

	VIETTEL AI RACE	TD069
	Thu Thập và Quản Lý Dữ Liệu Chuỗi Thời Gian	Lần ban hành: 1

Dữ liệu chuỗi thời gian (time series) xuất hiện trong hầu hết các lĩnh vực: từ tín hiệu cảm biến công nghiệp, log máy chủ, giao dịch tài chính cho đến dữ liệu khí tượng. Khác với dữ liệu dạng bảng thông thường, chuỗi thời gian đòi hỏi các phương pháp thu thập và quản lý đặc thù để đảm bảo tính liên tục, đồng bộ và chất lượng. Dưới đây là các phương pháp và nguyên tắc triển khai chi tiết, đủ để xây dựng một hệ thống bền vững và có khả năng mở rộng.

1. Xác định Nguồn Dữ Liệu và Yêu Cầu Kỹ Thuật

- Nguồn dữ liệu cảm biến công nghiệp: cảm biến nhiệt độ, áp suất, độ ẩm, rung động từ dây chuyền sản xuất.
- Nguồn dữ liệu số hóa khác: log hệ thống, dữ liệu giao dịch thời gian thực từ ngân hàng, sàn giao dịch.
- Yêu cầu quan trọng:
 - Độ chính xác của thời gian ghi nhận (timestamp) phải thống nhất múi giờ và đồng bộ với đồng hồ chuẩn (NTP).
 - Đảm bảo tần suất ghi nhận ổn định (1 Hz, 10 Hz hoặc cao hơn tùy ứng dụng).
 - Khả năng mở rộng khi số lượng cảm biến hoặc thiết bị tăng gấp nhiều lần.

	VIETTEL AI RACE	TD069
	Thu Thập và Quản Lý Dữ Liệu Chuỗi Thời Gian	Lần ban hành: 1

Các tổ chức cần lập kế hoạch ngay từ đầu về băng thông mạng, chuẩn truyền dữ liệu (MQTT, OPC-UA, gRPC) và phương án chống mất dữ liệu khi đường truyền gặp sự cố.

2. Thu Thập Dữ Liệu Thời Gian Thực

- Hạ tầng thu thập
 - Sử dụng gateway IoT để gom dữ liệu từ nhiều cảm biến, thực hiện lọc và chuẩn hóa trước khi gửi về máy chủ.
 - Dùng giao thức truyền tin gọn nhẹ như MQTT để giảm độ trễ.
- Đảm bảo tính toàn vẹn
 - Thêm checksum hoặc chữ ký số trong gói dữ liệu.
 - Cơ chế retry tự động nếu mất gói, lưu tạm (buffer) tại thiết bị biên để tránh mất mát.
- Đồng bộ thời gian
 - Thiết lập NTP cho toàn bộ thiết bị.
 - Kiểm tra lệch giờ định kỳ, đặc biệt quan trọng trong các phân tích tương quan giữa nhiều nguồn dữ liệu.

3. Tiền Xử Lý và Làm Giàu Dữ Liệu

Trước khi lưu trữ, dữ liệu cần được xử lý để phục vụ phân tích lâu dài:

	VIETTEL AI RACE	TD069
	Thu Thập và Quản Lý Dữ Liệu Chuỗi Thời Gian	Lần ban hành: 1

- Làm sạch dữ liệu:
 - Loại bỏ các điểm bất thường do lỗi cảm biến (outlier).
 - Điền giá trị thiếu bằng phương pháp nội suy tuyến tính hoặc spline.
- Chuẩn hóa đơn vị: tất cả cảm biến phải quy về cùng hệ đo (°C, bar, v.v.) để dễ so sánh.
- Tạo đặc trưng (feature engineering):
 - Tính toán trung bình trượt, phương sai, các chỉ số theo khung thời gian (5 phút, 1 giờ).
 - Sinh nhãn sự kiện (ví dụ: “máy dừng”, “nhiệt độ vượt ngưỡng”) cho các bài toán dự đoán.
- Nén và mã hóa: Dùng định dạng như Parquet hoặc Protobuf để giảm dung lượng mà vẫn đảm bảo truy vấn nhanh.

4. Hệ Thống Lưu Trữ và Quản Lý

- Cơ sở dữ liệu chuỗi thời gian (TSDB)
 - Lựa chọn InfluxDB, TimescaleDB hoặc QuestDB cho khả năng ghi dữ liệu hàng triệu điểm/giây.
 - Hỗ trợ truy vấn theo khoảng thời gian, downsampling tự động.
- Quản lý phiên bản và phân vùng
 - Thiết lập chính sách TTL (time-to-live) để lưu dữ liệu chi tiết trong thời gian nhất định rồi nén hoặc tổng hợp.

	VIETTEL AI RACE	TD069
	Thu Thập và Quản Lý Dữ Liệu Chuỗi Thời Gian	Lần ban hành: 1

- Dùng sharding hoặc partitioning theo ngày/tháng để tăng tốc độ truy vấn.
- Bảo mật và quyền truy cập
 - Mã hóa dữ liệu khi truyền và khi lưu.
 - Kiểm soát truy cập theo vai trò (Role-Based Access Control) để giới hạn người xem hoặc chỉnh sửa.

5. Giám Sát, Bảo Trì và Khả Năng Mở Rộng

- Giám sát luồng dữ liệu:
 - Dùng các công cụ như Prometheus, Grafana để theo dõi tốc độ ghi, độ trễ, số gói lỗi.
 - Thiết lập cảnh báo khi có bất thường, ví dụ lưu lượng giảm đột ngột.
- Chiến lược mở rộng
 - Kiến trúc microservices cho phép thêm node ghi/đọc mà không ảnh hưởng dịch vụ hiện tại.
 - Dùng cloud-native (Kubernetes) để tự động cân bằng tải.
- Sao lưu và phục hồi
 - Sao lưu định kỳ sang các vùng địa lý khác nhau.
 - Thử nghiệm phục hồi để đảm bảo RPO (Recovery Point Objective) và RTO (Recovery Time Objective) đạt yêu cầu.

	VIETTEL AI RACE	TD069
	Thu Thập và Quản Lý Dữ Liệu Chuỗi Thời Gian	Lần ban hành: 1

- Giám sát luồng dữ liệu:
 - Dùng các công cụ như Prometheus, Grafana để theo dõi tốc độ ghi, độ trễ, số gói lỗi.
 - Thiết lập cảnh báo khi có bất thường, ví dụ lưu lượng giảm đột ngột.
- Chiến lược mở rộng
 - Kiến trúc microservices cho phép thêm node ghi/đọc mà không ảnh hưởng dịch vụ hiện tại.
 - Dùng cloud-native (Kubernetes) để tự động cân bằng tải.
- Sao lưu và phục hồi
 - Sao lưu định kỳ sang các vùng địa lý khác nhau.
 - Thử nghiệm phục hồi để đảm bảo RPO (Recovery Point Objective) và RTO (Recovery Time Objective) đạt yêu cầu.