Formal approaches to:

Coordination-free query evaluation and multi-query optimization in parallel and distributed systems

Bas Ketsman





Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Context: Declarative Networking, where Datalog based languages are used for parallel and distributed computing in clusters with disordered communication.

CALM-conjecture: No-coordination $\stackrel{?}{=}$ Monotonicity

[Hellerstein, 2010]

Context: Declarative Networking, where Datalog based languages are used for parallel and distributed computing in clusters with disordered communication.

CALM-conjecture: No-coordination $\stackrel{?}{=}$ Monotonicity

[Hellerstein, 2010]

[Ameloot, Neven, Van den Bussche, 2011]: TRUE

Context: Declarative Networking, where Datalog based languages are used for parallel and distributed computing in clusters with disordered communication.

CALM-conjecture: No-coordination $\stackrel{?}{=}$ Monotonicity

[Hellerstein, 2010]

[Ameloot, Neven, Van den Bussche, 2011]: TRUE

[Zinn, Green, Ludäscher, 2012]: FALSE

Context: Declarative Networking, where Datalog based languages are used for parallel and distributed computing in clusters with disordered communication.

CALM-conjecture: No-coordination $\stackrel{?}{=}$ Monotonicity

[Hellerstein, 2010]

[Ameloot, Neven, Van den Bussche, 2011]: TRUE

▶ for a setting where nodes have **no** information about the horizontal-distribution of records

[Zinn, Green, Ludäscher, 2012]: FALSE

Context: Declarative Networking, where Datalog based languages are used for parallel and distributed computing in clusters with disordered communication.

CALM-conjecture: No-coordination $\stackrel{?}{=}$ Monotonicity

[Hellerstein, 2010]

[Ameloot, Neven, Van den Bussche, 2011]: TRUE

▶ for a setting where nodes have **no** information about the horizontal-distribution of records

[Zinn, Green, Ludäscher, 2012]: FALSE

 for settings where nodes have information about the horizontal-distribution of record **Goal:** To clarify the relation between **monotonicity** and **coordination** in asynchronous systems and to reveal the more complete picture

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Definition

A query Q is monotone if $Q(\mathbf{I}) \subseteq Q(\mathbf{I} \cup \mathbf{J})$ for all database instances \mathbf{I} and \mathbf{J} .

Notation: M = class of monotone queries

Definition

A query Q is monotone if $Q(\mathbf{I}) \subseteq Q(\mathbf{I} \cup \mathbf{J})$ for all database instances \mathbf{I} and \mathbf{J} .

Notation: M = class of monotone queries

Example

- ▶ Q_Δ : Select triangles in a graph $\in \mathcal{M}$
- ▶ $Q_{<}$: Select open triangles in a graph $\notin M$

Definition

A query Q is monotone if $Q(\mathbf{I}) \subseteq Q(\mathbf{I} \cup \mathbf{J})$ for all database instances \mathbf{I} and \mathbf{J} .

Notation: M = class of monotone queries

Example



Definition

A query Q is monotone if $Q(\mathbf{I}) \subseteq Q(\mathbf{I} \cup \mathbf{J})$ for all database instances \mathbf{I} and \mathbf{J} .

Notation: \mathcal{M} = class of monotone queries

Example

▶ Q_Δ : Select triangles in a graph $\in \mathcal{M}$



Relational Transducer Networks

[Ameloot, Neven, Van den Bussche, 2011]

- Network $\mathcal{N} = \{x, y, u, z\}$
- ► Transducer Π
- messages can be arbitrarily delayed but never get lost



Relational Transducer Networks

[Ameloot, Neven, Van den Bussche, 2011]

- Network $\mathcal{N} = \{x, y, u, z\}$
- ► Transducer Π
- messages can be arbitrarily delayed but never get lost



Semantics defined in terms of runs over a transition system

Eventual Consistent Query Evaluation

Definition

A transducer Π computes a query Q if

- ▶ for all networks \mathcal{N} , ← Network independent
- ► for all databases I, ____ Distribution independent
- for all horizontal distributions H, and
- for every run of Π ,

 $out(\Pi) = Q(\mathbf{I}).$

Consistency requirement

Example: \mathcal{Q}_{Δ} : select all triangles







Example: Q_{Δ} : select all triangles



Example: \mathcal{Q}_{Δ} : select all triangles







Example: Q_{Δ} : select all triangles



Example: Q_{Δ} : select all triangles



Extremely naive, but works .. and is coordination-free!

















Coordination is needed to reason about the absence of records.

Coordination-freeness

Goal: separate data-communication from coordination-communication

Coordination-freeness

Goal: separate data-communication from coordination-communication

Definition

IT is coordination-free if for all inputs I there is a distribution on which IT computes Q(I) without having to do communication.

[Ameloot, Neven, Van den Bussche, 2011]

Example: Ideal Distribution

 Q_{Δ} : select all triangles







Example: Ideal Distribution

 Q_Δ : select all triangles





Algorithm:

- Broadcast all data
- Output triangles whenever new data arrives



Example: Ideal Distribution

 Q_Δ : select all triangles





Algorithm:

- ▶ (Broadcast all data)
- Output triangles whenever new data arrives



CALM-conjecture

[Ameloot, Neven, Van den Bussche, 2011]

A query has a coordination-free and eventually consistent execution strategy iff the query is monotone Theorem $\mathcal{F}_0 = \mathcal{M}$

Definition

 $\mathcal{F}_0 = \text{set of queries which are distributedly computed by coordination-free transducers}$

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization









Deduction rules

- in local database \Rightarrow in global database
- not in local database + in scope \Rightarrow not in global database
- ▶ not in local database + not in scope \Rightarrow unknown


Policy-aware Transducers



Policy-aware Transducers

[Zinn, Green, Ludäscher, 2012]

Definition

A distribution policy **P** for σ and \mathcal{N} is a total function from $facts(\sigma)$ to the power set of \mathcal{N} .

Definition

A policy-aware transducer is a transducer with access to ${\bf P}$ restricted to its active domain

Definition

 \mathcal{F}_1 = set of queries which are distributedly computed by policy-aware coordination-free transducers

Definition

A fact **f** is domain distinct from instance **I** when $adom(\mathbf{f}) \not\subseteq adom(\mathbf{I})$.



Definition

A query \mathcal{Q} is domain-distinct-monotone if $\mathcal{Q}(\mathbf{I}) \subseteq \mathcal{Q}(\mathbf{I} \cup \mathbf{J})$ for all \mathbf{I} and \mathbf{J} , with \mathbf{J} having only domain-distinct facts

Notation: $M_{distinct}$ = domain-distinct-monotone queries



Remark

 $\mathcal{M}_{distinct}$: class of queries preserved under extensions



Select open triangles in graph $\in \mathcal{M}_{distinct}$.



Ι





Select open triangles in graph $\in \mathcal{M}_{distinct}$.





Not domain-distinct from ${\bf I}$

Revised CALM-conjecture



Definition

 \mathcal{F}_1 = set of queries which are distributedly computed by policy-aware coordination-free transducers

For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

Example: $Q_{<}$: Select open triangles in a graph



For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

Example: $Q_{<}$: Select open triangles in a graph



For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

Example: $Q_{<}$: Select open triangles in a graph



For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

Example: $Q_{<}$: Select open triangles in a graph



For domain-distinct-monotone queries:

- ► broadcast all present and deduced absent facts
- ► Evaluate query on complete sets

Example: $Q_{<}$: Select open triangles in a graph



Complete Picture



[Ameloot, Ketsman, Neven, Zinn PODS 2014 best paper; TODS 2016]

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Goal: To clarify the relation between **monotonicity** and **coordination** in asynchronous systems and to reveal the more complete picture

- A four-level quantification of coordination exists in terms of the amount of information needed for the query to become coordination-free
- The CALM conjecture reveals a very robust relation between coordination and non-monotonic behaviour of queries

Open Questions & Future Work

- ► How to compute queries without coordination and with less communication?
- ► How to compute queries with "some" coordination?

Open Questions & Future Work

- ► How to compute queries without coordination and with less communication?
- ► How to compute queries with "some" coordination?

Related Work

- Oblivious broadcasting algorithms = broadcast fragment of local database [Ketsman, Neven ICDT 2015; ToCS 2016]
- A worst-case optimal load algorithm for join evaluation [Ketsman, Suciu, PODS 2017]

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Motivation

- Many systems rely on coordination for communication
 For example: MapReduce-like systems
- ► Avoiding coordination completely is not possible
- Minimize the number of communication steps
 Ideally: one round

(1-Round MPC model [Koutris & Suciu 2011])

Input = query Q



(1-Round MPC model [Koutris & Suciu 2011])

Input = query Q





(1-Round MPC model [Koutris & Suciu 2011])

Input = query Q



(1-Round MPC model [Koutris & Suciu 2011])

Input = query Q



Output = union of output at each server

Multi-Query Evaluation & optimization

Workload: Q_1, Q_2, \ldots, Q_n + fixed database



. . .

Goal: To formally reason about **single-round** query evaluation and **multi-query optimization** in systems where communication implies a synchronization barrier

Focus: conjunctive queries





Output = union of output at each server

Main question 1

Given target query and distribution policy:

Do we need to reshuffle?

Main question 1

Given target query and distribution policy:

Do we need to reshuffle?

Parallel-Correctness

Main question 1

Given target query and distribution policy:

Do we need to reshuffle?

Parallel-Correctness

Main question 2

Given target query and previously computed query:

Do we need to reshuffle?

Main question 1

Given target query and distribution policy:

Do we need to reshuffle?

Parallel-Correctness

Main question 2

Given target query and previously computed query:

Do we need to reshuffle?

Transferability

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Given target query and distribution policy:

Do we need to reshuffle?

Given target query and distribution policy:

Do we need to reshuffle?

Semantical correctness of simple evaluation algorithm

Given target query and distribution policy:

Do we need to reshuffle?

Semantical correctness of simple evaluation algorithm

Definition

Q is parallel-correct w.r.t. P, iff for every database I

$$\mathcal{Q}(I) = \bigcup_{\kappa \in \mathcal{N}} \mathcal{Q}(dist_{\mathbf{P},I}(\kappa))$$

Given target query and distribution policy:

Do we need to reshuffle?

Semantical correctness of simple evaluation algorithm

Definition

Q is parallel-correct w.r.t. P, iff for every database I

$$\mathcal{Q}(I) = \bigcup_{\kappa \in \mathcal{N}} \mathcal{Q}(dist_{\mathbf{P},I}(\kappa))$$
$$\supseteq \text{ by monotonicity}$$
Parallel-Correctness Complexity



(needs policy representation)

Other use of Parallel-Correctness



Parallel-correct for \mathcal{Q}_1

Other use of Parallel-Correctness



Parallel-correct for \mathcal{Q}_1

Possible Problems

- ▶ Reasoning about distribution policies is complex
- ► Not every distribution policy is equally efficient
- ► Choice of policy may be hidden behind abstraction layer
- ► Reasoning about query order before policies are known

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Given target query and previously computed query:

Do we need to reshuffle?

Given target query and previously computed query:

Do we need to reshuffle?

Definition

 $\mathcal{Q} \to_T \mathcal{Q}'$ iff \mathcal{Q}' is parallel-correct on every P where \mathcal{Q} is parallel-correct on

Given target query and previously computed query:

Do we need to reshuffle?

Parallel-correct for Q_1

Given target query and previously computed query:

Do we need to reshuffle?

Parallel-correct for Q_1

Parallel-correct for \mathcal{Q}_2

Given target query and previously computed query:

Do we need to reshuffle?

Parallel-correct for \mathcal{Q}_1

Parallel-correct for \mathcal{Q}_2

Very strong property

like query containment, but for parallel and distributed setting

Complexity

	Parallel-correctness*	Transferability
CQ	Π_p^2 -C	Π_p^3 -C
UCQ	Π_p^2 -c	Π_p^3 -C
UCQ≠	Π_p^2 -C	Π_p^3 -C
FO	undecidable	undecidable
sm-CQ	Х	NP-c

(*needs policy representation)

[Ameloot,Geck,Ketsman,Neven,Schwentick PODS 2015 best paper; Sigmod record 2016; CACM 2017; JACM (accepted)]

Outline

CALM Formalization

CALM Revision 1

Conclusion

Parallel-Correctness

Transferability

Conclusion

Coordination-free evaluation

Multi-Query optimization

Goal: To formally reason about **single-round** query evaluation and **multi-query optimization** in systems where communication implies a synchronization barrier

- A formal framework for reasoning about the correctness of single-round query evaluation and query optimization via distribution policies.
- ► Parallel-correctness: semantical correctness
- Transferability: like "containment" but for parallel and distributed query evaluation

Open Questions & Future Work

- ► How do parallel-correctness and transferability relate to query evaluation in practice?
- ► How much data needs to be reshuffled?

Open Questions & Future Work

- ► How do parallel-correctness and transferability relate to query evaluation in practice?
- ► How much data needs to be reshuffled?

Related Work

▶ Parallel-correctness for CQs with negation

[Geck, Ketsman, Neven, Schwentick ICDT 2016]

 Extension of Parallel-correctness to reason about multi-round evaluation with with Datalog

[Ketsman, Koutris, Albarghouti, submitted]

► Bag-semantics

ongoing

Thank you!