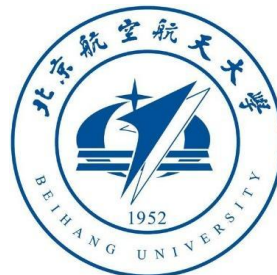# MiniClean: A Single-Machine System for Cleaning Big Graphs
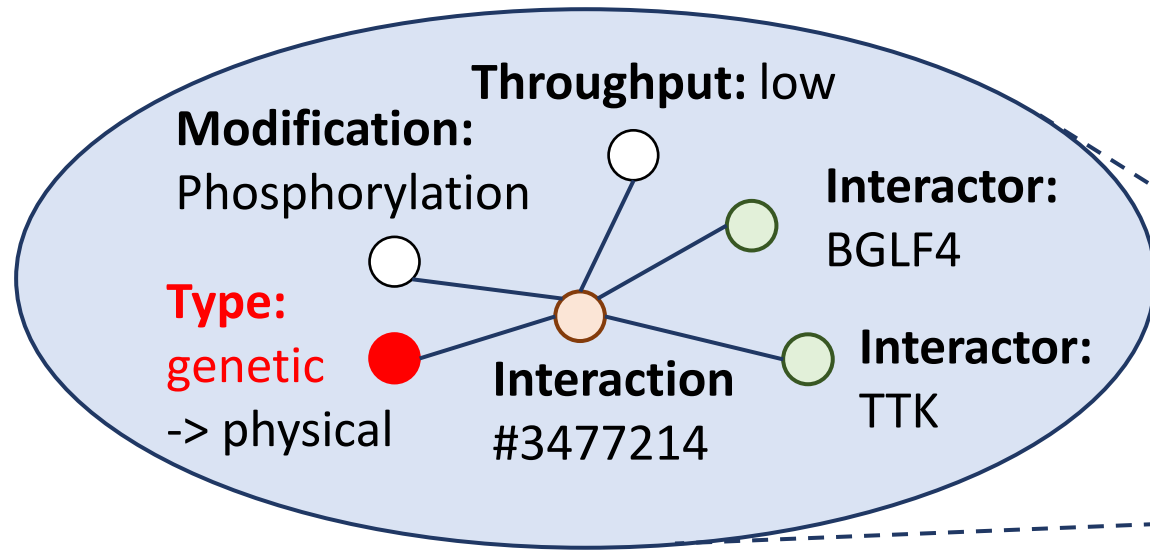
Wenchao Bai, Wenfei Fan, Jiahui Jin, Daji Li, Jian Li, Shuhao Liu, Mingliang Ouyang, Qiang Yuan

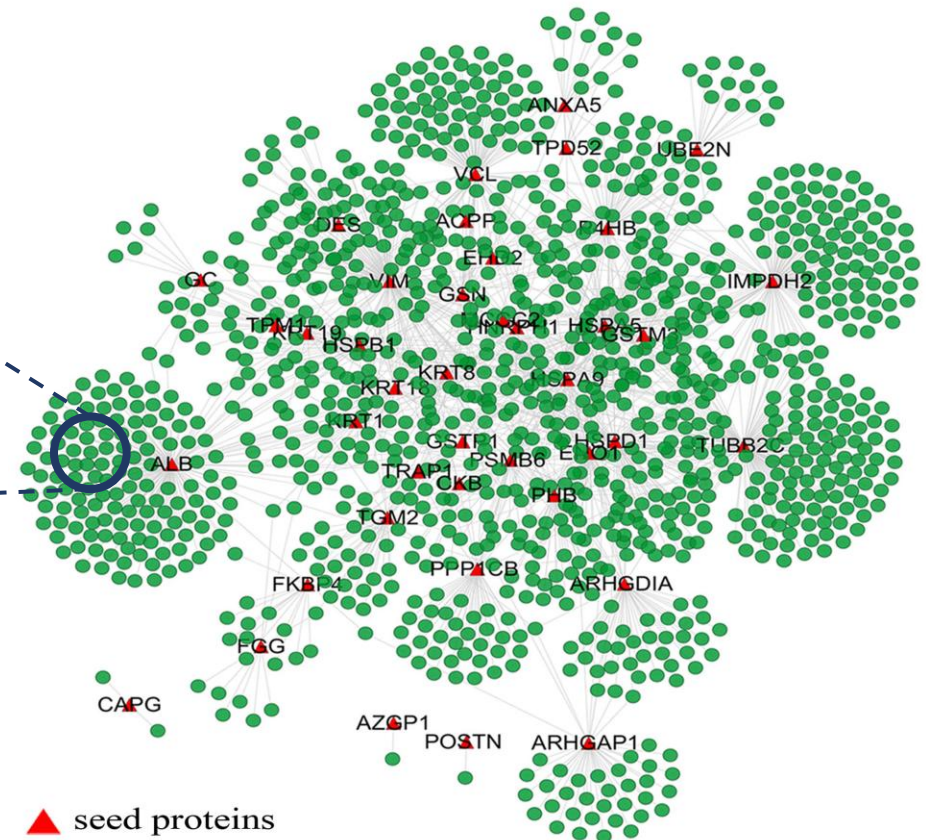Shenzhen Institute of Computing Science, Southeast University, University of Edinburgh, Beihang University

- **Introduction**:   00'00 – 08'55
- **Live Demo**:   08'55 – 12'21
- **Evaluation**:   12'21 – 13'38

**Throughput:** low

**Modification:** Phosphorylation

**Type:** genetic -> physical

**Interaction** #3477214

**Interactor:** BGLF4

**Interactor:** TTK

**A mislabeled instance in BioGRID[2]**

▲ seed proteins

[1] Chen, Chen, et al. "Construction and analysis of protein-protein interaction networks based on proteomics data of prostate cancer." *International journal of molecular medicine* 37.6 (2016): 1576-1586.
[2] BioGRID | Database of Protein, Chemical, and Genetic Interactions

**Modification:** Phosphorylation

**Throughput:** low

**Interactor:** BGLF4

**Type:** genetic -> physical

**Interaction** #3477214

**Interactor:** TTK

**A mislabeled instance in BioGRID[2]**

- *Errors (e.g., duplicates, conflicts) are common.*
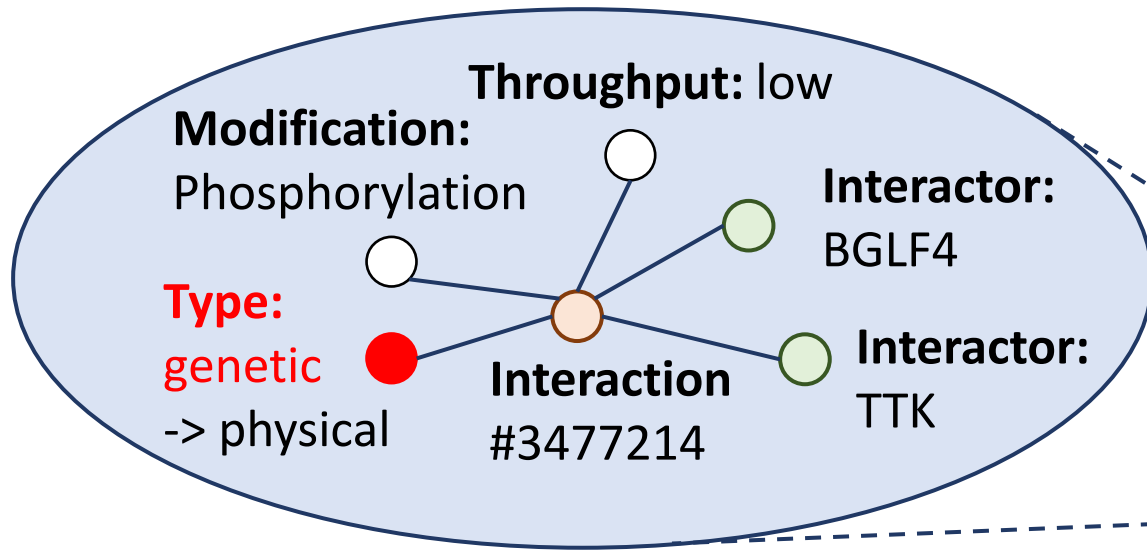- *Errors can misguide data-driven decision making.*

▲ seed proteins

[1] Chen, Chen, et al. "Construction and analysis of protein-protein interaction networks based on proteomics data of prostate cancer." *International journal of molecular medicine* 37.6 (2016): 1576-1586.
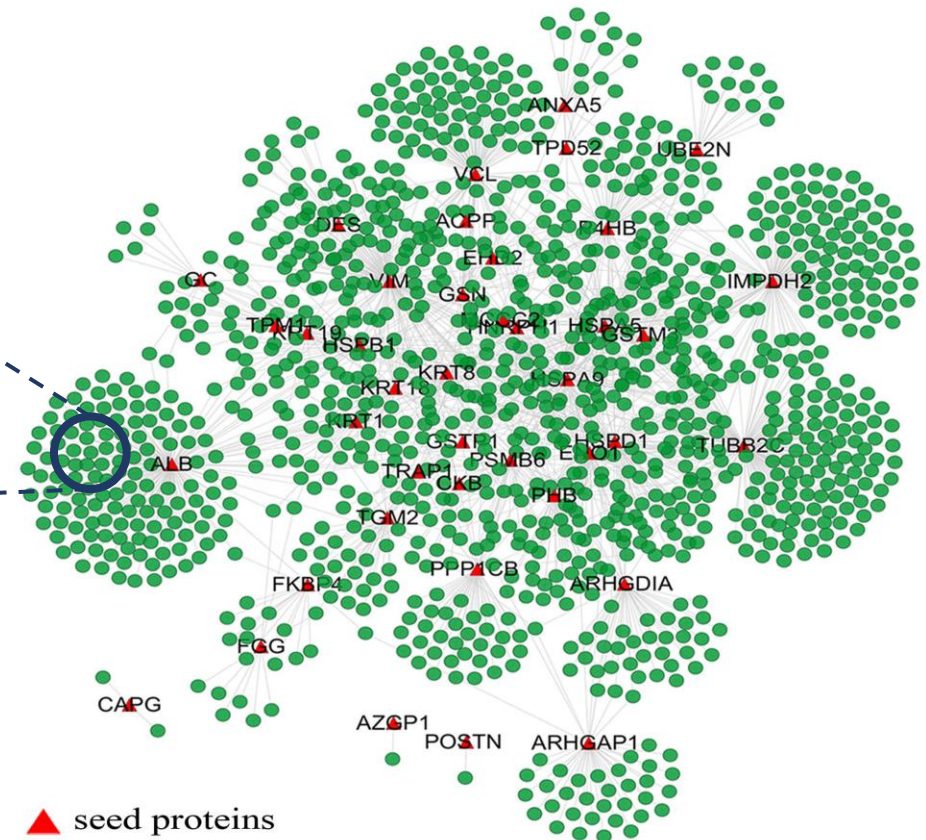[2] BioGRID | Database of Protein, Chemical, and Genetic Interactions

*Graph cleaning is crucial*

- **Strawman 1**: ML models (e.g., Ditto, KGClean).
  - Transform graphs into embeddings.
  - Detect errors (duplicates/conflicts) via binary classification.
  - Predictions are probabilistic and hard to explain.
- **Strawman 2**: Logic rules.
  - Apply rules via pattern matching and predicate verification.
  - Deduce dependencies to detect and correct errors.
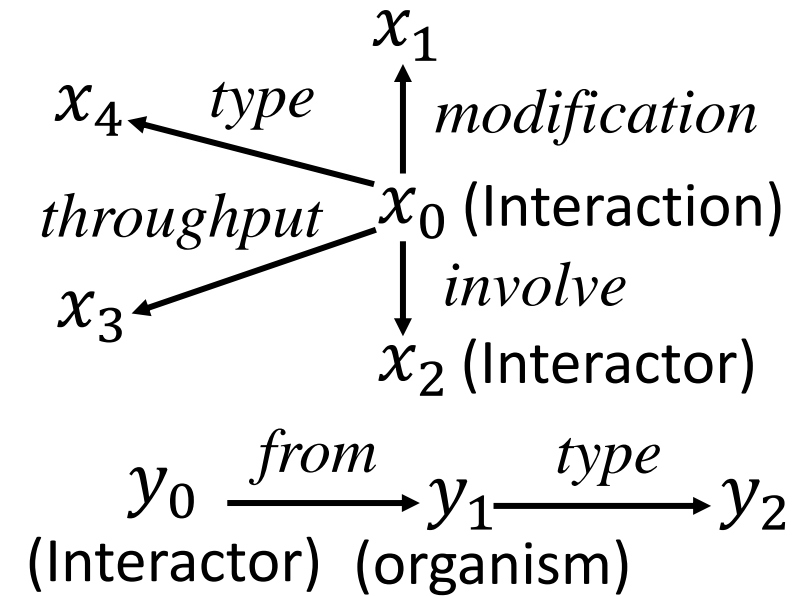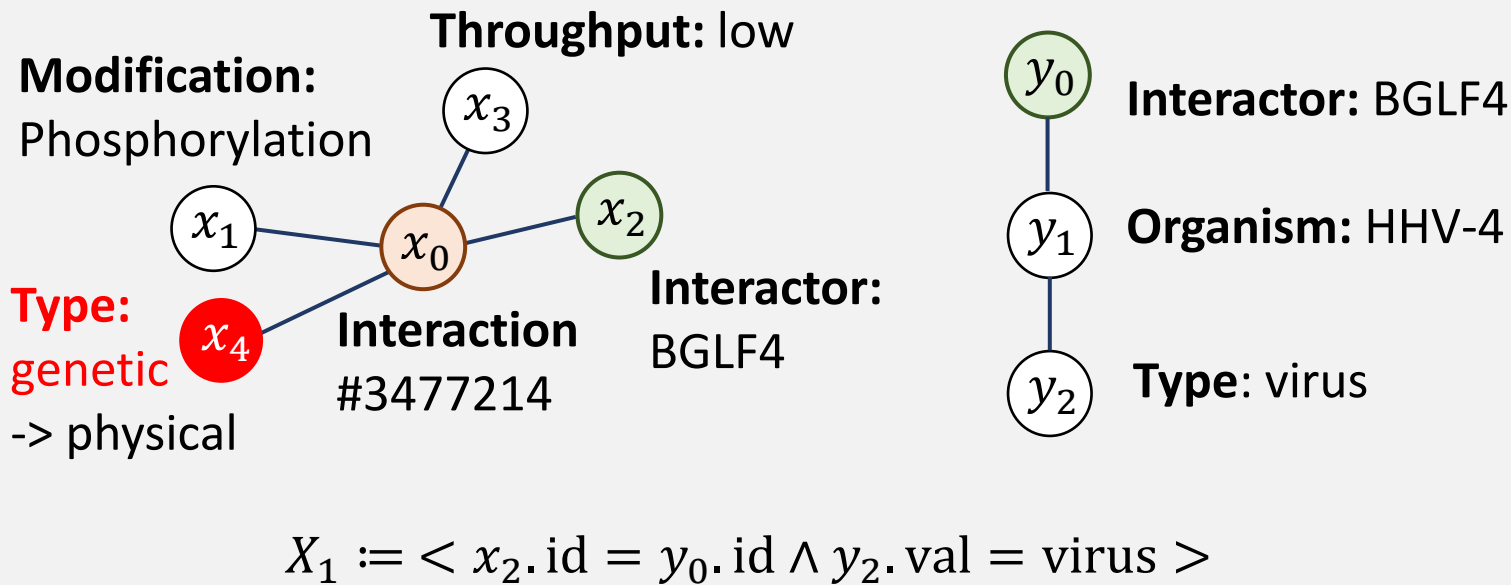  - Difficult to find rules that cover all the cases.

- **Strawman 1**: ML models (e.g., Ditto, KGClean).

  - Transform graphs into embeddings.

  - Detect errors (duplicates/conflicts) via binary classification.

  - Predictions are probabilistic and hard to explain.

- **Strawman 2**: Logic rules.

  - Apply rules via pattern matching and predicate verification.

  - Deduce dependencies to detect and correct errors.

  - Difficult to find rules that cover all the cases.

- **Our Solution**: ML models + Logic rules.

  - $M_1(x_0, x_1) > \delta_1 \land X_1 \rightarrow x_4.\text{val} = \text{physical}$

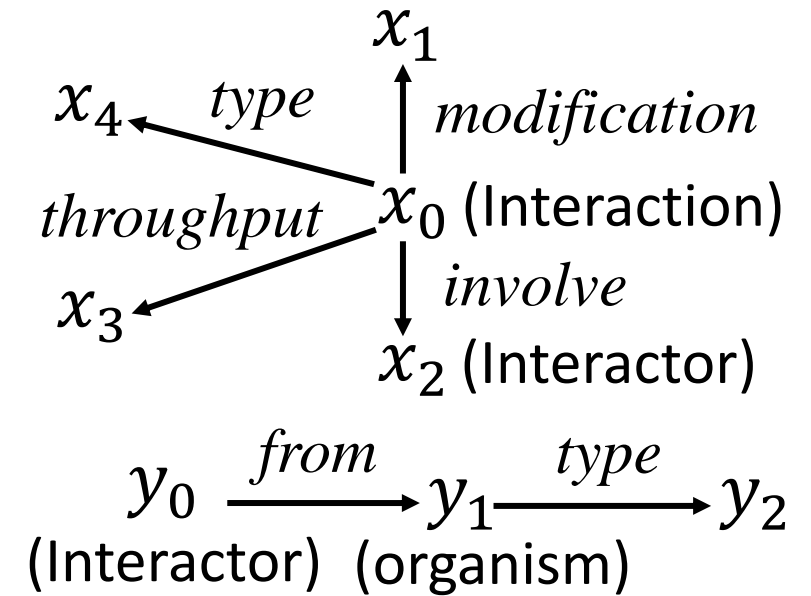  - $x_4.\text{val} = \text{genetic} \land X_2 \rightarrow M_2(x_0, x_3) > \delta_2$

**Modification:**
Phosphorylation

**Throughput:** low

$x_3$

$x_1$

$x_0$

$x_2$

**Interactor:**
BGLF4

**Type:**
genetic
-> physical

$x_4$

**Interaction**
#3477214

$y_0$

**Interactor:** BGLF4

$y_1$

**Organism:** HHV-4

$y_2$

**Type:** virus

$$X_1 := \; < x_2.\mathrm{id} = y_0.\mathrm{id} \land y_2.\mathrm{val} = \mathrm{virus} >$$

$x_1$

$x_4 \xleftarrow{type}$

$modification$

$throughput$

$x_0$ (Interaction)

$x_3$

$involve$

$x_2$ (Interactor)

$y_0 \xrightarrow{from} y_1 \xrightarrow{type} y_2$
(Interactor)  (organism)

- **Our Solution**: ML models + Logic rules.

  - $M_1(x_0, x_1) > \delta_1 \land X_1 \rightarrow x_4.\mathrm{val} = \mathrm{physical}$

  - $x_4.\mathrm{val} = \mathrm{genetic} \land X_2 \rightarrow M_2(x_0, x_3) > \delta_2$

**Throughput:** low

**Modification**

$x_1$

**Type:** genetic

$y_2$

$x_3$

$x_0$

**Interaction** #343119

$x_2$

**Interactor:** srmB

$y_0$ **Interactor:** srmB

$y_1$ **Organism:** K12/MG1655

$y_2$ **Type:** coli

$$X_2 := < x_2.\text{id} = y_0.\text{id} \wedge y_2.\text{val} = \text{coli} >$$

$x_1$

$x_4 \xleftarrow{type}$ *modification*

*throughput* $x_0$ (Interaction)

$x_3$ *involve*

$x_2$ (Interactor)

$y_0 \xrightarrow{from} y_1 \xrightarrow{type} y_2$
(Interactor) (organism)

- **Our Solution**: ML models + Logic rules.

  - $M_1(x_0, x_1) > \delta_1 \wedge X_1 \rightarrow x_4.\text{val} = \text{physical}$

  - $x_4.\text{val} = \text{genetic} \wedge X_2 \rightarrow M_2(x_0, x_3) > \delta_2$
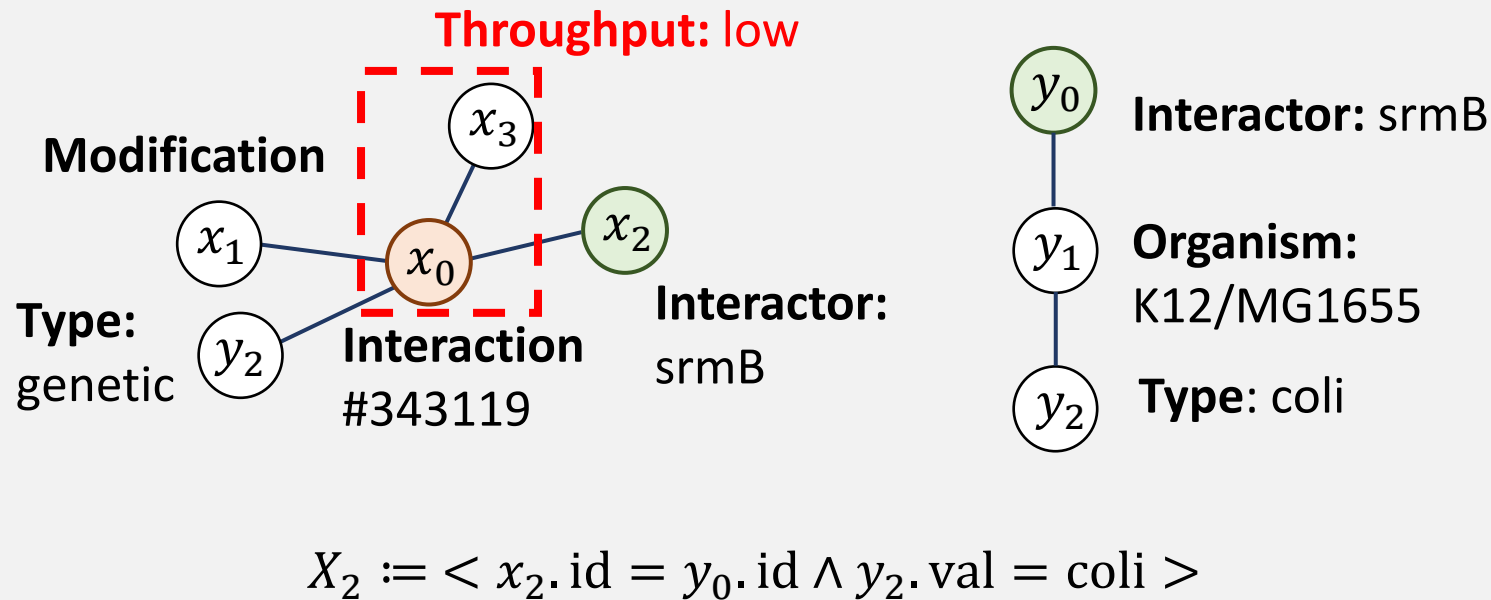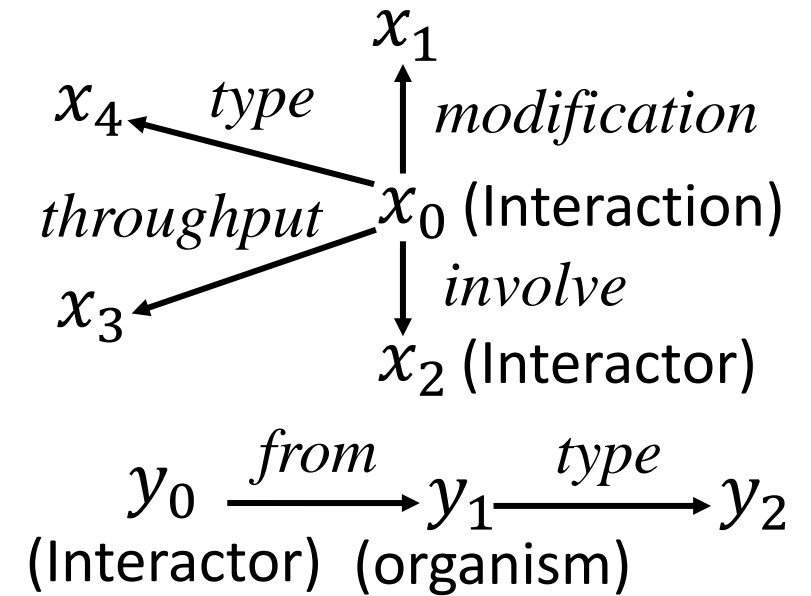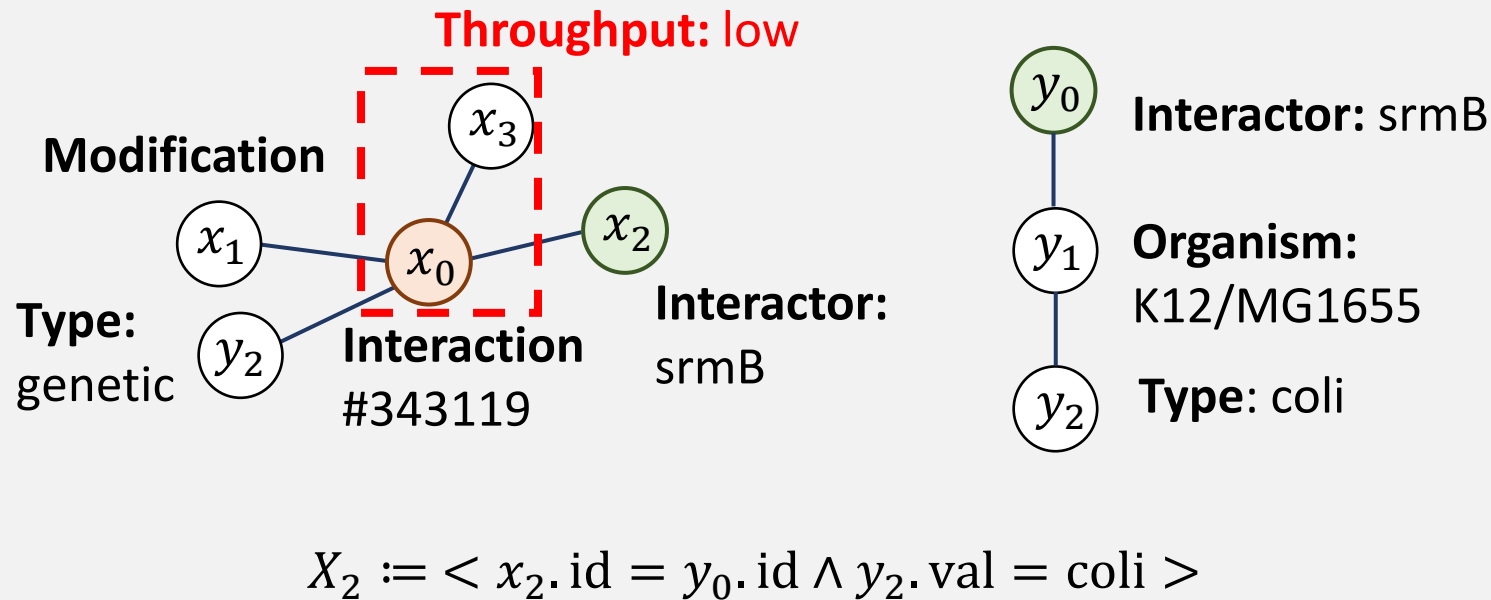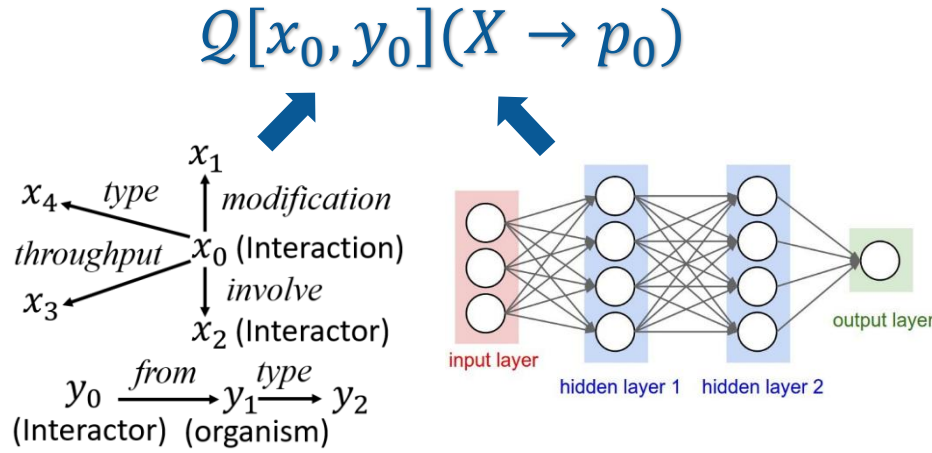
# Approaches to Graph Cleaning



$$X_2 := < x_2.\text{id} = y_0.\text{id} \land y_2.\text{val} = \text{coli} >$$

- **Our Solution**: ML models + Logic rules.
  - $M_1(x_0, x_1) > \delta_1 \land X_1 \rightarrow x_4.\text{val} = \text{physical}$
  - $x_4.\text{val} = \text{genetic} \land X_2 \rightarrow M_2(x_0, x_3) > \delta_2$

- *Generalizability (ML)*
- *Reliability (Logic)*
- *Explainability (Logic)*

$Q[x_0, y_0](X \rightarrow p_0)$

- ✓ $Q[x_0, y_0]$: dual star pattern
- ✓ $X \rightarrow p_0$: dependency
- ✓ $X$: precondition (conjunction of predicates)
- ✓ $p_0$: consequence (predicate)

Dependency deduction:

- ML predicates $M(x.\bar{A}, y.\bar{B})$: for ML classification & regression
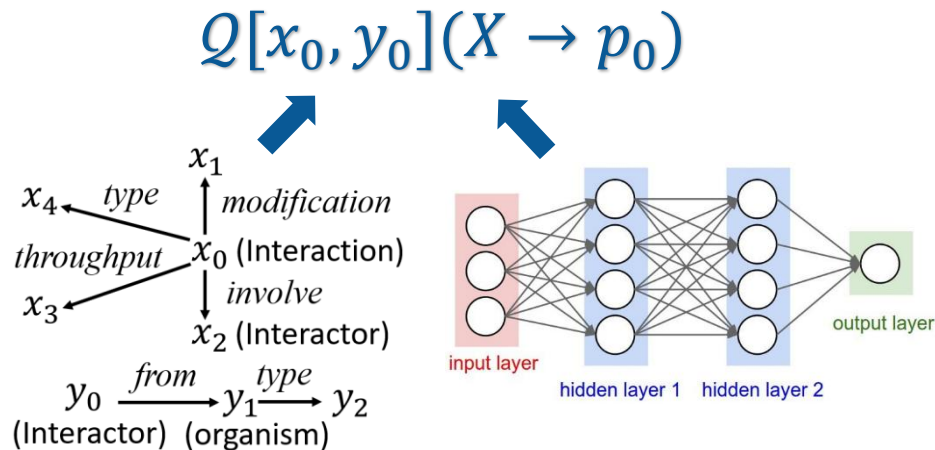- Variable and constant predicates $x.A \oplus y.B, x.A \oplus c$: for value associations

- ✓ $Q[x_0, y_0]$: dual star pattern
- ✓ $X \rightarrow p_0$: dependency
- ✓ $X$: precondition (conjunction of predicates)
- ✓ $p_0$: consequence (predicate)
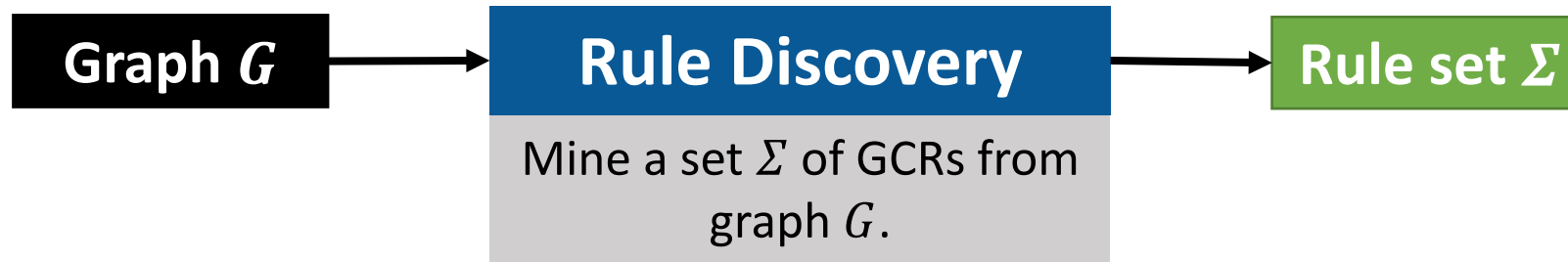
Dependency deduction:

- ML predicates $M(x.\bar{A}, y.\bar{B})$: for ML classification & regression

- Variable and constant predicates $x.A \oplus y.B, x.A \oplus c$: for value associations

**Functionalities**
- ✓ *Conflict resolution*
- ✓ *Entity resolution*
- ✓ *ML explanation*

**Graph $G$** → **Rule Discovery**

Mine a set $\Sigma$ of GCRs from graph $G$.

→ **Rule set $\Sigma$**

**Graph $G$** → **Rule Discovery**

Mine a set $\Sigma$ of GCRs from graph $G$.

→ **Rule set $\Sigma$**

**Error Detection**

Identify violations of any GCR of $\Sigma$ in graph $G$.

→ **Errors (duplicates /conflicts)**

**Graph** $G$ → **Rule Discovery**

Mine a set $\Sigma$ of GCRs from graph $G$.

→ **Rule set** $\Sigma$

**Corrections**

**Error Correction**

Correct identified violations based on GCRs in $\Sigma$

**Error Detection**

Identify violations of any GCR of $\Sigma$ in graph $G$.

**Cleaned Graph**

**No error is detected**

**Errors (duplicates /conflicts)**

## Single-machine solutions are often preferred:

- **Physical constraints on deployed hardware**.

  - A cluster on premise may not be feasible, e.g., at an edge location.

- **A cloud deployment is not viable**.

  - Prohibitive cost: network bandwidth at $0.05–0.09 per GB[1].

  - Privacy concerns.

[1] https://aws.amazon.com/pricing

- **Computation heavy**.

  - > 9h processing time on a 32-node cluster.

- **Excessive intermediate data**.

  - > 150GB intermediate result, way exceeding the memory capacity of a typical machine.

- **Parallel model for maximum resource utilization**.

  - Idle caused by I/O, data transfer, and task dependencies.

$$Q_l[x_0, \bar{x}_l](X_l, F_{Q_l})$$

$$Q_r[x_0, \bar{x}_r](X_r, F_{Q_r})$$

**Intersection** $\otimes$

$$Q_{l*r}[x_0, \bar{x}_l \cap \bar{x}_r](X_l \cap X_r, F_{Q_{l*r}})$$

**Union** $\oplus$

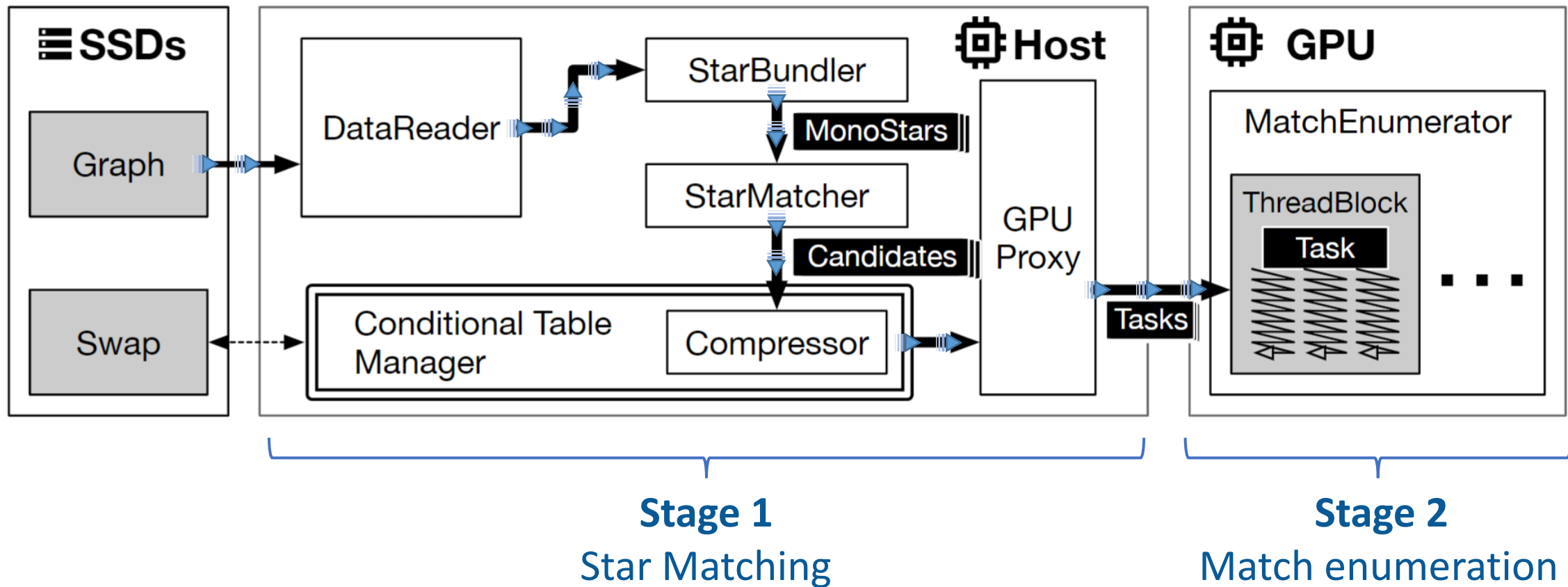$$Q_{l+r}[x_0, \bar{x}_l \cup \bar{x}_r](\emptyset, F_{Q_{l+r}})$$

✓ **Motivation**: Repetitive computation by matching common substructures.

✓ **Solution**: Bundle similar patterns into a group and match them together.

Conditional succinct matches

| Matches | | | Conditions | | |
|---|---|---|---|---|---|
| $x_0.\text{id}$ | $x_1.\text{val}$ | $x_2.\text{val}$ | $x_4.\text{val}$ | $x_0.\text{title}$ | $x_5.\text{val}$ |
| $u_3$ | OSDI | CS | {J, S} | MR | null |
| $u_4$ | OSDI | CS | {J, S} | GFS | null |
| $u_1$ | OSDI | null | {M, P, J, S} | TF | 2016 |
| $u_2$ | OSDI | null | {M, P, J, S} | tf | 2016 |

Unfolded matches

| $x_0.\text{id}$ | ... | $x_4.\text{val}$ | ... |
|---|---|---|---|
| $u_3$ | ... | J | ... |
| $u_3$ | ... | S | ... |
| ... | ... | ... | ... |

Filtered out by $x_2.\text{val} = \text{CS}$

✓ **Motivation**: Intermediate data has up to $O(|V|^{|Q|})$ materialized candidates.

✓ **Solution**: Conditional succinct table for compression.

**Stage 1**
Star Matching

**Stage 2**
Match enumeration

(1) Load the graph and configurations

**Stage 1**
Star Matching

**Stage 2**
Match enumeration

*Maximally overlapping I/O, CPU and GPU operations*

(2) Bundle star patterns and match them altogether

**Stage 1**
Star Matching

**Stage 2**
Match enumeration

(3) Compress the matches and swap them with SSDs if necessary.

**Stage 1**
Star Matching

**Stage 2**
Match enumeration

**Stage 1**
Star Match...

*(4) Launch GPU kernels for match enumeration*

**Stage 2**
Match enumeration

In our demonstration, we will walk the participants through

- **Data management**: visualize the input graph and GCRs;
- **GCR discovery**: configure and launch a rule discovery task;
- **GCR application**: inspect the errors detected/corrected by each GCR.

*We will use BioGRID, a real PPI graph in our demonstration.*

| System | BioGRID | | | SemScholar | | |
|---|---|---|---|---|---|---|
| | Time | ER-F1(%) | CR-F1(%) | Time | ER-F1(%) | CR-F1(%) |
| MiniClean | 312.7s | 97.5 | 97.2 | 4089.1s | 94.6 | 70.6 |
| Ditto | 7.7× | 91.0 | N/A | 4.3× | 90.3 | N/A |
| KGClean | 11.9× | N/A | 54.9 | 11.9× | N/A | 29.8 |

- **ER tasks**: outperforms ER model *Ditto* by 4.3%--6.5%.

- **CR tasks**: outperforms CR model *KGClean* by 40.4%--42.3%.

- **Efficiency**: 11.9× faster than ML baselines.

| System | BioGRID | | | SemScholar | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Time | ER-F1(%) | CR-F1(%) | Time | ER-F1(%) | CR-F1(%) |
| MiniClean | 312.7s | 97.5 | 97.2 | 4089.1s | 94.6 | 70.6 |
| Ditto | 7.7× | 91.0 | N/A | 4.3× | 90.3 | N/A |
| KGClean | 11.9× | N/A | 54.9 | 11.9× | N/A | 29.8 |

$$\Delta = 6.5\% \qquad \Delta = 4.3\%$$

- **ER tasks**: outperforms ER model *Ditto* by 4.3%--6.5%.

- **CR tasks**: outperforms CR model *KGClean* by 40.4%--42.3%.

- **Efficiency**: 11.9× faster than ML baselines.

| System | BioGRID | | | SemScholar | | |
|---|---|---|---|---|---|---|
| | Time | ER-F1(%) | CR-F1(%) | Time | ER-F1(%) | CR-F1(%) |
| MiniClean | 312.7s | 97.5 | 97.2 | 4089.1s | 94.6 | 70.6 |
| Ditto | 7.7× | 91.0 | N/A | 4.3× | 90.3 | N/A |
| KGClean | 11.9× | N/A | 54.9 | 11.9× | N/A | 29.8 |

$\Delta = 42.3\%$      $\Delta = 40.8\%$

- **ER tasks**: outperforms ER model *Ditto* by 4.3%--6.5%.

- **CR tasks**: outperforms CR model *KGClean* by 40.8%--42.3%.

- **Efficiency**: 11.9× faster than ML baselines.

| System | BioGRID | | | SemScholar | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Time | ER-F1(%) | CR-F1(%) | Time | ER-F1(%) | CR-F1(%) |
| MiniClean | 312.7s | 97.5 | 97.2 | 4089.1s | 94.6 | 70.6 |
| Ditto | 7.7× | 91.0 | N/A | 4.3× | 90.3 | N/A |
| KGClean | 11.9× | N/A | 54.9 | 11.9× | N/A | 29.8 |

- **ER tasks**: outperforms ER model *Ditto* by 4.3%--6.5%.

- **CR tasks**: outperforms CR model *KGClean* by 40.4%--42.3%.

- **Efficiency**: 11.9× faster than ML baselines.

| Single-machine (BioGRID) | | Scalability (SemScholar) | | Ablation study (SemScholar) | |
|---|---|---|---|---|---|
| **System** | **Time (s)** | **System** | **Time (s)** | **System** | **Time (s)** |
| Blaze | OOM | GCRClean-32-Node | 8.09× | noBundle | 1.35× |
| MiniGraph | 65.34× | MiniClean-2-GPU | 0.64× | noPipelinedPar | 1.26× |
| Hyperblocker | 11.92× | MiniClean-4-GPU | 0.48× | noIndPar | 1.12× |
| MiniClean-1-GPU | 259.6s | MiniClean-1-GPU | 4251.08s | MiniClean-2-GPU | 2993.7s |

- **Single-machine**: 65.34× faster than **MiniGraph**; 11.92× faster than **Hyperblocker**.

- **Multi-machine**: 8.09× faster than 32-node **GCRClean**.

- **Scalability**: speedup by 1.56× with **2 GPUs**; 2.08× with **4 GPUs**.

| Single-machine (BioGRID) | | Scalability (SemScholar) | | Ablation study (SemScholar) | |
|---|---|---|---|---|---|
| **System** | **Time (s)** | **System** | **Time (s)** | **System** | **Time (s)** |
| Blaze | OOM | GCRClean-32-Node | 8.09× | noBundle | 1.35× |
| MiniGraph | 65.34× | MiniClean-2-GPU | 0.64× | noPipelinedPar | 1.26× |
| Hyperblocker | 11.92× | MiniClean-4-GPU | 0.48× | noIndPar | 1.12× |
| MiniClean-1-GPU | 259.6s | MiniClean-1-GPU | 4251.08s | MiniClean-2-GPU | 2993.7s |

- **Single-machine**: 65.34× faster than **MiniGraph**; 11.92× faster than **Hyperblocker**.

- **Multi-machine**: 8.09× faster than 32-node **GCRClean**.

- **Scalability**: speedup by 1.56× with **2 GPUs**; 2.08× with **4 GPUs**.

| Single-machine (BioGRID) | | Scalability (SemScholar) | | Ablation study (SemScholar) | |
|---|---|---|---|---|---|
| **System** | **Time (s)** | **System** | **Time (s)** | **System** | **Time (s)** |
| Blaze | OOM | GCRClean-32-Node | 8.09× | noBundle | 1.35× |
| MiniGraph | 65.34× | MiniClean-2-GPU | 0.64× | noPipelinedPar | 1.26× |
| Hyperblocker | 11.92× | MiniClean-4-GPU | 0.48× | noIndPar | 1.12× |
| MiniClean-1-GPU | 259.6s | MiniClean-1-GPU | 4251.08s | MiniClean-2-GPU | 2993.7s |

- **Single-machine**: 65.34× faster than **MiniGraph**; 11.92× faster than **Hyperblocker**.

- **Multi-machine**: 8.09× faster than 32-node **GCRClean**.

- **Scalability**: speedup by 1.56× with **2 GPUs**; 2.08× with **4 GPUs**.

| Single-machine (BioGRID) | | Scalability (SemScholar) | | Ablation study (SemScholar) | |
|---|---|---|---|---|---|
| **System** | **Time (s)** | **System** | **Time (s)** | **System** | **Time (s)** |
| Blaze | OOM | GCRClean-32-Node | 8.09× | noBundle | 1.35× |
| MiniGraph | 65.34× | MiniClean-2-GPU | 0.64× | noPipelinedPar | 1.26× |
| Hyperblocker | 11.92× | MiniClean-4-GPU | 0.48× | noIndPar | 1.12× |
| MiniClean-1-GPU | 259.6s | MiniClean-1-GPU | 4251.08s | MiniClean-2-GPU | 2993.7s |

- **Single-machine**: 65.34× faster than **MiniGraph**; 11.92× faster than **Hyperblocker**.

- **Multi-machine**: 8.09× faster than 32-node **GCRClean**.

- **Scalability**: speedup by 1.56× with **2 GPUs**; 2.08× with **4 GPUs**.

# Welcome to play with MiniClean!