# FIRST: Fast Interactive Attributed Subgraph Matching

Presented by: Boxin Du

Arizona State University

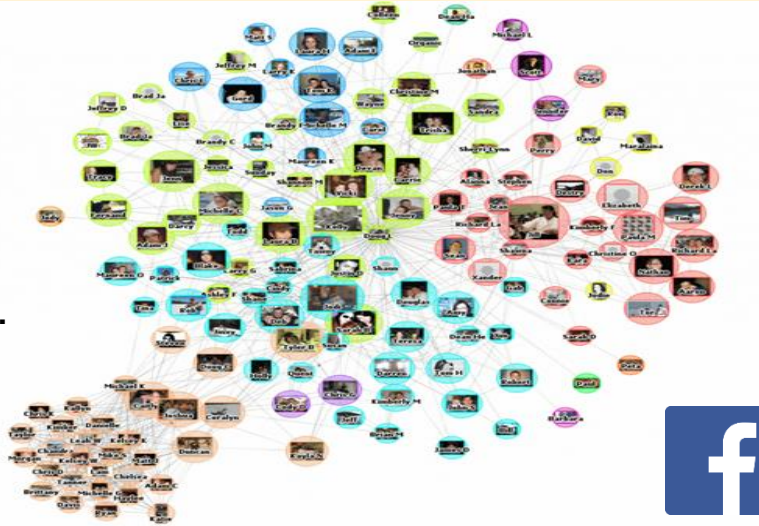Boxin Du      Si Zhang      Nan Cao      Hanghang Tong

DATA Lab

# Attributed Networks Are Everywhere

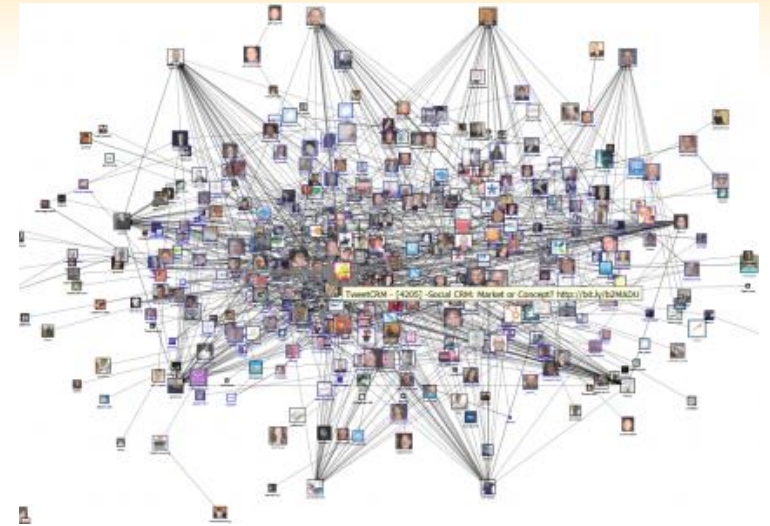**Node Attribute:**
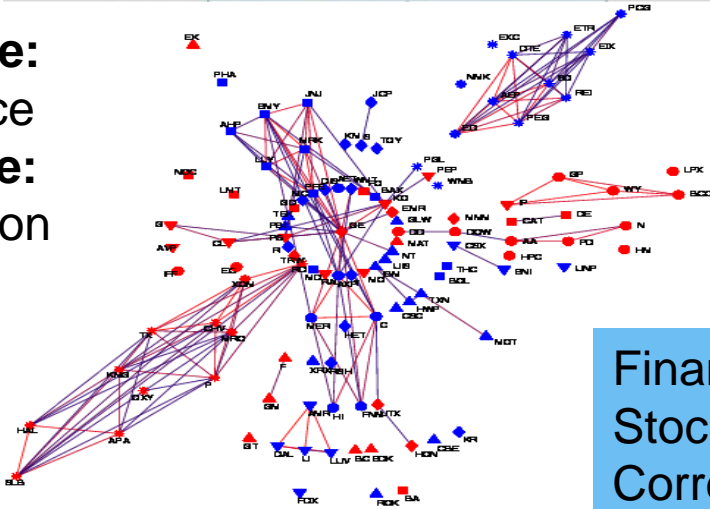Age/gender/
interests…
**Edge Attribute:**
Friendship/like…

**Node Attribute:**
Education/skills/
occupation…
**Edge Attribute:**
Follow/inmail…

**Node Attribute:**
Stock type/price
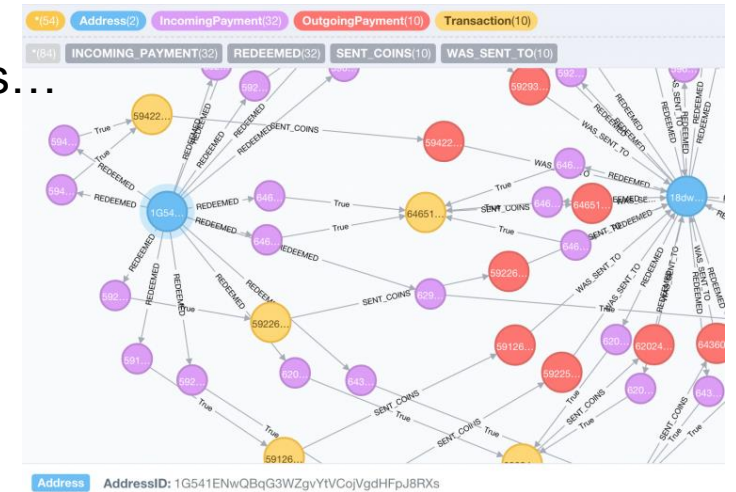**Edge Attribute:**
Stock correlation

**Node Attribute:**
Address/payments…
**Edge Attribute:**
Transaction types
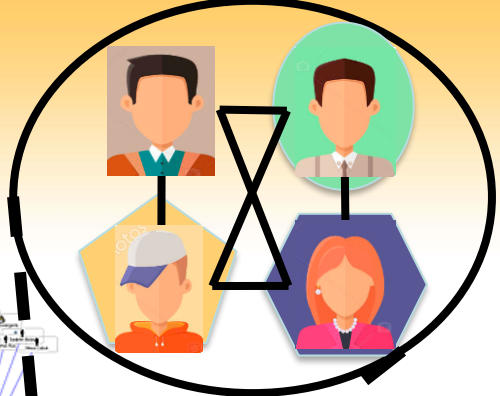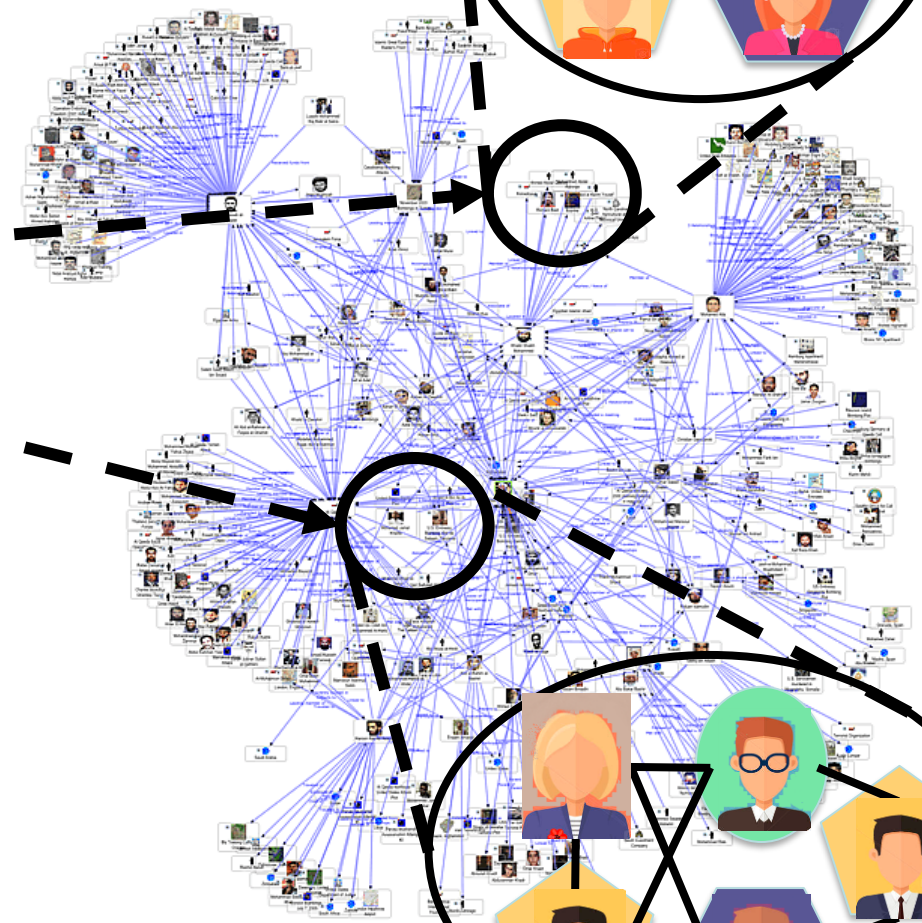
Financial:
Stock
Correlation

Bitcoin
Transaction

- **Q:** How to explore attributed networks?

DATA Lab

# A: Attributed Subgraph Matching

- To find structures to your interest.



**Query Graph:**

Accountant · Manager

Engineer · Market Investigator

**Exact Matching**

**Data Graph**

**Inexact Matching**

User

DATA Lab

# Attributed Subgraph Matching Applications

**Finding Chemical Substructures/ Similarities [OEChem TK library]**

Substructure -> chemical property

**Protein-to-protein Interaction Network [A Vinayagam' 14]**

Pathways in PPI -> Protein function

**Financial Fraud Pattern Detection [Neo4j]**

Subgraph Pattern -> Fraud Pattern

**Social Media Analysis [W Fan' 12]**

Graph patterns -> social patterns

and so on…

**DATA Lab**

# Existing Methods for Attributed Subgraph Matching (MANY!)

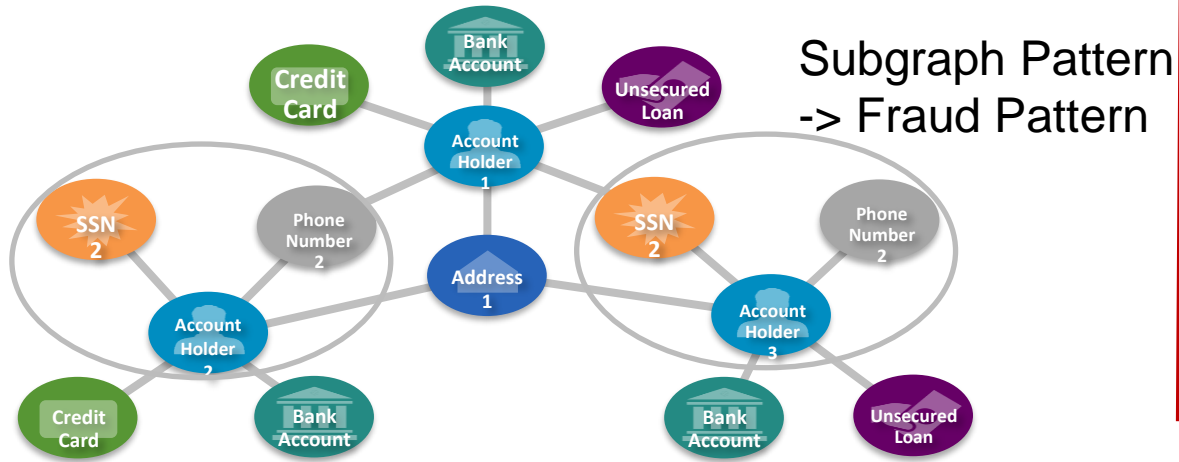| Algorithm | Author & Conference | Exact Matching | Inexact Matching | Node Attribute | Edge Attribute | Require no Index | Accurate query |
|-----------|---------------------|:--------------:|:----------------:|:--------------:|:--------------:|:----------------:|:--------------:|
| *R-WAG/I-WAG/S-WAG* | S Roy et al. TKDE' 15 | ✔ | ✔ | ✔ | ✘ | ✔ | ✘ |
| *MAGE* | R. Pienta et al. IEEE BigData' 14 | ✔ | ✔ | ✔ | ✔ | ✔ | ✘ |
| *NeMa* | A. Khan et al. VLDB' 13 | ✔ | ✔ | ✔ | ✘ | ✔ | ✘ |
| *IncMatch* | W. Fan et al. SIGMOD'11 | ✔ | ✘ | ✔ | ✘ | ✔ | ✘ |
| *SIGMA* | M.Mongiovi et al. CSB' 09 | ✘ | ✔ | ✔ | ✘ | ✘ | ✘ |
| *TALE* | Y. Tian et al. ICDE' 08 | ✔ | ✔ | ✔ | ✔ | ✘ | ✘ |
| *G-Ray* | H. Tong et al. KDD' 07 | ✔ | ✔ | ✔ | ✘ | ✔ | ✘ |
| ... | ... | ... | ... | ... | ... | ... | ... |

- **Obs:** The User Needs to Provide the **Accurate** Query Graph.

- **Q:** What if the user does not know exactly what s/he is looking for?

DATA Lab

# Interactive Attributed Subgraph Matching

- An illustrative example:
  - **Given:**
    - a social network with node attributes and edge attribute;
    - an initial query with attributes;
  - **Find:** the best matching subgraph(s) with **query revision on-the-fly**.

**Initial Query:**

**Data Network:**

**Attributes:**

⬠ Programming    🟡 Machine Learning

⬡ Data Mining    🟦 Visualization

——— Face2Face

- - - E-mail    〜〜 Phone

DATA Lab

# Interactive Attributed Subgraph Matching (cont'd)

- An illustrative example (cont'd): Revising and matching process.



- **Challenge**: How to respond **efficiently**?
  - w/o re-running algorithm or re-building Graph indexes

DATA Lab

# Roadmap

➢ Motivation

➢ Problem Definition

➢ Proposed Solution: FIRST family
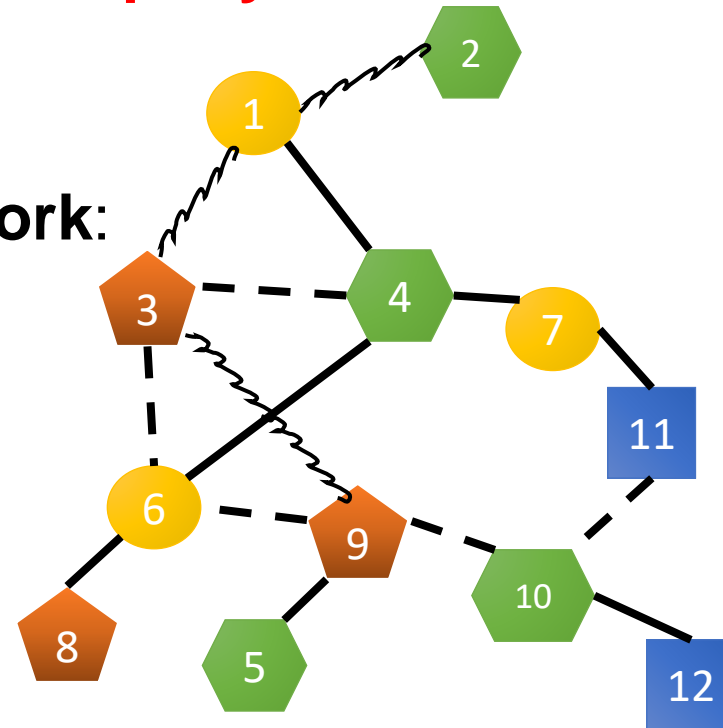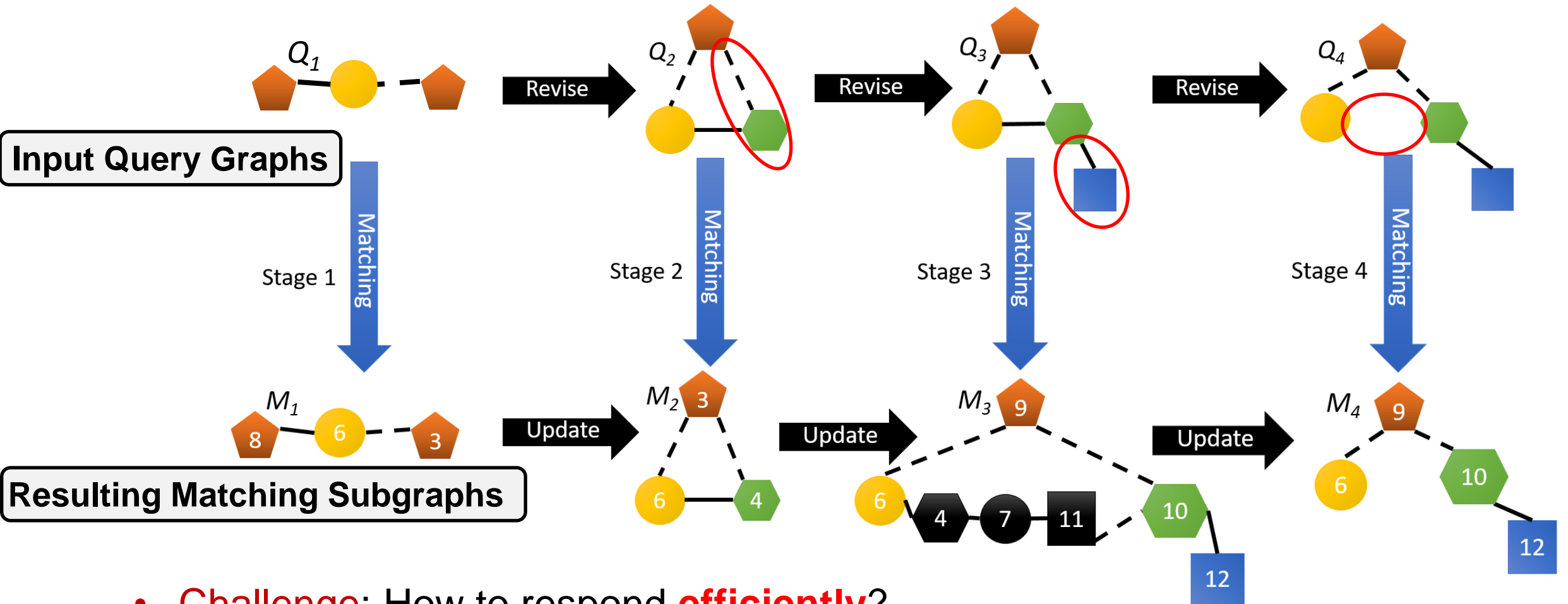
➢ Experiments

➢ Conclusions

DATA Lab

# Problem Definition

**Given**:

**❶** an undirected attributed network $G = (A, \ N_A, E_A)$, $[A: n \times n]$

**❷** an undirected initial query graph $Q = (A_q, N_q, E_q)$, $[A_q: k \times k]$

**❸** the initial matching graph $M$

**❹** the revised query graph $\widetilde{Q}$;

**Output:** the updated matching subgraph $\widetilde{M}$.

**Input:**

$n = 12$



Data Network

❷ Initial Query $Q$

❸ Initial Matching $M$

❹ Revised Query $\widetilde{Q}$

$k = 3$

**Output:**

$\widetilde{M}$

DATA Lab

# Roadmap:

➢ Motivation

➢ Problem Definition

➢ **Proposed Solution: FIRST family**

    ➢ Key ideas

    ➢ Details

➢ Experiments

➢ Conclusions

**DATA Lab**

# Key ideas

- <span style="color:red">Key Idea #1: Matching as cross-network node similarity</span>
- <span style="color:gold">Potential Benefit:</span>
  - Encodes both topology and attribute during matching
  - Major computation: Sylvester equation

- <span style="color:red">Key Idea #2: Explore the smoothness of query graphs</span>
- <span style="color:gold">Potential Benefit:</span>
  - View the revised query as a perturbation of previous query
  - Incrementally solve Sylvester equation for fast response

DATA Lab

# Key Idea #1: Matching as **cross-network node similarity**

- Step 1: Find Similarity Matrix (**S**): *FIRST-Q/N/E*
  - **Intuition:** cross-network node similarity [Zhang, et al KDD'16]
  - **Major Computation:** to solve the Sylvester Equation

$$\mathbf{s} = \alpha \mathbf{W}\mathbf{s} + (1 - \alpha)\mathbf{h}$$

$\mathbf{W}$ : Kronecker graph of **G** and **Q** (filtered by node/edge attribute);
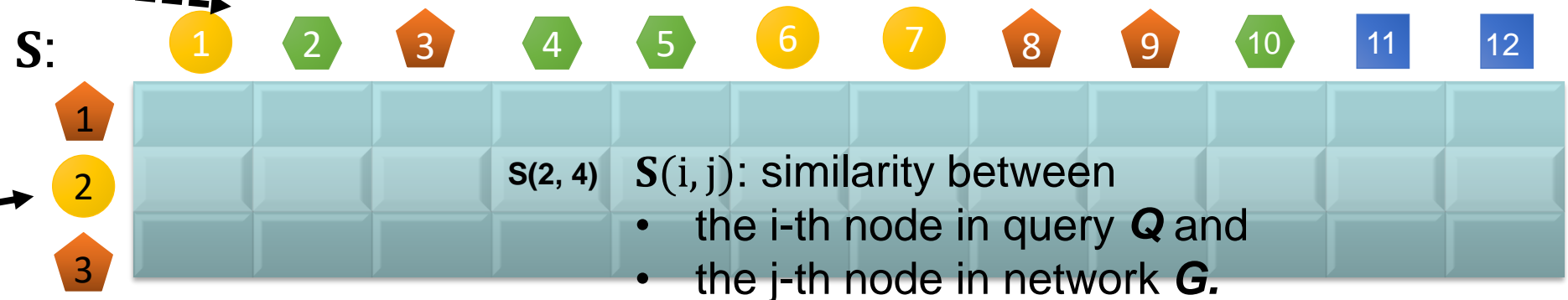
$\mathbf{s}$: Vectorized similarity matrix $\mathbf{S}$;

$\mathbf{h}$: Vectorized preference matrix $\mathbf{H}$.



S(2, 4)   **S**(i, j): similarity between
  - the i-th node in query **Q** and
  - the j-th node in network **G**.

# Key Idea #1 (cont'd):

- Step 2: Find Matching Subgraph: *Sim2Sub*
  - **Intuition:**
    1. From similarity matrix to permutation matrix **X**;
    2. From permutation matrix to subgraph.
  - **Major Computation:**
    - ➢ Calculate Matching Indicator Matrix **X** (k by n).



X(2,4)  X(i, j): 1 if the i-th node in **Q** matches the j-th node in **G**, 0 otherwise.

    - ➢ **How to**: Use 'goodness' function $g(\mathbf{X})$:

$$\mathbf{X}^* = argmax\big(g(\mathbf{X})\big)$$

$$= argmax[-\big\|\mathbf{XAX}' - \mathbf{A_q}\big\|_F^2 + a * trace(\mathbf{SX}') - b * \big\|\mathbf{XX}' - \mathbf{I}\big\|_F^2]$$

Matching subgraph Connectivity

Quality of individual matching nodes

Permutation matrix



(1)  *G*

Data Network

(2) Initial Query  *Q*

# Summary of steps in Key Idea #1

- Procedure: Input graphs ➡ cross-network similarity vector ➡ matching subgraph



Filtered by attributes

Sylvester Equation

$argmax(g(\mathbf{X}))$

Sim2Sub

$\mathbf{W}$

$\mathbf{s}$

$\mathbf{X}$

Initial Query $Q$

$G$

Data Network

Initial Matching $M$

$\mathbf{W} = \mathbf{N}[\mathbf{E} \odot (\mathbf{A} \otimes \mathbf{A_q})]\mathbf{N}$

$\mathbf{s} = \alpha \mathbf{W}\mathbf{s} + (1-\alpha)\mathbf{h}$

Local search+ matching nodes connecting

Kronecker Graph filtered by attributes

Similarity Vector

Matching Indicator Matrix

Step 1: FIRST-Q/N/E

Step 2: Sim2Sub

DATA Lab

# Details:

| Scenarios | W |
|---|---|
| Topology only | $\mathbf{W} = \mathbf{A} \otimes \mathbf{A_q}$ |
| Topology + node attribute | $\mathbf{W} = \mathbf{N}\left(\mathbf{A} \otimes \mathbf{A_q}\right)\mathbf{N}$ |
| Topology + node attribute + edge attribute | $\mathbf{W} = \mathbf{N}[\mathbf{E}\odot\left(\mathbf{A} \otimes \mathbf{A_q}\right)]\mathbf{N}$ |

$\mathbf{N}$: the node attribute matrix of input networks ($G \& Q$):

$\mathbf{N} = \sum_{p=1}^{K} N_A^p \otimes N_q^p$, K: the number of distinct node labels.

$\mathbf{E}$: the edge attribute matrix of input networks ($G \& Q$):

$\mathbf{E} = \sum_{l=1}^{L} E_A^l \otimes E_q^l$, L: the number of distinct edge labels.

DATA Lab

# FIRST-Q: Handle Topology Revision

**Scenario 1**: During query revision, only graph topology is changed.

- **Goal**:
    Fast computation of similarity vector after topology revision.

- **Observation**:
    1. We already have: $\mathbf{s} = \alpha \mathbf{W} \mathbf{s} + (1 - \alpha)\mathbf{h}$, $\widetilde{\mathbf{W}} = \mathbf{W} + \Delta \mathbf{W}$,
    2. The approximated similarity matrix:
$$\hat{\mathbf{s}} = (1 - \alpha)(\mathbf{I} - \alpha \widehat{\mathbf{W}})^{-1}\mathbf{h}$$

- **Solution**:
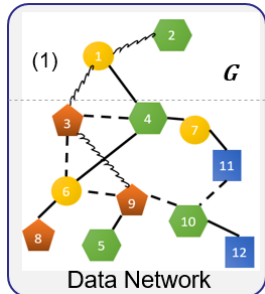    1. Calculate $\mathbf{W}$ in pre-computing stage, only $\Delta \mathbf{W}$ in interactive stage.
    2. Low-rank approx. & matrix inverse lemma for fast computation.

DATA Lab

# FIRST-Q: Handle Topology Revision

- **How to (Details):**

**Pre-computing Stage:**
1. Low-rank approximation for data network:

$$\mathbf{A} \approx \mathbf{U_A} \mathbf{\Lambda_A} \mathbf{U_A^T}$$

2. Store $\mathbf{U_A}$, $\mathbf{\Lambda_A}$.

**Interactive Stage:**
1. $\widetilde{\mathbf{A}}_\mathbf{q} \approx \mathbf{U_Q} \mathbf{\Lambda_Q} \mathbf{U_Q^T}$
2. Compute: node attribute matrix $\mathbf{N}$, diagonal degree matrix $\mathbf{D}$.
3. $\widehat{\mathbf{W}} = \mathbf{L} \; \mathbf{\Lambda} \; \mathbf{R}$

Construct $\mathbf{L}, \mathbf{\Lambda}, \mathbf{R}$ from the approximation of $\mathbf{A}$ & $\widehat{A}_q$.

4. $\hat{\mathbf{s}} = (1 - \alpha)\mathbf{P^{-1}} [\mathbf{D_1^{-1}} + \alpha \mathbf{D_1^{-1}} \mathbf{L} (\mathbf{\Lambda^{-1}} - \alpha \mathbf{R} \mathbf{D_1^{-1}} \mathbf{L})^{-1}]\mathbf{P^{-1}} \mathbf{h}$

- Result:
  - Time complexity: $O(r^2 t^2 k n + rtkn + K^2 kn)$
  - Space complexity: $O(k^2 rn + m_1)$

By Sherman-Morrison Lemma

Revised Query

Data Network

DATA Lab

# FIRST-N: Handle Node Attribute Revision

**Scenario 2**: During query revising, only node attribute is changed.

- **Goal**:
  Fast computation of similarity vector after node attribute revision.
- **Observation**:
  - We already have: $\hat{\mathbf{s}} = (1 - \alpha)(\mathbf{I} - \alpha\widehat{\mathbf{W}})^{-1}\mathbf{h}, \widetilde{\mathbf{N}_q} = \mathbf{N}_q + \Delta\mathbf{N},$
  - The topology keeps unchanged.
  - Low-rank approx. of both $\mathbf{A}$ and $\mathbf{A}_q$ (pre-compute).
- **Solution**:
  - Calculate $\mathbf{W}$ in pre-computing stage, only $\Delta\mathbf{W}$ in interactive stage.
  - Low-rank approx. & matrix inverse lemma for fast computation.

Arizona State University

# FIRST-N: Handle Node Attribute Revision

- **How to (Details):**

Pre-computing Stage:

1. Low-rank approximation:

$$\mathbf{A} \approx \mathbf{U_A \Lambda_A U_A^T}$$
$$\mathbf{A_q} \approx \mathbf{U_q \Lambda_q U_q^T}$$

2. Construct approximated $\mathbf{W}$ :

$$\widehat{\mathbf{W}} = \boxed{\mathbf{L}} \boxed{\mathbf{\Lambda}} \boxed{\mathbf{R}}$$

3. Store $\mathbf{L, R, \Lambda}$.

> This can further speed up this algorithm compared to FIRST-Q.

Interactive Stage:

1. Major Computation: Compute $\mathbf{N}$, diagonal degree matrix $\mathbf{D}$ with $\widetilde{\mathbf{N_q}}$.

2. Intermediate matrix:
$$\mathbf{P} = \mathbf{D^{-1/2}}\widetilde{\mathbf{N_q}}, \mathbf{D_1} = \mathbf{P^{-1}P^{-1}}.$$

3. Calculate similarity vector:
$$\hat{\mathbf{s}} = (1 - \alpha)\mathbf{P^{-1}}[\mathbf{D_1}^{-1} + \alpha\mathbf{D_1}^{-1}\mathbf{L}$$
$$(\mathbf{\Lambda^{-1}} - \alpha\mathbf{R}\mathbf{D_1}^{-1}\mathbf{L})^{-1}]\mathbf{P^{-1}h}$$

- Result:
➢ Time complexity: $O(r^2t^2kn + rtkn + K^2kn)$
➢ Space complexity: $O(k^2rn + m_1)$

> By Sherman-Morrison Lemma

# FIRST-E: Handle Edge Attribute Revision

**Scenario 3**: During query revising, only edge attributes are changed.

- **Goal**:

    Fast computation of similarity vector after edge attribute revision.

- **Observation**:

    1. Pre-compute: The low-rank approximation of the <u>edge attributed</u> adjacency matrix ($\mathbf{E_A^l} \odot \mathbf{A}$);

    2. Interactive: Only approx. of the revised <u>edge attributed</u> adjacency matrix.

- **Solution keys:**

    1. Low-rank approximation;

    2. matrix inverse lemma;

    3. Block matrix property.

# FIRST-E: Handle Edge Attribute Revision (cont'd)

- **How to (Details):**

**Pre-computing Stage:**

1. The edge attribute batch 

$$\mathbf{E}_A^l \odot \mathbf{A} \approx$$

> This can further speed up computation.

2. Store $\mathbf{U}_A^l, \mathbf{\Lambda}_A^l$;

3. If the index **l'** of changed edge attribute is available:

$$\mathbf{E}_q^k \odot \mathbf{A}_q \approx \mathbf{U}_q^k \, \mathbf{\Lambda}_q^k \, (\mathbf{U}_q^k)^T \quad \textit{(all } k)$$

4. Store $\mathbf{U}_q^k, \mathbf{\Lambda}_q^k$;

**Interactive Stage:**

1. Compute: $\mathbf{N}, \mathbf{D}$. $\mathbf{E}_q^k \odot \mathbf{A}_q \approx \mathbf{U}_q^k \, \mathbf{\Lambda}_q^k \, (\mathbf{U}_q^k)^T$

$$(k \in \mathbf{l'})$$

2. Construct block matrix $\mathbf{U}, \mathbf{\Lambda}$;

$$\mathbf{U} = \begin{pmatrix} \| \| \cdots \| \end{pmatrix}, \mathbf{\Lambda} = \begin{pmatrix} \mathbf{Y}_1 & & & \\ & \mathbf{Y}_2 & & \\ & & \ddots & \\ & & & \mathbf{Y}_L \end{pmatrix}$$

$$\mathbf{V}_1, \mathbf{V}_2, \ldots, \mathbf{V}_L$$

$$\mathbf{V}_l = \mathbf{U}_A^l \otimes \mathbf{U}_q^k \qquad \mathbf{Y}_j = \mathbf{\Lambda}_A^l \otimes \mathbf{\Lambda}_q^k$$

3. $\hat{\mathbf{s}} = (1 - \alpha) \left[ \mathbf{I} + \alpha \mathbf{L} (\mathbf{\Lambda}^{-1} - \alpha \mathbf{R} \mathbf{L})^{-1} \mathbf{R} \right] \mathbf{h}$;

$$(\mathbf{L} = \mathbf{D}^{-1/2} \mathbf{N} \mathbf{U}, \mathbf{R} = \mathbf{U}^T \mathbf{N} \mathbf{D}^{-1/2}).$$

- Time complexity: $O(r^2 t^2 kn + Lrtkn + K^2 Lkn)$
- Space complexity: $O(Lrtkn + m_1)$

**DATA Lab**

# Roadmap:

➢ Motivation

➢ Problem Definition

➢ Proposed Solution: FIRST family

   ➢ Key ideas

   ➢ Three scenarios

➢ Experiments

➢ Conclusions

**DATA Lab**

# Experimental Setup

- Datasets Summary

| Name | # of Nodes | # of Edges | Node/Edge Attribute |
|------|-----------|-----------|---------------------|
| DBLP | 9,143 | 16,338 | Node attribute only |
| Flickr | 12,974 | 16,149 | Node attribute only |
| LastFm | 136,421 | 1,685,524 | Node attribute only |
| ArnetMiner | 1,274,360 | 4,756,194 | Node & Edge attribute |
| LinkedIn | 6,726,290 | 19,360,690 | Node attribute only |

- Baseline methods:
  - G-ray [Tong et al. KDD' 07]
  - MAGE [Pienta et al. IEEE BigData' 14]
  - FINAL (and its variants) [Zhang et al. KDD'16]

# Experiment Result - Effectiveness (Nodes)

Query samples:

## % Exact Matching Nodes (Higher is Better)

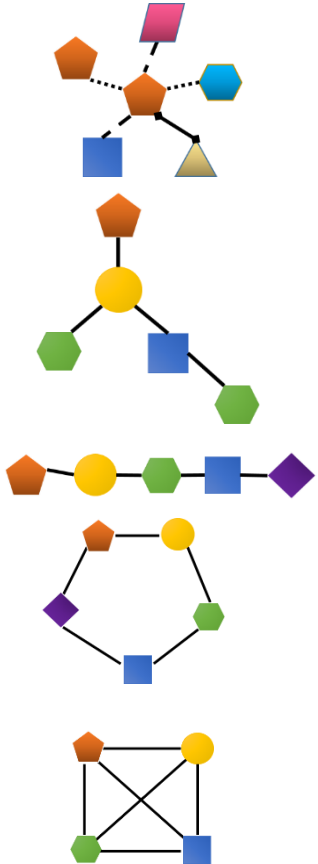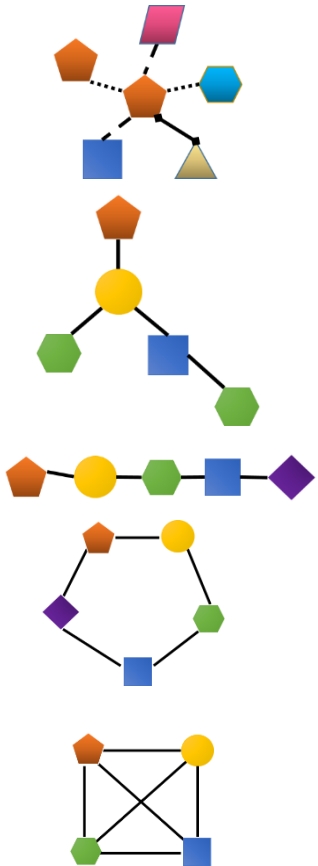| Algorithms | G-Ray | MAGE | FIRST |
|---|---|---|---|
| Star(N) | 37.5 | * | **75.0** |
| E-Star(N) | **83.3** | * | 71.4 |
| Line(N) | 50.0 | * | **83.3** |
| Loop(N) | 27.3 | * | **71.4** |
| Clique(N) | 25.0 | * | **57.1** |
| Star(NE) | * | 30.0 | **40.0** |
| E-Star(NE) | * | 33.3 | **41.7** |
| Line(NE) | * | 33.3 | **62.5** |
| Loop(NE) | * | 27.3 | **33.3** |
| Clique(NE) | * | 60.0 | **66.7** |

* = Not applicable

Observation: Generally FIRST family generates more accurate results.

# Experiment Result - Effectiveness (Nodes)

Query samples:

## % of Extra Nodes (Lower is better)



| Algorithms | G-Ray | MAGE | FIRST |
|---|---|---|---|
| Star(N) | 62.5 | * | **0.0** |
| E-Star(N) | **0.0** | * | **0.0** |
| Line(N) | 50.0 | * | **0.0** |
| Loop(N) | **0.0** | * | **0.0** |
| Clique(N) | 25.0 | * | **0.0** |
| Star(NE) | * | 50.0 | **0.0** |
| E-Star(NE) | * | **0.0** | **0.0** |
| Line(NE) | * | 33.3 | **0.0** |
| Loop(NE) | * | **27.3** | 44.4 |
| Clique(NE) | * | 40.0 | **0.0** |

\* = Not applicable

Observation: Generally FIRST family generates more accurate results.

DATA Lab

# Experiment Result - Effectiveness (Edges)
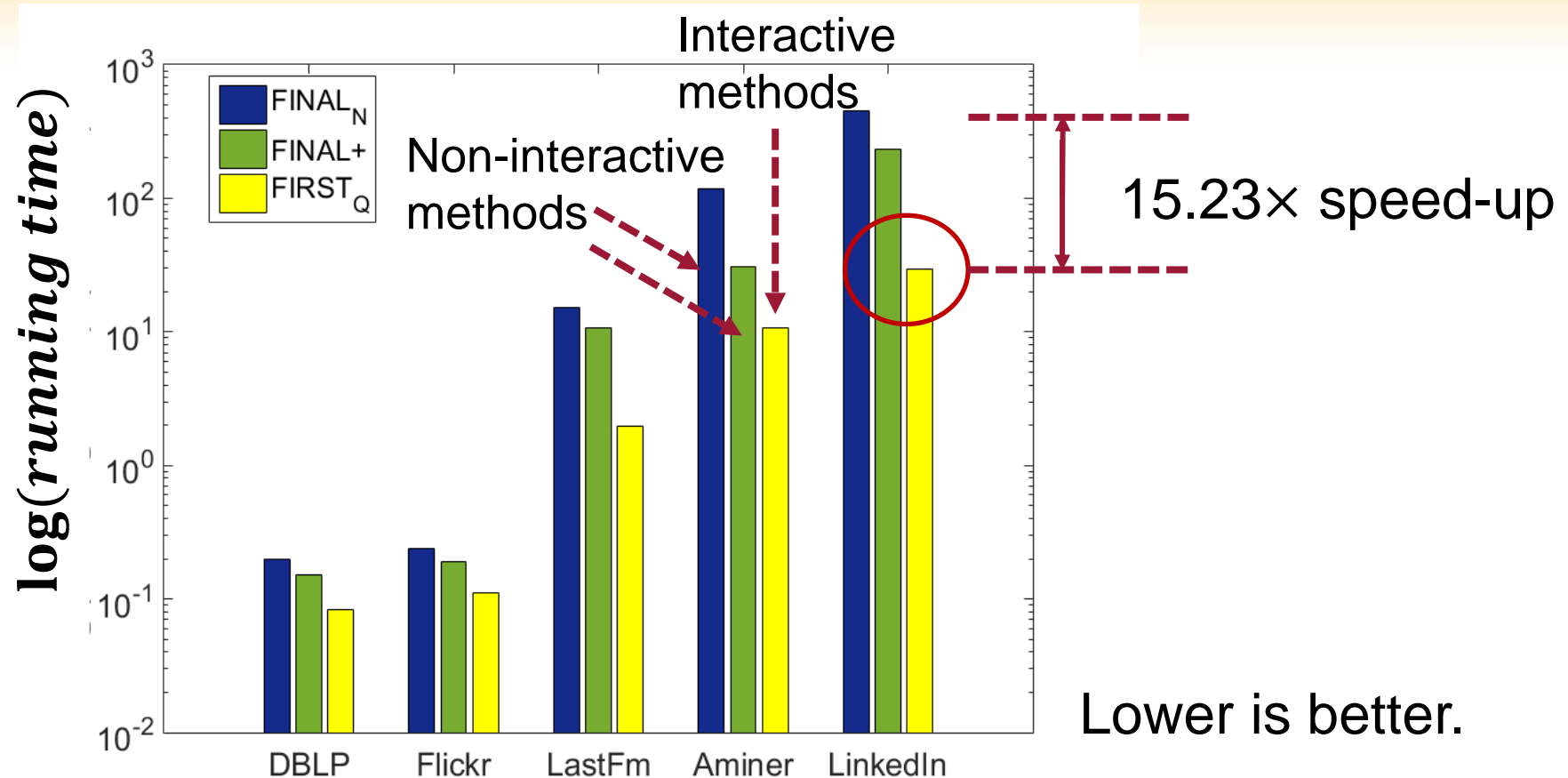
Query samples:

% Exact Matching Edges (Higher is better)

| Algorithm | G-Ray | MAGE | FIRST |
|-----------|-------|------|-------|
| Star(N) | 33.3 | * | **57.1** |
| E-Star(N) | **60.0** | * | 50.0 |
| Line(N) | 40.0 | * | **60.0** |
| Loop(N) | 8.3 | * | **42.9** |
| Clique(N) | 7.1 | * | **12.5** |
| Star(NE) | * | 0.0 | **14.3** |
| E-Star(NE) | * | 7.1 | **9.0** |
| Line(NE) | * | 0.0 | **14.3** |
| Loop(NE) | * | 0.0 | **14.3** |
| Clique(NE) | * | 0.0 | **12.5** |

\* = Not applicable

Observation: Generally FIRST family generates more accurate results.
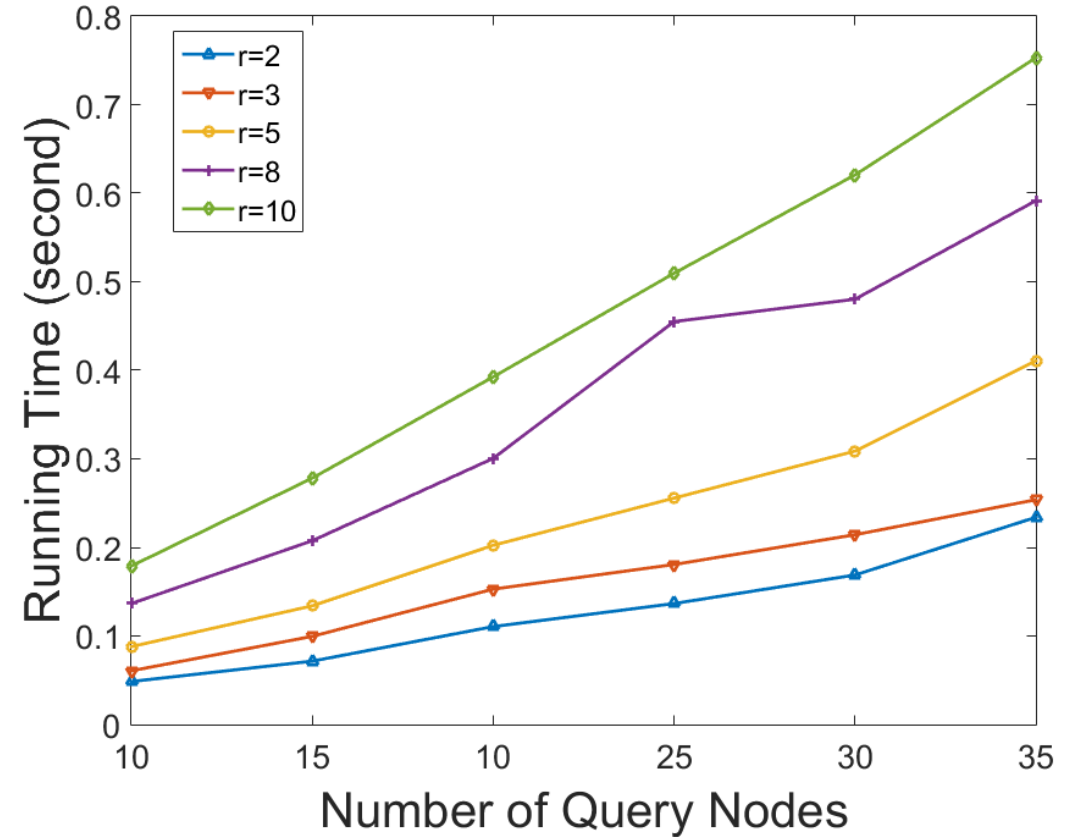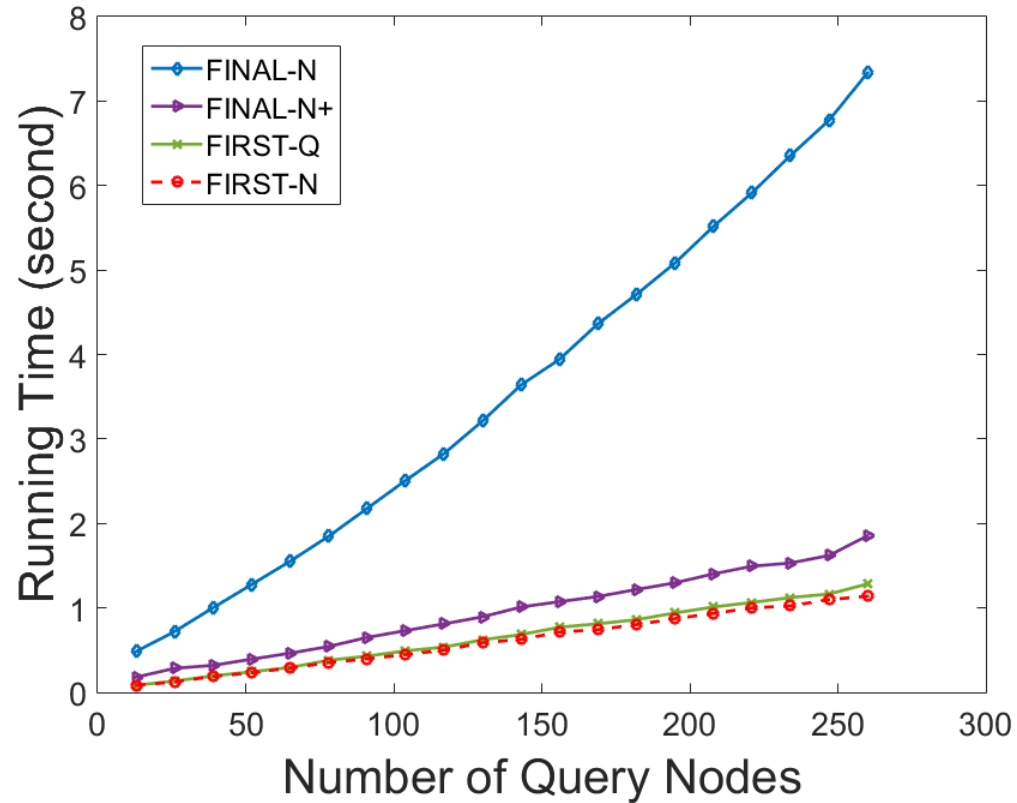
# Experiment Result - Effectiveness

Query samples:                % Extra Matching Edges (Lower is better)

| Algorithm | G-Ray | MAGE | FIRST |
|-----------|-------|------|-------|
| Star(N) | 66.7 | * | **0.0** |
| E-Star(N) | **0.0** | * | **0.0** |
| Line(N) | 60.0 | * | **0.0** |
| Loop(N) | 8.3 | * | **0.0** |
| Clique(N) | 35.7 | * | **0.0** |
| Star(NE) | * | 42.9 | **0.0** |
| E-Star(NE) | * | **0.0** | **0.0** |
| Line(NE) | * | 27.3 | **0.0** |
| Loop(NE) | * | **30.0** | 42.9 |
| Clique(NE) | * | 33.3 | **0.0** |

\* = Not applicable

Observation: Generally FIRST family generates more accurate results.

**DATA Lab**

- Observation: >15× speedup with 6,726,290-node data network.

- Observation: FIRST family is more efficient than baseline methods.
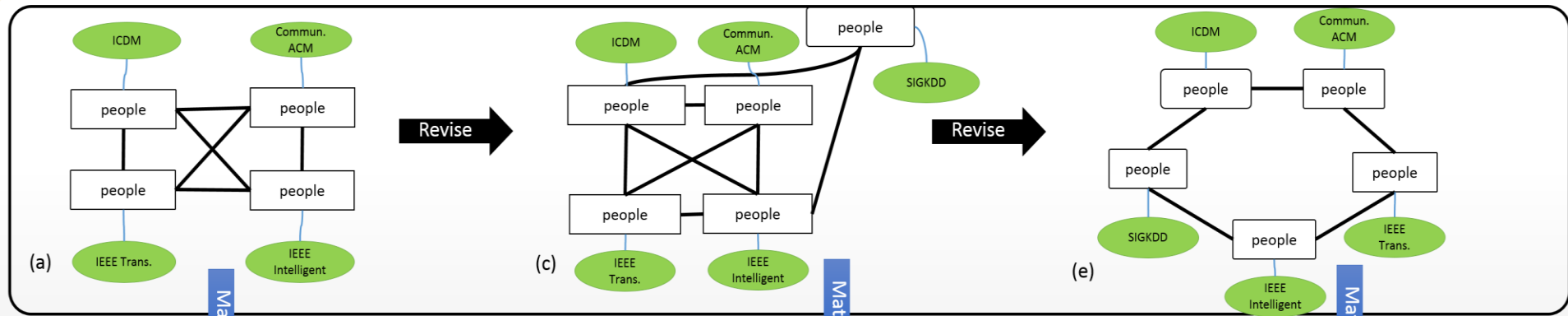
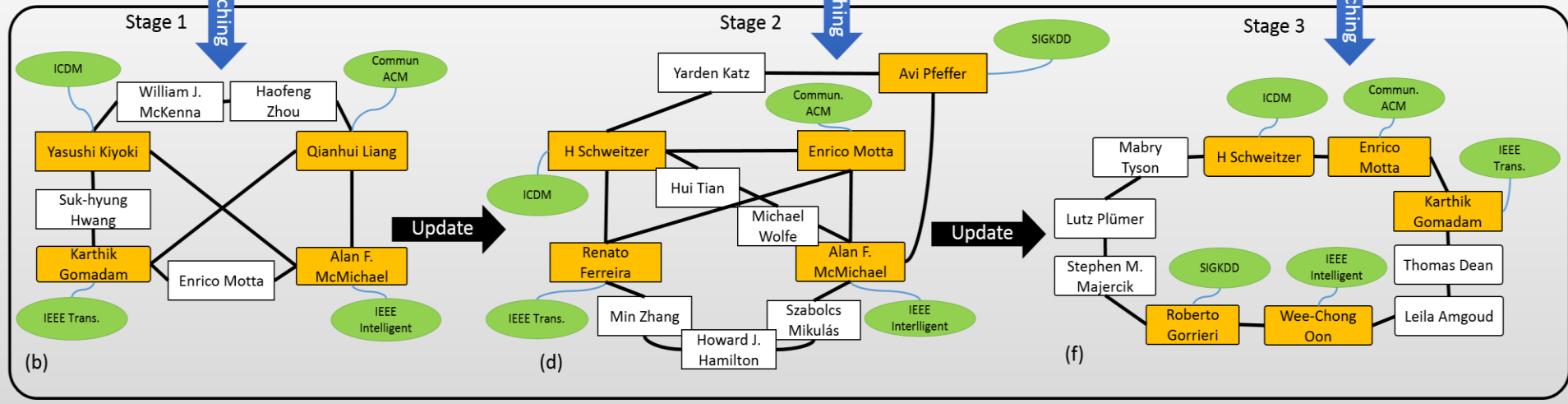# Experiment Result- Efficiency



- Observation: FIRST family scales linearly with regard to size of query graph.

# Experiments – Case Studies (on DBLP)

**Query graphs:**



**Matching subgraphs:**



Node attribute (conferences)

Extra/intermediate node (in matching subgraph)

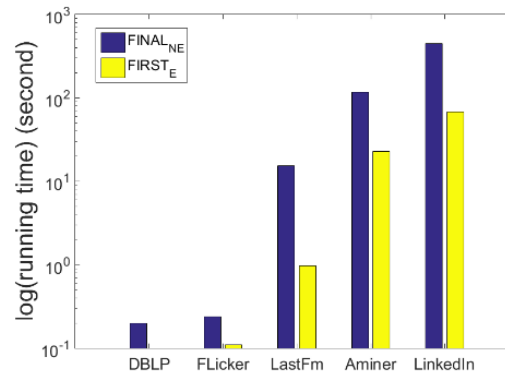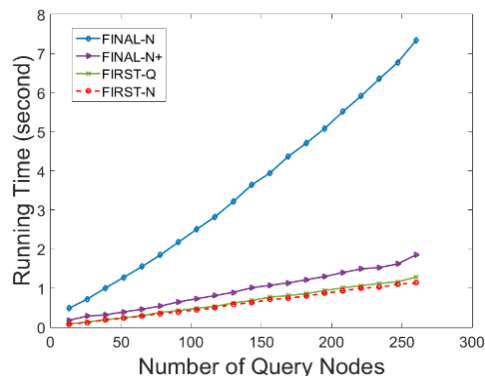people — Query node

Name — Matching node

# Roadmap:

- ➢ Motivation
- ➢ Problem Definition
- ➢ Proposed Solution: FIRST family
  - ➢ Key ideas
  - ➢ Three scenarios
- ➢ Experiments
  - ➢ Setup
  - ➢ Results
  - ➢ Case Study
- ➢ Conclusions

DATA Lab

# Conclusions

- **Goal**: **Efficient** Methods for **Interactive** Attributed Subgraph Matching.

- **Solution**: **FIRST family**
  - ➢ Key Idea #1: Subgraph matching as cross-network node similarity
  - ➢ Key Idea #2: Explore the smoothness of query graphs

- **Results:**
  - ➢ Linear scalability w.r.t the size of data network/query;
  - ➢ Better quality of matching subgraph against baselines.



| | % Exact Matching Nodes | | |
|---|---|---|---|
| Algorithm | G-Ray | MAGE | FIRST |
| Star(N) | 37.5 | * | **75.0** |
| E-Star(N) | **83.3** | * | 71.4 |
| Line(N) | 50.0 | * | **83.3** |
| Loop(N) | 27.3 | * | **71.4** |
| Clique(N) | 25.0 | * | **57.1** |
| Star(NE) | * | 30.0 | **40.0** |
| E-Star(NE) | * | 33.3 | **41.7** |
| Line(NE) | * | 33.3 | **62.5** |
| Loop(NE) | * | 27.3 | **33.3** |
| Clique(NE) | * | 60.0 | **66.7** |