

CƠ SỞ DỮ LIỆU PHÂN TÁN

NGUYỄN THỊ THANH THỦY

1

Nội dung môn học

- Giới thiệu
- Thiết kế CSDL phân tán
- Điều khiển dữ liệu phân tán
- Xử lý truy vấn phân tán
- Xử lý giao dịch phân tán

- Nhân bản dữ liệu
- Tích hợp CSDL – Các hệ thống đa CSDL
- Các hệ thống CSDL song song
- Quản lý dữ liệu ngang hàng (Peer-to-Peer)
- Xử lý dữ liệu lớn
- NoSQL, NewSQL và Polystores
- Quản lý dữ liệu Web

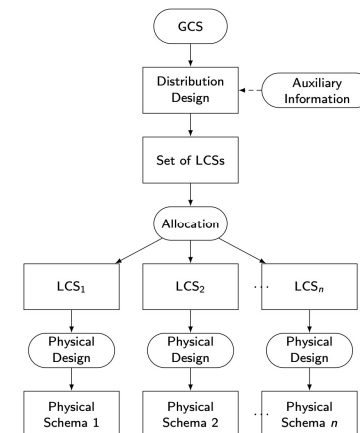
2

Nội dung

- Thiết kế CSDL phân tán
 - Phân mảnh
 - Phân tán dữ liệu
 - Các phương pháp kết hợp

3

Thiết kế phân tán



4

Nội dung

- Thiết kế CSDL phân tán
 - Phân mảnh
 - Phân tán dữ liệu
 - Các phương pháp kết hợp

5

Phân mảnh

- Có thể phân tán các quan hệ hay không?
 - Quan hệ không phải là một đơn vị phân tán thích hợp
 - Các khung nhìn (view) của ứng dụng thường là các tập con của các quan hệ
 - Nếu khung nhìn được định nghĩa dựa trên các quan hệ cho trước tại các trạm khác nhau thì:
 - Hoặc quan hệ không được nhân bản và được lưu trữ tại một trạm: cần phải thực hiện rất nhiều các thủ tục truy nhập từ xa.
 - Hoặc quan hệ được nhân bản tại tất cả hoặc một số trạm có chạy ứng dụng: khó khăn khi thực thi cập nhật, lãng phí không gian lưu trữ.

6

Phân mảnh

- Đơn vị phân tán hợp lý là gì?
 - Các mảnh của các quan hệ (quan hệ con)
 - Thực hiện đồng thời một số giao dịch mà có thể truy nhập đến các phần khác nhau của một quan hệ
 - Có thể thực thi song song truy vấn đơn (thao tác trên các mảnh) → nội truy vấn đồng thời
 - Nhược điểm:
 - Nếu các khung nhìn không thể được xác định trên một mảnh đơn sẽ được yêu cầu xử lý thêm → giảm khả năng thực thi
 - Kiểm soát dữ liệu ngữ nghĩa (đặc biệt là thực thi tính toàn vẹn) khó khăn hơn do có thể phải thực hiện tìm kiếm dữ liệu tại nhiều trạm

7

CSDL ví dụ

EMP			ASG			
ENO	ENAME	TITLE	ENO	PNO	RESP	DUR
E1	J. Doe	Elect. Eng.	E1	P1	Manager	12
E2	M. Smith	Syst. Anal.	E2	P1	Analyst	24
E3	A. Lee	Mech. Eng.	E2	P2	Analyst	6
E4	J. Miller	Programmer	E3	P3	Consultant	10
E5	B. Casey	Syst. Anal.	E3	P4	Engineer	48
E6	L. Chu	Elect. Eng.	E4	P2	Programmer	18
E7	R. Davis	Mech. Eng.	E5	P2	Manager	24
E8	J. Jones	Syst. Anal.	E6	P4	Manager	48
			E7	P3	Engineer	36
			E8	P3	Manager	40

PROJ				PAY	
PNO	PNAME	BUDGET	LOC	TITLE	SAL
P1	Instrumentation	150000	Montreal	Elect. Eng.	40000
P2	Database Develop.	135000	New York	Syst. Anal.	34000
P3	CAD/CAM	250000	New York	Mech. Eng.	27000
P4	Maintenance	310000	Paris	Programmer	24000

8

Phân mảnh theo chiều ngang

PROJ₁: các dự án với ngân sách thấp hơn \$200,000

PROJ₂: các dự án với ngân sách lớn hơn hoặc bằng \$200,000

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris

9

Phân mảnh theo chiều ngang

PROJ₁: các dự án với ngân sách thấp hơn \$200,000

PROJ₂: các dự án với ngân sách lớn hơn hoặc bằng \$200,000

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris

PROJ₁

PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York

PROJ₂

PNO	PNAME	BUDGET	LOC
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris

10

Phân mảnh theo chiều dọc

PROJ₁: thông tin về ngân sách dự án

PROJ₂: thông tin về tên và vị trí của dự án

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris

11

Phân mảnh theo chiều dọc

PROJ₁: thông tin về ngân sách dự án

PROJ₂: thông tin về tên và vị trí của dự án

PROJ			
PNO	PNAME	BUDGET	LOC
P1	Instrumentation	150000	Montreal
P2	Database Develop.	135000	New York
P3	CAD/CAM	250000	New York
P4	Maintenance	310000	Paris

PROJ₁

PNO	BUDGET
P1	150000
P2	135000
P3	250000
P4	310000

PROJ₂

PNO	PNAME	LOC
P1	Instrumentation	Montreal
P2	Database Develop.	New York
P3	CAD/CAM	New York
P4	Maintenance	Paris

12

Tính đúng đắn của việc phân mảnh

- Tính đầy đủ
 - Việc phân rã quan hệ R thành các mảnh R_1, R_2, \dots, R_n được coi là hoàn thành khi và chỉ khi mỗi mục dữ liệu trong R đều có thể được tìm thấy trong một số R_i nào đó.
- Tính phục hồi
 - Nếu quan hệ R được phân rã thành các mảnh R_1, R_2, \dots, R_n , thì cần tồn tại một toán tử quan hệ nào đó ∇ sao cho:

$$R = \nabla_{1 \leq i \leq n} R_i$$
- Tính tách biệt
 - Nếu quan hệ R được phân rã thành các mảnh R_1, R_2, \dots, R_n , và mục dữ liệu d_i thuộc R_j , thì d_i không thuộc về bất kỳ mảnh R_k ($k \neq j$) nào khác.

13

Các phương pháp cấp phát (/định vị)

- Không nhân bản
 - Khi phân chia: Mỗi mảnh chỉ có tại một trạm
- Nhân bản
 - Nhân bản đầy đủ: mỗi mảnh có tại một trạm
 - Nhân bản một phần: mỗi mảnh có tại một số trạm
- Quy tắc:
 - Nếu số lượng truy vấn chỉ đọc lớn hơn nhiều số lượng truy vấn cập nhật, thì nhân bản là tốt
 - Ngược lại, nhân bản sẽ gây ra rất nhiều vấn đề.

14

So sánh các phương pháp nhân bản

	Nhân bản hoàn toàn	Nhân bản một phần	Phân mảnh
XỬ LÝ TRUY VẤN	Dễ	Cùng mức độ khó khăn	
QUẢN LÝ THU MỤC	Dễ hoặc không tồn tại	Cùng mức độ khó khăn	
ĐIỀU KHIỂN ĐỒNG THỜI	Vừa phải	Khó	Dễ
ĐỘ TIN CẬY	Rất cao	Cao	Thấp
TÍNH THỰC TẾ	Có thể áp dụng	Thực tế	Có thể áp dụng

15