

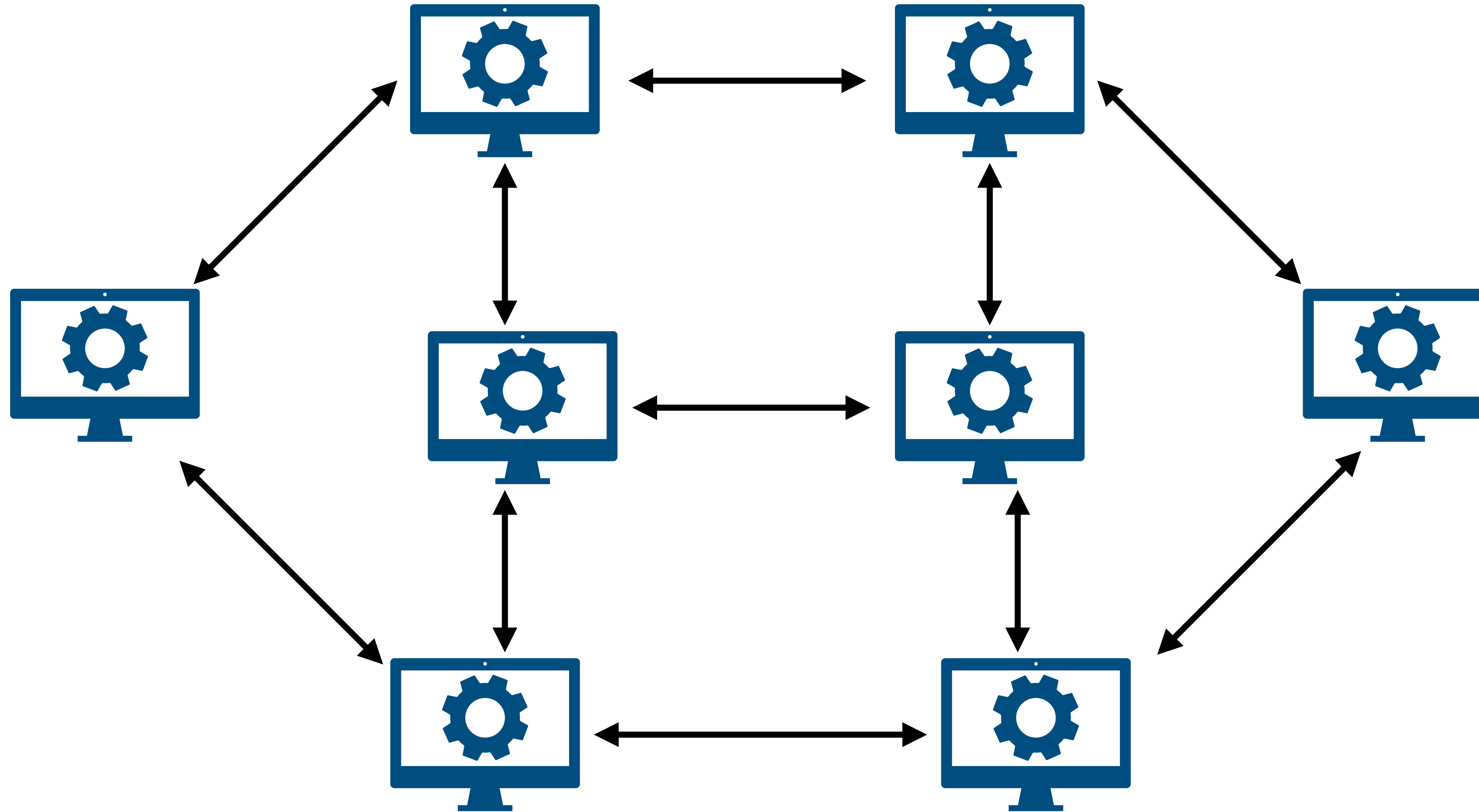
Optimizing Stratified Datalog with Count

Fernando Hechavarría Fajardo, Heba Aamer, and Bas Ketsman

Vrije Universiteit Brussel, Belgium

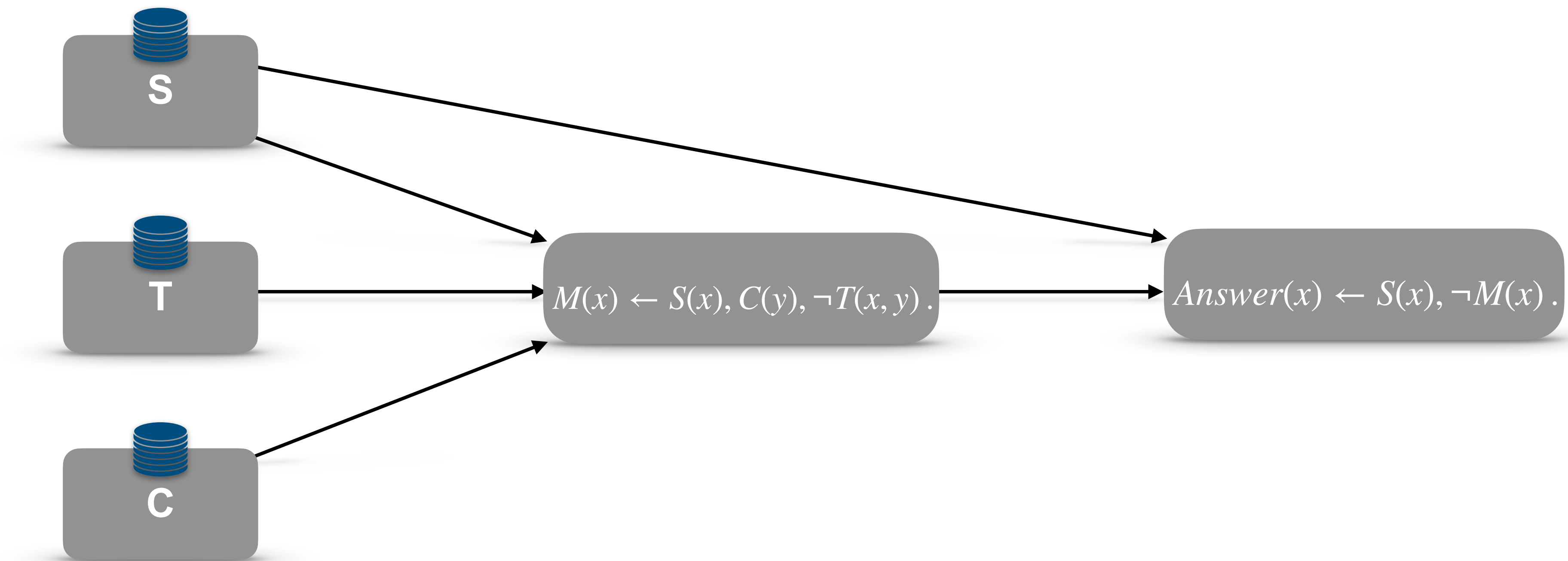
22 November 2024

Distributed Systems

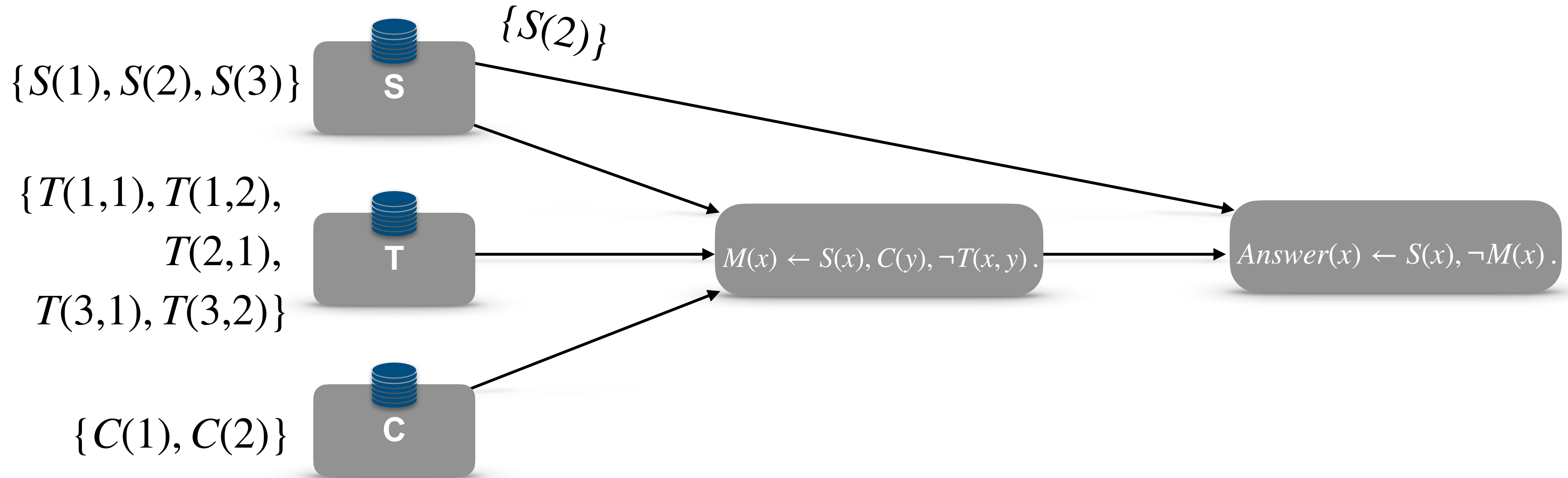


Distributed Datalog Programs

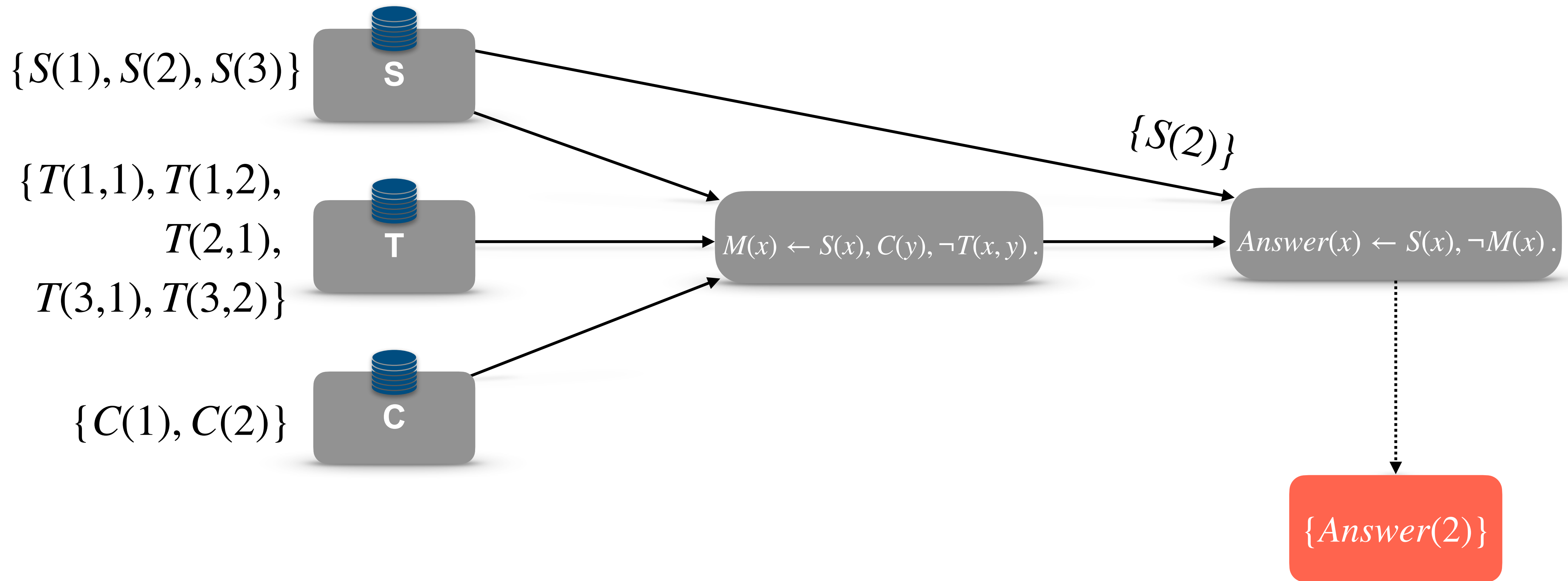
Find all *students (S) taking (T) all courses (C)*.



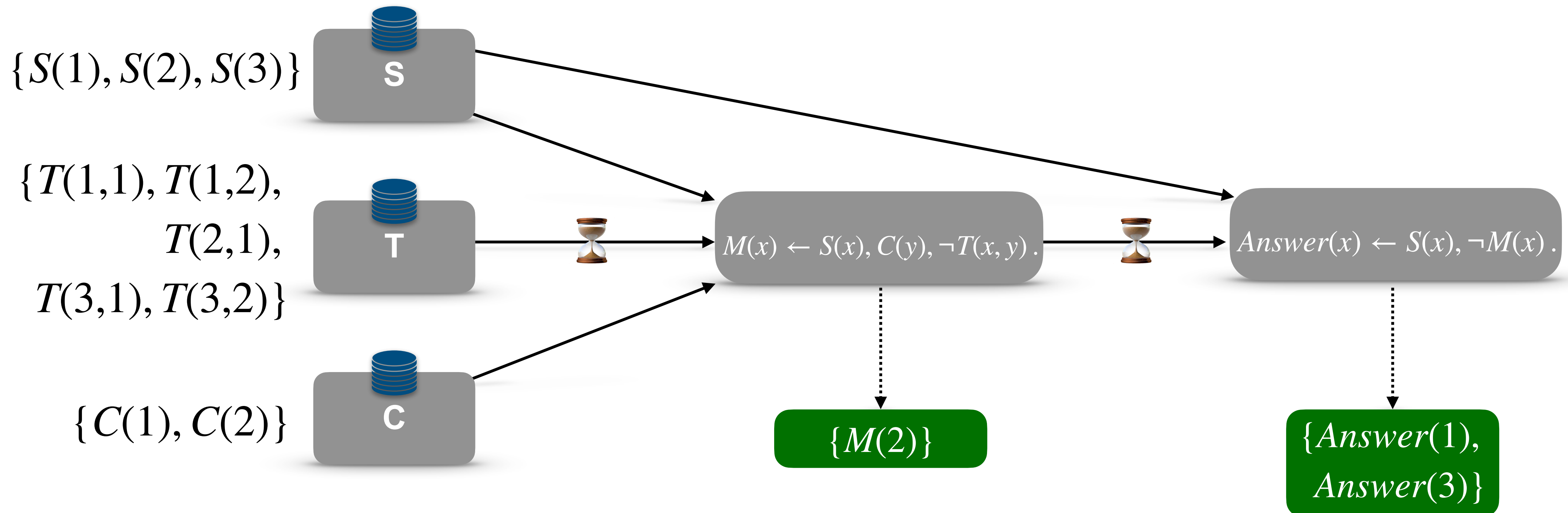
Distributed Datalog Programs



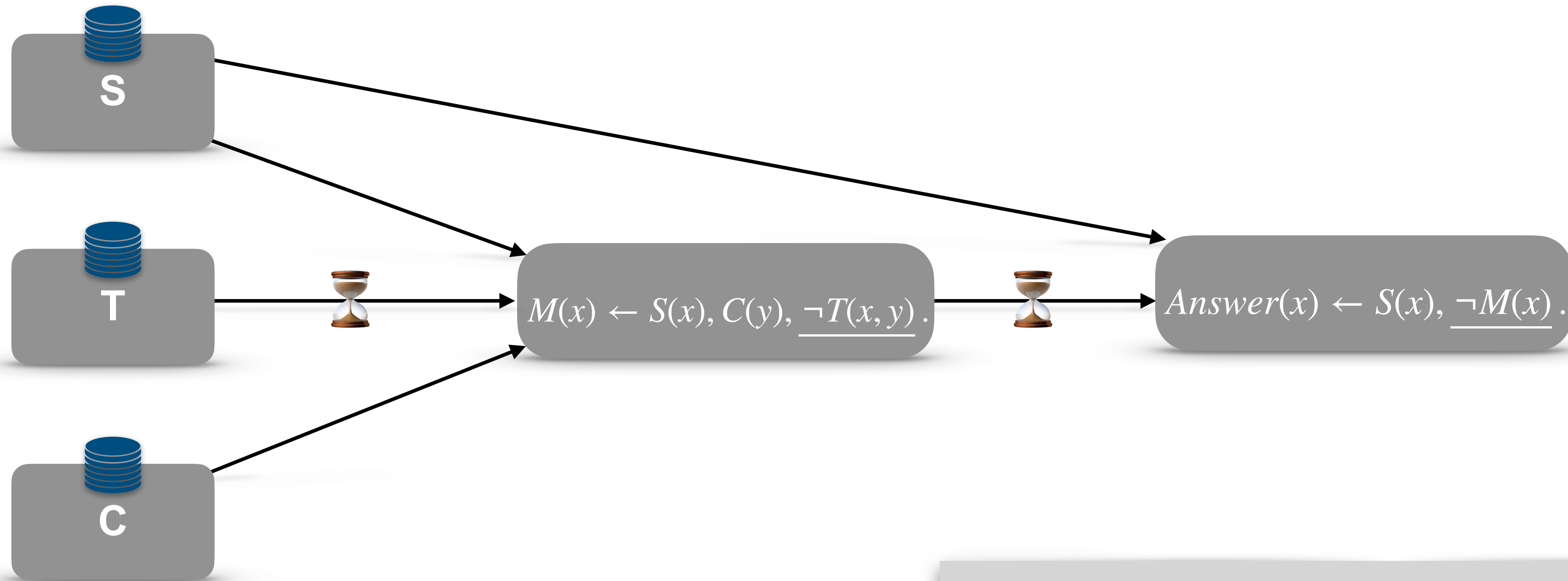
Distributed Datalog Programs



Distributed Datalog Programs



Synchronization \approx Strata



$$\frac{M(x) \leftarrow S(x), C(y), \neg T(x, y)}{Answer(x) \leftarrow S(x), \neg M(x)}$$

Goal

How can we reduce synchronization overhead?

Positive Programs

- Positive Datalog programs are confluent.
- Positive Datalog programs can only express monotone queries.

Confluence is undecidable for Datalog programs with negation.

Negation \rightarrow Incremental Count

$$M(x) \leftarrow S(x), C(y), \neg T(x, y).$$

$$\text{Answer}(x) \leftarrow S(x), \neg M(x).$$

$$T_{\#}(x; \text{count}\langle \rangle) \leftarrow C(y), T(x, y).$$

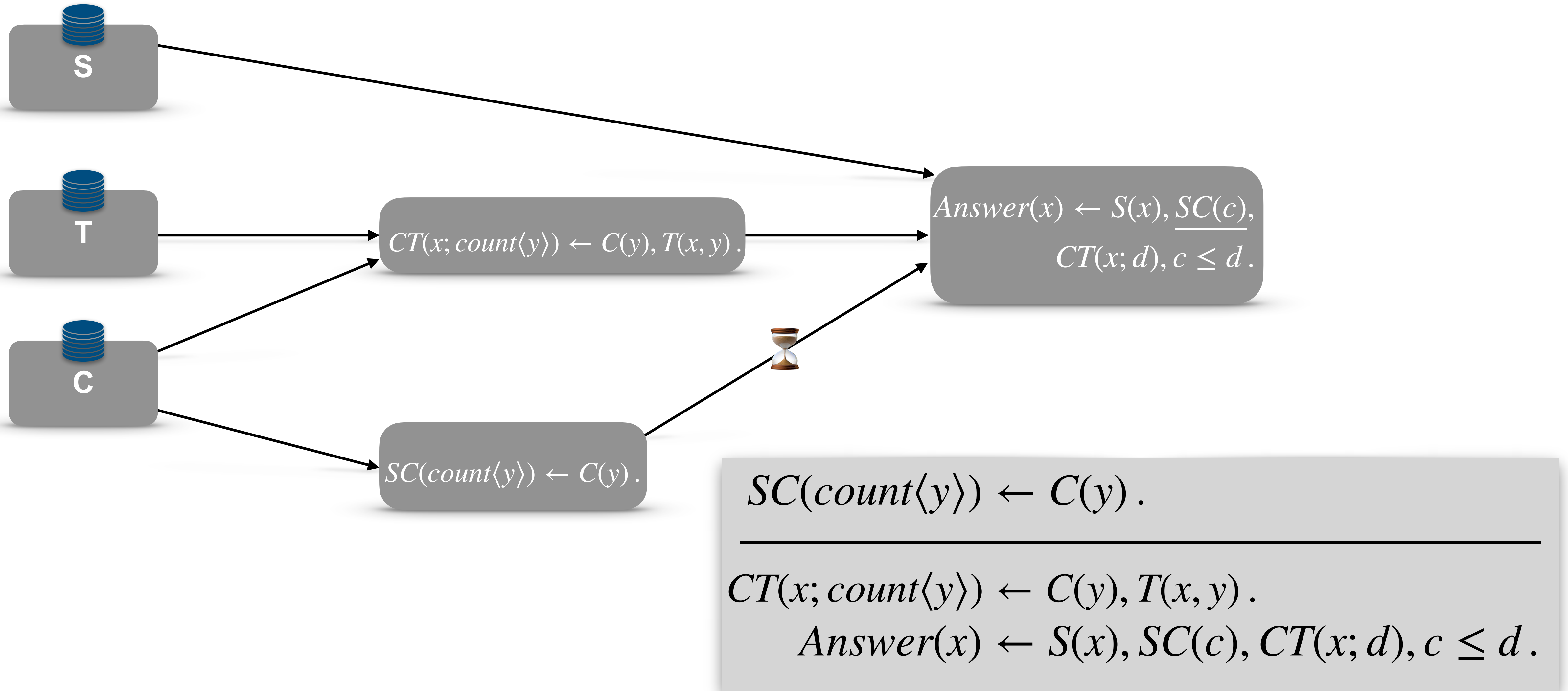
$$M(x) \leftarrow S(x), C(y), T_{\#}(x; c), c \leq 0.$$
$$M_{\#}(x; \text{count}\langle \rangle) \leftarrow M(x).$$

$$\text{Answer}(x) \leftarrow S(x), M_{\#}(x; c), c \leq 0.$$

Replacing Negation by Count

- The rewriting does not lead directly to a program with a reduced number of synchronization steps.
- Even Positive Datalog programs with count are not always confluent.

Fewer synchronization steps



Local Optimization through Counting

$$M(x) \leftarrow S(x), C(y), \neg T(x, y).$$

$$\text{Answer}(x) \leftarrow S(x), \neg M(x).$$

$$SC(\text{count}\langle y \rangle) \leftarrow C(y).$$

$$CT(x; \text{count}\langle y \rangle) \leftarrow C(y), T(x, y).$$

$$\text{Answer}(x) \leftarrow S(x), SC(c), CT(x; d), c \leq d.$$

Local Optimization through Counting

In a program with precisely one QF-generating rule of the form

$$\text{QF}(\mathbf{y};) \leftarrow_{\{\top\}} \text{R}(\mathbf{x};), \neg \text{T}(\mathbf{z};). \quad (1)$$

with $\emptyset \subsetneq \mathbf{y} \subsetneq \mathbf{x} \cap \mathbf{z}$, every rule of the form

$$\text{F}(\mathbf{u};) \leftarrow_{\{\text{QF}\}} \text{S}(\mathbf{v};), \neg \text{QF}(\mathbf{w};). \quad (2)$$

can be replaced by the rules

$$\text{Count}_{\top}(\mathbf{y}; 0) \leftarrow \text{Adom}(\mathbf{y}_1), \dots, \text{Adom}(\mathbf{y}_k).$$

$$\text{Count}_{\top}(\mathbf{y}; \text{count}\langle \mathbf{x} \cap \mathbf{z} \rangle) \leftarrow \text{R}(\mathbf{x};).$$

$$\text{Count}_{\leq}(\mathbf{y}; 0) \leftarrow \text{Adom}(\mathbf{y}_1), \dots, \text{Adom}(\mathbf{y}_k).$$

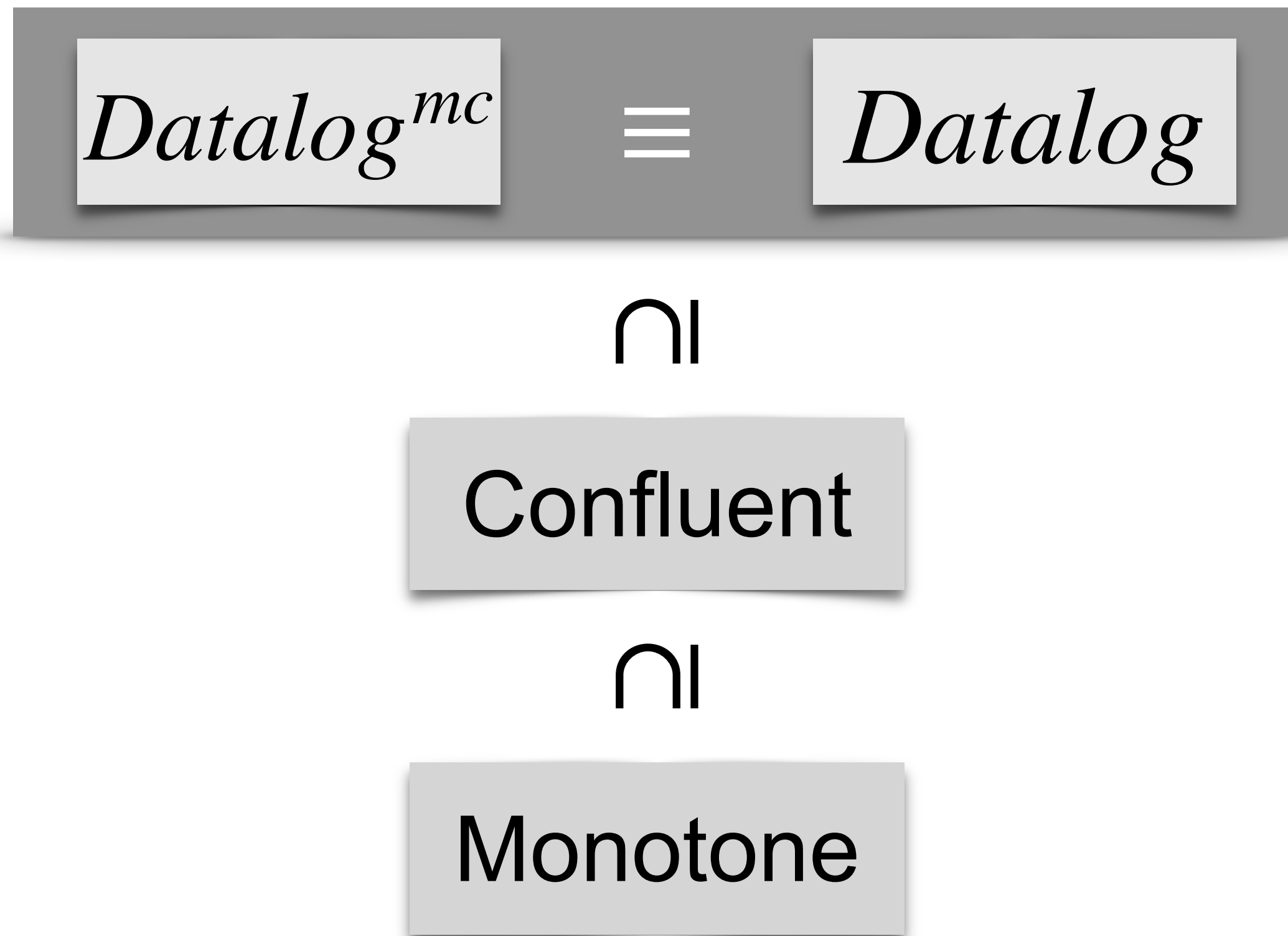
$$\text{Count}_{\leq}(\mathbf{y}; \text{count}\langle \mathbf{x} \cap \mathbf{z} \rangle) \leftarrow \text{R}(\mathbf{y};), \text{T}(\mathbf{z};).$$

$$\text{F}(\mathbf{u};) \leftarrow_{\{\text{Count}_{\top}\}} \text{S}(\mathbf{v};), \text{Count}_{\top}(\mathbf{w}; c), \text{Count}_{\leq}(\mathbf{w}; d), c \leq d.$$

Additionally, rule (1) can be removed if rule (2) is the only rule using QF. Note that we assume that relation names Count_{\top} and Count_{\leq} do not exist in the original program.

Results

We provide a rewriting technique that can reduce the *required* synchronization up to half.



Future Work

- Broader identification of confluent programs.
- Use of database-dependent constants to optimize further.
- Experimental validation to prove whether the rewriting technique improves performance in practice.