

Spivak - Database Migration Theory

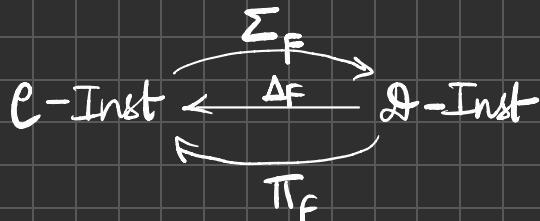
- * \mathcal{C} -instance is a functor $I: \mathcal{C} \rightarrow \text{Set}$ where \mathcal{C} is a schema
- * An **instance homomorphism** between 2 instance $I, J: \mathcal{C} \rightarrow \text{Set}$ is a natural transformation $I \xrightarrow{\sim} J$ (also 'schema mapping')
- * If \mathcal{C}, \mathcal{D} are 2 categories, $F: \mathcal{C} \rightarrow \mathcal{D}$, is a functor for any set-value functor $I: \mathcal{D} \rightarrow \text{Set}$ we refer its composite $F \circ I: \mathcal{C} \rightarrow \text{Set}$ as pullback of I along F .

$$\mathcal{C} \xrightarrow{F} \mathcal{D} \xrightarrow{I} \text{Set} \quad \mathcal{C} \xrightarrow{F \circ I} \text{Set}$$

Given $F: \mathcal{C} \rightarrow \mathcal{D}$, a data migration functor Δ_F turns \mathcal{D} -Inst to \mathcal{C} -Inst

Note:

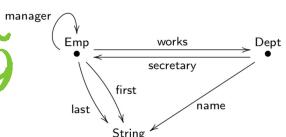
* A 'database schema' \mathcal{C} consist of a pair $\mathcal{C} = (G, \cong)$ where G is a graph & \cong is a congruence on G .



Diversion →

Schemas and Instances

Schema (graph)



$[manager.works] = [works]$ $[secretary.works] = []$

Inst.

Emp				
ID	mgr	works	first	last
101	103	q10	Al	Akin
102	102	x02	Bob	Bo
103	103	q10	Carl	Cork

Dept		
ID	sec	name
q10	101	CS
x02	102	Math
...		

String	
ID	
1	name
2	age
3	Int

CQL code for the schema:

entities

Emp
Dept

foreign keys

manager : Emp → Emp
works : Emp → Dept
secretary : Dept → Emp

attributes

first last : Emp → string
name : Dept → string

path equations

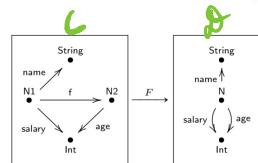
manager.works = works
secretary.works = Department

$\Delta \rightarrow$ Duplicate / Destroy tables / columns
 $\Sigma \rightarrow$ Union Data
 $\Pi \rightarrow$ Pair & Query Data (Join)

Example: (from a conf. video)

A schema mapping $F : S \rightarrow T$ is an equation-preserving function:

$nodes(S) \rightarrow nodes(T)$



$edges(S) \rightarrow paths(T)$

Not edges necessarily

$C \xrightarrow{F} \emptyset$

$F(Int) = Int \quad F(String) = String$
 $F(N1) = N \quad F(N2) = N$
 $F(name) = [name] \quad F(age) = [age] \quad F(salary) = [salary]$
 $F(f) = []$

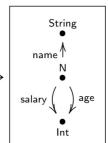
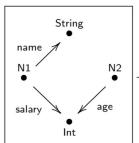
N1			N2		N			
ID	name	salary	ID	age	ID	name	salary	age
1	Alice	\$100	4	20	a	Alice	\$100	20
2	Bob	\$250	5	20	b	Bob	\$250	20
3	Sue	\$300	6	30	c	Sue	\$300	30

$\leftarrow \Delta_F$

$C\text{-inst}$ $\emptyset\text{-inst}$

What Δ_F did here?

It duplicates / divides N to N_1 & N_2 for whatever column they have (here very simple) → It 'projects'



C-inst

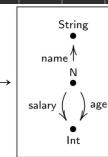
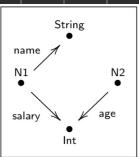
N1			N2	
ID	Name	Salary	ID	Age
1	Alice	\$100	4	20
2	Bob	\$250	5	20
3	Sue	\$300	6	30

D-inst

N			
ID	Name	Salary	Age
a	Alice	\$100	null ₁
b	Bob	\$250	null ₂
c	Sue	\$300	null ₃
d	null ₄	null ₅	20
e	null ₆	null ₇	20
f	null ₈	null ₉	30

→ unions the data
(Note:

In every row of N either data of N₁ or N₂ is only stored no mixed happened!
Like "outer - union"



N1

N2

Σ_F

N			
ID	name	salary	age
a	Alice	\$100	20
b	Alice	\$100	20
c	Alice	\$100	30
d	Bob	\$250	20
e	Bob	\$250	20
f	Bob	\$250	30
g	Sue	\$300	20
h	Sue	\$300	20
i	Sue	\$300	30

→ $\Pi_F \rightarrow$ products the data

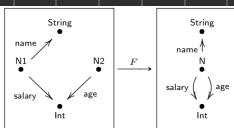
⇒ Queries in CORL : Target Schema

$Q : S \rightarrow T$
 source ↴ schema (F : S → X)
 ↴ schema (G : T → X) Intermediate 'X'
 (query's result)

Evaluation : $\text{eval}_Q \stackrel{\approx}{=} \Delta_Q \circ \Pi_F$

Query execution is done via this.

{ coeval_Q is how the query results to source schema }



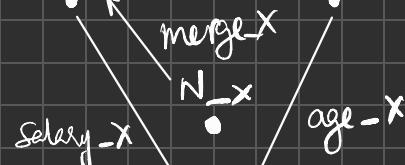
N1			N2		
ID	name	salary	ID	age	
1	Alice	\$100	4	20	
2	Bob	\$250	5	20	
3	Sue	\$300	6	30	

Π_F

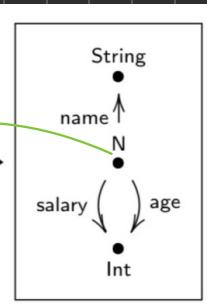
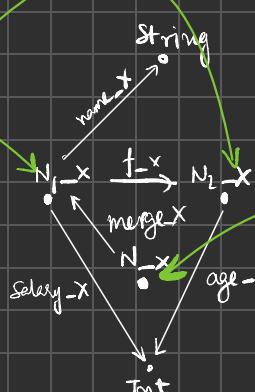
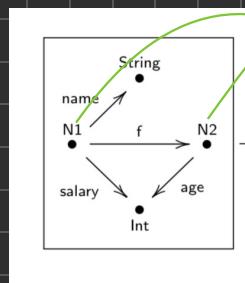
N			
ID	name	salary	age
a	Alice	\$100	20
b	Alice	\$100	20
c	Alice	\$100	30
d	Bob	\$250	20
e	Bob	\$250	20
f	Bob	\$250	30
g	Sue	\$300	20
h	Sue	\$300	20
i	Sue	\$300	30

for this, string

$X :$



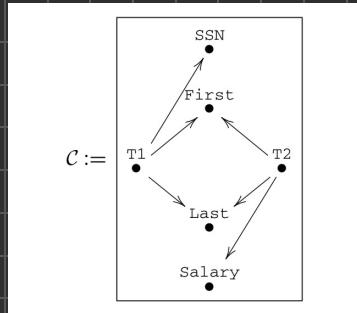
age-X



$$\zeta \xrightarrow{F} X \xleftarrow{T} \theta$$

$$\text{eval}_Q = \Delta_\theta \circ \Pi_F$$

other e.g.s. from FDM paper
C-Inst



T1			
ID	SSN	First	Last
T1-001	115-234	Bob	Smith
T1-002	122-988	Sue	Smith
T1-003	198-877	Alice	Jones

T2			
ID	First	Last	Salary
T2-A101	Alice	Jones	\$100
T2-A102	Sam	Miller	\$150
T2-A104	Sue	Smith	\$300
T2-A110	Carl	Pratt	\$200

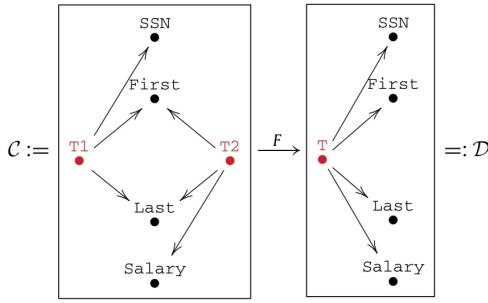
SSN	
ID	SSN
	115-234
	118-334
	122-988
	198-877
	342-164

First	
ID	First
	Adam
	Alice
	Bob
	Carl
	Sam
	Sue

Last	
ID	Last
	Jones
	Miller
	Pratt
	Richards
	Smith

Salary	
ID	Salary
	\$100
	\$150
	\$200
	\$250
	\$300

$C \xrightarrow{F} \mathcal{D}$



Example 2.21 (Pullback). Let $F : C \rightarrow \mathcal{D}$ be the translation given in diagram (7). In this example, we explore the data migration functor $\Delta_F : \mathcal{D}\text{-Inst} \rightarrow C\text{-Inst}^4$ by applying it to a \mathcal{D} -instance J . Note that even though our translation F points forwards (from C to \mathcal{D}), our migration functor Δ_F points "backwards" (from \mathcal{D} -instances to C -instances). We will see why it works that way, but first we bring the discussion down to earth by working with a particular \mathcal{D} -instance.

Consider the instance J , on schema \mathcal{D} , defined by the table

T				
ID	SSN	First	Last	Salary
XF6671	115-234	Bob	Smith	\$250
XF8911	122-988	Sue	Smith	\$300
XF2211	198-877	Alice	Jones	\$100

$\mathcal{D}\text{-Inst}$

Δ_F

and having the four leaf tables from Example 2.11, display (5). Pulling back along the translation F , we are supposed to get an instance $\Delta_F(J)$ on schema C , which we must describe. But the description is easy: $\Delta_F(J)$ splits up the columns of table T according to the translation F . The four leaf tables will be exactly the same as above (i.e. the same as in Example 2.11, (5)), and the two fact tables will be something like?

T1			
ID	SSN	First	Last
T1-001	115-234	Bob	Smith
T1-002	122-988	Sue	Smith
T1-003	198-877	Alice	Jones

T2			
ID	First	Last	Salary
A21	Alice	Jones	\$100
A67	Bob	Smith	\$250
A91	Sue	Smith	\$300

After Σ_F (left push f.f.)

After (Rightpush forward functor)
 Π_F

T				
ID	SSN	First	Last	Salary
T1-002	122-988	Sue	Smith	\$300
T1-003	198-877	Alice	Jones	\$100

T				
ID	SSN	First	Last	Salary
T1-001	115-234	Bob	Smith	T1-001.Salary
T1-002	122-988	Sue	Smith	T1-002.Salary
T1-003	198-877	Alice	Jones	T1-003.Salary
T2-A101	T2-A101.SSN	Alice	Jones	\$100
T2-A102	T2-A102.SSN	Sam	Miller	\$150
T2-A104	T2-A104.SSN	Sue	Smith	\$300
T2-A110	T2-A110.SSN	Carl	Pratt	\$200