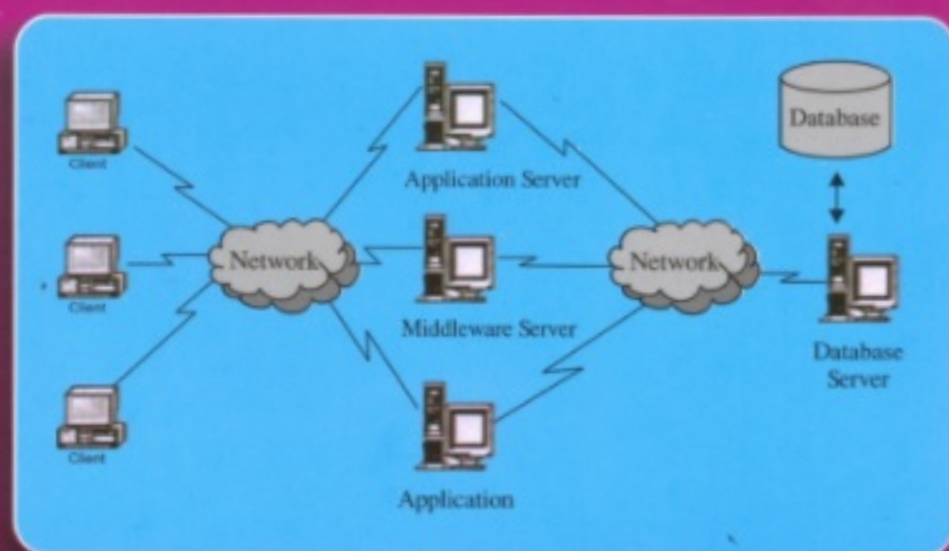


NGUYỄN VĂN HUÂN - PHẠM VIỆT BÌNH

# GIÁO TRÌNH HỆ CƠ SỞ DỮ LIỆU

## PHÂN TÁN & SUY DIỄN

Lý thuyết và thực hành



NHÀ XUẤT BẢN KHOA HỌC VÀ KỸ THUẬT

NGUYỄN VĂN HUÂN - PHẠM VIỆT BÌNH

**GIÁO TRÌNH**

**HỆ CƠ SỞ DỮ LIỆU  
PHÂN TÁN VÀ SUY DIỄN**

**NHÀ XUẤT BẢN KHOA HỌC VÀ KỸ THUẬT  
Hà Nội - 2009**



# MỤC LỤC

|   |    |
|---|----|
| <b>MỞ ĐẦU</b>   | 3  |
| <b>MỤC LỤC</b>  | 5  |
| <b>CHƯƠNG I. GIỚI THIỆU CHUNG VỀ CƠ SỞ DỮ LIỆU</b>              | 9  |
| 1.1. KHÁI NIỆM CƠ BẢN VỀ CÁC HỆ CƠ SỞ DỮ LIỆU                   | 9  |
| 1.1.1. Cơ sở dữ liệu là gì                                      | 9  |
| 1.1.2. Sự cần thiết của các hệ cơ sở dữ liệu                    | 10 |
| 1.1.3. Mô hình kiến trúc tổng quát cơ sở dữ liệu 3 mức          | 11 |
| 1.1.4. Mục tiêu của các hệ cơ sở dữ liệu                        | 14 |
| 1.1.5. Hệ quản trị CSDL & người quản trị CSDL                   | 15 |
| 1.1.6. Ràng buộc dữ liệu  | 17 |
| 1.1.7. Các mô hình truy xuất dữ liệu                            | 18 |
| 1.1.7.1. Mô hình cơ sở dữ liệu Client Server                    | 18 |
| 1.1.7.2. Mô hình Client/Server nhiều lớp                        | 19 |
| 1.1.7.3. Kỹ thuật lập trình cơ sở dữ liệu - Web động            | 20 |
| 1.1.7.4. Kiến trúc hệ thống Server (Server System Architecture) | 21 |
| 1.1.7.5. Các mô hình kiến trúc ứng dụng                         | 23 |
| 1.2. CÁC MÔ HÌNH CƠ SỞ DỮ LIỆU                                  | 23 |
| 1.2.1. Mở đầu   | 24 |
| 1.2.2. Mô hình dữ liệu (Data Model)                             | 24 |
| 1.2.2.1. Phân biệt giữa các mô hình dữ liệu                     | 25 |
| 1.2.2.2. Các hệ thống CSDL đối tượng và tri thức                | 25 |
| 1.2.3. Mô hình CSDL phân cấp (Hierarchy Data Model)             | 25 |
| 1.2.3.1. Cấu trúc biểu diễn dữ liệu phân cấp                    | 25 |
| 1.2.3.2. Ngôn ngữ thao tác trên CSDL phân cấp                   | 26 |
| 1.2.4. Mô hình CSDL mạng (Network Data Model)                   | 28 |
| 1.2.4.1. Cấu trúc biểu diễn dữ liệu mạng                        | 28 |
| 1.2.4.2. Ngôn ngữ dữ liệu thao tác trên CSDL mạng               | 30 |
| 1.2.5. Cách tiếp cận mô hình CSDL quan hệ                       | 30 |
| <b>CHƯƠNG II. CƠ SỞ DỮ LIỆU PHÂN TÁN</b>                        | 33 |
| 2.1. HỆ CƠ SỞ DỮ LIỆU PHÂN TÁN                                  | 33 |
| 2.1.1. Định nghĩa CSDL phân tán                                 | 33 |
| 2.1.2. Phân loại cơ sở dữ liệu phân tán                         | 35 |
| 2.1.3. Các đặc điểm chính của cơ sở dữ liệu phân tán            | 36 |
| 2.1.5. Xử lý dữ liệu phân tán                                   | 38 |
| 2.1.5. Ưu nhược điểm của việc sử dụng cơ sở dữ liệu phân tán    | 39 |
| 2.1.6. Cơ sở dữ liệu phân tán và cơ sở dữ liệu tập trung        | 40 |
| 2.1.7. Kiến trúc cơ bản của CSDL phân tán                       | 41 |
| 2.1.8. Hệ quản trị CSDL phân tán                                | 42 |
| 2.2. CÁC MÔ HÌNH XỬ LÝ PHÂN TÁN                                 | 45 |
| 2.2.1. Mô hình xử lý Master – Slave                             | 45 |
| 2.2.2. Các hệ khách/đại lý                                      | 45 |
| 2.2.3. Các hệ phân tán ngang hàng                               | 47 |
| 2.2.4. Môi trường đa tầng                                       | 47 |
| 2.3. THIẾT KẾ CƠ SỞ DỮ LIỆU PHÂN TÁN                            | 48 |
| 2.3.1. Các chiến lược thiết kế                                  | 48 |

|   |  |     |
|---|--|-----|
| 2.3.1.1.                                  | Quá trình thiết kế từ trên xuống (top-down)              | 48  |
| 2.3.1.2.                                  | Quá trình thiết kế từ dưới lên (bottom-up)               | 49  |
| 2.3.2.                                    | Các vấn đề thiết kế                                      | 50  |
| 2.3.2.1.                                  | Lý do phân mảnh  | 50  |
| 2.3.2.2.                                  | Các quy tắc phân mảnh đúng đắn                           | 50  |
| 2.3.2.3.                                  | Các yêu cầu thông tin                                    | 51  |
| 2.3.3.                                    | Phân mảnh ngang  | 51  |
| 2.3.3.1.                                  | Hai kiểu phân mảnh ngang                                 | 51  |
| 2.3.3.2.                                  | Yêu cầu thông tin của phân mảnh ngang                    | 51  |
| 2.3.3.3.                                  | Phân mảnh ngang nguyên thủy                              | 54  |
| 2.3.3.4.                                  | Phân mảnh ngang dẫn xuất                                 | 60  |
| 2.3.3.5.                                  | Kiểm định tính đúng đắn                                  | 62  |
| 2.3.4.                                    | Phân mảnh dọc  | 62  |
| 2.3.5.                                    | Phân mảnh hỗn hợp  | 71  |
| 2.3.6.                                    | Cấp phát   | 71  |
| 2.3.6.1.                                  | Bài toán cấp phát  | 72  |
| 2.3.6.2.                                  | Cách tiếp cận 1  | 72  |
| 2.3.6.3.                                  | Cách tiếp cận 2  | 75  |
| 2.4.                                      | XỬ LÝ VẤN TIN  | 80  |
| 2.4.1.                                    | Bài toán xử lý vấn tin                                   | 81  |
| 2.4.2.                                    | Phân rã vấn tin  | 84  |
| 2.4.3.                                    | Cục bộ hóa dữ liệu phân tán                              | 91  |
| 2.4.4.                                    | Tối ưu hoá vấn tin phân tán                              | 101 |
| 2.5.                                      | QUẢN LÝ GIAO DỊCH  | 104 |
| 2.5.1.                                    | Giao dịch (Transaction)                                  | 104 |
| 2.5.2.                                    | Giao dịch phân tán                                       | 109 |
| 2.5.3.                                    | Tính khả tuần tự của các lịch biểu và việc sử dụng chung | 115 |
| 2.5.4.                                    | Các kỹ thuật điều khiển tương tranh bằng khóa            | 116 |
| 2.5.4.1.                                  | Mô hình khóa cơ bản                                      | 119 |
| 2.5.4.2.                                  | Mô hình khóa đọc và khóa ghi                             | 121 |
| 2.5.4.3.                                  | Thuật toán điều khiển tương tranh bằng nhãn thời gian    | 123 |
| <b>CHƯƠNG III. CƠ SỞ DỮ LIỆU SUY DIỄN</b> |  | 129 |
| 3.1.                                      | GIỚI THIỆU CHUNG   | 129 |
| 3.2.                                      | CƠ SỞ DỮ LIỆU SUY DIỄN                                   | 129 |
| 3.2.1.                                    | Mô hình cơ sở dữ liệu suy diễn                           | 129 |
| 3.2.2.                                    | Lý thuyết mô hình đối với cơ sở dữ liệu quan hệ          | 131 |
| 3.2.2.1.                                  | Nhìn nhận cơ sở dữ liệu theo quan điểm logic             | 131 |
| 3.2.2.2.                                  | Nhìn lại cơ sở dữ liệu quan hệ                           | 131 |
| 3.2.3.                                    | Nhìn nhận cơ sở dữ liệu suy diễn                         | 132 |
| 3.2.4.                                    | Các giao tác trên cơ sở dữ liệu suy diễn                 | 133 |
| 3.3.                                      | CƠ SỞ DỮ LIỆU DỰA TRÊN LOGIC                             | 133 |
| 3.3.1.                                    | Cú pháp  | 133 |
| 3.3.2.                                    | Ngữ nghĩa  | 134 |
| 3.3.3.                                    | Cấu trúc cơ bản  | 134 |
| 3.3.4.                                    | Cấu trúc của câu hỏi                                     | 137 |
| 3.3.5.                                    | So sánh DATALOG với đại số quan hệ                       | 138 |
| 3.3.6.                                    | Các hệ cơ sở dữ liệu chuyên gia                          | 142 |
| 3.4.                                      | MỘT SỐ VẤN ĐỀ KHÁC                                       | 142 |

|   |     |
|---|-----|
| <b>CHƯƠNG IV. CƠ SỞ DỮ LIỆU HƯỚNG ĐỐI TƯỢNG</b>   | 145 |
| 4.1. NGUYÊN TẮC CỦA CÁC MÔ HÌNH HƯỚNG ĐỐI TƯỢNG   | 145 |
| 4.1.1. Mô hình hóa các đối tượng  | 145 |
| 4.1.2. Phương pháp  | 146 |
| 4.1.3. Lớp (Class)  | 146 |
| 4.1.4. Các liên kết thừa kế giữa các lớp  | 147 |
| 4.1.5. Lược đồ lớp  | 148 |
| 4.2. TÍNH BỀN VỮNG CỦA CÁC ĐỐI TƯỢNG  | 148 |
| 4.2.1. Cơ sở dữ liệu hướng đối tượng  | 148 |
| 4.2.2. Quản lý tính bền vững  | 148 |
| <b>CHƯƠNG V. THỰC HÀNH MỘT SỐ ỨNG DỤNG</b>  | 151 |
| 5.1. THIẾT KẾ MỘT HỆ CƠ SỞ DỮ LIỆU KẾ TOÁN  | 151 |
| 5.1.1. Đặt vấn đề bài toán  | 151 |
| 5.1.2. Chiến lược   | 158 |
| 5.1.3. Phân tích  | 158 |
| 5.1.4. Thiết kế   | 171 |
| 5.1.5. Mô tả thiết kế hệ cơ sở dữ liệu phân tán cho hệ thống kế toán                            | 172 |
| 5.1.6. Sơ đồ phân cấp chức năng của hệ thống  | 176 |
| 5.2. THỰC HÀNH VỚI MỘT SỐ THUẬT TOÁN ĐIỀU KHIỂN TƯƠNG<br>TRANH TRONG QUẢN LÝ GIAO DỊCH PHÂN TÁN | 178 |
| <b>MỘT SỐ ĐỀ ĐÃ THI QUA CÁC NĂM</b>   | 189 |
| <b>TÀI LIỆU THAM KHẢO</b>   | 205 |



# CHƯƠNG I

## GIỚI THIỆU CHUNG VỀ CƠ SỞ DỮ LIỆU

---

### 1.1. KHÁI NIỆM CƠ BẢN VỀ CÁC HỆ CƠ SỞ DỮ LIỆU

Trong chương này trình bày những khái niệm cơ bản về các hệ cơ sở dữ liệu do E.F Codd đề xuất. Những khái niệm này bao gồm mục tiêu của một hệ cơ sở dữ liệu. Sự cần thiết phải tổ chức dữ liệu dưới dạng cơ sở dữ liệu. Tính độc lập của dữ liệu thể hiện mô hình kiến trúc 3 mức. Vì vậy có thể nói cơ sở dữ liệu phản ánh tính trung thực, khách quan của thế giới dữ liệu. Không dư thừa thông tin và cũng không thiếu thông tin. Nội dung của chương bao gồm các phần:

- Cơ sở dữ liệu là gì;
- Sự cần thiết của các hệ cơ sở dữ liệu;
- Mô hình kiến trúc 3 mức cơ sở dữ liệu;
- Mục tiêu của các hệ cơ sở dữ liệu;
- Hệ quản trị CSDL & người quản trị CSDL;
- Tổ chức lưu trữ dữ liệu;
- Các mô hình truy xuất.

#### 1.1.1. Cơ sở dữ liệu là gì

Cơ sở dữ liệu là một bộ sưu tập rất lớn về các loại dữ liệu tác nghiệp, bao gồm các loại dữ liệu âm thanh, tiếng nói, chữ viết, văn bản, đồ họa, hình ảnh tĩnh hay hình ảnh động... được mã hoá dưới dạng các chuỗi bit và được lưu trữ dưới dạng File dữ liệu trong các bộ nhớ của máy tính. Cấu trúc lưu trữ dữ liệu tuân theo các quy tắc dựa trên lý thuyết toán học. Cơ sở dữ liệu phản ánh trung thực thế giới dữ liệu hiện thực khách quan.

*Cơ sở dữ liệu là tài nguyên thông tin dùng chung cho nhiều người:* Cơ sở dữ liệu (CSDL) là tài nguyên thông tin chung cho nhiều người cùng sử dụng. Bất kỳ người sử dụng nào trên mạng máy tính, tại các thiết bị đầu cuối, về nguyên tắc có quyền truy nhập khai thác toàn bộ hay một phần dữ liệu theo chế độ trực tuyến hay tương tác mà không phụ thuộc vào vị trí địa lý của người sử dụng với các tài nguyên đó.

*Cơ sở dữ liệu được các hệ ứng dụng khai thác bằng ngôn ngữ con dữ liệu hoặc bằng các chương trình ứng dụng để xử lý, tìm kiếm, tra cứu, sửa đổi, bổ sung hay loại bỏ dữ liệu.* Tìm kiếm và tra cứu thông tin là một trong những chức năng quan trọng và phổ biến nhất của các dịch vụ cơ sở dữ liệu. Hệ quản trị CSDL - HQTCSDL (DataBase Management System - DBMS) là phần mềm điều khiển các chiến lược truy nhập CSDL. Khi người sử dụng đưa ra yêu cầu truy nhập bằng một ngôn ngữ con dữ liệu nào đó, HQTCSDL tiếp nhận và thực hiện các thao tác trên CSDL lưu trữ.



Đối tượng nghiên cứu của CSDL là các thực thể và mối quan hệ giữa các thực thể. Thực thể và mối quan hệ giữa các thực thể là hai đối tượng khác nhau về căn bản. Mối quan hệ giữa các thực thể cũng là một loại thực thể đặc biệt. Trong cách tiếp cận CSDL quan hệ, người ta dựa trên cơ sở lý thuyết đại số quan hệ để xây dựng các quan hệ chuẩn, khi kết nối không tổn thất thông tin và khi biểu diễn dữ liệu là duy nhất. Dữ liệu được lưu trữ trong bộ nhớ của máy tính không những phải tính đến yếu tố về tối ưu không gian lưu trữ, mà phải đảm bảo tính khách quan, trung thực của dữ liệu hiện thực. Nghĩa là phải đảm bảo tính nhất quán của dữ liệu và giữ được sự toàn vẹn của dữ liệu.

### 1.1.2. Sự cần thiết của các hệ cơ sở dữ liệu

Tổ chức lưu trữ dữ liệu theo lý thuyết cơ sở dữ liệu có những ưu điểm:

*Giảm bớt dư thừa dữ liệu trong lưu trữ:* Trong các ứng dụng lập trình truyền thống, phương pháp tổ chức lưu trữ dữ liệu vừa tốn kém, lãng phí bộ nhớ và các thiết bị lưu trữ, vừa dư thừa thông tin lưu trữ. Nhiều chương trình ứng dụng khác nhau cùng xử lý trên các dữ liệu như nhau, dẫn đến sự dư thừa đáng kể về dữ liệu. Ví dụ trong các bài toán nghiệp vụ quản lý "Cước thuê bao điện thoại" và "Doanh thu & sản lượng", tương ứng với mỗi một chương trình là một hay nhiều tệp dữ liệu được lưu trữ riêng biệt, độc lập với nhau. Trong cả 2 chương trình cùng xử lý một số thuộc tính của một cuộc đàm thoại như "số máy gọi đi", "số máy gọi đến", "hướng cuộc gọi", "thời gian bắt đầu" và "thời gian kết thúc"... Nhiều thuộc tính được mô tả và lưu trữ nhiều lần độc lập với nhau. Nếu tổ chức lưu trữ theo lý thuyết CSDL thì có thể hợp nhất các tệp lưu trữ của các bài toán trên, các chương trình ứng dụng có thể cùng chia sẻ tài nguyên trên cùng một hệ CSDL.

*Tổ chức lưu trữ dữ liệu theo lý thuyết CSDL sẽ tránh được sự không nhất quán trong lưu trữ dữ liệu và bảo đảm được tính toàn vẹn của dữ liệu:* Nếu một thuộc tính được mô tả trong nhiều tệp dữ liệu khác nhau và lặp lại nhiều lần trong các bản ghi, khi thực hiện việc cập nhật, sửa đổi, bổ sung sẽ không sửa hết nội dung các mục đó. Nếu dữ liệu càng nhiều thì sự sai sót khi cập nhật, bổ sung càng lớn. Khả năng xuất hiện mâu thuẫn, không nhất quán thông tin càng nhiều, dẫn đến không nhất quán dữ liệu trong lưu trữ. Tất yếu kéo theo sự dị thường thông tin, thừa, thiếu và mâu thuẫn thông tin.

Thông thường, trong một thực thể, giữa các thuộc tính có mối quan hệ ràng buộc lẫn nhau, tác động ảnh hưởng lẫn nhau. Cước của một cuộc đàm thoại phụ thuộc vào khoảng cách và thời gian cuộc gọi, tức là phụ thuộc hàm vào các thuộc tính máy gọi đi, máy gọi đến, thời gian bắt đầu và thời gian kết thúc cuộc gọi. Các trình ứng dụng khác nhau cùng xử lý cước đàm thoại trên các thực thể lưu trữ tương ứng khác nhau chưa hẳn cho cùng một kết quả về sản lượng phút và doanh thu. Điều này lý giải tại sao trong một doanh nghiệp, cùng xử lý trên các chỉ tiêu quản lý mà số liệu báo cáo của các phòng ban, các công ty con lại cho các kết quả khác nhau, thậm chí còn trái ngược nhau. Như vậy, có thể khẳng định, nếu dữ liệu không tổ chức theo lý thuyết cơ sở dữ liệu, tất yếu không thể phản ánh thế giới hiện thực dữ liệu, không phản ánh đúng bản chất vận động của dữ liệu.

*Sự không nhất quán dữ liệu trong lưu trữ làm cho dữ liệu mất đi tính toàn vẹn của nó.* Tính toàn vẹn dữ liệu đảm bảo cho sự lưu trữ dữ liệu luôn luôn đúng. Không thể có mã vùng ngoài quy định của cơ quan quản lý, hoặc ngày sinh của một nhân viên không thể xảy ra sau ngày tốt nghiệp ra trường của nhân viên đó...

*Tô chức lưu trữ dữ liệu theo lý thuyết CSDL có thể triển khai đồng thời nhiều ứng dụng trên cùng một CSDL:* Điều này có nghĩa là các ứng dụng không chỉ chia sẻ chung tài nguyên dữ liệu mà còn trên cùng một CSDL có thể triển khai đồng thời nhiều ứng dụng khác nhau tại các thiết bị đầu cuối khác nhau.

*Tô chức dữ liệu theo lý thuyết cơ sở dữ liệu sẽ thống nhất các tiêu chuẩn, thủ tục và các biện pháp bảo vệ, an toàn dữ liệu:* Các hệ CSDL sẽ được quản lý tập trung bởi một người hay một nhóm người quản trị CSDL, bằng các hệ quản trị CSDL. Người quản trị CSDL có thể áp dụng thống nhất các tiêu chuẩn, quy định, thủ tục chung như quy định thống nhất về mẫu biểu báo cáo, thời gian bổ sung, cập nhật dữ liệu. Điều này làm dễ dàng cho công việc bảo trì dữ liệu. Người quản trị CSDL có thể bảo đảm việc truy nhập tới CSDL, có thể kiểm tra, kiểm soát các quyền truy nhập của người sử dụng. Ngăn chặn các truy nhập trái phép, sai quy định từ trong ra hoặc từ ngoài vào...

### **1.1.3. Mô hình kiến trúc tổng quát cơ sở dữ liệu 3 mức**

Mô hình kiến trúc 3 mức của hệ CSDL gồm: *Mức trong, mức mô hình dữ liệu (Mức quan niệm) và mức ngoài.* Giữa các mức tồn tại các ánh xạ quan niệm trong và ánh xạ quan niệm ngoài. Trung tâm của hệ thống là mức quan niệm, tức là mức mô hình dữ liệu. Ngoài ra còn có khái niệm người sử dụng, hệ quản trị CSDL và người quản trị CSDL.

*Người sử dụng:* Là những người tại thiết bị đầu cuối truy nhập vào các hệ CSDL theo chế độ trực tuyến hay tương tác bằng các chương trình ứng dụng hay bằng các ngôn ngữ con dữ liệu. Thường là các chuyên viên kỹ thuật tin học, có trình độ thành thạo biết lập trình và biết sử dụng ngôn ngữ con thao tác dữ liệu (SQL Server, Oracle...). Người sử dụng có thể truy nhập toàn bộ hay một phần CSDL mà họ quan tâm, phụ thuộc vào quyền truy nhập của họ. Cách nhìn CSDL của người sử dụng nói chung là triu tượng. Họ nhìn CSDL bằng mô hình ngoài; gọi là mô hình con dữ liệu. Chẳng hạn người sử dụng là một nhân viên của phòng kế toán tài chính, chỉ nhìn thấy tập các xuất hiện kiểu bản ghi ngoài về doanh thu, sản lượng trong tháng, không thể nhìn thấy các xuất hiện kiểu bản ghi lưu trữ về các chỉ tiêu kỹ thuật của đường thông, mạng lưới...

*Mô hình ngoài:* Mô hình ngoài là nội dung thông tin của CSDL dưới cách nhìn của người sử dụng. Là nội dung thông tin của một phần dữ liệu tác nghiệp được một người hoặc một nhóm người sử dụng quan tâm. Nói cách khác, mô hình ngoài mô tả cách nhìn dữ liệu của người sử dụng và mỗi người sử dụng có cách nhìn dữ liệu khác nhau. Nhiều mô hình ngoài khác nhau có thể cùng tồn tại trong một hệ CSDL, nghĩa là có nhiều người sử dụng chia sẻ chung cùng một cơ sở dữ liệu. Hơn nữa, có thể mô hình ngoài quan hệ, mô hình ngoài phân cấp hay mô hình ngoài kiểu mạng cũng có thể tồn tại trong một cơ sở dữ liệu. Sơ đồ ngoài không làm "hiện" mà được nhúng vào trong logic một đơn tác có liên quan.

- Mô hình ngoài gồm nhiều xuất hiện kiểu bản ghi ngoài, nghĩa là mỗi một người sử dụng có một sơ đồ dữ liệu riêng, một khung nhìn dữ liệu riêng. Bản ghi ngoài của người sử dụng có thể khác với bản ghi lưu trữ và bản ghi quan niệm.

- Mô hình ngoài được xác định bởi một sơ đồ ngoài bao gồm các mô tả về kiểu bản ghi ngoài như tên các trường, kiểu dữ liệu các trường, độ rộng của trường...

- Ngôn ngữ con dữ liệu của người sử dụng thao tác trên các bản ghi ngoài.