



Comparing the Decompositions Produced by Software Clustering Algorithms using Similarity Measurements

*2001 IEEE International Conference on
Software Maintenance (ICSM'01).*

Brian S. Mitchell & Spiros Mancoridis
Math & Computer Science, Drexel University

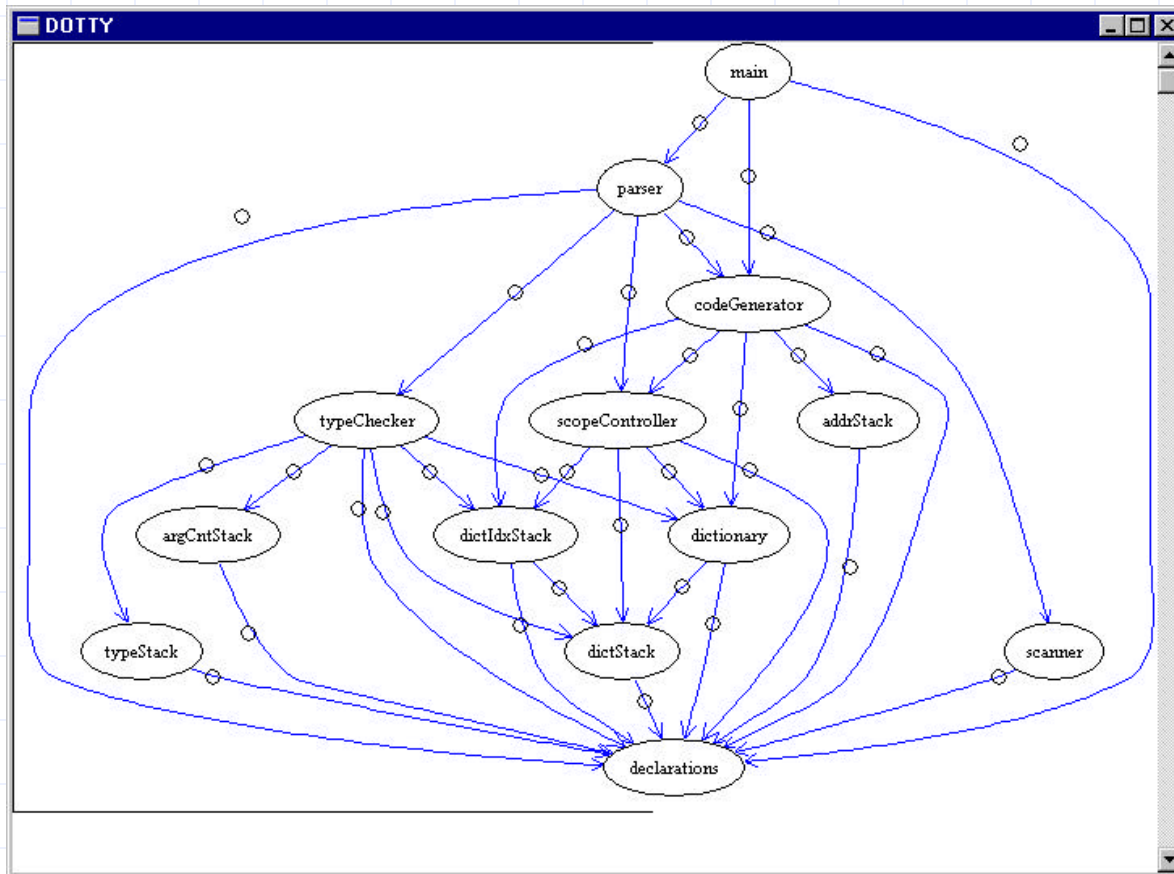
Motivation

*Using module dependencies
when determining the
similarity between two
decompositions
is a good idea...*



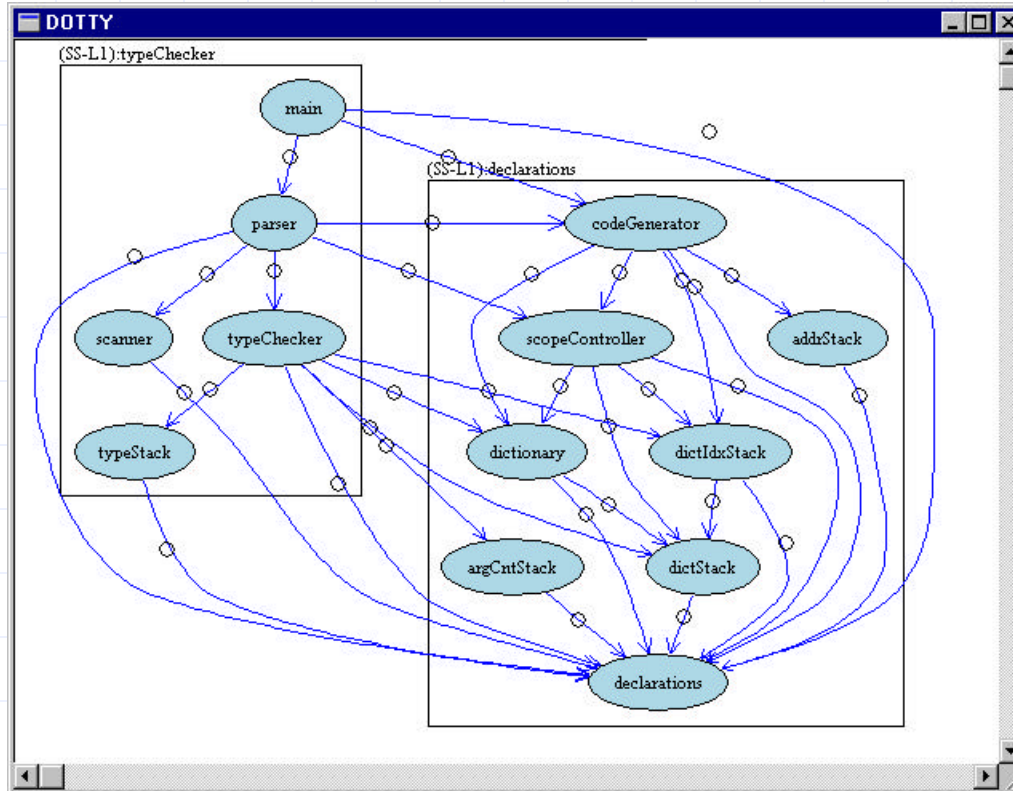
Clustering the Structure of a System (1)

Given the structure of a system...



Clustering the Structure of a System (2)

The goal is to partition the system structure graph into clusters...

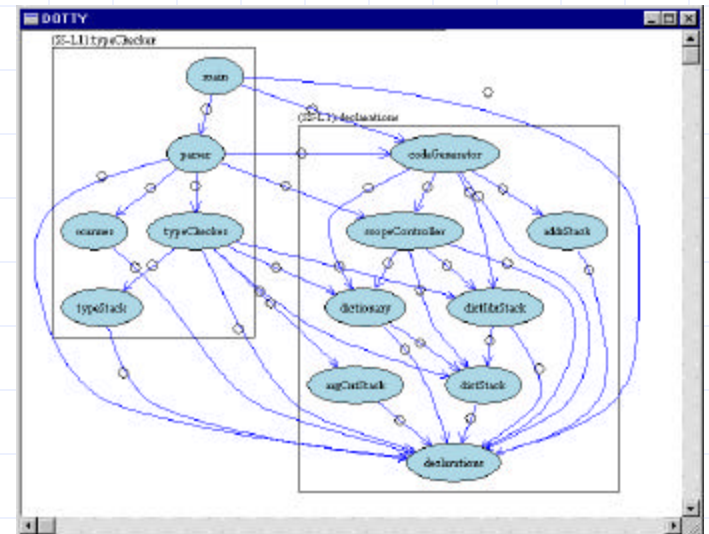
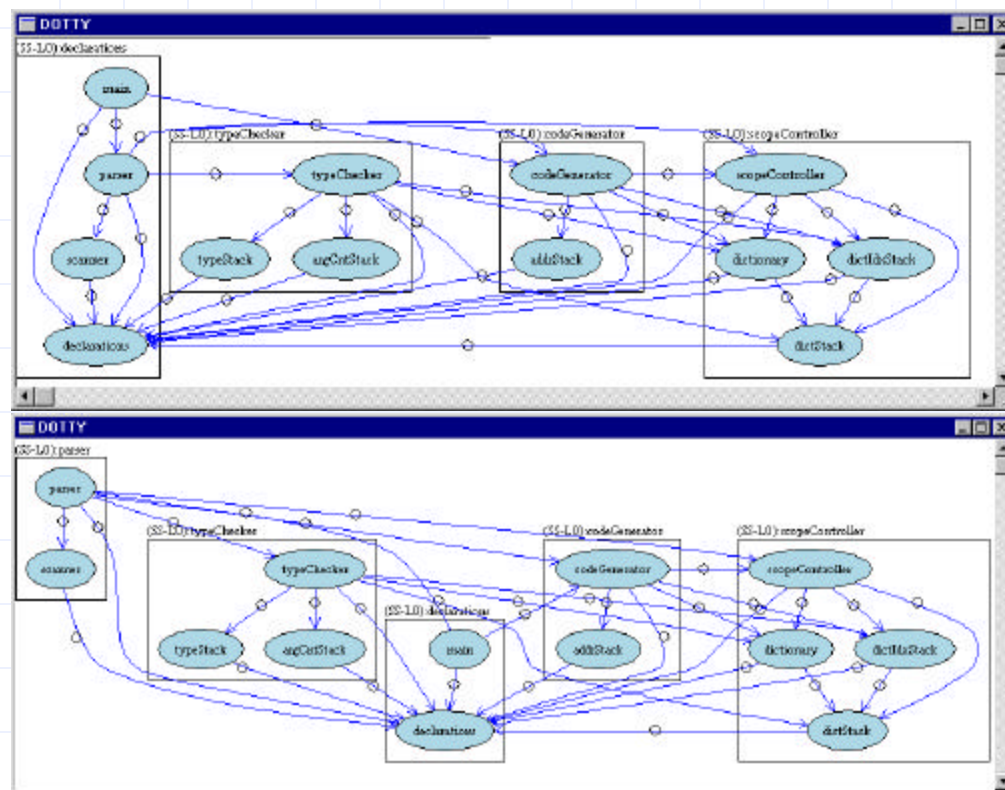


The clusters should represent the subsystems



Clustering the Structure of a System (3)

But how do we know that the clustering result is good?



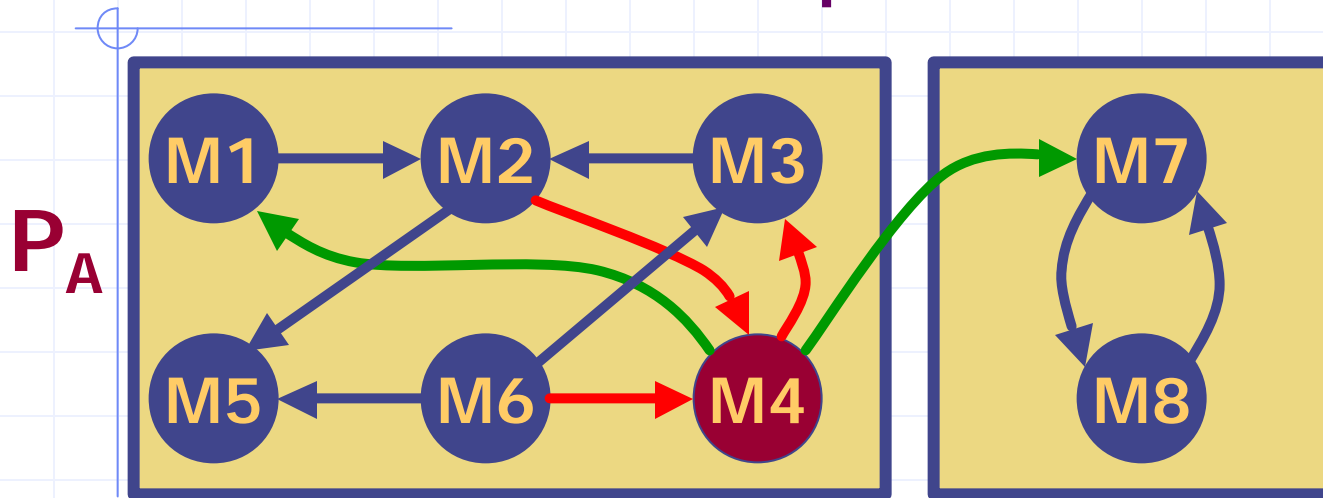
Ways to Evaluate Software Clustering Results...

Given a software clustering result, we can:

- ◆ Assess it against a mental model
- ◆ Assess it against a benchmark standard
- ◆ Techniques:
 - Subjective Opinions
 - Similarity Measurements



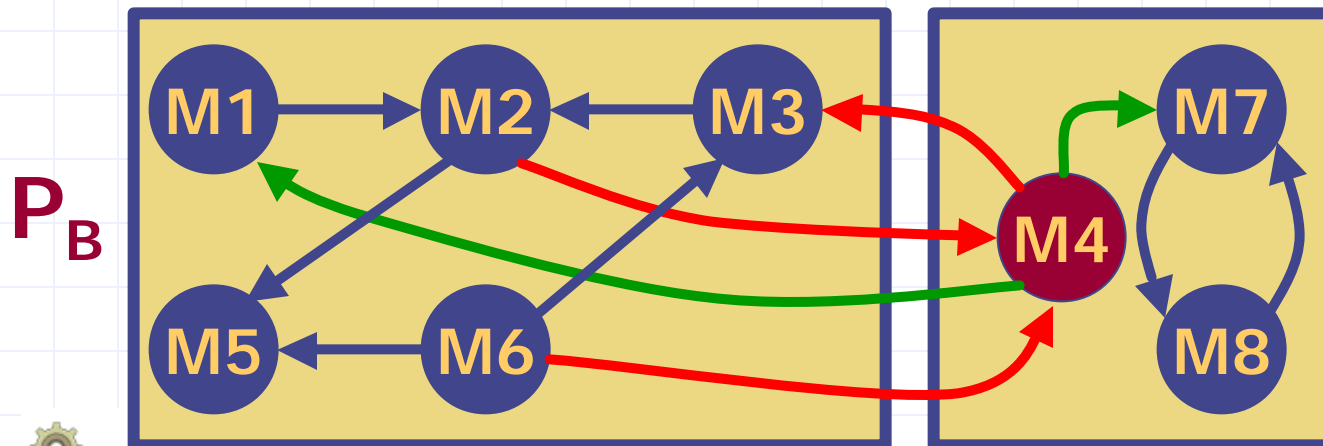
Example: How “Similar” are these Decompositions?



Blue Edges:
Similarity still the same...

Green Edges:
Similarity still the same...

Red Edges:
Not as similar...



Conclusions:
Once we add the **red** edges the similarity between P_A and P_B decreases



Observations

- ◆ Edges are important for determining the similarity between decompositions
- ◆ Existing measurements don't consider edges:
 - Precision / Recall (*similarity*)
 - MoJo (*distance*)
- ◆ **Our idea:** Use the edges to determine similarity

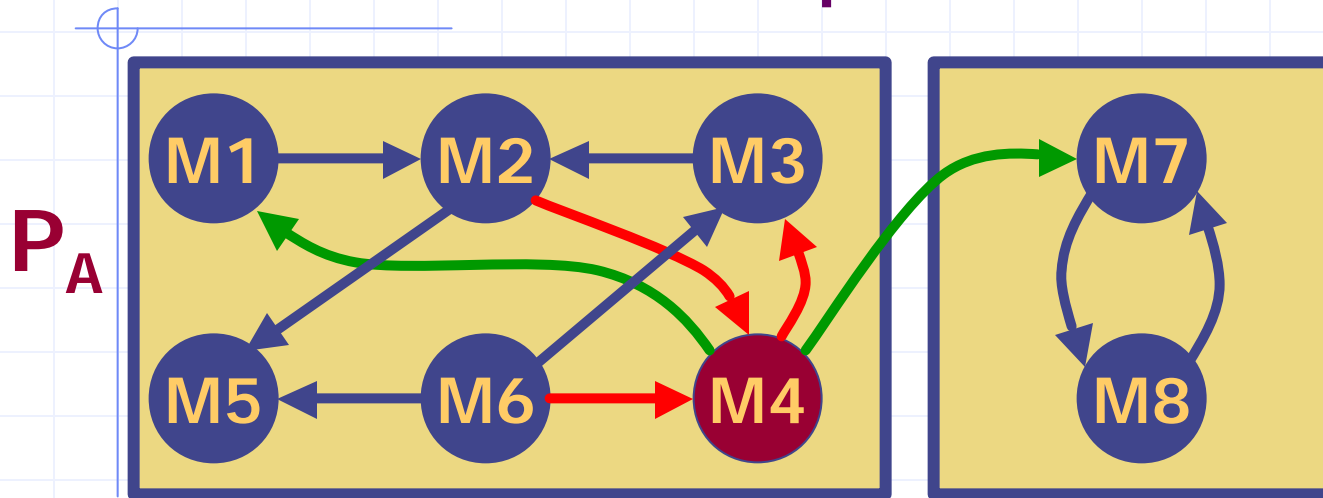


Research Objectives

- ◆ Create new similarity measurements that use dependencies (edges)
 - **EdgeSim** (*similarity*)
 - **MeCI** (*distance*)
- ◆ Evaluate the new similarity measurements against MoJo & Precision/Recall
- ◆ Use similarity measurements to support evaluation of software clustering results (see our WCRE'01 paper)

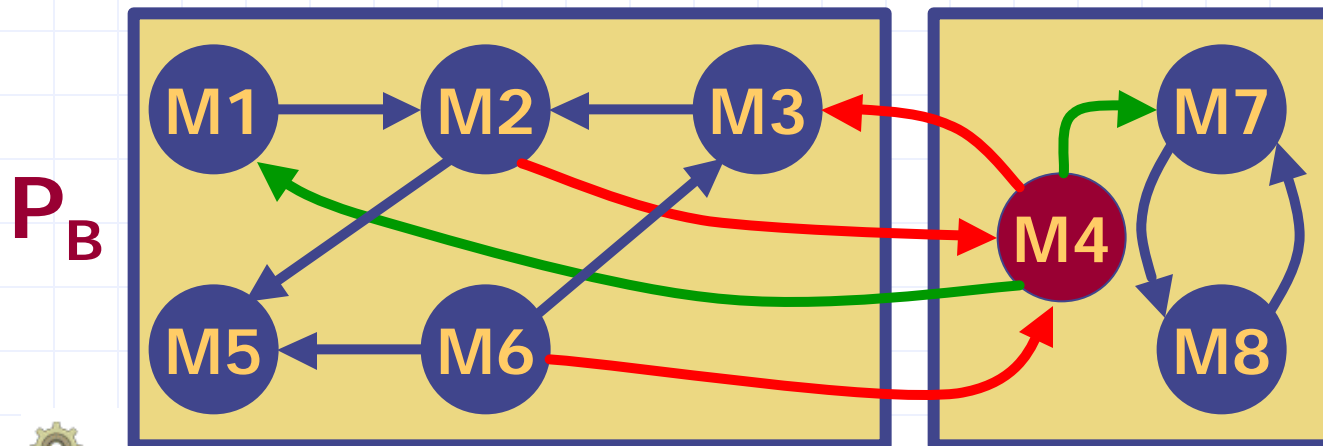


Example: How “Similar” are these Decompositions?



Add
Blue Edges:
PR, MoJo, MeCI &
EdgeSim **unchanged**.

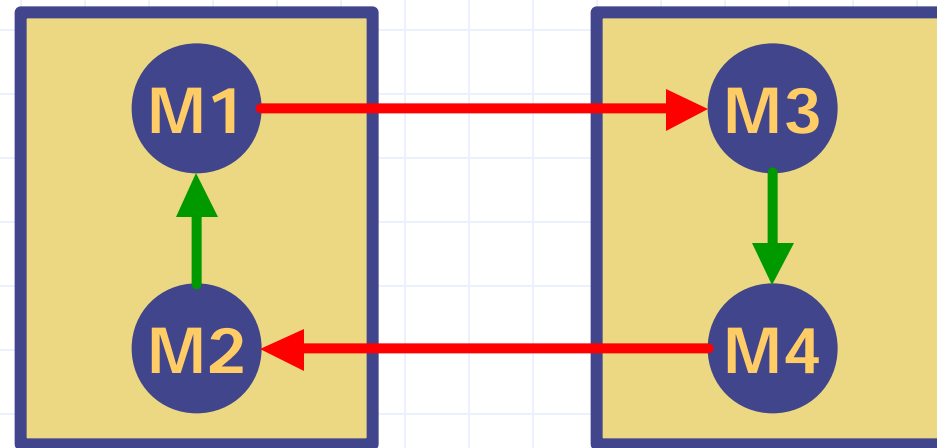
Add
Green Edges:
PR, MoJo, MeCI &
EdgeSim **unchanged**.



Add
Red Edges:
PR, MoJo **unchanged**.
EdgeSim, MeCI
reduced.



Definitions



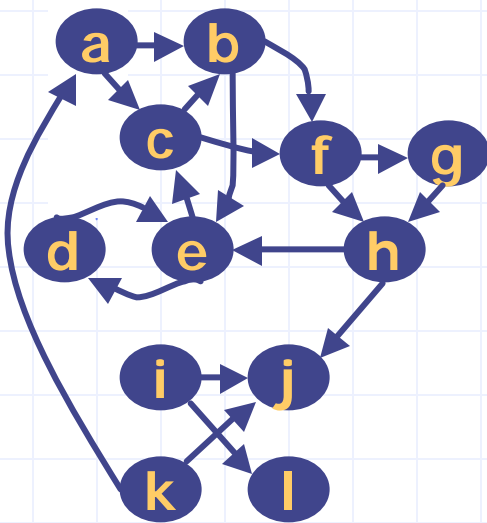
Internal/Intra-Edge: Edge within a cluster

External/Inter-Edge: Edge between two clusters

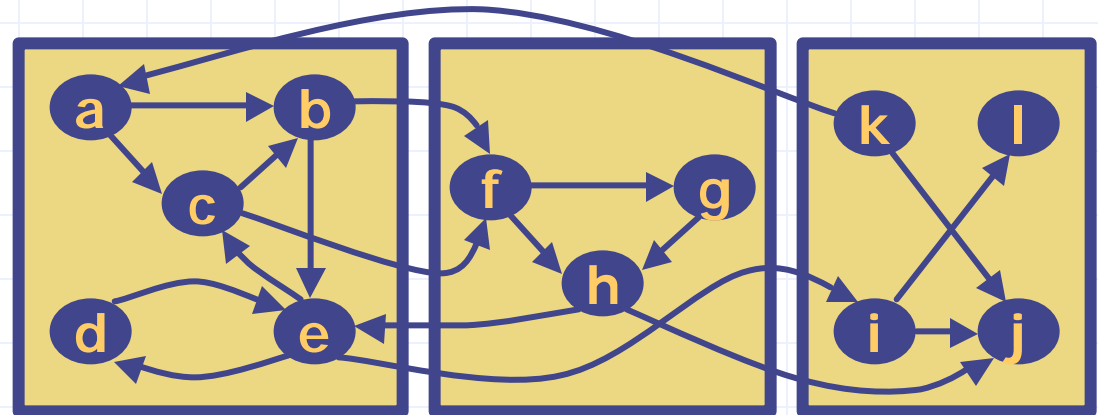


EdgeSim Example

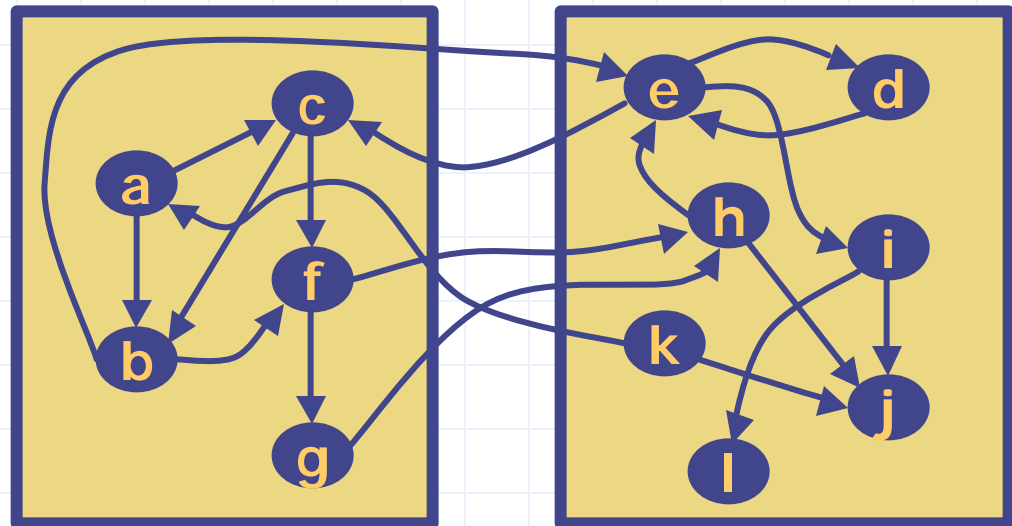
MDG



P_A



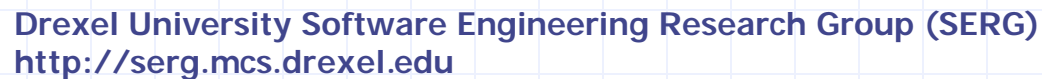
P_B



MDG

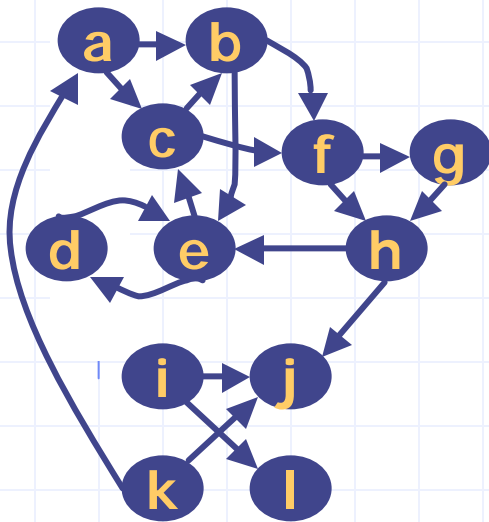


Find Common Inter- and Intra-Edges



EdgeSim Example

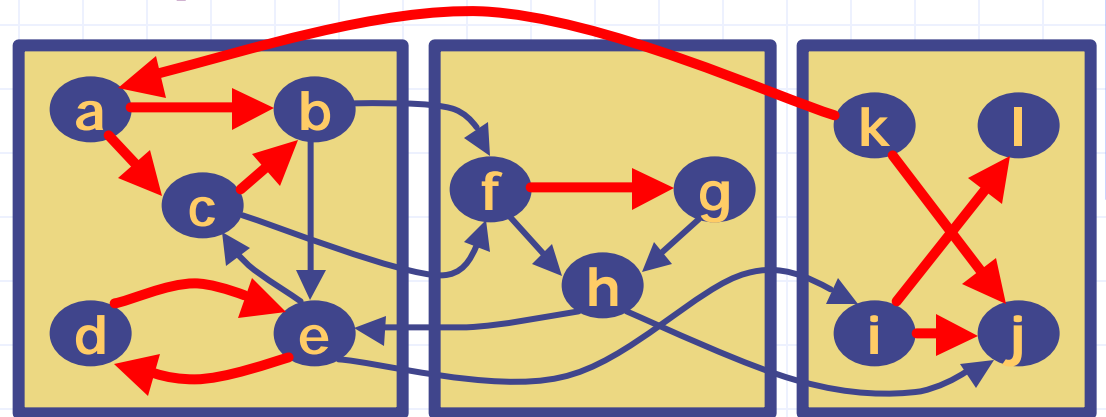
MDG



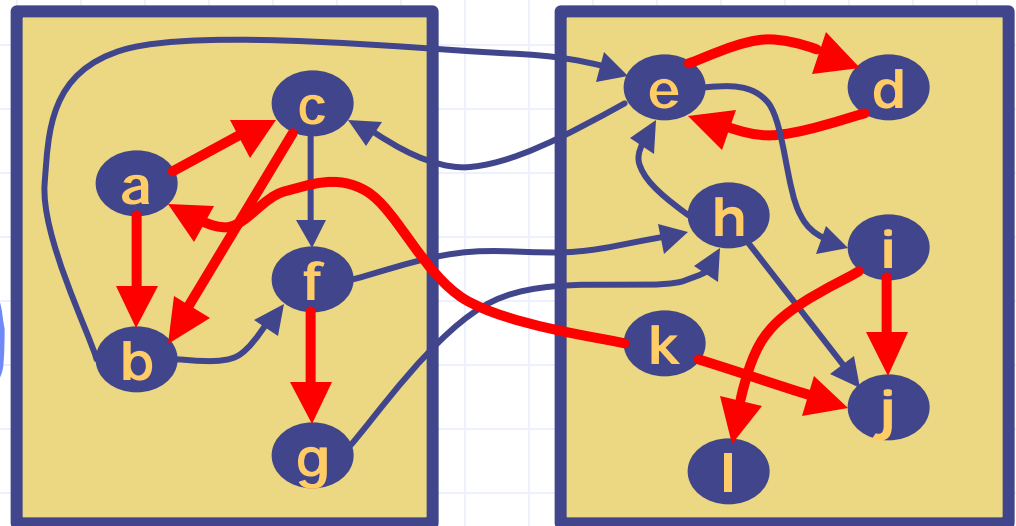
$$\frac{\text{Common Edge Weight}}{\text{Total Edge Weight}} = \frac{10}{19} = 53\%$$



P_A

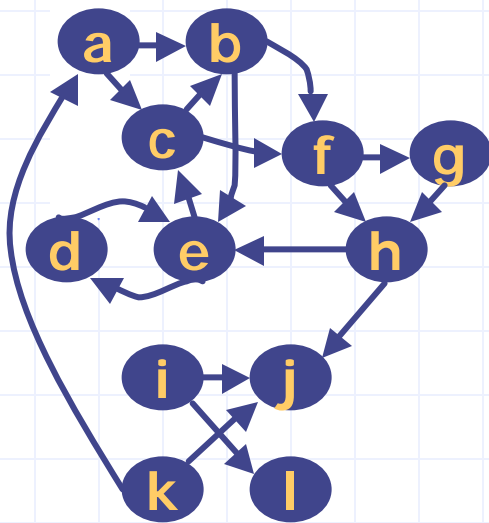


P_B

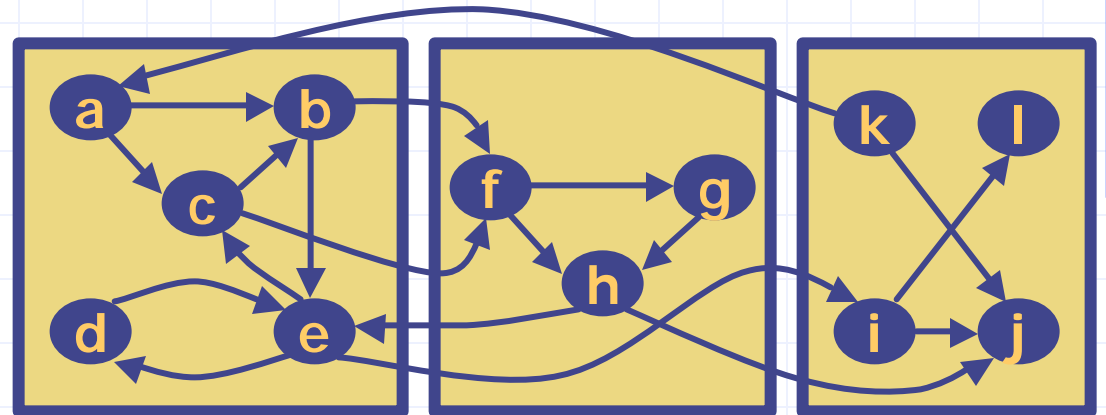


MeCI Example

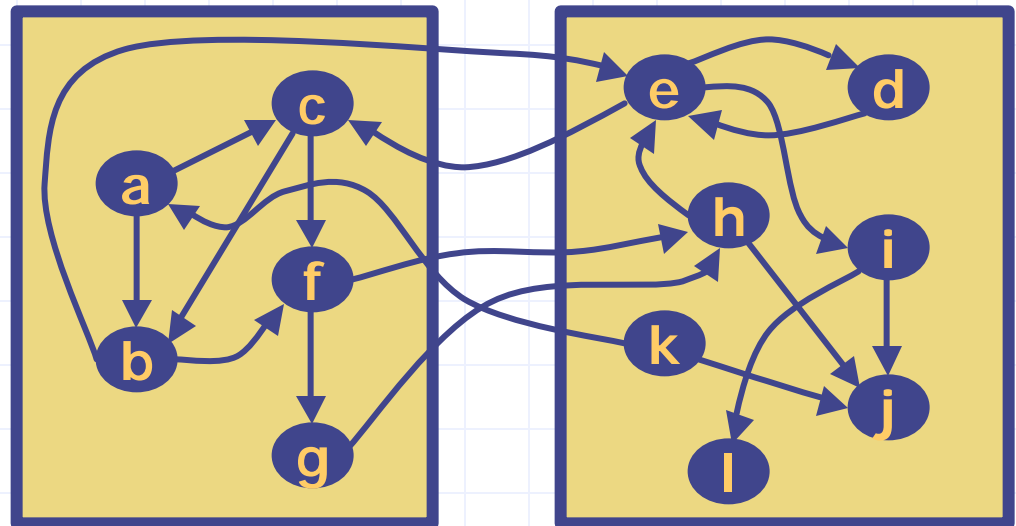
MDG



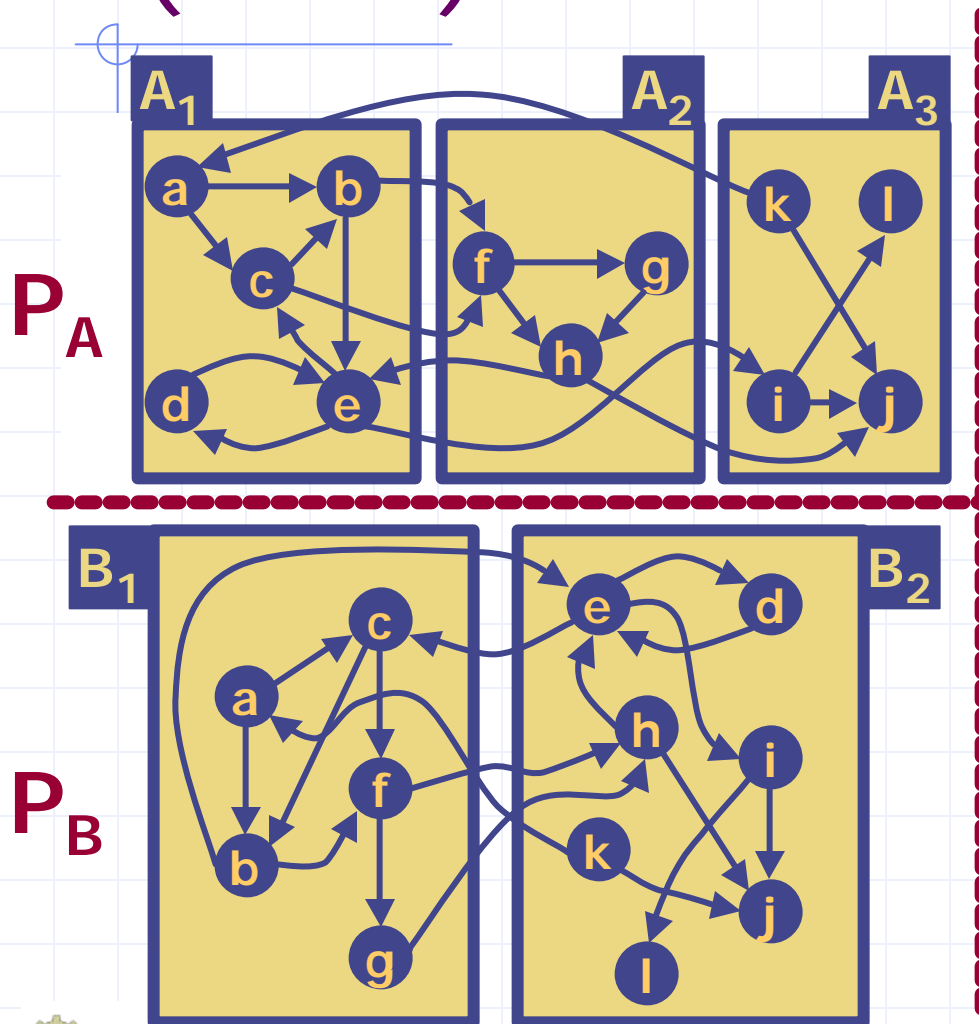
P_A



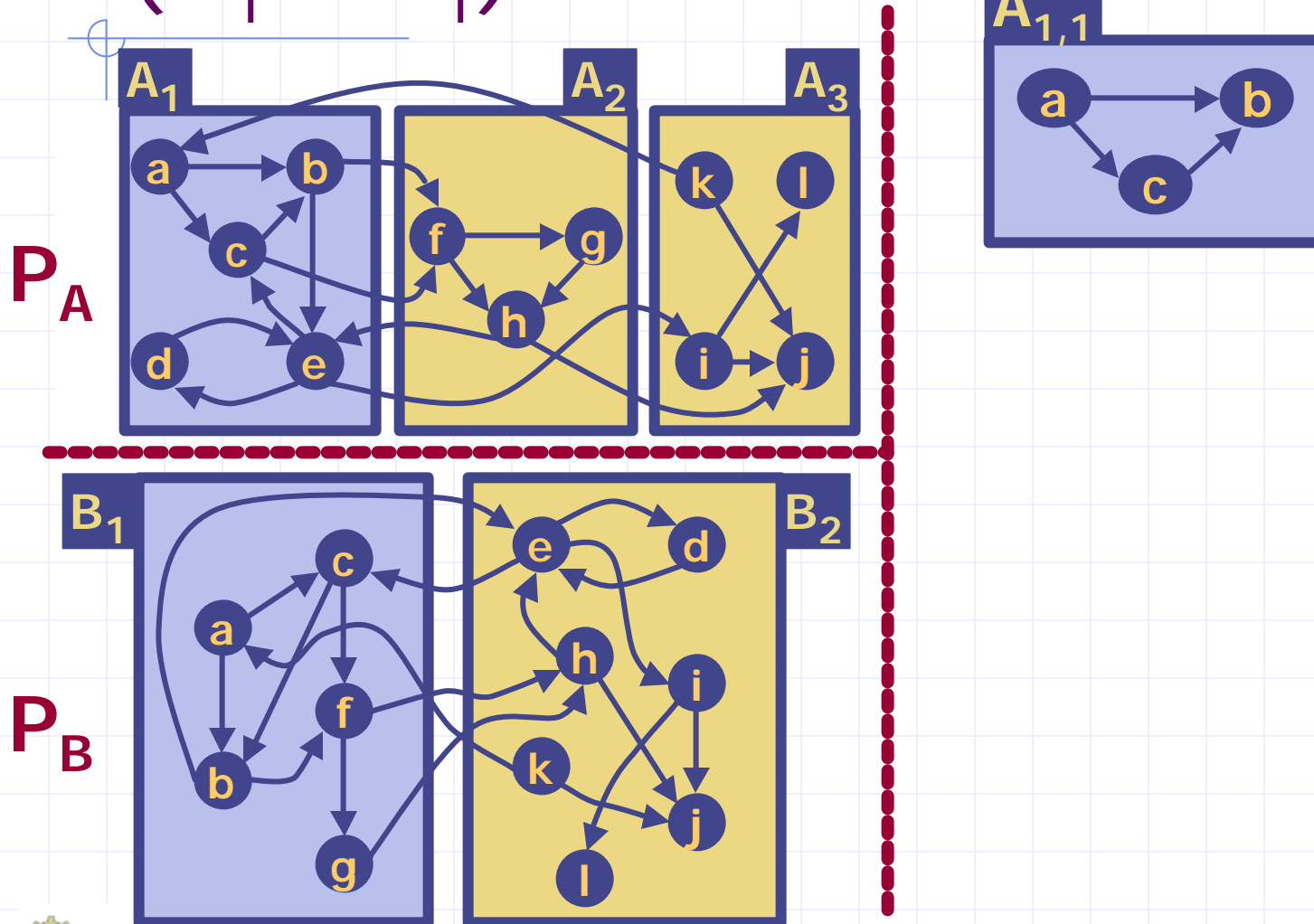
P_B



MeCI Example ($A \rightarrow B$)

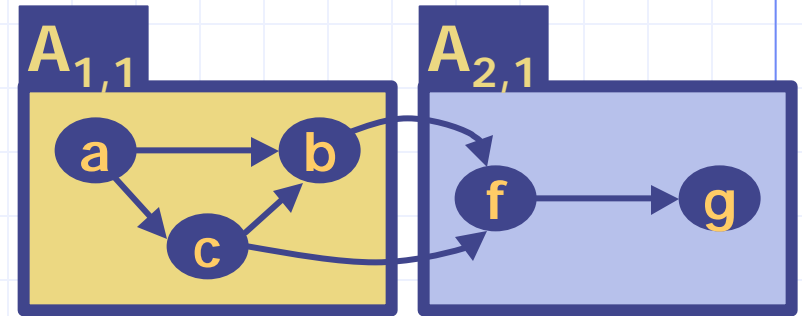
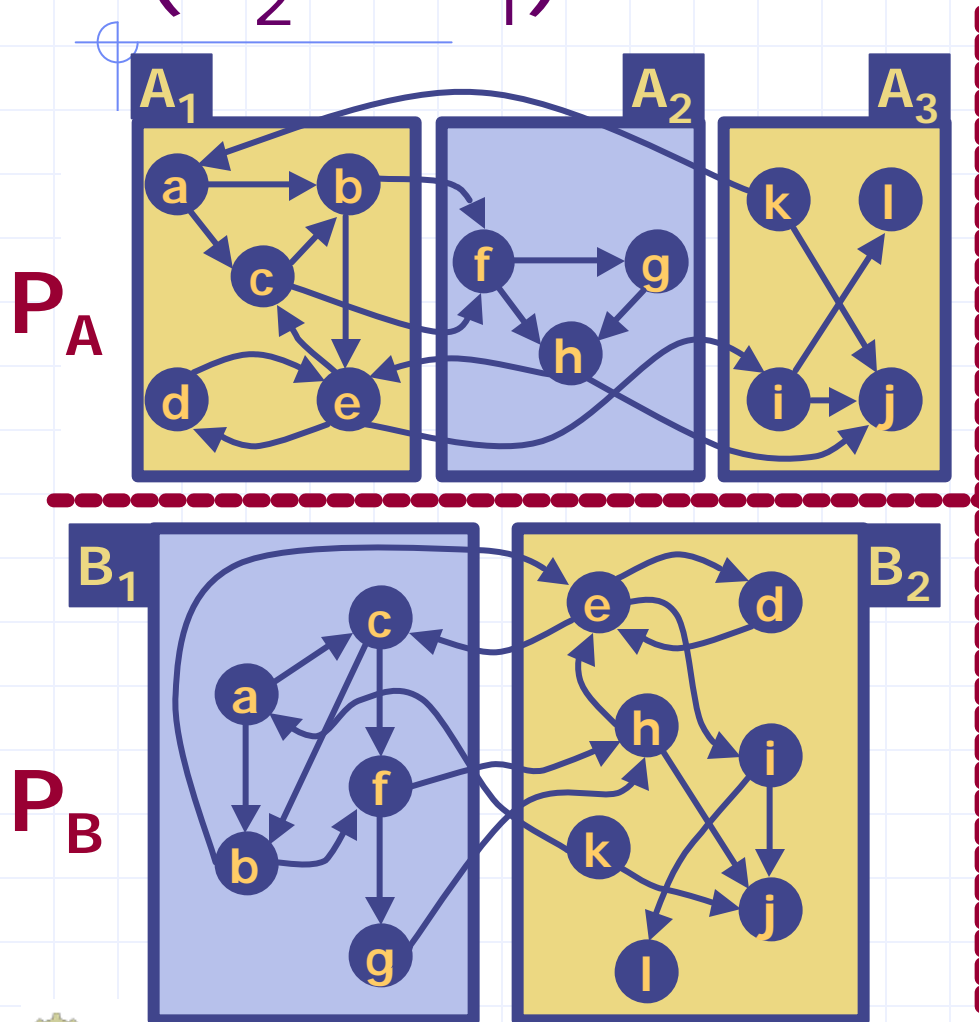


MeCI Example ($A_1 \cap B_1$)

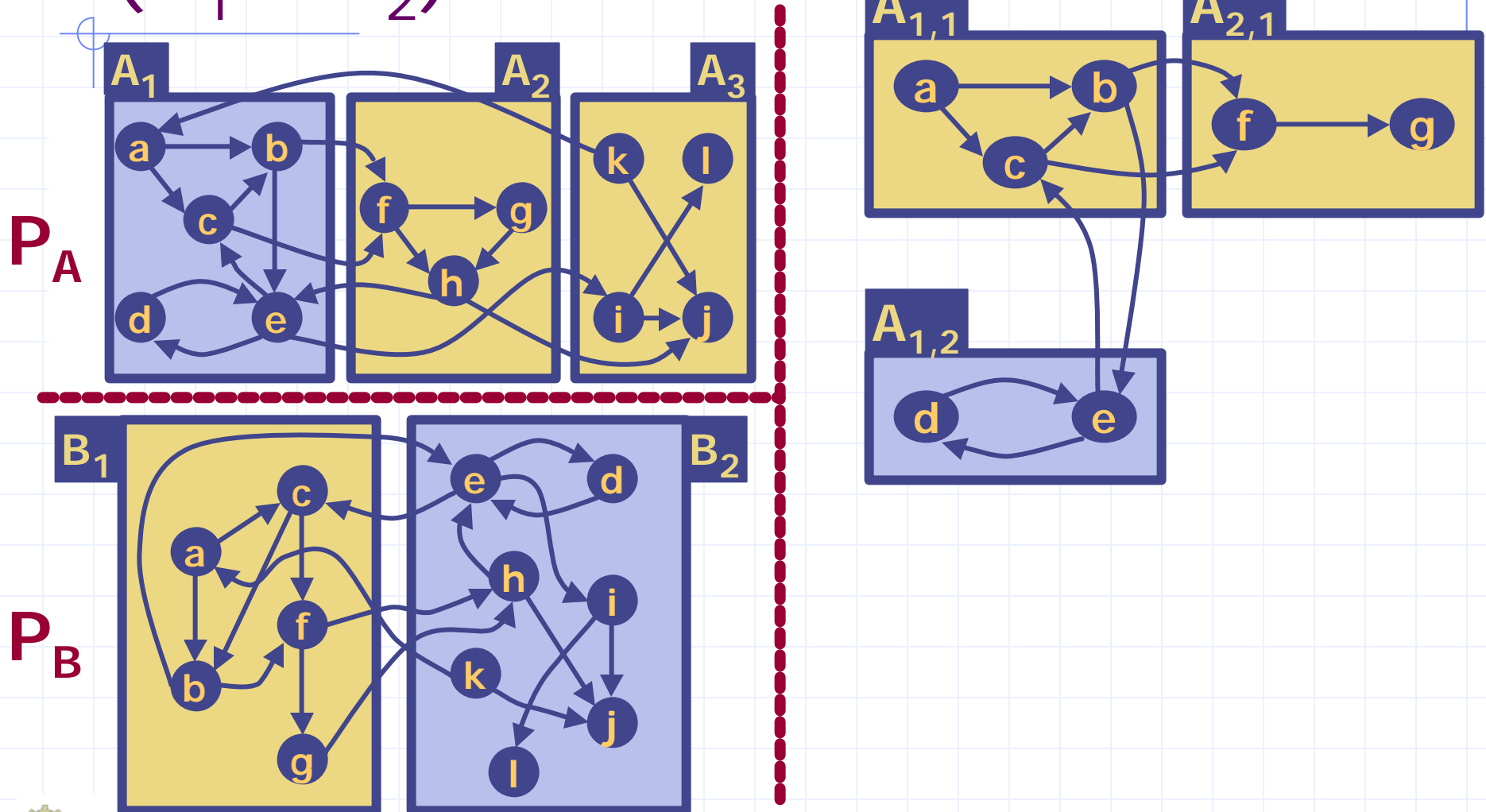


MeCI Example

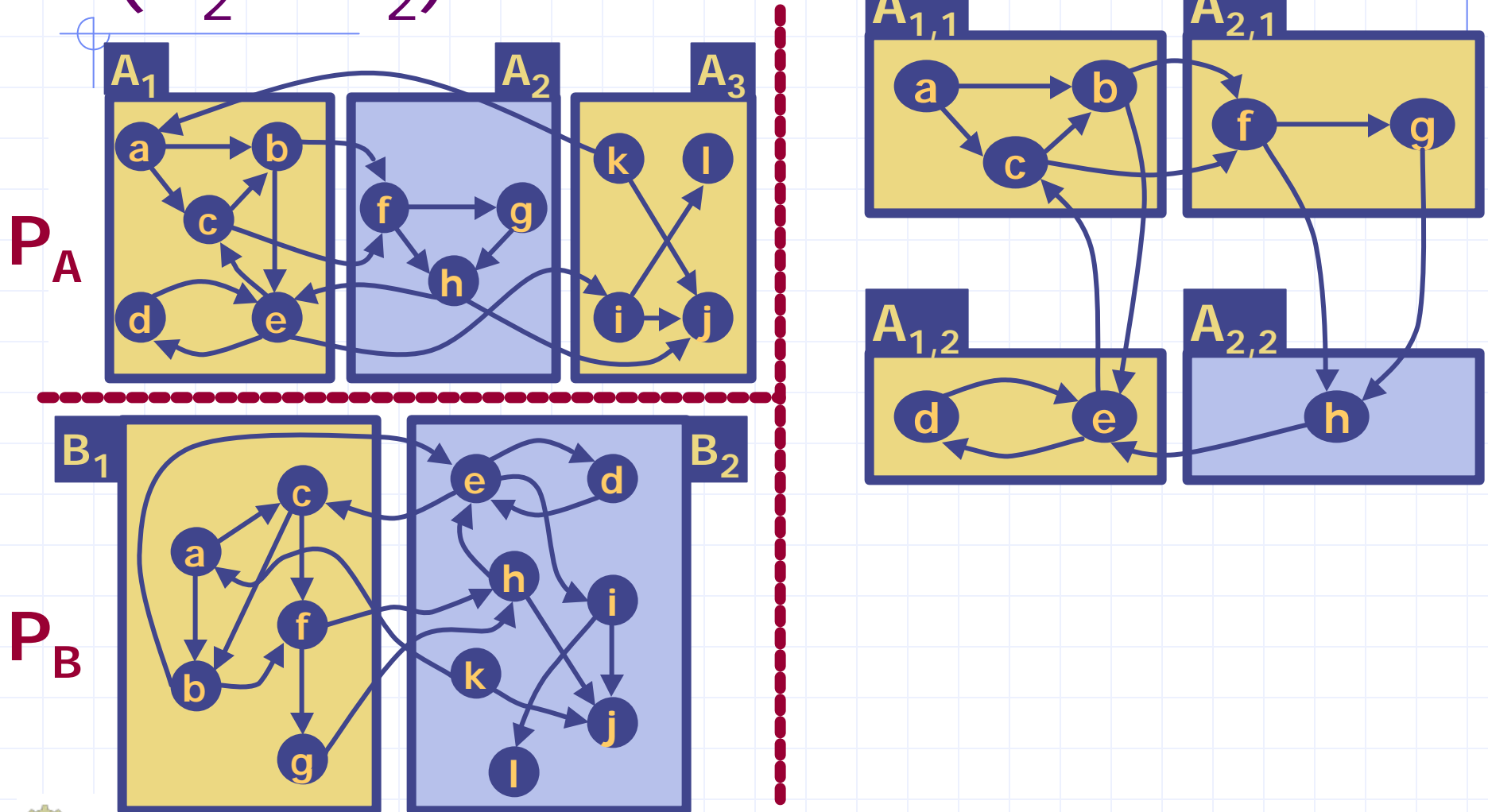
$(A_2 \cap B_1)$



MeCI Example ($A_1 \cap B_2$)

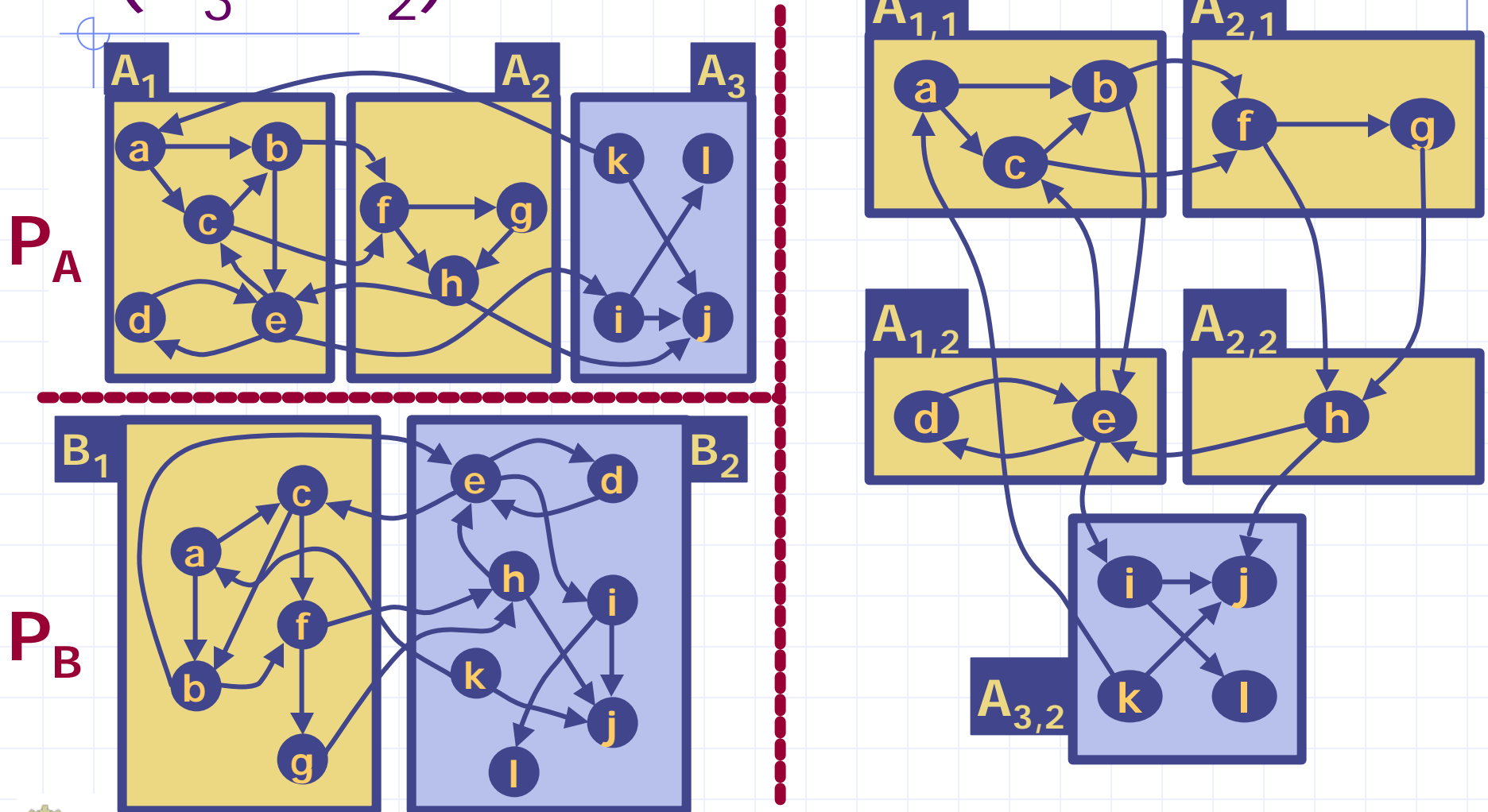


MeCI Example ($A_2 \cap B_2$)

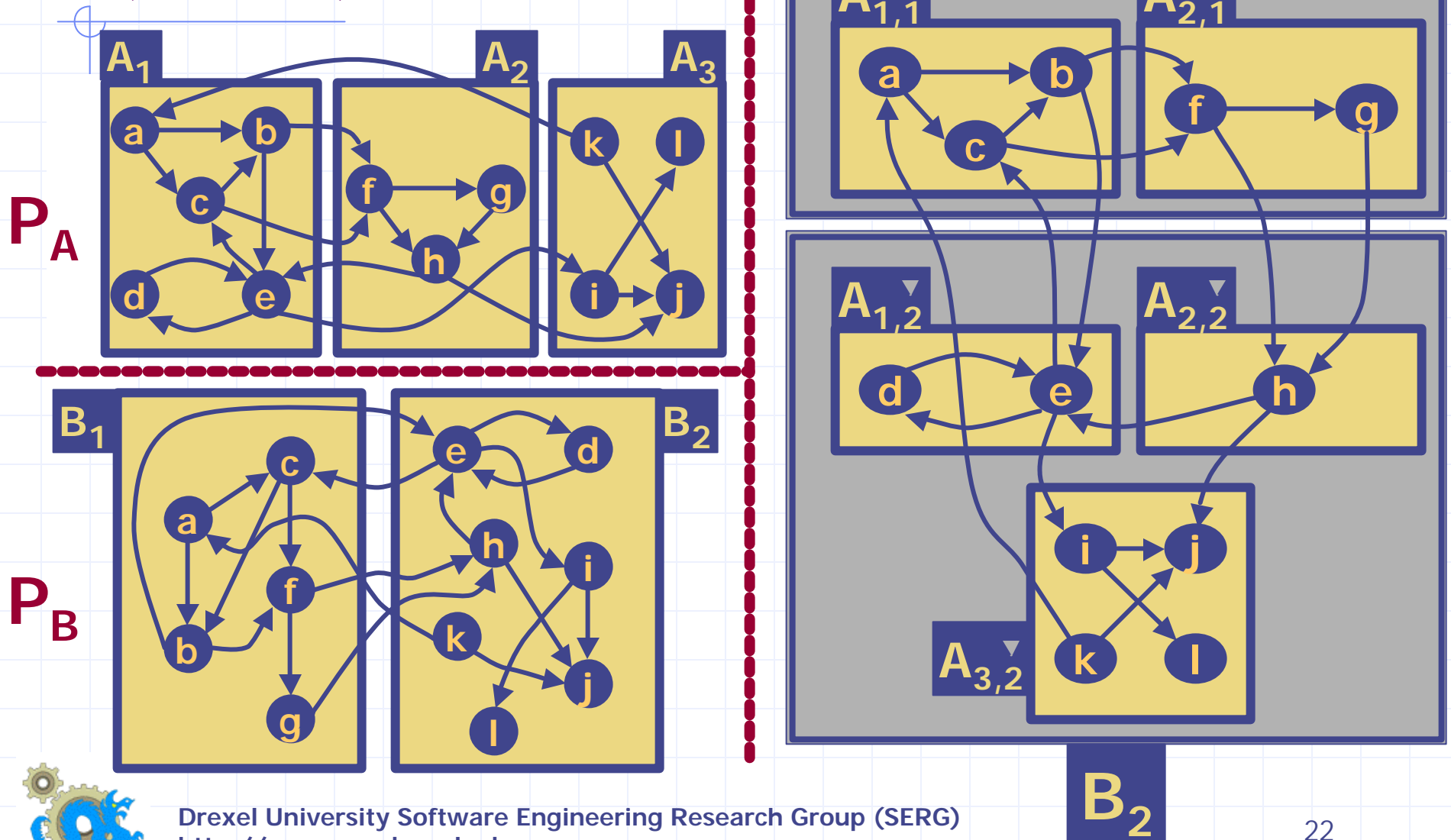


MeCI Example

$(A_3 \cap B_2)$

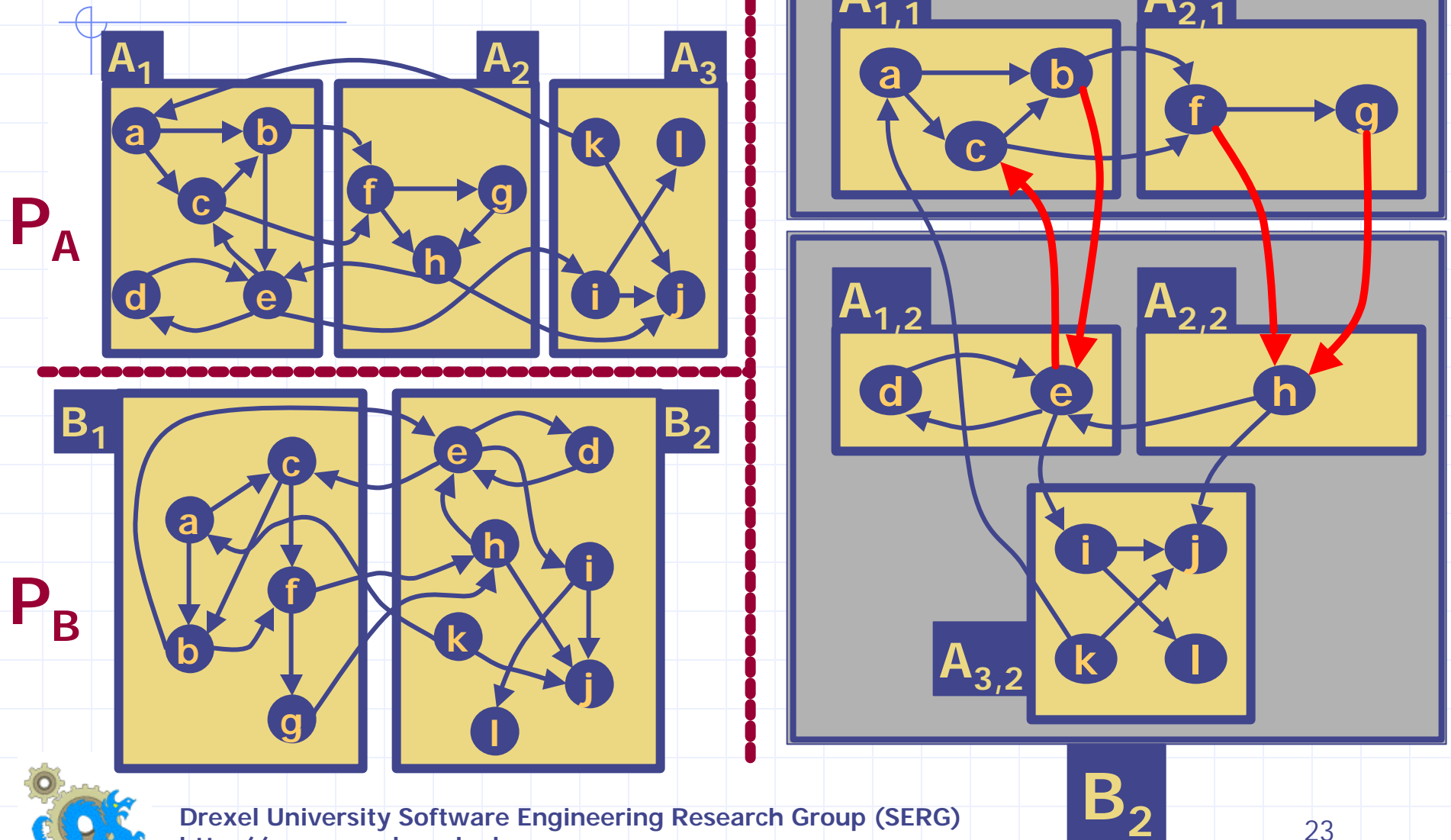


MeCI Example ($A \rightarrow B$)

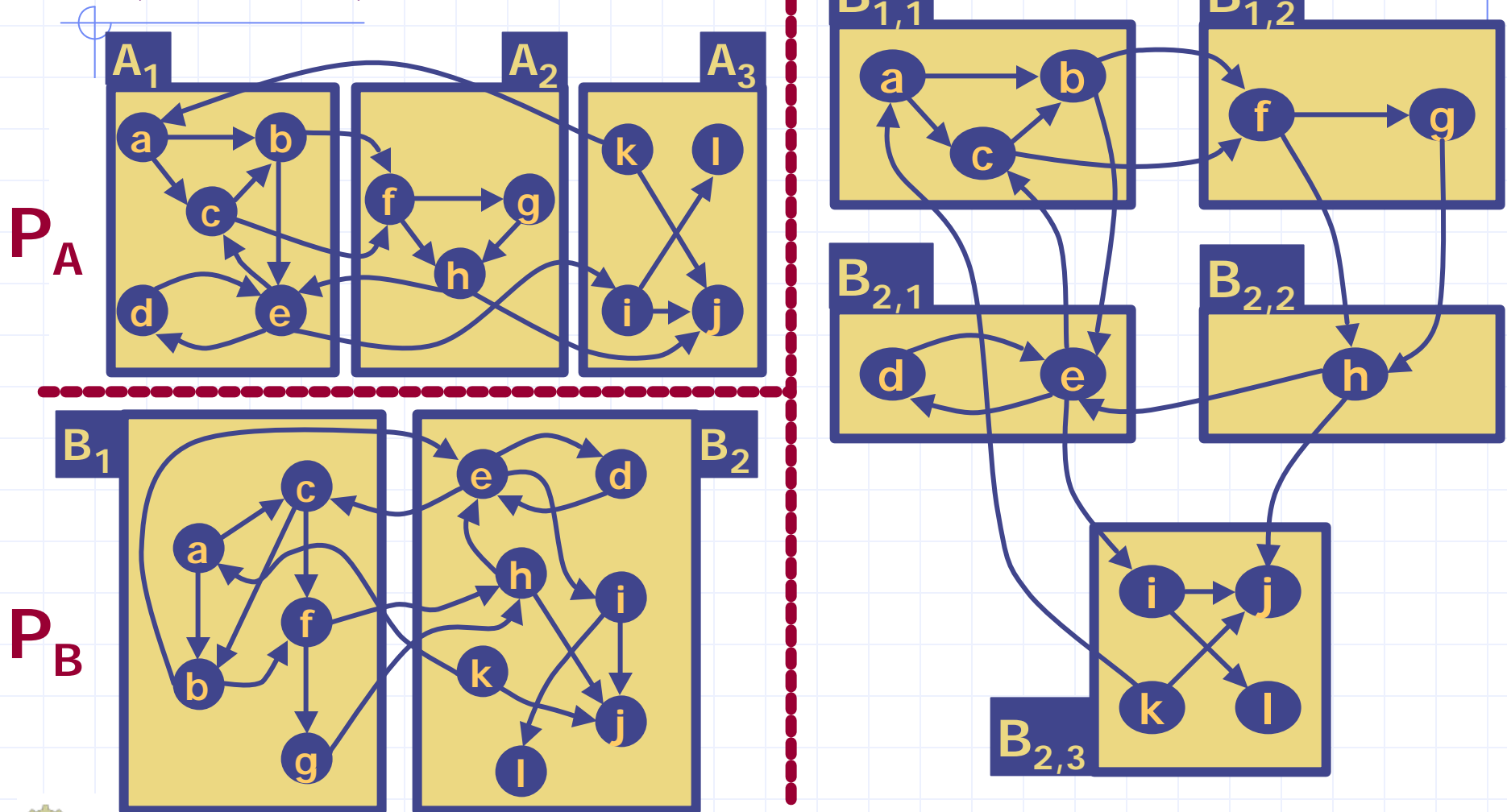


MeCI Example ($A \rightarrow B$)

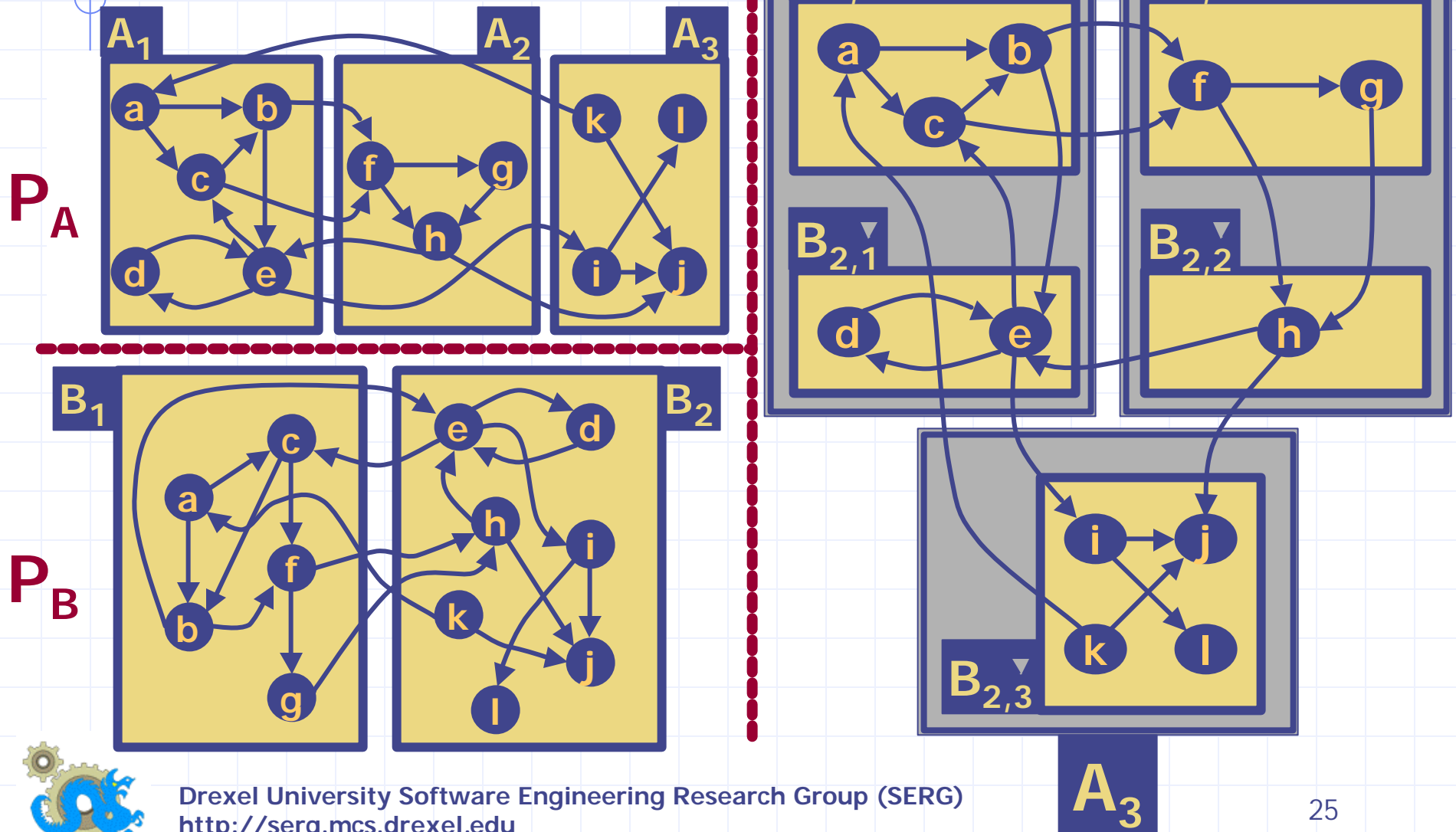
Newly Introduced Inter-Edges



MeCI Example ($B \rightarrow A$)

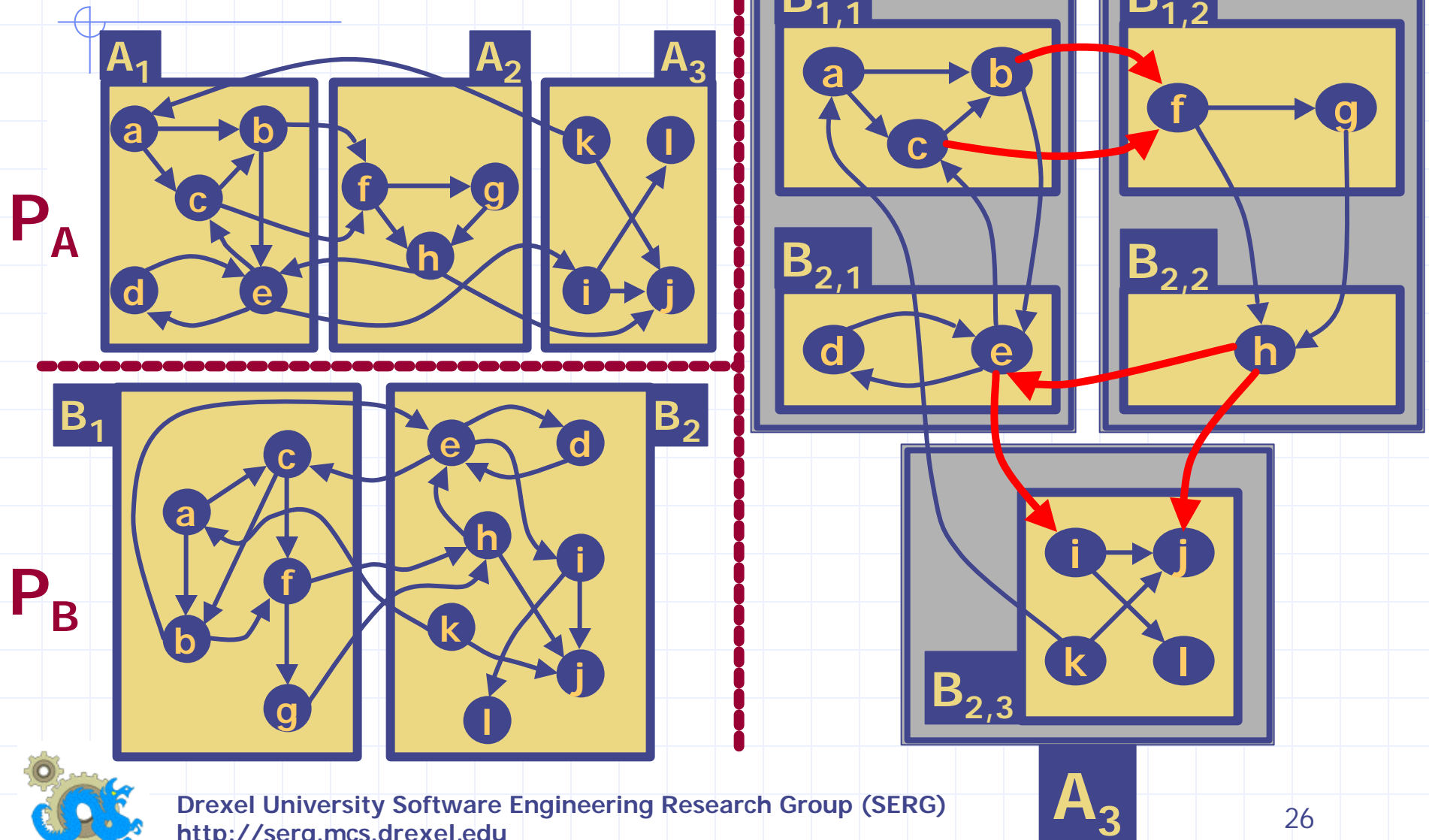


MeCI Example ($B \rightarrow A$)

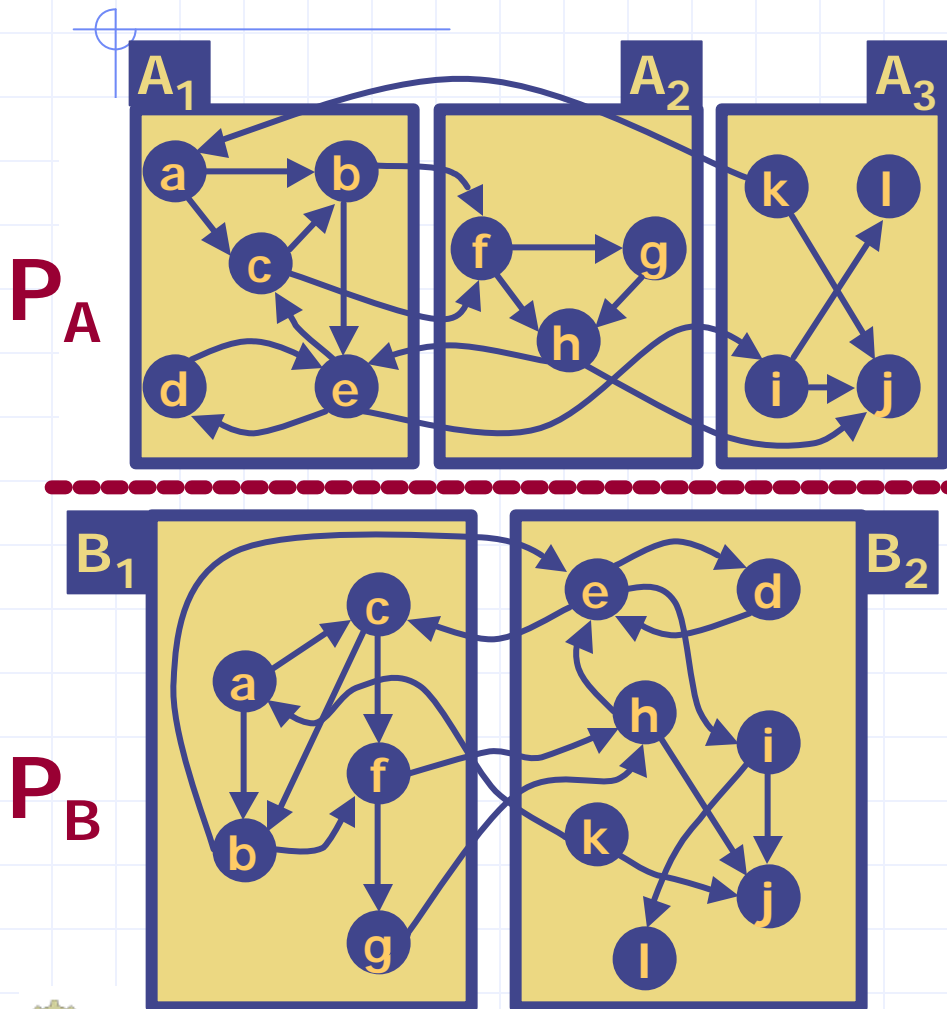


MeCI Example ($B \rightarrow A$)

Newly Introduced Inter-Edges



MeCI Calculation



Inter-Edges Introduced

MeCI(A→B):
 ({b,e},{e,c},{g,h},{f,h})

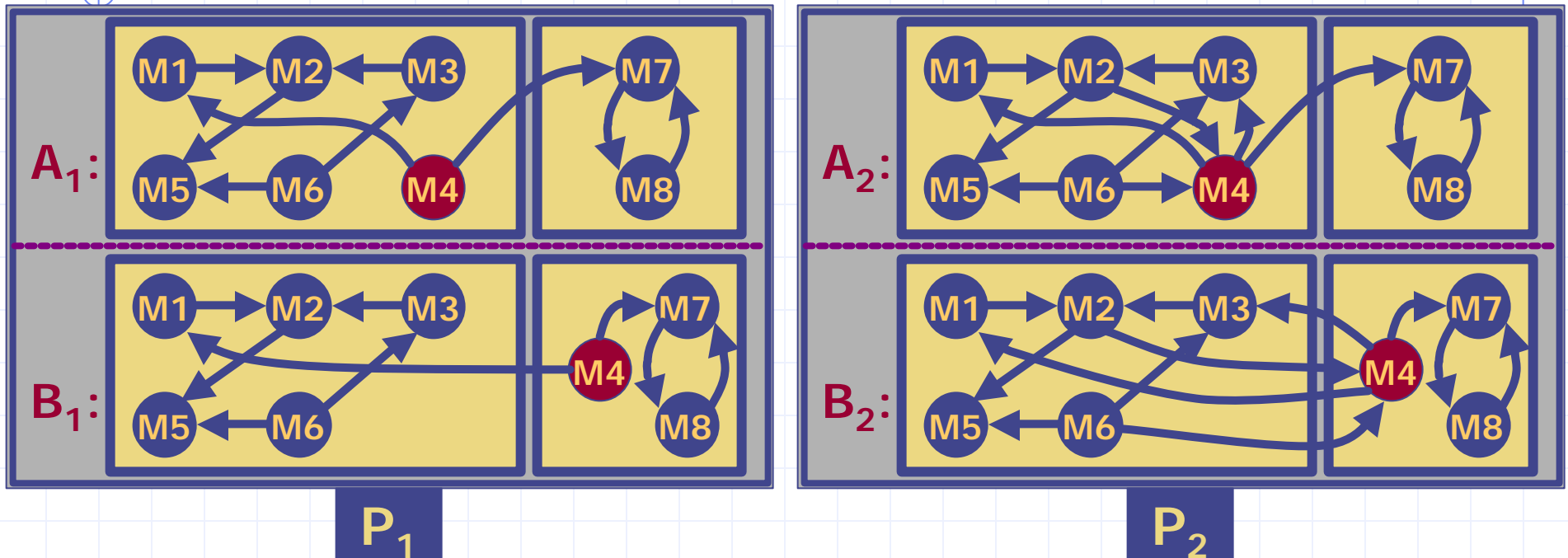
MeCI(B→A):
 ({e,i},{h,j},{b,f},{c,f},{h,e})

$$\text{MeCI} = \left[1 - \frac{\max_W(M_{A \rightarrow B}, M_{B \rightarrow A})}{\text{Total Edge Weight}} \right]$$

$$\text{MeCI} = \left[1 - \frac{5}{19} \right] = 73.7\%$$



Similarity Measurement Recap



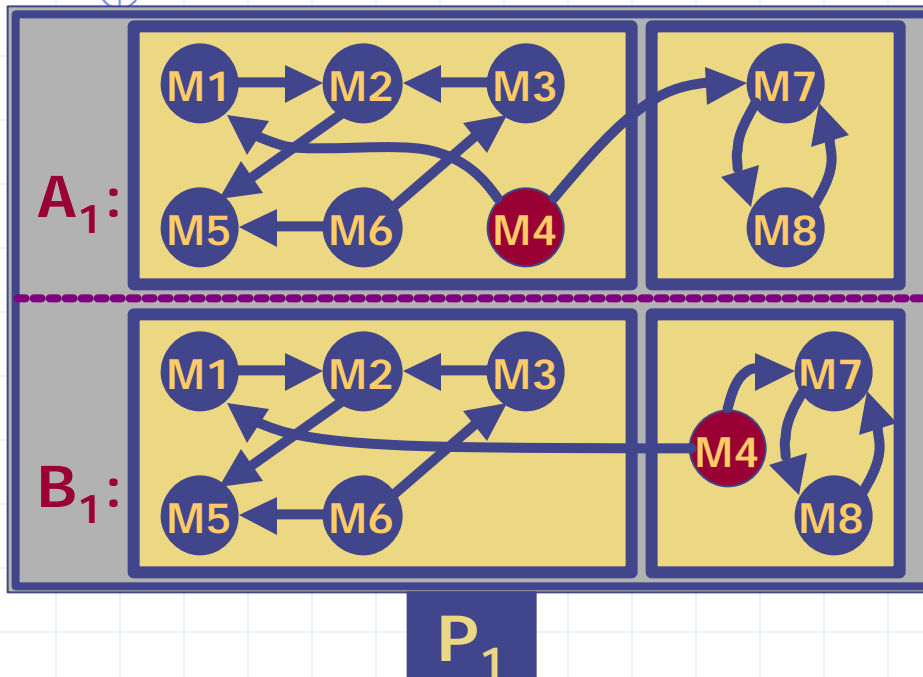
$$\text{MoJo}(P_1) = \text{MoJo}(P_2) = 87.5\%$$

$$\text{PR}(P_1) = \text{PR}(P_2) = \text{P:}84.6\%, \text{R:}68.7\%, \text{AVG}_{\text{PR}}=76.7\%$$

Conclusion... P_1 is equally similar to P_2

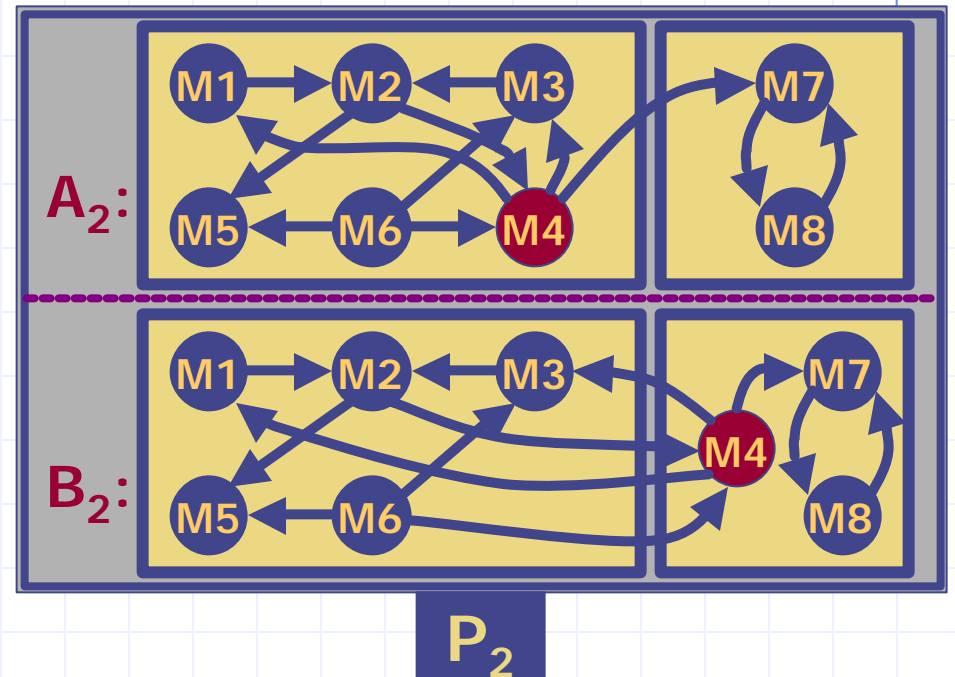


Similarity Measurement Recap



EdgeSim(P₁) = 77.8%

MeCI(P₁) = 88.9%



EdgeSim(P₂) = 58.3%

MeCI(P₂) = 66.7%

Conclusion... P₁ is more similar than P₂



Summary: EdgeSim & MeCI

◆ EdgeSim:

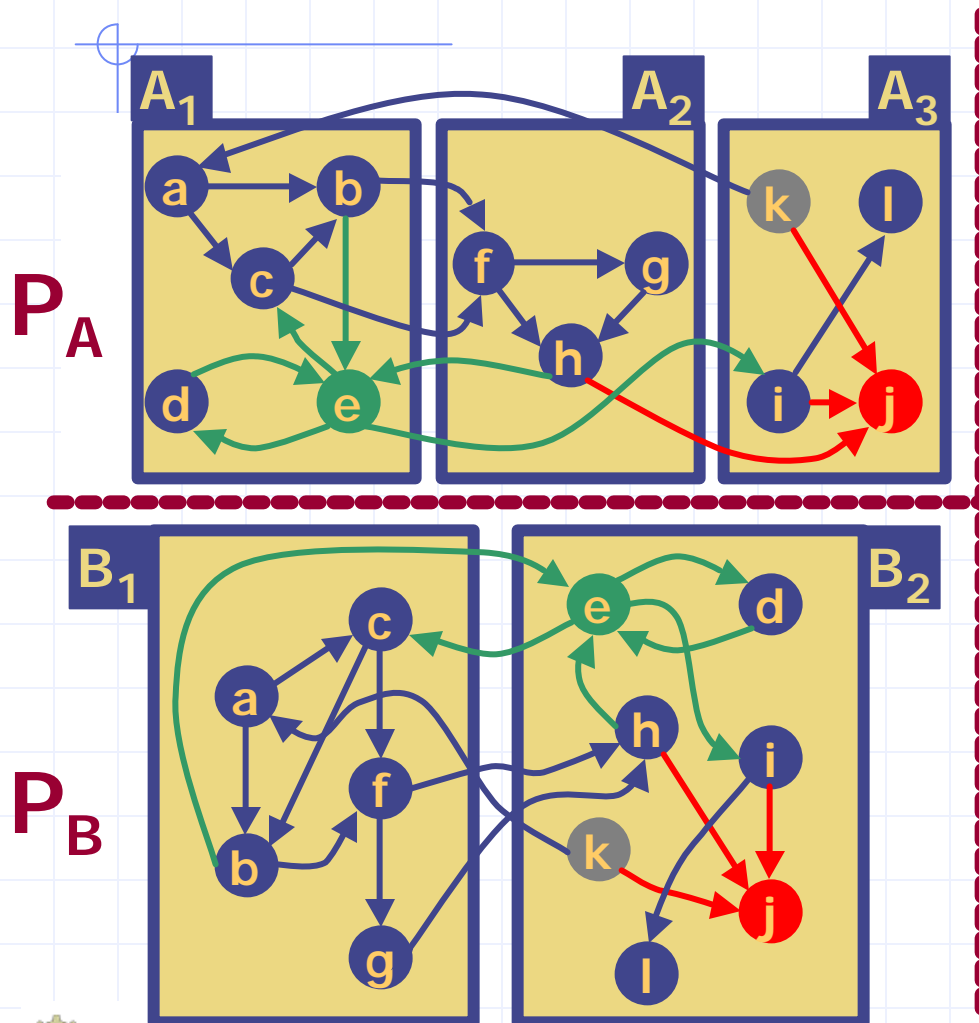
- Rewards clustering algorithms for preserving the edge types
- Penalizes clustering algorithms for changing the edge types

◆ MeCI:

- Rewards the clustering algorithm for creating cohesive “subclusters”



Special Modules



Omnipresent Modules:

"Strong" Connection to other Modules

Library Modules:

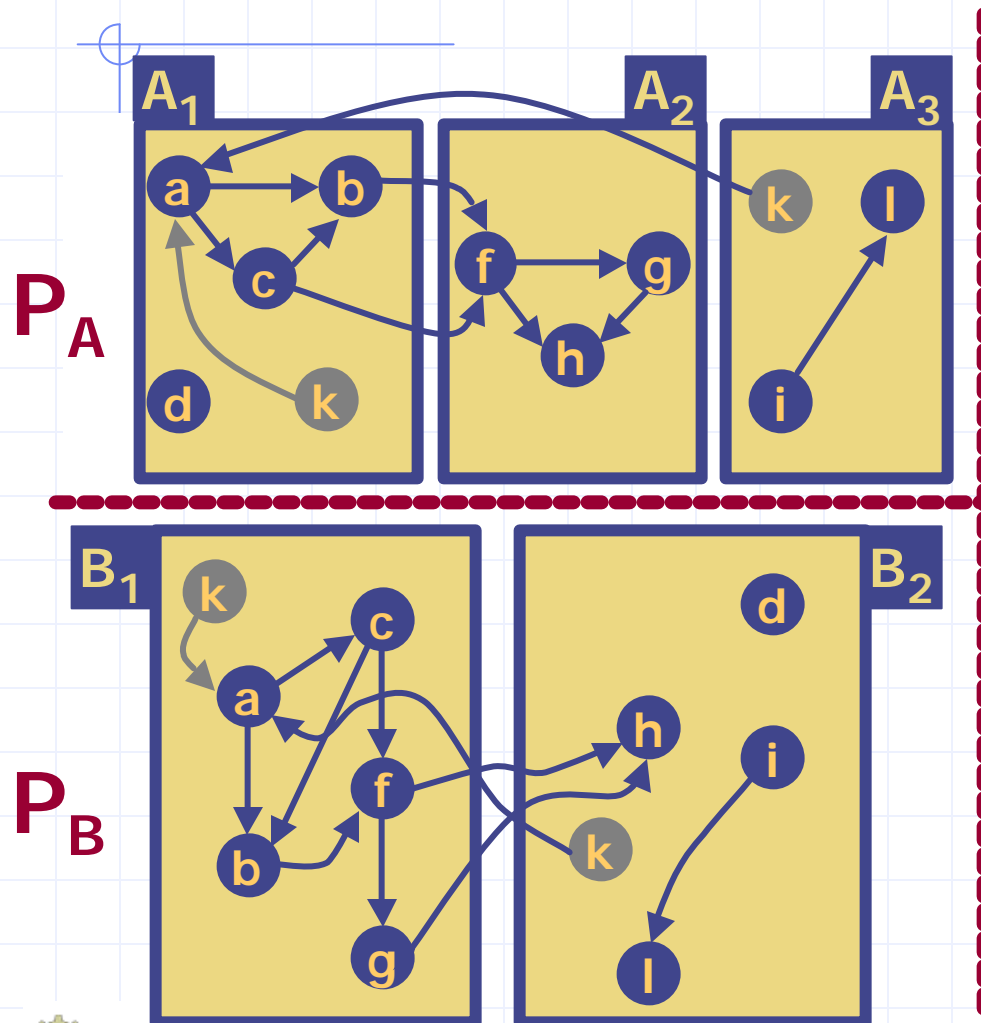
Always used by other modules, never use other modules

Isomorphic Modules:

Modules equally connected to other subsystems



Special Modules



Special Treatment of Special Modules helps to determine the Similarity

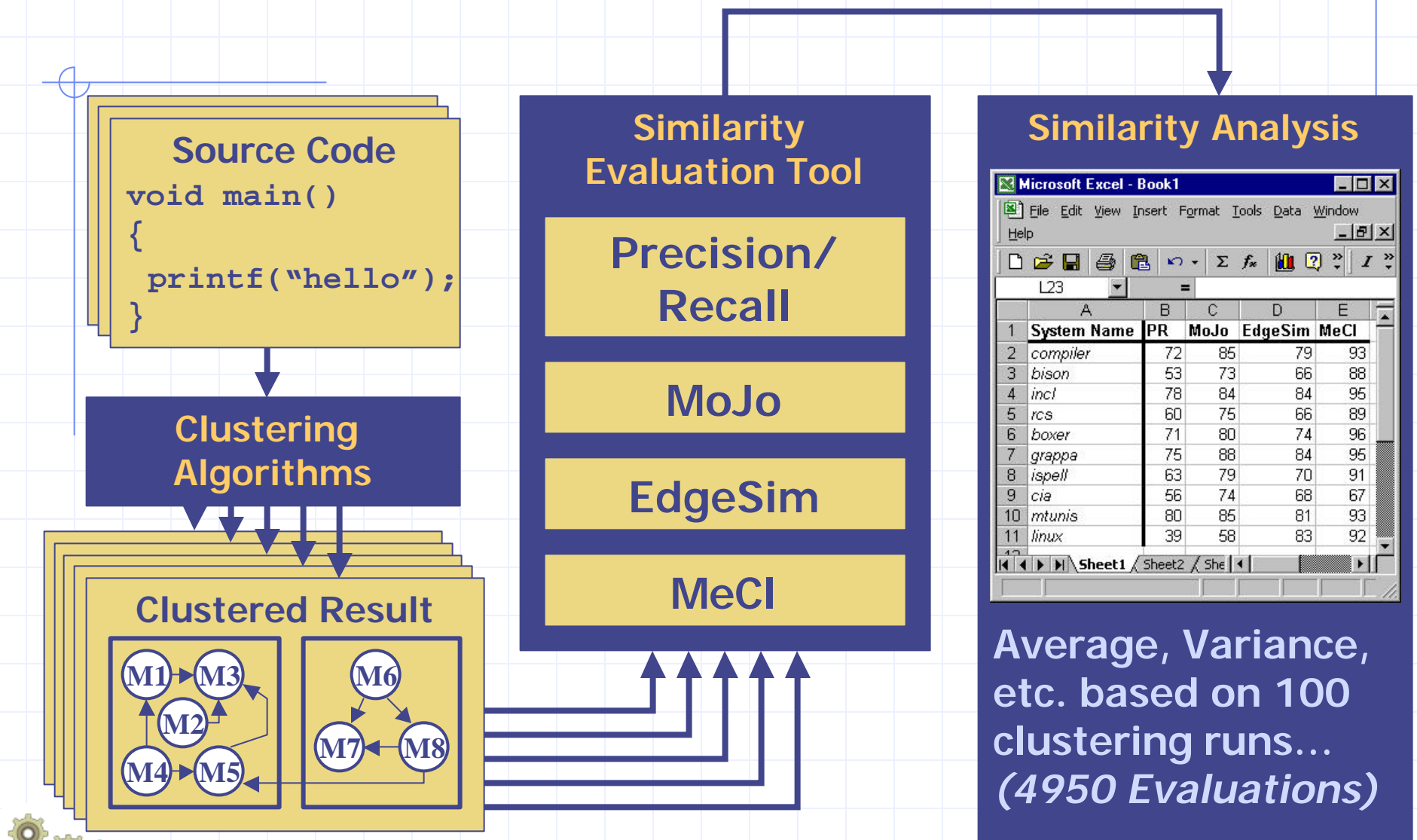
Omnipresent Modules:
Removed

Library Modules:
Removed

Isomorphic Modules:
Replicated



Case Study Overview



Case Study Observations

- ◆ All similarity measurements exhibit consistent behavior for the systems studied

For all systems examined:

**If $\text{MeCl}(S_A) < \text{MeCl}(S_B)$ then $\text{MoJo}(S_A) < \text{MoJo}(S_B)$,
 $\text{PR}(S_A) < \text{PR}(S_B)$, and $\text{EdgeSim}(S_A) < \text{EdgeSim}(S_B)$**

- ◆ Removal of “special” modules improved all similarity measurements
- ◆ Treating isomorphic modules specially only improved similarity slightly
- ◆ EdgeSim and MeCl produced higher and less variable similarity values than Precision/Recall and MoJo



Questions



◆ Special Thanks To:

- AT&T Research
- Sun Microsystems
- DARPA
- NSF
- US Army

