# Semantic information retrieval
## A Description logics based approach

Naouel Karam

karam@isima.fr

LIMOS   NEOLIS

# Context

**Classical information retrieval**

- Statistical methods
- Considerable number of results
- Non appropriate results
- Lot of work to the user

**Semantic methods**

- Promising approach
- NLP

# The proposed approach

- **Input**
  - Query Q as a natural language description
  - A set of documents {D1,…,Dn}
- **Problem**
  - Sort {D1,…,Dn} according to their semantic distance w.r.t. Q
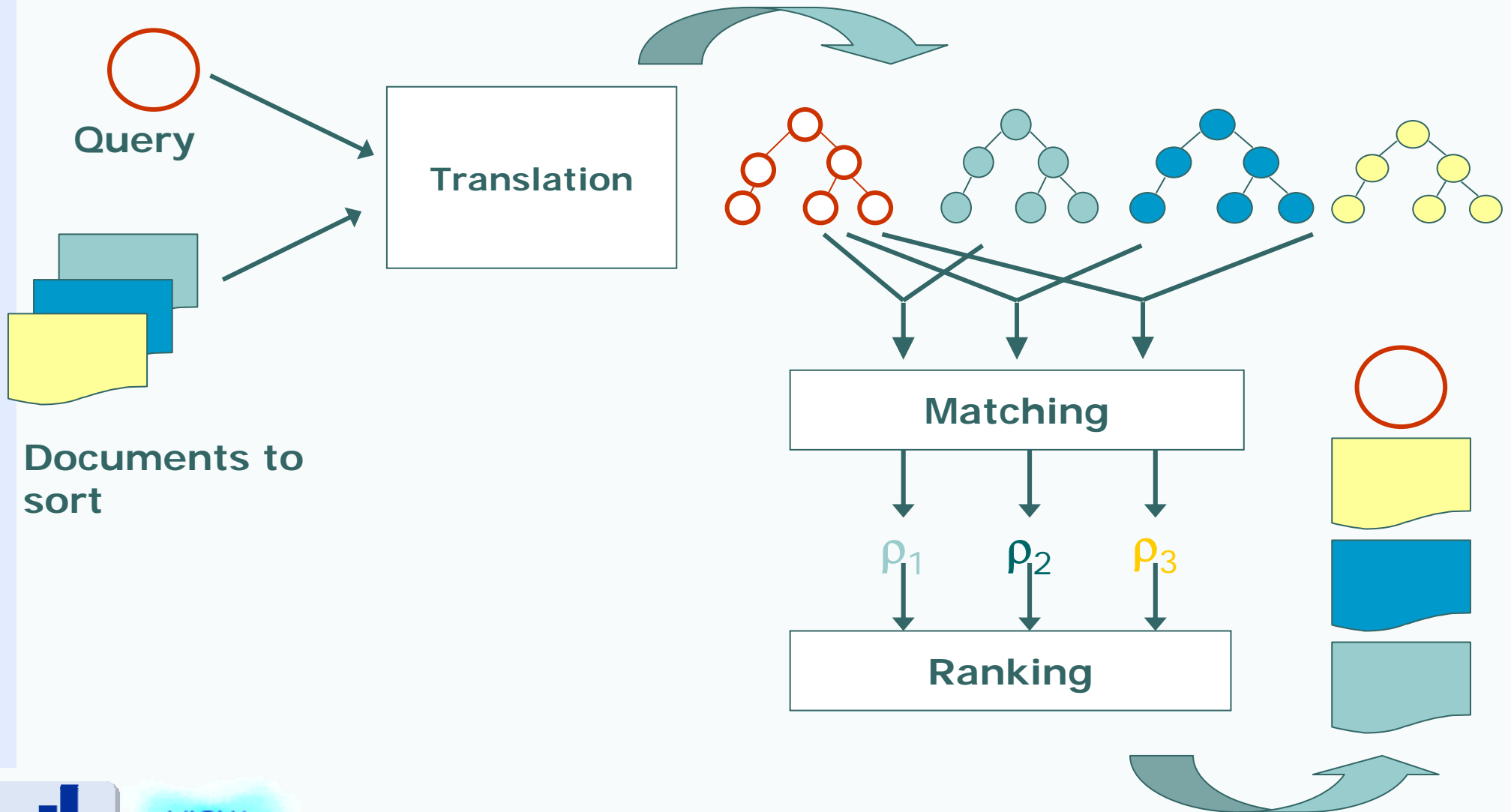  - Need to detect the related information between the query and the documents

Matching Problem

## The proposed approach - An example

**Researchers** are **embedded in** a **laboratory**. They examine guinea pigs and discover **factors** that **give rise** of protein receptors. They study only mice.
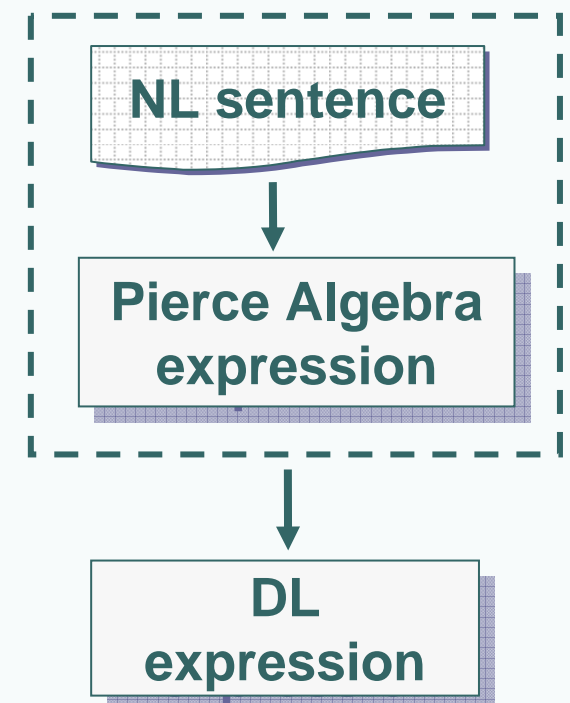
**Scientists** are **attached to** a **research laboratory**. They discover **genes** which produce **specialized** protein receptors. These genes are found in cells.

Query

Translation

Documents to sort

Matching

$\rho_1$ $\rho_2$ $\rho_3$

Ranking

# The translation step
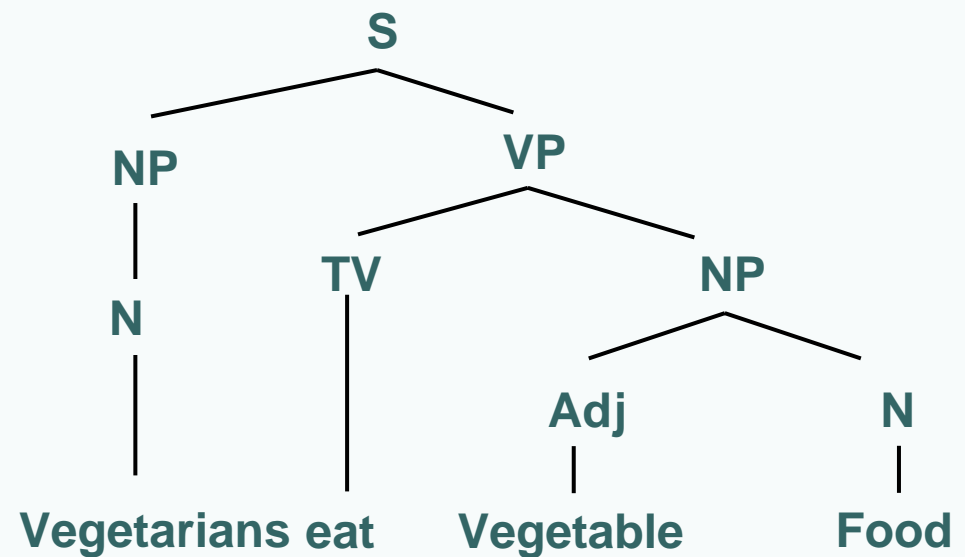
- Formal representation

  - Description Logics (DL)

  - 2 reasons:

    - Well defined semantics / correct algorithms

    - Link with Natural language already established

- Based on existent work [Schmidt 92, 96]

  - Correspondence between syntactic constructions and a semantic representation

  - Connection: Pierce algebras

**NL sentence**

$\downarrow$

**Pierce Algebra expression**

$\downarrow$

**DL expression**

# The translation step – NL and Pierce algebras

- Syntax

| | Production rules |
|---|---|
| (i) | S → NP + VP |
| (ii) | VP → TV + NP |
| (iii) | NP → Adj + N |
| (iv) | NP → N |

# The translation step – NL and Pierce algebras

- **Semantics**

  **[.]** : {N, Adj, NP} → sets

  {TV} → binary relation

| | Production rule | Semantic association |
|---|---|---|
| (i) | S → NP + VP | [NP] ⊆ [VP] |
| (ii) | VP → TV + NP | [TV]:[NP] |
| (iii) | NP → Adj + N | [Adj] ∩ [N] |
| (iv) | NP → N | [N] |

**P: [vegetarian] ⊆ [eat] : ([vegetable] ∩ [food])**

**NP: [vegetarian]**

**VP : [eat] : ([vegetable] ∩ [food])**

**N: [vegetarian]**

**TV: [eat]**

**NP: [vegetable] ∩ [food]**

**Adj: [vegetable]**

**N: [food]**

**vegetarians: [vegetarian]**  **eat : [eat]**  **vegetable: [vegetable]**  **food: [food]**

# The translation step – Pierce algebra and DL

| Pierce algebra | Description logics |
|---|---|
| set | concept |
| Binary relation | role |
| Subset relation ( $\subseteq$ ) | subsumption ( $\sqsubseteq$ ) |
| intersection ( $\cap$ ) | conjunction ( $\sqcap$ ) |
| Pierce product ( : ) | Existential quantification ( $\exists$ ) |

$$[\text{Vegetarian}] \subseteq [\text{eat}] : ([\text{Vegetarian}] \cap [\text{Food}])$$

$$\text{Vegetarian} \sqsubseteq \exists \text{eat}.(\text{Vegetarian} \sqcap \text{Food})$$

## The translation step

- Restricted framework: sentences with complements, quantifiers (all, some, only), number restrictions, negation, passive form.

- Quantifiers

  - « Some persons eat fruit »    Person ◆ ∃ eat.Fruit ⁄B ⊥

  - « All persons eat fruit »    Person ❖ ∃ eat.Fruit

  - « No persons eat fruit »    Person ◆ ∃ eat.Fruit B ⊥

- Number restrictions

  - « John loves more than 3 girls »    John ❖ ≥3 love. Girl

  - « John loves at most 2 girls »    John ❖ ≤2 love. Girl

  - « John loves exactly 1 girl »    John ❖ ≤1 love.Girl ◆ ≥1loves.Girl

# The translation step

- **Relational nouns**
  - « A Father has sons »  Father ❖ ∃ son. T
- **Negation**
  - « is not comfortable »  ¬ comfortable
- **Passive form**
  - « is teached by »  ∃ teach $^{-1}$

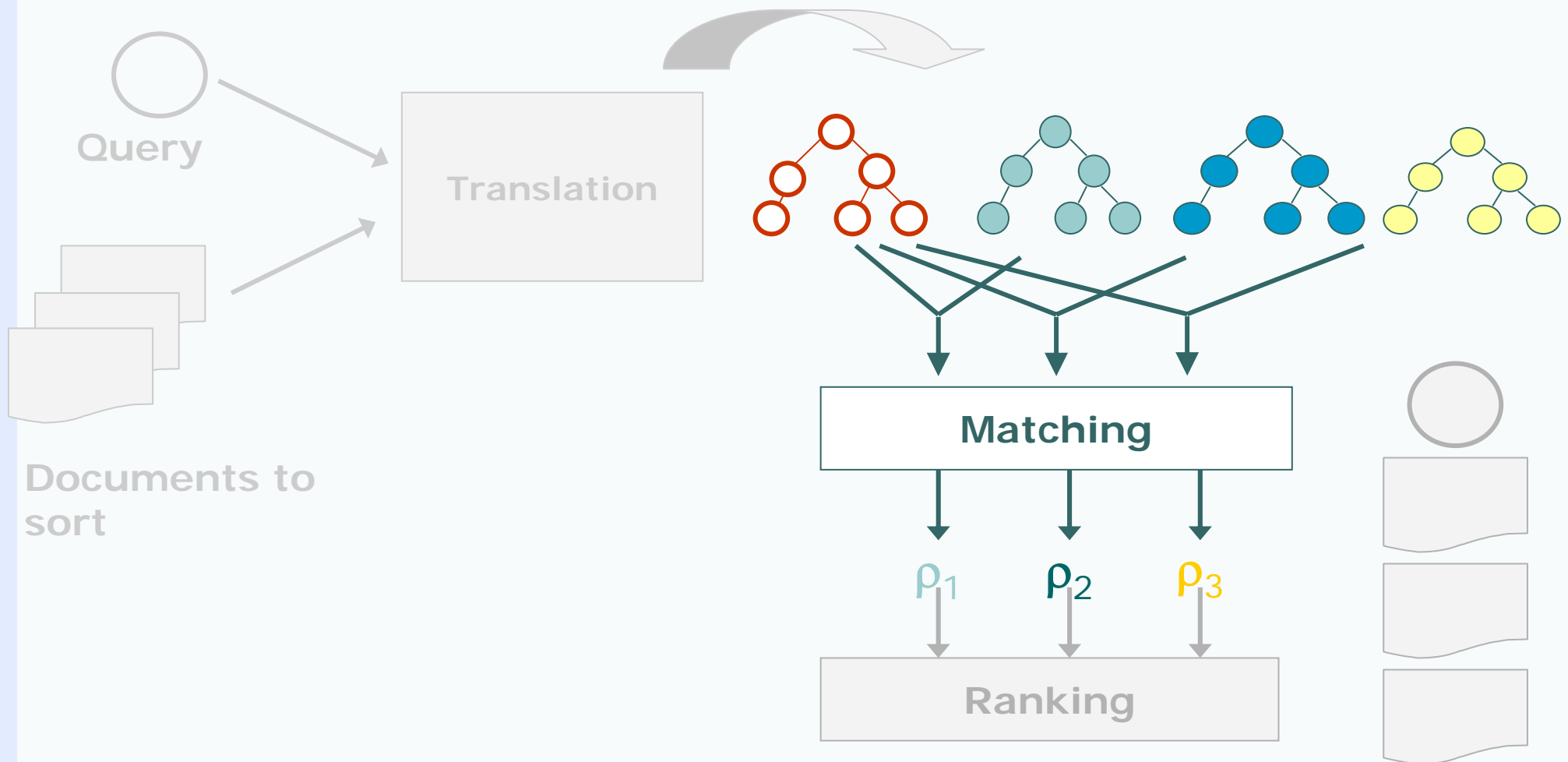# The translation step – An example

Q

French scientists are attached to a research laboratory. They discover genes. These genes are found in cells and produce specialized protein receptors.

$T_Q$

Scientist $\doteq$ French $\sqcap \exists$ attached-to.Research laboratory
$\sqcap \exists$ discover.(Gene $\sqcap \exists$ found-in.Cell $\sqcap \exists$ produce.(Specialized $\sqcap$ Protein receptor)) $\sqcap \overline{\text{Scientist}}$

Gene $\doteq$ $\exists$ found-in.Cell $\sqcap \exists$ produce.(Specialized $\sqcap$ Protein receptor) $\sqcap \overline{\text{Gene}}$
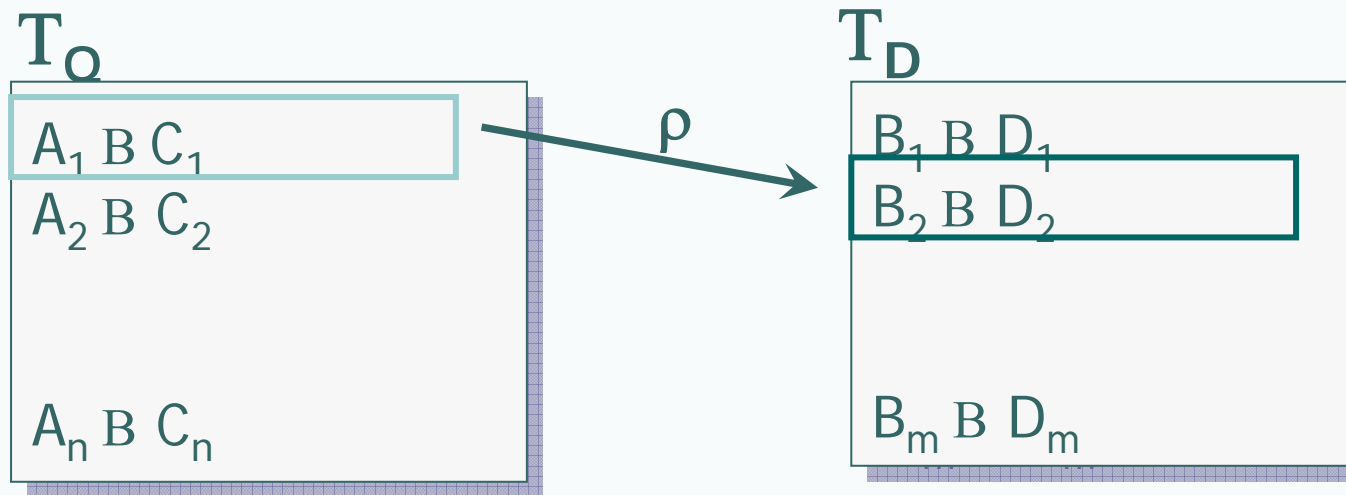
# The global process

# The matching step

- Similar to schema matching problems (Databases, XML,…)

- Existing approaches: schema = tree structure

- Framework : description logics

  - Schema = Terminology

# Matching definition

- Operation that takes two schema as input and returns a correspondence between elements from the two schemas

- Correspondence is a pair of related elements

- Matching terminologies

  - Elements to relate: defined concepts in the terminologies

$T_Q$

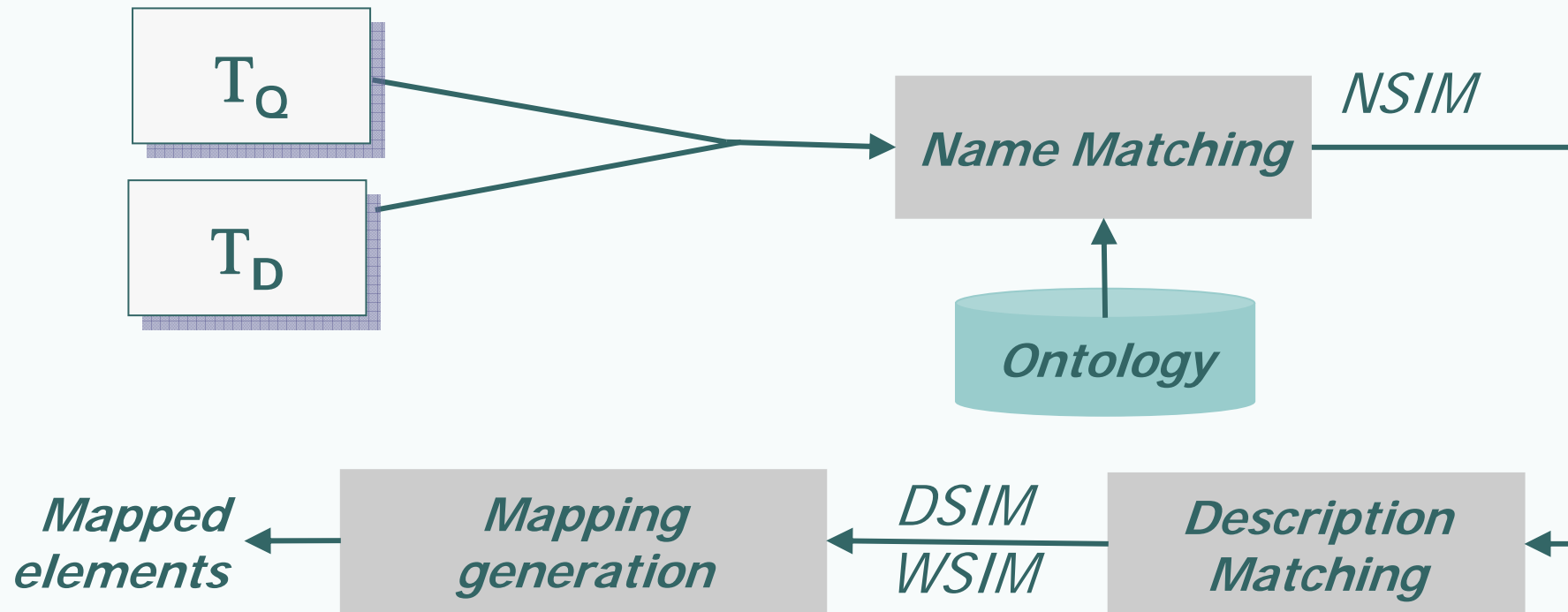| $A_1$ B $C_1$ |
| $A_2$ B $C_2$ |
| |
| $A_n$ B $C_n$ |

$\rho$

$T_D$

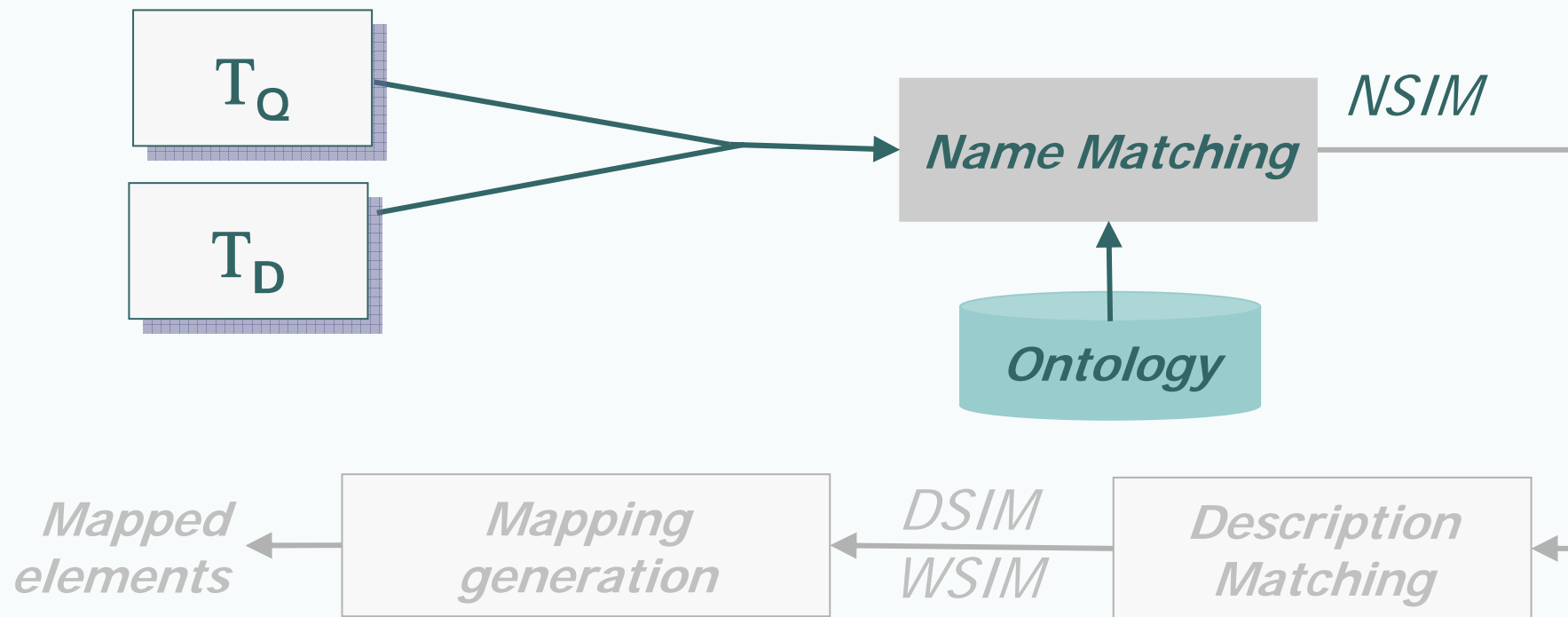| $B_1$ B $D_1$ |
| $B_2$ B $D_2$ |
| |
| $B_m$ B $D_m$ |

Names $A_1$ and $B_2$ are similar

Descriptions $C_1$ et $D_2$ are similar

- 2 steps:
  - **Name matching**
  - **Description matching**

# The matching step

# The matching step



$T_Q$

$T_D$

Name Matching

Ontology

NSIM
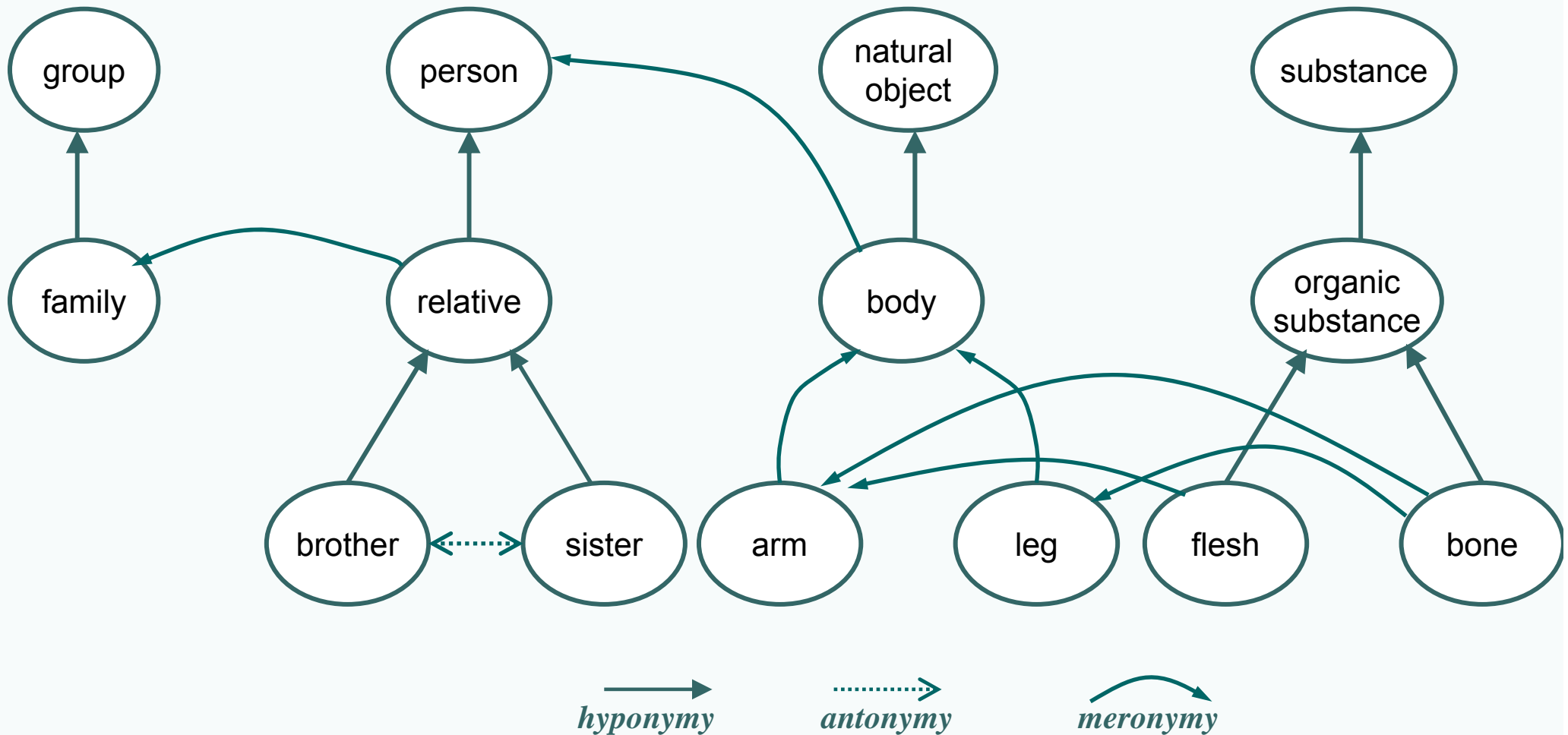
Mapped elements

Mapping generation

DSIM WSIM

Description Matching

# The name matching

- Computes name similarity coefficients *NSIM* between concept names

- Based on the notion of "semantic relatedness" (`rel`)
  - Degree of semantic similarity between two lexically expressed concepts
  - Based on the semantic relations of `WorNet`

# The name matching - WordNet



Legend:
- → hyponymy
- ⇢ antonymy
- ⌒→ meronymy

Nodes: group, person, natural object, substance, family, relative, body, organic substance, brother, sister, arm, leg, flesh, bone
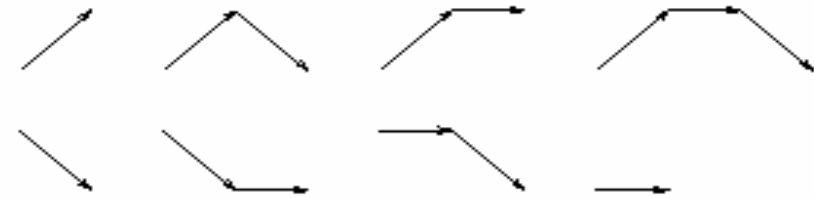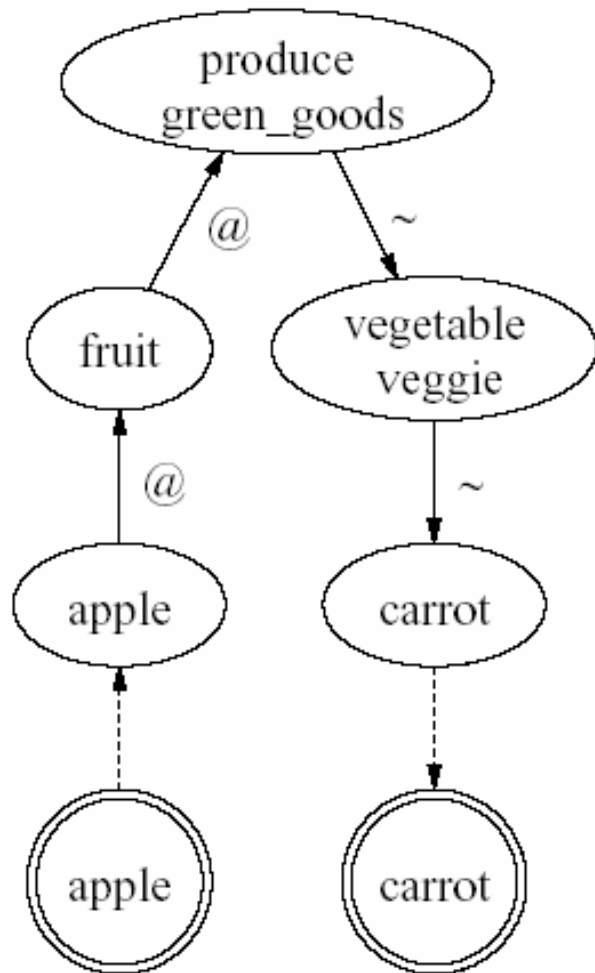
# The name matching

- Semantic relatedness
- 2 concepts are semantically close if:
  - Path not long
  - Path does not change direction too often

- **NSIM** ™ [0,1]
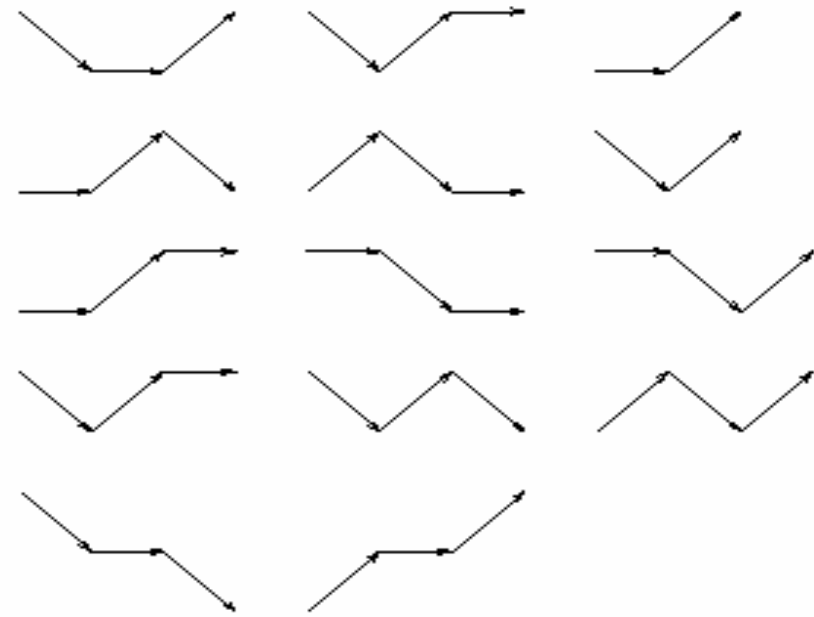
$$\text{rel}(C_1, C_2) = \lambda - \text{len}(C_1, C_2) - k * \text{turns}(C_1, C_2)$$

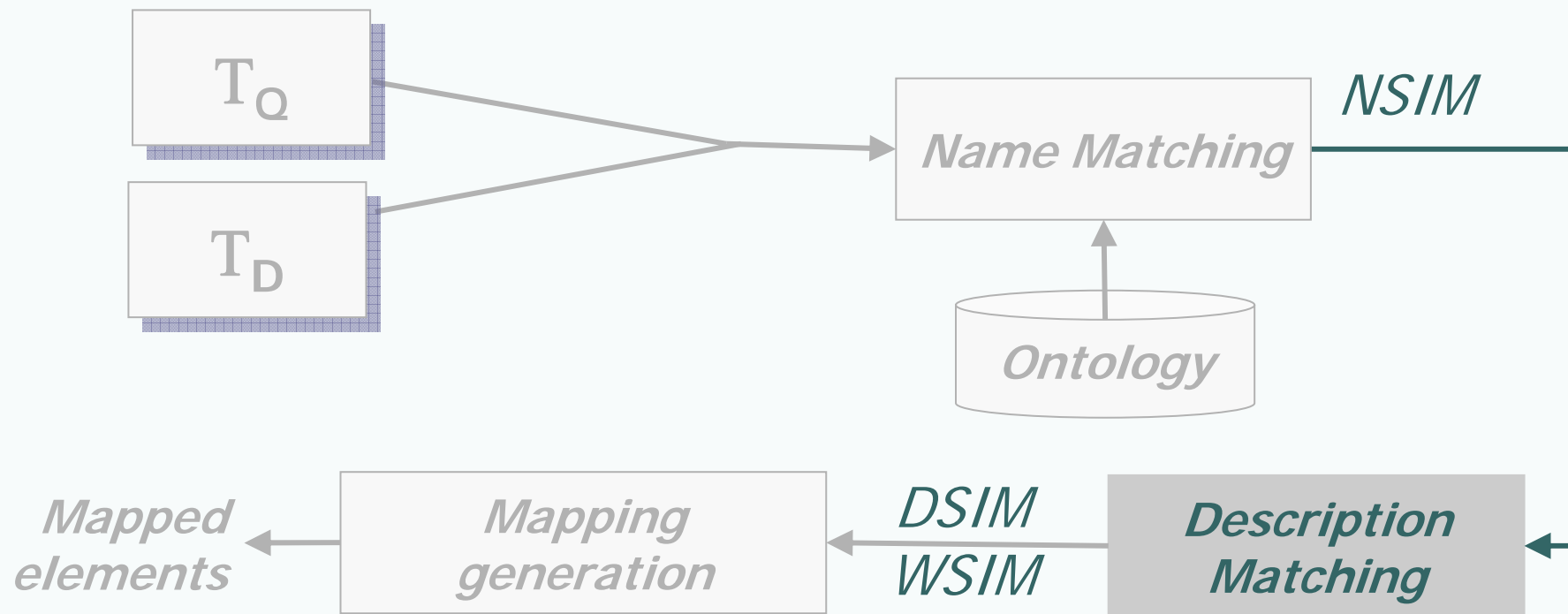$$\text{NSIM}(C_1, C_2) = \text{rel}(C_1, C_2) / \lambda$$

(a)

(b)

# The matching step

# The description matching

- Intuition



$$DSIM\,(C,D) = 1 - \frac{|\ C - D\ |}{|\ C\ |}$$

Difference operator

# The difference operator

- Allows to remove from a given description the information contained in another description

- Take into account linguistic relations (semantic relatedness) between concept and role names when computing the difference → "**Similarity difference**"

The difference algorithm based on the notion of **subsumption**

- **Goal :** define a subsumption taking into account linguistic relations between concept and role names

→ based on hierarchies

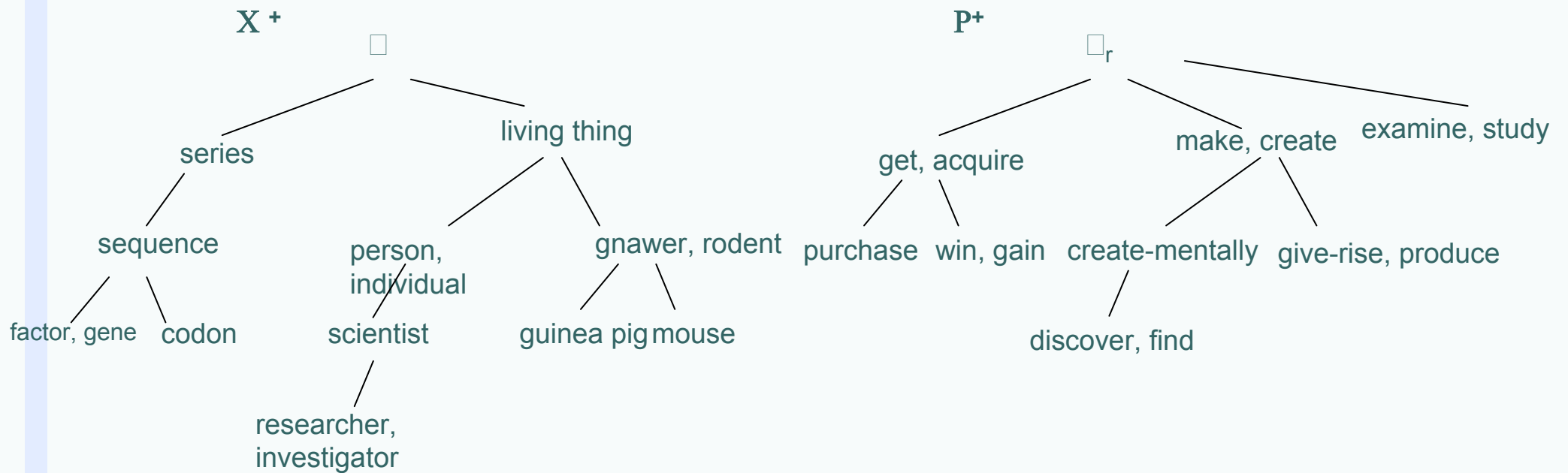→ based on similarities "**Similarity subsumption**"

# Difference based on concept and role hierarchies

# Hierarchies

- A support $\mathcal{S} = (\mathcal{C}^+, \mathcal{R}^+)$

  $\mathcal{C}^+ = (\mathcal{N}_A, \leq_C)$ where $\leq_C$ is a partial order relation defined on $\mathcal{N}_A$.

  $\mathcal{R}^+ = (\mathcal{N}_R, \leq_R)$ where $\leq_R$ is a partial order relation defined on $\mathcal{N}_R$.



Logic under consideration: $\mathcal{ALEH_S}$

# A structural subsumption algorithm for ALEH$_S$

- Based on graphs
- 3 steps
    - Concept descriptions are turned into a normal form
    - Normal forms represented by tree descriptions
    - Subsumption est caracterized in term of tree homomorphism

# Normalisation rules for ALEH$_S$

$$\forall r.C \sqcap \forall r.D \;\rightarrow\; \forall r.(C \sqcap D)$$

$$\forall s.C \sqcap \forall r.D \;\rightarrow\; \forall s.C \sqcap \forall r.(C \sqcap D), \text{si } r <_R s$$

$$\forall s.C \sqcap \exists r.D \;\rightarrow\; \forall s.C \sqcap \exists r.(C \sqcap D), \text{si } r \leq_R s$$

$$\forall r.\top \;\rightarrow\; \top$$

$$C \sqcap \top \;\rightarrow\; C$$

$$P \sqcap \neg Q \;\rightarrow\; \bot, \text{pour tout } P, Q \in N_C \text{ tel que } P \leq_C Q$$
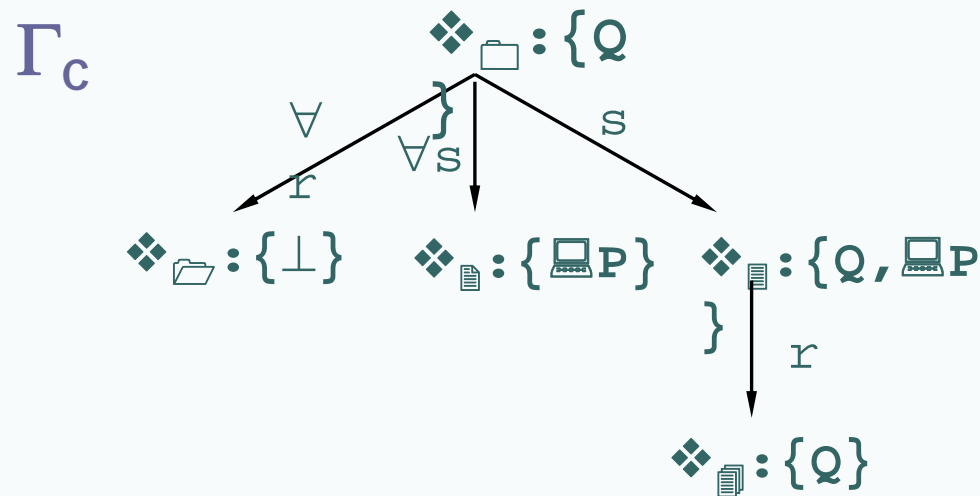
$$\exists r.\bot \;\rightarrow\; \bot$$

$$C \sqcap \bot \;\rightarrow\; \bot$$

# Decription trees

$$C \doteq Q \sqcap \forall r.P \sqcap \forall s.\neg P \sqcap \exists s.(Q \sqcap \exists r.Q)$$
$$C' \doteq Q \sqcap \forall r.\bot \sqcap \forall s.\neg P \sqcap \exists s.(Q \sqcap \neg P \sqcap \exists r.Q)$$

$$r \leq_R s$$

$$C \doteq Q \sqcap \forall r'. \bot \sqcap \forall s. \neg P \sqcap \exists s. (Q \sqcap \neg P \sqcap \exists r. Q)$$
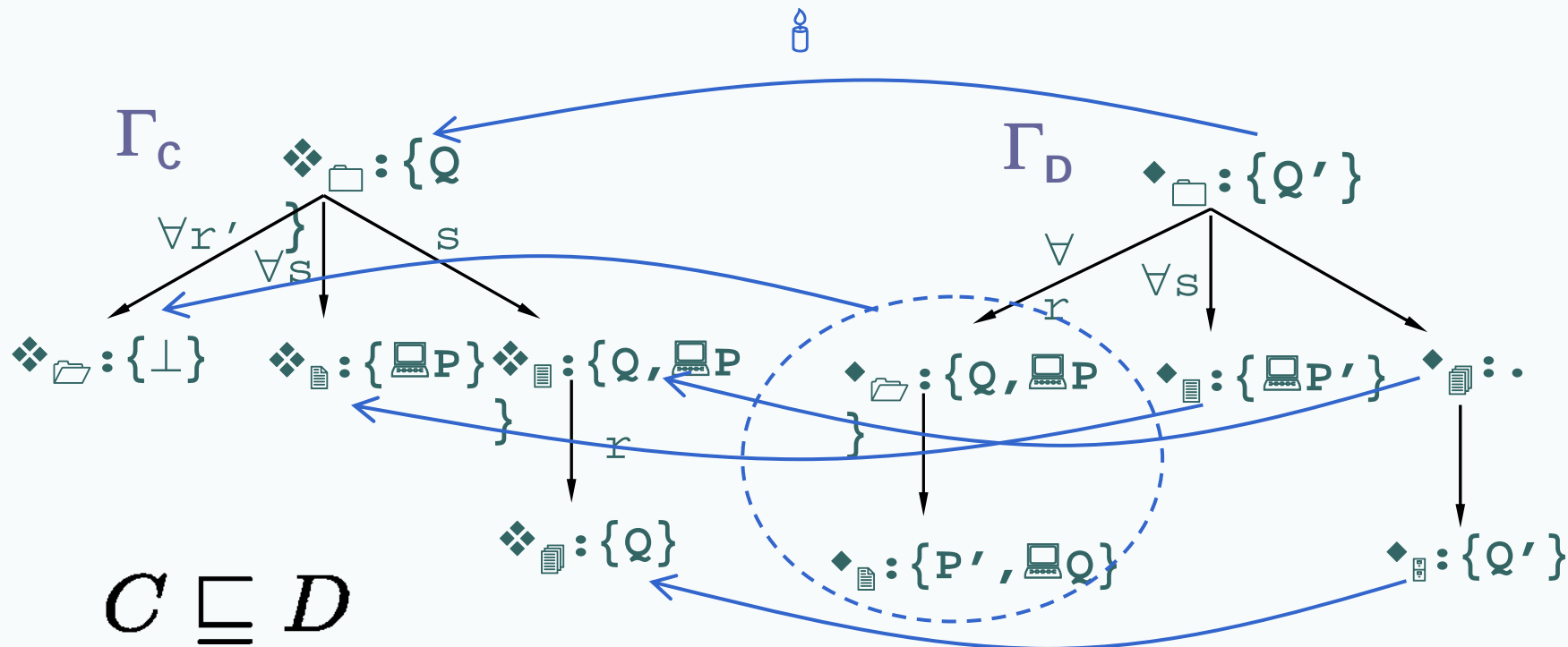
$$D \doteq Q' \sqcap \forall r. (Q \sqcap \exists s. (P' \sqcap \neg Q)) \sqcap \forall s. \neg P' \sqcap \exists s. (\exists s. Q')$$
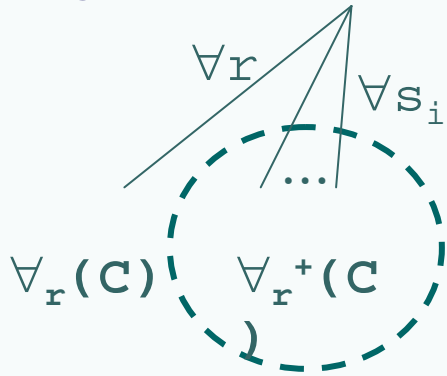
$$P \leq_C P' \quad r \leq_R r' \leq_R s$$



$$C \sqsubseteq D$$

# ALEH$_S$ Difference

$$P' \leq_C P$$
$$r \leq_R s_i$$

$\Gamma_C$: $\qquad \{\cancel{P_1}, P_2, \cancel{P_3}\}$

$\forall r$ $\qquad \forall s_i$

$\ldots$

$\forall_r(C)$ $\qquad \forall_r^+(C$

$\Gamma_D$: $\qquad \{P'_1, Q_1, Q_2, P_3\}$

$\forall s_i$

$\ldots$

$\forall_r^*(D$

$$diff_s(C,D) = P_2 \sqcap \forall r.\ diff_s(\forall_r(C), \forall_r^+(C) \sqcap \forall_r^*(D))$$

$$\sqcap \bigsqcap_{E^{™}E} \exists r.E$$

LIMOS

$$s_i \leq_R r$$

$\Gamma_C:$ $\{P_1, P_2, P_3\}$     $\Gamma_D:$ $\{P_1', Q_1, Q_2, P_3\}$



$$diff_s(C,D) = P_2 \sqcap \forall r. \; diff_s(\forall_r(C), \forall_r^+(C) \sqcap \forall_r^*(C))$$

$$\sqcap \prod_{E^{TM}E} \exists r.E$$

$$s_i \leq_R r$$

$\Gamma_C$: $\{P_1, P_2, P_3\}$

$\Gamma_D$: $\{P_1', Q_1, Q_2, P_3\}$

$s_i$ $r$

$s_i$ $r$

$C_{k+1}$ $C_n$ $C_1$ ... $C_k$

$D_1$ $D_m$

$\lceil \forall_r^*(C) \lceil$

$\forall_r^*(D)$

$\mathrm{E}$ $\mathcal{S}^{\sqsupseteq}$

$\exists_r(C)$

$\exists_r^-(D)$

$\exists_r^-(C)$

$$diff_s(C,D) = P_2 \lceil \forall r. \; diff_s(\forall_r(C), \forall_r^+(C) \lceil \forall_r^*(C))$$

$$\lceil \lceil_{E^{\mathrm{TM}}E} \exists r.E$$

# Difference based on similarities between concept and role names

# Similarity subsumption

- $\Sigma_1, \Sigma_2$ the sets of symbols of two terminologies $T_1$ and $T_2$

$$\alpha(c) = \{c' \in \Sigma_2 \mid nsim(c, c') > TH\}$$
$$\alpha(\text{human}) = \{\text{person, individual, someone}\}$$

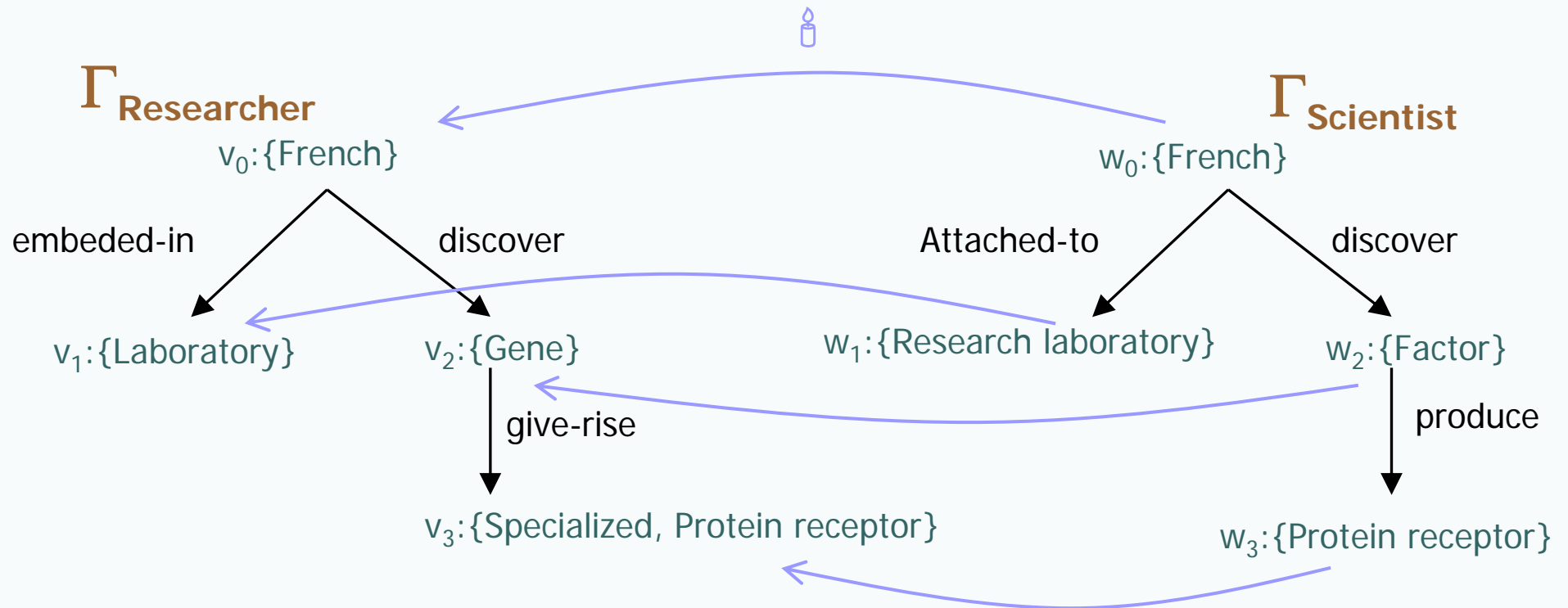- $\sigma$ Substitution: can replace a symbol c by an element of $\alpha(c)$

$\sigma(\exists\text{has-child.human}) = \exists \text{ has-offspring.person}$

$\sigma(\text{human}\sqcap\exists\text{has-child.human}) = \text{individual}\exists\text{has-offspring.person}$

$C \sqsubseteq_\alpha D$, iff there exists a substitution $\sigma$ w.r.t. $\alpha$ such that $C \sqsubseteq \sigma(D)$.

# Similarity subsumption

- **Homomorphism**



$\Gamma_{\textbf{Researcher}}$

$v_0$:{French}

embeded-in      discover

$v_1$:{Laboratory}    $v_2$:{Gene}

give-rise

$v_3$:{Specialized, Protein receptor}

$\Gamma_{\textbf{Scientist}}$

$w_0$:{French}

Attached-to      discover

$w_1$:{Research laboratory}    $w_2$:{Factor}

produce

$w_3$:{Protein receptor}

Resercher $\sqsubseteq_\alpha$ Scientist

$\sigma$ over $\alpha$ such that: Researcher $\sqsubseteq \sigma$ (Scientist)

# Similarity difference
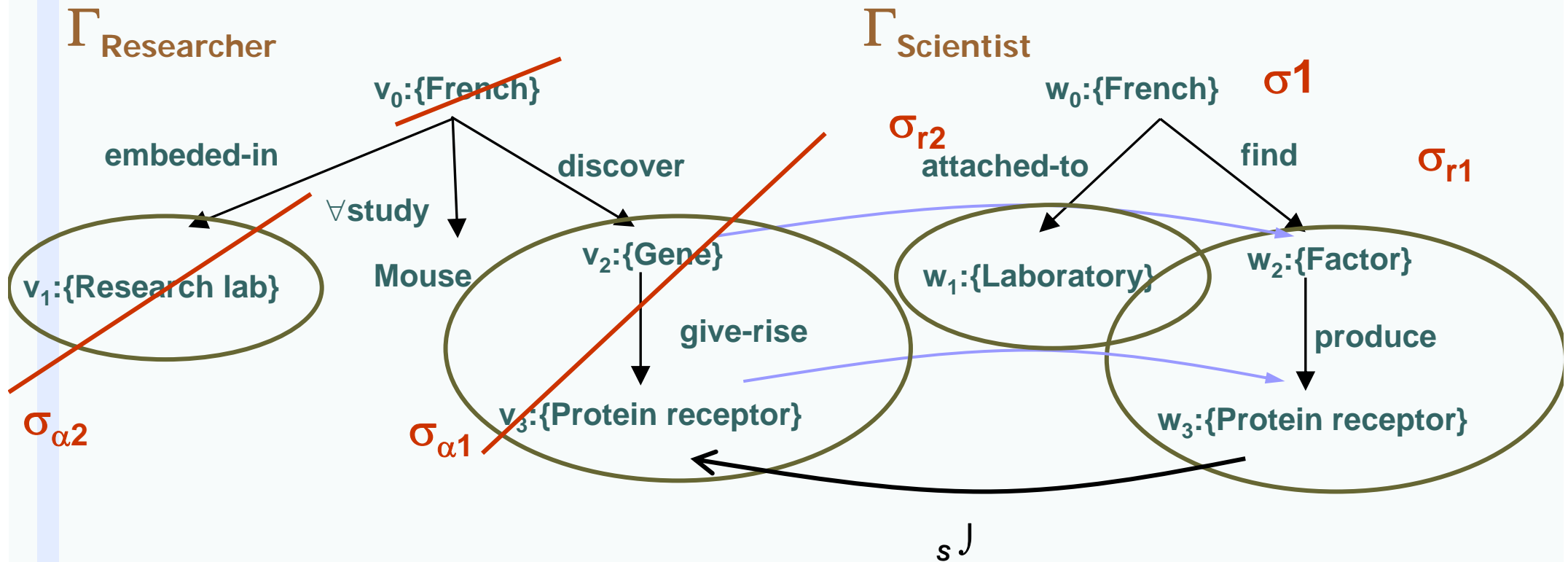
An expansion $\delta_{\mathcal{S}}$ w.r.t. $\mathcal{S}$

$$\delta_{\mathcal{S}}(C) = \prod_{\substack{\sigma_i \in \mathcal{S} \\ C_i \in C}} \sigma_i(C_i)$$

$E = C -_{\mathcal{S}} D$, iff there exists two sets of substi-tutions $\mathcal{S}_1$ and $\mathcal{S}_2$ w.r.t. $\alpha$ such that

$$\delta_{\mathcal{S}_1}(C) \sqcap \delta_{\mathcal{S}_2}(D) \equiv \delta_{\mathcal{S}_1}(E) \sqcap \delta_{\mathcal{S}_2}(D).$$

# Similarity difference



The sets of substitutions $\Sigma_1 = \{\sigma_{\alpha 1}, \sigma_{\alpha 2}\}$ $\Sigma_2 = \{\sigma_1, \sigma_{r1}, \sigma_{r2}\}$

s-diff (Resercher, Scientist) = $\forall$study. Mouse

## Similarity difference

$$\delta_{\mathcal{S}_1}(C) \sqcap \delta_{\mathcal{S}_2}(D) \equiv \delta_{\mathcal{S}_1}(E) \sqcap \delta_{\mathcal{S}_2}(D).$$

French $\sqcap$ $\exists$discover.($\sigma_{\alpha_1}$(Gene) $\sqcap$ $\exists\sigma_{\alpha_1}$(give-rise).$\sigma_{\alpha_1}$(Protein receptor)) $\sqcap$
$\exists$embedded-in.$\sigma_{\alpha_2}$(Reaserach lab)) $\sqcap$ $\forall$study.Mouse $\sqcap$ $\sigma_1$(French) $\sqcap$
$\exists\sigma_{\alpha_{r1}}$(find).(Factor $\sqcap$ $\exists$produce. Protein receptor)) $\sqcap$
$\exists\sigma_{\alpha_{r1}}$(attached-to).laboratory

$\equiv$

$\forall$study.Mouse $\sqcap$ $\sigma_1$(French) $\sqcap$ $\exists\sigma_{\alpha_{r1}}$(find).(Factor $\sqcap$
$\exists$produce. Protein receptor)) $\sqcap$ $\exists\sigma_{\alpha_{r1}}$(attached-to).laboratory

$$\textbf{\textit{DSIM}} (\text{Researcher, Scientist}) = 1 - \frac{|\text{ s-diff (Researcher, Scientist) }|}{|\text{ Researcher }|} = 0.77$$

# The description matching
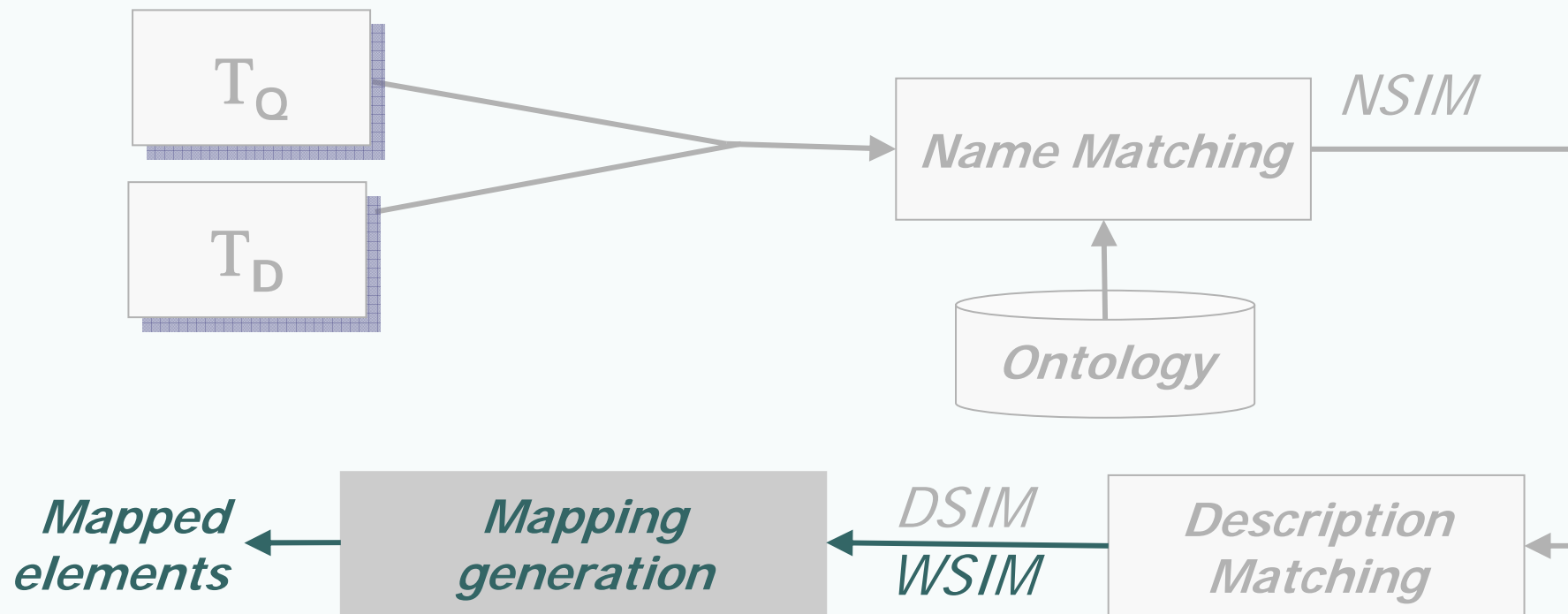
○ **WSIM** is mean of **NSIM** and **DSIM**

$$WSIM = w * NSIM + ( 1 - w ) * DSIM$$

**NSIM** (Researcher, Scientist) = 1
**DSIM** (Reasearcher, Scientist) = 0.77

=> **WSIM** (Reasearcher, Scientist) = 0.83

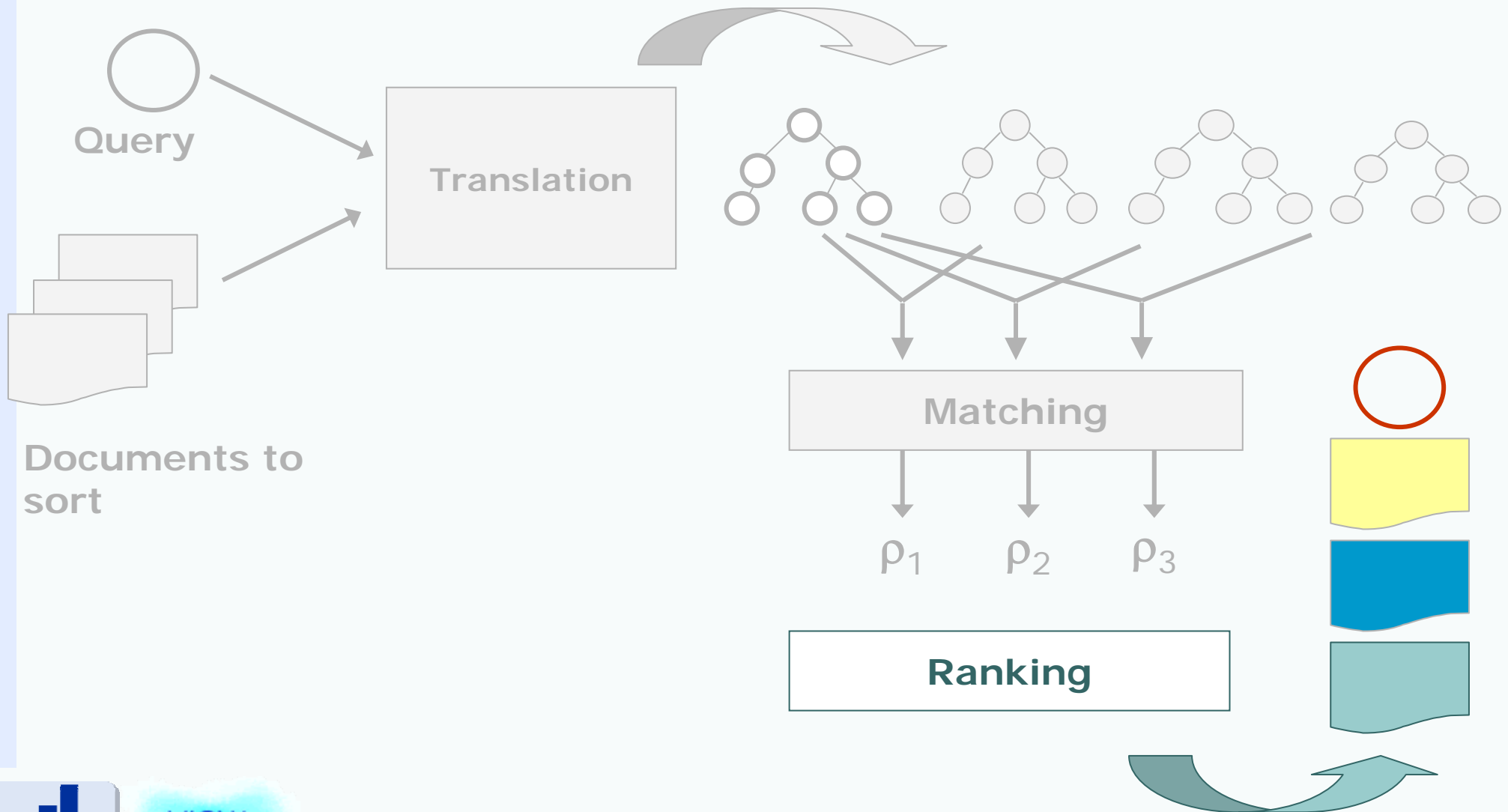($w$ = 0.3)

# The matching step

## The mapping generation

- A mapping is returned between elements having a weighted similarity greater than $th_{map}$

$$th_{map} = 0.75$$
$$WSIM \text{ (Reasearcher, Scientist)} = 0.83 > th_{map}$$

$$\rho \text{ (Reasearcher)} = \text{Scientist}$$

# The global process



Query

Documents to sort

Translation

Matching

$\rho_1$   $\rho_2$   $\rho_3$

Ranking

# The ranking step

- The ranking function

  - Based on the matching result
  - Computes the non covered part of the query by each document
  - Ranks the documents according to the size of this part
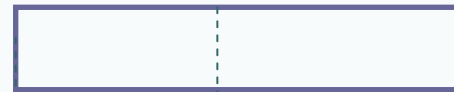
# The ranking step



$T_Q$

A_1 B ...
A_2 B ...

A_n B ...

$T_D$

B_1 B ...
B_2 B ...

B_m B ...

$A_1$

$\rho (A_1)$

s-diff$(A_1, \rho(A_1))$

$A_n$

$\rho (A_n)$

s-diff$(A_n, \rho(A_n))$

**Part of Q non covered by D**

$$\text{diff}_s (\circledast_Q, \circledast_D) = \int_{i=1..n} \text{s-diff} (A_i, \rho (A_i))$$

# Future work

- Approximate matching

- Application of matching to other type of data: web services

  - Representation / Adaptation to needs

- Extension of the method to

  - Structural subsumption algorithm $\mathcal{ALEN}$

# Implementation



**text**

**Input**

**Extraction module**

**Syntactic analysis**

$\mathcal{ALE}$-terminologies

**Input**

**Ontology**
(WordNet 2.0)

**Matching module**

**Subsumption test**

**difference**