



# **An Approach to Information Integration Based on the AMN Formalism**

M. G. Manukyan

*Yerevan State University, Armenia*

First Workshop on Programming the Semantic Web

November 11, 2012, Boston, USA



# Goals

- To use two or more databases (data sources) in the frame of a large database, possibly virtual, containing information from all sources, so the data can be queried as a unit.

# Agenda

- Problems of Information Integration
- State of the Art
- Canonical Model and Principle of its Extension
- AMN Formalization of Data Model
- Mapping from RDM into Canonical Model
- Summary & Future Work

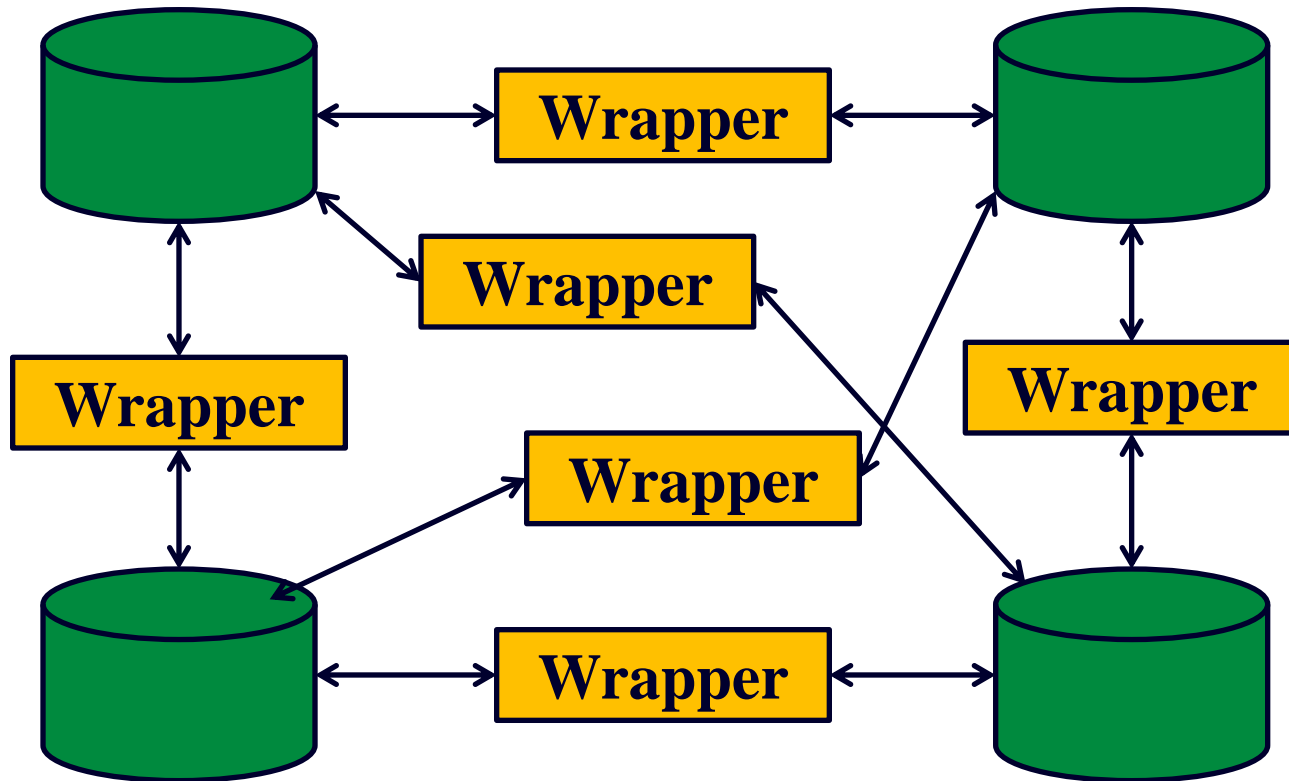
# Problems of Information Integration

- **Data Type Differences** The same property is presented by different types;
- **Value Differences** The value of the same property have different encodings;
- **Semantic Differences** Terms may be given different interpretations;
- **Missing Values** A source might not record information of a type that all or most of the other sources provide.

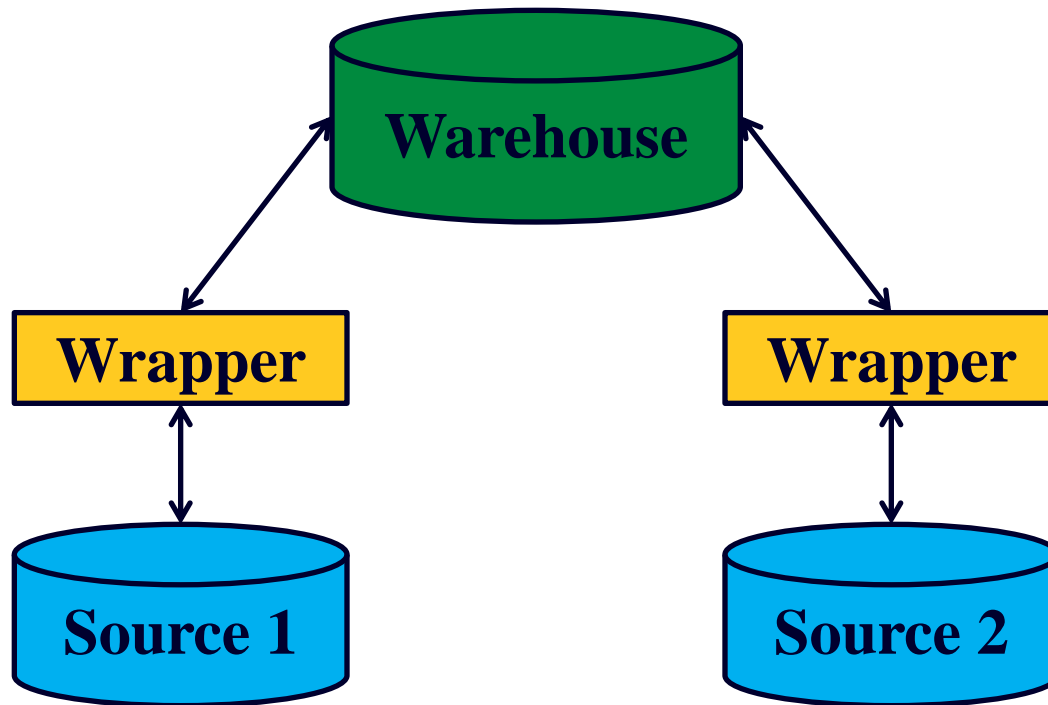
# Common Approaches

- **Federated Databases** Everybody talks directly to everyone else;
- **Warehousing** Sources are translated from their local schema to a global schema and copied to a central DB;
- **Mediation** Transforms a user query into a sequence of source queries.

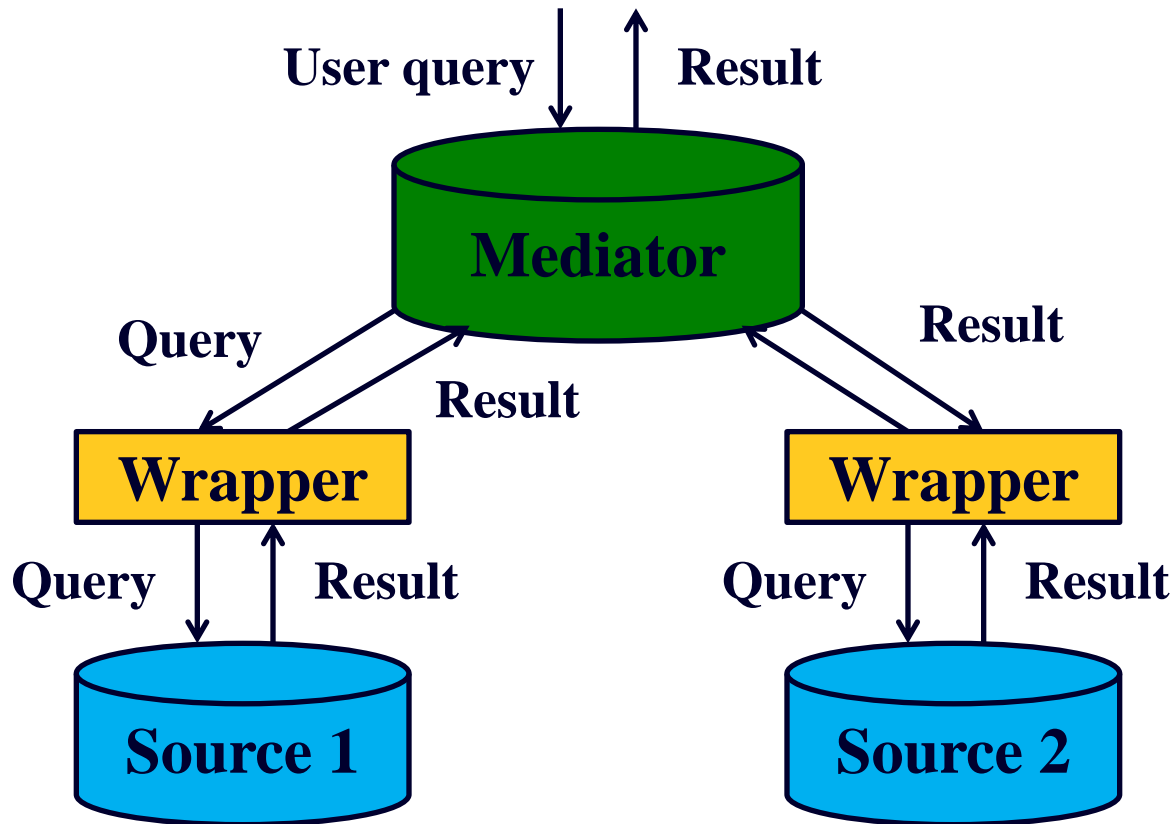
# Federated Database Systems



# Data Warehouses



# Mediators







# Two Mediation Approaches

- **Global as View** Global schema is defined as a view over the data sources;
- **Local as View** Sources are defined in terms of global schema; mediator finds all ways to build query from views.

# Example: Local View (1)

The mediator supports a **par(c,p)** relation:

- Source 1 provides some parent facts:

$$V1(c,p) \leftarrow \text{par}(c,p)$$

- Source 2 supports only grandparent facts:

$$V2(c,g) \leftarrow \text{par}(c,p) \text{ AND } \text{par}(p,g)$$

# Example: Local View (2)

- **Query** (great-grandparents):

$ggp(c,x) \leftarrow par(c,u) \text{ AND } par(u,v)$   
 $\text{AND } par(v,x)$

$Sol1(c,x) \leftarrow V1(c,u) \text{ AND } V1(u,v)$   
 $\text{AND } V1(v,x)$

$Sol2(c,x) \leftarrow V1(c,u) \text{ AND } V2(u,x)$

$Sol3(c,x) \leftarrow V2(c,u) \text{ AND } V1(u,x)$



# LAV and GAV Comparison

- **GAV** approach has better query processing capabilities (simpler to implement);
- **LAV** is more extensible (new data sources easier integrated).



# GLAV

**GLAV** approach combines the **GAV** and **LAV** approaches:

- Provides direct access to data sources;
- Provides capability to define global schema as a view over the data sources.

# Theoretical Basis (1)

Works of SYNTHESIS Group (<http://synthesis.ipi.ac.ru>):

- **AMN** is used to formalize data models;
- The canonical model should be extensible;
- A kernel of the canonical model is fixed. For each specific data model **M** of the environment an appropriate kernel extension is defined axiomatically;
- The canonical model is synthesized as a union of extensions, constructed for models **M** of the environment.

# Theoretical Basis (2)

In the frame of the considered formalism is offered to construct:

- A mapping from  $\mathbf{M}_s$  into an extension  $\mathbf{M}_t$ ;
- **AMN** semantics of  $\mathbf{M}_s$ ;
- **AMN** semantics of the extended  $\mathbf{M}_t$ .

After that the **B** - technology is applied to proof that  $\mathbf{M}_s$  is a refinement of the extension of  $\mathbf{M}_t$ . We say that specification  $\mathbf{S}_1$  refines specification  $\mathbf{S}_2$ , if it is possible to use  $\mathbf{S}_1$  instead of  $\mathbf{S}_2$  so that the user of  $\mathbf{S}_2$  does not notice this substitution.

# XDM as Kernel of Canonical Model

The following kinds of **XDM** objects are considered:

- **Basic Objects** constants of atomic types, variables of different sorts, etc.
- **Compound Objects** are defined in terms of *application* and *binding* in  $\lambda$  - calculus.



# Type System of XDM (1)

- To define metadata;
- To formalize the signature of **XDM** symbols;
- Formalism:
  - term calculus;
  - inference rules of judgments.

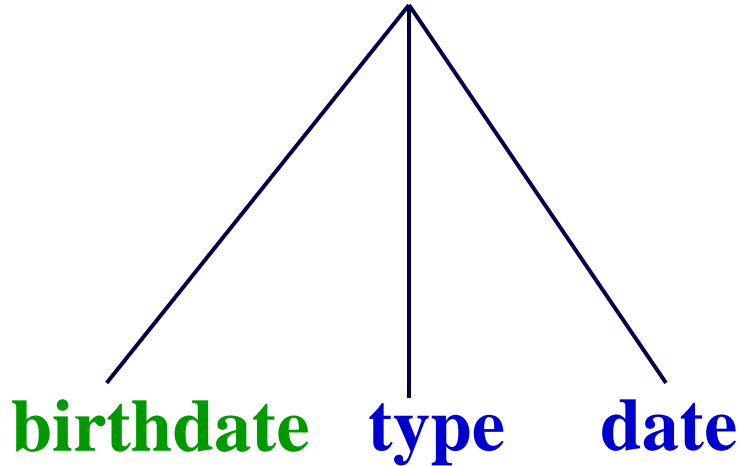
# Type System of XDM (2)

- **Basis** consists of conventional atomic types, for example, **int**, **string**, **boolean**, etc.
- *Attribution*. If **v** is a basic object variable and **t** is a typed object, then **attribution(v, type t)** is a typed object.

# Type System of XDM (3)

(Example)

**attribution**



**<xdmattr>**

**<xdmatp>**

**<xdms name="type" cd="ecc"/>**

**<xdms name="date" cd="ecc"/>**

**</xdmatp>**

**<xdmelm name="birthdate"/>**

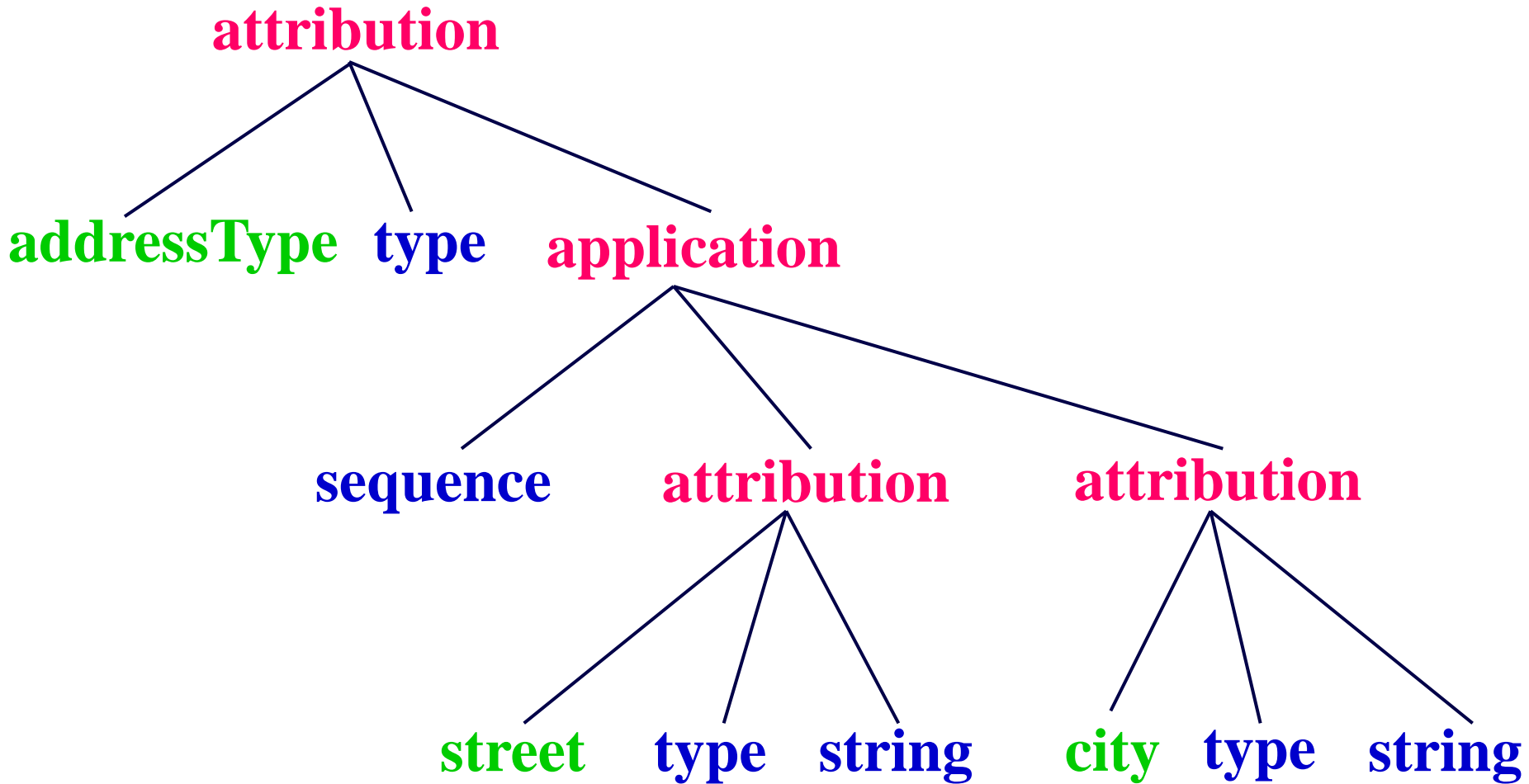
**</xdmattr>**

# Type System of XDM (4)

- *Abstraction*. If  $v$  is a basic object variable and  $t$ ,  $A$  are typed objects, then  $\text{binding}(\text{lambda}, \text{attribution}(v, \text{type } t), A)$  is typed object;
- *Application*. If  $F$  and  $A$  are typed objects, then  $\text{application}(F, A)$  is typed object;
- *Function Space*. If  $t$  and  $u$  are typed objects, and  $v$  is a basic variable object, then  $\text{binding}(\text{PiType}, \text{attribution}(v, \text{type } t), u)$  is typed object.

# Type System of XDM (5)

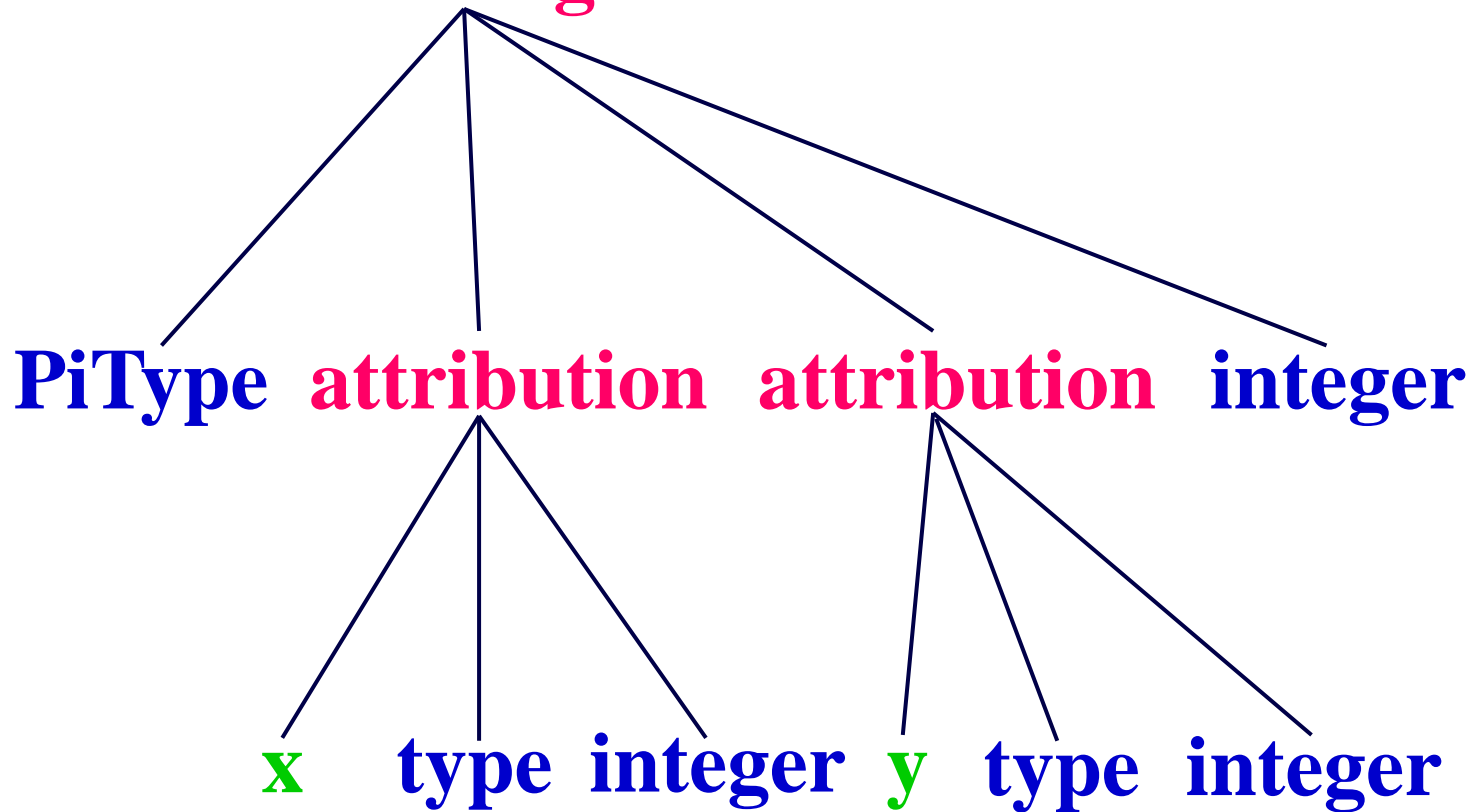
(Example)



# Type System of XDM (6)

(Example)

**binding**



**integer × integer → integer**

# Kernel Extension Principle (1)

- The extension of canonical model is formed by consideration of each new data model by adding new symbols to its DDL (symbols are used for data models concepts representation);
- For applying a concept on the canonical model level the following rule is proposed:

**Concept ← Symbol Context Definition**

# Kernel Extension Principle (2)

(Example)

- To support the concepts of **referential integrity** and **key's** of relational data model, we have expanded the kernel with the following symbols:  
**key, unique, foreign key, constraint, on update, on delete, cascade, and set null.**
- $S = \{\underline{Snum}, Sname, Status, City\}$
- $P = \{\underline{Pnum}, Pname, Color, Weight, City\}$
- $SP = \{\underline{Snum}, \underline{Pnum}, Qty\}$



# Kernel Extension Principle (3)

(Example)

**S** ← **attribution**(**S**, type TypeContext,  
constraint ConstContext)

TypeContext ← **application**(sequence, AppContext)

AppContext ← **attribution**(**Snum**, type int),  
**attribution**(**Sname**, type string),  
**attribution**(**Status**, type int),  
**attribution**(**City**, type string)

ConstContext ← **attribution**(**ConstName**, key Snum)

# Kernel Extension Principle (4)

(Example)

**SP** ← **attribution**(**SP**, **type** TypeContext,  
**constraint** KConstContext,  
**constraint** FKConstContext1,  
**constraint** FKConstContext2)

...

KConstContext ← **attribution**(ConstName, **key**  
**application**(list, **Snum**, **Pnum**))

FKConstContext1 ← **attribution**(ConstName,  
**foreign key** **application**(ref, ConstNameOFS),  
**on update cascade, on delete cascade**)

# AMN as Formal Language of Specification

**AMN** is a state-oriented formalism for software development and consists the following notations:

- the logical notation;
- the basic set notation;
- the relational notation;
- the mathematical object notation;
- the generalized substitution notation.



# Basic Constructions of AMN

- Abstract Machine;
- Refinement;
- Implementation.

# Abstract Machine Notation

**MACHINE** *MachineDeclaration*

**USES** *MachineList*

**SEES** *MachineList*

**INCLUDES** *InstantiatedMachineList*

**EXTENDS** *InstantiatedMachineList*

**SETS** *SetList*

**CONSTANTS** *ConstantList*

**PROPERTIES** *Predicate*

**VARIABLES** *VariableList*

**INVARIANTS** *Predicate*

**INITIALISATION** *Substitution*

**OPERATIONS** *OperationList*

**END**

# Refinement

REFINEMENT  $r$

REFINES  $m$

SEES  $sm$

INCLUDES  $im$

SETS  $s$

CONSTANTS  $c$

PROPERTIES  $P(s,c)$

VARIABLES  $x$

INVARIANTS  $I(x)$

INITIALISATION  $S$

OPERATIONS  $O_1; O_2; \dots ; O_n$

END

# Formalization of the Concept Refinement in AMN (1)

REFINEMENT **M**

REFINES **K**

CONSTANTS  $c_M$

PROPERTIES  $P_M$

VARIABLES **v**

INVARIANTS  $I_M$

INITIALISATION  $Init_M$

OPERATIONS

```
y ← op(x) =  
  PRE  $Pre_{op,M}$   
  THEN  $Def_{op,M}$   
  END
```

...

END

REFINEMENT **N**

REFINES **M**

CONSTANTS  $c_N$

PROPERTIES  $P_N$

VARIABLES **w**

INVARIANTS  $I_N$

INITIALISATION  $Init_N$

OPERATIONS

```
y ← op(x) =  
  PRE  $Pre_{op,M}$   
  THEN  $Def_{op,M}$   
  END
```

...

END

# Formalization of the Concept Refinement in AMN (2)

**Definition** We say that **M** refines **N** if holds the following *proof obligations*:

- $P_M \wedge P_N \Rightarrow \exists(v,w) \bullet (I_M \wedge I_N)$
- $P_M \wedge P_N \Rightarrow [\text{Init}_N] \neg ([\text{Init}_M] \neg I_N)$
- $P_M \wedge P_N \wedge I_M \wedge I_N \wedge \text{Pre}_{\text{op},M} \Rightarrow$   
 $\text{Pre}_{\text{op},N} \wedge [\text{Def}_{\text{op},N}\{y \rightarrow y'\}] \neg$   
 $([\text{Def}_{\text{op},M}] \neg (I_N \wedge y' = y))$



# AMN Formalization of Data Model (1)

Our approach to information integration assumes to create:

- **AMN** semantics for **XDM** (mapping from **XDM** to **AMN**);
- **AMN** semantics for each source data model (mapping from source data model to **AMN**);
- The canonical model extension by concepts of the source data model which are not yet present in target model.

After that the **B** - technology is applied to proof that source data model is a refinement of the extension of target model.

# AMN Formalization of Data Model (2)

First of all we introduce an abstract machine in which some common concepts of data models are included.

## MACHINE

**DatamodelContext**

## SETS

**OBJ, Types**

## ABSTRACT\_CONSTANTS

**atomicTypes, extent, typeOf**

## PROPERTIES

**$\text{atomicTypes} \in \mathcal{P}(\text{OBJ}) \wedge \text{extent} \in \text{Types} \rightarrow \mathcal{P}(\text{OBJ}) \wedge$   
 $\text{typeOf} \in \text{OBJ} \rightarrow \text{Types}$**

## END

# AMN Formalization of Data Model (3)

## (AMN semantics for XDM)

For each data model an abstract machines hierarchy is created which represents **AMN** semantics for that model. We use the following definition of **XDM** schema to create the **AMN** semantics for **XDM**:

**Definition** We say that **S** is a schema, if

$S = \langle \text{name}, \text{atomicity}, f \rangle$  or

$S = \langle \text{name}, \text{typeOp}(S_1, S_2, \dots, S_n), f \rangle$ , and  $S_i$  is a schema, where  $\text{typeOp} \in \{\text{sequence}, \text{choice}, \text{all}\}$ ,  $f \in \{?, *, +, \perp\}$ ,  $1 \leq i \leq n$ .

# AMN Formalization of Data Model (4)

## (AMN semantics for XDM)

**MACHINE** XDM\_Schema

(p\_name, p\_operationType, p\_argOfOperation, p\_frequency)

**CONSTANTS** Frequency, OperationTypes, FunctionSpace

**PROPERTIES**

Frequency = {?, \*, +,  $\perp$ }  $\wedge$

OperationTypes = {sequence, choice, all, empty}  $\wedge$

FunctionSpace = seq(Types) x Types

**VARIABLES** name, operationType, argOfOperation, frequency

**INVARIANT**

name  $\in$  STRING  $\wedge$  operationType  $\in$  OperationTypes  $\wedge$

argOfOperation  $\in$  P(*extent*(XDM\_Schema))  $\wedge$

frequency  $\in$  Frequency

**CONSTRAINTS** p\_name  $\in$  STRING  $\wedge$  p\_operationType  $\in$  operationTypes  $\wedge$  ...

**INITIALISATION** name := p\_name  $\wedge$  operationType := p\_operationType  $\wedge$  ...

**END**

# Extension of the Canonical Model

```
MACHINE MchT
SETS ST
CONSTANTS CT
PROPERTIES PT
VARIABLES VT
INVARIANTS IT
INITIALISATION S
END
```

```
MACHINE ConceptsST
SETS SST
CONSTANTS CST
PROPERTIES PST
END
```

```
MACHINE MchST
EXTENDS MchT
SEES ConceptsST
END
```

```
REFINEMENT Mchs
REFINES MchST
SEES ConceptsST
SETS Ss
CONSTANTS Cs
```

```
PROPERTIES Ps
VARIABLES Vs
INVARIANTS Is ∧ Conds
INITIALISATION S
END
```

# AMN Semantics of Relational Data Model (1)

**MACHINE**

**REL\_Schema**

**VARIABLES**

**schemaName, state**

**INVARIANT**

**schemaName  $\in$  STRING  $\wedge$**

**state  $\in$  P(STRING  $\times$  atomicTypes)**

**END**

# AMN Semantics of Relational Data Model (2)

**MACHINE**

**REL\_DB\_Schema**

**VARIABLES**

**dbSchemaName, schemas**

**INVARIANT**

**dbSchemaName**  $\in$  **STRING**  $\wedge$   
**schemas**  $\in$  **P**(*extent*(**REL\_Schema**))

**END**

# Mapping from Relational Data Model into Canonical Model (1)

**MACHINE**

**REL\_Canonical\_Mapping**

**USES**

**REL\_Schema, REL\_DB\_Schema, XDM\_Schema**

**OPERATIONS**

**xdm\_Schema**  $\leftarrow$  **REL\_attr\_to\_XDM(attr) =**

**PRE**

**attr**  $\in$  **string**  $\times$  **atomicTypes**  $\wedge$  **xdmSchema**  $\in$   
*extent*(**XDM\_Schema**)

**THEN**

**xdmSchema** = **XDM\_Schema**(*dom*(**attr**), *ran*(**attr**),  $\emptyset$ ,  $\perp$ )

**END**



# Mapping from Relational Data Model into Canonical Model (2)

$\text{xdm\_Schema} \leftarrow \text{REL\_Schema\_to\_XDM}(\text{relSchema}) =$

**PRE**

$\text{relSchema} \in \textit{extent}(\text{REL\_Schema}) \wedge$

$\text{xdmSchema} \in \textit{extent}(\text{XDM\_Schema})$

**THEN**

$\text{xdmSchema} = \text{XDM\_Schema}(\textit{name}, \textit{sequence}, \textit{recSequence}, \perp)$

**END**

$\text{xdm\_Schema} \leftarrow \text{REL\_DB\_Schema\_to\_XDM}(\text{relDbSchema}) =$

**PRE**

$\text{relDbSchema} \in \textit{extent}(\text{REL\_DB\_Schema}) \wedge$

$\text{xdmSchema} \in \textit{extent}(\text{XDM\_Schema})$

**THEN**

$\text{xdmSchema} = \text{XDM\_Schema}(\textit{name}, \textit{sequence}, \textit{recSequence}, \perp)$

**END**

**END**

# DML Formalization Principles in AMN

**MACHINE**

**RelationalAlgebra**

**USES**

**Relation**

**OPERATIONS**

$s \leftarrow \text{union}(r, q) =$

**PRE**

$r \in \text{extent}(\text{Relation}) \wedge q \in \text{extent}(\text{Relation}) \wedge$   
 $s \in \text{extent}(\text{Relation}) \wedge \text{isUnionCompatibility}(r, q)$

**THEN**

$s.\text{schema} := r.\text{schema} \parallel s.\text{tuples} := r.\text{tuples} \cup q.\text{tuples}$

**END**

# Summary & Future Work

## ■ Canonical Model

- The XDM has been formalized in the frame of the AMN,
- Principle of extension of canonical model has been proposed and formalized in the frame of AMN.

## ■ Example of Mappings

- The RDM has been formalized in the frame of AMN,
- Mapping from AMN representation of RDM to canonical model has been proposed.

## ■ In the future we plan to develop

- Algorithms for generation of mappings from the source models into canonical model,
- Temporal extension of canonical model.



**THANK YOU FOR YOUR ATTENTION!**