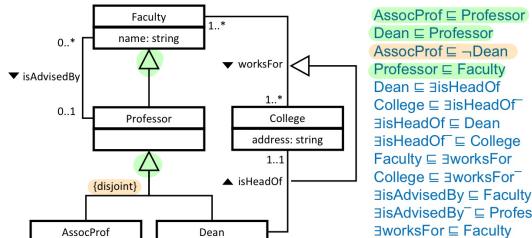


Example of DL-Lite_A TBox

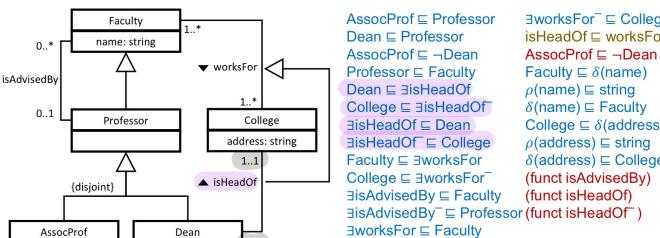


$\exists \text{worksFor} \sqsubseteq \text{College}$
 $\exists \text{isHeadOf} \sqsubseteq \text{worksFor}$
 $\exists \text{AssocProf} \sqsubseteq \text{Professor}$
 $\exists \text{Dean} \sqsubseteq \text{Professor}$
 $\exists \text{AssocProf} \sqsubseteq \neg\text{Dean}$
 $\exists \text{Professor} \sqsubseteq \text{Faculty}$
 $\exists \text{Dean} \sqsubseteq \neg\text{Dean}$
 $\exists \text{AssocProf} \sqsubseteq \text{Faculty}$
 $\exists \text{Faculty} \sqsubseteq \delta(\text{name})$
 $\rho(\text{name}) \sqsubseteq \text{string}$
 $\delta(\text{name}) \sqsubseteq \text{Faculty}$
 $\text{College} \sqsubseteq \delta(\text{address})$
 $\rho(\text{address}) \sqsubseteq \text{string}$
 $\delta(\text{address}) \sqsubseteq \text{College}$
 $\text{Faculty} \sqsubseteq \exists \text{worksFor}$
 $\exists \text{worksFor} \sqsubseteq \text{College}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Faculty}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Professor}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Dean}$
 $\exists \text{worksFor} \sqsubseteq \text{Faculty}$

NOTE: DL-Lite cannot capture completeness of a hierarchy. This would require disjunction (OR).

college
 \downarrow
 $\exists \text{isHeadOf}(o_1, o_2) + \text{esclusione}$
 $\neg \exists \text{isHeadOf}(o_1, o_2)$
 $\exists \text{isHeadOf}(o_1, o_2) \wedge \exists \text{isHeadOf}(o_2, o_1)$

Example of DL-Lite_A TBox



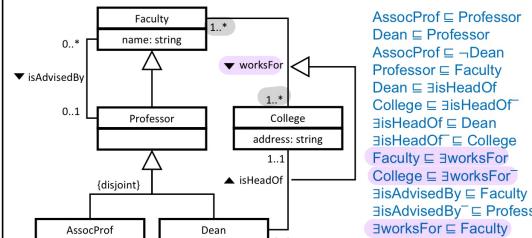
NOTE: DL-Lite cannot capture completeness of a hierarchy. This would require disjunction (OR).

il ruolo ha una corrispondenza esatta con i concetti

- Generalizzazione/specializzazione: esiste una superclasse e le proprie sottoclassi (ereditare)
- nome associazione e direzione (come leggere) -> viene trattato come un ruolo
- nome dell'attributo e il suo tipo primitivo
- molti elementi: numero di oggetti che possono essere associati esattamente a un oggetto del lato opposto

- facciamo in questo modo perché tutte le tuple in *isAdvisedBy* utilizzando sicuramente una faculty e un professore
- l'insieme Faculty possiede tante istanza a cui non è associato alcun professore => tramite il ruolo
- lo stesso ragionamento vale per i professori
- sono elementi non coinvolti nel ruolo

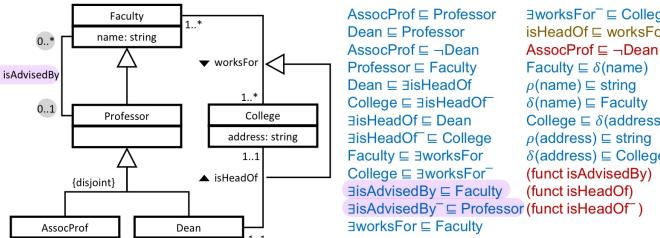
Example of DL-Lite_A TBox



$\exists \text{worksFor} \sqsubseteq \text{College}$
 $\exists \text{isHeadOf} \sqsubseteq \text{worksFor}$
 $\exists \text{AssocProf} \sqsubseteq \text{Professor}$
 $\exists \text{Dean} \sqsubseteq \text{Professor}$
 $\exists \text{AssocProf} \sqsubseteq \neg\text{Dean}$
 $\exists \text{Professor} \sqsubseteq \text{Faculty}$
 $\exists \text{Dean} \sqsubseteq \neg\text{Dean}$
 $\exists \text{AssocProf} \sqsubseteq \text{Faculty}$
 $\exists \text{Faculty} \sqsubseteq \delta(\text{name})$
 $\rho(\text{name}) \sqsubseteq \text{string}$
 $\delta(\text{name}) \sqsubseteq \text{Faculty}$
 $\text{College} \sqsubseteq \delta(\text{address})$
 $\rho(\text{address}) \sqsubseteq \text{string}$
 $\delta(\text{address}) \sqsubseteq \text{College}$
 $\text{Faculty} \sqsubseteq \exists \text{worksFor}$
 $\exists \text{worksFor} \sqsubseteq \text{College}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Faculty}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Professor}$
 $\exists \text{isHeadOf} \sqsubseteq \text{Dean}$
 $\exists \text{worksFor} \sqsubseteq \text{Faculty}$

NOTE: DL-Lite cannot capture completeness of a hierarchy. This would require disjunction (OR).

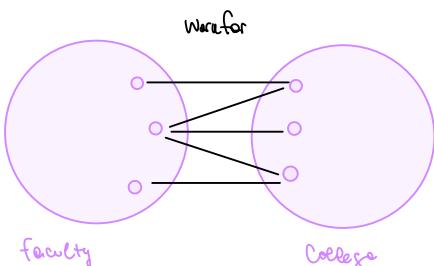
Example of DL-Lite_A TBox



NOTE: DL-Lite cannot capture completeness of a hierarchy. This would require disjunction (OR).

- i due insiemi generati dall'interpretazione non cambiano: Professor e faculty
- quello che dipende dall'associazione è il ruolo
- parto sempre dal ruolo e lo definisco

worksFor



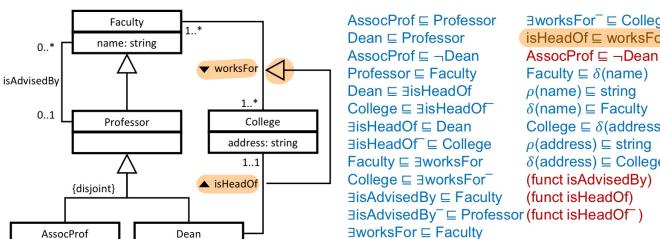
- tutti gli elementi sono coinvolti del ruolo *workFor*
- utilizzo tutte le istanze del concetto *faculty* nel ruolo
- essendo degli insiemi anche se si ripetono le facoltà la tiro fuori solo una volta all'interno dell'associazione
- quindi sicuramente abbiamo una coincidenza a livello insiemistico

tutte le facoltà

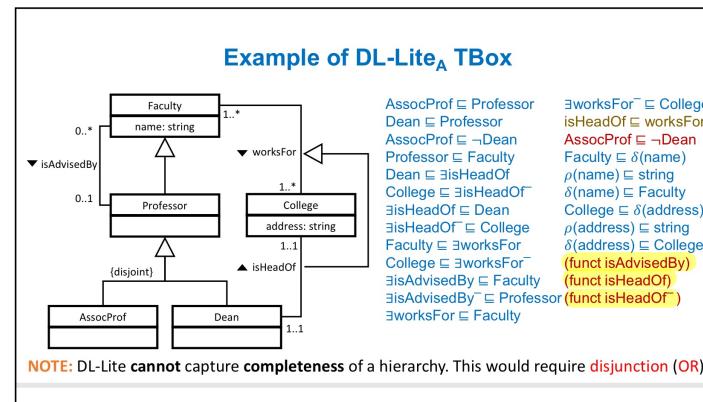
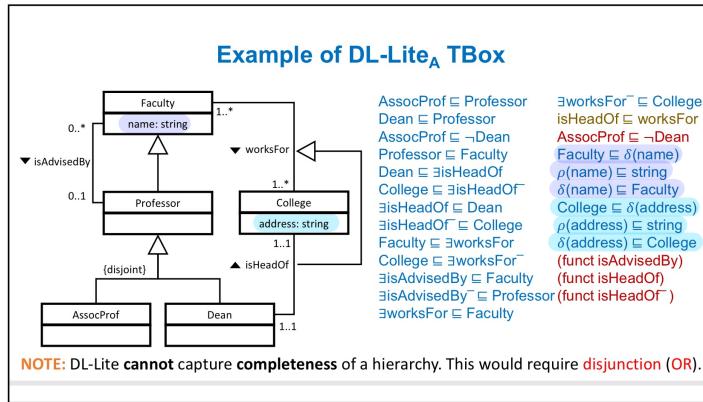
\downarrow
 $\exists \text{workFor}(o_1, o_2)$

tutti i college

Example of DL-Lite_A TBox



NOTE: DL-Lite cannot capture completeness of a hierarchy. This would require disjunction (OR).



$$\text{Sem}(\mathcal{T}, \mathcal{O}) = \{ \mathcal{B} \mid q_i(\mathcal{O}) \subseteq p_G(\mathcal{S}) \text{ for each } \langle q_i, p_G \rangle \in \mathcal{M} \}$$

$$\text{Inv}(\mathcal{T}, \mathcal{O}) = \{ \mathcal{B} \mid q_i(\mathcal{O}) \supseteq p_G(\mathcal{B}) \wedge \langle p_G, p_G \rangle \in \mathcal{M} \}$$

≡

$$\text{FA} \vdash \langle q_S, p \rangle \quad \forall x. \psi_S(x) \simeq g(x)$$

$$(\text{A} \cup \langle s, p_G \rangle) \quad \exists \bar{x}. S(\bar{x}) \rightarrow \psi_S(\bar{x})$$

$$\mathcal{M}(\mathcal{O}) = \{ \mathcal{B}(\bar{c}) \mid \langle q_j, p \rangle \in \mathcal{M} \text{ s.t. } \bar{c} \in q_S(\mathcal{O}) \}$$

- asserzioni funct sono per ruoli e attributi
- isHeadOf:**
 - Ogni Dean ha esattamente un college (perché è funzionale e obbligatorio)
 - Ogni College ha esattamente un Dean (per l'inverso funzionale)
- isAdvisedBy:**
 - per ogni faculty c'è al più Professor
 - per un professor possono esserci più faculty => no funzionale

more repetitive
1a.1

Ex 1)

Exercise

Consider the following information integration specification $J = \langle G, M, S \rangle$, where:

- Global schema $G = \{ \text{Author}(id, name), \text{Book}(code, title, authorName) \}$
- Source schema $S = \{ \text{bestSeller}(bookCode, bookTitle, authorId), \text{award}(authorName, bookTitle, date) \}$
- Mapping $M = \{ m_1: \forall x, y, z, k, w. \text{bestSeller}(x, y, z) \wedge \text{award}(k, y, w) \rightarrow \text{Book}(x, y, k), m_2: \forall x, y, z, k, w. \text{bestSeller}(x, y, z) \wedge \text{award}(k, y, w) \rightarrow \text{Author}(z, k) \}$

- Tell if M is a GAV, LAV, or GLAV mapping.
- Provide the formal definition of retrieved global database for a GAV information integration specification w.r.t. a database.
- Compute the retrieved global database for J w.r.t. D , where D is as follows:
 $D = \{ \text{bestSeller}(1, "Less", 10), \text{bestSeller}(2, "Il cardellino", 20), \text{award}("Greer", "Less", 2018), \text{award}("Tartt", "Il cardellino", 2014), \text{award}("Desiati", "Spatriati", 2020) \}$.
- Given the following CQ over G : $q = \{ () \mid \exists x, y, k, t. \text{Author}(x, y) \wedge \text{Book}(k, t, y) \}$, compute $\text{cert}(q, J, D)$.
- Write at least two LAV mapping between S and G .

1) GAV con operatore GLAV

3) $M(D) = \{ \text{Book}(1, "Less", "Greer"), \text{Book}(2, "Il cardellino", "Tartt"), \text{Author}(1, "Greer"), \text{Author}(2, "Tartt") \}$

9) $\text{cert}(q, J, D) = q(M(D)) = \{ () \} \sqcup \emptyset$

5) $\exists x, y, z. \text{bestSeller}(x, y, z) \rightarrow \exists h. \text{Book}(x, y, h)$

$\wedge \text{Author}(z, h)$

$\exists x, y, z. \text{award}(x, y, z) \rightarrow \exists h, w. \text{Book}(h, w, x)$

$\wedge \text{Author}(z, h)$

Exercise

Consider the following TBox T :

$$T = \{ \text{book} \sqsubseteq \text{product}, \text{e-book} \sqsubseteq \text{product}, \text{product} \sqsubseteq \text{book} \sqcup \text{e-book}, \text{book} \sqsubseteq \exists \text{printedBy}, \exists \text{printBy}^- \sqsubseteq \text{PrintCompany}, \text{e-book} \sqsubseteq \neg \text{book} \}$$

- Tell, motivating your answer, whether TBox T can be expressed in DL-LiteA. If it cannot, write a TBox T' containing only the axioms of T that are expressible in DL-LiteA.
- Given the ABox $A = \{ \text{product}(p1), \text{book}(b1), \text{e-book}(eb1), \text{printedBy}(b2, pc1) \}$, tell whether the ontology (T, A) is satisfiable (consistent), and if so, show an interpretation that is a model for it.
- Given the following conjunctive query $q: \{ (x) \mid \exists y. \text{product}(x) \wedge \text{printedBy}(x, y) \wedge \text{PrintCompany}(y) \}$, compute the perfect rewriting of q w.r.t. T' , where T' is the DL-LiteA TBox provided by you above.
- Given the ABox A above, compute the certain answers of the query q over the ontology (T, A) .
- Add to the TBox T the DL-Lite assertions needed for formalizing the following statement "every e-book is distributed by at least one e-book store" (consider to have the role distributedBy).

logical programming rules (PIs)

$$\text{product}(z) \leftarrow \text{book}(z)$$

$$\text{product}(z) \leftarrow \neg \text{book}(z)$$

$$\text{printedBy}(z, f(z)) \leftarrow \text{book}(z)$$

$$\text{printedBy}(z, f(z)) \leftarrow \text{printedBy}(f(z), z)$$

1) In DL-Lite you can't specialize on obtain concept and also uses conjunction for computational purposes.

$$T' = T / \{ \text{product} \sqsubseteq \text{book} \sqcup \text{e-book} \}$$

2) $\exists N \in T, N: \text{e-book} \sqsubseteq \text{book}$

$$q_N = \{ () \mid \exists x. \text{e-book}(x) \wedge \text{book}(x) \}$$

$$\text{Evol}_{\text{CWA}}(\text{SQL}(q_N), \text{DB}(A)) = \emptyset \Rightarrow O = \langle T, A \rangle$$

addio facilife

$$I = (\Delta^I, \cdot^I) \quad \Delta^I = \{ p_1, b_1, eb_1, b_2, pc_1 \}$$

$$\text{book}^I = \{ b_1 \}$$

$$\text{e-book}^I = \{ eb_1 \}$$

$$\text{product}^I = \{ p_1, b_1, eb_1 \}$$

$$\text{printedBy}^I = \{ (b_2, pc_1), (b_1, pc_1) \}$$

$$\text{printCompany}^I = \{ pc_1 \}$$

$I \models O = \langle T, A \rangle$
fatto le istanze
non si mette

3) logical programming rules (PIs)

$T = \{ \text{book} \sqsubseteq \text{product},$
 $\text{e-book} \sqsubseteq \text{product},$
 ~~$\text{product} \sqsubseteq \text{book} \sqcup \text{e-book}$~~ ,
 $\text{book} \sqsubseteq \exists \text{printedBy},$
 $\exists \text{printeBy}^- \sqsubseteq \text{PrintCompany},$
 ~~$\text{e-book} \sqsubseteq \neg \text{book}$~~ }

$\text{product}(z) \leftarrow \text{book}(z)$
 $\text{product}(z) \leftarrow \text{e-book}(z)$
 $\text{printedBy}(z, f(z)) \leftarrow \text{book}(z)$
 $\text{printCompany}(z_1) \leftarrow \text{printed}(z_2, z_1)$

$$r_{q,T} = \text{perfectRef}(q, T_p) = \{ (x) | \exists y. \text{product}(x) \wedge \text{printedBy}(x, y) \wedge \text{printCompany}(y) \} \cup$$

$$\{ (x) | \exists y. \text{book}(x) \wedge \text{printedBy}(x, z) \wedge \text{printCompany}(z) \} \cup$$

$$\{ (x) | \exists y. \text{e-book}(x) \wedge \text{printedBy}(x, z) \wedge \text{printCompany}(z) \} \cup$$

$$\{ (x) | \exists y, h. \text{product}(x) \wedge \text{printedBy}(x, z) \wedge \text{printedBy}(h, y) \} \cup$$

$$\{ (x) | \exists y, h. \text{book}(x) \wedge \text{printedBy}(x, z) \wedge \text{printedBy}(h, y) \} \cup$$

$$\{ (x) | \exists y, h. \text{e-book}(x) \wedge \text{printedBy}(x, z) \wedge \text{printedBy}(h, y) \} \cup$$

$$\{ (x) | \exists y. \text{product}(x) \wedge \text{printedBy}(x, z) \} \cup$$

$$\{ (x) | \exists y. \text{book}(x) \wedge \text{printedBy}(x, z) \} \cup$$

$$\{ (x) | \exists y. \text{e-book}(x) \wedge \text{printedBy}(x, z) \} \cup$$

$$\{ (x) | \exists y. \text{book}(x) \}$$

$$4) \text{cert}(q_1, \langle T, A \rangle) = \text{Eder}_{\text{RWA}}(\text{RL}(r_{q_1, T}), \text{D}(A)) \\ = \{ (b_1) \}$$

$$5) \text{e-book} \sqsubseteq \exists \text{distributed}$$

$$\exists \text{distributed} \sqsubseteq \text{e-book}$$

$$\exists \text{distributed}^- \sqsubseteq \text{e-store}$$

ESF 2)

Exercise

Consider the following information integration specification $J = \langle G, M, S \rangle$, where:

- Global schema $G = \{ \text{Researcher}(rid, name), \text{Paper}(pid, pTitle, authorName) \}$
- Source schema $S = \{ \text{conferencePaper}(pid, pTitle, rid), \text{bestPaper}(authorName, pTitle, year) \}$
- Mapping $M = \{ m_1: \forall x, y, z. \text{conferencePaper}(x, y, z) \rightarrow \exists w. \text{Paper}(x, y, w) \wedge \text{Researcher}(z, w) \}$
 $m_2: \forall x, y, z. \text{bestPaper}(x, y, z) \rightarrow \exists v. \text{Paper}(v, y, x) \}$

- Tell, motivating your answer, if M is a GAV, LAV, or GLAV mapping.
- Provide the formal definition of universal solution for an information integration specification w.r.t to a source database D .
- Compute a global instance K that is an universal solution for J w.r.t. D , where D is as follows:
 $D = \{ \text{conferencePaper}(1, "title1", 10), \text{conferencePaper}(2, "title2", 20),$
 $\text{bestPaper}("Bob", "title1", 2010),$
 $\text{bestPaper}("Sam", "title3", 2012), \text{bestPaper}("Bill", "title4", 2021) \}$
- Given the following CQ over G : $q = \{ (y, z) \mid \exists x. \text{Paper}(x, y, z) \}$, compute, motivating your answer, $\text{cert}(q, J, D)$.
- Write at least two (possibly new) GAV mapping between S and G .

1) LAV con Spectre di GLAV

3) $n = \text{ch}(M, D) = \{ \text{Paper}(1, "title1", 1), \text{Researcher}(10, 21)$
 $\text{Paper}(2, "title2", 2), \text{Researcher}(20, 22),$
 $\text{Paper}(23, "title3", "Sam"),$
 $\text{Paper}(24, "title4", "Bill") \}$

dove 2:
non label
necessario

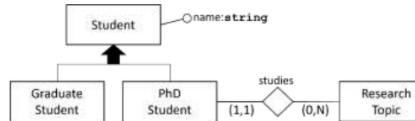
4) $\text{cert}(q, J, D) = q^+(ch(M, D)) = \{ ("title1", "Bob"),$
 $("title3", "Sam"),$
 $("title4", "Bill) \}$

nessuna variabile nella
tuple

5) $\exists x, y, z. \text{conference}(x, y, z) \wedge \text{best}(x, y, z) \rightarrow \text{Paper}(x, y, z)$
 $\exists u, v.$

Exercise

Consider the following Entity-Relationship diagram modelling a portion of the university domain:



- Provide an equivalent formalization with a TBox T in Description Logics of the represented domain.
- Tell, motivating your answer, whether the domain can be formalized through a TBox T' in DL-LiteA only the portion of the domain that can be captured.
- Given the ABox $A = \{ \text{PhdStudent}(sam), \text{Student}(bill), \text{studies}(bill, OBDM) \}$, tell whether the ontology (T, A) is satisfiable (consistent), and if so, show an interpretation that is a model for it.
- Given the following conjunctive query $q: \{ (x) \mid \exists y. \text{Student}(x) \wedge \text{studies}(x, y) \}$, compute the perfect rewriting of q w.r.t. T' , where T' is the DL-LiteA TBox provided by you above.
- Let A be the ABox above, compute the certain answers of the query q over the ontology (T', A) .

1) $T = \{ gs \sqcup phds \subseteq \text{student}$
 $\text{student} \subseteq gs \sqcup phds$
 $phds \subseteq \exists \text{studies}$
 $\exists \text{studies} \subseteq \text{phds}$
 $\exists \text{studies} \subseteq \text{researchT}$
 $\exists \text{name} \subseteq \text{student}$
 $\text{student} \subseteq \delta(\text{name})$
 $\text{phds} \subseteq \text{xsd:string}$
 (funct studies)
 $(\text{funct name}) \}$

2) DL-Lite non escludono le generalizzazioni e i constraints

$$T' = T \setminus \{ \text{student} \subseteq gs \sqcup phds \} \\ \cup \{ ps \subseteq \text{student} \} \\ \cup \{ phds \subseteq \text{student} \}$$

aggiungere anche:
 $gs \subseteq \exists phds$

3) soddisfacibile:

$$\exists f \in T, f_1: (\text{funct studies})$$

$$q_{f_1} = \{ () \mid \exists o, i, j. \text{studies}(o, i) \wedge \text{studies}(o, j) \rightarrow i \neq j \}$$

$$\text{Else}_{\text{QWA}}(\text{SQL}(q_{f_1}), \text{DB}(A)) = \emptyset$$

$$\exists f \in T, f_2: (\text{funct name}) \rightarrow q_{f_2} = \{ () \mid$$

$$q_{f_2} = \{ () \mid \exists o, i, j. \text{name}(o, i) \wedge \text{name}(o, j) \rightarrow i \neq j \}$$

$$\text{Else}_{\text{QWA}}(\text{SQL}(q_{f_2}), \text{DB}(A)) = \emptyset$$

$$\Rightarrow \exists f: f \in O - \langle T, A \rangle$$

da aggiungere qua.

considering $\Sigma = (\Delta^I, \Gamma)$ with $\Delta^I = \{\text{Sau}, \text{Lille}, \text{OBOS}\}$

$$\text{student}^I = \{\text{Lille}, \text{Sau}\}$$

$$gS^I = \{\}$$

$$\text{phdS}^I = \{\text{Sau}, \text{Lille}\}$$

$$\text{research}^I = \{\text{OBOS}\}$$

$$\text{name}^I = \{(\text{Lille}, \text{"Lille"}), (\text{Sau}, \text{"Sau"})\}$$

$$\text{studies}^I = \{(\text{Lille}, \text{OBOS}), (\text{Sau}, \text{OBOS})\}$$

$$\Rightarrow I = \emptyset$$

9) QISD:

$$\text{student}(z) \leftarrow gS(z)$$

$$\text{student}(z) \leftarrow \text{phdS}(z)$$

$$\text{studies}(z_1, f(z)) \leftarrow \text{phdS}(z)$$

$$\text{phdS}(z_1) \leftarrow \text{studies}(z_1, z_2)$$

$$\text{research}(z_1) \leftarrow \text{studies}(z_1, z_2)$$

$$\begin{aligned} \text{perfectRow}(q_1, T_P) = r_{q_1, T} &= \{(x) \mid \exists y. \text{student}(x) \wedge \text{studies}(x, y)\} \cup \\ &= \{(x) \mid \exists y. gS(x) \wedge \text{studies}(x, y)\} \cup \\ &= \{(x) \mid \exists y. \text{phdS}(x) \wedge \text{studies}(x, y)\} \cup \\ &= \{(x) \mid \exists y. \text{student}(x) \wedge \text{phdS}(x)\} \cup \\ &= \{(x) \mid gS(x) \wedge \text{phdS}(x)\} \cup \\ &= \{(x) \mid \text{phdS}(x) \wedge \text{phdS}(\neg x)\} = \{(x) \mid \text{phdS}(x)\} \end{aligned}$$

$$5) \text{ cert}(q_1, \langle T_1, \Delta \rangle) = \text{Ev}_{CNA}(\text{SEL}(r_{q_1, T}), DB(A)) = \{(Sau), (Lille)\}$$

ESE 3

Exercise

Consider the following information integration specification $J = (G, M, S)$, where:

- Global schema $G = \{ \text{Patient}(pid, pName, pAge), \text{Hospitalization}(pid, doctor, date) \}$
- Source schema $S = \{ \text{MediacalRecord}(pid, pName, hospDate, doctoid) \}$
- Mapping $M = \{ m_1: \forall x, y, z, k. \text{MediacalRecord}(x, y, z, k) \rightarrow \exists w. \text{Patient}(x, y, w), m_2: \forall x, y, z. \text{MediacalRecord}(x, y, z, k) \rightarrow \text{Hospitalization}(x, k, z) \}$

- Tell, motivating your answer, if M is a GAV, LAV, or GLAV mapping.
- Describe the Naïve Chase Algorithm
- Compute a $\text{ch}(M, D)$, where M is the above mapping and D is as follows:
 $D = \{ \text{MediacalRecord}(10, "Rossi", 06/08/2010, 35), \text{MediacalRecord}(20, "Bianchi", 09/11/2012, 35), \text{MediacalRecord}(30, "Verdi", 15/01/2013, 45) \}$.
- Given the following CQ over G : $q = \{(y, x, w) \mid \exists z, k. \text{Hospitalization}(x, y, z) \wedge \text{Hospitalization}(w, y, k)\}$, compute the perfect UCQ rewriting of q w.r.t. J .
- Compute $\text{cert}(q, J, D)$.

1) LAV special case of GLAV

$$3) h = \text{ch}(M, D) = \{ \text{Patient}(10, "Rossi", 31), \\ \text{Patient}(20, "Bianchi", 32), \\ \text{Patient}(30, "Verdi", 43), \\ \text{Hosp.}(10, 35, 06/08/2010), \\ \text{Hosp.}(20, 35, 09/11/2012), \\ \text{Hosp.}(30, 45, 15/01/2013) \}$$

$\exists i$ Jersi
labeled nulls

$$4) \begin{array}{c} (q, J) \\ \text{Rew}_{UCQ} = q_M = \{ (y, x, w) \mid \exists r_1, r_2, r_3. \text{MedicR}(x, r_1, r_2, r_3) \wedge \\ \exists r_4, r_5, r_6. \text{MedicR}(w, r_4, r_5, r_6) \} \end{array}$$

$$9) \begin{array}{c} \text{cert}(q, J, D) = q^{\perp}(\text{ch}(M, D)) = q_M(D) \\ = \{ (35, 10, 10), (35, 20, 20), (45, 30, 30), \\ (35, 10, 20), (35, 20, 10) \} \end{array}$$

considero
I rebello di
MedicR con le
dellese type

Exercise

Consider the following domain description about electronic devices:

"We know that both phones and computers are devices, and that a device that is both a phone and a computer is a smartphone. As we also know, a smartphone is both a phone and a computer. Furthermore, we know that each phone has at least one number associated with it".

- Provide a Description Logic TBox T that formalizes all the aspects in the domain description.
- Tell, giving motivation for the answer, whether the domain can be formalized through a TBox T' in DL-LiteA or not. If it cannot, express in DL-LiteA only the portion of the domain that can be captured.
- Given the ABox $A = \{ \text{Phone}(d1), \text{Computer}(d1) \}$, tell whether the ontology (T, A) is satisfiable (consistent), and if so, show an interpretation that is a model for it.
- Given the following conjunctive query $q: \{ (x) \mid \exists y. \text{Device}(x) \wedge \text{hasNumber}(x, y) \}$, compute the perfect rewriting of q w.r.t. T' , where T' is the DL-LiteA TBox provided by you above.
- Let A' be the ABox $\{ \text{Smartphone}(sm), \text{Computer}(co) \}$, $\text{hasNumber}(de, num) \}$, compute the certain answers of the query q over the ontology (T', A') .

$$1) T = \{ \text{phone} \sqcup \text{pc} \sqsubseteq \text{devices}, \\ \text{phone} \sqcap \text{pc} \sqsubseteq \text{smart}, \\ \text{smart} \sqsubseteq \text{phone} \sqcap \text{pc}, \\ \text{phone} \sqsubseteq \exists \text{hasNumber}, \\ \exists \text{hasNumber} \sqsubseteq \text{phone}, \\ \exists \text{hasNumber} \sqsubseteq \text{Number}, \\ \text{phone} \sqsubseteq \forall \text{pc}, \\ (\text{funct hasNumber}) \}$$

2) No.

phone \sqsubseteq devices
 pc \sqsubseteq devices
 smart \sqsubseteq phone
 smart \sqsubseteq pc

3)

Advanced Data Management & Laboratory

Date: 13/02/2024

9 CFU

Exercise 1

Describe the procedure for answering queries over ontologies through query rewriting. Moreover, given a query language Q and an ontology language L, provide the formal definition of Q-rewritability.

Exercise 2

Consider the following information integration specification $J = \langle G, M, S \rangle$, where:

- Global schema $G = \{ \text{Resources}(eId, \text{surname}), \text{Projects}(\text{name}, \text{budget}), \text{Costs}(\text{project}, \text{resource}, \text{cost}) \}$
- Source schema $S = \{ \text{employees}(eId, \text{surname}, \text{salary}), \text{activities}(aId, \text{name}, \text{budget}), \text{worksIn}(eId, aId) \}$
- Mapping $M = \{ m_1: \forall x, y, w \exists z, k, t. \text{employees}(x, y, z) \wedge \text{worksIn}(x, k) \wedge \text{activities}(k, w, t) \rightarrow \text{Resources}(x, y, w), m_2: \forall x, z, w \exists y, t. \text{employees}(x, y, z) \wedge \text{worksIn}(x, k) \wedge \text{activities}(k, w, t) \rightarrow \text{Costs}(w, x, z), m_3: \forall y, z \exists x. \text{activities}(x, y, z) \rightarrow \text{Projects}(y, z) \}$

- Tell whether M is a GAV, LAV, or GLAV mapping.
- Give the formal definition of Retrieved Global Database and describe its properties with respect to $\text{sem}(J, D)$, where D is a database for S .
- Compute the perfect rewriting w.r.t. J of the query $q = \{(y, z, k) \mid \exists x. \text{Resources}(x, y, z) \wedge \text{Projects}(z, k)\}$
- Compute $\text{cert}(q, J, D)$, where D is as follows: $D = \{ \text{employees}(1, \text{"white"}, 30k), \text{employees}(2, \text{"brown"}, 40k), \text{employees}(3, \text{"smith"}, 50k), \text{activities}(p1, \text{"prj1"}, 100k), \text{activities}(p2, \text{"prj2"}, 200k), \text{activities}(p3, \text{"prj3"}, 300k), \text{worksIn}(1, p1), \text{worksIn}(2, p2) \}$.
- Write at least one (possibly new) LAV mapping between S and G .

Exercise 3

Consider the following DL-LiteA TBox T :

$$T = \{ \begin{array}{l} \text{car} \sqsubseteq \text{bike} \sqsubseteq \text{vehicle}, \\ \text{vehicle} \sqsubseteq \text{car} \sqcup \text{bike}, \\ \text{car} \sqsubseteq \exists \text{ownerOf}, \\ \exists \text{ownerOf} \sqsubseteq \text{person}, \\ \exists \text{ownerOf} \sqsubseteq \text{car}, \\ \text{vehicle} \sqsubseteq \delta(\text{numberOfWheels}), \\ \text{car} \sqsubseteq \neg \text{bike}, \\ (\text{funct } \exists \text{ownerOf}) \end{array} \quad \begin{array}{l} \text{car} \sqsubseteq \text{vehicle}, \\ \text{bike} \sqsubseteq \text{vehicle}, \\ \text{car} \sqsubseteq \exists \text{ownerOf}, \\ \exists \text{ownerOf} \sqsubseteq \text{person}, \\ \exists \text{ownerOf} \sqsubseteq \text{car}, \\ \text{vehicle} \sqsubseteq \delta(\text{numberOfWheels}), \\ \text{car} \sqsubseteq \neg \text{bike}, \\ (\text{funct } \exists \text{ownerOf}) \end{array}$$

- Tell, motivating your answer, whether TBox T can be expressed in DL-LiteA. If it cannot, write a TBox T' containing only the axioms of T that are expressible in DL-LiteA.
- Given the ABox $A = \{ \text{car}(c1), \text{vehicle}(v), \text{ownerOf}(p, v2) \}$, tell whether the ontology (T, A) is satisfiable (consistent), and if so, show an interpretation that is a model for it.
- Given the following conjunctive query $q: \{ (x) \mid \exists y. \text{car}(x) \wedge \text{ownerOf}(y, x) \}$, apply the PerfectRef algorithm to the query q and TBox T' .
- Compute the certain answers of the query q over the ontology (T', A) .
- Add to the TBox T the DL-Lite assertions needed for formalizing the following statement "every bike has exactly a brand" (suppose to have the hasBrand attribute).

(separate sheet)

Exercise 4

Describe the role of l-diversity in data anonymization using k-anonymity.

Exercise 5

Classify the approaches for multidimensional indexing.

EFE 2)

1) GAV \Rightarrow GAV mapping because respects the following def.

$$3) \text{new}_{\text{L}}^{(q,J)} = \bigcup_{q \in J} q = \{ (y, z, k) \mid \exists x, t, s, w, x. \text{employees}(x, y, t) \wedge \text{worksIn}(x, r) \wedge \text{activities}(r, z, s) \wedge \text{activities}(w, z, u) \}$$

$$4) \text{new}_{\text{L}}^{(q,J)}(D) = \bigcup_{q \in J} q(D) = \text{cert}(q, J, D) = \{ ("white", "prj1", 100), ("brown", "prj2", 200) \}$$

$$5) \exists x, y, z. \text{employees}(x, y, z) \rightarrow \exists d, \beta. \text{Resources}(x, y, \alpha) \wedge \text{Costs}(d, \beta, z)$$

EFE 3)

1) T cannot be expressed in DL-LiteA

- an atomic concept cannot be specialized
- distinctions are not reflected \rightarrow computational complexity

$$\begin{aligned} T = \{ & \text{car} \sqsubseteq \text{vehicle}, \\ & \text{line} \sqsubseteq \text{vehicle}, \\ & \text{car} \sqsubseteq \exists \text{ownerOf}, \\ & \exists \text{ownerOf} \sqsubseteq \text{person}, \\ & \exists \text{ownerOf} \sqsubseteq \text{car}, \\ & \text{vehicle} \sqsubseteq \delta(\text{numberOfWheels}), \\ & \text{car} \sqsubseteq \neg \text{line}, \\ & (\text{funct } \exists \text{ownerOf}) \} \end{aligned}$$

2)

Theorem: Θ satisfiability depends only on N_J and functionality assertions.

$\exists N \in T, N : q_N = \{ | | \exists x. \text{car}(x) \wedge \text{line}(x) \}$

$\text{Elab}_{\text{CWA}}(\text{SQL}(q_N), \text{DB}(A)) = \emptyset \quad \leftarrow T_1, A \not\models q_N$

$\exists F \in T, f : q_F = \{ | | \exists o_1, o_2. \text{ownerOf}(o_1, o_2) \wedge \text{ownerOf}(o_2, o_1) \wedge o_1 \neq o_2 \}$

$\text{Elab}_{\text{CWA}}(q_F, \text{DB}(A)) = \emptyset \quad A \not\models q_F$

\Rightarrow doesn't exist objects violating the ontology assertions $\Rightarrow \exists I \models O - \langle T, A \rangle$

$$g = (D^I, \cdot^I) \quad D^I = \{c_1, j_1, p, j_2\}$$

$$\text{vehicle}^I = \{c_1, j_1, j_2\}$$

$$\text{car}^I = \{c_1, j_1\}$$

$$\text{line}^I = \{\}$$

$$\text{ownerOf}^I = \{(p, c_1), (p, c_1)\}$$

$$\text{numOfWheels}^I = \{ (c_1, "0"), (j_1, "1"), (j_2, "1") \}$$

$$\text{person}^I = \{p\}$$

3) $T = \{ \text{car} \sqsubseteq \text{vehicle}, \quad \text{vehicle}(z) \leftarrow \text{car}(z)$
 $\text{line} \sqsubseteq \text{vehicle}, \quad \text{vehicle}(z) \leftarrow \text{line}(z)$
 $\text{car} \sqsubseteq \text{ownerOf}, \quad \text{ownerOf}(f(z), z) \leftarrow \text{car}(z)$
 $\text{ownerOf} \sqsubseteq \text{person}, \quad \text{person}(z_1) \leftarrow \text{ownerOf}(z_1, z_2)$
 $\text{ownerOf} \sqsubseteq \text{car}, \quad \text{car}(z_1) \leftarrow \text{ownerOf}(z_1, z_2)$
 ~~$\text{vehicle} \sqsubseteq \text{numOfWheels}$~~ ,
 ~~$\text{car} \sqsubseteq \text{line}$~~ ,
 ~~$(\text{funct } \text{ownerOf})$~~

$\text{perfectRep}(q, T_p) = f_{q,T} = \{ (x) | \exists y. \text{car}(x) \wedge \text{ownerOf}(y, x) \} \cup$

$= \{ (x) | \exists y, z. \text{ownerOf}(y, x) \wedge \text{ownerOf}(z, y) \} \cup$

$= \{ (x) | \text{car}(x) \wedge \text{car}(x) \} \cup$

$= \{ (x) | \exists y. \text{ownerOf}(y, x) \} \cup$

$= \{ (x) | \text{car}(x) \}$

4) $\text{cert}(q, \langle T, A \rangle) = \text{Elab}_{\text{CWA}}(\text{SQL}(f_{q,T}), \text{DB}(A))$

$= \{ (c_1), (j_2) \}$

5)

$\text{line} \sqsubseteq \text{hasBrand}$

$\text{hasBrand} \sqsubseteq \text{line} ?$

$(\text{funct } \text{hasBrand})$

realte von
ein Objekt zu
einer spezifischen
Marke

Esercizio

Exercise

Consider the following information integration specification $J = \langle G, M, S \rangle$, where:

- Global schema $G = \{ \text{Author}(id, name), \text{Book}(code, title, authorName) \}$
- Source schema $S = \{ \text{bestSeller}(bookCode, bookTitle, authorId), \text{award}(authorName, bookTitle, date) \}$
- Mapping $M = \{ m_1: \forall x,y,z,k,w. \text{bestSeller}(x, y, z) \wedge \text{award}(k, y, w) \rightarrow \text{Book}(x, y, k), \\ m_2: \forall x,y,z,k,w. \text{bestSeller}(x, y, z) \wedge \text{award}(k, y, w) \rightarrow \text{Author}(z, k) \}$

- Tell if M is a GAV, LAV, or GLAV mapping.
- Provide the formal definition of retrieved global database for a GAV information integration specification w.r.t. a database.
- Compute the retrieved global database for J w.r.t. D , where D is as follows:
 $D = \{ \text{bestSeller}(1, "Less", 10), \text{bestSeller}(2, "Il cardellino", 20), \text{award}("Greer", "Less", 2018), \\ \text{award}("Tartt", "Il cardellino", 2014), \text{award}("Desiati", "Spatriati", 2020) \}$.
- Given the following CQ over G : $q = \{() \mid \exists x,y,k,t. \text{Author}(x,y) \wedge \text{Book}(k,t,y) \}$, compute $\text{cert}(q, J, D)$.
- Write at least two LAV mapping between S and G .

Risoluzione

- GAV and a special case of GLAV
- Compute the retrieved global database for J w.r.t. D , where D is as follows:
 - Prendo tutte le mapping asssetions m_1, m_2 e le calcolo sull'istanza D
 - Il risultato verrà aggiunto come fatto a $M(D)$
 - Fisso un fatto in D e lo faccio scorrere su tutti gli altri fino a quando non viene rispettato il mapping su S e creo il nuovo fatto in G . Un volta terminato passo al prossimo fatto successivo in D (devo scorrere tutto fra tutti).

$$M(D) = \{ \text{Book}(1, "Less", "Greer"), \text{Book}(2, "Il cardellino", "Tartt") \\ \text{Author}(10, "Greer"), \text{Author}(20, "Tartt") \}$$

- Compute the $\text{cert}(q, J, D)$:

$$\text{cert}(q, J, D) = \{ \bar{c} \mid \bar{c} \in q(\mathbf{B}) \forall B \in \text{sem}(J, D) \}$$

$$\text{cert}(q, J, D) = \{ () \}$$

- Write at least two LAV mapping between S and G :

$$m_1 : \forall x, y, z. \text{bestSeller}(x, y, z) \rightarrow \exists k. \text{Book}(x, y, k) \wedge \text{Author}(z, k)$$

$$m_2 : \forall x, y, z. \text{award}(x, y, z) \rightarrow \exists k. \exists w. \text{Book}(k, y, x) \wedge \text{Author}(w, x)$$

Exercise

Consider the following DL-Lite TBox T :

$$T = \{ \text{doctor} \sqsubseteq \text{person}, \\ \text{doctor} \sqsubseteq \exists \text{worksFor}, \\ \exists \text{worksFor}^- \sqsubseteq \text{hospital}, \\ \text{medStudent} \sqsubseteq \text{doctor}, \\ \text{medStudent} \sqsubseteq \neg \exists \text{worksFor} \}$$

- Tell whether the TBox T is satisfiable, and if so, show a model for T .
- Given the ABox $A = \{ \text{medStudent}(bob) \}$, tell whether the ontology (T, A) is satisfiable (consistent), and if so, show an interpretation that is a model for it.
- Given the following conjunctive query $q: \{ (x) \mid \exists y. \text{person}(x) \wedge \text{worksFor}(x, y) \wedge \text{hospital}(y) \}$, compute the perfect rewriting of q w.r.t. T .
- Given the ABox $A' = \{ \text{doctor}(sam), \text{person}(bill), \text{worksFor}(bill, H3) \}$, compute the certain answers of the query q over the ontology (T, A') .
- Add to the TBox T the DL-Lite assertions needed for formalizing the following statement “every doctor has a specialization”.

Risoluzione

- Determine whether the TBox T is satisfiable:

T is satisfiable, if it admits at least one model such that $I \models T$.

$$I = (\Delta^I, \cdot^I) \quad \Delta^I = \{ Alice, H1 \}$$

| | |
|---|----------------------------|
| $\text{person}^I = \{ Alice \}$ | |
| $\text{doctor}^I = \{ Alice \}$ | 1 okay |
| $\text{worksFor}^I = \{ (Alice, H1) \}$ | 2 okay |
| $\text{hospital}^I = \{ H1 \}$ | 3 okay |
| $\text{medStudent}^I = \emptyset$ | 4, 5 okay (in modo banale) |

- | | |
|---|---|
| 1) $\text{doctor} \sqsubseteq \text{person}$ | $\text{doctor}^I \subseteq \text{person}^I$ |
| 2) $\text{doctor} \sqsubseteq \exists \text{worksFor}$ | $\text{doctor}^I \subseteq \{ o \mid \exists o'. (o, o') \in \text{worksFor}^I \}$ |
| 3) $\exists \text{worksFor}^- \sqsubseteq \text{hospital}$ | $\{ o' \mid \exists o. (o, o') \in \text{worksFor}^I \} \subseteq \text{hospital}^I$ |
| 4) $\text{medStudent} \sqsubseteq \text{doctor}$ | $\text{medStudent}^I \subseteq \text{doctor}^I$ |
| 5) $\text{medStudent} \sqsubseteq \neg \exists \text{worksFor}$ | $\text{medStudent}^I \cap \{ o \mid \exists o'. (o, o') \in \text{worksFor}^I \} = \emptyset$ |

I is a model for T so T is satisfiable.

- Is the ontology satisfiable?

We are given the ABox:

$$A = \{ \text{medStudent}(bob) \} \Rightarrow bob^I \in \text{medStudent}^I$$

Since $bob^I \in \text{medStudent}^I$, by (4) we also have $bob^I \in \text{doctor}^I$. Then (2) implies $bob^I \in \exists \text{worksFor}$. This contradicts axiom (5), which requires that $bob^I \notin \exists \text{worksFor}$.

Conclusion:

$\langle \mathcal{T}, \mathcal{A} \rangle$ is **unsatisfiable** (i.e., it admits no model).

There is a contradiction between axioms (2), (4), and (5) when applied to the individual `bob`.

- Compute the perfect rewriting of the query `q`:

| | |
|---|--|
| $doctor \sqsubseteq person$ | $person(z) \leftarrow doctor(z)$ |
| $doctor \sqsubseteq \exists worksFor$ | $worksFor(z, f(z)) \leftarrow doctor(z)$ |
| $\exists worksFor^- \sqsubseteq hospital$ | $hospital(z) \leftarrow worksFor(f(z), z)$ |
| $medStudent \sqsubseteq doctor$ | $doctor(z) \leftarrow medStudent(z)$ |

We apply the perfect rewriting algorithm:

$$\begin{aligned} \text{PerfectRef}(q, T) = r_{q,T} &= \{(x) \mid \exists y. person(x) \wedge worksFor(x, y) \wedge hospital(y)\} \\ &\vee \{(x) \mid \exists y, k. person(x) \wedge worksFor(x, y) \wedge \text{worksFor}(k, y)\} \\ &\vee \{(x) \mid \exists y. person(x) \wedge \text{worksFor}(x, y)\} \\ &\vee \{(x) \mid \exists y. \text{doctor}(x) \wedge worksFor(x, y)\} \\ &\vee \{(x) \mid doctor(x) \wedge \text{doctor}(x)\} = \{(x) \mid doctor(x)\} \end{aligned}$$

NB: per applicare la riduzione, il join viene fatto su una var esistenziale e ne introduco un'altra con l'unificazione.

- Since the q_N is not applicable, we can say that the ontology is satisfiable and the result is the following:

$$\text{cert}(q, \langle \mathcal{T}, \mathcal{A} \rangle) = \text{Eval}_{CWA}(r_{q,T}, \mathbf{DB}(\mathcal{A})) = \{(Sam), (bill)\}$$

- Add to the **TBox**:

every doctor has a specialization

| | |
|--|---|
| $doctor \sqsubseteq \delta(specilization)$ | ogni dottore é un dottore specializzato |
| $\rho(specilization) \sqsubseteq xsd : string$ | la specializzazione é espressa in stringa |
| $\delta(specilization) \sqsubseteq doctor$ | i specializzati sono dottori |

8.1 Perplessitá

Conjunctive Mapping – Examples (1)

- Conjunctive **LAV** mapping
 $\forall x, y. \text{Tab1}(x, y) \rightarrow \exists z. \text{Person}(x) \wedge \text{Parent}(y) \wedge \text{hasFather}(x, z)$
- Conjunctive **GAV** mapping
 $\forall x. \forall y. \exists z. \text{Tab1}(x, y) \wedge \text{Tab2}(y, z) \rightarrow \text{hasFather}(x, y)$
- Conjunctive **GLAV** mapping
 $\forall x. \forall y. \exists z. \text{Tab1}(x, y) \wedge \text{Tab2}(y, z) \rightarrow \exists v. \text{hasFather}(x, y) \wedge \text{hasFather}(y, v)$

Checking violations of negative inclusions

For each **NI** N in T we compute a boolean CQ q_N according to the following rules:

| N | q_N |
|---|---|
| $A_1 \sqsubseteq \neg A_2$ | $\Rightarrow \{(\) \exists x. A_1(x) \wedge A_2(x)\}$ |
| $\exists P \sqsubseteq \neg A \text{ or } A \sqsubseteq \neg \exists P$ | $\Rightarrow \{(\) \exists x, y. P(x, y) \wedge A(y)\}$ |
| $\exists P^- \sqsubseteq \neg A \text{ or } A \sqsubseteq \neg \exists P^-$ | $\Rightarrow \{(\) \exists x, y. P(x, y) \wedge A(y)\}$ |
| $\exists P_1 \sqsubseteq \neg \exists P_2$ | $\Rightarrow \{(\) \exists x, y, z. P_1(x, y) \wedge P_2(x, z)\}$ |
| $\exists P_1 \sqsubseteq \neg \exists P^-_2$ | $\Rightarrow \{(\) \exists x, y, z. P_1(x, y) \wedge P_2(z, x)\}$ |
| $\exists P^-_1 \sqsubseteq \neg \exists P_2$ | $\Rightarrow \{(\) \exists x, y, z. P_1(x, y) \wedge P_2(y, z)\}$ |
| $\exists P^-_1 \sqsubseteq \neg \exists P^-_2$ | $\Rightarrow \{(\) \exists x, y, z. P_1(x, y) \wedge P_2(z, y)\}$ |
| $P_1 \sqsubseteq \neg P_2 \text{ or } P^-_1 \sqsubseteq \neg P^-_2$ | $\Rightarrow \{(\) \exists x, y. P_1(x, y) \wedge P_2(x, y)\}$ |
| $P^-_1 \sqsubseteq \neg P_2 \text{ or } P_1 \sqsubseteq \neg P^-_2$ | $\Rightarrow \{(\) \exists x, y. P_1(x, y) \wedge P_2(y, x)\}$ |

Checking violations of functionality assertions

For each **functionality assertions** F in T we compute a Boolean FOL query q_F according to the following rules:

| <i>functionality assertion</i> | q_F |
|--------------------------------|---|
| $(\text{funct } P)$ | $\Rightarrow \{(\) \exists x, y, z. P_1(x, y) \wedge P_2(x, z) \wedge y \neq z\}$ |
| $(\text{funct } P^-)$ | $\Rightarrow \{(\) \exists x, y, z. P_1(y, x) \wedge P_2(z, x) \wedge y \neq z\}$ |

Le associazioni

- Molteplicità:



Un'Azienda impiega molte Persone

Una Persona lavora per un'unica Azienda

Esempio di vincolo di cardinalità

- Ad ogni impiegato sono assegnati da 1 a 5 incarichi
- Ogni incarico è assegnato ad al più 50 impiegati



Notes about the skolem terms:

- Every PI is a logical programming rule (**only concepts and roles inclusions**).
- Skolem terms are used when the TBox contains an existential quantification on the right-hand side of a concept inclusion
- Skolem term that stands in for the anonymous individual
- Skolem terms are only used to satisfy existential quantifiers over roles, not concepts directly.
- You apply Skolem terms in rules where a concept implies the existence of a role

$\text{Professor} \sqsubseteq \exists \text{teaches}$

$\text{teaches}(z, f(z)) \leftarrow \text{Professor}(z)$

- You do not apply Skolem terms when reasoning about a concept inclusion like

$\exists \text{teaches}^- \sqsubseteq \text{Course}$

$\text{Course}(z_1) \leftarrow \text{teaches}(z_2, z_1)$

When can't you apply rules with Skolem terms? Even if the rule uses a Skolem term, you can't always apply it during rewriting. You **can't** apply a rule when:

- The Skolem term can't be unified with a constant
- The skolem term can't be unified with a distinguished variable (output variable)
- The skolem term can't be unified with a join variable

- The disjunction is not allowed in D-Lite in fact can be reformulated, but the conjunction is permitted since is in \mathcal{AL}

$$B \sqsubseteq C_1 \sqcap C_2$$

$$B \sqsubseteq C_1$$

$$B \sqsubseteq C_2$$

$$C_1 \sqcup C_2 \sqsubseteq B$$

$$C_1 \sqsubseteq B$$

$$C_2 \sqsubseteq B$$

To preserve DL-Lite's good computational properties, the disjunction in the original TBox must be removed or reformulated using only allowed constructs and also.

- Captures all the basic constructs of UML Class Diagrams and of the ER Model except covering **constraints in generalizations** \Rightarrow atomic concepts can't be specialized

Generalizzazioni/specializzazioni

- si riempie il corpo della freccia freccia se è totale (altrimenti parziale)
- si riempie la punta della freccia nel caso sia esclusiva (non esclusiva)
- **totale:** unione degli insiemi delle istanze delle entità figlie dà l'insieme delle istanze dell'entità genitore
- **esclusiva:** l'intersezione degli insiemi delle istanze delle entità figlie dà l'insieme vuoto

- La maximally-sound rewriting deve avere al più n distinti atomi
- **NB:** la maximally-sound CQ-rewriting della query q_G deve contenere \bar{x} come risposta della query.

The Perfect Rewriting Algorithm

Input: An I.I.S. $J = \langle G, M, S \rangle$ where M is a Conjunctive LAV mapping, and a CQ $q_g = \{ \bar{x} \mid \varphi(\bar{x}) \}$ where φ has n conjuncts.

Output: A UCQ Q over S

```

 $Q \leftarrow \emptyset$ 
For each CQ  $q' = \{ \bar{x} \mid \psi(\bar{x}) \}$  over  $S$  with at most  $n$  conjuncts
  If  $q'$  is a sound rewriting of  $q_g$ , i.e., if  $\exp^{(q',M)} \sqsubseteq q_g$ 
     $Q \leftarrow Q \cup \{q'\}$ 
  End If;
End For;
return  $Q$ ;
  
```

9 Domande Savo

- Let $\langle J, D \rangle$ be information integration system, where $J = \langle G, M, S \rangle$. Provide the formal definition of when a global database B satisfies M w.r.t. D in case of sound, complete, and exact mappings semantics.

Given an $\langle J, D \rangle$, where $J = \langle G, M, S \rangle$ and D is a source database for S :

- Sound-mapping semantics:** the sources contain *partial* but *correct* information regarding the global schema:

$$\text{sem}(J, D) = \{B \mid q_S(D) \subseteq q_G(B) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

- Complete-mapping semantics:** the sources contain *complete* but possibly *incorrect* information:

$$\text{sem}(J, D) = \{B \mid q_S(D) \supseteq q_G(B) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

- Exact-mapping semantics:** the sources contain *exact* information regarding the global schema:

$$\text{sem}(J, D) = \{B \mid q_S(D) \equiv q_G(B) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

Exact and complete mapping are the most reasonable semantics for IIS. However, complete and exact mappings can exhibit very counterintuitive behaviour and generate inconsistent systems (with no semantics).

• Tell if M is a GAV, LAV, or GLAV mapping.

This is a **GAV** mapping M which means is composed by a finite set of **GAV** mapping assertions of this form:

$$\forall m = \langle q_S, g \rangle \in M$$

where q_S is a FOL query over the source schema S and $g \in G$ is a single atom of the global schema. The logical formalization of GAV mapping assertions are the following:

$$\forall \bar{x}. \varphi_S(\bar{x}) \rightarrow g(\bar{x})$$

We can also consider this mapping, in a more generic way as **GLAV** mapping M , since **GAV** \subseteq **GLAV**.

• Provide the formal definition of retrieved global database for a GAV information integration specification w.r.t. a database.

Consider an information integration system $\langle J, D \rangle$, where $J = \langle G, M, S \rangle$ is a **GAV** information integration specification, and D is a source database for S . We call the **retrieved global database** for J with respect to D , denoted by $M(D)$, the global database obtained by applying the queries in all mapping assertions in M and transferring to the elements of G the corresponding tuples retrieved from D . Formally:

$$M(D) = \{g(\bar{c}) \mid \langle q_S, g \rangle \in M \text{ and } \bar{c} \in q_S(D)\}.$$

Note that, since mappings are of type **GAV**, the tuples to be transferred to the global schema G are always definite, i.e., they do not contain existentially quantified elements. Also it is easy to prove the following properties:

- Property 1:** $M(D) \in \text{sem}(J, D)$.

- Property 2:** $M(D) \subseteq B$ for each $B \in \text{sem}(J, D)$.

Thus the retrieved global database is the smallest global database in $\text{sem}(J, D)$.

• Provide the formal definition of universal solution for an information integration specification w.r.t to a source database D .

Consider an information integration system $\langle J, D \rangle$, where $J = \langle G, M, S \rangle$ is a **(G)LAV** information integration specification and D is a source database. A global instance K over G is like a global database over G except that it may contain *labeled nulls*.

Definition: a global instance K is a **solution** for J w.r.t. D if $q_S(D) \subseteq q_G(K)$ for each mapping assertion $m = \langle q_S, q_G \rangle$ in M , i.e., $(K, D) \models M$.

Definition: a global instance K is a **universal solution** for J w.r.t. D if K is a **solution** for J w.r.t. D such that, for each $B \in \text{sem}(J, D)$, there exists a homomorphism from K to B .

• Describe the Naïve Chase Algorithm.

Definition: a mapping assertion $\langle q_S, q_G \rangle \in M$ is **triggerable** over D for a tuple of constants \bar{a} (called the trigger) if $\bar{a} \in q_S(D)$.

The definition of Universal Solutions **does not guarantee existence**. Indeed, the definition is not constructive and does not provide a procedure to compute it. **Universal solutions may not be unique**. For this reasons we will use the **chase algorithms** to compute universal solutions. In the literature, there are different versions of this algorithm; we will see the simplest which uses **existential rules**. The algorithm is based on the following intuitions:

- we need to materialize all triggerable mappings
- we can use variables (labeled nulls) to deal with existentially quantified attributes
- we have to materialize in the most general way

The output of the algorithm is denoted with $\text{ch}(M, D)$ and according to the **C theorem** is also a universal solution of J w.r.t. D .

• Let $J = \langle G, M, S \rangle$ be an information integration specification and let L be a class of queries. Provide the formal definition of maximally-sound L-rewriting of a query q w.r.t. J .

Definition: a **sound L-rewriting** of q w.r.t. M is a query q' in L such that, for every database D for S , the following condition holds:

$$\text{if } \bar{a} \in q'(D) \text{ then } \bar{a} \in \text{cert}(q, J, D)$$

Definition: a **sound L-rewriting** q' of q w.r.t. J is a **maximally-sound L-rewriting** of q w.r.t. J if for every database D for S there exists no **sound L-rewriting** q'' of q w.r.t. J such that:

$$q'(D) \subseteq q''(D)$$

- Provide the definition of the query containment problem. Given the queries:

$$q = \{(y) \mid P(x, y) \wedge P(y, 'c')\}$$

$$q' = \{(w) \mid P(w, k)\}$$

Tell, motivating your answer, if there is a **containment relationship** between them.

Definition: let q and q' be two queries of the same arity over a schema \mathbf{S} . Query containment is the problem of checking whether $q(\mathbf{D}) \subseteq q'(\mathbf{D})$, for every database \mathbf{D} for \mathbf{S} . We write $q \sqsubseteq q'$ to denote that q is contained in q' .

Theorem: let q and q' be two CQs of the same arity over \mathbf{S} . We have that $q \sqsubseteq q'$ if and only if there exists a homomorphism h from q' to q such that $h(\bar{x}_1) = \bar{x}_2$.

We can tell that: $q \sqsubseteq q'$ because it exists a homomorphism from q' to q such that:

$$h(w) = y \quad h(k) = c'$$

- Let Q be a query language and L be an ontology language. Provide the definition of Q -rewritability in L .

Query answering can be thought as done in two phases:

- **Perfect rewriting:** produce from \mathbf{q} and the **TBox** \mathcal{T} a new query $r_{q,T}$ (called the perfect rewriting of \mathbf{q} w.r.t. \mathcal{T}).
- **Query evaluation:** evaluate $r_{q,T}$ over the **ABox** \mathcal{A} seen as a complete database (without considering the **TBox** \mathcal{T}). The result of such evaluation is $\text{cert}(\mathbf{q}, \langle \mathcal{T}, \mathcal{A} \rangle)$

Definition [Q-rewritability]: query answering is Q -rewritable if for every \mathcal{T} in \mathcal{L} and query q , the perfect rewriting $r_{q,T}$ of q w.r.t. \mathcal{T} can be expressed in Q .

10 Domande Paraboschi

- Describe the advantages of the multidimensional model.
- Describe the principles used for the application of k-anonymity over large data collections.
- Describe the main ideas of the Mondrian algorithm and the requirements guaranteed on the produced output
- Describe the role of ℓ -diversity in data anonymization using k-anonymity
- Describe the difference between k-anonymity and ℓ -diversity
- Describe the use of Randomized response to support differential privacy
- Classify the approaches for multidimensional indexing
- Multidimensional index per poche dimensioni
- Describe the main motivations behind the development of the key-value model.
- Describe the advantages of the map-reduce paradigm.
- Talk about the CAP theorem

10.1 Describe the advantages of the multidimensional model

- **Conceptual expressiveness:** The multidimensional model provides a natural and semantically rich framework for organizing data in decision support systems. It is based on the separation between facts, representing quantitative measures of interest, and dimensions, which define the contextual axes along which the facts are analyzed. This separation facilitates the interpretation and formulation of analytical queries by end users, aligning with the cognitive process involved in multidimensional reasoning.
- **Support for OLAP operations:** The model is designed to natively support Online Analytical Processing (OLAP), enabling a wide range of operations that are essential for interactive data exploration. These operations include:

- *Slice* — restricts the data to a single value of one dimension. **projection**
- *Dice* — selects a subcube by filtering on multiple dimensions. **selection**
- *Roll-up* — aggregates data by ascending a hierarchy in a dimension. **aggregate**
- *Drill-down* — refines the analysis by descending a hierarchy to more detailed levels.
- *Pivot* — reorients the dimensional view to examine data under different perspectives. **change dimensional orientation**

These operations are expressive, efficient, and essential for multidimensional data analysis in business intelligence contexts.

- **Efficient query performance:** By leveraging precomputed aggregates and materialized views, often organized as multidimensional cubes, the model supports high-performance analytical queries even over large volumes of data. Indexing structures and storage techniques (e.g., MOLAP, ROLAP, HOLAP) further optimize the execution of aggregation queries.
- **Schema modularity and maintainability:** The strict distinction between facts and dimensions enhances the modularity of the data schema. This allows for more straightforward maintenance and evolution of the data warehouse, including the addition of new measures or dimensions without impacting existing queries.
- **Facilitation of user interaction:** The multidimensional model is well aligned with the structure of visual analytics tools and dashboards. This compatibility enables intuitive navigation and manipulation of data by users who may not possess technical expertise in query languages or relational modeling.
- **Temporal and hierarchical analysis:** The model naturally accommodates temporal dimensions and hierarchical relationships within dimensions, which are crucial for trend analysis, time-series comparisons, and aggregated reporting at multiple levels of granularity.
- **Integration with data warehousing:** The model fits seamlessly within the architecture of data warehouses, where data is often preprocessed, cleaned, and transformed into a format optimized for multidimensional access. It supports star and snowflake schemas, which are standard modeling techniques in the warehousing domain.

10.2 Describe the principles used for the application of k-anonymity over large data collections

- **Privacy-preserving motivation:** The application of k -anonymity in large-scale datasets arises from the need to mitigate re-identification risks in the release or analysis of microdata. When datasets contain quasi-identifiers—attributes that, while not directly identifying, can be linked to external information—they pose a significant threat to individual privacy. k -anonymity seeks to prevent such threats by guaranteeing that each record is indistinguishable from at least $k - 1$ other records with respect to the quasi-identifiers.
 - data are de-identified (identifier is removed)
 - combining QI + identifier
 - equivalence class: group of records in a database that share the same values for the QI attributes
- **Definition and theoretical basis:** A dataset satisfies k -anonymity if every equivalence class formed by identical values of quasi-identifiers contains at least k records. This ensures that an adversary cannot uniquely associate a record with an individual with a confidence greater than $1/k$, assuming no access to sensitive attributes.
- **Techniques for achieving k -anonymity:** The enforcement of k -anonymity involves modifying quasi-identifying attributes to form sufficiently large equivalence classes. The two principal techniques are:
 - **Generalization** increases data indistinguishability by replacing specific attribute values with more general categories, thereby coarsening the resolution of the data space.
can be applied at different levels of granularity (cell, row, column)

- **Suppression** removes values or entire records that cannot be adequately generalized without violating utility constraints or resulting in excessive information loss.

These transformations are applied selectively to minimize utility degradation while satisfying the anonymity condition.

- **Domain generalization hierarchies:** To guide the generalization process, domain generalization hierarchies (DGHs) are constructed for each quasi-identifier. A DGH defines permissible abstraction levels, allowing the algorithm to select the minimal generalization level that satisfies the k -anonymity constraint. These hierarchies are typically domain-specific and must be designed in accordance with semantic coherence and analytical needs.
- **Partitioning-based algorithms:** Efficient implementation of k -anonymity over large datasets relies on scalable partitioning strategies. Algorithms such as Mondrian perform multidimensional recursive partitioning, where the dataset is iteratively split along chosen attributes until all resulting partitions satisfy the minimum cardinality requirement. Each partition is then generalized independently to ensure consistency and privacy.
- **Scalability and performance concerns:** Applying k -anonymity to large data collections requires computationally efficient methods due to the combinatorial nature of the generalization space. Heuristics, indexing structures, and distributed processing are often employed to ensure scalability. Nevertheless, the trade-off between anonymity guarantees and data utility remains a critical design consideration.

- **Limitations and context of applicability:** While k -anonymity provides a clear and interpretable privacy guarantee, it is limited in scope. It does not address attribute disclosure or homogeneity attacks, where sensitive values within an equivalence class are uniform or nearly so. Consequently, k -anonymity is often used in combination with more robust models such as ℓ -diversity or t -closeness to achieve stronger privacy protection.

- grants identity disclosure/data privacy
- but no attribute disclosure: homogeneity and background knowledge attacks

10.3 Describe the main ideas of the Mondrian algorithm and the requirements guaranteed on the produced output

- **General overview and purpose:** The Mondrian algorithm is a multidimensional, recursive partitioning method designed to enforce k -anonymity on tabular datasets. It generalizes quasi-identifying attributes such that each record becomes indistinguishable from at least $k - 1$ others within a partition. The key objective is to preserve data utility while satisfying privacy constraints, using a heuristic approach that balances granularity and anonymity.

how it works:

- Each attribute in QI represents a dimension
- Each tuple in PT represents a point in the space defined by QI
- Tuples with the same QI value are represented by giving a multiplicity value to points
- The multi-dimensional space is partitioned by splitting dimensions such that each area contains at least k occurrences of point values
- All the points in a region are generalized to a unique value
- The corresponding tuples are substituted by the computed generalization

- **Handling multidimensional generalization:** Unlike global recoding approaches that generalize values across the entire dataset, Mondrian applies local generalization within each partition. Once the recursion halts for a given partition, all records in that group are generalized uniformly according to the attribute ranges within the partition. This leads to more adaptive and information-preserving generalizations, particularly effective in high-dimensional spaces.
- **Output guarantees:** The algorithm ensures that the resulting dataset satisfies the k -anonymity property by construction. Each equivalence class, defined by the generalized values of quasi-identifiers, contains at least k records. The quasi-identifiers in each class are made identical through generalization, thus achieving indistinguishability among the records in the class.
- **Efficiency and scalability:** Mondrian is designed to be computationally efficient and scalable to large datasets. Its recursive structure and dimension-driven splitting strategy make it suitable for real-world scenarios where data volume and dimensionality pose significant challenges to privacy-preserving transformations.
- **Heuristic nature and information loss:** Although Mondrian does not guarantee minimal information loss, it offers a favorable trade-off by producing high-utility anonymized datasets. The heuristic dimension selection and median split criteria lead to balanced partitions that often yield better data quality than global approaches. Nevertheless, the algorithm is not optimal and its performance may vary depending on data distribution and parameter settings.

10.4 Describe the role of ℓ -diversity in data anonymization using k -anonymity

- identity disclosure: l'attaccante riesce a identificare l'identità della persona
- attribute disclosure: attaccante non sa chi la persona ma scopre un attributo sensibile

- **Motivation for extended protection:** While k -anonymity ensures that each individual is indistinguishable from at least $k - 1$ others based on quasi-identifiers, it does not inherently protect against attribute disclosure. Specifically, if all records in a k -anonymous equivalence class share the same sensitive attribute value, an adversary can infer that value with certainty, defeating the purpose of anonymization. ℓ -diversity was introduced to address this fundamental limitation.

applied to a generalized table T altriimenti posso avere blocchi i-diverse
- **Definition and conceptual foundation:** The ℓ -diversity principle strengthens k -anonymity by requiring that each equivalence class contains at least ℓ “well-represented” distinct values for the sensitive attribute. This condition increases the uncertainty of an adversary attempting to infer sensitive information, even after identifying the equivalence class to which a target individual belongs.
- **Variants of ℓ -diversity:** To account for different data distributions and adversarial models, several formulations of ℓ -diversity have been proposed:
 - **Distinct ℓ -diversity** requires that each equivalence class contains at least ℓ distinct sensitive values.

per ogni classe di equivalenza la distribuzione dei valori deve essere superiore a una soglia log()
 - **Entropy ℓ -diversity** imposes a minimum entropy threshold, ensuring that the distribution of sensitive values in each class is sufficiently diverse. **Eagle**
 - **Recursive (c, ℓ) -diversity** restricts the dominance of the most frequent values by bounding the cumulative frequency of the top values relative to the others.

These variants provide increasing levels of protection in contexts with skewed or semantically clustered sensitive attributes.

- **Integration with k -anonymity:** ℓ -diversity operates as an additional constraint layered on top of k -anonymity. After a dataset is partitioned into k -anonymous equivalence classes, each class is examined for compliance with the ℓ -diversity criterion. If a class fails, further generalization or suppression is applied until the required diversity level is achieved.

- **Privacy guarantees and adversarial resistance:** By requiring multiple distinct values for the sensitive attribute within each group, ℓ -diversity limits the adversary's confidence in correctly guessing the sensitive value of a given individual. This provides resistance to homogeneity attacks and improves robustness against background knowledge.

- **Limitations and further extensions:** Despite its strengths, ℓ -diversity is not without limitations. In datasets where the sensitive attribute has an inherently skewed distribution, satisfying ℓ -diversity may lead to excessive generalization or even infeasibility. Furthermore, it does not account for the semantic similarity between distinct sensitive values. These limitations have motivated the development of more refined models such as t -closeness, which considers the distance between the distribution of sensitive values in each class and their distribution in the overall dataset.

10.5 Describe the difference between k -anonymity and ℓ -diversity

- **Objective and threat model:** k -anonymity and ℓ -diversity address distinct aspects of privacy in microdata publishing. k -anonymity focuses on preventing identity disclosure by ensuring that each individual is indistinguishable from at least $k - 1$ others based on quasi-identifiers. ℓ -diversity extends this protection to mitigate attribute disclosure, ensuring that sensitive attributes remain uncertain even after successful record linkage.

- **Privacy definition:** k -anonymity requires that the projection of the dataset on quasi-identifiers yields equivalence classes of size at least k . ℓ -diversity requires, in addition, that each such class contains at least ℓ well-represented sensitive values, thus increasing the diversity of information an attacker could observe.

- **Vulnerability coverage:** k -anonymity is vulnerable to homogeneity and background knowledge attacks. These occur when all records in an equivalence class share the same sensitive value or when external knowledge can help narrow down the true value. ℓ -diversity addresses these attacks by ensuring variation in the sensitive attributes within each group.

- **Implementation strategy:** Both models typically rely on generalization and suppression to achieve the required anonymity levels. However, ℓ -diversity places stricter requirements on equivalence classes, often necessitating more aggressive generalization or restructuring of the data, particularly in skewed distributions.

- **Data utility trade-off:** k -anonymity generally preserves more data utility because it only enforces indistinguishability on quasi-identifiers. ℓ -diversity introduces additional constraints that may require greater distortion of the data, potentially leading to increased information loss in exchange for stronger privacy guarantees.
- **Formal hierarchy:** ℓ -diversity is a strictly stronger privacy notion than k -anonymity. Every dataset that satisfies ℓ -diversity also satisfies k -anonymity (with $k \geq \ell$), but the converse does not hold. ℓ -diversity is therefore considered a refinement that addresses limitations inherent in k -anonymity.
- **Limitations and extensions:** While ℓ -diversity improves on k -anonymity, it still suffers from limitations such as vulnerability to semantic similarity and difficulty handling skewed sensitive value distributions. These challenges have motivated the development of more nuanced models like t -closeness, which measures the statistical distance between sensitive value distributions.

10.6 Describe the use of Randomized Response to support Differential Privacy

- **introduction to randomized response:** randomized response is a privacy-preserving technique originally proposed by Warner (1965) to allow individuals to respond truthfully to sensitive binary questions while retaining plausible deniability. It introduces intentional random noise into individual responses to obscure the true values, thus protecting respondents' privacy.
- **mechanism definition:** given a population of n individuals, each with a binary value $x_i \in \{0, 1\}$ representing their truthful answer to a sensitive question, the randomized response mechanism produces an output y_i for each individual, according to the following probabilistic rule:

$$y_i = \begin{cases} x_i & \text{with probability } \frac{e^\varepsilon}{1+e^\varepsilon} \\ \bar{x}_i & \text{with probability } \frac{1}{1+e^\varepsilon} \end{cases} \quad \begin{array}{l} \text{rumore che introduce} \\ \text{più epsilon è grande meno rumore introduce} \\ \text{perché la probabilità che } y_i = x_i \text{ è altissima} \end{array}$$

where $\varepsilon \geq 0$ is the privacy parameter and $\bar{x}_i = 1 - x_i$ is the negated value of the true response. The mechanism is denoted by $\text{RR}_\varepsilon(x_1, \dots, x_n) = (y_1, \dots, y_n)$. the noise responses

- **differential privacy guarantee:** this mechanism satisfies ε -differential privacy. The key property is that the probability distribution over outputs does not significantly change when a single input x_i is altered. In particular, for neighboring inputs $x \simeq x'$ and for all sets E of possible outputs,

Intuition behind the definition:
 • The neighbouring relation \simeq captures what is protected
 • The probability bounds capture how much protection we get

$$\Pr[M(x) \in E] \leq e^\varepsilon \Pr[M(x') \in E]$$

Differential Privacy (Dwork et al., 2006, Dwork, 2006)
 - Input space X (with symmetric neighboring relation \simeq)
 - Output space Y (with symmetric neighboring relation \simeq)
 - Privacy parameter $\varepsilon \geq 0$
 - Differential Privacy: $\Pr[M(x) \in E] \leq e^\varepsilon \Pr[M(x') \in E]$ for all sets of outputs E : $y \in Y$ have
 $y \simeq x$ and for all sets of outputs E : $y \in Y$ have
 $y \simeq x'$ and for all sets of outputs E : $y \in Y$ have
 $y \simeq x$

which ensures that the presence or absence of an individual's data has a limited impact on the final output, providing strong privacy guarantees.

- **plausible deniability:** since the reported response y_i may differ from the true value x_i , and this discrepancy is governed by a known probability distribution, it is impossible for an observer to determine with certainty whether a given answer reflects the respondent's actual data. This provides each individual with plausible deniability.
- **utility and unbiased estimation:** despite the introduced noise, statistical utility can be preserved. By averaging the observed responses and applying a correction

factor, an unbiased estimator for the true mean $\frac{1}{n} \sum_{i=1}^n x_i$ can be obtained. the expected error in estimating the mean using the randomized responses \tilde{y}_i satisfies, with high probability:

$$\left| \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n \tilde{y}_i \right| \leq \mathcal{O}\left(\frac{1}{\varepsilon\sqrt{n}}\right)$$

I can apply a correction term to the observed values because I know their distribution

- la formula mostra che possiamo ottenerci la media delle risposte vere con un errore minore anche se ogni singola risposta ha aggiunto rumore alla propria risposta per proteggere la privacy
- l'errore diminuisce all'aumentare del numero di partecipanti aumenta se voglio più privacy cioè epsilon piccolo

indicating that the accuracy improves with larger datasets and stronger privacy parameters (larger ε). ↗ Privacy

- **role in local differential privacy:** randomized response is a canonical example of a *local differential privacy* (LDP) mechanism, where each individual perturbs their data before sharing it with the curator. this eliminates the need for a trusted central party and supports large-scale data collection with formal privacy guarantees.

10.7 Classify the approaches for multidimensional indexing

Multidimensional indexing techniques are essential for supporting efficient content-based retrieval in large data collections, especially when data are represented in terms of feature vectors or spatial descriptors. These techniques can be classified according to the dimensionality of the data they are designed to manage, the data model (vector-based vs. metric-based), and the strategy used for space partitioning.

- **Low-dimensional indexing:** These approaches are typically employed for spatial data applications such as Geographic Information Systems (GIS) or Computer-Aided Design (CAD), where the data have two to four dimensions. The most common indexing methods include:
 - geometric entities
 - normalized

- **k-d Trees:** Binary trees that alternate the splitting dimension at each level, effective for range queries.
- **Point Quadtrees:** Space is recursively partitioned into quadrants; suitable for dynamic insertion of spatial points.
- **R Trees and variants:** Designed for managing geometric objects (e.g., polygons, rectangles), using Minimum Bounding Regions (MBRs) to organize data in a balanced tree. Variants include:
 - * **R*** Trees: Incorporate forced reinsertion to reduce overlap.
 - * **R+ Trees:** Avoid overlap by allowing objects to appear in multiple leaves.

- **High-dimensional indexing:** Used in multimedia applications where feature vectors can have 64 to 500 dimensions. Traditional spatial indexes tend to degrade in performance due to the “curse of dimensionality”. Key structures include:
 - performance degrade as number of dimensionality increases

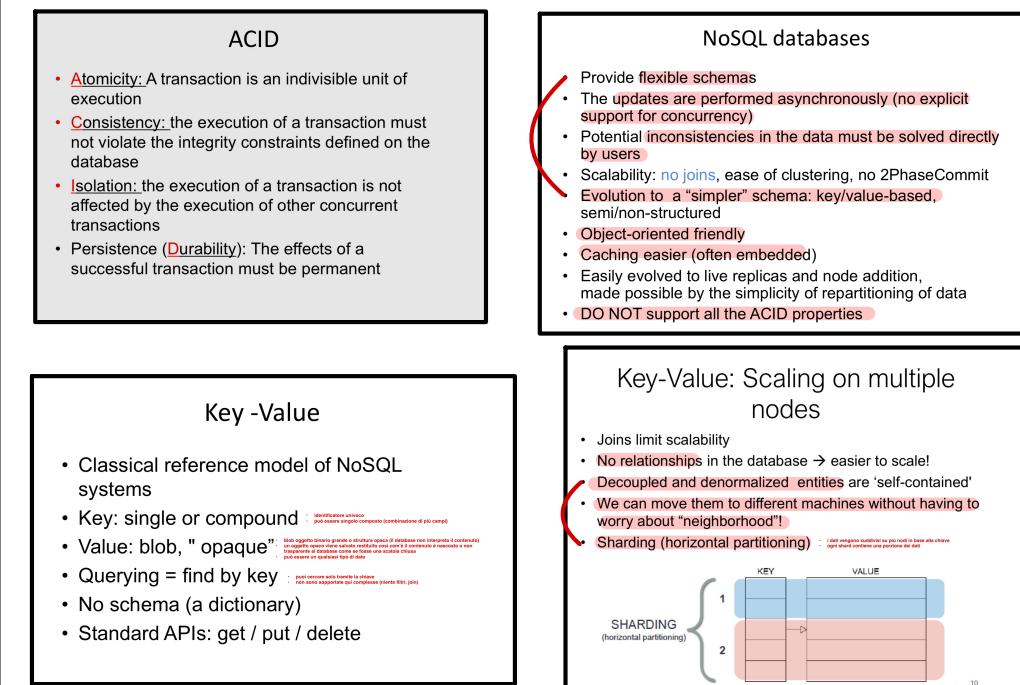
- **SS Trees:** Use Minimum Bounding Spheres (MBS) to form isotropic clusters with fewer dimensions per entry.
- **SR Trees:** Combine bounding spheres and rectangles to balance space efficiency and overlap minimization.

- **Very high-dimensional or metric-based indexing:** These structures rely on metric properties (e.g., triangle inequality) rather than explicit coordinate-based partitioning. They are suitable for domains where distance functions are expensive or non-vectorial.

- **M Trees:** Balanced metric trees that cluster data using routing objects and radii, supporting pruning via metric constraints.
- **Distance-based indexes:** These include general-purpose metric trees that organize data solely based on inter-object distances.
- **Comparison with unidimensional indexing:** While structures like B+ Trees can simulate multidimensional indexing by concatenating multiple features into a composite key, this results in a linear ordering that fails to preserve spatial proximity. In contrast, true multidimensional indexes explicitly exploit geometric or metric nearness.

10.8 Describe the main motivations behind the development of the key-value model

The development of the key-value model within the NoSQL paradigm is primarily motivated by the necessity to handle the increasing complexity and scale of modern data management, particularly in the context of Big Data applications. Classical relational database management systems (RDBMS) offer strong transactional guarantees (ACID properties), rigid schemas, and support for complex queries and relationships. However, these features often impose significant limitations in terms of scalability, flexibility, and performance when applied to large, dynamic, and heterogeneous datasets.



10.9 Describe the advantages of the map-reduce paradigm

AI Overview

MapReduce is a programming model and a software framework for processing and generating large datasets with a parallel, distributed algorithm on a cluster of computers. It's designed to handle big data by breaking down complex tasks into smaller, manageable parts that can be processed concurrently. The core idea is to distribute the data processing across multiple nodes, enabling faster computations compared to traditional sequential processing.

Here's a breakdown of the key concepts:

1. Map Function:
 - The map function takes input data and transforms it into a set of intermediate key-value pairs. 
 - It processes each input element independently, allowing for parallel execution on different nodes.
 - For example, in a word counting application, the map function might take a document as input and output (word, 1) pairs for each word in the document.
2. Reduce Function:
 - The reduce function takes the intermediate key-value pairs generated by the map function and aggregates them based on their keys.
 - It combines the values associated with the same key to produce a smaller set of output key-value pairs.
 - In the word counting example, the reduce function would sum up the counts for each word, resulting in (word, total_count) pairs.
3. Parallel Processing:
 - MapReduce distributes the input data across multiple nodes in a cluster.
 - The map and reduce functions are executed in parallel on these nodes, significantly speeding up the overall processing time.
 - The framework handles the complexities of data distribution, task scheduling, and communication between nodes.
4. Key-Value Pairs:

- MapReduce relies on the concept of key-value pairs for both input and output data.
- This structured format allows for efficient grouping and aggregation of data during the reduce phase.

5. Advantages:

Scalability:

MapReduce can handle massive datasets by distributing the workload across a large number of machines.

Fault Tolerance:

The framework can automatically recover from node failures by rescheduling tasks on other available nodes.

Parallel Processing:

MapReduce enables parallel execution of map and reduce tasks, leading to faster processing times.

Simplified Programming:

The framework handles the complexities of distributed computing, allowing developers to focus on the logic of their map and reduce functions.

6. Evolution:

- While MapReduce was a foundational technology for big data processing, newer frameworks like Apache Spark have emerged with improved performance and features.
- Spark offers in-memory processing, which can significantly reduce data access times compared to MapReduce's disk-based approach.
- However, MapReduce is still used in some legacy applications and remains relevant for understanding the fundamentals of distributed data processing.

- Support for iterative enhancements:** Systems like Apache Spark extend the map-reduce paradigm by enabling in-memory computation and support for advanced features such as directed acyclic graphs (DAGs), thereby improving efficiency in iterative workloads.

In summary, the map-reduce paradigm offers a robust, scalable, and flexible framework for processing massive datasets, making it a cornerstone in the design of modern data-intensive applications in both industry and research.

10.10 Describe the CAP theorem

The CAP theorem, also known as Brewer's theorem, formalizes the trade-offs inherent in the design of distributed data management systems. It states that in a distributed system, it is impossible to simultaneously guarantee all three of the following properties:

- Consistency (C):** All nodes in the system see the same data at the same time. This means that after an update operation completes, all clients accessing the system receive the updated value, regardless of the node they contact.
- Availability (A):** Every request receives a (non-error) response, even if some nodes in the system are unreachable. The system continues to operate and serve client requests under normal and degraded conditions.
- Partition Tolerance (P):** The system continues to function even when there are arbitrary delays or failures in communication between nodes. This includes tolerance to network partitions that may divide the system into disjoint subsets unable to communicate.

According to the CAP theorem, a distributed system can guarantee at most two of these three properties at any given time. This leads to three archetypal system designs:

- CP (Consistency + Partition Tolerance):** These systems maintain consistency and can tolerate network partitions but may sacrifice availability during a partition (e.g., Google BigTable, HBase).
- AP (Availability + Partition Tolerance):** These systems remain available and partition-tolerant but may serve stale or inconsistent data (e.g., Amazon DynamoDB, Cassandra).
- CA (Consistency + Availability):** These systems provide both consistency and availability but cannot function correctly if a network partition occurs. Such configurations are feasible only in non-partitioned or single-node environments.

It is important to note that the CAP theorem applies under the assumption of a network partition. In the absence of partitions, it is possible for a system to exhibit all three properties. Moreover, modern systems often provide tunable consistency and allow developers to choose trade-offs at the granularity of operations, data items, or users.

In conclusion, the CAP theorem highlights the fundamental limitations of distributed systems and guides the design of NoSQL databases and cloud storage infrastructures by encouraging explicit trade-offs among consistency, availability, and partition tolerance.

11 Theory notes

- $\text{sem}(J, \mathbf{D}) = \{\mathbf{B} \mid \mathbf{B} \text{ is a global database such that } (\mathbf{B}, \mathbf{D}) \models M\}$
- $\text{cert}(q, J, \mathbf{D}) = \{(c_1, \dots, c_n) \mid (c_1, \dots, c_n) \in q(\mathbf{B}) \text{ for each } \mathbf{B} \in \text{sem}(J, \mathbf{D})\}$
- **1)** **Theorem:** let $J = \langle G, M, S \rangle$ be an information integration specification, \mathbf{D} be a source database for S , and q be a FOL query over G . Computing $\text{cert}(q, J, D)$ is **undecidable**.

- **Definition:** a mapping assertion m is of the form: $m = \langle q_S, q_G \rangle$ where:

- $q_S = \{\bar{x} \mid \varphi_S(\bar{x})\}$ is a FOL query over S .
- $q_G = \{\bar{x} \mid \varphi_G(\bar{x})\}$ is a FOL query over G .

These mappings have a natural **logical formulation**:

$$\rho : \forall \bar{x}. \varphi_S(\bar{x}) \rightarrow \varphi_G(\bar{x})$$

- Given an $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ and \mathbf{D} is a source database for S :

- **Sound-mapping semantics:** the sources contain *partial* but *correct* information regarding the global schema:

$$\text{sem}(J, \mathbf{D}) = \{\mathbf{B} \mid q_S(\mathbf{D}) \subseteq q_G(\mathbf{B}) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

- **Complete-mapping semantics:** the sources contain *complete* but possibly *incorrect* information:

$$\text{sem}(J, \mathbf{D}) = \{\mathbf{B} \mid q_S(\mathbf{D}) \supseteq q_G(\mathbf{B}) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

- **Exact-mapping semantics:** the sources contain *exact* information regarding the global schema:

$$\text{sem}(J, \mathbf{D}) = \{\mathbf{B} \mid q_S(\mathbf{D}) \equiv q_G(\mathbf{B}) \text{ for each } \langle q_S, q_G \rangle \in M\}.$$

- Inside sound conjunctive mappings, we distinguish three important and widely used classes:

| | | |
|------------------------------|---|--|
| Global-As-View (GAV) | $\forall \bar{x}. \varphi_S(\bar{x}) \rightarrow g(\bar{x}).$ | $\exists a \text{ sx}$ |
| Global-As-View (LAV) | $\forall \bar{x}. s(\bar{x}) \rightarrow \varphi_G(\bar{x}).$ | $\exists a \text{ dx}$ |
| Global-As-View (GLAV) | $\forall \bar{x}. \varphi_S(\bar{x}) \rightarrow \varphi_G(\bar{x}).$ | $\exists a \text{ sx} \wedge a \text{ dx}$ |

in the \forall goes the variables defined on the both part of the mapping.

- **Definition (complete database):** the answers of evaluating the query q over a **database instance** \mathbf{D} is the set $q(\mathbf{D})$, defined as:

$$q = \{(x_1, \dots, x_n) \mid \varphi(x_1, \dots, x_n)\}$$

$$q(\mathbf{D}) = \{\langle c_1, \dots, c_n \rangle \in \mathbf{D}^n \mid \mathbf{D}, \langle c_1, \dots, c_n \rangle \models \varphi(x_1, \dots, x_n)\}$$

- **Definition (missing values):** consider a first-order logic (FOL) query

$$q = \{(x_1, \dots, x_n) \mid \varphi(x_1, \dots, x_n)\}$$

A tuple $\langle c_1, \dots, c_n \rangle \in \text{dom}^n$ is a **certain answer** for q over a database \mathbf{D} if:

$$\text{cert}(q, \mathbf{D}) = \{\langle c_1, \dots, c_n \rangle \mid \langle c_1, \dots, c_n \rangle \in q(v(\mathbf{D})) \text{ for every valuation } v\}.$$

In questo caso utilizzo \exists per le variabili non incluse nella risposta della query, ossia le bounded variables. Infatti se voglio rispondere con tutte le var avrò:

$$q = \{(x_1, \dots, x_n) \mid \varphi(x_1, \dots, x_n)\}$$

$$q_S = \{(n) \mid \exists d, r. \text{Customers}(n, d, r)\}$$

$$q_G = \{(n) \mid \exists i. \text{Person}(i, n)\}$$

Nel caso della formalizzazione logica utilizzo \forall per indicare quelle variabili che compaiono in entrambe le parti della mappatura:

$$m_1 = \langle q_S, q_G \rangle$$

$$\rho : \forall \bar{x}. \varphi_S(\bar{x}) \rightarrow \varphi_G(\bar{x})$$

$$m_1 : \forall n. \exists d. \exists r. \text{Customers}(n, d, r) \rightarrow \exists i. \text{Person}(i, n)$$

11.1 Materialization in GAV

- We want to produce a database $\mathbf{M}(\mathbf{D})$ to represent the result of the integration.
- We want to use the materialization $\mathbf{M}(\mathbf{D})$ to compute certain answers to queries.
- $\mathbf{M}(\mathbf{D})$ is computed by evaluating the source-side queries in each GAV mapping over the source database \mathbf{D} , and inserting the results into the global schema \mathbf{G} .
- We call **retrieved global database** for J w.r.t. \mathbf{D} , denoted by $\mathbf{M}(\mathbf{D})$:

$$\mathbf{M}(\mathbf{D}) = \{g(\bar{c}) \mid \langle q_S, g \rangle \in M \text{ and } \bar{c} \in q_S(\mathbf{D})\}$$

- Since GAV mappings have no existential variables on the right-hand side, all tuples in $\mathbf{M}(\mathbf{D})$ are definite.

Properties:

- **(P1):** $\mathbf{M}(\mathbf{D}) \in \text{sem}(J, \mathbf{D})$
- **(P2):** $\mathbf{M}(\mathbf{D}) \subseteq \mathbf{B}$, for every $\mathbf{B} \in \text{sem}(J, \mathbf{D})$

Hence, $\mathbf{M}(\mathbf{D})$ is the *smallest* global database in $\text{sem}(J, \mathbf{D})$.

- **Theorem:** consider an information integration system $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **CQ-GAV** information integration specification, and \mathbf{D} is a source database. For any **Conjunctive Query (CQ)** q over G , we have:

$$\text{cert}(q, J, \mathbf{D}) = q(\mathbf{M}(\mathbf{D})).$$

- **Theorem:** consider an information integration system $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **FOL-GAV** information integration specification, and \mathbf{D} is a source database. For any **Conjunctive Query (CQ)** q over G , we have:

$$\text{cert}(q, J, \mathbf{D}) = q(\mathbf{M}(\mathbf{D})).$$

11.2 Virtualization in GAV

- Rewrite a user query over the global schema \mathbf{G} into a query over the source schema \mathbf{S} and compute the certain answers by evaluating the rewritten query directly over the source database.
- Virtualization in GAV is very simple:** we replace atoms in the user query with their definition in the mapping \rightarrow *unfolding*.
- Assume a class of queries \mathbf{L} (a language of queries), e.g., **CQs**, **UCQs**, **FOL**, etc.
- Definition:** given a query q over G , a **perfect L-rewriting** of q with respect to J is a query $\text{rew}_{\mathbf{L}}^{(q,J)}$ in the language \mathbf{L} such that the following condition holds for every database \mathbf{D} for S :

$$\bar{a} \in \text{cert}(q, J, \mathbf{D}) \text{ if and only if } \bar{a} \in \text{rew}_{\mathbf{L}}^{(q,J)}(\mathbf{D}).$$

- The **unfolding** of a query q over G with respect to M (denoted by $\text{unf}_{q,M}$) is the query over S obtained from q by substituting every atom $R(\bar{z})$ in q with the corresponding definition $\psi_S(\bar{z})$, such that $\langle \psi_S, R \rangle \in M$.
- The unfolding of a **UCQ** is obtained by unfolding each disjunct individually and then combining the results with a union.
- Theorem:** given an information integration specification $J = \langle G, M, S \rangle$, where M is a **CQ-GAV** mapping, and given a **UCQ** q over G , we have that $\text{unf}_{q,M}$ is a **perfect UCQ-rewriting** of q with respect to J .
- Corollary:** for any **UCQ** q over G , and source database \mathbf{D} for S , we have:

$$\text{cert}(q, J, \mathbf{D}) = \text{unf}_{q,M}(\mathbf{D}).$$

- Theorem:** given an information integration specification $J = \langle G, M, S \rangle$, where M is a **FOL-GAV** mapping, and given a **UCQ** q over G , we have that $\text{unf}_{q,M}$ is a **perfect FOL-rewriting** of q with respect to J .
- Observe:** in this case, $\text{unf}_{q,M}$ is a **FOL** query instead of a **UCQ**.
- Theorem: UCQ-QA-IIIS(CQ-GAV)** is polynomial in **data complexity**, actually **LogSpace**. This result holds also in the case of **UCQ-QA-IIIS(FOL-GAV)** which is polynomial (in **LogSpace**) in **data complexity**.

11.3 Query Answering in (G)LAV Materialization

- Thm:** **UCQ-QA(CQ-GLAV)** is decidable: via materialization and virtualization.
- As for **GAV**, we want to produce a database $\mathbf{M}(\mathbf{D})$ to represent the result of the integration, i.e., a global instance representing all the global databases in $\text{sem}(J, \mathbf{D})$.
- With (G)LAV we cannot simply materialize views from the sources:
 - There is no direct correspondence between definitions and global predicates.
 - Missing information (existential quantification) in the global schema definition
- Intuitively, to materialize GLAV mappings we need **null** values
- A global instance \mathbf{K} over G is similar to a global database over G , but it may contain **labeled nulls**, (placeholders for unknown constants). We denote the set of constants in \mathbf{K} as **Const(\mathbf{K})**, and the set of labeled nulls in \mathbf{K} as **Nulls(\mathbf{K})**. For convenience, we use the following notation for the domain of \mathbf{K} :

$$\text{dom}(\mathbf{K}) = \text{Const}(\mathbf{K}) \cup \text{Nulls}(\mathbf{K})$$

- Definition:** let \mathbf{A} and \mathbf{B} be two database instances over the same schema. A **homomorphism** from \mathbf{A} to \mathbf{B} is a function h from $\text{dom}(\mathbf{A})$ to $\text{dom}(\mathbf{B})$ (where $\text{dom}(\mathbf{A})$ and $\text{dom}(\mathbf{B})$ can both contain constants, variables/labeled nulls) satisfying the following properties:

$$h : \text{dom}(\mathbf{A}) \rightarrow \text{dom}(\mathbf{B})$$

- For every constant $c \in \text{Const}(\mathbf{A})$, $h(c) = c \Rightarrow \text{const}(\mathbf{A}) \subseteq \text{dom}(\mathbf{B})$
- For every relational atom $R(\bar{a}) \in \mathbf{A}$, we have $R(h(\bar{a})) \in \mathbf{B} \Rightarrow h(\mathbf{A}) \subseteq \mathbf{B}$

$$h(\mathbf{A}) = \{R(h(\bar{a})) \mid R(\bar{a}) \in \mathbf{A}\}.$$

Note: For any $x \in \text{Nulls}(\mathbf{A})$, the homomorphism h may map x to either a constant or a labeled null in $\text{dom}(\mathbf{B})$; that is, $h(x) \in \text{dom}(\mathbf{B})$, with no requirement that $h(x) = x$.

- Thm:** the composition $h' \circ h$ is itself a homomorphism h'' from \mathbf{A} to \mathbf{C} .
- Definition:** a global instance \mathbf{K} is a **solution** for J w.r.t. \mathbf{D} if $q_S(\mathbf{D}) \subseteq q_G(\mathbf{K})$ for each mapping assertion $m = \langle q_S, q_G \rangle$ in M , i.e., $(\mathbf{K}, \mathbf{D}) \models M$.
- Definition:** a global instance \mathbf{K} is a **universal solution** for J w.r.t. \mathbf{D} if \mathbf{K} is a **solution** for J w.r.t. \mathbf{D} such that, for each $\mathbf{B} \in \text{sem}(J, \mathbf{D})$, there exists a homomorphism from \mathbf{K} to \mathbf{B} .
- Theorem U:** let \mathbf{K} be a universal solution for J w.r.t. \mathbf{D} . For any **CQ** q , we have

$$\text{cert}(q, J, \mathbf{D}) = q^\perp(\mathbf{K}).$$

- The definition of Universal Solutions **does not guarantee existence**. Indeed, the definition is not constructive and does not provide a procedure to compute it. **Universal solutions may not be unique**.

- For GAV mappings, the retrieved global database is a universal solution.
- The **chase** is a family of **algorithms** that can compute universal solutions. In the literature, there are different versions of this algorithm; we will see the simplest.
- The chase can be used to perform reasoning tasks over **existential rules**.
- **Definition:** an **existential rule** is a first-order logic (FOL) formula of the following form, where φ and ψ are conjunctions of atoms:

$$\forall \bar{x}. \forall \bar{y}. \varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z}. \psi(\bar{x}, \bar{z}).$$

- **Definition:** a mapping assertion $\langle q_S, q_G \rangle \in M$ is **triggerable** over \mathbf{D} for a tuple of constants \bar{a} (called the trigger) if $\bar{a} \in q_S(\mathbf{D})$.

Example: m_1 is triggerable for $\langle a, a \rangle$ but not for $\langle b, c \rangle$

- The chase is based on a set of **intuitions**
 - we need to materialize all triggerable mappings \rightarrow as facts into \mathbf{K}
 - we can use variables (labeled nulls) to fill up the existentially quantified attributes
 - we have to materialize in the most general way

The Naïve Chase Algorithm

Input: $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a GLAV information integration specification and \mathbf{D} is a **source database**

Output: A **global instance** \mathbf{K} over G (i.e., an incomplete database over G)

```

 $\mathbf{K} \leftarrow \emptyset$ 
For each  $m : \forall \bar{x}. \varphi(\bar{x}) \rightarrow \exists \bar{y}. \psi(\bar{x}, \bar{y})$  occurring in  $M$  do:
  For each  $\bar{c} \in q_S(\mathbf{D})$ , where  $q_S = \{\bar{x} \mid \varphi(\bar{x})\}$  do:
    For any  $y_i \in \bar{y}$  pick a fresh labeled null  $n_i \in \bar{n}$ ;
    Add to  $\mathbf{K}$  all the atoms occurring in  $\psi(\bar{c}, \bar{n})$ ;
  End For;
End For;
return  $\mathbf{K}$ ;

```

We denote by $ch(M, \mathbf{D})$ the output \mathbf{K} of the algorithm

- **Theorem C:** given $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **GLAV** information integration specification and \mathbf{D} is a source database, we have that $ch(M, \mathbf{D})$ is a universal solution of J w.r.t. \mathbf{D} .
- **Theorem:** Given $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **GLAV** information integration specification and \mathbf{D} is a source database, and a **CQ** q , we have that

$$cert(q, J, \mathbf{D}) = q^\perp(ch(M, \mathbf{D})).$$

- **Thm:** UCQ-QA(CQ-GLAV) is decidable (polynomial in data complexity).

- **Proof:** for a **GLAV** mapping M and database \mathbf{D} , we can compute a universal solution \mathbf{K} of \mathbf{D} w.r.t. M via the naïve chase (in polynomial time w.r.t. the size of \mathbf{D}) and use the naïve evaluation to answer the query q over \mathbf{K} (in polynomial time w.r.t. the size of \mathbf{D}).

11.4 Query Answering in (G)LAV Virtualization

- We want to produce a query over the sources that captures the results of a user query over the global schema.
- We cannot use simple unfolding techniques with **LAV** mappings (as for **GAV**), since **LAV** does not define the global schema atoms directly $\rightarrow \exists$ quantifier on the right
- **Definition:** given a query q over G , a **perfect L-rewriting** of q w.r.t. J is a query $rew_L^{(q, J)}$ in the language \mathbf{L} such that the following condition holds for every database \mathbf{D} for S :

$$\bar{a} \in cert(q, J, \mathbf{D}) \iff \bar{a} \in rew_L^{(q, J)}(\mathbf{D})$$

- Fix **UCQ** as global query language (L_G) and source query language (L_S).
- **Definition:** we can show that **perfect UCQ-rewriting always exists** and is computable, and we can provide an algorithm for its computation. The algorithm is based on the following **intuitions**:
 - **Sound rewritings:** rewritings that capture only portion of the certain answers
 - **Maximally-sound rewritings:** rewritings that capture as much as possible in a certain class
 - **Query containment:** checking whether the answers for a query are always contained in the answers for another.

- **Definition:** a **sound L-rewriting** of q w.r.t. M is a query q' in \mathbf{L} such that, for every database \mathbf{D} for S , the following condition holds:

$$\bar{a} \in q'(\mathbf{D}) \Rightarrow \bar{a} \in cert(q, J, \mathbf{D})$$

Sound rewritings may not always produce the same results as the original query, they ensure that they **never return incorrect results**, instead it is guaranteed to be a **valid result** for the original query but not the **completeness** of the answer.

- **Observation:** the unsatisfiable query $\{\{\} \mid \perp\}$, (denoted by \perp , representing an empty set) is always a **sound rewriting**, because it doesn't return any incorrect results, there are simply no results at all.
- **Definition:** a **sound L-rewriting** q' of q w.r.t. J is a **maximally-sound L-rewriting** of q w.r.t. J if for every database \mathbf{D} for S there exists no **sound L-rewriting** q'' of q w.r.t. J such that:

$$q'(\mathbf{D}) \subseteq q''(\mathbf{D}) \equiv q' \sqsubseteq q''$$

maximally-sound rewritings are one of the **best possible approximations** of certain answers that one can obtain in a language. Of course, they may not be unique.

- **Definition:** let q_m be the **UCQ** defined as the union of all the **maximally-sound CQ rewritings** of q w.r.t. J .

- Thm: q_m is a **perfect UCQ-rewriting** of q w.r.t. J .
 - **Observation:** in order to build a perfect rewriting we can
 - Compute all maximally sound rewritings of the given query
 - Output the union of all these rewritings
- This is still not an algorithm:
- How can we generate the maximally sound rewritings?
 - How many of them are there?

To answer these questions we must first introduce the following notions: **Query containment** and **Query Expansion**.

- **Definition:** the **expansion** $\exp^{(qs, M)}$ of qs w.r.t. a mapping M is the query over G obtained by replacing each atom $R(\bar{z})$ in qs with the corresponding view q_R over G , where $\langle R, q_R \rangle \in M$.
- **Observe:** $\exp^{(qs, M)}$ is a **CQ**.
- Intuitively, we use M as a reverse mapping from the source schema to the global schema and unfold qs over M in a **GAV** fashion.
- **Definition:** let q and q' be two queries of the same arity over a schema S . **Query containment** is the problem of checking whether $q(D) \subseteq q'(D)$, for every database D for S . We write $q \sqsubseteq q'$ to denote that q is contained in q' .
- **Theorem:** let q and q' be two **CQs** of the same arity over S . We have that $q \sqsubseteq q'$ if and only if there exists a homomorphism h from q' to q such that $h(\bar{x}_1) = \bar{x}_2$.
- **Theorem (check sound rewritings):** qs is a sound rewriting of q_G if and only if $\exp^{(qs, M)} \sqsubseteq q_G$, i.e., $\exp^{(qs, M)}(D) \subseteq q_G(D)$, for every database D for G .
- **Theorem:** let $q_G = \{\bar{x} \mid \varphi(\bar{x})\}$ be a **CQ** over G such that φ consists of n distinct atoms. If a **CQ** $qs = \{\bar{x} \mid \psi(\bar{x})\}$ is a **maximally-sound CQ-rewriting** of q_G w.r.t. J , then ψ is logically equivalent to an existential conjunction of at most n distinct atoms.

The Perfect Rewriting Algorithm

Input: An I.I.S. $J = \langle G, M, S \rangle$ where M is a Conjunctive LAV mapping, and a CQ $q_g = \{\bar{x} \mid \varphi(\bar{x})\}$ where φ has n conjuncts.
Output: A UCQ Q over S

```

 $Q \leftarrow \emptyset$ 
For each CQ  $q' = \{\bar{x} \mid \psi(\bar{x})\}$  over  $S$  with at most  $n$  conjuncts
  If  $q'$  is a sound rewriting of  $q_g$ , i.e., if  $\exp^{(q', M)} \sqsubseteq q_g$ 
     $Q \leftarrow Q \cup \{q'\}$ 
  End If;
End For;
return  $Q$ ;

```

- **Theorem:** let $J = \langle G, M, S \rangle$ be a **LAV** information integration specification and q_G be a **CQ** over G . For input J and q_G , the **perfect rewriting algorithm** computes a perfect **UCQ-rewriting** of q_G w.r.t. J .

- **Proof:** the algorithm computes the union of all the **maximally-sound CQ-rewritings** of the input query.
- **Theorem:** **UCQ-QA(CQ-GLAV)** is decidable, **LOGSPACE** in data complexity.
- **Proof:** For a **GLAV** mapping M , a query q over G , and database D for S :
 - We can compute the perfect **UCQ-rewriting** q' of q w.r.t. J (independently from the size of D)
 - Then we can evaluate q' over D in polynomial time w.r.t. the size of D (actually in **LOGSPACE**)

Materializzazione:

- Calcolare la materializzazione: ossia per ogni mappatura $\langle q_S, q_G \rangle$ importare tutte le tuple $q_S(D)$ all'interno della materialized seguendo la definizione della mappatura (inserisco il fatto q_G), ottenendo $M(D)$ il retrieved global db, oppure K global instance

$$\begin{aligned} & \langle q_S, q_G \rangle \\ & \forall \bar{x}. \varphi_S(\bar{x}) \rightarrow g(\bar{x}) \\ & \bar{c} \in q_S(D) : \forall \bar{x}. \varphi_S(\bar{x}) \\ & K = M(D) = \{g(\bar{c})\} \end{aligned}$$

- Per generare K generico applicare il chase naive algorithm dove $K = ch(M, D)$
- Calcolare la certain:

$$cert(q, JD) = q((M(D)) = q^\perp(ch(M, D))$$

Vitualizzazione:

- Riscrivo le query q su G in query su S impiegando le mapping assertion in M
- In **GAV** basta effettuare $unf_{q, M}$ perché vie una definizione diretta senza **nulls**
- In **(G)LAV** devo computare la q_m impiegando il **the perfect rewriting algorithm**

$$cert(q, J, D) = unf_{q, M}(D) = q_m(D)$$

11.5 Ontology perfect rewriting

Algorithm *PerfectRef(Q; T_P)*
Input: union of conjunctive queries Q , set of DL-Lite_A Pls T_P
Output: union of conjunctive queries PR
 $PR := Q;$
repeat
 $PR' := PR;$
 $\text{for each } q \in PR' \text{ do}$
 $\quad \text{for each atom } g \in q \text{ do}$
 $\quad \quad \text{for each } PI \in T_P \text{ do}$
 $\quad \quad \quad \text{if } PI \text{ is applicable to } g \text{ then } PR := PR \cup \{ \text{ApplyPI}(q, g, PI) \};$
 $\quad \quad \quad \text{for each pair of atoms } g_1, g_2 \in q \text{ do}$
 $\quad \quad \quad \quad \text{if } g_1 \text{ and } g_2 \text{ unify then } PR := PR \cup \{ \text{Reduce}(q, g_1, g_2) \};$
 $\quad \quad \quad \text{until } PR' = PR;$
 $\text{return } PR;$

Observations: 1) Termination follows from having only finitely many different rewritings.
 2) Nls or functionalities do not play any role in the rewriting of the query.

12 Observations

- The architecture of an **IIS** can be defined using **FOL**
- Consider to have instance **D** of a schema **S** over a domain **dom** we can see **D** as a FOL interpretation.
- For each quantifier variable I have to specify it such as: $\exists x. \exists y.$ or $\forall x. \forall y.$
- **cert(q,J,D)** is, the intersection of the answers for q over all the global database in $sem(J, D)$.
- Querying $\langle J, D \rangle$ simply means posing queries over the global schema G of the information integration specification J .
- Complete and exact mappings can exhibit very counterintuitive behavior: may give rise to **inconsistent** systems (i.e., systems with no semantics $sem(J, D) = \emptyset$).
- **1)** The theorem above tells us that we cannot expect to answer queries over an Information integration system with arbitrarily expressive mappings. We need to impose further restrictions on the languages involved \Rightarrow The restriction that we will adopt is to assume mappings that only use **conjunctive queries (CQ)**.
- GAV mappings provide direct information on what data populates the relations of the global schema.
- LAV mappings do not provide direct information on what data populates the global schema. Single elements of the global schema may not be explicitly defined
- We can think of **GLAV mappings** as the *composition* of a **LAV** and a **GAV** mapping. Indeed, each **GLAV mapping** $\langle q_s, q_g \rangle$ can be rewritten as:

$$\langle q_s, p \rangle \cup \langle p, q_g \rangle$$

- From an algorithmic perspective, we can consider LAV and GLAV the same.
- The answer of a query over a database instance **D** is defined as $q(\mathbf{D})$
- we use **null** to represent incomplete information. The common assumption is that **NULL** represents an unknown value. SQL uses this notion of NULL. SQL engines evaluate conditions using three truth values. Tuples evaluating to true are then returned.
- We call a database instance **complete** if it does not contain nulls.
- A more reasoned approach to incomplete information in relational databases is the use of multiple null values. We often refer to this model of nulls as **Marked Nulls**.
- To capture the notion of incomplete information, we use the notion of **valuation**.
- G is an empty theory over the alphabet \mathbf{A}_G , meaning G is a database schema without integrity constraints.
- Assume an (incomplete) database **D** and a query q . We denote by $q(\mathbf{D})$ the evaluation of q over **D**.

- Both a *database instance* and a *query* are defined as sets of *atoms*, also called facts in the case of databases.

A is an instance of a database schema **S**

$$A = \{R(a, b), S(t, v)\}$$

13 Theorems

- **Theorem:** let $J = \langle G, M, S \rangle$ be an information integration specification, **D** be a source database for S , and q be a **FOL** query over G . Computing $cert(q, J, \mathbf{D})$ is **undecidable**. \Rightarrow We need restrictions for the mapping.
- **Theorem:** let q be a **FOL** query. Checking whether $\bar{a} \in cert(q, \mathbf{D})$ can be done in **EXPTIME** in combined complexity and is **NP-Complete** in data complexity.
- **Theorem:** if q is a **Union of Conjunctive Queries (UCQ)**, then:

$$cert(q, \mathbf{D}) = q^\perp(\mathbf{D}).$$

This theorem states that the set of **certain answers** for a UCQ q over **D** is equivalent to the result of the naïve evaluation of q over **D**.

- **Theorem:** let q be a **Union of Conjunctive Queries (UCQ)**. Checking whether $\bar{a} \in cert(q, \mathbf{D})$ is **NP-Complete** in combined complexity and is in **LOGSPACE** in data complexity.
- **Corollary:** query answering for **Union of Conjunctive Queries (UCQs)** is in **PTIME** in data complexity.
- **Definition:** **QA-IIS** is the following decision problem:
 - **Input:** an I.I.S. $J = \langle G, M, S \rangle$, a first-order logic (FOL) query $q(\bar{x})$, a source database **D**, and a tuple of constants \bar{c} .
 - **Question:** is $\bar{c} \in cert(q, J, \mathbf{D})$?
 - **Theorem:** the problem **QA-IIS** is undecidable.
- **Definition (restrictions applied):** **UCQ-QA-IIS(\mathcal{L})**, where \mathcal{L} is a class of mappings, is the following problem:
 - **Input:** an I.I.S. $J = \langle G, M, S \rangle$, where G and S are empty theories, M is a mapping in \mathcal{L} , a source database **D**, a **UCQ** q , and a tuple of constants \bar{c} .
 - **Question:** Is $\bar{c} \in cert(q, J, \mathbf{D})$?
 - **Reply:** we will see, **UCQ-QA-IIS(\mathcal{L})** is decidable in several interesting cases: **CQ-GAV**, **CQ-LAV** and **CQ-GLAV** mappings.
- **Theorem:** consider an information integration system $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **CQ-GAV** information integration specification, and **D** is a source database. For any **Conjunctive Query (CQ)** q over G , we have:

$$cert(q, J, \mathbf{D}) = q(M(\mathbf{D})).$$

- **Theorem:** consider an information integration system $\langle J, \mathbf{D} \rangle$, where $J = \langle G, M, S \rangle$ is a **FOL-GAV** information integration specification, and \mathbf{D} is a source database. For any **Conjunctive Query (CQ)** q over G , we have:

$$\text{cert}(q, J, \mathbf{D}) = q(\mathbf{M}(\mathbf{D})).$$

- **Theorem:** given an information integration specification $J = \langle G, M, S \rangle$, where M is a **CQ-GAV** mapping, and given a **UCQ** q over G , we have that $\text{unf}_{q,M}$ is a **perfect UCQ-rewriting** of q with respect to J .

- **Corollary:** for any **UCQ** q over G , and source database \mathbf{D} for S , we have:

$$\text{cert}(q, J, \mathbf{D}) = \text{unf}_{q,M}(\mathbf{D}).$$

- **Theorem:** given an information integration specification $J = \langle G, M, S \rangle$, where M is a **FOL-GAV** mapping, and given a **UCQ** q over G , we have that $\text{unf}_{q,M}$ is a perfect **FOL-rewriting** of q with respect to J .

- **Theorem:** **UCQ-QA-IIS(CQ-GAV)** is polynomial in ***data complexity***, actually **LogSpace**. This result holds also in the case of **UCQ-QA-IIS(FOL-GAV)** which is polynomial (in **LogSpace**) in ***data complexity***.