

NUEVAS DIRECCIONES EN EL ACCESO E INTEGRACIÓN DE DATOS

Simone Santini

Escuela Politécnica Superior
Universidad Autónoma de Madrid

1975 (datos bancários)	2005 (datos <i>imaging</i>)
almacenados en ficheros, en formatos "ad-hoc"	almacenados en ficheros, en varios formatos estándar, directos a la visualización
para el acceso, se escriben procedimientos en FORTRAN o Cobol	para el acceso, se usan procedimientos en C, Java, Perl; si los datos se encuentran en un sitio web, las procedimientos se escriben como <i>servlets</i>
es necesario conocer no sólo la estructura lógica de los datos, sino también la manera en que están puestos en el disco	es necesario conocer no sólo la estructura lógica de los datos, sino también el formato de datos y acaso la estructura del sitio web donde se encuentren.

En el caso de los datos bancários: la situación cambia hacia 1980 con las bases de datos relacionales, que permiten un acceso a la estructura *lógica* de los datos independientemente de su estructura física.

```
select dept, avg(age)
from employees
where salary > 30.000
groupby dept
```

En el caso de datos iconográficos (*imaging*), todavía no existe un formalismo que permita una tratación tan general. El problema no es exclusivamente técnico: hay que analizar los asuntos a la base del modelo relacional:

- i) una base de datos es formada de datos que pertenecen a ciertos *tipos base* definidos *a priori*, agrupados en estructuras homogéneas y completamente expresadas;
- ii) tipos de datos más complejos (p.ej. objetos) son declarados explícitamente en sus componentes estructurales, y estas componentes son datos "macroscópicos" con un valor semántico definido;
- iii) cada dato tiene una representación única en el sistema;
- iv) los datos se agrupan en estructuras de un tipo definido (conjuntos de tuplas o, más recientemente, árboles XML).

Los datos iconográficos no respetan estos supuestos.

Ejemplo:

Volúmen: demasiado complejo para considerarse un tipo base y, además es necesario acceder a sus componentes en el curso de una interrogación.

No se puede considerar como una estructura (relacion) definida explícitamente. Imagínese una definición de relación:

```
Vol(X:int, Y:int, Z:int, value:float)
```

y una interrogación sobre el tamaño de la región conectada más grande en la cual una cierta condición sobre el valor está verificada, tal que el cálculo de la región se tenga que expresar en SQL:

```
select count(v)
from Vol
where ...
```

Datos como volúmenes se representan con características.

Definición de estructuras: muchos almacenes de datos interesantes no publican una estructura relacional. Ejemplos:

base de datos
orientadas a
objetos

data banks de
biología

ontologías

archivos de
citaciones

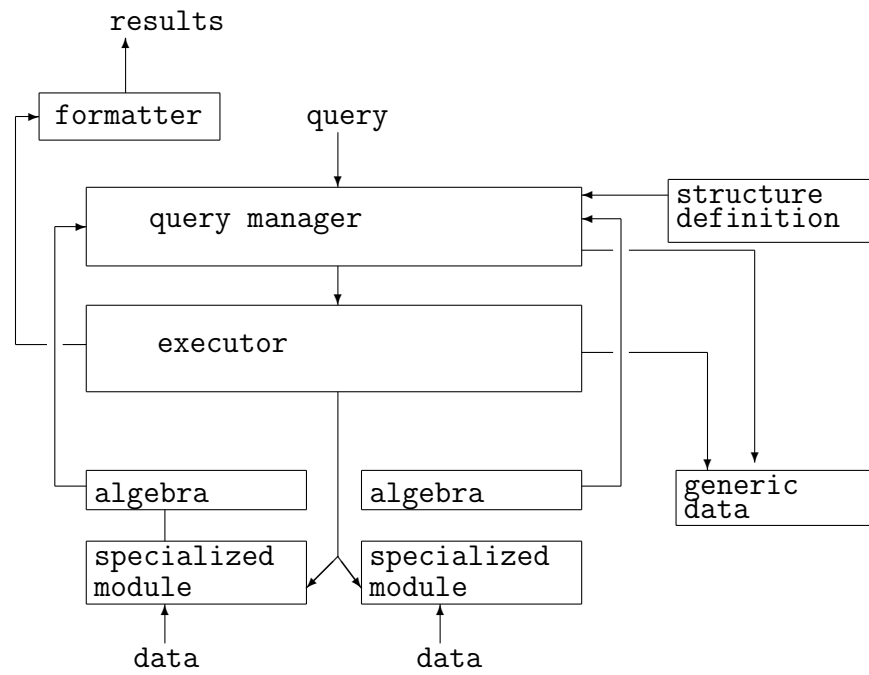
En estos casos se podrán utilizar en pleno todas las posibilidades que la fuente de datos ofrece sólo si el sistema de almacenamiento central o de integración de datos permite de expresar las mismas estructuras que se encuentran en estas fuentes.

Por ejemplo, la banca de datos *gene ontology* es un grafo acíclico de 23.000 nodos generados por las relaciones *is-a* y *part-of*.

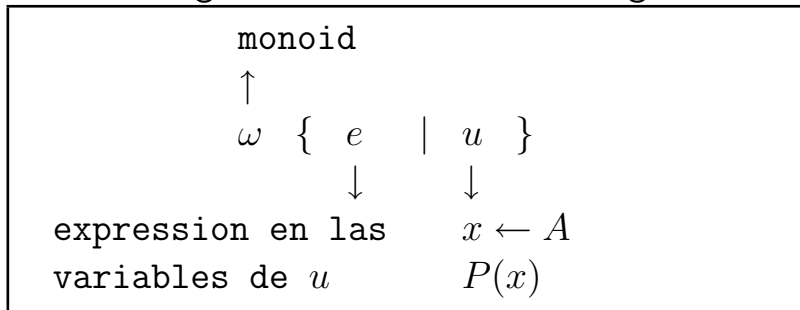
La "research intellectual property network" es una red de citaciones de patentes que es también estructurada como un grafo acíclico.

Los problemas generales que se deben enfrentar para definir una base de datos de imaging data:

- i) una integración en la base de datos de estructuras complejas y variadas que incluyan conjuntos, listas, árboles, y grafos;
- ii) una representación de tipos de datos abstractos, así que se pueda acceder a la estructura de datos iconográficos (como volúmenes) sin que sea necesario representarlos explícitamente en el sistema;
- iii) la posibilidad de usar representaciones diferentes para los datos que se almacenan en la base de datos.



forma general de una interrogación



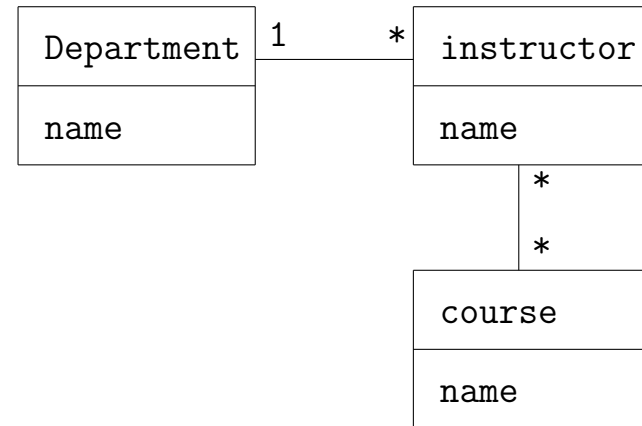
$\omega\{e x \leftarrow A\}$	<pre> r := null foreach x in A do r := r ω e(x) od </pre>	$\cup\{x^2 x \leftarrow [1, 2, 3]\}$	<pre> r := ∅ foreach x ∈ A do r := r ∪ {x^2} od </pre>	$\{1, 3, 9\}$
$\omega\{e x \leftarrow A, P(x)\}$	<pre> r := null foreach x in A do if P(x) then r := r ω e(x) fi od </pre>	$\vee\{\text{true} x \leftarrow [1, 2, 3], x > 2\}$	<pre> r := false foreach x ∈ A do if x > 2 then r := r ∨ true fi od </pre>	true

Esta sintaxis puede expresar toda la algebra relacional (SQL), y las más comunes algebras para bases de datos orientadas a objetos (OQL). Y más...

```
class Department {
  name : string;
  instructors : set(instructor);
}
```

```
class instructor {
  name : string;
  teaches : set(course);
}
```

```
class course {
  name : string;
}
```



$$\cup\{e.name \mid el \leftarrow \cup\{d.instructors \mid d \leftarrow Depts, d.name = CSE\}, \\ e \leftarrow el, \forall\{c.name = cse101 \mid c \leftarrow e.teaches\}\}$$

Busca todos los profesores del departamento CSE que enseñan cse101---se note que la relación entre profesores (instructors) y cursos es "many to many".

La primera interrogación saca el conjunto de todos los profesores de CSE que, de echo, es un conjunto de conjuntos. La expresión $e \leftarrow el$ hace un *unnesting* de estos conjuntos. La expresión final busca, entre todos estos profesores, los que enseñan cse101.

Para insertar grafos, se usan expresiones tipo patrones, de la misma familia de los patrones que se usan en XPath para hacer interrogaciones en árboles.

Las expresiones que yo uso permiten no solo la definición de patrones que partes de grafos pueden satisfacer, sino aún la definición de *variables* a las cuales partes del grafo que encaja con el patron se pueden asignar.

$$g \vdash \pi(x_1 : \pi_1, \dots, x_n : \pi_n)$$

Esta expresión genera un conjunto de tuplas $\langle x_1, \dots, x_n \rangle$, una para cada combinación de variables correspondientes a una parte del grafo que encaja con el patrón.

grafos: $\Gamma(\nu)$
 nodos: $\nu = (\text{id} : \text{skolem}, v : \text{int})$

patron: $\omega\{e|x \vdash \pi, x \leftarrow G\}$

$G : \{\Gamma\alpha\}, G : \{\Gamma\alpha\}, \dots$
 $x : \Gamma(\alpha).$

ejemplos: $\cup\{\{x\}|x \leftarrow G, x \vdash (v = 3)\}$ grafos con un nodo tal que $v = 3$;
 $\cup\{n|x \leftarrow G, x \vdash n:(v = 3)\}$ nodos con $v = 3$
 nota: $n : \{\nu\}.$

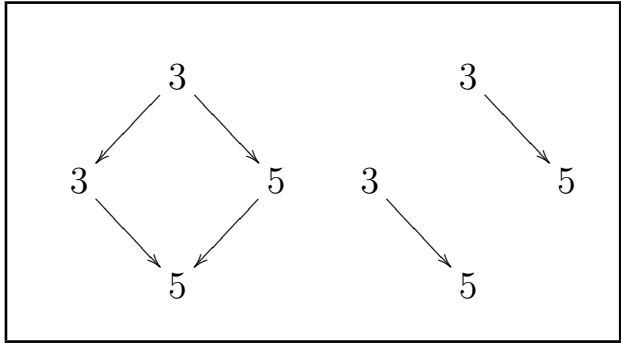
satisfacción de un patron: $g \models \pi$

Nota: para simplificar la notación, desde este momento escribiré las condiciones del tipo $(v = 3)$ como 3.

Composición serial:

$g_1 \models \pi_1, g_2 \models \pi_2 \Rightarrow g \models \pi_1 \text{---} \pi_2$: g es el "sink to source" de g_1, g_2 .

$\uplus\{n|x \leftarrow G, x \vdash n:(3 \text{---} 5)\}$ \Longrightarrow



$\uplus\{(n, m)|x \leftarrow G, x \vdash (n:3 \text{---} m:5)\}$ \Longrightarrow

$\{(3, 5), (3, 5)\}$

repetición de patrones:

$$\pi_1[-\pi_2](n) \equiv \pi_1 \overbrace{-\pi_2-\pi_2-\dots-\pi_2}^n$$

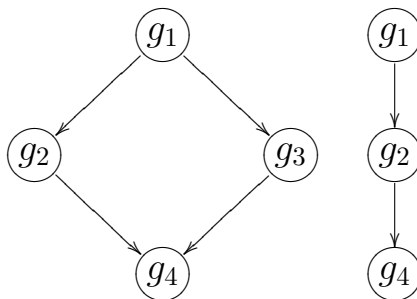
$$\pi_1[-\pi_2](n, m)$$

$$\pi_1[-\pi_2]^* \equiv \pi_1[-\pi_2](0, \infty)$$

$$\pi_1[-\pi_2]^+ \equiv \pi_1[-\pi_2](1, \infty)$$

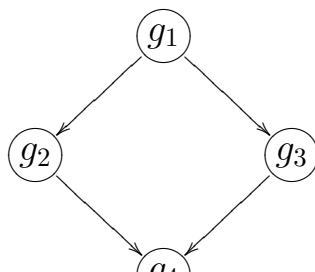
composición paralela:

$$\pi_1\{-\pi_2, -\pi_3\}-\pi_4$$



este patrón no requiere que los dos caminos sean distintos.

$$\pi_1! \{-\pi_2, -\pi_3\} - \pi_4$$



este si!

$\omega \{ e \mid u \}$
↓ ↓

expressions para
construir
resultados

Todas la variables
definidas en u pueden
ser mencionadas en e

Types: ν : nodos (*nodes*);
 $\Pi(\nu)$: caminos (*paths*);
 $\Gamma(\nu)$: grafos (*graphs*)

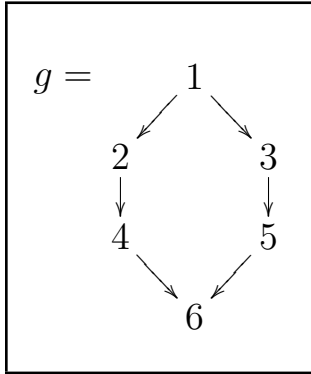
$\text{dom}(\alpha) \subseteq \text{dom}(\Pi(\alpha)) \subseteq \text{dom}(\Gamma(\alpha))$

gmin, gmax : monoid min y max para grafos:

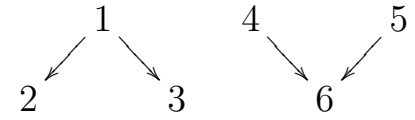
$$g_1 \leq g_2 \Leftrightarrow N(g_1) \subseteq N(g_2) \wedge E(g_1) \subseteq E(g_2)$$

$a, b, : \Gamma(\alpha)$: $a \text{---} b$ conecta b a a ;

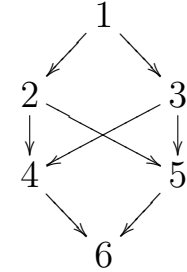
"merge": une dos grafos identificando los nodos con valores iguales.



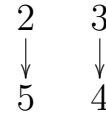
$$\cup \{(p, q) \mid g \vdash p : (1\{ \text{---} 2, \text{---} 3\}), \\ g \vdash q : (\{4 \text{---}, 5 \text{---}\}6) \\ \}$$



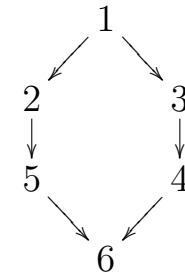
$$\cup \{p \text{---} q \mid g \vdash p : (1\{ \text{---} 2, \text{---} 3\}), \\ g \vdash q : (\{4 \text{---}, 5 \text{---}\}6) \\ \}$$



$$\cup \{(r \text{---} s, t \text{---} w) \mid g \vdash p : (1\{ \text{---} r : 2, \text{---} t : 3\}), \\ g \vdash q : (\{w : 4 \text{---}, s : 5 \text{---}\}6) \\ \}$$



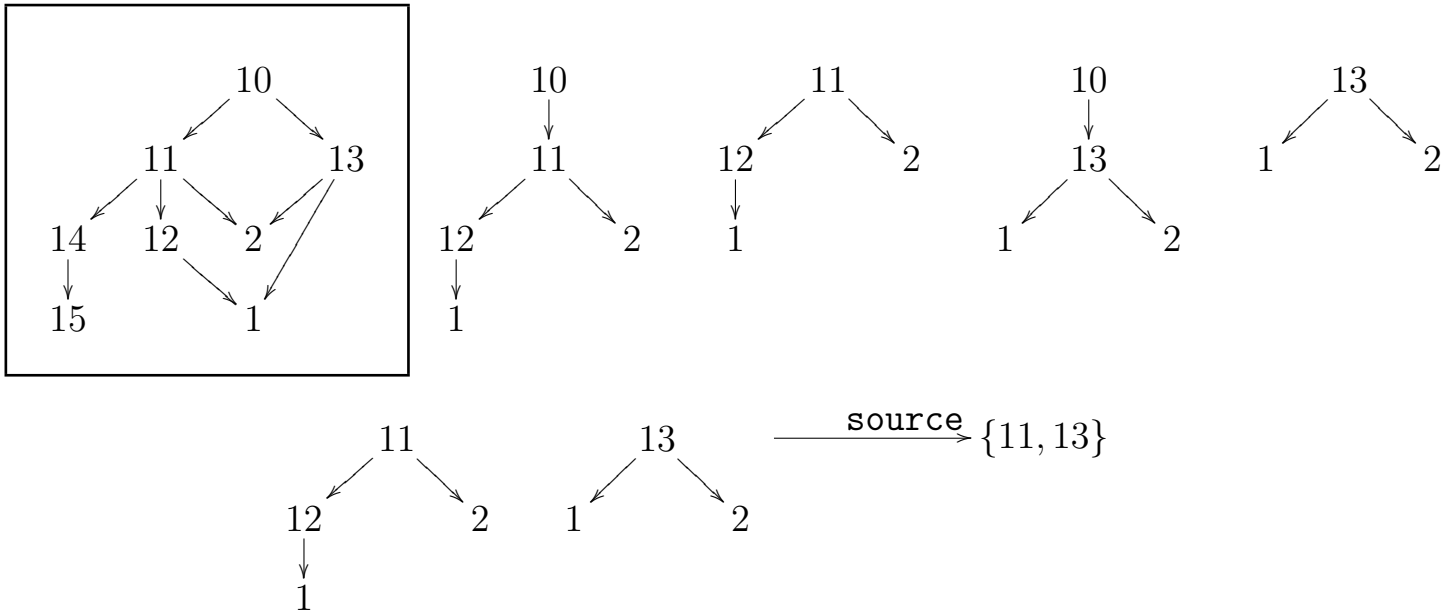
$$\cup \{\text{merge}\{p, q, r, s, t, w\} \mid g \vdash p : (1\{ \text{---} r : 2, \text{---} t : 3\}), \\ g \vdash q : (\{w : 4 \text{---}, s : 5 \text{---}\}6) \\ \}$$



$\text{merge}\{p, q, r, s, t, w\} = \text{merge}(p, \text{merge}(q, \text{merge}(r, \text{merge}(s, \text{merge}(t, w))))))$

$$\cup \{p \mid g \vdash p : (1\{ \text{---} 2 \text{---} 5, \text{---} 3 \text{---} 4\} \text{---} 6)\}$$

Antepasado más proximo (*least common ancestor*) de los nodos con id=1 e id=2
 ($DB = \{\Gamma(\alpha)\}$)

$$\cup \{ \text{source}(x) \mid x \leftarrow \text{gmin} \{ a \mid \begin{array}{l} x \vdash a : ((\text{true})\{[\neg(\text{true})] * \neg(id = 1), [\neg(\text{true})] * \neg(id = 2)\}), \\ x \leftarrow DB \\ \} \} \end{array} \}$$


Unos ejemplos de una red de patentes:

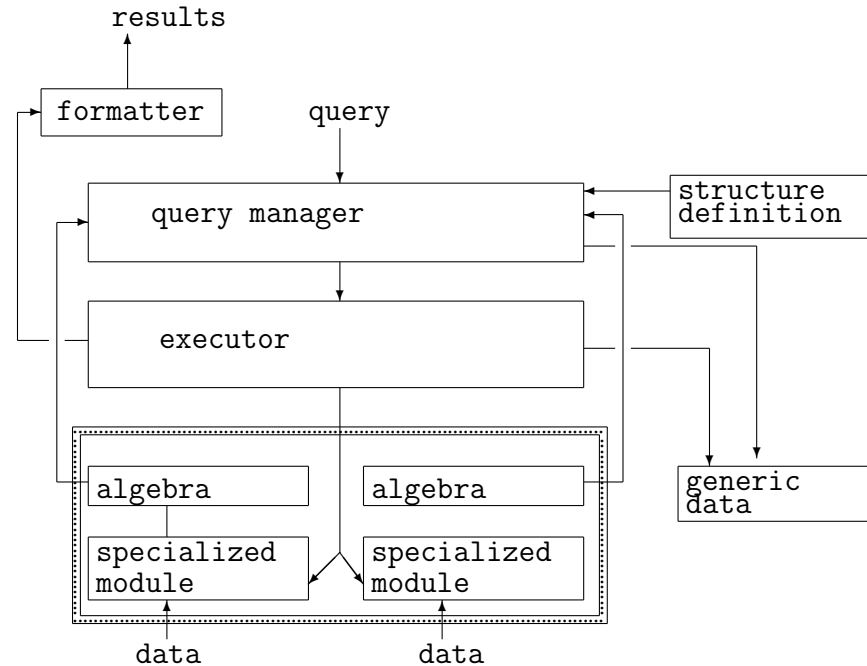
find all patents that can be traced to patent X

$$\cup\{a|g \leftarrow DB, g \vdash a : t[\text{---}t]^* \text{---}(name = X)\}$$

"Find all patent that have two distinct citation paths to X"

$$\cup\{a|g \leftarrow DB, g \vdash a : t!\{[\text{---}t]^*, [\text{---}t]^*\} \text{---}(name = X)\}$$

Tipos de datos complejos: la estructura de los datos iconográficos (imaging) es demasiado fina, y sus elementos no tienen bastante valor semántico para que se puedan almacenar directamente. En este caso, el sistema que se está desarrollando prevee su definición como *tipos de datos abstractos*, definidos por sus *algebras*.



En el caso específico, la idea es de proceder a una definición completa y general de algebras de características (*feature algebras*) para traducir interrogaciones en la *algebra interfaz* del tipo de datos a un *query plan* expresado en las algebras de características.

Ejemplo: definición de datos de imágenes.

Abstract data type:

```
type: image, region, color;  
constraint image < region;  
function: sim : image, image -> float;  
           segment : image -> {region};  
           area : region -> float;  
           carea : image, color -> float;  
           creg : image, region, color -> float;  
  
end
```

Se asuma que las imagenes son representadas con dos características: un histograma de color y un conjunto de regiones, cada una con su histograma de color. En el sistema se añade la siguiente definición

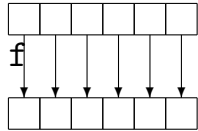
```
feature: hist, sphist;  
representation: h : image -> hist;  
                 sh : image : sphist;  
functions: hval : hist, color -> float;  
            sval : sphist, region, color -> float;  
equations: carea(img, col) == hval(h(img), col);  
            carea(img, col) == sum(r in segment(img))  
                               (sval(sh(img), r, color));  
            creg(img, r, color) = sval(sh(img), r, color);  
  
end
```

Para permitir este tipo de traducción entre tipos de datos y características, se está procediendo a una definición y clasificación general de los tipos de características de imágenes, y a una definición de algebras adecuadas para los varios tipos de imágenes.

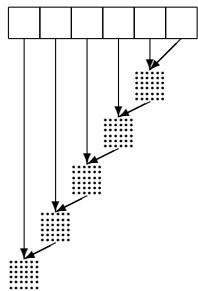
Para dar un ejemplo, se han empezado a estudiar los datos organizados en array. Se consideran sólo array unidimensionales, ya que un array de dimensionalidad n se puede considerar como un array unidimensional de arrays de dimensionalidad $n-1$.

La algebra de array solo consiste en dos operadores:

i) la aplicación de función m : $m.f.[u_1, \dots, u_n] = [f(u_1), \dots, f(u_n)]$;

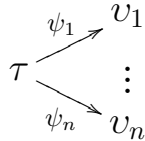


ii) la reducción r : $m.f.[u_1, \dots, u_n] = f(u_1, f(u_2, \dots, f(u_n, 0) \dots))$;



Basandose en estos operadores, y en funciones f especializadas, se ha definido, por ejemplo, una algebra de histograma, aplicándola a histogramas de color, de direcciones de *edges*, etc.

características:



cada característica permite una implementación eficiente de una parte del algebra del tipo.

Constants:

identity: $\iota : \alpha \rightarrow \alpha : x \mapsto x$;

projections: $\pi_i : \alpha_1 \times \dots \times \alpha_n \rightarrow \alpha_i : (x_1, \dots, x_n) \mapsto x_i$;

diagonal: $\Delta : \alpha \rightarrow \alpha \times \alpha : x \mapsto (x, x)$.

Algebra de funciones:

evaluación:

$$\frac{f : \alpha \rightarrow \beta \quad a : \alpha}{f.a : \beta}$$

composición:

$$\frac{f : \delta \rightarrow \beta \quad g : \alpha \rightarrow \delta}{f \circ g : \alpha \rightarrow \beta}$$

producto Cartesiano:

$$\frac{f : \alpha \rightarrow \beta \quad g : \delta \rightarrow \gamma}{f \times g : \alpha \times \delta \rightarrow \beta \times \gamma}$$

homomorphismos:

$$\frac{f : \alpha \rightarrow \nu(\beta)}{[\tau, \nu](f) : \tau(\alpha) \rightarrow \nu(\beta)}$$

$$[\tau, \nu] : (\alpha \rightarrow \nu(\beta)) \times \tau(\alpha) \rightarrow \nu(\beta);$$

condicional:

$$\frac{f : \alpha \rightarrow \beta \quad g : \alpha \rightarrow \beta \quad P : \alpha \rightarrow \mathbf{2}}{\text{if}(P, f, g) : \alpha \rightarrow \beta}$$

El futuro...

Incluir sistemas de base de datos en sistemas de flujo de datos: en el caso de tipo de datos imaging, muy a menudo el acceso a una base de datos se efectua como parte de un proceso mas complejo que incluye el procesamiento de los datos que se buscan. Estos procesos muchas veces se expresan como flujos de datos, y la integración de bases de datos en sistemas de flujo de datos ofrece muchas oportunidades de optimización.

Ejemplo: base de datos de pacientes con (entre otros) la edad del paciente y un volumen MRI del cerebro.

