\lambda_1 \mathbf{S}_F + \lambda_2 \mathbf{I}_F)}_{\mathbf{Q}} \mathbf{f} + \mathbf{f}^T \underbrace{\lambda_3 \mathbf{h}}_{\mathbf{r}} \end{aligned}$$

- The general form: MIQP & MILP with slack variable z

$$\text{minimize}_{\mathbf{f}} \quad \mathbf{f}^T \mathbf{Q} \mathbf{f} + \mathbf{r}^T \mathbf{f}$$

$$\begin{aligned} \text{subject to} \quad & 0 \leq \sum_i^F \mathbf{f}(i) \leq F \\ & 0 \leq \mathbf{f}(i) \leq 1, \quad i = 1, \dots, F. \end{aligned}$$

Slack variable z

$$\text{minimize}_{\mathbf{f}, z} \quad z + \mathbf{r}^T \mathbf{f}$$

$$\begin{aligned} \text{subject to} \quad & 0 \leq \sum_i^F \mathbf{f}(i) \leq F \\ & 0 \leq \mathbf{f}(i) \leq 1, \quad i = 1, \dots, F. \\ & \mathbf{f}^T \mathbf{Q} \mathbf{f} - z \leq 0, \quad z \geq 0 \end{aligned}$$

“0-pit” Problem

$$\begin{aligned} & \arg \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \lambda_1 \mathbf{f}^T \mathbf{S}_F \mathbf{f} + \lambda_2 \mathbf{f}^T \mathbf{I}_F \mathbf{f} + \lambda_3 \mathbf{f}^T \mathbf{h} \\ & = \min_{\mathbf{f} \in \{0,1\}^{F \times 1}} \mathbf{f}^T \underbrace{(\lambda_1 \mathbf{S}_F + \lambda_2 \mathbf{I}_F)}_Q \mathbf{f} + \mathbf{f}^T \underbrace{\lambda_3 \mathbf{h}}_r \end{aligned}$$

- “0-pit” problem:
 - All the terms are positive
 - Optimal solution: for \mathbf{f} = all-0 vector
- Solution:
 - **EAGLE-Fix**: Explicitly set the # of selected features in \mathbf{f}
 - **EAGLE-Flex**: Set negative value to the normalization term

Experiments: data

○ Feature space (28 in total):

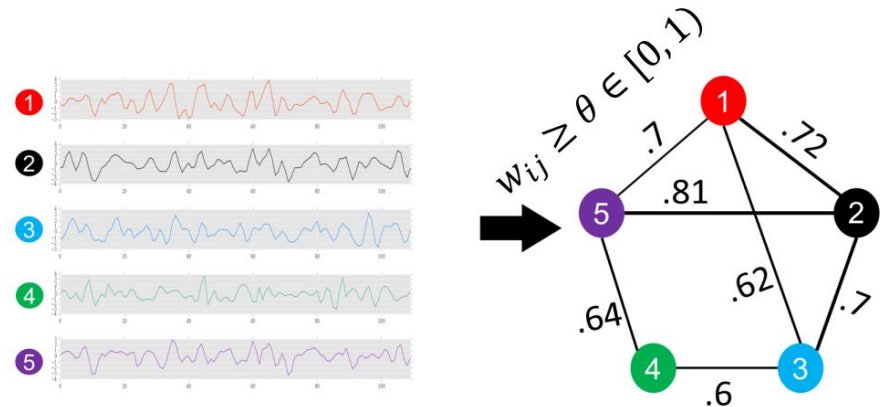
■ Node-specific:

in/out-degree, PageRank, hubs, authorities, roles, ...

■ Structure-specific:

in/out neighbors and # in/out edges of egonets, ..., distribution of communities, motifs, ...

Domain	Name	Nodes	Edges	Description
Connectomics	Brain-Voxel1	3 789	399 069	directed unweighted
	Brain-Voxel2	3 789	148 648	directed unweighted
Citation networks	HepTh	27 770	352 807	directed unweighted
	HepPh	34 546	421 578	directed unweighted
Social science	Epinions	75 879	508 837	directed unweighted
	Slashdot0811	77 360	905 468	directed unweighted
	Slashdot0922	82 168	948 464	directed unweighted



Experiments: baseline methods

- Random selection
- Surprising selection

No existing domain-specific summarization method!

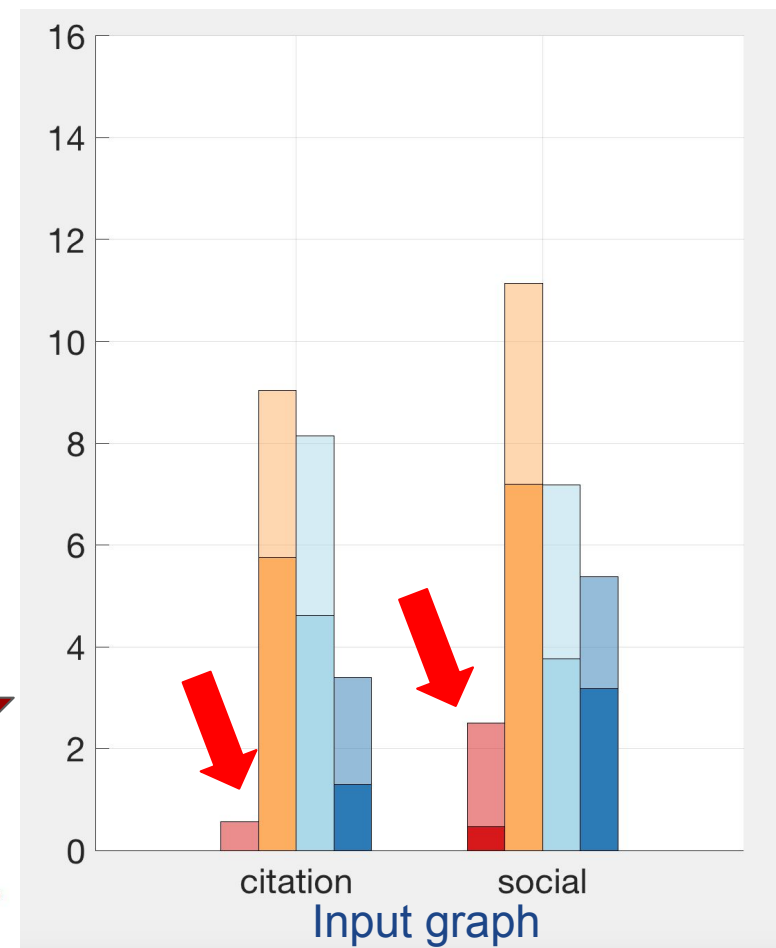
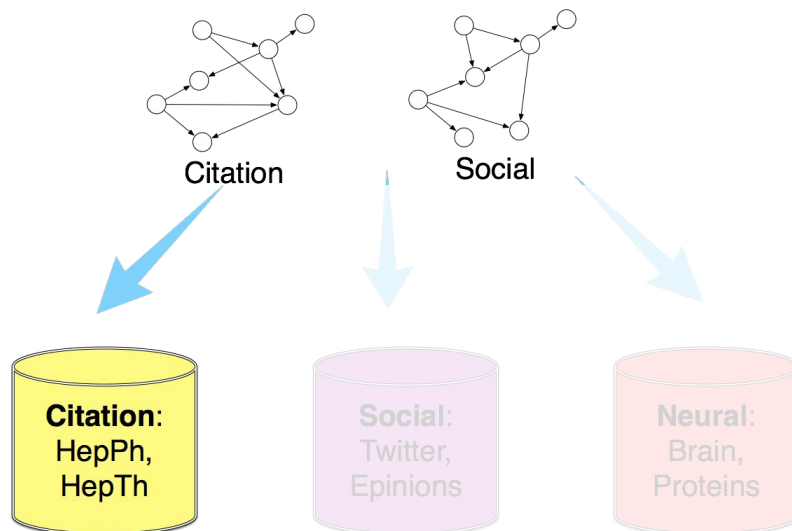
$$\arg \min_{\mathbf{f}} \lambda_1 \underbrace{\mathbf{f}^T \mathbf{S}_F \mathbf{f}}_{\text{Diversity}} + \lambda_2 \underbrace{\|\mathbf{f}\|_0}_{\text{Conciseness}} + \lambda_3 \cdot \underbrace{\mathbf{f}^T \mathbf{h}}_{\text{Domain-specificity}}$$

- SCAGNOSTICS
 - Pick each feature *independently* based on its anomalies in density, shape and trend.
 - 9 scores: stringiness, skewness, skinniness, etc.

Diversity & Domain-specificity

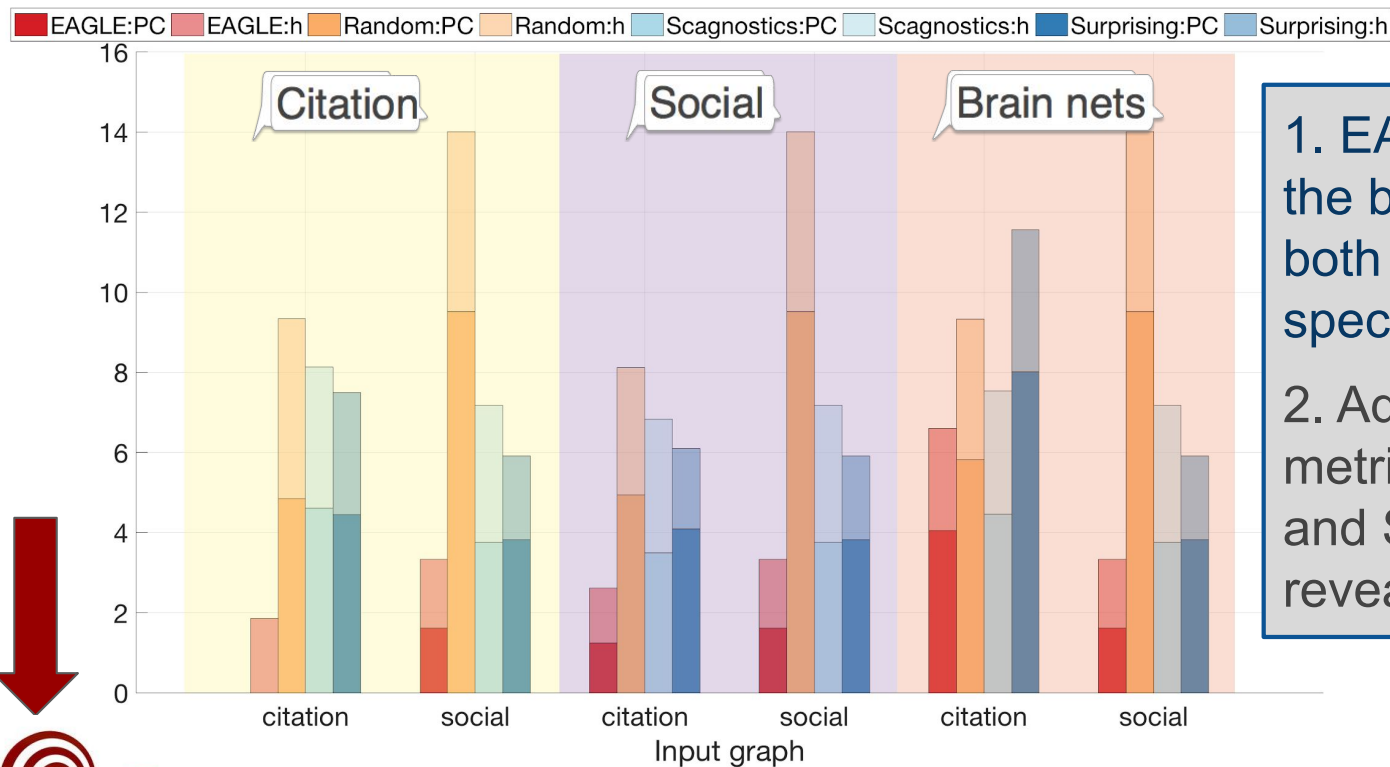
- Metric: Pearson correlation (PC)

■ EAGLE:PC
 ■ EAGLE:h
 ■ Random:PC
 ■ Random:h
 ■ Scagnostics:PC
 ■ Scagnostics:h
 ■ Surprising:PC
 ■ Surprising:h



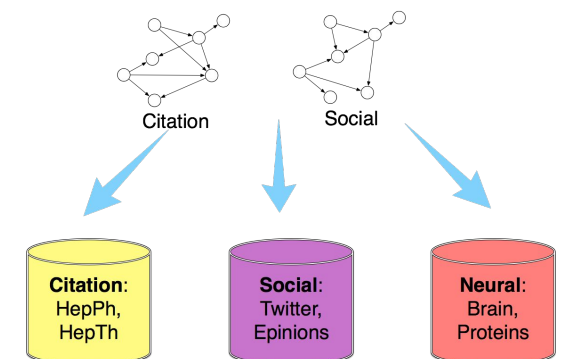
Diversity & Domain-specificity

- Metric: Pearson correlation (PC)



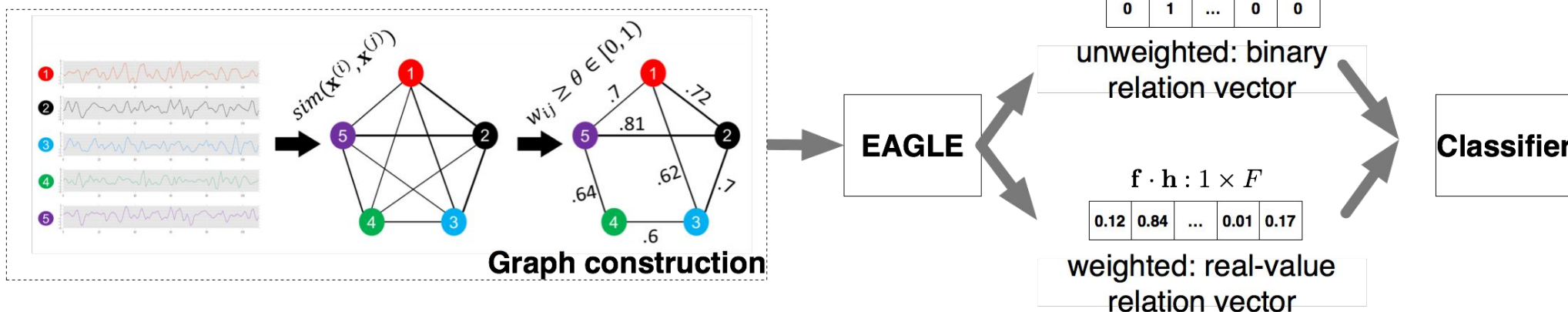
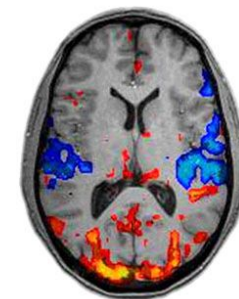
1. EAGLE outperforms the baselines in terms of both diversity and domain specificity.

2. Additional evaluation metrics (Kendall's Tau and Spearman's Rank) reveal similar patterns.



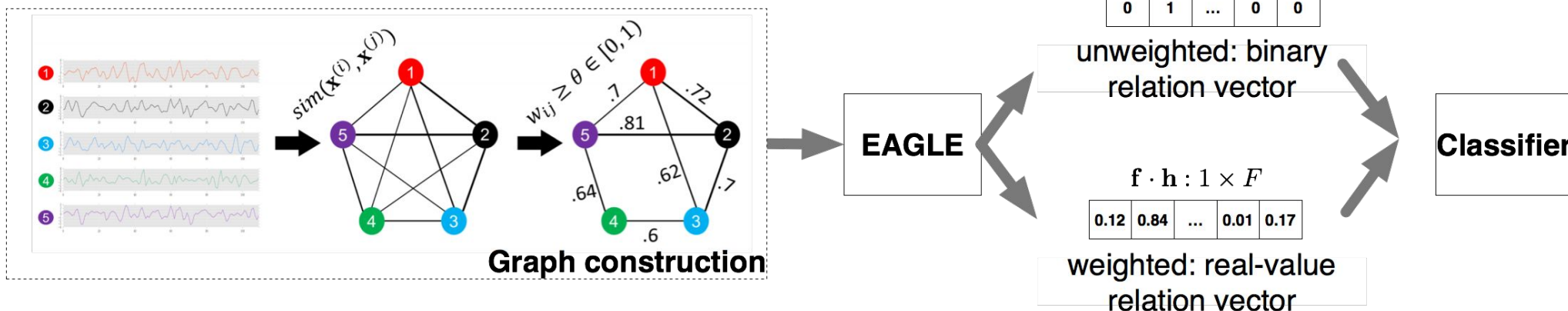
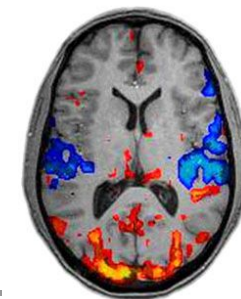
Case study: brain graph classification

- Setup (EAGLE-Fix & -Flex)
 - COBRE dataset: 72 patients with schizophrenia and 76 healthy controls, 1166 fMRI time series.
 - Threshold: 0.6
 - Feature space: 11 (degree, clustering coeff, betweenness, ...)



Case study: brain graph classification

- Setup (EAGLE-Fix & -Flex)
 - COBRE dataset: 72 patients with schizophrenia and 76 healthy controls, 1166 fMRI time series.
 - Threshold: 0.6
 - Feature space: 11 (degree, clustering coeff, betweenness, ...)
- Baselines
 - Baseline 1: average feature values
 - Baseline 2: “flatten” adjacency matrix



Case study: brain graph classification

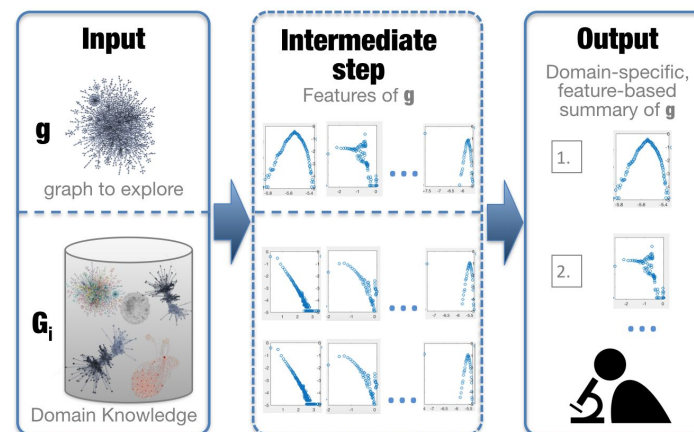
Method Category	Unweighted		Weighted		
	Ordinary	Surprising	Ordinary	Surprising	
EAGLE-FLEX	0.6893	0.5499	0.7096	0.7296	←
EAGLE-FIX: 6	0.5114	0.5445	0.6961	0.7371	Classification on COBRE: AUC scores per method
EAGLE-FIX: 8	0.6795	0.5904	0.7216	0.7079	
EAGLE-FIX: 10	0.5003	0.4989	0.7032	0.6807	
Full	-	-	0.6681	0.7147	←
Baselines	Baseline 1:	0.7028	Baseline 2:	0.1099	←

Although not designed explicitly for this, features selected by EAGLE can be applied to specific tasks such as classification with at least as good performance.

EAGLE-Flex improves performance by effectively eliminating noise from the data.

Conclusion & Contributions

- EAGLE: a **novel graph summarization** technique that *learns an unknown graph from known ones*.
- Informative graph features that satisfy:
 - Diversity
 - Conciseness
 - Domain-specificity
 - Interpretability
 - Efficiency
- Formulation of graph exploration as **constrained optimization**
 - Two efficient solutions: Eagle-Fix and Eagle-Flex.
 - Applications.



Thank you! Questions?

- EAGLE: a **novel graph summarization** technique that *learns an unknown graph from known ones*.
- Informative graph features that satisfy:
 - Diversity
 - Conciseness
 - Domain-specificity
 - Interpretability
 - Efficiency
- Formulation of graph exploration as **constrained optimization**
 - Two efficient solutions: Eagle-Fix and Eagle-Flex.
 - Applications.

