CATEGORICAL CALCULUS AND ALGEBRA FOR MULTI-MODEL DATA

Jiaheng Lu Department of Computer Science University of Helsinki, Flmand



- Multi-model Databases
- Categorical Algebra and Calculus
- Algebraic Transformation Rules
- Conclusion
- Note: this talk will involve more in database knowledge, but only use basic knowledge in category theory, such as limit and thin category.



• Data come from different sources and have different formats







Smart phone

Camera

Social media

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI Acknowledgement: Icons created by Freepik - Flaticon





An example of different data and query in databases

Social network



Table

Customer_ID	Name	Credit_limits
1	Mary	5,000
2	John	3,000
3	William	2,000

Persons made the order:

"1"-->"34e5e759" "2"-->"0c6df508"

Order information:



An example of different data and query





One application with different models of data

Relational data: customer databases

- •Graph data: social networks
- •Hierarchical data: catalog, product
- •Key-value data: orders by customers

How to integrate those heterogenous data to provide a unified service?



Multi-Model Databases System

• One unified database system for multi-model data





What is DBMS?

 A Database Management System (DBMS) is software designed to efficiently manage data, with traditional systems storing data in the form of tables (RDBMS).

Student ID	First name	Last name	Department
001	John	Smith	Biology
002	Emily	Johnson	Physics
003	Michael	Brown	History
004	Sarah	Davis	English

Students Relational Table

What is multi-model database management system

• A multi-model database management system (MMDBMS) is designed to support multiple data models against a single, integrated backend.

• Document, graph, relational and key-value models are examples of data models that may be supported by a multi-model database.



Four models of data

Graph:



Relation:

Customer_ID	Name	Credit_limits
1	Mary	5,000
2	John	3,000
3	William	2,000

Key-value:

"1"-->"34e5e759" "2"-->"0c6df508"

Document:



Advantages of MMDBMS over the traditional relational database

- Handling diverse data types
 - Handle various types of data, such as graph, relation, document and key-value data and more models
- Enhanced query capabilities
 - Support content-based search for multi-model data or spatial queries for geospatial data.
- Flexible schema
 - Greater flexibility in schema design and evolution. Relational DBMS has the fixed database schema definitions.



Multi-model databases products







• Supporting graph, document, key/value and object models.





Query: Return all products which are ordered by a friend of a customer whose credit_limit>4000 Answer: John is a friend of Mary (the credit_limit of Mary > 4000)







Query language of OrientDB:

SELECT

expand(out("Knows").Orders.orderlines.Product_no)
FROM Customers
WHERE CreditLimit > 4000

Recommendation query: Return all products which are ordered by any friend of a customer whose credit_limit>4000



Challenge: a new theory foundation

Research goal: Call for a unified model and theory for multi-model data!

The theory of traditional relations is not adequate to mathematically describe modern database systems.

One possible theory foundation: Category Theory

- Introduced to mathematics world by Samuel Eilenberg and Sauders MacLane in 1944
- Developed for a unified language of topology and algebra





HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI Samuel Eilenberg

Sauders MacLane



- Databases are set categories:
 - Objects are sets and morphisms are functions
- We assume that it is a thin Category (or Posetal Category)
 - Given a pair of objects X and Y in a category C, and any two morphisms f, g: $X \rightarrow Y$, we say that C is a thin category if and only if the morphisms f and g are equal.



An example of Categorical Unification



Relational algebra and relational calculus

 In the field of relational databases, relational algebra and relational calculus are developed as two formal languages for query databases.

 Similarly, categorical algebra and categorical calculus are developed to query category databases.



Relational algebra

- Operators:
 - Selection: σ (sigma)
 - Projection: Π
 - Union: \cup
 - Intersection : \cap
 - Difference: -
 - Cartesian Product: \times
- Derived operators:
 - Joins (equi-join) 🖂



HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



Categorical algebra

- Set operators:
 - Unary operator:
 - Map: f
 - Selection: $\boldsymbol{\sigma}$
 - Projection: Π
 - Binary operator:
 - Division: ÷
 - getParent(D₁,D₂) (tree data)
 - getAncestor(D₁,D₂) (tree data)
 - Tenary operator:
 - getReach(S,T,E) (graph data)
 - getNHop(S,T,E) (graph data)

Categorical algebra

- Category operators:
 - Sets and Functions to Category:
 - Cat(S₁,...,S_n, f₁ : S_{i1} \rightarrow S_{j1} ,..., f_m : S_{im} \rightarrow S_{jm})
 - This operator, called **Categorification**, constructs a category using a given set of objects and morphisms.
 - Category to Set
 - Limit which converts a category into a relational object (set).
 - Lim(Cat(S₁,...,S_n, f₁ : S_{i1} \rightarrow S_{j1} ,..., f_m : S_{im} \rightarrow S_{jm}))



Example of categorical algebra: Selection



Query: Find all courses taken by "Smith" S1:= $\sigma_{\text{student} \cdot \text{Last_name}=\text{"Smith"}}$ (Registration) S2:= S1 · Course · Title Return S2

Example of categorical algebra: Division





Two examples of query plan with categorical algebra



Figure 4: Two holistic query plans involving three types of data



- Categorical calculus, a declarative language for describing results in the category; Categorical algebra, a procedural language for listing operations in the category.
- The formulae of the Categorical Calculus

Formulae with range terms	Safe variables
$x_1 \in O_1$	<i>x</i> ₁
$x_1 \in O_1 \land \neg (x_1 \in O_2)$	<i>x</i> ₁
Formulae with function and range terms	Safe variables
$((f_1: x_1 \to x_2) = f_2 \circ g_1) \land (x_1 \in S_1) \land (x_2 \in S_2)$	x_1, x_2
$(\pi_1:(x_1,x_2)\to x_1)\wedge(x_1\in S_1)\wedge(x_2\in S_2)$	x_1, x_2
Formulae with predicate, range and function terms	Data model
$(x_1 \in S_1) \land (x_2 \in S_2) \land (x_1 \leadsto^E x_2) \land (x_1 \cdot \text{Name} = "\text{John"})$	Graph
$(x_1 \in D_1) \land (x_2 \in D_2) \land (x_1 \text{ isAncestor } x_2)$	Tree
Formulae with unsafe terms	Unsafe variable
$x_2 \in S_2, x_3 \in S_3, \exists x_1(x_1 > x_3 \land x_2 = 6)$	x_1
$\forall x_1 \exists x_2 \in S_2(x_1 > x_2)$	x_1
$(x_1 \in S_1) \lor f(x_1) = a_1$	<i>x</i> ₁



• The formulae of the Categorical Calculus

Reachable predicate from node x_1 to x_2

Formulae with range terms	Safe variables
$x_1 \in O_1$	x_1
$x_1 \in O_1 \land \neg (x_1 \in O_2)$	x_1
Formulae with function and range terms	Safe variables
$((f_1: x_1 \to x_2) = f_2 \circ g_1) \land (x_1 \in S_1) \land (x_2 \in S_2)$	x_1, x_2
$(\pi_1:(x_1,x_2)\to x_1)\wedge(x_1\in S_1)\wedge(x_2\in S_2)$	x_1, x_2
Formulae with predicate, range and function terms	Data model
$(x_1 \in S_1) \land (x_2 \in S_2) \land (x_1 \rightsquigarrow^E x_2) \land (x_1 \cdot \text{Name} = "\text{John"})$	Graph
$(x_1 \in D_1) \land (x_2 \in D_2) \land (x_1 \text{ isAncestor } x_2)$	Tree
Formulae with unsafe terms	Unsafe variable
$x_2 \in S_2, x_3 \in S_3, \exists x_1(x_1 > x_3 \land x_2 = 6)$	x_1
$\forall x_1 \exists x_2 \in S_2(x_1 > x_2)$	x_1
$(x_1 \in S_1) \lor f(x_1) = a_1$	x_1



• The formulae of the Categorical Calculus

Ancestor predicate to determine the relationship between two nodes in trees.

Safe variables
<i>x</i> ₁
x_1
Safe variables
x_1, x_2
x_1, x_2
Data model
Graph
Tree
Unsafe variable
x_1
x_1
x_1



• The formulae of the Categorical Calculus

Formulae with range terms	Safe variables
$x_1 \in O_1$	<i>x</i> ₁
$x_1 \in O_1 \land \neg (x_1 \in O_2)$	<i>x</i> ₁
Formulae with function and range terms	Safe variables
$((f_1: x_1 \to x_2) = f_2 \circ g_1) \land (x_1 \in S_1) \land (x_2 \in S_2)$	x_1, x_2
$(\pi_1:(x_1,x_2)\to x_1)\wedge (x_1\in S_1)\wedge (x_2\in S_2)$	x_1, x_2
Formulae with predicate, range and function terms	Data model
$(x_1 \in S_1) \land (x_2 \in S_2) \land (x_1 \rightsquigarrow^E x_2) \land (x_1 \cdot \text{Name} = "\text{John"})$	Graph
$(x_1 \in D_1) \land (x_2 \in D_2) \land (x_1 \text{ isAncestor } x_2)$	Tree
Formulae with unsafe terms	Unsafe variable
$x_2 \in S_2, x_3 \in S_3, \exists x_1(x_1 > x_3 \land x_2 = 6)$	
$\forall x_1 \exists x_2 \in S_2(x_1 > x_2)$	
$(x_1 \in S_1) \lor f(x_1) = a_1$	<i>x</i> ₁

Unsafe

variables refer to a variable that has possibly infinite number of values or is unbounded.



Categorical calculus and categorical algebra are equivalent (I)





Categorical calculus and categorical algebra are equivalent (II)



Query: Find the titles of courses taken by all students

S1:= Registration[Student] ÷ Student

 $S2:=S1 \cdot Course \cdot Title$

Equivalent calculus:

{ $x \mid x \in Title, \forall y \in Student, \exists r \in Registration, r \cdot Student = y \land r \cdot$ Course \cdot Title= x}



- Multi-model Databases
- Categorical Algebra and Calculus
- Algebraic Transformation Rules (Query optimization)
- Conclusion

Algebraic transformation rules (I)

- Rewrite the algebraic operators for query optimization.
- Limit and Projection:

 $\pi_{S_1}(Lim(Cat(S_1, S_2, f: S_1 \to S_2))) \equiv S_1$

 $\pi_{S_2}(Lim(Cat(S_1, S_2, f: S_1 \to S_2))) \subseteq S_2$

 \bullet Pushing σ to one or multiple objects in lim

 $\sigma_C(Lim(S_1,S_2)) \equiv Lim(\sigma_C(S_1),S_2)$

 $\sigma_C(Lim(S_1,S_2)) \equiv Lim(\sigma_{C_1}(S_1),\sigma_{C_2}(S_2))$

Algebraic transformation rules (II)

• Pushing σ to one or multiple objects in getReach:

 $getReach(\sigma_{C_1}(S), \sigma_{C_2}(T), E) \equiv \sigma_{C_1 \wedge C_2}(getReach(S, T, E))$

Commuting function mapping with the product operator.

 $(f \otimes g)(S_1 \times S_2) \equiv f(S_1) \times g(S_2)$

• Commuting π with the Lim operation.

 $\pi_L(Lim(Cat(R_1, R_2, f_1 : R_1 \to R_2))) \equiv Lim(Cat(\pi_{L_1}(R_1), \pi_{L_2}(R_2), f_2 : \pi_{L_1}(R_1) \to \pi_{L_2}(R_2)))$

Algebraic transformation rules (III)

- Commuting g with the Lim operation..
- The following diagram holds:

$$\begin{array}{ccc} S_1 & \stackrel{f_1}{\longrightarrow} & S_2 \\ g_1 & & & \downarrow g_2 \\ S'_1 & \stackrel{f_2}{\longrightarrow} & S'_2 \end{array}$$

then the two operators g and Lim can be commuted as follows:

 $g(Lim(Cat(S_1, S_2, f_1 : S_1 \to S_2))) \equiv Lim(Cat(g_1(S_1), g_2(S_2), f_2 : g_1(S_1) \to g_2(S_2)))$

Algebraic transformation rules (III)

- Commuting g with the Lim operation..
- The following diagram holds:

$$\begin{array}{ccc} S_1 & \stackrel{f_1}{\longrightarrow} & S_2 \\ g_1 & & & \downarrow g_2 \\ S'_1 & \stackrel{f_2}{\longrightarrow} & S'_2 \end{array}$$

then the two operators g and limit can be commuted as follows:

 $g(Lim(Cat(S_1, S_2, f_1 : S_1 \to S_2))) \equiv Lim(Cat(g_1(S_1), g_2(S_2), f_2 : g_1(S_1) \to g_2(S_2)))$



An optimization query plan with algebraic operators transformation



Figure 4: Two holistic query plans involving three types of data

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

 $\sigma_C(Lim(S_1,S_2)) \equiv Lim(\sigma_{C_1}(S_1),\sigma_{C_2}(S_2))$



- Theorem 13. Categorical calculus and categorical algebra can express all of the following:
 - Relational calculus and algebra queries;
 - Graph pattern matching and graph reachability queries;
 - XML twig pattern queries.



Theorem 14. Consider a category C with p objects and q morphisms. Let the maximum number of elements in any object of C be n.

1. The data time complexity of computations for categorical calculus and categorical algebra in \mathscr{C} is bounded by $O(q \cdot n^p)$.

2. The space complexity of these computations is bounded by NSPACE[log n].



- Previous works use category theory on relational databases, but this paper focuses on multi-model databases.
 - Libkin and Wong (1997) showcase the connection between database operations and the categorical notion of a monad.
 - Schultz and Spivak (2016) introduce a categorical query language that serves as a data integration scripting language
- There are existing algebra and calculus for relational data, graph data, and object-oriented data, but not multi-model data.

.



- Model multi-model databases as thin set category
- Define categorical algebra and calculus, and their equivalence
- Develop the algebraic transformation rules for query optimization

Applied category theory (ACT) here can contribute to **practical query processing and optimization** of multi-model databases.



- Jeremy Gibbons, Fritz Henglein, Ralf Hinze & Nicolas Wu (2018): Relational algebra by way of adjunctions. Proc. ACM Program. Lang. 2(ICFP), pp. 86:1–86:28.
- Leonid Libkin & Limsoon Wong (1997): Query Languages for Bags and Aggregate Functions. J. Comput. Syst. Sci. 55(2), pp. 241–272.
- Patrick Schultz, David I. Spivak, Christina Vasilakopoulou & Ryan Wisnesky (2016): Algebraic Databases. arXiv:1602.03501.
- Allen Van Gelder & Rodney W. Topor (1991): Safety and translation of relational calculus. ACM Trans Database Syst. 16(2), p. 235–278, doi:10.1145/114325.103712. Available at <u>https://doi.org/10.1145/114325.103712</u>.
- Jiaheng Lu, Irena Holubová: Multi-model Databases: A New Journey to Handle the Variety of Data. ACM Comput. Surv. 52(3): 55:1-55:38 (2019)
- Jiaheng Lu: A Categorical Unification for Multi-Model Data: Part I Categorical Model and Normal Forms. CoRR abs/2502.19131 (2025)