

# CLASSIFYING $\mathcal{ELH}$ ONTOLOGIES IN SQL DATABASES

Vincent Delaitre and Yevgeny Kazakov  
(Presented by Rob Shearer)

Oxford University Computing Laboratory

October 24, 2009





# OUTLINE

## 1 INTRODUCTION

## 2 PROCEDURE OUTLINE

## 3 PROBLEMS AND SOLUTIONS

## 4 RESULTS

 $\mathcal{ELH}$  AND  $\mathcal{OWL} 2 \text{ EL}$ 

OWL 2 Syntax	DL Syntax
Class expressions:	
ObjectIntersectionOf( $C D$ )	$C \sqcap D$
ObjectSomeValuesFrom( $r C$ )	$\exists r.C$
Axioms:	
SubClassOf( $C D$ )	$C \sqsubseteq D$
EquivalentClasses( $C D$ )	$C \equiv D$
SubObjectPropertyOf( $r s$ )	$r \sqsubseteq s$

- $\mathcal{ELH}$  is a simple sub-fragment of  $\mathcal{OWL} 2 \text{ EL}$
- Has a very simple polynomial-time classification procedure [Baader et al., IJCAI 2003, 2005]
- Sufficiently expressive for many ontologies such as SNOMED, FMA, NCI, GO and large part of GALEN
- Has a potential of scaling to even larger ontologies



# ARE WE READY FOR ONTOLOGIES WITH MILLIONS OF CLASSES?

- **SNOMED CT** contains over **300,000** classes—probably the largest ontology available so far
- Can be classified **in minutes** using many existing reasoners:

	Time
CEL	21min.42s.
FaCT++	16min.05s.
RACER	19min.30s.
SNOROCKET	1min.06s.
CB	0min.45s.



# ARE WE READY FOR ONTOLOGIES WITH MILLIONS OF CLASSES?

- SNOMED CT contains over 300,000 classes—probably the largest ontology available so far
- Can be classified in minutes using many existing reasoners:

	Time	Memory
CEL	21min.42s.	700MB*
FaCT++	16min.05s.	320MB*
RACER	19min.30s.	900MB
SNOROCKET	1min.06s.	2GB
CB	0min.45s.	400MB

- But memory consumption could be a problem for ontologies 10x larger.



# SECONDARY MEMORY ONTOLOGY REASONING

- The main idea: use a DBMS for processing of ontologies
- Advantages:
  - 1 Low main memory footprint
  - 2 Persistence: can save / restore computations
  - 3 Transactions and fault tolerance
  - 4 Possible to adapt to multi-user environments
- Disadvantage:
  - 1 Slow (because of the secondary memory characteristics)



# SECONDARY MEMORY ONTOLOGY REASONING

- The main idea: use a DBMS for processing of ontologies
- Advantages:
  - 1 Low main memory footprint
  - 2 Persistence: can save / restore computations
  - 3 Transactions and fault tolerance
  - 4 Possible to adapt to multi-user environments
- Disadvantage:
  - 1 Slow (because of the secondary memory characteristics)
- Our main results:
  - It is possible to classify *ELH* ontologies in SQL databases
  - Naive approach has poor performance
  - Optimizations (caching) improve performance significantly
  - **Able to classify SNOMED CT in 20min using 32MB of RAM.**



## (UN)RELATED WORKS

- Conjunctive query answering in  $\mathcal{EL}$  using relational databases [Lutz, Toman, Wolter, IJCAI 2009]
  - large instance data
  - medium-size schema (60,000 classes)
  - main focus is on query response
- “DB-backed” module in the IBM SHER system
  - Uses a Datalog engine
  - Presumably can work with  $\mathcal{EL}^+$  ontologies
  - Cannot classify SNOMED CT(?)
- RDF databases:
  - Can query large triple stores
  - Can use custom rules
  - Cannot classify OWL ontologies(?)





# OUTLINE

1 INTRODUCTION

2 PROCEDURE OUTLINE

3 PROBLEMS AND SOLUTIONS

4 RESULTS



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 **Normalization** to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 **Normalization** to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r. \underline{(B \sqcap C)} \rightsquigarrow$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 **Normalization** to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \quad \rightsquigarrow \quad A \sqsubseteq \exists r.\underline{D} \quad \underline{D} \sqsubseteq B \sqcap C$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 **Normalization** to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \quad \rightsquigarrow \quad A \sqsubseteq \exists r.D \quad D \sqsubseteq \underline{B \sqcap C}$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 **Normalization** to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \quad \rightsquigarrow \quad A \sqsubseteq \exists r.D \quad D \sqsubseteq B \quad D \sqsubseteq C$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 Normalization to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \quad \rightsquigarrow \quad A \sqsubseteq \exists r.D \quad D \sqsubseteq B \quad D \sqsubseteq C$$

2 Deriving consequences using the rules [Brandt, ECAI 2004]:



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

1 Normalization to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \quad \rightsquigarrow \quad A \sqsubseteq \exists r.D \quad D \sqsubseteq B \quad D \sqsubseteq C$$

2 Deriving consequences using the rules [Brandt, ECAI 2004]:

$$\text{IR1} \quad \frac{}{A \sqsubseteq A} \qquad \text{IR2} \quad \frac{}{A \sqsubseteq \top} \quad (\text{tautologies})$$



 $\mathcal{ELH}$  CLASSIFICATION PROCEDURE

**1** Normalization to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \rightsquigarrow A \sqsubseteq \exists r.D \quad D \sqsubseteq B \quad D \sqsubseteq C$$

**2** Deriving consequences using the rules [Brandt, ECAI 2004]:

$$\text{IR1} \frac{}{A \sqsubseteq A} \quad \text{IR2} \frac{}{A \sqsubseteq \top} \quad (\text{tautologies})$$

$$\text{CR1} \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O} \quad \text{CR2} \frac{A \sqsubseteq B \quad A \sqsubseteq C}{A \sqsubseteq D} : B \sqcap C \sqsubseteq D \in \mathcal{O}$$



# $\mathcal{ELH}$ CLASSIFICATION PROCEDURE

**1** Normalization to simple axioms of five forms:

$$(1) A \sqsubseteq B \quad (2) A \sqcap B \sqsubseteq C \quad (3) A \sqsubseteq \exists r.B \quad (4) \exists r.B \sqsubseteq C \quad (5) r \sqsubseteq s$$

## EXAMPLE

$$A \sqsubseteq \exists r.(B \sqcap C) \rightsquigarrow A \sqsubseteq \exists r.D \quad D \sqsubseteq B \quad D \sqsubseteq C$$

**2** Deriving consequences using the rules [Brandt, ECAI 2004]:

$$\text{IR1} \frac{}{A \sqsubseteq A} \quad \text{IR2} \frac{}{A \sqsubseteq \top} \quad (\text{tautologies})$$

$$\text{CR1} \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O} \quad \text{CR2} \frac{A \sqsubseteq B \quad A \sqsubseteq C}{A \sqsubseteq D} : B \sqcap C \sqsubseteq D \in \mathcal{O}$$

$$\text{CR3} \frac{A \sqsubseteq B}{A \sqsubseteq \exists r.C} : B \sqsubseteq \exists r.C \in \mathcal{O} \quad \text{CR4} \frac{A \sqsubseteq \exists r.B}{A \sqsubseteq \exists s.B} : r \sqsubseteq s \in \mathcal{O}$$

$$\text{CR5} \frac{A \sqsubseteq \exists r.B \quad B \sqsubseteq C}{A \sqsubseteq D} : \exists r.C \sqsubseteq D \in \mathcal{O}$$



# DATABASE ORGANIZATION

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan (type 1)

Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem (type 3)

- Use two tables to assign ids to classes and object properties

class	id	object property	id
Heart	1	isPartOf	1
MuscularOrgan	2		
CirculatorySystem	3		

- Use five tables to store normalized axioms of each type

ax_t1	ax_t2	ax_t3	ax_t4	ax_t5
$A \sqsubseteq B$	$A \sqcap B \sqsubseteq C$	$A \sqsubseteq \exists r.B$	$\exists r.B \sqsubseteq C$	$r \sqsubseteq s$
1 2	. . .	1 1 3	. . .	. .



## COMPLETION USING SQL QUERIES

ax_t1	ax_t2	ax_t3	ax_t4	ax_t5
$A \sqsubseteq B$	$A \cap B \sqsubseteq C$	$A \sqsubseteq \exists r.B$	$\exists r.B \sqsubseteq C$	$r \sqsubseteq s$
1 2	2 5 1	1 1 3	3 4 5	1 2

- Use two tables to output the results of inferences:

s_t1	s_t2
$A \sqsubseteq B$	$A \sqsubseteq \exists r.B$
1 1	1 1 3

- Use SQL commands to perform inferences:

$$\text{IR1} \frac{}{A \sqsubseteq A}$$

```
INSERT INTO s_t1
SELECT class.id, class.id;
```

$$\text{CR1} \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```



# OUTLINE

1 INTRODUCTION

2 PROCEDURE OUTLINE

3 PROBLEMS AND SOLUTIONS

4 RESULTS



## PROBLEM 1: ASSIGNMENT OF THE IDS

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan    Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem

class

id

object property

id

 $A \sqsubseteq B$  $A \sqcap B \sqsubseteq C$  $A \sqsubseteq \exists r. B$  $\exists r. B \sqsubseteq C$  $r \sqsubseteq s$



# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan

Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



class	id
Heart	1

object	property	id

$A \sqsubseteq B$
1

$A \sqcap B \sqsubseteq C$

$A \sqsubseteq \exists r.B$

$\exists r.B \sqsubseteq C$

$r \sqsubseteq s$



# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan

Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



class	id
Heart	1
MuscularOrgan	2

object	property	id

$A \sqsubseteq B$
1 2

$A \sqcap B \sqsubseteq C$

$A \sqsubseteq \exists r.B$

$\exists r.B \sqsubseteq C$

$r \sqsubseteq s$





# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan      Heart  $\sqsubseteq$   $\exists$ isPartOf.CirculatorySystem



class	id
Heart	1
MuscularOrgan	2

object	property	id

$A \sqsubseteq B$
1   2

$A \sqcap B \sqsubseteq C$

$A \sqsubseteq \exists r.B$
1

$\exists r.B \sqsubseteq C$

$r \sqsubseteq s$



# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan

Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



class	id
Heart	1
MuscularOrgan	2

object	property	id
	isPartOf	1

$A \sqsubseteq B$
1 2

$A \sqcap B \sqsubseteq C$

$A \sqsubseteq \exists r. B$
1 1

$\exists r. B \sqsubseteq C$

$r \sqsubseteq s$



# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan      Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



class	id
Heart	1
MuscularOrgan	2
CirculatorySystem	3

object	property	id
	isPartOf	1

$A \sqsubseteq B$
1   2

$A \cap B \sqsubseteq C$

$A \sqsubseteq \exists r. B$
1   1   3

$\exists r. B \sqsubseteq C$

$r \sqsubseteq s$



# PROBLEM 1: ASSIGNMENT OF THE IDs

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan      Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem

class	id
Heart	1
MuscularOrgan	2
CirculatorySystem	3

object	property	id
	isPartOf	1

$A \sqsubseteq B$
1   2

$A \cap B \sqsubseteq C$

$A \sqsubseteq \exists r. B$
1   1   3

$\exists r. B \sqsubseteq C$

$r \sqsubseteq s$

- On-disc table lookup is **too slow!**
- Making a query for every occurrence of a class is impractical due to overheads (connection + parsing + transaction)

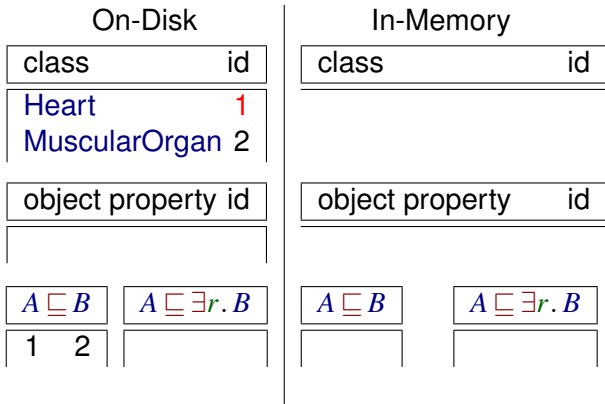


## SOLUTION: IN-MEMORY CACHING

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan    Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem

- Insert into in-memory tables with fresh ids





## SOLUTION: IN-MEMORY CACHING

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan      Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



- Insert into in-memory tables with **fresh ids**

On-Disk

class	id
Heart	1
MuscularOrgan	2

object	property	id
A	$\sqsubseteq B$	1
A	$\sqsubseteq \exists r. B$	2

A  $\sqsubseteq$  B

1

A  $\sqsubseteq \exists r. B$ 

2

In-Memory

class	id
Heart	3

object	property	id
A	$\sqsubseteq B$	
A	$\sqsubseteq \exists r. B$	3

A  $\sqsubseteq$  BA  $\sqsubseteq \exists r. B$ 

3



## SOLUTION: IN-MEMORY CACHING

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan    Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem



- Insert into in-memory tables with **fresh ids**

On-Disk

class	id
Heart	1
MuscularOrgan	2

object	property	id

$A \sqsubseteq B$
1   2

$A \sqsubseteq \exists r. B$

1	2
---	---

--	--

In-Memory

class	id
Heart	3

object	property	id

object	property	id
	isPartOf	4

object	property	id

$A \sqsubseteq B$

$A \sqsubseteq \exists r. B$
3   4

--	--

3	4
---	---



## SOLUTION: IN-MEMORY CACHING

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan    Heart  $\sqsubseteq$   $\exists$ isPartOf.CirculatorySystem



- Insert into in-memory tables with **fresh ids**

On-Disk

class	id
Heart	1
MuscularOrgan	2

Heart	1
MuscularOrgan	2

object	property	id

--	--	--

$A \sqsubseteq B$
-------------------

$A \sqsubseteq \exists r. B$
------------------------------

1	2
---	---

--	--

In-Memory

class	id
Heart	3
CirculatorySystem	5

Heart	3
CirculatorySystem	5

object	property	id

isPartOf	4
----------	---

$A \sqsubseteq B$
-------------------

$A \sqsubseteq \exists r. B$
------------------------------

--	--

3	4	5
---	---	---





## SOLUTION: IN-MEMORY CACHING

## EXAMPLE

Heart  $\sqsubseteq$  MuscularOrgan    Heart  $\sqsubseteq \exists$ isPartOf.CirculatorySystem

- Insert into in-memory tables with **fresh ids**
- **Resolve uniqueness of ids using SQL queries when the tables are large enough**

## On-Disk

class	id
Heart	1
MuscularOrgan	2

object	property	id

$A \sqsubseteq B$	$A \sqsubseteq \exists r. B$
1 2	

## In-Memory

class	id
Heart	3
CirculatorySystem	5

object	property	id
	isPartOf	4

$A \sqsubseteq B$	$A \sqsubseteq \exists r. B$
	3 4 5



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

s_t1		ax_t1	
$A \sqsubseteq B$		$A \sqsubseteq B$	
1	1	1	2
2	2	1	4
3	3	2	3

- Repeated application of joins are necessary to compute the closure



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure

s_t1		ax_t1	
$A \sqsubseteq B$		$A \sqsubseteq B$	
1	1	1	2
2	2	1	4
3	3	2	3
1	2		
1	4		
2	3		



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure

s_t1		ax_t1	
$A \sqsubseteq B$		$A \sqsubseteq B$	
1	1	1	2
2	2	1	4
3	3	2	3
1	2		
1	4		
2	3		
1	3		



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure
- Instead one can compute the closure for a part of the table in-memory

s_t1	ax_t1																
<table border="1"><thead><tr><th colspan="2"><math>A \sqsubseteq B</math></th></tr></thead><tbody><tr><td>1</td><td>1</td></tr><tr><td>2</td><td>2</td></tr><tr><td>3</td><td>3</td></tr></tbody></table>	$A \sqsubseteq B$		1	1	2	2	3	3	<table border="1"><thead><tr><th colspan="2"><math>A \sqsubseteq B</math></th></tr></thead><tbody><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>4</td></tr><tr><td>2</td><td>3</td></tr></tbody></table>	$A \sqsubseteq B$		1	2	1	4	2	3
$A \sqsubseteq B$																	
1	1																
2	2																
3	3																
$A \sqsubseteq B$																	
1	2																
1	4																
2	3																
<table border="1"><thead><tr><th colspan="2">tmp</th></tr></thead><tbody><tr><td>1</td><td>1</td></tr></tbody></table>	tmp		1	1													
tmp																	
1	1																



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure
- Instead one can compute the closure for a part of the table in-memory

s_t1	ax_t1																
<table border="1"><thead><tr><th colspan="2"><math>A \sqsubseteq B</math></th></tr></thead><tbody><tr><td>1</td><td>1</td></tr><tr><td>2</td><td>2</td></tr><tr><td>3</td><td>3</td></tr></tbody></table>	$A \sqsubseteq B$		1	1	2	2	3	3	<table border="1"><thead><tr><th colspan="2"><math>A \sqsubseteq B</math></th></tr></thead><tbody><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>4</td></tr><tr><td>2</td><td>3</td></tr></tbody></table>	$A \sqsubseteq B$		1	2	1	4	2	3
$A \sqsubseteq B$																	
1	1																
2	2																
3	3																
$A \sqsubseteq B$																	
1	2																
1	4																
2	3																
<table border="1"><thead><tr><th colspan="2">tmp</th></tr><tr><th colspan="2"><math>A \sqsubseteq B</math></th></tr></thead><tbody><tr><td>1</td><td>1</td></tr><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>4</td></tr></tbody></table>	tmp		$A \sqsubseteq B$		1	1	1	2	1	4							
tmp																	
$A \sqsubseteq B$																	
1	1																
1	2																
1	4																



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure
- Instead one can compute the closure for a part of the table in-memory

s_t1	
$A \sqsubseteq B$	
1	1
2	2
3	3

ax_t1	
$A \sqsubseteq B$	
1	2
1	4
2	3

tmp	
$A \sqsubseteq B$	
1	1
1	2
1	4
1	3



## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure
- Instead one can compute the closure for a part of the table in-memory
- **And output the result into the main table**

s_t1		ax_t1	
$A \sqsubseteq B$		$A \sqsubseteq B$	
1	1	1	2
2	2	1	4
3	3	2	3
1	2		
1	4		
1	3		





## PROBLEM 2: SLOW JOINS

$$\text{CR1 } \frac{A \sqsubseteq B}{A \sqsubseteq C} : B \sqsubseteq C \in \mathcal{O}$$

```
INSERT IGNORE INTO s_t1
SELECT s_t1.A, ax_t1.C
FROM s_t1 JOIN ax_t1
ON s_t1.B = ax_t1.A;
```

- Repeated application of joins are necessary to compute the closure
- Instead one can compute the closure for a part of the table in-memory
- And output the result into the main table
- Repeat similarly for the other parts

s_t1	ax_t1																						
<table border="1"><thead><tr><th colspan="2">A <math>\sqsubseteq</math> B</th></tr></thead><tbody><tr><td>1</td><td>1</td></tr><tr><td>2</td><td>2</td></tr><tr><td>3</td><td>3</td></tr><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>4</td></tr><tr><td>1</td><td>3</td></tr></tbody></table>	A $\sqsubseteq$ B		1	1	2	2	3	3	1	2	1	4	1	3	<table border="1"><thead><tr><th colspan="2">A <math>\sqsubseteq</math> B</th></tr></thead><tbody><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>4</td></tr><tr><td>2</td><td>3</td></tr></tbody></table>	A $\sqsubseteq$ B		1	2	1	4	2	3
A $\sqsubseteq$ B																							
1	1																						
2	2																						
3	3																						
1	2																						
1	4																						
1	3																						
A $\sqsubseteq$ B																							
1	2																						
1	4																						
2	3																						
<table border="1"><thead><tr><th colspan="2">tmp</th></tr></thead><tbody><tr><td>2</td><td>2</td></tr></tbody></table>	tmp		2	2																			
tmp																							
2	2																						



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced

s\_t1

$A \sqsubseteq B$	
1	1
2	2
3	3
1	2
2	3
1	3



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced
- Can be done using one self join and marking the result as non-direct subsumptions.

s\_t1

$A \sqsubseteq B$		
1	1	×
2	2	×
3	3	×
1	2	
2	3	
1	3	×



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced
- Can be done using one self join and marking the result as non-direct subsumptions.
- This results in many on-disk updates since the number of non-direct subsumptions is large

s\_t1

$A \sqsubseteq B$		
1	1	×
2	2	×
3	3	×
1	2	
2	3	
1	3	×



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced
- Can be done using one self join and marking the result as non-direct subsumptions.
- This results in many on-disk updates since the number of non-direct subsumptions is large
- **Instead, transitive reduction can be performed for parts of the table in-memory, marking only direct subsumptions on the disk**

s_t1	
$A \sqsubseteq B$	
1	1
2	2
3	3
1	2
2	3
1	3

tmp

$A \sqsubseteq B$	
1	1
1	2
1	3



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced
- Can be done using one self join and marking the result as non-direct subsumptions.
- This results in many on-disk updates since the number of non-direct subsumptions is large
- **Instead, transitive reduction can be performed for parts of the table in-memory, marking only direct subsumptions on the disk**

s\_t1

$A \sqsubseteq B$	
1	1
2	2
3	3
1	2
2	3
1	3

tmp

$A \sqsubseteq B$		
1	1	×
1	2	
1	3	×



# PROBLEM 3: TRANSITIVE REDUCTION

- To produce the taxonomy, the output table needs to be transitively reduced
- Can be done using one self join and marking the result as non-direct subsumptions.
- This results in many on-disk updates since the number of non-direct subsumptions is large
- **Instead, transitive reduction can be performed for parts of the table in-memory, marking only direct subsumptions on the disk**

s_t1	
$A \sqsubseteq B$	
1	1
2	2
3	3
1	2
2	3
1	3

○

tmp	
$A \sqsubseteq B$	
1	1
1	2
1	3

× ○ ×



# OUTLINE

- 1 INTRODUCTION
- 2 PROCEDURE OUTLINE
- 3 PROBLEMS AND SOLUTIONS
- 4 RESULTS**





# TIMINGS FOR DIFFERENT STAGES (TIME IN SECONDS)

Action	NCI	GO	Galen <sup>-</sup>	Snomed
Loading/Preprocessing	17.85	5.99	23.41	127.51
Completion	5.78	7.29	53.13	783.30
Transitive reduction	10.32	6.10	21.44	249.23
Formating output	1.56	0.98	2.88	23.76
<b>Total</b>	<b>35.51</b>	<b>20.36</b>	<b>100.86</b>	<b>1183.80</b>

- NCI ([www.cancer.gov](http://www.cancer.gov)) contains 27,652 classes
- GO ([www.geneontology.org](http://www.geneontology.org)) contains 20,465 classes
- Galen<sup>-</sup> ([www.co-ode.org/galen](http://www.co-ode.org/galen)) contains 23,136 classes (functionality, inverses, and transitivity removed)
- Snomed ([www.ihtsdo.org](http://www.ihtsdo.org)) contains 315,489 classes

available at:

[db-reasoner.googlecode.com](http://db-reasoner.googlecode.com)



# COMPARISON WITH IN-MEMORY REASONERS (TIME IN SECONDS)

Reasoner	NCI	GO	Galen <sup>-</sup>	Snomed
CB	7.64	1.23	3.36	45.17
CEL v.1.0	3.60	1.02	169.23	1302.18
FaCT++ v.1.3.0	4.60	10.50	—	965.84
HermiT v.0.9.3	70.23	92.76	—	—
<b>DB</b>	<b>35.51</b>	<b>20.36</b>	<b>100.86</b>	<b>1183.80</b>

- CB ([cb-reasoner.googlecode.com](http://cb-reasoner.googlecode.com))
- CEL ([lat.inf.tu-dresden.de/systems/cel/](http://lat.inf.tu-dresden.de/systems/cel/))
- FaCT++ ([owl.man.ac.uk/factplusplus](http://owl.man.ac.uk/factplusplus))
- HermiT ([hermit-reasoner.com](http://hermit-reasoner.com))



## CONCLUSIONS

- $\mathcal{ELH}$  classification is implementable in SQL databases
- Not as simple as it might first seem
- Optimizations are achieved using in-memory processing
- Performance is comparable to existing in-memory reasoners



## CONCLUSIONS

- $\mathcal{ELH}$  classification is implementable in SQL databases
- Not as simple as it might first seem
- Optimizations are achieved using in-memory processing
- Performance is comparable to existing in-memory reasoners
- Does it scale to millions of classes? –Not on a laptop:
  - Snomed x 1 = 20min
  - Snomed x 5 = 4h.30min
  - Snomed x 10 = did not finish overnight



## CONCLUSIONS

- *ELH* classification is implementable in SQL databases
- Not as simple as it might first seem
- Optimizations are achieved using in-memory processing
- Performance is comparable to existing in-memory reasoners
- Does it scale to millions of classes? –Not on a laptop:
  - Snomed x 1 = 20min
  - Snomed x 5 = 4h.30min
  - Snomed x 10 = did not finish overnight

### Future work

- Extension to *OWL 2 EL*
- Tuning the DB engine / testing on a real DB server



# CONCLUSIONS

- *ELH* classification is implementable in SQL databases
- Not as simple as it might first seem
- Optimizations are achieved using in-memory processing
- Performance is comparable to existing in-memory reasoners
- Does it scale to millions of classes? –Not on a laptop:
  - Snomed x 1 = 20min
  - Snomed x 5 = 4h.30min
  - Snomed x 10 = did not finish overnight

## Future work

- Extension to *OWL 2 EL*
- Tuning the DB engine / testing on a real DB server

Please be kind and not ask too difficult questions!

or send them to: [yevgeny.kazakov@comlab.ox.ac.uk](mailto:yevgeny.kazakov@comlab.ox.ac.uk)