

SDMv2 Status

CALIM-2011, Manchester Jul., 2011

F. Viallefond.



Laboratoire d'Étude du Rayonnement et de la Matière en Astrophysique

Outline

1. What is a Data Model?

- Domain \rightarrow Language
- Structures
- Data Model

2. Formalization

- Categories, Functors
- Sketches, Models and Theories

3. Examples of diagrams

- Primitive Data Types
- Initial and Final objects
- Products, Coproducts, Direct limits

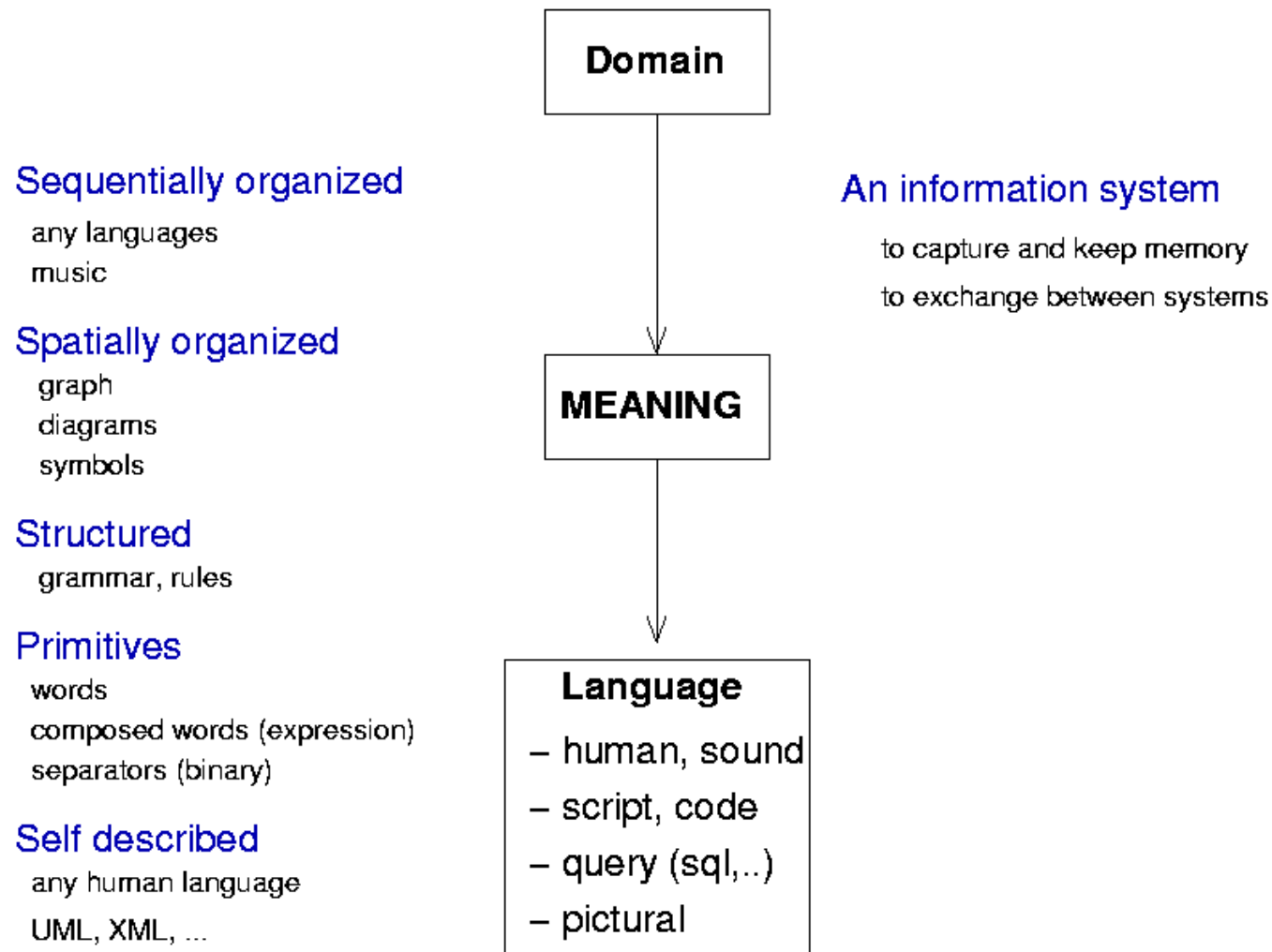
4. Relation, Recursive constructs, topology of the MSDM (SDMv2)

5. Monads, lambda calculus

6. Type specification: the measurements and their topology

7. Conclusions and perspectives

Domain and Language



To represent a set of measurements

Examples of words (*physical quantities*):

- Length, Area, Angle, Solid angle, Aperture efficiency, Rotation measure
- Speed
- Angular rate
- Noise equivalent power
- FluxDensity (*Jy which is not SI...*)
- ...

Note that:

1. All these have units.
2. Dimensioned, dimensionless and mixed case units!
3. They may have units which uses powers of rational numbers!
4. Physical expressions are composition of such words

Measurements in context

We assign domain specific meaning to words:

- Station
- Antenna
- Spectral window
- Feed
- Configuration description
- ...

Note that:

- We conceptualize
- We compose
- There are context-independencies and context dependencies
- ...

Motivations to have a data model

A measurement set is a set of concrete concepts at different levels,
a) words such as the physical quantities (universal concepts),
b) compositions of words giving rise to relations.

We must all share a common understanding of these concepts.

These must be easily usable in information systems (data reduction packages, DBMS, ...).

These must be *a)* concisely described to insure reliability and *b)* properly understood.

Very important for optimization of calculus

(architecture: structure, factorization, localization, slicing, ... i.e. geometry),

The measurement set to become a data model

The mathematicians

- 1/ have developed all the abstract concepts that we use (often implicitly...)
- 2/ give a methodology defining what is a **model** and a **theory**!

The theory of categories: used in fundamental computer science.

Implementation of the SDMV2 is based on **generic programming techniques**.

The work of conceptualization which was required fits (to my surprise) very well with that theory!

Status and prospect: The 'SDMV2' will be fully explained by **a theory**.

From 2008 to now...

1. 2008: devel concepts of phys. quan. and several important structures, proof that it was doable demonstrated by a small proto
2. 2009: more generics (EMBRACE context), radiotelescope generic
3. 2010/11 mathematics (very steep learning curve...), converter ASDM to SDMV2
4. 2011 'reverse engineering' to elaborate the theory, more algebraic types in the implementation.

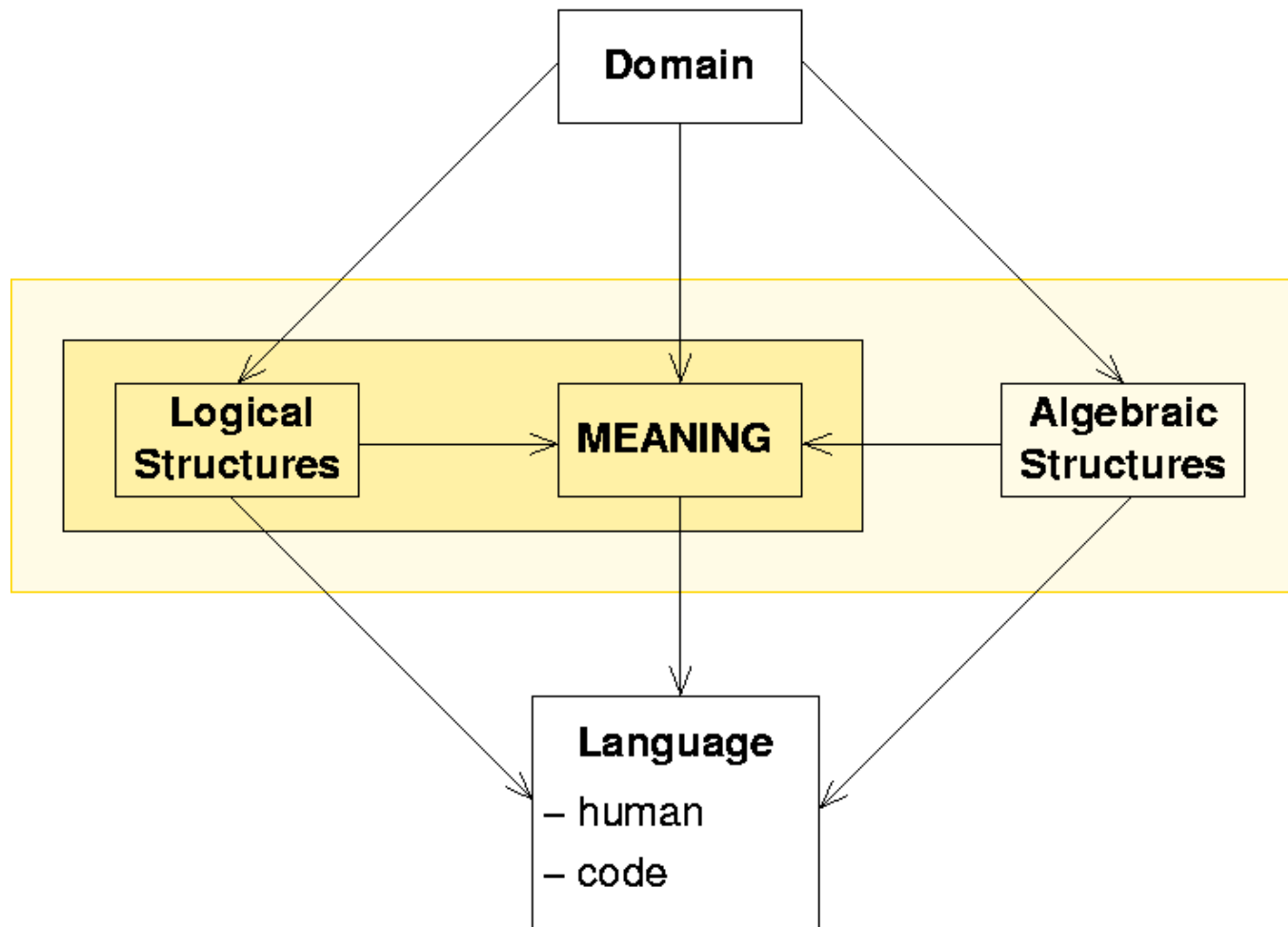
What is a data model?

A model is the composition of a structure (mathematical logic) with algebra.

Example: the relational data model.

- The semantic is captured through constraints.
- The structure gives the meaning of things in a formal language.

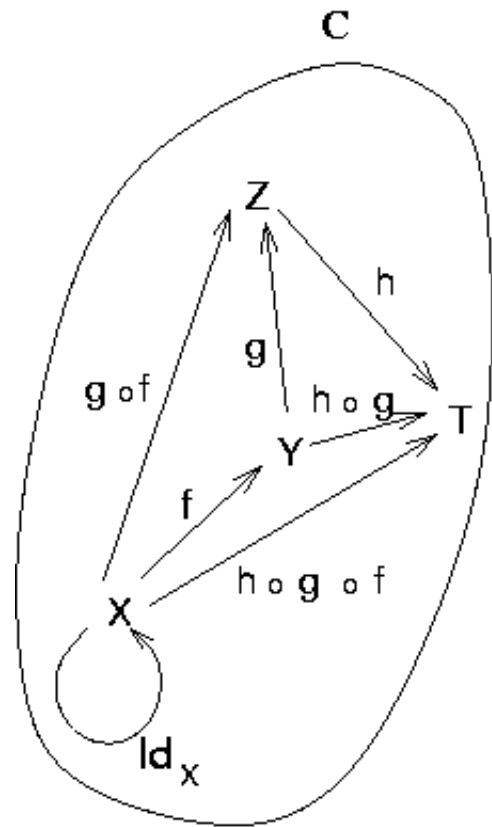
4 commutable triangles



Formalization

- Category
- Functor
- Natural transform
- Product and coproduct:
example of diagrams, a cone and a cocone
- Direct limit
- Sketches, Models and Theories

Category C



Collection of objects:

X, Y, Z, T

Morphisms of objects:

f, g, h

- Identity:

$$\forall X \in \mathbf{C} \exists \text{Id}_X \in \mathbf{C}$$

- Transitive composition:

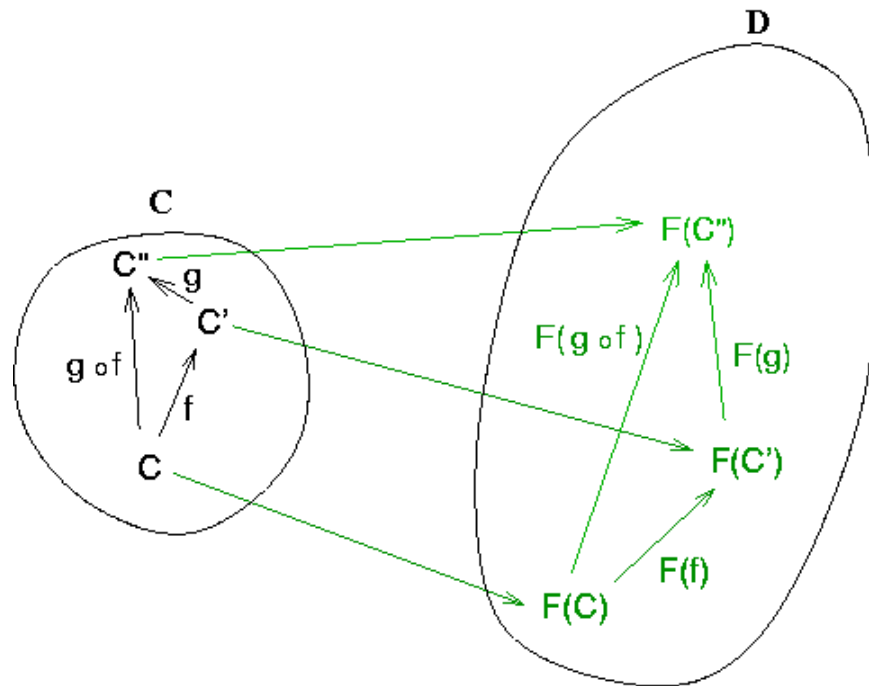
$$\begin{array}{c}
 X \xrightarrow{f} Y \xrightarrow{g} Z \\
 \searrow \quad \quad \nearrow \\
 \quad \quad \quad g \circ f
 \end{array}$$

- Associativity:

$$\begin{array}{ccc}
 & (h \circ g) \circ f = h \circ (g \circ f) & \\
 X & \xrightarrow{\quad} & Z \\
 f \downarrow & \searrow \quad \quad \nearrow & \uparrow h \\
 & g \circ f & h \circ g \\
 Y & \xrightarrow{\quad} & T \\
 & g &
 \end{array}$$

Functor

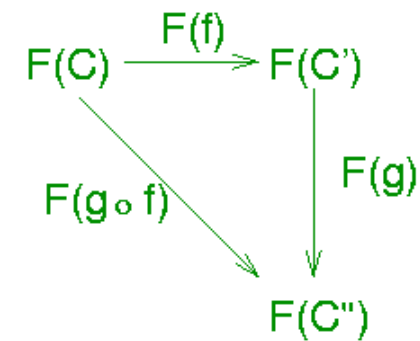
$$F: \mathbf{C} \rightarrow \mathbf{D}$$



Two categories \mathbf{C} and \mathbf{D}

The morphism $F: \mathbf{C} \rightarrow \mathbf{D}$
is a functor if:

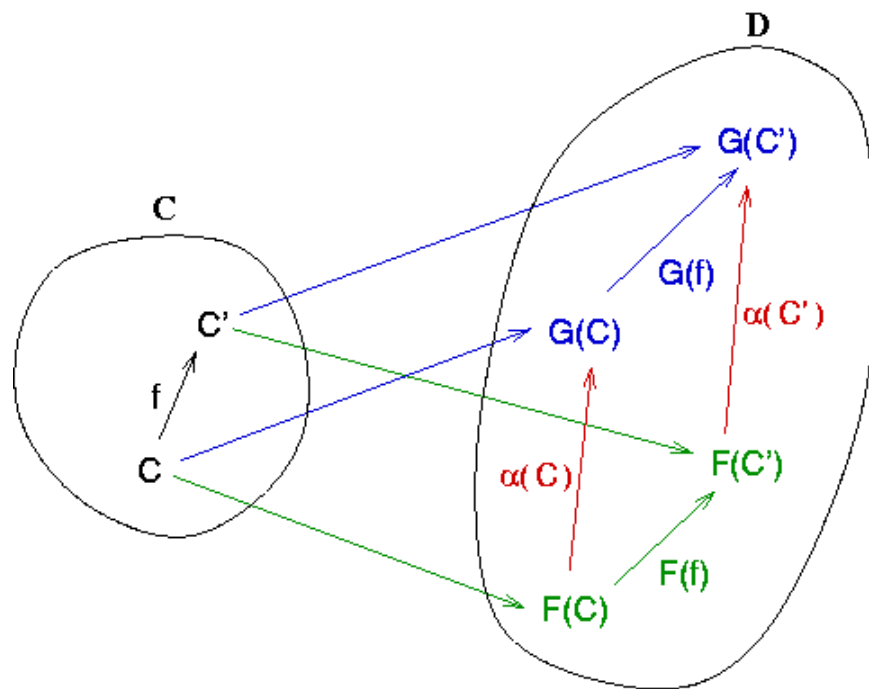
- $\forall C \in \mathbf{C} \exists F(C) \in \mathbf{D}$
- $F(\text{Id}_C) = \text{Id}_{F(C)}$
- and the diagram



is commutative

Natural transformation

$$\alpha: F \rightarrow G$$



Two functors:

$$F: C \rightarrow D$$

$$G: C \rightarrow D$$

The morphism

$$\alpha: F \rightarrow G$$

is natural if the diagram

$$\begin{array}{ccc} F(C) & \xrightarrow{F(f)} & F(C') \\ \alpha(C) \downarrow & & \downarrow \alpha(C') \\ G(C) & \xrightarrow{G(f)} & G(C') \end{array}$$

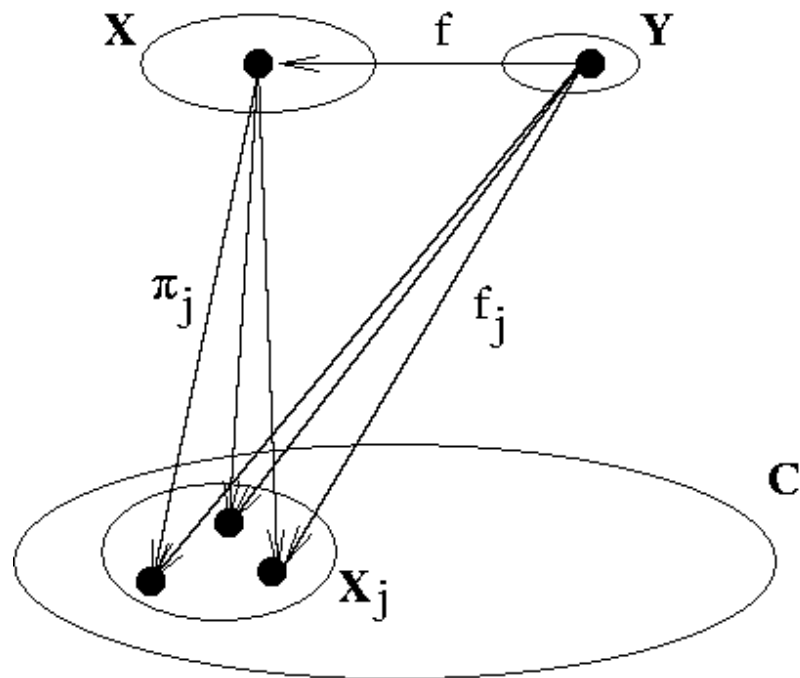
is commutative

Product

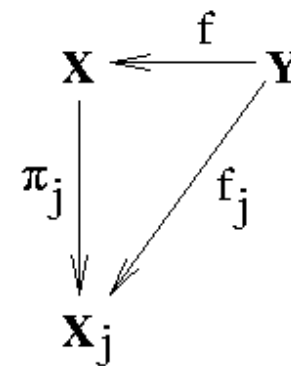
an projective cone

$$\mathbf{X} = \prod_{j \in J} \mathbf{X}_j$$

$f_j = f \circ \pi_j$ is unique



a commutative diagram



a product of morphisms

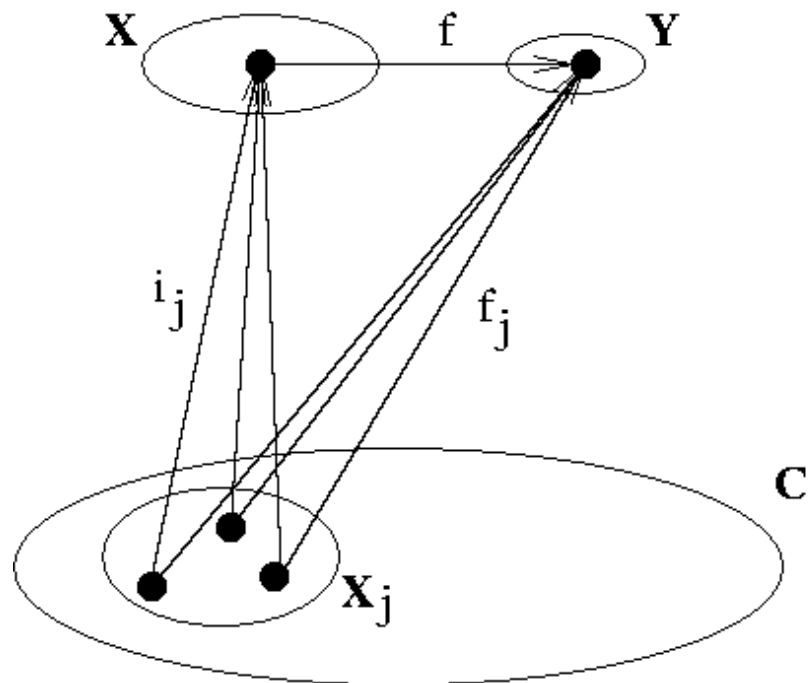
$$\langle f_1 \dots f_n \rangle$$

Coproduct

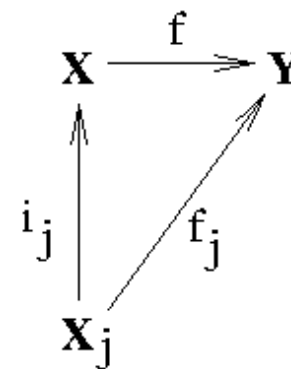
an injective cone

$$\mathbf{X} = \bigoplus_{j \in J} \mathbf{X}_j$$

$f_j = f \circ i_j$ is unique



a commutative diagram



Coproduct = colimit
an inductive limit

$$\mathbf{X} = \varinjlim \mathbf{X}_i$$

a direct set $\langle I, \leq \rangle$

a direct system $\langle \mathbf{X}_i, f_{ij} \rangle$

a disjoint union $\mathbf{X} = \varinjlim \mathbf{X}_i$

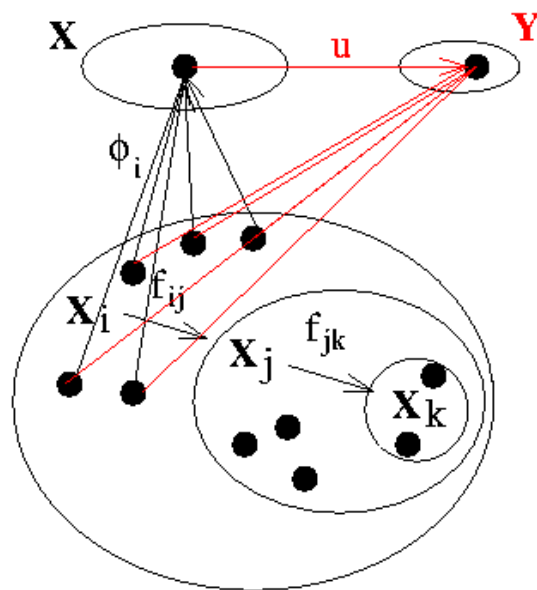
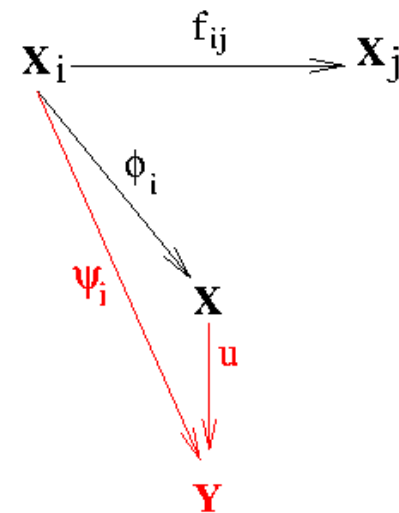


diagram commutes $\forall i, j$



$u: \mathbf{X} \rightarrow \mathbf{Y}$ is unique $\forall i, j$

Direct limit = colimit
an inductive limit

$$\mathbf{X} = \varinjlim \mathbf{X}_i$$

a direct set $\langle I, \leq \rangle$

a direct system $\langle \mathbf{X}_i, f_{ij} \rangle$

a disjoint union $\mathbf{X} = \varinjlim \mathbf{X}_i = \bigoplus_i \mathbf{X}_i / \sim$

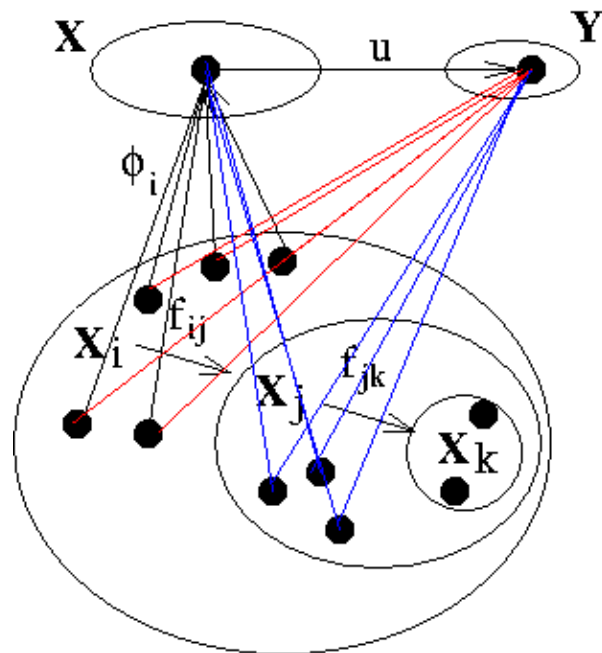
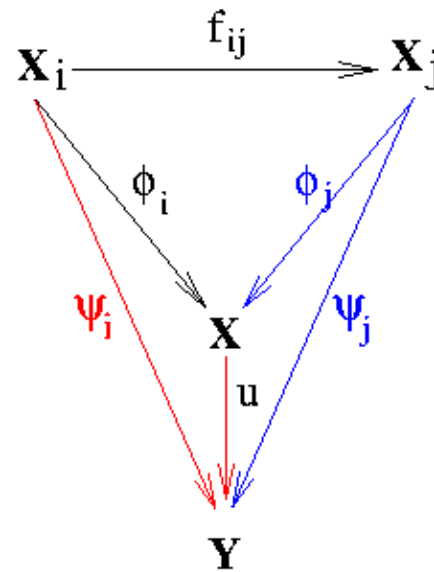


diagram commutes $\forall i, j$



$u: \mathbf{X} \rightarrow \mathbf{Y}$ is unique $\forall i, j$

Data base axes

- Relation

Relational model: $\text{Hom}_{\text{Set}^{\text{Rel}}}$ is empty

Example: Table<Antenna>

- Recursion

Use-cases require a cone or co-cone (projection or induction)

Example: Table<Antenna> (in case of APA and PAF)

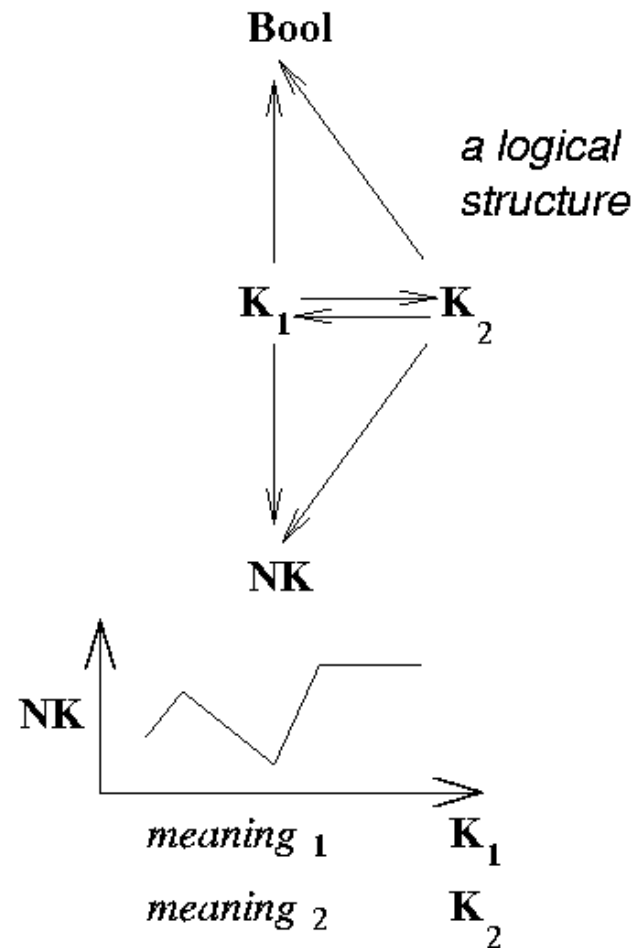
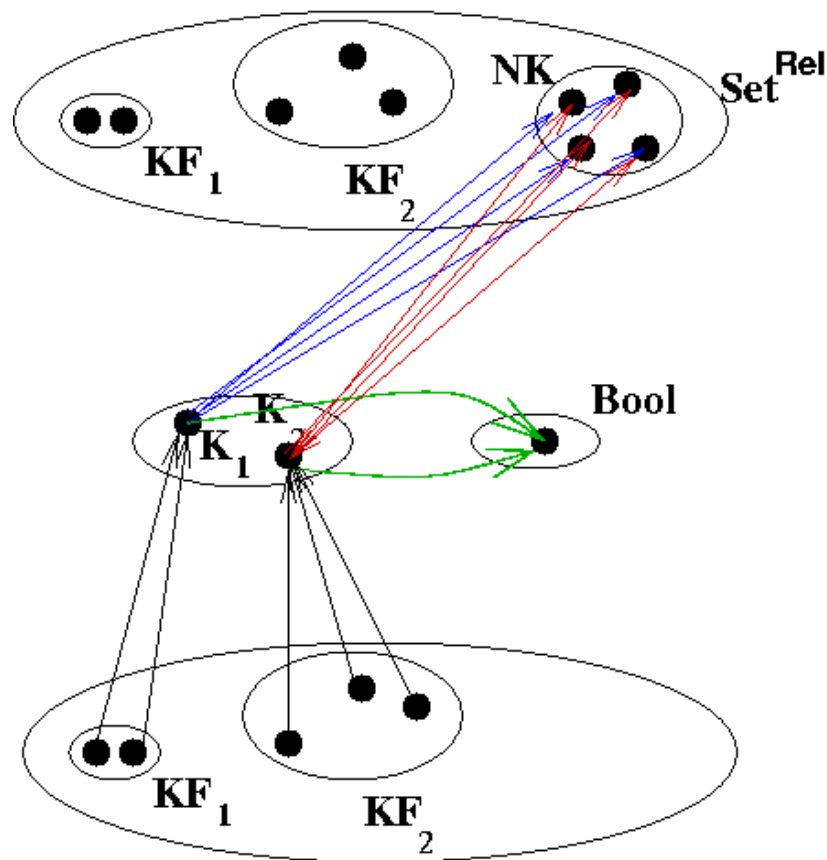
- Monoids: A table is a data base

$\text{Mon}(\text{Set}^{\text{Rel}}, +)$

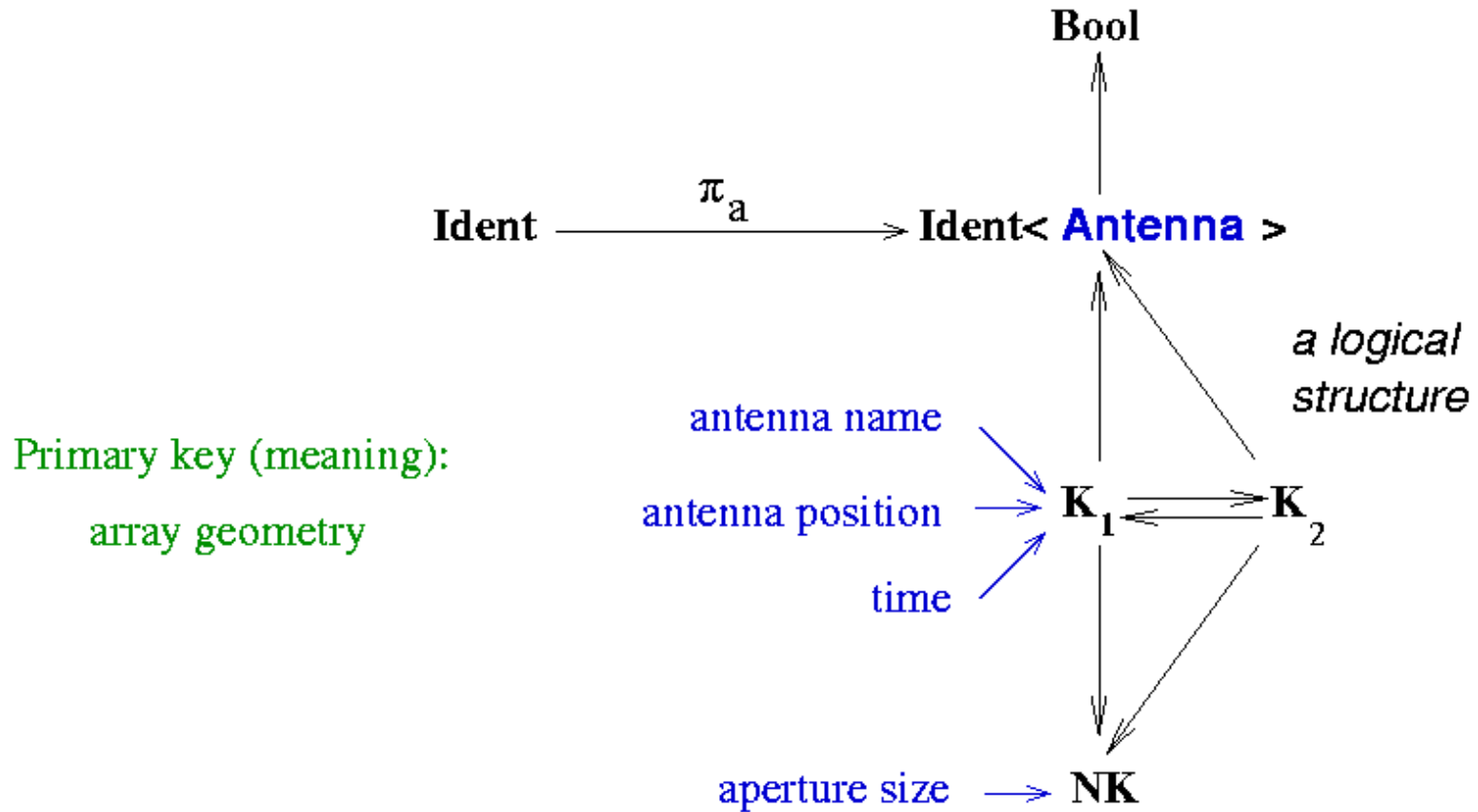
Relation the category Set^{Rel}

Defined using projective and injective cones

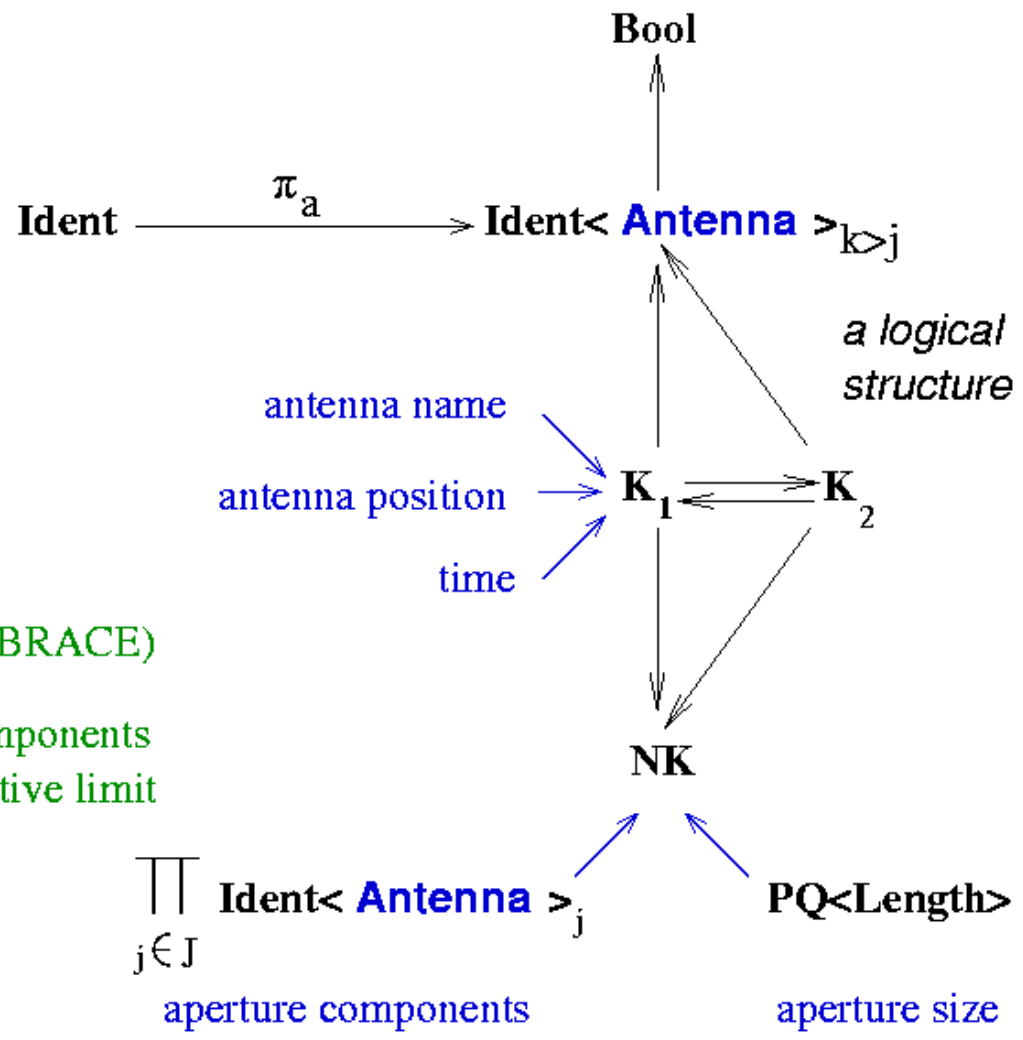
Must conform to the normal forms (logic based)



Relation Table < **Antenna** >



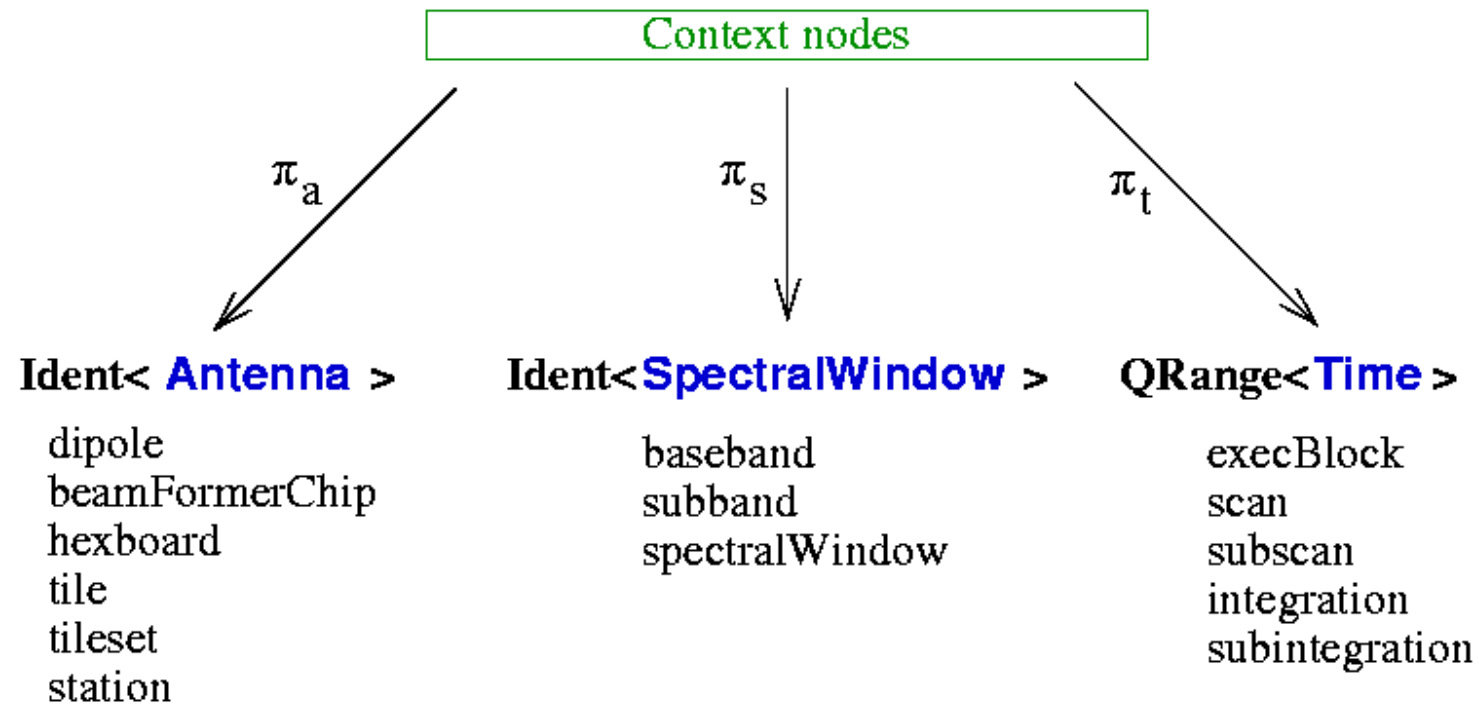
Relation Table < Antenna >



Use-case of APA (e.g. EMBRACE)

Antenna \longrightarrow Antenna components
 A recursive object, a projective limit

Topological space axis basis



↔ ● ↔
antennaProcessor

→ ● ←
downConverter
polyPhaseFilter
tunableFilter
correlator

↔ ● ↔
obsExecutor

→ ● ←
integrator

Processors

Monad

Let consider a functor $T: C \longrightarrow C$ and 2 natural transformations:

the induction $\eta_t: \text{ident} \longrightarrow T(\text{ident})$

the direct system $\mu_t: T(T(\text{ident})) \longrightarrow T(\text{ident})$

We have:

the monad (T, η_t, μ_t)

$$\begin{array}{ccc}
 T^3 & \xrightarrow{\mu_t T} & T^2 \\
 \mu_t T \downarrow & & \downarrow \mu_t \\
 T^2 & \xrightarrow{\mu_t} & T
 \end{array}
 \qquad
 \begin{array}{ccc}
 T & \xrightarrow{T \eta_t} & T^2 \\
 \eta_t T \downarrow & \parallel & \downarrow \mu_t \\
 T^2 & \xrightarrow{\mu_t} & T
 \end{array}$$

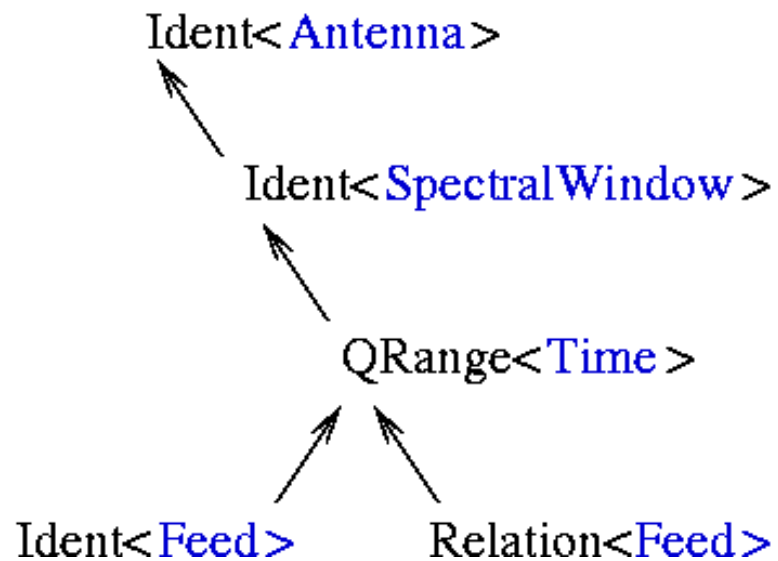
its T -algebra (t, h)

$$\begin{array}{ccc}
 T^2 t & \xrightarrow{Th} & T t \\
 \mu_t \downarrow & & \downarrow h \\
 T t & \xrightarrow{h} & t
 \end{array}
 \qquad
 \begin{array}{ccc}
 t & \xrightarrow{\eta_t} & T t \\
 \downarrow l_t & & \downarrow h \\
 & & t
 \end{array}$$

Application: Table< Feed >

Proposition: A table is a monad which has for its algebraic structure the vector space, the directed set

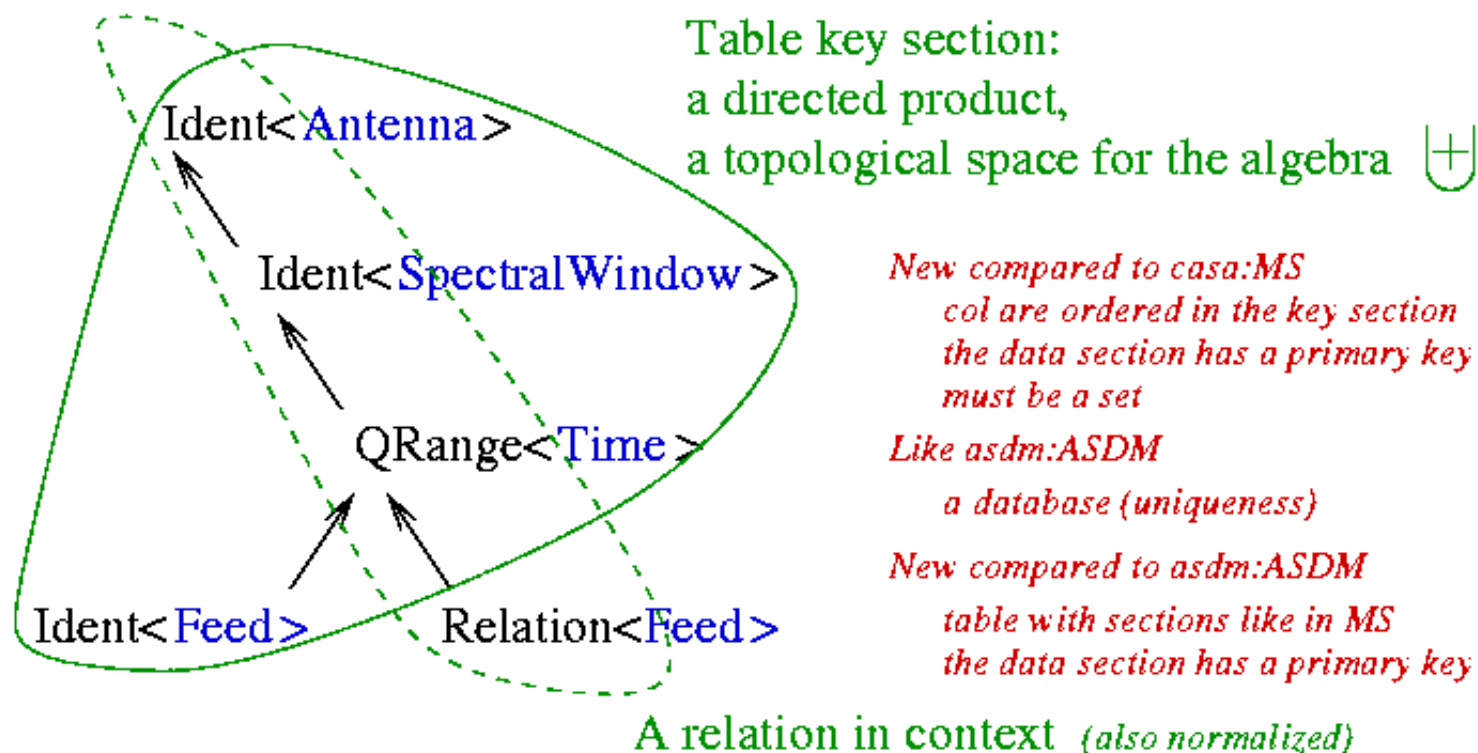
Ident<Antenna> \longrightarrow Ident<SpectralWindow> \longrightarrow QRange<Time>



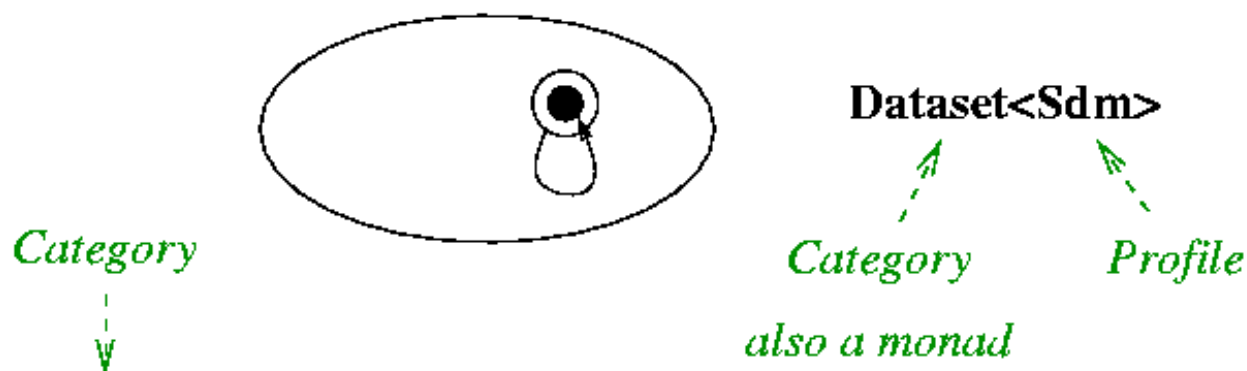
Application: Table< Feed >

Proposition: A table is a monad which has for its algebraic structure the vector space, the directed set

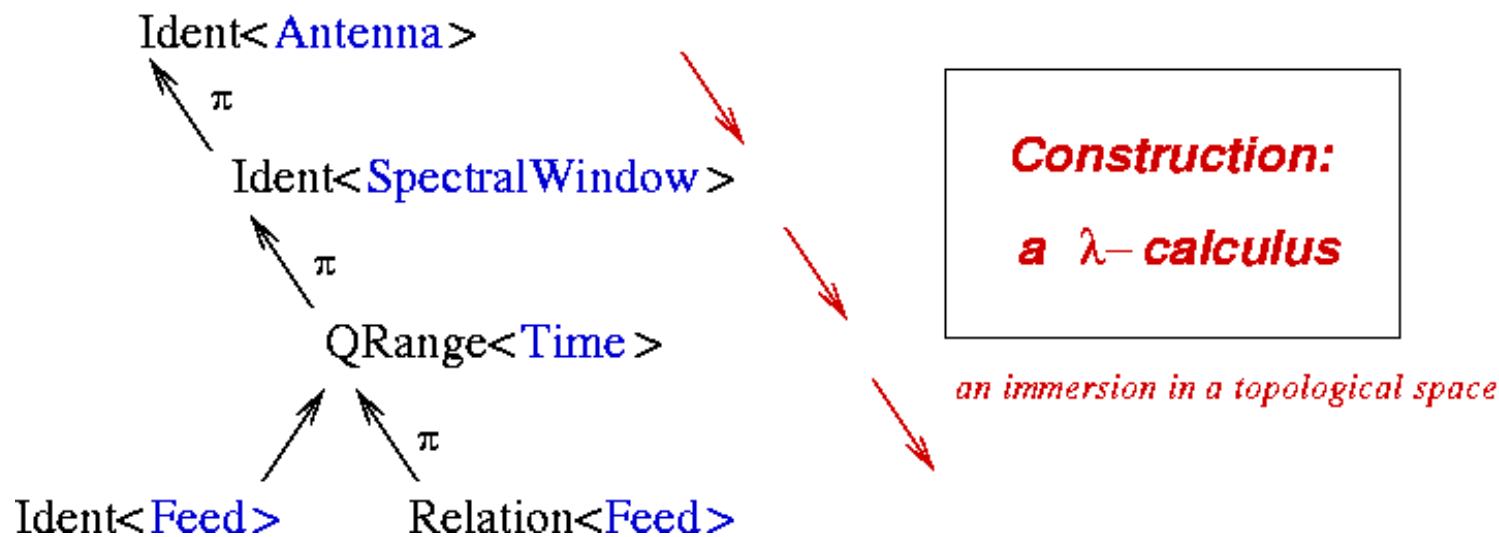
Ident<Antenna> \longrightarrow Ident<SpectralWindow> \longrightarrow QRange<Time>



Tables, singletons, are monads in a category Dataset



Example: Table< Feed > in the SDMv2



Where is the meaning?

We have seen that a relation is the composition of attributes.

Relations are named and their meanings are bound to constraints, rules of constructions.

An attribute: association “column name,type” (*NB: relational model: type=PDT*).

Names are words, a subset of the set **String**.

In practice very convenient but at the cost of several deficiencies:

1. **Not efficient**: parsing string is expensive. Solutions to this are:

carry the meaning via names of variables and use PDT

but

2. **Not type-safe**: codes not robust
3. **Meanings embedded in codes**: semantic embedded in the codes
4. **Tendency to use implicit words**: information 'hidden' in documentation

Proposition: Attribute (*in relation*) must rely on data models.

at the cost of departures from the relational data model

Type specification

A type is determined by a domain.

a set which gives all the possible values:

basis of a type.

SDMv2 data types have data models:

basis + logical structure + algebraic structure

Type specification: QValue(s) data model:

1. Abstract QValue data model

- A value is parametrized by a primitive data type
- A value may be *undefined* or *actual* or *virtual*
virtual values are expressions parametrized by a function type:
POLYNOMIAL, STAIRCASE, CHEBYSHEV

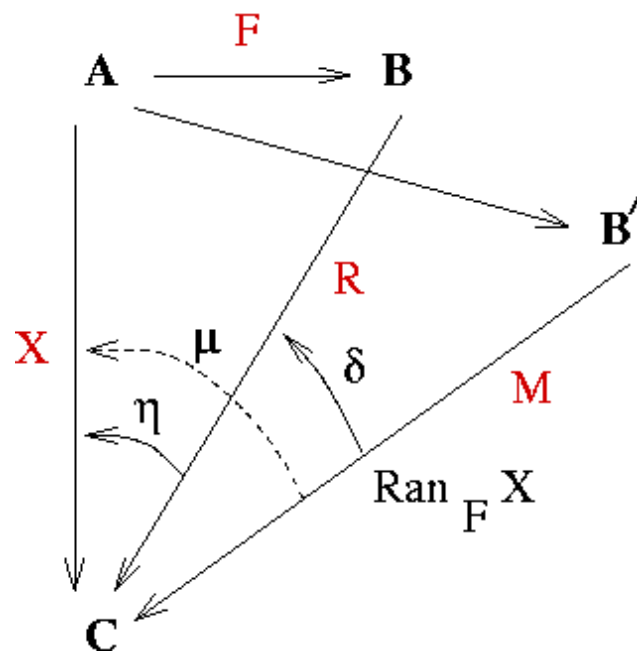
2. Abstract QValues data model

- A collection of values has meta-data to describe its structure
Any collection of values is parametrized by a dimensionality to discriminate
scalars from **vectors**
- The actual content of a collection is independent of its representation.
there is a **polymorphic representation to support data compression.**
- Collections have iterators

Type specification: Simple use-case: PQuantity

1. A categorical construct: the Kan extension
2. Handling multiple system of units
3. Its topology: a 2-category on a vector space
4. Its algebraic structure to get its categorical theory

Kan extension



3 categories, A, B, C

3 functors, F, M, X

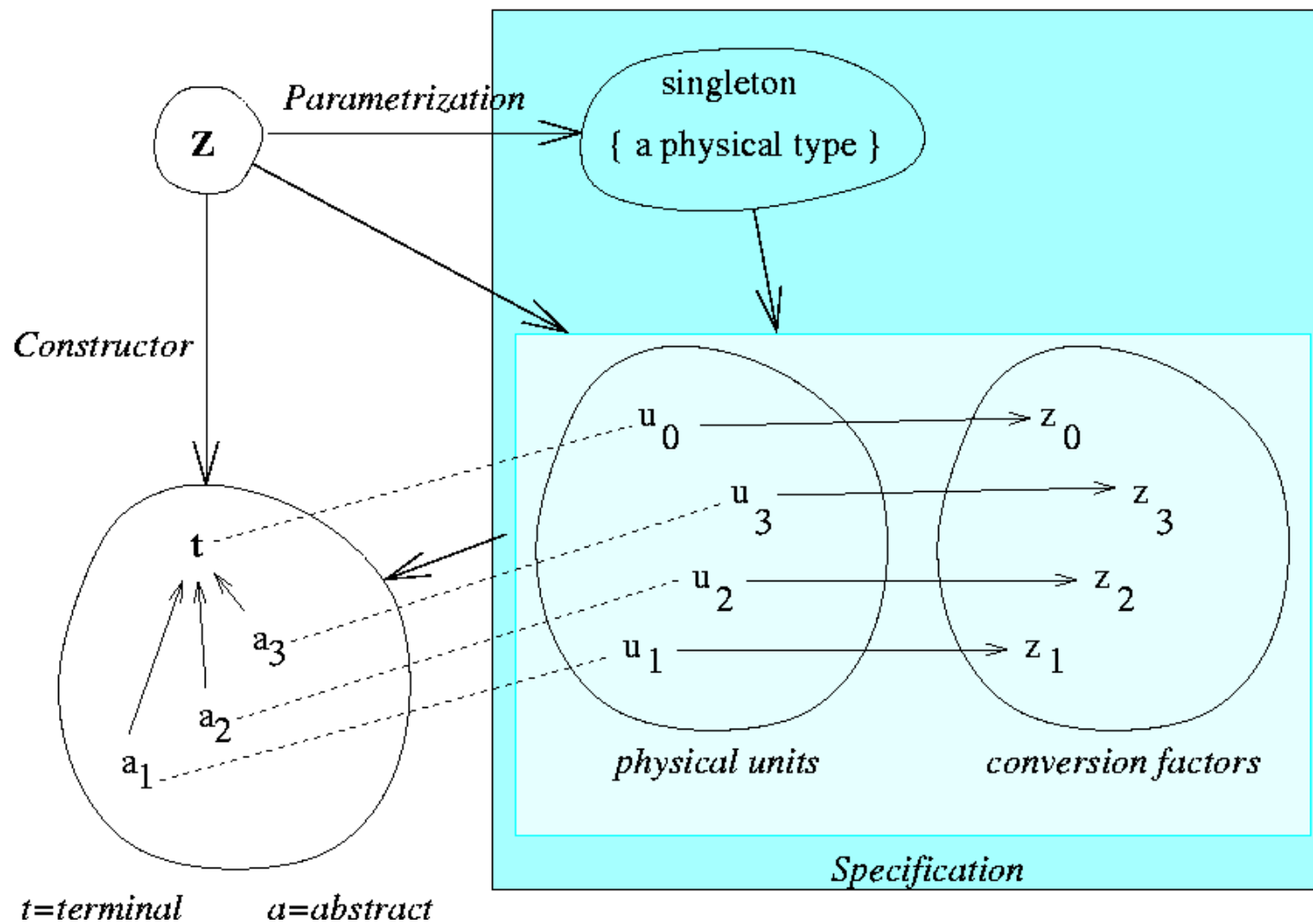
2 natural transformations η, δ

$$\forall M: B \longrightarrow C$$

Let μ be a natural transformation

$$\forall \mu: MF \longrightarrow X$$

$$\implies \delta: M \longrightarrow R \text{ is unique}$$

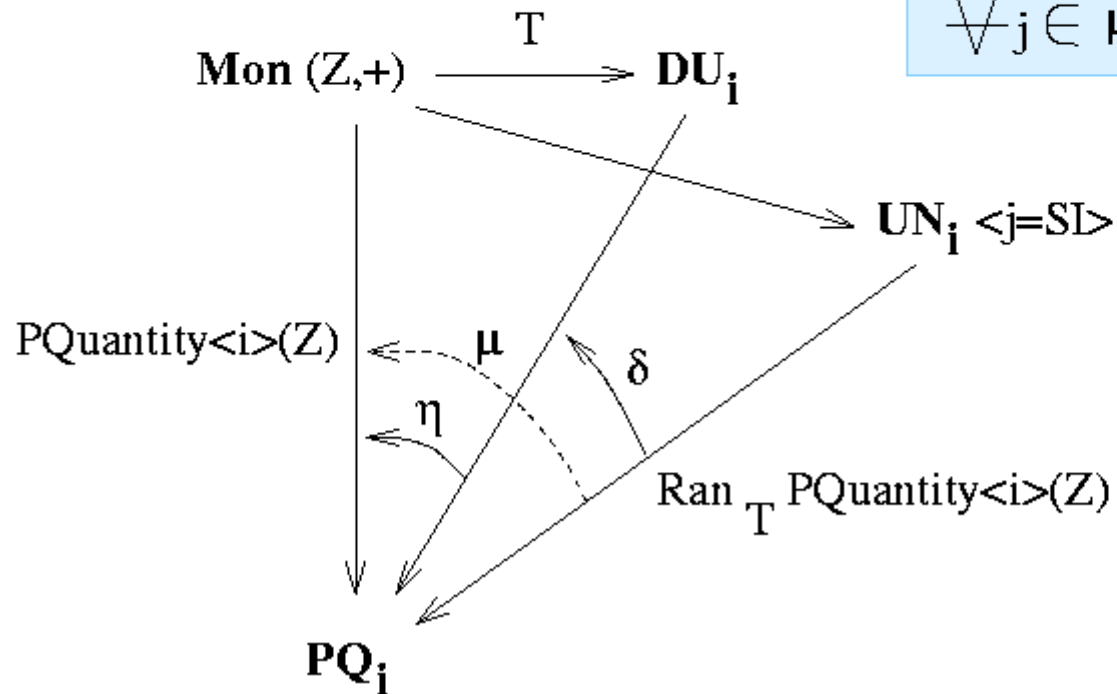


PQuantity data type

$\delta_i : \text{units}_i \rightarrow \text{canonic units}_i$

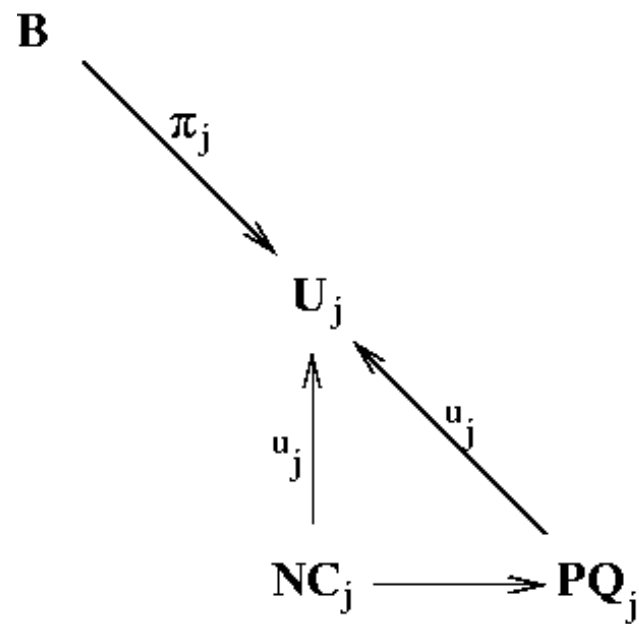
$\forall i \in \{ \text{PQname} \}$

$\forall j \in \mu_{i,j} : \text{units}_i \{j\} \rightarrow Z$



PQuantity

Let U_j a unit element along the axis j in a vector space B an object $\in \mathbf{Vect}$
Let $NC_j \in \mathbf{R}$ the dimension unit along the axis j of a point in that space



Examples with $j=0$

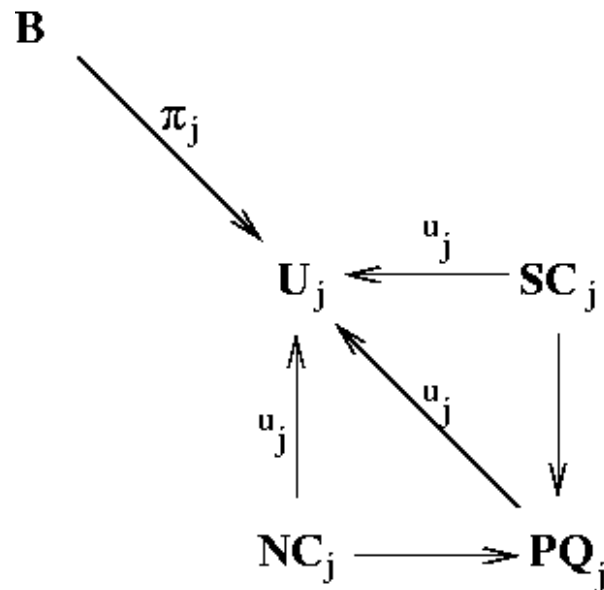
<i>Phys quan. k</i>	$NC_{0,k}$
<i>Length</i>	<i>1</i>
<i>Area</i>	<i>2</i>
<i>SpatialFrequency</i>	<i>-1</i>

PQuantity

Let U_j a unit element along the axis j in a vector space \mathbf{B} an object $\in \mathbf{Vect}$

Let $NC_j \in \mathbf{R}$ the dimension unit along the axis j of a point in that space

Let $SC_j \in \mathbf{R}$ the dimension unit ratio along the axis j of a point in that space



Examples with $j=0$

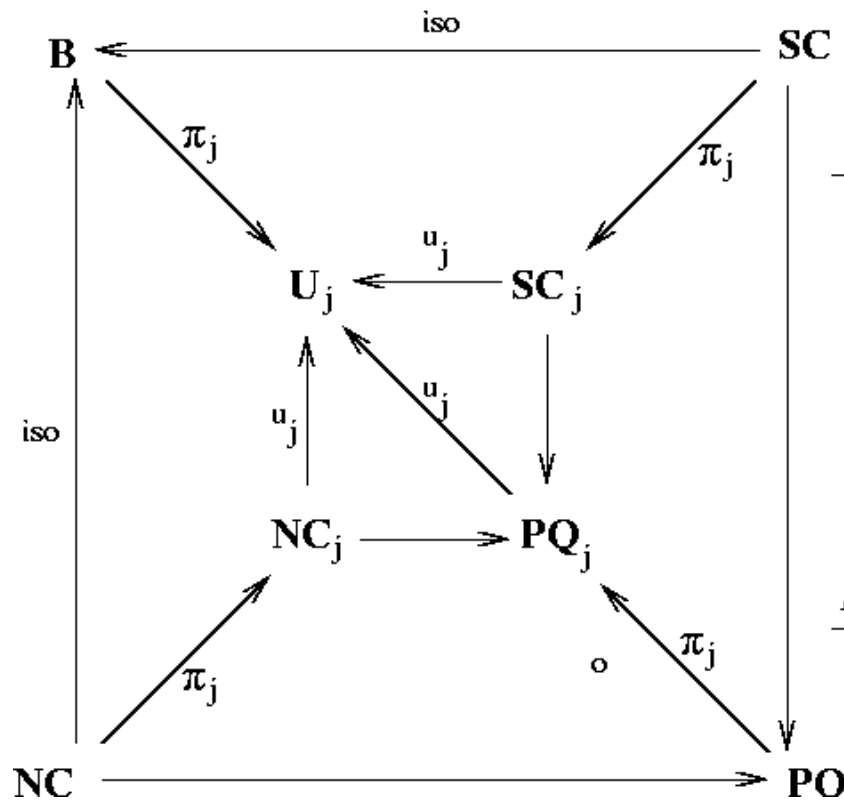
<i>Phys quan. k</i>	$NC_{0,k}$	$SC_{0,k}$
<i>Length</i>	<i>1</i>	<i>0</i>
<i>Area</i>	<i>2</i>	<i>0</i>
<i>SpatialFrequency</i>	<i>-1</i>	<i>0</i>
<i>Angle</i>	<i>0</i>	<i>1</i>
<i>SolidAngle</i>	<i>0</i>	<i>2</i>
<i>RotationMeasure</i>	<i>-2</i>	<i>1</i>

PQuantity

Let U_j a unit element along the axis j in a vector space \mathbf{B} an object $\in \mathbf{Vect}$

Let $NC_j \in \mathbf{R}$ the dimension unit along the axis j of a point in that space

Let $SC_j \in \mathbf{R}$ the dimension unit ratio along the axis j of a point in that space



Examples with $j=0$

<i>Phys quan. k</i>	$NC_{0,k}$	$SC_{0,k}$
<i>Length</i>	1	0
<i>Area</i>	2	0
<i>SpatialFrequency</i>	-1	0
<i>Angle</i>	0	1
<i>SolidAngle</i>	0	2
<i>RotationMeasure</i>	-2	1

Examples with $j=0$ and 2

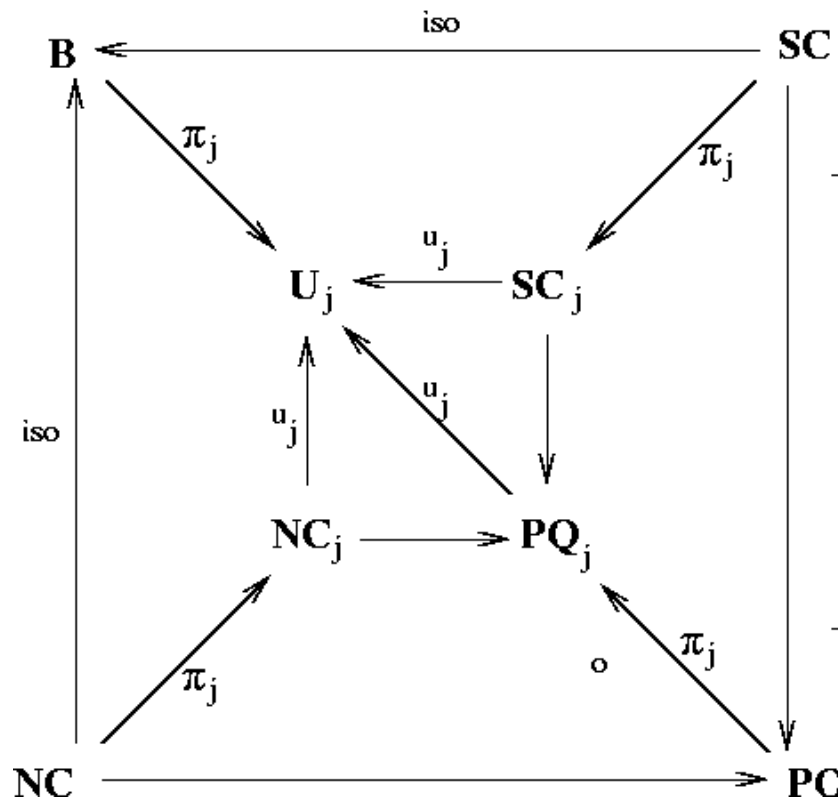
<i>Phys quan. k</i>	$NC_{0,k}$	$SC_{0,k}$	$NC_{2,k}$
<i>Speed</i>	1	0	-1
<i>AngularRate</i>	0	1	-1

PQuantity

Let U_j a unit element along the axis j in a vector space \mathbf{B} an object $\in \mathbf{Vect}$

Let $NC_j \in \mathbf{Q}$ the dimension unit along the axis j of a point in that space

Let $SC_j \in \mathbf{Q}$ the dimension unit ratio along the axis j of a point in that space



Examples with $j=0$

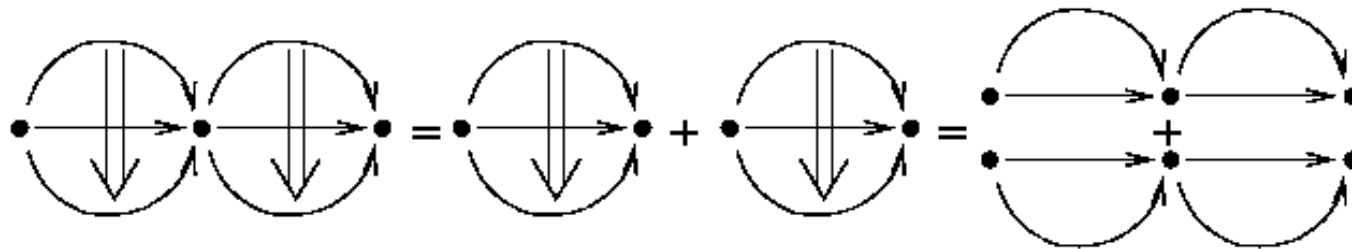
<i>Phys quan. k</i>	$NC_{0,k}$	$SC_{0,k}$
<i>Length</i>	1	0
<i>Area</i>	2	0
<i>SpatialFrequency</i>	-1	0
<i>Angle</i>	0	1
<i>SolidAngle</i>	0	2
<i>RotationMeasure</i>	-2	1

Examples with $j=0$ and 2

<i>Phys quan. k</i>	$NC_{0,k}$	$SC_{0,k}$	$NC_{2,k}$
<i>Speed</i>	1	0	-1
<i>AngularRate</i>	0	1	-1
<i>NEP</i>			$-1/2$

PQuantity

The topology of PQuantity is a 2–category on a vector space.



this vector space is of dimension 7:

- 0 Length
- 1 Mass
- 2 Time
- 3 Temperature
- 4 LuminousIntensity
- 5 MolarConcentration
- 6 ElectricCurrent

the horizontal composition along the fundamental physical unit basis
the vertical composition for the dimension,dimensionless property

PQuantity is a monoid for the addition.

Which diagram we need to add to get the theory?

Physical expressions require the addition and the multiplication (ring)

Let consider the following composite graph:

2 monoids PQ_a and PQ_b and the functor m such that $PQ_c = PQ_a \times PQ_b$

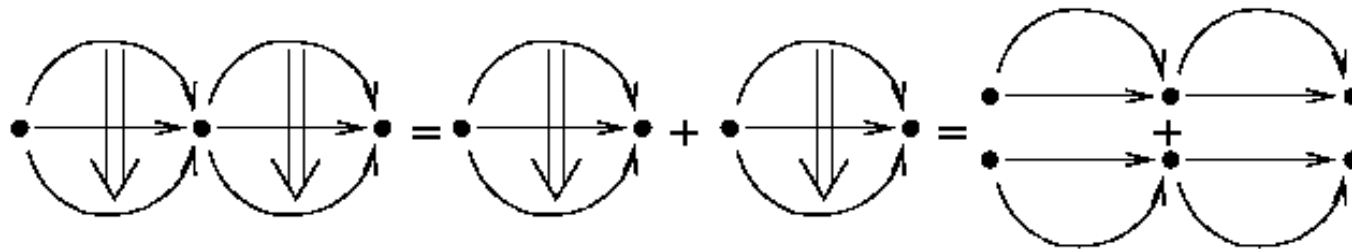
This functor requires to define an other functor for the DU algebra:
 $DU_c = DU_a \times DU_b$, a product in the topologic space.

Proposition: This functor is the product of 7 4-simplices (pyramides)

NB: this requires using the ternary operator “?:” to be constructible.

PQuantity

The topology of PQuantity is a 2–category on a vector space.

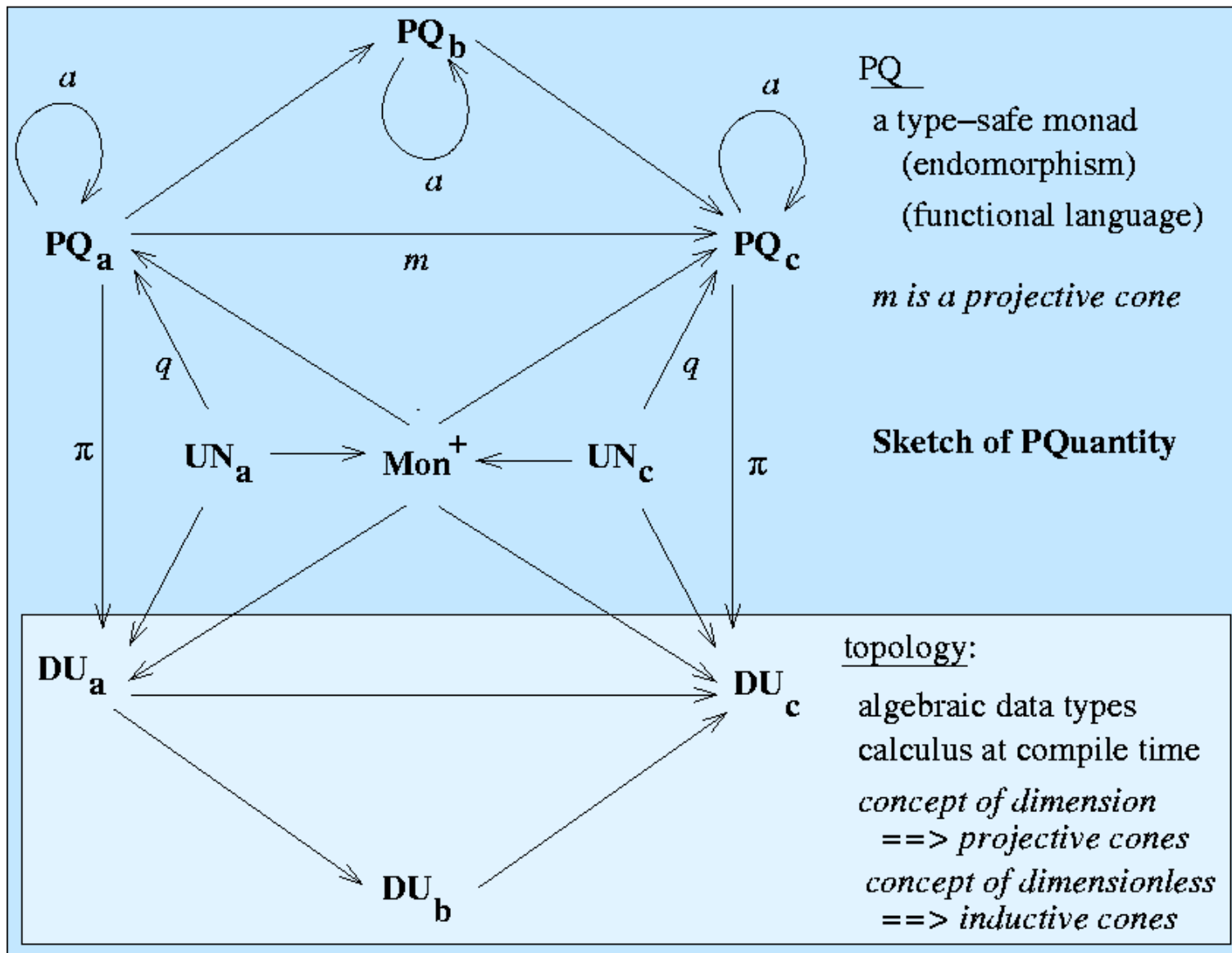


this vector space is of dimension 7:

- 0 Length
- 1 Mass
- 2 Time
- 3 Temperature
- 4 LuminousIntensity
- 5 MolarConcentration
- 6 ElectricCurrent

the horizontal composition along the fundamental physical unit basis
the vertical composition for the dimension,dimensionless property

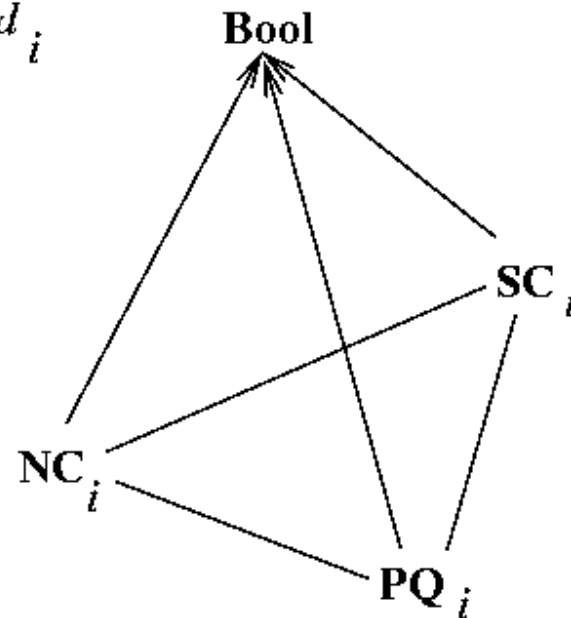
PQuantity is a monoid for the addition.



The Multiplier functor of PQuantity

An algebraic type with a closure: $SC_i / SC_i = Id_i$

Identity element: $Id = \bigoplus_{j \in J} Id_j$



A 3-simplex:

Space	Regions in the DSL
2D facette NC,PQ,Bool	sub-category of the dimensionned PQ
2D facette SC,PQ,Bool	sub-category of the dimensionless PQ
3D volume	category PQ: general case

Conclusions

1. The theory of the measurement set has been mostly developed
2. It is a data model 'above' the relational model
3. It encapsulates the relational model
4. Tables are sets containing a subset of their powersets, allow recursive definitions
5. Tables are topos
6. Tables are monads
7. The measurements are the object of a category
(*allows functional programming*)
8. The DM provides a rigorous DSL allowing expressive codes
9. The DM implies very robust codes (minimum of exceptions, type-safety i.e. validation at compile time).

10. The DM implies efficient code
11. The theory of categories provides a useful language for clean codes and should help when optimizing for parallel processing
12. Generic programming in C++ allows to express the mathematical formalism

Prospects

1. Consider implementing “Concepts” although not in C++11?
2. Deliver an implementation by the end of 2011 (ALMA CIPT)
3. Fully test the implementation with EMBRACE, edit profiles for other instruments
4. Produce a DM for the image domain using this approach
5. Think about an algebraic tree to work directly with these DMs
6. Generate queries in the context of DBMSs (test of the logic part)