DEPARTMENT OF COMPUTER SCIENCE





Ontology Based Data Access (OBDA)

Ian Horrocks Information Systems Group Department of Computer Science University of Oxford

Motivation

- Huge quantity of data increasing at an exponential rate
- Identifying & accessing relevant data is of critical importance
- Handling data variety & complexity often turns out to be main challenge
- Semantic Technology can seamlessly integrate heterogeneous data sources















Semantic Technology

Rich **conceptual schemas** used to integrate heterogeneous sources

- User Centric
 - Schema modelled according to user intuitions
 - Independent of physical structure/storage of data
- Declarative
 - Improved understandability
 - Easier design, maintenance and evolution
- Logic-based semantics
 - Precise and formally specified meaning
 - Machine processable
- Used at both design and query time
 - Check validity and consequences of design
 - Easier query formulation and enriched query answers









Semantic Technology: Scalability Challenge



How can we provide (empirically) scalable query answering?











OWL Profiles

OWL 2 defines language subsets, aka **profiles** that can be "more simply and/or efficiently implemented"

- OWL 2 QL
 - Based on DL-Lite
 - Efficiently implementable via rewriting into relational queries (OBDA)





































(:p1, rdf:type, :Pipeline) (:p1, :fromFacility, :f1) (:f1, rdf:type, :OilFacility) (:p2, rdf:type, :OilPipeline) (:p2, :fromFacility, :f2) (:f2, rdf:type, :OilFacility) (:p3, rdf:type, :OilPipeline)











































































































Given QL ontology \mathcal{O} query \mathcal{Q} and mappings \mathcal{M} :











Given QL ontology \mathcal{O} query \mathcal{Q} and mappings \mathcal{M} :

Use O to rewrite Q → Q' s.t. answering Q' without O is equivalent to answering Q w.r.t. O for any dataset













Given QL ontology \mathcal{O} query \mathcal{Q} and mappings \mathcal{M} :

- Use O to rewrite Q → Q' s.t. answering Q' without O is equivalent to answering Q w.r.t. O for any dataset
- Map ontology queries → DB queries (typically SQL) using mappings *M* to rewrite *Q*' into a DB query











Given QL ontology \mathcal{O} query \mathcal{Q} and mappings \mathcal{M} :

- Use O to rewrite Q → Q' s.t. answering Q' without O is equivalent to answering Q w.r.t. O for any dataset
- Map ontology queries → DB queries (typically SQL) using mappings *M* to rewrite *Q*' into a DB query
- Evaluate (SQL) query against DB





Information Systems Group









Information Systems Group









Information Systems Group











Information Systems Group









Information Systems Group







Query Rewriting — Issues

1 Rewriting

- May be large (worst case exponential in size of ontology)
- Queries may be hard for existing DBMSs

2 Mappings

May be difficult to develop and maintain

3 Expressivity

- OWL 2 QL (necessarily) has (very) restricted expressive power, e.g.:
 - No functional or transitive properties
 - No universal (for-all) restrictions
 - ...









OWL Profiles – Beyond QL?

OWL 2 defines language subsets, aka **profiles** that can be "more simply and/or efficiently implemented"

- OWL 2 QL
 - Based on DL-Lite
 - Efficiently implementable via rewriting into relational queries (OBDA)
- OWL 2 RL
 - Based on "**Description Logic Programs**" (\approx DL \cap Datalog)
 - Implementable via Datalog query answering
- OWL 2 EL
 - Based on *EL*⁺⁺
 - Implementable via Datalog query answering plus "filtration"









RL/Datalog Query Ans. via Materialisation

Given (RDF) data DB, RL/Datalog ontology \mathcal{O} and query \mathcal{Q} :











RL/Datalog Query Ans. via Materialisation

Given (RDF) data DB, RL/Datalog ontology \mathcal{O} and query \mathcal{Q} :

• Materialise (RDF) data DB \rightarrow DB' s.t. evaluating Q w.r.t. DB' equivalent to answering Q w.r.t. DB and O

nb: Closely related to chase procedure used with DB dependencies

Research Council





RL/Datalog Query Ans. via Materialisation

Given (RDF) data DB, RL/Datalog ontology \mathcal{O} and query \mathcal{Q} :

• Materialise (RDF) data DB \rightarrow DB' s.t. evaluating Q w.r.t. DB' equivalent to answering Q w.r.t. DB and O

nb: Closely related to chase procedure used with DB dependencies

Evaluate Q against DB'











 $\mathcal{O} \left\{ \mathsf{OilEquip} \sqsubseteq \forall \mathsf{hasPart}.\mathsf{OilEquip} \right.$



Information Systems Group





EPSRC Optique



 $\mathcal{O} \left\{ \mathsf{OilEquip} \sqsubseteq \forall \mathsf{hasPart}.\mathsf{OilEquip} \right.$















$Q_1 \quad Q(x) \leftarrow \mathsf{OilEquip}(x)$



























 $\mathcal{O} \left\{ \begin{array}{ll} \mathsf{OilEquip} \sqsubseteq \forall \mathsf{hasPart.OilEquip} & \mathsf{OilEquip}(x) \land \mathsf{hasPart}(x,y) \to \mathsf{OilEquip}(y) \end{array} \right.$

$$\mathsf{DB} \left\{ \begin{array}{l} \mathsf{OilEquip}(\mathsf{a}) \\ \mathsf{hasPart}(\mathsf{a},\mathsf{b}) \\ \mathsf{hasPart}(\mathsf{b},\mathsf{c}) \end{array} \right.$$

$Q_1 \quad Q(x) \leftarrow \mathsf{OilEquip}(x)$























 $\mathsf{DB} \left\{ \begin{array}{l} \mathsf{OilEquip}(\mathsf{a}) \\ \mathsf{hasPart}(\mathsf{a},\mathsf{b}) \\ \mathsf{hasPart}(\mathsf{b},\mathsf{c}) \end{array} \right.$

$$\mathsf{DB'} \left\{ \begin{array}{l} \mathsf{OilEquip(a)} \\ \mathsf{hasPart(a,b)} \\ \mathsf{hasPart(b,c)} \\ \mathsf{OilEquip(b)} \end{array} \right.$$

$$\mathcal{Q}_1 \quad Q(x) \leftarrow \mathsf{OilEquip}(x)$$











 $\mathsf{DB}\left\{\begin{array}{l}\mathsf{OilEquip(a)}\\\mathsf{hasPart(a,b)}\\\mathsf{hasPart(b,c)}\end{array}\right.$

$$\mathsf{DB'} \left\{ \begin{array}{l} \mathsf{OilEquip}(a) \\ \mathsf{hasPart}(a,b) \\ \mathsf{hasPart}(b,c) \\ \mathsf{OilEquip}(b) \\ \mathsf{OilEquip}(c) \end{array} \right.$$

$$\mathcal{Q}_1 \quad Q(x) \leftarrow \mathsf{OilEquip}(x)$$











$Q_1 \quad Q(x) \leftarrow \mathsf{OilEquip}(x)$

 $\rightsquigarrow \{a, b, c\}$









Materialisation — Issues

1 Scalability

- Ptime complete
- Efficiently implementable in practice?

2 Updates

- Additions relatively easy (continue materialisation)
- But what about retraction?

3 Migrating data to RDF

- Materialisation assumes data in "special" (RDF triple) store
- How can legacy data be migrated?

4 Expressivity

• $\mathsf{QL} \not\subseteq \mathsf{RL}$; in particular, no invention of new individuals











Materialisation: Scalability

- Efficient Datalog/RL engine is critical
- Existing approaches mainly target distributed "shared-nothing" architectures, often via map reduce
 - High communication overhead
 - Typically focus on small fragments (e.g., RDFS), so don't really address expressivity issue
 - Even then, query answering over (distributed) materialized data is nontrivial and may require considerable communication









RDFox Datalog Engine

- Targets SOTA main-memory, mulit-core architecture
 - Optimized in-memory storage with 'mostly' lock-free parallel inserts
 - Memory efficient: commodity server with 128 GB can store >10⁹ triples
 - Exploits multi-core architecture: 10-20 x speedup with 32/16 threads/cores
 - LUBM 120K (>10¹⁰ triples) in 251s (20M t/s) on T5-8 (4TB/1024 threads)











RDFox Datalog Engine

- Incremental addition and retraction of triples
 - Retraction via novel FBF "view maintenance" algorithm
 - Retraction of 5,000 triples from materialised LUBM 50k in less than 1s
- Many other novel features
 - Handles more general (than RL) Dalalog and SWRL rules
 - SPARQL features such as BIND and FILTER in rule bodies
 - Native equality handling (owl:sameAs) via rewriting
 - Stratified negation as failure (NAF)









Materialisation: Data Migration

- Need to specify a suitable **migration** process
 - Use R2RML mappings to extract data and transform into RDF
 - But where do these mappings come from?
- Recall query rewriting:
 - Mappings \mathcal{M} are R2RML mappings
 - Run mappings in reverse to extract and transform data

• "Lazy ETL"

- Deploy query rewriting (OBDA) system
- Extend ${\mathcal O}$ and ${\mathcal M}$ as needed
- Use ${\boldsymbol{\mathcal{M}}}$ to ETL data into RDF store













Materialisation: Expressivity

- RL is more powerful than QL, but $QL \not\subseteq RL$
 - In particular, no "individual creation" (RHS existentials)
 - Can't express, e.g., OilPipeline
 Pipeline
 IfromFacility.OilFacility
- Recall OWL 2 EL
 - Based on *EL*⁺⁺
 - Implementable via Datalog query answering plus "filtration"









Given (RDF) Data Set, EL ontology \mathcal{O} and query \mathcal{Q} :





Information Systems Group









Research Council

Given (RDF) Data Set, EL ontology \mathcal{O} and query \mathcal{Q} :

 Over-approximate O into Datalog program D





Information Systems Group









Research Council

Given (RDF) Data Set, EL ontology \mathcal{O} and query \mathcal{Q} :

- **Over-approximate** \mathcal{O} into Datalog program D
- Evaluate Q over D + Data Set (via materialisation)













Given (RDF) Data Set, EL ontology \mathcal{O} and query \mathcal{Q} :

- Over-approximate O into Datalog program D
- Evaluate Q over D + Data Set (via materialisation)
- Use (polynomial) Filtering Procedure to eliminate spurious answers













Discussion

- **QL-Rewriting** has many advantages
 - Data can be left untouched and in legacy storage
 - Exploits existing DB infrastructure and scalability
 - • • •
- But what if more expressiveness/flexibility is needed?
 - Query answering for EL and RL still tractable (polynomial)
 - Critically depend on Datalog scalability RDFox to the rescue!
 - Easy migration path from QL-rewriting via "lazy ETL"









Future Work

- Piloting, evaluation and tuning
- Porting to other large-scale architectures
- Semantic (data) partitioning for distributed architectures
- (Incremental maintenance of) aggregations
- Improved query planning
- Stream reasoning
- Hybrid rewriting/materialisation (on demand) approach
- Expressiveness beyond RL/EL via PAGOdA techniques









Acknowledgements







































Engineering and Physical Sciences Research Council













Information Systems Group







Thank you for listening











Thank you for listening

THREE LOGICIANS WALK INTO A BAR...



Any questions?



Information Systems Group









Research Council