

Query Processing for Non-traditional Applications

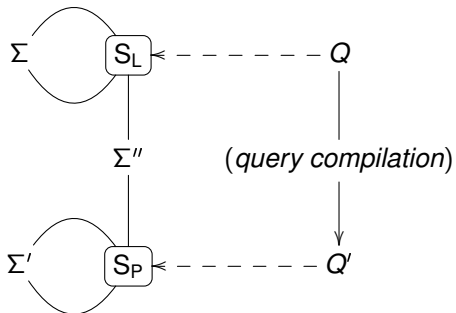
CS848 Spring 2013

Cheriton School of CS

Logical and Physical Schemas

Physical Design and Query Compilation: Overview

$$\Sigma = (\Sigma \cup \Sigma' \cup \Sigma'')$$



Standard Design: Discussion

Standard (relational) Physical Design

CREATE TABLE `employee` DDL command causes

- 1 a *logical symbol* `employee` to be created;
- 2 a (disk-based) *file* of (appropriate) records to be created; and
- 3 a link between these two objects to be *recorded* (where?)

- Multiple indices for the same table
- Horizontal partitioning
- Views and Materialized views

Points to consider

- 1 How are the above options *recorded* in a RDBMs? (speculation is ok)
- 2 Is there a uniform and compact way to describe *all* of the above options?

First-order Logic

- 1 First-order signatures for S_L and (most of) S_P ,
- 2 First-order sentences for Σ (Σ' , Σ''),
- 3 First order formulae for Q and (most of) Q' .

First-order Logic

- 1 First-order signatures for S_L and (most of) S_P ,
- 2 First-order sentences for Σ (Σ' , Σ''),
- 3 First order formulae for Q and (most of) Q' .

⇒ we'll need some additional *auxiliary* information in the cases of

S_P : which attributes are *input parameters*?

what are the symbol's *performance characteristics*?

Q' : how are formulae *mapped* to imperative programs?

are there plan features *not* captured by formula syntax?

LOGICAL DESIGN

Syntax of FOL: Signatures

Vocabularies are called *signatures* in FOL.

Syntax of FOL: Signatures

Vocabularies are called *signatures* in FOL.

The *non-logical parameters* in FOL consist of infinite disjoint collections $\{P_1, P_2, \dots\}$ and $\{f_1, f_2, \dots\}$ of *predicate symbols* and *function symbols*, respectively.

The *arity* of each symbol is a non-negative integer n , denoted $Ar(P_i)$ or $Ar(f_j)$.

- P/i denotes predicate symbol P where $Ar(P) = i$.
- f/j denotes function symbol f where $Ar(f) = j$.

Syntax of FOL: Signatures

Vocabularies are called *signatures* in FOL.

The *non-logical parameters* in FOL consist of infinite disjoint collections $\{P_1, P_2, \dots\}$ and $\{f_1, f_2, \dots\}$ of *predicate symbols* and *function symbols*, respectively.

The *arity* of each symbol is a non-negative integer n , denoted $\text{Ar}(P_i)$ or $\text{Ar}(f_j)$.

- P/i denotes predicate symbol P where $\text{Ar}(P) = i$.
- f/j denotes function symbol f where $\text{Ar}(f) = j$.

Predicate symbols of arity 0 are also called *propositions*,
Function symbols of arity 0 are also called *constants*.

Syntax of FOL: Signatures

Vocabularies are called *signatures* in FOL.

The *non-logical parameters* in FOL consist of infinite disjoint collections $\{P_1, P_2, \dots\}$ and $\{f_1, f_2, \dots\}$ of *predicate symbols* and *function symbols*, respectively.

The *arity* of each symbol is a non-negative integer n , denoted $Ar(P_i)$ or $Ar(f_i)$.

- P/i denotes predicate symbol P where $Ar(P) = i$.
- f/j denotes function symbol f where $Ar(f) = j$.

Predicate symbols of arity 0 are also called *propositions*,
Function symbols of arity 0 are also called *constants*.

A *signature* S in FOL is a (possibly infinite) selection of non-logical parameters.

- S^P denotes all predicate symbols in S .
- S^F denotes all function symbols in S .

ACME Case: Signatures for PAYROLL

OPTION 1

- $S_L^P = \{\text{employee}/3\}$
- $S_L^F = \emptyset$

ACME Case: Signatures for PAYROLL

OPTION 1

- $S_L^P = \{\text{employee}/3\}$
- $S_L^F = \emptyset$

Fewest non-logical parameters: a single 3-ary predicate symbol.

- 1st arg: an employee number
- 2nd arg: an employee name
- 3rd arg: an employee salary

ACME Case: Signatures for PAYROLL

OPTION 1

- $S_L^P = \{\text{employee}/3\}$
- $S_L^F = \emptyset$

Fewest non-logical parameters: a single 3-ary predicate symbol.

- 1st arg: an employee number
- 2nd arg: an employee name
- 3rd arg: an employee salary

1st arg serves a special role: the set of employee number values is identified with the set of employees.

ACME Case: Signatures for PAYROLL

OPTION 1

- $S_L^P = \{\text{employee}/3\}$
- $S_L^F = \emptyset$

Fewest non-logical parameters: a single 3-ary predicate symbol.

- 1st arg: an employee number
- 2nd arg: an employee name
- 3rd arg: an employee salary

1st arg serves a special role: the set of employee number values is identified with the set of employees.

Each 3-tuple in $(\text{employee})^I$ suggests two things.

- 1 The employee number is a *visible object identifier* of some employee.
- 2 The remaining two components of the 3-tuple express two facts about the employee.

Signatures for PAYROLL

OPTION 2

- $S_L^P = \{\text{employee}/1\}$
- $S_L^F = \{\text{employee-number}/1, \text{name}/1, \text{salary}/1\}$

Signatures for PAYROLL

OPTION 2

- $S_L^P = \{\text{employee}/1\}$
- $S_L^F = \{\text{employee-number}/1, \text{name}/1, \text{salary}/1\}$

Trades the need to remember the role of argument positions with the need to learn and remember additional non-logical parameters.

Signatures for PAYROLL

OPTION 2

- $S_L^P = \{\text{employee}/1\}$
- $S_L^F = \{\text{employee-number}/1, \text{name}/1, \text{salary}/1\}$

Trades the need to remember the role of argument positions with the need to learn and remember additional non-logical parameters.

Introduce

- unary predicates to capture the various kinds of entities, and
- unary functions to capture entity attributes.

Signatures for PAYROLL

OPTION 2

- $S_L^P = \{\text{employee}/1\}$
- $S_L^F = \{\text{employee-number}/1, \text{name}/1, \text{salary}/1\}$

Trades the need to remember the role of argument positions with the need to learn and remember additional non-logical parameters.

Introduce

- unary predicates to capture the various kinds of entities, and
- unary functions to capture entity attributes.

Advantages:

- Separates entity classification from entity description: an entity e in a given interpretation \mathcal{I} is an employee exactly when $e \in (\text{employee})^{\mathcal{I}}$.
- All information about entities, such as a `name` or `salary`, is captured by unary functions.

Signatures for PAYROLL

OPTION 1 versus OPTION 2

Latter allows the possibility that more than one employee can have the same *combination* of values for attributes `employee-number`, `name` and `salary`.

Signatures for PAYROLL

OPTION 1 versus OPTION 2

Latter allows the possibility that more than one employee can have the same *combination* of values for attributes `employee-number`, `name` and `salary`.

Replacing an n -ary predicate symbol with one unary predicate symbol and n unary function symbols is called *reification*.

Signatures for PAYROLL

OPTION 1 versus OPTION 2

Latter allows the possibility that more than one employee can have the same *combination* of values for attributes `employee-number`, `name` and `salary`.

Replacing an n -ary predicate symbol with one unary predicate symbol and n unary function symbols is called *reification*.

Disadvantages:

- Requires all entities to have a value for all attributes.

Signatures for PAYROLL

OPTION 1 versus OPTION 2

Latter allows the possibility that more than one employee can have the same *combination* of values for attributes `employee-number`, `name` and `salary`.

Replacing an n -ary predicate symbol with one unary predicate symbol and n unary function symbols is called *reification*.

Disadvantages:

- Requires all entities to have a value for all attributes.
- Therefore requires simulating partial functions (e.g., “null inapplicable” values).

Signatures for PAYROLL

OPTION 3

- $S_L^P = \{\text{employee}/1, \text{employee-number}/2, \text{name}/2, \text{salary}/2\}$
- $S_L^F = \emptyset$

Overcomes disadvantages of OPTION 2: replaces each unary function symbol with a new binary predicate symbol.

Signatures for PAYROLL

OPTION 3

- $S_L^P = \{\text{employee}/1, \text{employee-number}/2, \text{name}/2, \text{salary}/2\}$
- $S_L^F = \emptyset$

Overcomes disadvantages of OPTION 2: replaces each unary function symbol with a new binary predicate symbol.

Makes it possible for an entity (including employees) to have any number of employee numbers, names or salaries, including none.

Signatures for PAYROLL

OPTION 3

- $S_L^P = \{\text{employee}/1, \text{employee-number}/2, \text{name}/2, \text{salary}/2\}$
- $S_L^F = \emptyset$

Overcomes disadvantages of OPTION 2: replaces each unary function symbol with a new binary predicate symbol.

Makes it possible for an entity (including employees) to have any number of employee numbers, names or salaries, including none.

Replacing function symbols with new predicate symbols is always possible when a function free signature is desired.

Variables and Well-Formed Formulae

Denoted V , the *variables* in FOL are a countably infinite collection of symbols

$$\{x_1, x_2, \dots\}$$

disjoint from the set of non-logical parameters.

Variables and Well-Formed Formulae

Denoted V , the *variables* in FOL are a countably infinite collection of symbols

$$\{x_1, x_2, \dots\}$$

disjoint from the set of non-logical parameters.

The following grammars define the *terms*, *atoms* and *well formed formulae* induced by S , denoted $\text{TERM}(S)$, $\text{ATOM}(S)$ and $\text{WFF}(S)$, respectively.

- $\text{Term} ::= x$ (where $x \in V$) | $f(\text{Term}_1, \dots, \text{Term}_n)$ (where $f/n \in S^F$)
- $\text{Atom} ::= \text{Term}_1 \approx \text{Term}_2$ | $P(\text{Term}_1, \dots, \text{Term}_n)$ (where $P/n \in S^P$)
- $\phi, \psi ::= \text{Atom}$ | $\neg \phi$ | $(\phi \wedge \psi)$ | $\exists x. \phi$ (where $x \in V$)

Variables and Well-Formed Formulae

Denoted V , the *variables* in FOL are a countably infinite collection of symbols

$$\{x_1, x_2, \dots\}$$

disjoint from the set of non-logical parameters.

The following grammars define the *terms*, *atoms* and *well formed formulae* induced by S , denoted $\text{TERM}(S)$, $\text{ATOM}(S)$ and $\text{WFF}(S)$, respectively.

- $\text{Term} ::= x$ (where $x \in V$) | $f(\text{Term}_1, \dots, \text{Term}_n)$ (where $f/n \in S^F$)
- $\text{Atom} ::= \text{Term}_1 \approx \text{Term}_2$ | $P(\text{Term}_1, \dots, \text{Term}_n)$ (where $P/n \in S^P$)
- $\phi, \psi ::= \text{Atom}$ | $\neg \phi$ | $(\phi \wedge \psi)$ | $\exists x. \phi$ (where $x \in V$)

Assumes S denotes an FOL signature; we omit S when clear from context.

Variables and Well-Formed Formulae (cont'd)

The *logical parameters* in FOL:

- (*equality*) \approx
- (*negation*) \neg
- (*conjunction*) \wedge
- (*existential quantification*) \exists

Variables and Well-Formed Formulae (cont'd)

The *logical parameters* in FOL:

- (*equality*) \approx
- (*negation*) \neg
- (*conjunction*) \wedge
- (*existential quantification*) \exists

Convenient to have additional logical parameters as syntactic shorthand:

- (*disjunction*) \vee : “ $(\phi \vee \psi)$ ” \rightsquigarrow “ $\neg(\neg\phi \wedge \neg\psi)$ ”
- (*implication*) \rightarrow : “ $(\phi \rightarrow \psi)$ ” \rightsquigarrow “ $(\neg\phi \vee \psi)$ ”
- (*equivalence*) \equiv : “ $(\phi \equiv \psi)$ ” \rightsquigarrow “ $((\phi \rightarrow \psi) \wedge (\psi \rightarrow \phi))$ ”
- (*universal quantification*) \forall : “ $\forall x.\phi$ ” \rightsquigarrow “ $\neg\exists x.\neg\phi$ ”.

Variables and Well-Formed Formulae (cont'd)

The *logical parameters* in FOL:

- (*equality*) \approx
- (*negation*) \neg
- (*conjunction*) \wedge
- (*existential quantification*) \exists

Convenient to have additional logical parameters as syntactic shorthand:

- (*disjunction*) \vee : “ $(\phi \vee \psi)$ ” \rightsquigarrow “ $\neg(\neg\phi \wedge \neg\psi)$ ”
- (*implication*) \rightarrow : “ $(\phi \rightarrow \psi)$ ” \rightsquigarrow “ $(\neg\phi \vee \psi)$ ”
- (*equivalence*) \equiv : “ $(\phi \equiv \psi)$ ” \rightsquigarrow “ $((\phi \rightarrow \psi) \wedge (\psi \rightarrow \phi))$ ”
- (*universal quantification*) \forall : “ $\forall x.\phi$ ” \rightsquigarrow “ $\neg\exists x.\neg\phi$ ”.

Also common practice to omit parenthesis in well-formed formulae when intentions are clear, e.g.:

“($\phi_1 \wedge \phi_2 \wedge \phi_3$)” instead of “($\phi_1 \wedge (\phi_2 \wedge \phi_3)$)”

Free Variables

Given $t \in \text{TERM}$ or $\phi \in \text{WFF}$: $\text{Fv}(t)$ and $\text{Fv}(\phi)$ denote the *free variables* of a term and of a well formed formula, respectively.

$$\text{Fv}(t) = \begin{cases} \{x\} & \text{if } t = "x", \text{ and} \\ \bigcup_{1 \leq i \leq n} \text{Fv}(t_i) & \text{when } t = "f(t_1, \dots, t_n)" \text{ otherwise.} \end{cases}$$

$$\text{Fv}(\phi) = \begin{cases} \bigcup_{1 \leq i \leq n} \text{Fv}(t_i) & \text{if } \phi = "P(t_1, \dots, t_n)", \\ \text{Fv}(t_1) \cup \text{Fv}(t_2) & \text{if } \phi = "t_1 \approx t_2", \\ \text{Fv}(\psi) & \text{if } \phi = "\neg \psi", \\ \text{Fv}(\psi_1) \cup \text{Fv}(\psi_2) & \text{if } \phi = "(\psi_1 \wedge \psi_2)", \text{ and} \\ \text{Fv}(\psi) - \{x\} & \text{when } \phi = "\exists x.\psi" \text{ otherwise.} \end{cases}$$

A well-formed formula ϕ is *closed* if $\text{Fv}(\phi) = \emptyset$. A closed well-formed formula is also called a *sentence*.

Semantics of FOL: Interpretations

Assume S denotes a signature in FOL.

An *interpretation* $\mathcal{I}(S)$ of S is a pair $\langle \Delta^{\mathcal{I}(S)}, (\cdot)^{\mathcal{I}(S)} \rangle$.

- 1 $\Delta^{\mathcal{I}(S)}$ is a non-empty *domain* of entities.
- 2 $(\cdot)^{\mathcal{I}(S)}$ is an *interpretation function*.

Semantics of FOL: Interpretations

Assume S denotes a signature in FOL.

An *interpretation* $\mathcal{I}(S)$ of S is a pair $\langle \Delta^{\mathcal{I}(S)}, (\cdot)^{\mathcal{I}(S)} \rangle$.

- 1 $\Delta^{\mathcal{I}(S)}$ is a non-empty *domain* of entities.
- 2 $(\cdot)^{\mathcal{I}(S)}$ is an *interpretation function*.

For each $P/n \in S_P$, $(P/n)^{\mathcal{I}(S)}$ is a subset of $(\Delta^{\mathcal{I}(S)})^n$.

Semantics of FOL: Interpretations

Assume S denotes a signature in FOL.

An *interpretation* $\mathcal{I}(S)$ of S is a pair $\langle \Delta^{\mathcal{I}(S)}, (\cdot)^{\mathcal{I}(S)} \rangle$.

- 1 $\Delta^{\mathcal{I}(S)}$ is a non-empty *domain* of entities.
- 2 $(\cdot)^{\mathcal{I}(S)}$ is an *interpretation function*.

For each $P/n \in S_P$, $(P/n)^{\mathcal{I}(S)}$ is a subset of $(\Delta^{\mathcal{I}(S)})^n$.

For each $f/n \in S_F$, $(f/n)^{\mathcal{I}(S)}$ is a total function: $(\Delta^{\mathcal{I}(S)})^n \rightarrow \Delta^{\mathcal{I}(S)}$.

Semantics of FOL: Interpretations

Assume S denotes a signature in FOL.

An *interpretation* $\mathcal{I}(S)$ of S is a pair $\langle \Delta^{\mathcal{I}(S)}, (\cdot)^{\mathcal{I}(S)} \rangle$.

- 1 $\Delta^{\mathcal{I}(S)}$ is a non-empty *domain* of entities.
- 2 $(\cdot)^{\mathcal{I}(S)}$ is an *interpretation function*.

For each $P/n \in S_P$, $(P/n)^{\mathcal{I}(S)}$ is a subset of $(\Delta^{\mathcal{I}(S)})^n$.

For each $f/n \in S_F$, $(f/n)^{\mathcal{I}(S)}$ is a total function: $(\Delta^{\mathcal{I}(S)})^n \rightarrow \Delta^{\mathcal{I}(S)}$.

Write $\langle e_1, \dots, e_n \rangle$ to denote an *n-tuple*, an element of $(\Delta^{\mathcal{I}(S)})^n$.

Valuations

Assume \mathcal{I} is an interpretation of signature S .

Valuations

Assume \mathcal{I} is an interpretation of signature S .

A *valuation* over \mathcal{I} is written $\mathcal{V}(\mathcal{I})$ (or as \mathcal{V} when \mathcal{I} is clear from context) and is a total function: $V \rightarrow \Delta^{\mathcal{I}}$.

Valuations

Assume \mathcal{I} is an interpretation of signature S .

A *valuation* over \mathcal{I} is written $\mathcal{V}(\mathcal{I})$ (or as \mathcal{V} when \mathcal{I} is clear from context) and is a total function: $V \rightarrow \Delta^{\mathcal{I}}$.

For a given $x \in V$ and $e \in \Delta^{\mathcal{I}}$, the valuation $\mathcal{V}[x \mapsto e]$ is defined as follows:

$$\mathcal{V}[x_1 \mapsto e](x_2) = \begin{cases} e & \text{if “}x_1\text{” = “}x_2\text{”, and} \\ \mathcal{V}(x_2) & \text{otherwise.} \end{cases}$$

Valuations

Assume \mathcal{I} is an interpretation of signature S .

A *valuation* over \mathcal{I} is written $\mathcal{V}(\mathcal{I})$ (or as \mathcal{V} when \mathcal{I} is clear from context) and is a total function: $V \rightarrow \Delta^{\mathcal{I}}$.

For a given $x \in V$ and $e \in \Delta^{\mathcal{I}}$, the valuation $\mathcal{V}[x \mapsto e]$ is defined as follows:

$$\mathcal{V}[x_1 \mapsto e](x_2) = \begin{cases} e & \text{if “}x_1\text{” = “}x_2\text{”, and} \\ \mathcal{V}(x_2) & \text{otherwise.} \end{cases}$$

A valuation \mathcal{V} is extended to apply to any $t \in \text{TERM}$ in *the* way that satisfies

$$\mathcal{V}(t) = (f)^{\mathcal{I}}(\mathcal{V}(t_1), \dots, \mathcal{V}(t_n))$$

whenever $t = “f(t_1, \dots, t_n)”$.

Models

Assume S is a signature in FOL and also that $\phi \in \text{WFF}(S)$.

Models

Assume S is a signature in FOL and also that $\phi \in \text{WFF}(S)$.

An interpretation \mathcal{I} of S and valuation \mathcal{V} over \mathcal{I} is a *model* of ϕ , written

$$\mathcal{I}, \mathcal{V} \models \phi,$$

iff one of the following conditions apply:

- $\phi = "P(t_1, \dots, t_n)"$ and $\langle \mathcal{V}(t_1), \dots, \mathcal{V}(t_n) \rangle \in (P)^{\mathcal{I}}$,
- $\phi = "t_1 \approx t_2"$ and $\mathcal{V}(t_1) = \mathcal{V}(t_2)$,
- $\phi = "\neg\psi"$ and $\mathcal{I}, \mathcal{V} \not\models \psi$,
- $\phi = "(\psi_1 \wedge \psi_2)"$, $\mathcal{I}, \mathcal{V} \models \psi_1$ and $\mathcal{I}, \mathcal{V} \models \psi_2$, or
- $\phi = "\exists x.\psi"$ and $\mathcal{I}, \mathcal{V}[x \mapsto e] \models \psi$ for some $e \in \Delta^{\mathcal{I}}$.

Satisfiability and Logical Consequence

Assume Σ is a theory (over signature S).

Satisfiability and Logical Consequence

Assume Σ is a theory (over signature S).

We say the following.

- 1 The pair \mathcal{I}, \mathcal{V} is a *model* of Σ if $\mathcal{I}, \mathcal{V} \models \psi$ for all $\psi \in \Sigma$.

Satisfiability and Logical Consequence

Assume Σ is a theory (over signature S).

We say the following.

- 1 The pair \mathcal{I}, \mathcal{V} is a *model* of Σ if $\mathcal{I}, \mathcal{V} \models \psi$ for all $\psi \in \Sigma$.
- 2 Σ is *satisfiable* if it has a model and *unsatisfiable* otherwise.

Satisfiability and Logical Consequence

Assume Σ is a theory (over signature S).

We say the following.

- 1 The pair \mathcal{I}, \mathcal{V} is a *model* of Σ if $\mathcal{I}, \mathcal{V} \models \psi$ for all $\psi \in \Sigma$.
- 2 Σ is *satisfiable* if it has a model and *unsatisfiable* otherwise.
- 3 ϕ is a *logical consequence* of Σ , written

$$\Sigma \models \phi,$$

iff $\mathcal{I}, \mathcal{V} \models \phi$ for any model \mathcal{I}, \mathcal{V} of Σ .

Satisfiability and Logical Consequence

Assume Σ is a theory (over signature S).

We say the following.

- 1 The pair \mathcal{I}, \mathcal{V} is a *model* of Σ if $\mathcal{I}, \mathcal{V} \models \psi$ for all $\psi \in \Sigma$.
- 2 Σ is *satisfiable* if it has a model and *unsatisfiable* otherwise.
- 3 ϕ is a *logical consequence* of Σ , written

$$\Sigma \models \phi,$$

iff $\mathcal{I}, \mathcal{V} \models \phi$ for any model \mathcal{I}, \mathcal{V} of Σ .

The fundamental problem of reasoning in a given FOL theory $\Sigma(S)$ is the problem of *logical implication*: establishing which $\phi \in \text{WFF}(S)$ are logical consequences of $\Sigma(S)$.

ACME Case: Logical Constraints for PAYROLL

On identification.

Assume S_L is given by OPTION 1.

ACME Case: Logical Constraints for PAYROLL

On identification.

Assume S_L is given by OPTION 1.

The condition that *employees can be identified by their employee number* can be expressed as the FOL sentence

$$\forall x_1, x_2, y_1, y_2. (\exists z. (\text{employee}(z, x_1, x_2) \wedge \text{employee}(z, y_1, y_2)) \rightarrow ((x_1 \approx y_1) \wedge (x_2 \approx y_2))).$$

ACME Case: Logical Constraints for PAYROLL

On identification.

Assume S_L is given by OPTION 1.

The condition that *employees can be identified by their employee number* can be expressed as the FOL sentence

$$\forall x_1, x_2, y_1, y_2. (\exists z. (\text{employee}(z, x_1, x_2) \wedge \text{employee}(z, y_1, y_2)) \rightarrow ((x_1 \approx y_1) \wedge (x_2 \approx y_2))).$$

Ensures that each employee is associated with a single 3-tuple in $(\text{employee})^{\mathcal{I}}$ in any interpretation \mathcal{I} for ACME's PAYROLL system.

ACME Case: Logical Constraints for PAYROLL

On identification.

Assume S_L is given by OPTION 1.

The condition that *employees can be identified by their employee number* can be expressed as the FOL sentence

$$\forall x_1, x_2, y_1, y_2. (\exists z. (\text{employee}(z, x_1, x_2) \wedge \text{employee}(z, y_1, y_2)) \rightarrow ((x_1 \approx y_1) \wedge (x_2 \approx y_2))).$$

Ensures that each employee is associated with a single 3-tuple in $(\text{employee})^{\mathcal{I}}$ in any interpretation \mathcal{I} for ACME's PAYROLL system.¹

¹Remember introductory comments: the collection of all data corresponds to an interpretation \mathcal{I} .

ACME Case: Logical Constraints for PAYROLL

On identification.

Assume S_L is given by OPTION 1.

The condition that *employees can be identified by their employee number* can be expressed as the FOL sentence

$$\forall x_1, x_2, y_1, y_2. (\exists z. (\text{employee}(z, x_1, x_2) \wedge \text{employee}(z, y_1, y_2)) \rightarrow ((x_1 \approx y_1) \wedge (x_2 \approx y_2))).$$

Ensures that each employee is associated with a single 3-tuple in $(\text{employee})^{\mathcal{I}}$ in any interpretation \mathcal{I} for ACME's PAYROLL system.¹

Called a *functional dependency* in relational schema.

¹Remember introductory comments: the collection of all data corresponds to an interpretation \mathcal{I} .

Logical Constraints for PAYROLL (identification cont'd)

For S_L given by OPTION 2:

$$\forall x, y. ($$
$$\text{employee}(x) \wedge \text{employee}(y)$$
$$\wedge \text{employee-number}(x) \approx \text{employee-number}(y))$$
$$\rightarrow x \approx y).$$

Logical Constraints for PAYROLL (identification cont'd)

For S_L given by OPTION 2:

$$\forall x, y. (\\ \text{employee}(x) \wedge \text{employee}(y) \\ \wedge \text{employee-number}(x) \approx \text{employee-number}(y)) \\ \rightarrow x \approx y).$$

For S_L given by OPTION 3:

$$\forall x, y. (\\ \exists z. (\text{employee}(x) \wedge \text{employee}(y) \\ \wedge \text{employee-number}(x, z) \wedge \text{employee-number}(y, z)) \\ \rightarrow x \approx y).$$

Logical Constraints for PAYROLL (identification cont'd)

For S_L given by OPTION 2:

$$\forall x, y. (\\ \text{employee}(x) \wedge \text{employee}(y) \\ \wedge \text{employee-number}(x) \approx \text{employee-number}(y)) \\ \rightarrow x \approx y).$$

For S_L given by OPTION 3:

$$\forall x, y. (\\ \exists z. (\text{employee}(x) \wedge \text{employee}(y) \\ \wedge \text{employee-number}(x, z) \wedge \text{employee-number}(y, z)) \\ \rightarrow x \approx y).$$

More accurately, latter states that: *no pair of distinct employees may have any employee number at all in common* (becomes possible in OPTION 3 for employees to have any number of employee numbers).

Logical Constraints for PAYROLL

On property functionality.

Assume OPTION 3 chosen by ACME's APS department.

Logical Constraints for PAYROLL

On property functionality.

Assume OPTION 3 chosen by ACME's APS department.

Then necessary to disallow the number of possible values for an `employee-number`, `name`, or `salary` attribute for a given employee to exceed one.

Must add constraints to the logical constraints Σ to ensure the attributes are *partial functions*.

$$\forall x, y. (\exists z. (\text{employee-number}(z, x) \wedge \text{employee-number}(z, y)) \rightarrow (x \approx y))$$

$$\forall x, y. (\exists z. (\text{name}(z, x) \wedge \text{name}(z, y)) \rightarrow (x \approx y))$$

$$\forall x, y. (\exists z. (\text{salary}(z, x) \wedge \text{salary}(z, y)) \rightarrow (x \approx y))$$

Logical Constraints for PAYROLL

On typing.

The additional unary predicates in the PAYROLL signature can be used to ensure that attribute values are of appropriate types.

Logical Constraints for PAYROLL

On typing.

The additional unary predicates in the PAYROLL signature can be used to ensure that attribute values are of appropriate types.

For S_L given by OPTION 1:

$$\forall x, y, z. (\text{employee}(x, y, z) \rightarrow (\text{integer}(x) \wedge \text{string}(y) \wedge \text{integer}(z)))$$

Logical Constraints for PAYROLL

On typing.

The additional unary predicates in the PAYROLL signature can be used to ensure that attribute values are of appropriate types.

For S_L given by OPTION 1:

$$\forall x, y, z. (\text{employee}(x, y, z) \rightarrow (\text{integer}(x) \wedge \text{string}(y) \wedge \text{integer}(z)))$$

For S_L given by OPTION 2:

$$\forall x. (\text{employee}(x) \rightarrow (\text{integer}(\text{employee-number}(x)) \wedge \text{string}(\text{name}(x)) \wedge \text{integer}(\text{salary}(x))))$$

Logical Constraints for PAYROLL

For S_L given by OPTION 3:

$$\forall x.(\text{employee}(x) \rightarrow \exists y, z, w.(\text{employee-number}(x, y) \wedge \text{integer}(y) \\ \wedge \text{name}(x, z) \wedge \text{string}(z)) \\ \wedge \text{salary}(x, w) \wedge \text{integer}(w)))$$

Logical Constraints for PAYROLL

For S_L given by OPTION 3:

$$\forall x.(\text{employee}(x) \rightarrow \exists y, z, w.(\text{employee-number}(x, y) \wedge \text{integer}(y) \\ \wedge \text{name}(x, z) \wedge \text{string}(z)) \\ \wedge \text{salary}(x, w) \wedge \text{integer}(w)))$$

OPTION 3 also makes it possible to say that *only employees have employee numbers*.

$$\forall x.(\exists y.\text{employee-number}(x, y) \rightarrow \text{employee}(x))$$

PHYSICAL DESIGN

(take 1)

Access Paths and Simple Scanning

Assume ACME's DBA department selects a very simple physical design for PAYROLL: all employee information is recorded in a main-memory array.

```
array emp-array [1 to  $n$ ] of
  integer emp-num
  integer emp-salary
  string  emp-name
```


Access Paths and Simple Scanning

Assume ACME's DBA department selects a very simple physical design for PAYROLL: all employee information is recorded in a main-memory array.

```
array emp-array [1 to  $n$ ] of
  integer emp-num
  integer emp-salary
  string  emp-name
```

- 1 The salary, employee-number and name of each employee is recorded at some position in the array (in corresponding fields).
- 2 DBA ensures array entries are ordered by a *major sort* on emp-num values.

Access Paths and Simple Scanning

To capture a physical design in FOL, think in terms of *capabilities* attached to new predicate symbols (that become part of a physical vocabulary).

Access Paths and Simple Scanning

To capture a physical design in FOL, think in terms of *capabilities* attached to new predicate symbols (that become part of a physical vocabulary).

The organization of array `emp-array` suggests two capabilities in particular.

Access Paths and Simple Scanning

To capture a physical design in FOL, think in terms of *capabilities* attached to new predicate symbols (that become part of a physical vocabulary).

The organization of array `emp-array` suggests two capabilities in particular.

- 1 Scanning all entries: `emp-array0/3`.

Access Paths and Simple Scanning

To capture a physical design in FOL, think in terms of *capabilities* attached to new predicate symbols (that become part of a physical vocabulary).

The organization of array `emp-array` suggests two capabilities in particular.

- 1 Scanning all entries: `emp-array0/3`.
- 2 Scanning all entries with the first field matching a given `num` value: `emp-array1/3`.

Access Paths and Simple Scanning

To capture a physical design in FOL, think in terms of *capabilities* attached to new predicate symbols (that become part of a physical vocabulary).

The organization of array `emp-array` suggests two capabilities in particular.¹

- 1 Scanning all entries: `emp-array0/3`.
- 2 Scanning all entries with the first field matching a given `num` value: `emp-array1/3`.

¹The DBA department must provide the code that *implements* these capabilities in a library or at runtime, e.g., code that performs a binary search of `emp-array` in the case of `emp-array0` and `emp-array1`.

Access Paths and Simple Scanning

The *access paths* of S are a distinguished subset S_A of the predicate symbols S_P that comprise the physical vocabulary of S .

Access Paths and Simple Scanning

The *access paths* of S are a distinguished subset S_A of the predicate symbols S_P that comprise the physical vocabulary of S .

The *binding pattern* of an access path $P \in S_A$ is denoted $Bp(P)$ and is a non-negative integer satisfying $0 \leq Bp(P) \leq Ar(P)$.

Access Paths and Simple Scanning

The *access paths* of S are a distinguished subset S_A of the predicate symbols S_P that comprise the physical vocabulary of S .

The *binding pattern* of an access path $P \in S_A$ is denoted $Bp(P)$ and is a non-negative integer satisfying $0 \leq Bp(P) \leq Ar(P)$.

Write $P/n/m$ to indicate that P is a predicate symbol with arity n and that P is also an access path with binding pattern m (and write $P/n/m \in S_A$).

Access Paths and Simple Scanning

The *access paths* of S are a distinguished subset S_A of the predicate symbols S_P that comprise the physical vocabulary of S .

The *binding pattern* of an access path $P \in S_A$ is denoted $Bp(P)$ and is a non-negative integer satisfying $0 \leq Bp(P) \leq Ar(P)$.

Write $P/n/m$ to indicate that P is a predicate symbol with arity n and that P is also an access path with binding pattern m (and write $P/n/m \in S_A$).

The declaration of `emp-array` is captured by adding two new predicate symbols that are also access paths to S_A .

$$\{\text{emp-array}_0/3/0, \text{emp-array}_1/3/1\}$$

Access Paths and Simple Scanning

The *access paths* of S are a distinguished subset S_A of the predicate symbols S_P that comprise the physical vocabulary of S .

The *binding pattern* of an access path $P \in S_A$ is denoted $Bp(P)$ and is a non-negative integer satisfying $0 \leq Bp(P) \leq Ar(P)$.

Write $P/n/m$ to indicate that P is a predicate symbol with arity n and that P is also an access path with binding pattern m (and write $P/n/m \in S_A$).

The declaration of `emp-array` is captured by adding two new predicate symbols that are also access paths to S_A .

$$\{\text{emp-array}_0/3/0, \text{emp-array}_1/3/1\}$$

This new set of predicate symbols is now the *physical vocabulary* S_P of ACME's PAYROLL system.

ACME Case: Mapping Constraints

How do we know (ensure) that `employee` is properly represented by `emp-array` (and in turn by the access paths `emp-array0` and `emp-array1`)?

ACME Case: Mapping Constraints

How do we know (ensure) that `employee` is properly represented by `emp-array` (and in turn by the access paths `emp-array0` and `emp-array1`)?

ACME's DBA group must specify *mapping* or *correspondence* constraints Σ'' over the signature $(S_L \cup S_P)$.

ACME Case: Mapping Constraints

How do we know (ensure) that `employee` is properly represented by `emp-array` (and in turn by the access paths `emp-array0` and `emp-array1`)?

ACME's DBA group must specify *mapping* or *correspondence* constraints Σ'' over the signature $(S_L \cup S_P)$.

Such constraints provide the necessary “connections” for all possible interpretations \mathcal{I} (encoding factual data) of the *logical vocabulary* S_L and the *physical vocabulary* S_P of PAYROLL.

ACME Case: Mapping Constraints

How do we know (ensure) that `employee` is properly represented by `emp-array` (and in turn by the access paths `emp-array0` and `emp-array1`)?

ACME's DBA group must specify *mapping* or *correspondence* constraints Σ'' over the signature $(S_L \cup S_P)$.

Such constraints provide the necessary “connections” for all possible interpretations \mathcal{I} (encoding factual data) of the *logical vocabulary* S_L and the *physical vocabulary* S_P of PAYROLL.

With OPTION 1 for S_L , DBA can add the following sentences to Σ .

- 1 $\forall x, y, z. (\text{employee}(x, y, z) \rightarrow \text{emp-array0}(x, z, y))$
- 2 $\forall x, y, z. (\text{emp-array0}(x, z, y) \rightarrow \text{employee}(x, y, z))$

ACME Case: Access Path so far

- 1 `emp-array0/3/0` allows us to *scan* all employees;
- 2 `emp-array1/3/1` allows us to *find* (all) employees given `enumber`.

ACME Case: Access Path so far

- 1 `emp-array0/3/0` allows us to *scan* all employees;
- 2 `emp-array1/3/1` allows us to *find* (all) employees given `enumber`.

What if we also want to find employees *by their name*?

- 1 use `emp-array0` and *filter out* non-matching employees (“selection”)

ACME Case: Access Path so far

- 1 `emp-array0/3/0` allows us to *scan* all employees;
- 2 `emp-array1/3/1` allows us to *find* (all) employees given `enumber`.

What if we also want to find employees *by their name*?

- 1 use `emp-array0` and *filter out* non-matching employees (“selection”)
- 2 *improve the physical design* to allow efficient search ⇒ “create index”

ACME Case: Access Path Code Templates

Note: Code templates for access paths must be provided by ACME's DBA department.

ACME Case: Access Path Code Templates

Note: Code templates for access paths must be provided by ACME's DBA department.

E.g., Pseudo-code templates realizing a *first/next* protocol for `emp-array0` might be given as follows (variables would be renamed for each occurrence of `emp-array0` in a query plan).

```
function emp-array0-first
   $i := 0$ 
  return emp-array0-next
```

```
function emp-array0-next
   $i := i + 1$ 
  if ( $i > n$ ) return false
   $x_1 := \text{emp-array}[i].\text{emp-salary}$ 
   $x_2 := \text{emp-array}[i].\text{emp-num}$ 
   $x_3 := \text{emp-array}[i].\text{emp-name}$ 
  return true
```

ACME Case: Access Path Code Templates

Note: Code templates for access paths must be provided by ACME's DBA department.

E.g., Pseudo-code templates realizing a *first/next* protocol for `emp-array0` might be given as follows (variables would be renamed for each occurrence of `emp-array0` in a query plan).

```
function emp-array0-first
   $i := 0$ 
  return emp-array0-next
```

```
function emp-array0-next
   $i := i + 1$ 
  if ( $i > n$ ) return false
   $x_1 := emp-array[i].emp-salary$ 
   $x_2 := emp-array[i].emp-num$ 
   $x_3 := emp-array[i].emp-name$ 
  return true
```

Assumes a global state recording bindings of (possible copies of) variables.

- 1 x_1 , x_2 and x_3 to communicate the contents of `emp-array`.
- 2 i and n to record scanning status and size of `emp-array`.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*. Alternatives to an array could also have served the same purpose: linked lists, simple search trees, and so on.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*. Alternatives to an array could also have served the same purpose: linked lists, simple search trees, and so on.

It is beyond the scope of this book to consider the synthesis of such basic data structures as part of the job of query compilation.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*. Alternatives to an array could also have served the same purpose: linked lists, simple search trees, and so on.

It is beyond the scope of this book to consider the synthesis of such basic data structures as part of the job of query compilation.

However, we will see how physical design based on more complex data structures can be usefully *decomposed* into such basic data structures using FOL.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*. Alternatives to an array could also have served the same purpose: linked lists, simple search trees, and so on.

It is beyond the scope of this book to consider the synthesis of such basic data structures as part of the job of query compilation.

However, we will see how physical design based on more complex data structures can be usefully *decomposed* into such basic data structures using FOL.¹

¹A decomposition of more complex data structures can enable compilation opportunities that would otherwise not be possible.

Access Path Code Templates

Examples of atomic query plans have so far been based on using an array as a *basic collection type*. Alternatives to an array could also have served the same purpose: linked lists, simple search trees, and so on.

It is beyond the scope of this book to consider the synthesis of such basic data structures as part of the job of query compilation.

However, we will see how physical design based on more complex data structures can be usefully *decomposed* into such basic data structures using FOL.¹

Main point: Once given `first/next` “black box” code templates for the basic data structures (such as records, arrays, linked lists and simple search trees) constraints can then be expressed in FOL that do the rest.

¹A decomposition of more complex data structures can enable compilation opportunities that would otherwise not be possible.

Summary: Data vs. Metadata in FOL

Metadata (database schema)

- 1 Signature S_L and constraints Σ for the *logical schema*,
- 2 Signature S_P and constraints Σ' for the *physical schema*,
- 3 Constraints Σ'' that relate S_L to S_P .

Data (database instance)

A *first-order structure (interpretation)* that

- 1 interprets symbols in S_L and S_P and
- 2 satisfies $\Sigma \cup \Sigma' \cup \Sigma''$.