

HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG
KHOA CÔNG NGHỆ THÔNG TIN 1



BÁO CÁO LẦN I
MÔN HỌC: THỰC TẬP TỐT NGHIỆP
Doanh nghiệp: Công ty cổ phần VCCORP
Giảng viên hướng dẫn: Đỗ Thị Liên
Nhóm: 25

Họ tên: Hoàng Văn An

Mã sinh viên: B20DCCN045

Số điện thoại: 0867865001

Email: vananhoang10052002@gmail.com

Hà Nội – 2024

[illegible]

[illegible]

DANH MỤC THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
API: Application Programming Interface	Giao diện lập trình ứng dụng
Aerospike	Aerospike là một hệ thống cơ sở dữ liệu NoSQL hiệu năng cao được thiết kế để xử lý khối lượng dữ liệu lớn và cung cấp tốc độ truy xuất nhanh, hỗ trợ lưu trữ vĩnh viễn và khả năng mở rộng cao
Hbase	HBase là một cơ sở dữ liệu phân tán, mã nguồn mở, được thiết kế để cung cấp khả năng lưu trữ và truy xuất dữ liệu có cấu trúc lớn với hiệu suất cao, chạy trên hệ thống file HDFS của Hadoop.
Kafka	Kafka là một nền tảng streaming phân tán mạnh mẽ, được thiết kế để xử lý và truyền tải lượng lớn dữ liệu theo thời gian thực với độ tin cậy cao và khả năng mở rộng linh hoạt.
Spark	Spark là một framework xử lý dữ liệu phân tán mạnh mẽ, hỗ trợ xử lý dữ liệu lớn với hiệu suất cao thông qua các API dễ sử dụng cho SQL, streaming, machine learning, và graph processing.
Docker	Docker là một nền tảng mã nguồn mở giúp tự động hóa việc triển khai ứng dụng bên trong các container phần mềm, đảm bảo tính nhất quán, di động và hiệu quả trong việc phát triển và vận hành phần mềm.

CHƯƠNG I: GIỚI THIỆU CHUNG ĐƠN VỊ THỰC TẬP

Chương này cung cấp một cái nhìn tổng quan về đơn vị thực tập, bao gồm lịch sử hình thành và phát triển, tầm nhìn và sứ mệnh của công ty. Chương này cũng sẽ giới thiệu các sản phẩm và dịch vụ chính mà công ty cung cấp, nhấn mạnh vào những giá trị cốt lõi và cam kết của công ty trong việc mang lại các giải pháp công nghệ tiên tiến và hiệu quả cho khách hàng.

1.1. Thông tin về đơn vị thực tập:

1.1.1. Giới thiệu về VCCorp:

- Được thành lập vào năm 2006, sau 18 năm phát triển, Công ty CP VCCorp (VCCorp) là một trong những công ty tiên phong trong lĩnh vực công nghệ và công nghệ cao ở Việt Nam. VCCorp đã xây dựng và ứng dụng thành công hạ tầng cloud computing có khả năng lưu trữ và tính toán lượng dữ liệu lớn (Big Data, Data mining) hỗ trợ cho nền tảng công nghệ quảng cáo trực tuyến Admicro, phục vụ trên 50 triệu người dùng (chiếm trên 90% tổng số người dùng Internet tại Việt Nam) với hơn 200 website uy tín, gần 30 website trong số đó thuộc sở hữu độc quyền, phục vụ hàng tỷ lượt xem mỗi tháng. Ngoài ra, VCCorp nằm trong top 3 công ty phát hành game trên nền tảng di động với thương hiệu Sohagame, và mới đây hệ sinh thái chuyển đổi số Bizfly đang được VCCorp phát triển và triển khai mạnh mẽ. Đồng thời VCCorp cũng là đơn vị xây dựng Lotus - mạng xã hội do người Việt Nam làm chủ.



Hình 1.1. Tổng quan về các sản phẩm của VCCorp

1.2. Cơ sở thực tập:

1.2.1. Thông tin thực tập

- Tên doanh nghiệp: Công ty cổ phần VCCORP
- Vị trí thực tập: Java Backend Developer
- Khối được phân: Khối Adtech của VCCORP
- Leader, mentor hướng dẫn: Leader Ngô Văn Vĩ
- Hình thức thực tập: Tại văn phòng công ty
- Mô tả công việc của team:
 - Xử lý các phần backend Server của các hệ thống phân tán, các hệ thống dữ liệu, data minning.
 - Tối ưu, xử lý dữ liệu.

1.2.2. Đội ngũ hướng dẫn thực tập sinh:

- Leader Ngô Văn Vĩ:
 - Chuyên gia với kiến thức chuyên sâu về các công nghệ mới nhất và xu hướng phát triển trong ngành.

1.2.3. Cơ sở vật chất:

- VCCorp có văn phòng được thiết kế theo tiêu chuẩn công nghệ cao, với không gian mở và các khu vực làm việc riêng biệt. Thiết kế văn phòng tối ưu hóa sự linh hoạt và sự giao tiếp giữa các bộ phận..
- Công nghệ, thiết bị, cập nhật liên tục, đảm bảo nhân viên luôn tiếp cận với những công nghệ mới nhất.
- Cơ sở vật chất của VCCorp được đảm bảo về mặt an ninh và an toàn, với hệ thống bảo vệ chuyên nghiệp và các biện pháp phòng cháy chữa cháy đầy đủ.

1.2.4. Dịch vụ hỗ trợ thực tập sinh:

- Đào tạo và hướng dẫn:
 - Cung cấp chương trình đào tạo ban đầu để giúp thực tập sinh làm quen với công việc và các công cụ cần thiết.
 - Được hướng dẫn bởi các chuyên gia có kinh nghiệm trong ngành, giúp thực tập sinh nâng cao kỹ năng chuyên môn.
- Môi trường làm việc chuyên nghiệp:
 - Cung cấp môi trường làm việc hiện đại, tiện nghi với đầy đủ trang thiết bị cần thiết.
 - Không gian làm việc thân thiện, khuyến khích sự sáng tạo và hợp tác giữa các nhân viên.
- Phản hồi và đánh giá:
 - Thường xuyên cung cấp phản hồi và đánh giá hiệu quả làm việc, giúp thực tập sinh nhận biết được điểm mạnh và điểm cần cải thiện.
 - Hỗ trợ phát triển cá nhân thông qua các buổi mentoring và coaching.
- Cơ hội phát triển nghề nghiệp:
 - Tạo điều kiện để thực tập sinh tham gia vào các dự án thực tế, tích lũy kinh nghiệm quý báu.
 - Cơ hội trở thành nhân viên chính thức của VCCorp sau khi hoàn thành chương trình thực tập xuất sắc.

1.2.5. Thời gian thực tập:

- Thực tập từ 26/6/2024 đến 11/8/2024
- Lên công ty theo lịch đã đăng kí, buổi sáng từ 9h đến 12h, buổi chiều từ 13h30 đến 18h

1. NỘI DUNG CÔNG VIỆC TUẦN 2 (1/7 – 7/7):

1.1. Tìm hiểu về các công nghệ thực tế được sử dụng trong công ty

- **Track 1: Tìm hiểu và viết báo cáo về Aerospike, sau đó thực hiện các task sử dụng công nghệ đã tìm hiểu đó, mentor giải đáp các thắc mắc về công nghệ**
 - Cài đặt Aerospike
 - Ghi 1 bản ghi vào Aerospike (Ghi vĩnh viễn/có expire time)
 - Đọc 1 bản ghi từ Aerospike

- Đọc nhiều bản ghi từ Aerospike
 - Sử dụng AQL
 - Tìm hiểu về các operation khác của Aerospike
 - Tìm hiểu về EventLoop và Callback trong Aerospike
 - Scan và Query
 - Tìm hiểu index trong Aerospike
 - So sánh giữa Aerospike với Redis
 - Tìm hiểu các chiến lược caching cho cả đọc và ghi
- **Track 2: Tìm hiểu và viết báo cáo về HBase, sau đó thực hiện các task sử dụng công nghệ đã tìm hiểu đó, mentor giải đáp các thắc mắc về công nghệ**
- Cài đặt HBase (có thể dùng docker)
 - Thao tác bằng HBase shell
 - Thao tác bằng Java API
 - Get
 - Exist bin
 - Exist row
 - Put
 - Delete bin
 - Delete Row
 - Bulk
 - Time To Live
 - Scan + Filter
 - Pagination using HBase Scan
 - Tìm hiểu vai trò của ZooKeeper trong HBase
 - Tìm hiểu về Compact trong Hbase
 - Tìm hiểu các thành phần cấu hình khi cài đặt Hbase

1.2. Những kết quả đạt được:

- Hoàn thành cài đặt Aerospike và HBase.
- Thực hiện thành công các thao tác cơ bản với Aerospike và HBase.
- Viết báo cáo chi tiết về các công nghệ, cách cài đặt, và các thao tác cơ bản.
- So sánh Aerospike và Redis, nắm rõ các ưu nhược điểm của từng công nghệ.
- Áp dụng các chiến lược caching hiệu quả.

1.3. Những điều chưa đạt được:

- Chưa tối ưu hóa toàn diện cho các thao tác với dữ liệu lớn trong Aerospike và HBase.
- Chưa triển khai các kịch bản phức tạp sử dụng EventLoop và Callback trong Aerospike.
- Cần tìm hiểu thêm về các mô hình dữ liệu phức tạp và các phương pháp tối ưu hóa hiệu suất trong HBase.

1.4. Bài học kinh nghiệm:

- Cần tìm hiểu kỹ lưỡng về cấu trúc và mô hình dữ liệu của từng công nghệ trước khi triển khai.
- Việc sử dụng Docker giúp giảm thiểu thời gian cài đặt và cấu hình môi trường phát triển.
- Cần thường xuyên tham khảo tài liệu và hướng dẫn chính thức của từng công nghệ để nắm bắt các tính năng mới và các phương pháp tối ưu hóa.
- Việc so sánh và lựa chọn công nghệ cần dựa trên yêu cầu cụ thể của dự án và khả năng mở rộng trong tương lai.

2. NỘI DUNG CÔNG VIỆC TUẦN 2 (7/7 – 14/7):

2.1. Tìm hiểu về các công nghệ thực tế được sử dụng trong công ty

- **Track 3: Tìm hiểu và viết báo cáo về Kafka, sau đó thực hiện các task sử dụng công nghệ đã tìm hiểu đó, mentor giải đáp các thắc mắc về công nghệ.**
 - Cài đặt Kafka
 - Produce message
 - Consume message
 - Consumer group
 - Tìm hiểu về partition và replicate trong Kafka
 - Tìm hiểu vai trò của ZooKeeper trong Kafka (tìm hiểu cả các loại khác thay thế zookeeper)
 - So sánh với rabbitmq
- **Track 4: Tìm hiểu và viết báo cáo về Spark, sau đó thực hiện các task sử dụng công nghệ đã tìm hiểu đó, mentor giải đáp các thắc mắc về công nghệ**
 - Tìm hiểu cách thức hoạt động của Spark
 - Tìm hiểu RDD
 - Tìm hiểu mô hình MapReduce
 - Tìm hiểu Spark Streaming
 - Tìm hiểu HDFS, Parquet file
 - Spark SQL
 - Spark SQL Optimization

2.2. Kết quả đạt được:

- Hoàn thành cài đặt và cấu hình Kafka và Spark.
- Thực hiện thành công các thao tác cơ bản với Kafka và Spark.
- Viết báo cáo chi tiết về các công nghệ, cách cài đặt, và các thao tác cơ bản.
- So sánh Kafka với RabbitMQ, nắm rõ các ưu nhược điểm của từng công nghệ.
- Hiểu cách tối ưu hóa truy vấn và xử lý dữ liệu với Spark.

2.3. Những điều chưa đạt được

- Chưa tối ưu hóa toàn diện cho các thao tác với dữ liệu lớn trong Kafka và Spark.
- Chưa triển khai các kịch bản phức tạp sử dụng Spark Streaming và các mô hình phức tạp trong Kafka.

- Cần tìm hiểu thêm về các phương pháp tối ưu hóa hiệu suất trong Kafka và Spark.

2.4. Bài học kinh nghiệm

- Cần tìm hiểu kỹ lưỡng về cấu trúc và mô hình dữ liệu của từng công nghệ trước khi triển khai.
- Việc sử dụng Docker giúp giảm thiểu thời gian cài đặt và cấu hình môi trường phát triển.
- Cần thường xuyên tham khảo tài liệu và hướng dẫn chính thức của từng công nghệ để nắm bắt các tính năng mới và các phương pháp tối ưu hóa.
- Việc so sánh và lựa chọn công nghệ cần dựa trên yêu cầu cụ thể của dự án và khả năng mở rộng trong tương lai.