



دانشگاه صنعتی امیرکبیر

(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

گزارش کارآموزی

محل کارآموزی: شرکت سامانه گستر سحاب پرداز

نام استاد کارآموزی

دکتر مسعود صبائی

نام دانشجو

امیرمحمد پیرحسینلو

۹۵۳۱۰۱۴

تابستان ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## چکیده

این دوره‌ی کارآموزی حول محور داده‌های حجیم برگزار شد و در آن با مفاهیم و ابزارهای پردازش داده‌ها آشنا شدیم، برنامه نویسی سیستم‌های توزیع شده را فراگرفتیم و با انجام دو پروژه بر مهارت‌های فنی خود افزودیم. جدا از نکات فنی، اجرای یک کار گروهی ۶ نفره از نکات برجسته‌ی این دوره بود که باعث شد کار در یک محیط واقعی با یک پروژه واقعی را تجربه کنیم. هدف از این گزارش، شرح فعالیت‌های انجام شده و مطالب فراگرفته در این دوره است که در هفت فصل به آن‌ها خواهیم پرداخت.

## واژه‌های کلیدی:

داده‌های حجیم، معماری توزیع شده، جاوا، پایگاه داده، موتور جستجو

# فهرست مطالب

د	فهرست اشکال
۱	۱ مقدمه
۳	۲ معرفی شرکت سحاب پرداز و روند کارآموزی در این شرکت
۴	۱-۲ مقدمه
۴	۲-۲ معرفی شرکت سحاب پرداز
۴	۳-۲ روند کارآموزی
۴	۱-۳-۲ نیمبو
۴	۲-۳-۲ مدت زمان کارآموزی
۴	۳-۳-۲ گروه بندی
۴	۴-۳-۲ پروژه اول
۵	۵-۳-۲ پروژه دوم
۵	۴-۲ C++ vs Java
۵	۱-۴-۲ نبود garbage collector در C++
۵	۲-۴-۲ ارث بری چندگانه
۵	۳-۴-۲ جاوا مستقل از سکو است
۶	۳ پروژه اول
۷	۱-۳ مقدمه
۷	۲-۳ تعریف پروژه اول
۷	۳-۳ ابزارهای پیاده سازی پروژه
۷	۱-۳-۳ زبان برنامه نویسی جاوا
۷	۲-۳-۳ git
۸	۳-۳-۳ maven
۸	۴-۳ مفاهیم و اصطلاحات
۸	۱-۴-۳ channel
۸	۲-۴-۳ publisher
۸	۳-۴-۳ subscriber
۸	۴-۴-۳ real time messaging (rtm)
۹	۵-۴-۳ satori rtm
۹	۵-۳ github channel
۱۰	۶-۳ satori API
۱۰	۷-۳ معماری نرم افزار

۱۱	۸-۳ نتیجه‌گیری
۱۱	۱-۸-۳ code
۱۱	۲-۸-۳ ساخت فایل jar و اجرای آن
۱۲	۴ پروژه دوم - آشنایی با ابزارها و مفاهیم جدید در حوزه‌ی داده‌های حجیم
۱۳	۱-۴ مقدمه
۱۳	۲-۴ تعریف پروژه دوم
۱۳	۳-۴ معرفی ابزارها
۱۳	۱-۳-۴ فریم‌ورک Hadoop
۱۵	۲-۳-۴ پایگاه داده HBase
۱۶	۳-۳-۴ موتور جستجوی Elasticsearch
۱۷	۴-۴ حالات نصب
۱۷	۱-۴-۴ Standalone mode
۱۷	۲-۴-۴ Distributed mode
۱۷	۵-۴ نصب ابزارها و اجرای کدهای قسمت‌های قبل
۱۸	۶-۴ نتیجه‌گیری
۱۹	۵ معماری پروژه دوم (موتور جستجو)
۲۰	۱-۵ مقدمه
۲۰	۲-۵ اجزای موتور جستجو
۲۱	۱-۲-۵ صف لینک‌ها
۲۱	۲-۲-۵ Fetcher
۲۲	۳-۲-۵ صف اسناد
۲۲	۴-۲-۵ Parser
۲۲	۵-۲-۵ Kafka
۲۲	۳-۵ نتیجه‌گیری
۲۳	۶ پیاده‌سازی پروژه دوم (موتور جستجو)
۲۴	۱-۶ مقدمه
۲۴	۲-۶ پکیج‌ها
۲۴	۱-۲-۶ crawler
۲۴	۲-۲-۶ index
۲۴	۳-۲-۶ storage
۲۴	۴-۲-۶ test
۲۵	۵-۲-۶ utils

۲۵	language ۶-۲-۶
۲۵	کلاس‌ها ۳-۶
۲۵	Crawler ۱-۳-۶
۲۵	Fetcher ۲-۳-۶
۲۶	LruCache ۳-۳-۶
۲۶	Parser ۴-۳-۶
۲۶	Elastic ۵-۳-۶
۲۷	HBase ۶-۳-۶
۲۷	Constants ۷-۳-۶
۲۷	Statistics ۸-۳-۶
۲۷	Pair ۹-۳-۶
۲۷	استثناها ۴-۶
۲۷	اجرای برنامه ۵-۶
۲۸	موارد باقی مانده ۶-۶
۲۸	نتیجه گیری ۷-۶
۲۹	بهبود عملکرد موتور جستجو ۷
۳۰	مقدمه ۱-۷
۳۰	بهبود جستجو با محاسبه تعداد لینک‌های ورودی ۲-۷
۳۰	بهبود جستجو با اعمال متن anchor link ۳-۷
۳۱	بهبود جستجو با Page Rank ۴-۷
۳۱	نتیجه گیری ۵-۷
۳۲	نتیجه گیری ۸
۳۴	واژه‌نامه‌ی انگلیسی به فارسی
۱۰	معماری نرم‌افزار پروژه اول ۱-۳
۲۰	معماری نرم‌افزار موتور جستجو ۱-۵

## فهرست اشکال

## فهرست اختصارات

فرم کامل	مخفف
<b>Hadoop Distributed File System</b>	<i>HDFS</i>
<b>HyperText Markup Language</b>	<i>HTML</i>
<b>HyperText Transfer Protocol</b>	<i>HTTP</i>
<b>JavaScript Object Notation</b>	<i>JSON</i>
<b>Java Virtual Machine</b>	<i>JVM</i>
<b>Real Time Messaging</b>	<i>RTM</i>
<b>Uniform Resource Locator</b>	<i>URL</i>
<b>Yet Another Resource Negotiator</b>	<i>YARN</i>

## فصل اول

### مقدمه



در این دوره‌ی کارآموزی، با ابزارها و مفاهیم جدیدی در حوزه‌ی داده‌های حجیم آشنا شدیم. هدف از این دوره یادگیری کار با ابزارهای برنامه نویسی، پایگاه داده‌های توزیع شده و سیستم عامل لینوکس بود که تا حد خیلی خوبی با آن‌ها آشنا شدیم. دو پروژه تعریف شد و در قالب تیم‌های ۶ نفره پیاده‌سازی شد. در هر کدام از پروژه‌ها ابتدا با ابزارهای مورد نیاز برای انجام پروژه آشنا شدیم، سپس طراحی معماری نرم‌افزار صورت گرفت و بعد از آن به پیاده‌سازی و برنامه نویسی پرداختیم.

در فصل دوم به معرفی شرکت سحاب پرداز و زمینه‌ی کاری آن می‌پردازیم، سپس به سراغ روند کارآموزی در این شرکت می‌رویم. در فصل سوم، سراغ پروژه اول و روند پیاده‌سازی آن می‌رویم.

در فصل‌های چهارم و پنجم، ابزارهای مورد نیاز برای ساخت پروژه دوم که یک موتور جستجو است، معماری نرم‌افزار آن و نحوه پیاده‌سازی و برنامه نویسی آن شرح داده خواهد شد.

در فصل ششم ایده‌هایی برای بهبود عملکرد موتور جستجوی ساخته شده در پروژه دوم مطرح می‌شود و نحوه پیاده‌سازی یکی از ایده‌ها توضیح داده می‌شد.

فصل آخر هم به نتیجه‌گیری اختصاص یافته است.

## فصل دوم

# معرفی شرکت سحاب پرداز و روند کارآموزی در این شرکت

## ۱-۲ مقدمه

در این فصل ابتدا به معرفی شرکت سحاب پرداز و زمینه‌ی کاری آن می‌پردازیم، سپس به سراغ روند کارآموزی در این شرکت می‌رویم و در انتها مطالبی که در هفته‌ی اول یادگرفتیم را بازگو می‌کنیم.

## ۲-۲ معرفی شرکت سحاب پرداز

این شرکت خصوصی در سال ۱۳۹۳ تاسیس شده است و هم اکنون در زمینه‌های مختلفی فعالیت می‌کند که مهم‌ترین آن **Big Data**<sup>۱</sup> است و در حال ایجاد محصولات خاص و تجربیات جدید در این زمینه است. ارزشمندترین سرمایه سحاب، اعضای خانواده و فرهنگی است که با آن زندگی می‌کنیم. اعضای این شرکت را ۷۵ نفر تشکیل می‌دهند که در بخش‌های منابع انسانی، زیرساخت، مدیریت و ... فعالیت می‌کنند.

آدرس: خیابان سهروردی، خیابان خرمشهر، پ ۲۴

شماره تماس: ۸۸۷۶۲۶۷۲ (۰۲۱)

ایمیل: [info@sahab.ir](mailto:info@sahab.ir)

## ۳-۲ روند کارآموزی

### ۱-۳-۲ نیمبو

نیمبو نام دوره کارآموزی فشرده است که تابستان امسال توسط شرکت سحاب پرداز برگزار می‌شود و مسیری است برای توسعه مهارت‌های فنی و غیر فنی که با آن می‌توان به دنیای حرفه‌ای‌ها وارد شد. در این دوره اهداف یادگیری متنوعی مبتنی بر پروژه در نظر گرفته شده است که در حوزه **Big Data** و **Data Science**<sup>۲</sup> است. از همه مهم‌تر مربی‌هایی در کنار ما هستند که در یادگیری به ما و تیم‌مان کمک می‌کنند. این مربی‌ها از با تجربه‌ترین همکاران سحاب انتخاب شده‌اند. آموزش و مهارت‌های دوره نیمبو طوری طراحی شده است که در آینده بتوانیم کاربردش را در مسیر شغلی و تحصیلی‌مان ببینیم.

### ۲-۳-۲ مدت زمان کارآموزی

کارآموزی در ۶ هفته برگزار می‌شود که در هر هفته از شنبه تا چهارشنبه از ساعت ۹ صبح تا ۶ بعد از ظهر باید در شرکت حضور داشته باشیم.

### ۳-۳-۲ گروه بندی

در این دوره ۱۸ نفر کارآموز حضور دارند که برای پروژه اول به ۹ تیم دو نفره و برای پروژه دوم به ۳ تیم ۶ نفره تقسیم می‌شوند.

### ۴-۳-۲ پروژه اول

پروژه اول در رابطه با آشنایی با **real time messaging (rtm)** است. در این پروژه قصد داریم به عنوان **Subscriber** از کانال **github-events** داده‌ها را دریافت کنیم و با پردازش داده‌ها افرادی که بیشترین فعالیت را در **github** داشته‌اند را بیابیم. بیشترین فعالیت بر اساس تعداد **commit** ها و **push** ها و چند متغیر دیگر تعریف می‌شود. جزئیات بیشتر درمورد روند کار با یک پروتکل **rtm**

<sup>۱</sup>داده‌ی حجیم  
<sup>۲</sup>علم داده

در سایت **satori**<sup>۳</sup> موجود است. زبان برنامه نویسی جاوا به همراه ابزار مدیریت پکیج **maven**<sup>۴</sup> برای پیاده سازی این پروژه انتخاب شده است.

## ۵-۳-۲ پروژه دوم

تعریف پروژه دوم پس از پایان پروژه اول در اختیار ما قرار می گیرد.

## ۴-۲ C++ vs Java

اولین کارگاه آموزشی در این هفته به مدت یک ساعت و نیم برگزار شد و محور آن بحث درمورد تفاوت های زبان های برنامه نویسی **Java** و **C++** بود. این کارگاه به این دلیل برگزار شد که تعدادی از کارآموزها در درس برنامه نویسی پیشرفته زبان **C++** را یاد گرفته بودند و با **Java** آشنایی نداشتند. مهم ترین تفاوت ها به شرح زیر است:

### ۱-۴-۲ نبود garbage collector در C++

زبان **C++** فاقد **garbage collector** است به همین دلیل برنامه نویس باید دقت کند هر جا که نیاز به یک **object**<sup>۵</sup> از بین می رود باید **destructor**<sup>۶</sup> آن را فراخوانی کند و هر جا نیاز به یک آرایه ندارد آن را مستقیماً **delete**<sup>۷</sup> کند در غیر این صورت فضای حافظه اختصاص داده شده به این موارد در حافظه سیستم آزاد نشده و در اختیار سیستم عامل قرار نمی گیرد.

### ۲-۴-۲ ارث بری چندگانه

یکی از ویژگی های زبان برنامه نویسی **C++**، ارث بری چندگانه<sup>۸</sup> است که در جاوا وجود ندارد. این ویژگی مزایا و معایب خود را دارد که از آن ها صرف نظر می کنیم اما برنامه نویس باید دقت کند که در صورت استفاده از ارث بری چندگانه ممکن است به مشکل **diamond**<sup>۹</sup> برخورد کند.

### ۳-۴-۲ جاوا مستقل از سکو است

یکی از مزایای عالی زبان **Java**، مستقل بودن آن از سیستم عامل<sup>۱۰</sup> یا سکو است. برنامه های جاوا در سیستم عاملی که **Java Virtual Machine (JVM)** در آن نصب است اجرا می شود و برنامه نویس را از نوشتن چند برنامه برای سیستم عامل های مختلف رها می سازد. اما زبان **C++** این گونه نیست و به دلیل وجود مترجم های مختلف مانند **g++** و **visual studio compiler** و ... برای هر سیستم عامل، ممکن است کدی که برای یک سیستم عامل نوشته شده باشد برای سیستم عامل دیگر قابل اجرا نباشد. تفاوت های دیگری هم وجود دارد که در گزارش به موارد بالا اکتفا می کنیم.

<sup>۳</sup><https://www.satori.com/docs/rtm-sdks/tutorials/java-sdk-quickstart>

<sup>۴</sup><https://maven.apache.org>

<sup>۵</sup>multiple inheritance

<sup>۹</sup><https://www.freecodecamp.org/news/multiple-inheritance-in-c-and-the-diamond-problem-7c12a9dddbec/>

<sup>۱۰</sup>platform independency

## فصل سوم

### پروژه اول

## ۱-۳ مقدمه

در هفته‌ی دوم، یک پروژه تعریف شد که باید در قالب تیم‌های دو نفره انجام می‌شد. تعریف پروژه در ادامه آمده است. هدف این پروژه آشنایی با رابط‌های برنامه نویسی و تمرین کار گروهی است. در انتها کدهای تیم‌ها مورد ارزیابی قرار گرفت و نواقص رفع شد. در این فصل ابتدا به تعریف پروژه می‌پردازیم سپس سراغ معرفی ابزارهای پیاده‌سازی پروژه می‌رویم. در ادامه، مفاهیم و اصطلاحاتی را که در حوزه پروژه است، تعریف و بیان می‌کنیم و در قسمت آخر نحوه‌ی اجرای کد پروژه اول را نشان می‌دهیم.

## ۲-۳ تعریف پروژه اول

در این پروژه قصد داریم با استفاده از رابط برنامه نویسی سایت [satori.com](https://satori.com) به کانال [github](https://github.com)<sup>۱</sup> دسترسی پیدا کرده و اطلاعاتی را از آن استخراج کنیم. این اطلاعات شامل نام کاربری افراد با بیشترین تعداد **commit** و **push** در سایت [github.com](https://github.com) در یک ساعت اخیر است که به طور تناوبی هر یک دقیقه یک‌بار این اطلاعات استخراج می‌شود.

## ۳-۳ ابزارهای پیاده‌سازی پروژه

برای پیاده‌سازی پروژه از ابزارهای زیر استفاده شده‌است:

## ۱-۳-۳ زبان برنامه نویسی جاوا

زبان برنامه نویسی جاوا یک زبان سطح بالا است و یکی از قدرتمندترین زبان‌های برنامه نویسی به شمار می‌رود که در اکثر زمینه‌ها کاربرد دارد. از ویژگی‌های آن می‌توان به موارد زیر اشاره کرد:

- **platform independency**<sup>۲</sup>
- **garbage collector**<sup>۳</sup>
- فراهم بودن تعداد بسیار زیادی از **library**<sup>۴</sup> ها و **API (Application Programming Interface)**<sup>۵</sup> ها
- مستندات عالی برای زبان
- تعداد بسیار زیادی ابزار مدیریت **package** و...

## ۲-۳-۳ git

**git** یک ابزار مدیریت نسخه است که در توسعه نرم‌افزارها از آن استفاده می‌شود. از کاربردهای آن می‌توان به نگهداری کدها، آسان‌سازی برنامه نویسی برای تیم‌ها و موارد بسیار زیاد دیگر اشاره کرد.

<sup>۱</sup><https://www.github.com>

<sup>۲</sup>مستقل از سکو است، یعنی در هر سیستم عاملی می‌توان از آن استفاده کرد.

<sup>۳</sup><https://www.geeksforgeeks.org/garbage-collection-java>

<sup>۴</sup>کتابخانه

<sup>۵</sup>رابط برنامه نویسی

## ۳-۳-۳ maven

**maven** یک ابزار مدیریت **package** است که برای ساخت نرم‌افزارهای جاوایی از آن استفاده می‌شود. هنگامی که تعداد کتابخانه‌های استفاده شده در پروژه زیاد می‌شود، به شدت روند ایجاد فایل **jar** را آسان می‌کند.

## ۴-۳ مفاهیم و اصطلاحات

پیش از شرح معماری پروژه بهتر است بعضی از اصطلاحات و مفاهیم توضیح داده شود تا ابهامات رفع شود و ابعاد پروژه روشن گردد.

## ۱-۴-۳ channel

**channel** یک **queue**<sup>۶</sup> است که همزمان چندین **agent**<sup>۸</sup> می‌توانند در آن داده **push** کنند و از طرف دیگر **agent** های دیگر می‌توانند داده‌ها را دریافت کنند. به دلیل اینکه چندین **agent** در حال کار با این صف هستند، این صف باید **thread-safe**<sup>۹</sup> باشد تا مشکل **race condition**<sup>۱۰</sup> پیش نیاید.

یک نکته مهم: عامل‌هایی که از صف داده‌ها را می‌خوانند می‌توانند به همه‌ی داده‌های قرار داده شده در کانال دسترسی داشته باشند.

## ۲-۴-۳ publisher

عامل‌هایی هستند که در صف (کانال) داده قرار می‌دهند.

برای مثال یک برنامه می‌تواند **log** های خروجی خود را داخل یک کانال قرار دهد. در اینجا **log** همان داده است که در کانال قرار می‌گیرد و برنامه هم همان عامل است

## ۳-۴-۳ subscriber

عامل‌هایی که داده‌ها را از صف (کانال) برداشته و مورد استفاده قرار می‌دهند. برای مثال یک برنامه ناظر که خروجی برنامه‌های دیگر را از کانال دریافت کرده و عمل نظارت را انجام می‌دهد. در صورتی که خطایی در یک برنامه رخ داده باشد یا میزان **load**<sup>۱۱</sup> در یک برنامه بالا رفته باشد، هشدار می‌دهد. در اینجا عامل همان برنامه‌ی ناظر است و داده‌های کانال خروجی برنامه‌ها هستند.

## ۴-۴-۳ real time messaging (rtm)

امروزه با پیشرفت تکنولوژی‌ها تعداد پیام‌رسان‌ها و شبکه‌های اجتماعی زیاد شده است. برای پیاده‌سازی این نرم‌افزارها معماری‌های زیادی وجود دارد. یکی از این معماری‌ها، معماری **publisher-subscriber**<sup>۱۲</sup> است.

در این معماری تعدادی از کاربرها (عامل‌ها) عضو یک کانال هستند پیام‌های آن کانال را دریافت می‌کنند و مورد استفاده قرار می‌دهند. به این کاربران **subscriber** گفته می‌شود. تعدادی از کاربران (عامل‌ها) هم در کانال داده (پیام) قرار می‌دهند که به آن‌ها **publisher** گفته می‌شود.

برای مثال در پیام‌رسان تلگرام کانال سایت **varzesh3.com** وجود دارد که تعداد زیادی کاربر دارد که هر روزه از اخبار جدید ورزشی

۶کانال  
۷صف  
۸عامل

<sup>۹</sup><https://www.baeldung.com/java-thread-safety>

<sup>۱۰</sup><https://stackoverflow.com/questions/34510/what-is-a-race-condition>

<sup>۱۲</sup>[https://en.wikipedia.org/wiki/Publish-subscribe\\_pattern](https://en.wikipedia.org/wiki/Publish-subscribe_pattern)

<sup>۱۱</sup>بارگیری

بهره‌مند می‌شوند. این کاربران نقش **subscriber** را دارند. ادمین‌های این کانال نقش **publisher** را ایفا می‌کنند چون آن‌ها مطالب را در کانال قرار می‌دهند.

فرایند ذکر شده به این صورت است که بلافاصله بعد از اینکه یک ادمین مطلبی را در کانال قرار داد، این مطلب در اختیار کاربران قرار می‌گیرد. به همین دلیل این فرایند یک پیام‌رسانی برخط (**rtm**) است.

### ۵-۴-۳ satori rtm

وبسایت [www.satori.com](http://www.satori.com) این امکان را فراهم می‌کند که با استفاده از یک رابط کاربری خوب بدون درگیر شدن با جزئیات ایجاد کانال به صورت نرم‌افزاری، با چند کلیک یک کانال ایجاد کنیم (برای استفاده‌های غیرتجاری ایجاد کانال رایگان است). این کانال بر روی سرورهای این شرکت ایجاد شده و می‌توانیم با استفاده از رابط برنامه نویسی از آن استفاده کنیم.

در ابتدا باید در وبسایت ثبت نام کنیم. بعد از ساخت کانال یک **endpoint** و یک **appkey** در اختیار ما قرار می‌گیرد که برای وصل شدن به کانال به آن‌ها نیاز داریم. پیکربندی‌های مختلفی برای یک کانال وجود دارد که به دو مورد آن در این جا می‌پردازیم.

- یک کانال می‌تواند **public**<sup>۱۳</sup> باشد یعنی هر کسی با دانستن نام کانال و **appkey** و **endpoint** می‌تواند به داده‌های آن دسترسی پیدا کند.

- یک کانال می‌تواند **private**<sup>۱۴</sup> باشد که در این صورت برای دسترسی به داده‌های آن باید احراز هویت به وسیله نام کاربری و رمز عبور صورت گیرد.

برای کانال‌ها پیکربندی‌های زیادی وجود دارد. برای مثال می‌توان تعیین کرد که چه داده‌هایی وارد کانال شود و جلوی ورود چه نوع داده‌هایی گرفته شود.

اعمال زیادی نیز می‌توان بر روی داده‌های کانال‌ها انجام داد. برای مثال:

- مشاهده **history**<sup>۱۵</sup> کانال که شامل داده‌های پیشین کانال است.

- فیلترکردن پیام‌های دریافتی بر اساس پارامترها

- دریافت پیام‌ها در یک بازه زمانی خاص به صورت **stream**<sup>۱۶</sup>

### ۵-۳ github channel

کانال **github** یک کانال عمومی است که رویدادهای اخیر سایت [github.com](http://github.com) در آن ثبت می‌شود. هر کاربری با داشتن **appkey** و **endpoint** آن می‌تواند به داده‌های این کانال دسترسی پیدا کند.

داده‌های این کانال شامل رویدادهای رخ داده مانند ایجاد یک **commit**، عملیات **push**، **tag** و ... است. از این کانال برای پیاده‌سازی پروژه استفاده می‌کنیم.

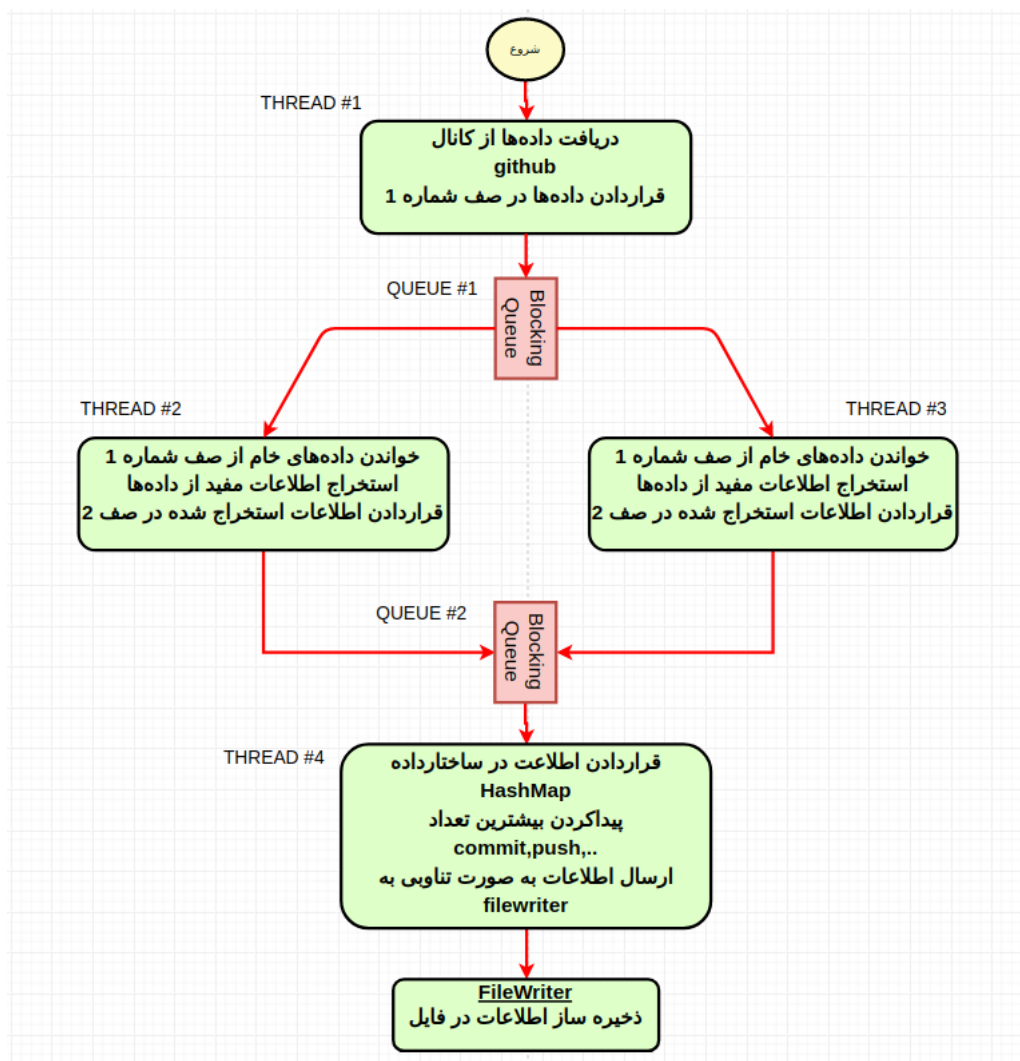
<sup>۱۳</sup> عمومی

<sup>۱۴</sup> خصوصی

<sup>۱۵</sup> تاریخچه

<sup>۱۶</sup> جریان





شکل ۳-۱: معماری نرم‌افزار پروژه اول

### ۶-۳ satori API

سایت [satori.com](https://satori.com) یک رابط برنامه نویسی در اختیار ما قرار می‌دهد که با استفاده از کلاس‌ها و فراخوانی توابع آن و در دست داشتن نام کانال، **appkey** و **endpoint** می‌توانیم به یک کانال وصل شده و از آن داده دریافت کنیم. برای ارسال داده هم باید علاوه بر موارد ذکر شده احراز هویت نیز انجام دهیم. در این پروژه به کانال **github** متصل شده و از آن داده دریافت می‌کنیم و از رویدادها مطلع می‌شویم.

### ۷-۳ معماری نرم‌افزار

شمای کلی نرم‌افزار در شکل ۳-۱ ترسیم شده است.

به صورت صریح چهار نخ در برنامه حضور دارند و از دو صف استفاده شده است.

ابتدا توسط نخ شماره یک پیام‌ها از کانال **github** توسط رابط برنامه نویسی سایت **satori** دریافت شده و داخل یک **blocking**

**queue**<sup>۱۷</sup> قرار داده می‌شود. حال سوالی که مطرح می‌شود این است که **blocking queue** چه فرقی با صف معمولی دارد؟ **blocking queue** صفی است که **thread-safe** است یعنی همزمان می‌توان در آن داده قرار داد و همزمان داده خواند. چون در این جا یک نخ در حال قرار دادن داده‌های کانال در صف و دو نخ در حال خواندن داده‌ها هستند از این ساختار داده استفاده کردیم. این ساختار داده در زبان جاوا پیاده‌سازی شده است.

از آن جا که سرور با سرعت زیادی داده در کانال قرار می‌دهد (تعداد کاربران **github** خیلی زیاد است به همین دلیل تعداد رخدادها زیاد است)، با یک نخ داده‌ها را در صف یک قرار داده و با دو نخ داده‌ها را از آن صف می‌خوانیم و پردازش می‌کنیم تا صف اشباع نشود. این تعداد نخ به صورت تجربی به دست آمده است.

پردازش به این شکل است که تنها داده‌هایی که شامل **commit** و **push** است را فیلتر کرده و سایر داده‌ها را دور می‌ریزیم. حال داده‌های به دست آمده را در صف شماره دو قرار می‌دهیم تا یک نخ دیگر عملیات نهایی را انجام دهد.

نخ شماره چهار داده‌های فیلتر شده را دریافت کرده و آن‌ها را در یک **HashMap**<sup>۱۸</sup> ذخیره می‌کند. در این ساختار داده، **Key** برابر نام کاربری و **value** شامل تعداد **commit** ها و **push** های شخص در یک ساعت اخیر است. هر یک دقیقه یک بار نام کاربری شخصی که بیشترین تعداد **commit** و **push** را دارد از **HashMap** استخراج شده و به مازول **FileWriter** فرستاده می‌شود تا در یک فایل ذخیره شود.

در انتها این فایل خروجی برنامه است که هر یک دقیقه یکبار نام شخصی که بیشترین فعالیت را در یک ساعت اخیر در **github** داشته است در انتهای آن ثبت می‌شود.

### ۸-۳ نتیجه‌گیری

برنامه‌ی مورد نظر در یک هفته آماده شد و نحوه‌ی **build**<sup>۱۹</sup> و اجرای آن در ادامه آمده است. با اجرای این پروژه، یک کار تیمی برنامه نویسی دو نفره را تجربه کردیم و با مفاهیم زیادی در حوزه معماری **publisher-subscriber** آشنا شدیم.

### code ۱-۸-۳

کد در آدرس زیر موجود است:

<https://github.com/Su6lime/GitHubTrends>

### ۲-۸-۳ ساخت فایل **jar** و اجرای آن

**mvn clean install**

**java -jar GitHubTrends-1.0-SNAPSHOT-jar-with-dependencies.jar**

**java -jar GitHubTrendsGUI.jar**

<sup>۱۷</sup><https://docs.oracle.com/javase/7/docs/api/java/util/concurrent/BlockingQueue.html>

<sup>۱۸</sup><https://docs.oracle.com/javase/8/docs/api/java/util/HashMap.html>

<sup>۱۹</sup>فرایند ساختن فایل **jar** را **build** می‌گویند.

## فصل چهارم

# پروژه دوم - آشنایی با ابزارها و مفاهیم جدید در حوزه‌ی داده‌های حجیم

## ۴-۱ مقدمه

در هفته سوم، تعریف پروژه دوم در اختیار ما قرار گرفت که هدف آن طراحی یک موتور جستجو است. برای طراحی چنین نرم‌افزاری، نیاز است با تعدادی از ابزارها و همچنین با مفاهیمی در حوزه‌ی داده‌های حجیم (**big data**) آشنا شویم. این ابزارها شامل موارد زیر می‌شوند:

- فریم ورک Hadoop<sup>۱</sup>

- پایگاه داده‌ی HBase<sup>۲</sup>

- موتور جستجوی Elasticsearch<sup>۳</sup>

- فریم ورک Spark<sup>۴</sup>

- پلتفرم Kafka<sup>۵</sup>

در هفته سوم با پایگاه داده‌های **ACID** و تئوری **CAP** آشنا شدیم. همچنین نحوه نصب و کار با ابزارها و فریم‌ورک‌های یاد شده را در سطح مقدماتی فراگرفتیم. در این فصل ویژگی‌های این ابزارها را توضیح خواهیم داد. طراحی معماری نرم‌افزار و برنامه نویسی در فصل‌های بعد بررسی می‌شوند.

## ۴-۲ تعریف پروژه دوم

هدف طراحی یک موتور جستجو است که باید حدود ۱۰۰ میلیون صفحه (**web page**) را بازدید و محتوای آن شامل لینک‌های خروجی و متن نوشته شده در صفحه را **index**<sup>۶</sup> کند. برای ذخیره‌سازی و پردازش داده‌ها دو سرور در اختیار ما قرار گرفته است. هر سرور ۸ ترابایت فضای ذخیره‌سازی از نوع دیسک سخت (**Hard Disk - HDD**)، ۱۶ گیگابایت حافظه **RAM** و یک **CPU** ۸ هسته دارد. فایل تعریف پروژه در کنار فایل گزارش قرار دارد.

## ۴-۳ معرفی ابزارها

در این بخش به معرفی جزئی چند تا از مهم‌ترین ابزارهای مورد استفاده در این پروژه می‌پردازیم.

### ۴-۳-۱ فریم‌ورک Hadoop

**Hadoop** یک فریم‌ورک برای ذخیره‌سازی داده‌های انبوه و پردازش آن‌ها است که بر بستر کامپیوترهای معمولی نصب و اجرا می‌شود. **Hadoop** با این فرض ساخته شده است که احتمال خراب شدن (**failure**) حافظه ذخیره‌سازی جانبی<sup>۷</sup> (معمولاً از نوع **HDD**) کم نیست، به همین دلیل فریم‌ورک طوری ساخته شده است که این خرابی را به صورت اتوماتیک تشخیص داده و با انجام عملیات‌هایی نظیر ایجاد افزونگی (**replication**) مانع از دست رفتن داده (**data loss**) می‌شود.

<sup>۱</sup>[https://en.wikipedia.org/wiki/Apache\\_Hadoop](https://en.wikipedia.org/wiki/Apache_Hadoop)

<sup>۲</sup>[https://en.wikipedia.org/wiki/Apache\\_HBase](https://en.wikipedia.org/wiki/Apache_HBase)

<sup>۳</sup><https://en.wikipedia.org/wiki/Elasticsearch>

<sup>۴</sup>[https://en.wikipedia.org/wiki/Apache\\_Spark](https://en.wikipedia.org/wiki/Apache_Spark)

<sup>۵</sup>[https://en.wikipedia.org/wiki/Apache\\_Kafka](https://en.wikipedia.org/wiki/Apache_Kafka)

<sup>۷</sup><https://www.komprise.com/glossary/terms/secondary-storage>

<sup>۶</sup>فهرست کردن

مدل پردازش داده‌ها در این فریم‌ورک، **Map/Reduce** است و یک موتور مخصوص این کار در دل **Hadoop** قرار دارد. از نسخه ۲ به بعد این موتور جای خود را به فریم‌ورک **YARN (Yet Another Resource Negotiator)** داده است، البته خود **YARN** از این موتور استفاده می‌کند.

این فریم‌ورک برای ذخیره‌سازی داده‌ها از یک سیستم مدیریت فایل به نام **HDFS (Hadoop Distributed File System)** استفاده می‌کند. داده‌ها به صورت بلوکی همراه با افزونگی توسط این **file system** در سرورها قرار می‌گیرند. این فریم‌ورک از پرونده‌های زیر تشکیل شده است:

### NameNode

مدیریت فایل سیستم **Hadoop** یا همان **HDFS** برعهده‌ی این پرونده است. متادیتا (**metadata**) ی فایل‌های موجود در خوشه (**cluster**) در اختیار این پرونده است. مثلاً اطلاعاتی مانند این که هر فایل به چند بلوک تقسیم شده است، هر بلوک در کدام سرور قرار دارد و برای هر بلوک چند افزونگی ایجاد شده و آن‌ها در کجا قرار دارند. تنها یک پرونده **NameNode** برای کل خوشه اجرا می‌شود. برای جلوگیری از **SPOF (Single Point Of Failure)** می‌توان یک **SecondaryNameNode** اجرا کرد که از **NameNode** اصلی به صورت دوره‌ای **snapshot**<sup>۸</sup> می‌گیرد و در صورت از کار افتادن **NameNode** اصلی، مانع از دست رفتن اطلاعات حیاتی و متادیتاها می‌شود.

### DataNode

داده‌های اصلی توسط پرونده‌های **DataNode** در سرورهای مختلف که در خوشه هستند نگهداری می‌شود. معمولاً به ازای هر سرور، یک پرونده **DataNode** اجرا می‌شود. این پرونده‌ها به صورت تناوبی در هر ۳ ثانیه صحت عملکرد خود را به **NameNode** گزارش می‌دهند.

### TaskTracker

داده‌ی مورد نظر برای انجام عملیات‌های **Map**<sup>۹</sup> و **Reduce**<sup>۱۰</sup> را از **DataNode** گرفته و پس از انجام عملیات، نتیجه را به **JobTracker** باز می‌گرداند. معمولاً در کنار هر **DataNode**، یک **TaskTracker** حضور دارد.

### JobTracker

برنامه‌ی **Map/Reduce** که توسط کاربر نوشته می‌شود به این پرونده تحویل داده می‌شود تا اجرا شود و نتایج توسط این پرونده به کاربر برگردانده می‌شود. برای انجام عملیات، این پرونده کارها را بین **TaskTracker** ها پخش می‌کند، نتیجه کار را از آن‌ها گرفته، ادغام کرده و نتیجه نهایی را به کاربر برمی‌گرداند. این پرونده در صورتی که یک **TaskTracker** با مشکل روبه‌رو شود یا نتواند وظایفش را انجام دهد، وظیفه را به یک **TaskTracker** دیگر سپرده و خطاهای پیش آمده را تا جایی که بتواند رفع می‌کند. اطلاعات مربوط به محل فایل‌ها در خوشه را از **NameNode** دریافت کرده و در اختیار **TaskTracker** ها قرار می‌دهد.

این‌طور می‌توان گفت که در این معماری **Master-Slave**<sup>۱۱</sup>، **NameNode** و **JobTracker** نقش ارباب (**master**) و **DataNode** ها و **TaskTracker** ها نقش برده (**slave**) را ایفا می‌کنند.

<sup>۸</sup> نمونه برداری  
<sup>۹</sup> نگاشت  
<sup>۱۰</sup> کاهش

<sup>۱۱</sup> [https://en.wikipedia.org/wiki/Master/slave\\_\(technology\)](https://en.wikipedia.org/wiki/Master/slave_(technology))

برای ارتباط با **Hadoop** می‌توان از **Java API** استفاده کرد. یک برنامه ساده که عملیات **MapReduce** را بر روی یک فایل کوچک انجام می‌دهد در کنار فایل گزارش قرار دارد. وظیفه این برنامه پیدا کردن فرکانس کلمات در فایل ورودی است. نحوه نصب فریم‌ورک و اجرای برنامه در ادامه آمده است.

#### ۲-۳-۴ پایگاه داده **HBase**

یک پایگاه‌داده‌ی توزیع شده<sup>۱۲</sup>، غیررابطه‌ای<sup>۱۳</sup>، ستون‌گرا<sup>۱۴</sup> و کلید-مقدار<sup>۱۵</sup> است که بر بستر **HDFS** اجرا می‌شود. برای ذخیره‌سازی مقدار بسیار زیادی از داده‌های تنک (**sparse**) همراه با ویژگی **fault tolerant** ساخته شده است. برای عملیات‌های خواندن و نوشتن سریع بر روی دیتاست (**dataset**) های بزرگ همراه با گذردهی (**throughput**) بالا و تاخیر (**latency**) کم بسیار مناسب است. به‌طور خلاصه از ویژگی‌های **HBase** می‌توان به موارد زیر اشاره کرد:

- فشرده‌سازی (**compression**)
- عملیات‌های درون حافظه‌ای<sup>۱۶</sup>
- سازگاری (**consistency**)
- تا حدودی قابلیت تحمل خطا<sup>۱۷</sup>
- ...

راه‌های دسترسی و فرمان‌دادن به **HBase** عبارتند از:

#### • **Java API**

#### • **REST API**<sup>۱۸</sup>

#### • **Apache Avro**<sup>۱۹</sup>

#### • **Apache Thrift**<sup>۲۰</sup>

در این پروژه از **Java API** استفاده می‌کنیم به همین دلیل از توضیح موارد دیگری که در بالا ذکر شد اجتناب می‌کنیم. یک نمونه کد ساده برای اتصال به **HBase** همراه گزارش ضمیمه شده است. در این کد چند جدول در داخل پایگاه‌داده ساخته شده و در آن داده قرار می‌دهیم و داده‌ها را می‌خوانیم و به روز رسانی می‌کنیم. نحوه نصب **HBase** و اجرای کد در ادامه آمده است.

<sup>۱۲</sup>distributed

<sup>۱۳</sup>non relational

<sup>۱۴</sup>column oriented

<sup>۱۵</sup>key-value store

<sup>۱۶</sup>in-memory operations

<sup>۱۷</sup>partially tolerable against failure

<sup>۱۸</sup>[https://en.wikipedia.org/wiki/Representational\\_state\\_transfer](https://en.wikipedia.org/wiki/Representational_state_transfer)

<sup>۱۹</sup>[https://en.wikipedia.org/wiki/Apache\\_Avro](https://en.wikipedia.org/wiki/Apache_Avro)

<sup>۲۰</sup>[https://en.wikipedia.org/wiki/Apache\\_Thrift](https://en.wikipedia.org/wiki/Apache_Thrift)

### ۳-۳-۴ موتور جستجوی Elasticsearch

یک موتور جستجوی متن توزیع‌شده، بلادرنگ (real time)، مقیاس‌پذیر، مستقل از سکو<sup>۲۱</sup> و کاراست که برپایه Apache Lucene<sup>۲۲</sup> ساخته شده است. امکان جستجوی بلادرنگ را فراهم می‌کند و مانند پایگاه داده‌ی MongoDB، سندگرا (document-based) است. انواع داده‌های ساخت‌یافته (structured) و بدون ساختار (unstructured) را پشتیبانی می‌کند. Elasticsearch موارد زیر را پشتیبانی می‌کند:

- توزیع‌پذیری (distribution)

- sharding<sup>۲۳</sup>

- تکثیر (replication)

- خوشه بندی (clustering)

- معماری چندگره‌ای<sup>۲۴</sup>

- عملیات‌های فله‌ای<sup>۲۵</sup>

- و...

اما موارد زیر را پشتیبانی نمی‌کند:

- عملیات‌های نگاشت-کاهش<sup>۲۶</sup>

- معاملات توزیع‌شده<sup>۲۷</sup>

برای برقراری ارتباط با Elasticsearch باید از پروتکل HTTP استفاده کرد. داده‌های ارسالی یا دریافتی در قالب ساختاردهی JSON در قسمت body درخواست (request) یا پاسخ (response) قرار می‌گیرند. برای مثال دستور زیر

```
curl --header "Content-Type: application/json" --request POST
--data '{"title": "Java 8 InDepth", "category": "Java"}'
http://localhost:9200/bookstore/books/1001
```

داده‌ی موجود در قالب JSON را در ایندکس bookstore با تاپ book و ایدی 1001 قرار می‌دهد. ساختار URL (Uniform Resource Locator) ارسالی به Elasticsearch به صورت زیر است:

---

<sup>۲۱</sup>cross-platform

<sup>۲۲</sup>[https://en.wikipedia.org/wiki/Apache\\_Lucene](https://en.wikipedia.org/wiki/Apache_Lucene)

<sup>۲۳</sup>[https://en.wikipedia.org/wiki/Shard\\_\(database\\_architecture\)](https://en.wikipedia.org/wiki/Shard_(database_architecture))

<sup>۲۴</sup>multinode architecture

<sup>۲۵</sup>bulk operations

<sup>۲۶</sup>Map-Reduce operations

<sup>۲۷</sup>distributed transactions

**http://server:port/index(lowercase)/Type/**

یک برنامه ساده برای ارسال داده به **ElasticSearch** و دریافت داده از آن همراه گزارش ضمیمه شده است که شامل کد جاوا و **curl**<sup>۲۸</sup> می‌شود. نحوه نصب **ElasticSearch** و اجرای برنامه در ادامه آمده است.

## ۴-۴ حالات نصب

ابزارهای نام برده شده در قسمت‌های قبل را می‌توان به دو صورت نصب کرد:

### ۱-۴-۴ Standalone mode

در این حالت، یک نمونه (**instance**) از هر کدام از ابزارها را بر روی تنها یک سرور اجرا می‌کنیم. نکته قابل توجه این است که این ابزارها برای سیستم‌های توزیع‌شده و خوشه‌ای ساخته شده‌اند و کارایی بالایی آن‌ها به خاطر استفاده از معماری توزیع‌شده است. با محدود کردن آن‌ها تنها به یک سرور نمی‌توانیم کارایی مطلوب را به دست آوریم. به همین دلیل شرکت دو سرور در اختیار ما قرار داده است تا آن‌ها را به صورتی که در ادامه توضیح می‌دهیم نصب کنیم.

### ۲-۴-۴ Distributed mode

در این حالت در هر سرور در خوشه یک پروسه از ابزار مورد نظر اجرا کرده و پیکربندی مربوط به ارباب (**master**) و برده‌ها (**slaves**) را به صورت مجزا انجام می‌دهیم.

## ۵-۴ نصب ابزارها و اجرای کدهای قسمت‌های قبل

کد کار با ابزارها در پوشه‌های مختلف ضمیمه شده است، منتهی قبل اجرای آن‌ها باید ابزارها نصب باشند. برای نصب آن‌ها به آدرس‌های زیر مراجعه کنید.

### • نصب Hadoop

<https://www.edureka.co/blog/install-hadoop-single-node-hadoop-cluster>

### • نصب HBase

<https://computingforgeeks.com/how-to-install-apache-hadoop-hbase-on-ubuntu>

### • نصب ElasticSearch

<https://www.elastic.co/guide/en/elasticsearch/reference/current/deb.html>

نحوه اجرای کدها در سیستم عامل لینوکس (باید جاوا از قبل نصب باشد):

---

<sup>۲۸</sup><https://curl.haxx.se>



- اجرای نمونه کد کار با **Hadoop**

```
java -jar test-1.0-SNAPSHOT-jar-with-dependencies.jar input.txt output
```

برای مشاهده خروجی:

```
cat output/part-r-00000
```

- اجرای نمونه کد کار با **HBase**

ابتدا باید **Hadoop** و **HBase** اجرا شده باشند:

```
start-all.sh
```

```
start-hbase.sh
```

حال کد را اجرا می‌کنیم:

```
java -jar 1-1.0-SNAPSHOT-jar-with-dependencies.jar
```

- اجرای نمونه کد کار با **ElasticSearch**

```
java -jar 1-1.0-SNAPSHOT-jar-with-dependencies.jar 2> output.txt
```

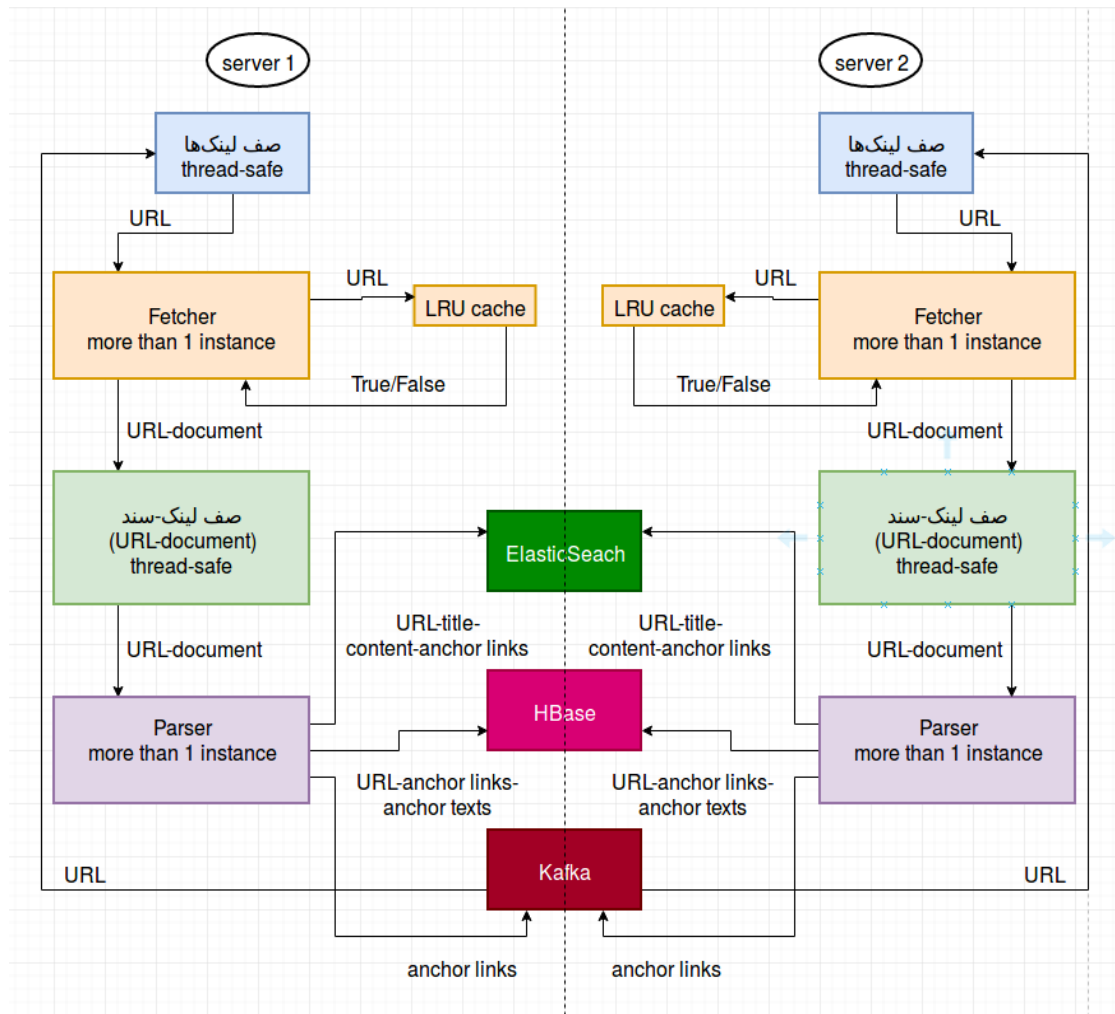
دستورات **curl** نیز در فایل **commands.txt** قرار دارد.

## ۴-۶ نتیجه گیری

در هفته سوم با ابزارها و مفاهیم جدیدی در حوزه داده‌های حجیم آشنا شدیم که از میان آن‌ها می‌توان به ابزارهای **HBase**، **Hadoop** و **ElasticSearch** اشاره کرد. نصب ابزارها و کار با آن‌ها را به صورت مقدماتی فراگرفتیم و معماری نرم‌افزار را در قدم بعدی انجام خواهیم داد.

## فصل پنجم

### معماری پروژه دوم (موتور جستجو)



شکل ۵-۱: معماری نرم‌افزار موتور جستجو

## ۵-۱ مقدمه

در هفته چهارم به طراحی معماری موتور جستجو پرداختیم. همچنین برنامه‌نویسی و پیاده‌سازی بخش‌های مستقل از معماری مانند بخش ارتباط با پایگاه‌داده‌ها انجام شد.

ابتدا هر کدام از افراد تیم معماری مدنظر خود را ارائه داد و در نهایت با ادغام طرح‌های موجود یک معماری جامع انتخاب شد. هر کدام از افراد تیم برنامه‌نویسی و پیاده‌سازی یک بخش از این معماری را بر عهده گرفتند. پیاده‌سازی بخش **HBase** به من واگذار شد. در ادامه اجزای این معماری معرفی و توضیح داده می‌شود.

## ۵-۲ اجزای موتور جستجو

شکل ۵-۱ نمای کلی از معماری موتور جستجو و یک معماری توزیع‌شده (**distributed**) را نشان می‌دهد. دو سرور داریم که دقیقاً عین هم هستند و پایگاه داده‌های آن‌ها به هم متصل است. طبیعتاً پایگاه داده‌ها مثل **Hbase** و موتور جستجوی **ElasticSearch** و پلتفرم **Kafka** را در حالت توزیع شده نصب و اجرا کرده‌ایم.

حال به شرح اجزا در شکل ۵-۱ می‌پردازیم:

## ۵-۲-۱ صف لینک‌ها

این صف شامل لینک‌ها (URL) هایی است که باید بازدید شده و محتوای آن‌ها مورد پردازش قرار گیرند. این صف در ابتدا با URL های زیر به عنوان مقادیر اولیه (seed اولیه) پر می‌شود:

- [https://en.wikipedia.org/wiki/Main\\_Page](https://en.wikipedia.org/wiki/Main_Page)
- <https://us.yahoo.com>
- <https://www.nytimes.com/>
- <https://www.msn.com/en-us/news>
- <http://www.telegraph.co.uk/news/>
- <http://www.alexas.com>
- <http://www.apache.org>
- [https://en.wikipedia.org/wiki/Main\\_Page/World\\_war\\_II](https://en.wikipedia.org/wiki/Main_Page/World_war_II)
- <http://www.news.google.com>
- <https://www.geeksforgeeks.org>
- <https://mvnrepository.com/>

به دلیل این‌که **object** های زیادی از این صف داده می‌خوانند، این صف باید بین نخ‌ها امن<sup>۱</sup> باشد.

## ۵-۲-۲ Fetcher

**object** ای است که به طور متناوب از صف لینک‌ها (در قسمت قبل توضیح داده شد) URL دریافت می‌کند سپس محتوای آن را دانلود می‌کند و URL و محتوا را در صف‌ها قرار می‌دهد.

پیش از دانلود محتوای یک لینک بررسی می‌شود که آیا درخواست به **host** این لینک مجاز است یا خیر. در صورت مجاز نبودن

درخواست، لینک در انتهای صف لینک‌ها قرار می‌گیرد. حال سوالی که مطرح می‌شود این است که معیار مجاز بودن چیست؟

اگر در ۳۰ ثانیه اخیر به یک **host** درخواست داده باشیم، دیگر مجاز نیستیم که به آن درخواست دهیم زیرا ممکن است آدرس IP ما را مسدود (**block**) کند (ممکن است فکر کند برنامه ما یک ربات مولد درخواست برای ایجاد حمله **DOS (Denial Of Service)**

است). برای بررسی این وضعیت از یک حافظه از نوع **LRU (Least Recently Used)** استفاده می‌کنیم.

چندین نمونه از **Fetcher** در برنامه اجرا می‌شوند تا سرعت پردازش و دریافت لینک‌ها افزایش یابد. البته تعداد آن‌ها به مواردی مانند

منابع سیستم (میزان **RAM**، قدرت **CPU**، میزان حجم دیسک ذخیره سازی جانبی (**HDD (Hard Disk Drive)**))، تعداد سایر

نخ‌ها (**thread** ها) و ... بستگی دارد.

<sup>۱</sup>Thread-safe

### ۵-۲-۳ صف اسناد

این صف شامل جفت (pair) های لینک-سند (url-document) است که باید محتوای بخش سند این جفت‌ها تجزیه (parse) شود و اطلاعات زیر از آن‌ها استخراج شود:

- عنوان سند (title)
  - محتوای پاراگراف‌ها یا برچسب (tag) های <p>
  - anchor link ها و عنوان‌های هر کدام
- به دلیل این که object های زیادی از این صف داده می‌خوانند، این صف باید امن<sup>۲</sup> باشد.

### ۵-۲-۴ Parser

object های Parser وظیفه دارند به صورت تناوبی از صف اسناد جفت‌های لینک-سند را خوانده و اطلاعاتی که در بخش صف اسناد توضیح داده شد را استخراج کنند. پس از استخراج عملیات‌های زیر انجام می‌شود:

- قرار دادن لینک سند، محتوای پاراگراف‌های سند و عنوان سند (title) در Elasticsearch
- قرار دادن لینک سند، anchor link ها و عنوان‌های هر کدام در HBase
- تحویل anchor link ها به Kafka

چندین پروسه Parser در برنامه داریم تا سرعت استخراج اطلاعات افزایش یابد. البته تعداد آن‌ها به مواردی مانند منابع سیستم (میزان RAM، قدرت CPU، میزان حجم دیسک ذخیره سازی جانبی (HDD))، تعداد سایر نخ‌ها (thread ها) و ... بستگی دارد.

### ۵-۲-۵ Kafka

Kafka مدیریت صف لینک‌ها را بر عهده دارد. anchor link ها را از Parser گرفته و آن‌هایی که تاکنون مشاهده نشده‌اند را در انتهای صف لینک‌ها قرار می‌دهد. این نوع معماری پردازش BFS (Breadth First Search) برای گراف وب را فراهم می‌کند.

## ۵-۳ نتیجه گیری

در هفته چهارم معماری نرم‌افزار طراحی شد و بخشی از کار برنامه‌نویسی و پیاده‌سازی انجام شد.

---

<sup>۲</sup> Thread-safe

## فصل ششم

### پیاده سازی پروژه دوم (موتور جستجو)

## ۶-۱ مقدمه

در هفته پنجم به پیاده سازی و برنامه نویسی موتور جستجو پرداختیم. اعضای تیم به برنامه نویسی بخش‌های مختلف پروژه پرداختند و در انتها با ترکیب کدها که در بستر **git** ذخیره شده بودند، برنامه‌ی واحدی که همان موتور جستجو است ساخته شد. این برنامه از دو بخش تشکیل شده است:

- خزنده: یک برنامه که با زبان جاوا نوشته شده است و وظیفه‌ی آن خزش (**crawl**) کردن صفحات وب و ذخیره اطلاعات آن است

و

- **search.sh**: یک اسکریپت است که از کاربر **query** گرفته و آن را در بین اسناد جستجو می‌کند و نتیجه را برمی‌گرداند. نتیجه شامل اسنادی است که شامل **query** هستند یا عبارات **query** در آن‌ها حضور دارد. نتایج با توجه به امتیازشان به صورت نزولی مرتب شده و نمایش داده می‌شوند (بهترین سند یافت شده اول نمایش داده می‌شود).

در ادامه به توضیح پکیج‌ها و کلاس‌های پروژه می‌پردازیم.

## ۶-۲ پکیج‌ها

پروژه از پکیج‌های زیر تشکیل شده است:

### ۶-۲-۱ crawler

این پکیج شامل تابع **main** است. در این تابع صف لینک‌ها با **URL** های اولیه پر می‌شود و نخ (**thread**) آمارگیری (کلاس **Statistics**) اجرا می‌شود. همچنین به تعداد نیاز از کلاس‌های **Fetcher** و **Parser** که خود نخ هستند نمونه ساخته شده و اجرا می‌شوند. این پکیج شامل کلاس‌های زیر است: (کلاس‌ها در ادامه توضیح داده خواهند شد).

• **Crawler**

• **Fetcher**

• **Parser**

• **LruCache**

### ۶-۲-۲ index

این پکیج شامل کلاس **Elastic** است که وظیفه برقراری ارتباط با **ElasticSearch** و تبادل داده با آن را برعهده دارد.

### ۶-۲-۳ storage

این پکیج شامل کلاس **HBase** است که وظیفه برقراری ارتباط با **HBase** و تبادل داده با آن را برعهده دارد.

### ۶-۲-۴ test

شامل کلاس‌هایی است که برای تست سایر کلاس‌های نرم‌افزار به کار می‌روند.

## utils ۵-۲-۶

کلاس‌های زیر در این پکیج قرار می‌گیرند:

• **Constant**

• **Pair**

• **Prints**

• **Statistics**

این کلاس‌ها در ادامه توضیح داده می‌شوند.

## language ۶-۲-۶

شامل کلاس‌های مورد نیاز برای تشخیص زبان یک متن است.

## ۳-۶ کلاس‌ها

برای پیاده سازی پروژه از کلاس‌های مختلفی استفاده شده است که در ادامه به تشریح آن‌ها می‌پردازیم:

## Crawler ۱-۳-۶

تابع **main** در این کلاس قرار دارد. عملیات انجام شده در این تابع در قسمت پکیج **crawler** توضیح داده شد.

## Fetcher ۲-۳-۶

این کلاس خود یک نخ است و وظیفه دارد که به طور متناوب از صف لینک‌ها (**URLs queue**)، لینک دریافت کرده و محتوای آن را دانلود کند سپس لینک و محتوایش (**document**) را در صف لینک-سند (**url-document**) قرار دهد. برای انجام این عملیات از توابع زیر استفاده شده است:

### fetchURL

یک لینک از سر صف لینک‌ها برداشته و برمی‌گرداند.

### getDomainIfLruAllowed

بررسی می‌کند که آیا به **host** لینک مورد نظر در ۳۰ ثانیه اخیر درخواستی داده شده است یا خیر. در صورت مثبت بودن، این لینک مجدداً در انتهای صف لینک‌ها قرار می‌گیرد و لینک دیگری اخذ می‌شود.

### fetch

محتوای یک لینک را دانلود کرده و برمی‌گرداند.

### putFetchedData

لینک و محتوای آن را در صف لینک-سند قرار می‌دهد.



### ۳-۳-۶ LruCache

در این کلاس پیاده سازی یک **LruCache** انجام شده است.

### ۴-۳-۶ Parser

این کلاس خود یک نخ است که به طور متناوب از صف لینک-سند داده (لینک و سند مربوط به آن) را استخراج می کند. محتوای لینک را **parse** کرده و موارد زیر را از آن استخراج می کند:

- عنوان لینک (**title**)

- **anchor link** ها و **anchor text** ها

- تگ های **<p>** (پاراگرافها)

پس از استخراج اطلاعات بالا، عملیات های زیر انجام می شود:

- قرار دادن لینک و عنوان آن و پاراگرافها در **ElasticSearch**

- قرار دادن لینک، **anchor link** ها و **anchor text** ها در **HBase**

- تحویل **anchor link** ها به **Kafka**

برای انجام عملیات های بالا از توابع زیر استفاده شده است:

**takeFetchedData**

از سر صف لینک-سند داده برداشته و برمی گرداند.

**extractLinkAnchors**

**anchor link** های یک سند را استخراج می کند و برمی گرداند.

**putToElastic**

لینک و عنوان آن و پاراگراف های سند مربوط به لینک را در **ElasticSearch** قرار می دهد.

**putAnchorsToHBase**

لینک، **anchor link** های درون سند مربوط به لینک و **anchor text** ها را در **HBase** قرار می دهد.

**putToKafka**

**anchor link** های مربوط به یک سند را به **Kafka** تحویل می دهد.

### ۵-۳-۶ Elastic

رابط برنامه و موتور جستجوی **ElasticSearch** است. مهم ترین تابع آن **indexData** است که لینک و عنوان آن و پاراگراف های سند مربوط به آن را دریافت کرده و داخل **ElasticSearch** قرار می دهد.

## ۶-۳-۶ HBase

رابط برنامه و پایگاه داده‌ی HBase است. مهم‌ترین تابع آن **insertLinks** است که لینک، **anchor link** های درون سند مربوط به لینک و **anchor text** ها را گرفته و داخل HBase قرار می‌دهد.

## ۷-۳-۶ Constants

مقادیر ثابت و پارامترهای برنامه در این کلاس قرار دارند.

## ۸-۳-۶ Statistics

کلاسی است که خود یک نخ است و وظیفه دارد به صورت تناوبی وضعیت سایر نخ‌ها ( **Parser** ها و **Fetcher** ها) را گزارش دهد.

## ۹-۳-۶ Pair

پیاده سازی یک ساختار داده برای نگهداری یک کلید و مقدار متناظر آن است.

## ۴-۶ استثناها

موتور جستجوی ساخته شده تا حد بسیار زیادی مطابق با اطلاعات داده شده در این گزارش و گزارش‌های قبلی است و فقط در چند مورد زیر تفاوت دارد که علت آن نیز ذکر شده است:

- عدم استفاده از پکیج تشخیص زبان: به دلیل خطای بالای تشخیص زبان و زمان زیادی هم که می‌گرفت، از کلاس‌های این پکیج استفاده نشد. در اسرع وقت راه‌حلی برای مشکل پیش آمده پیدا می‌کنیم.
- عدم استفاده از **Kafka**: به جای **Kafka** ، **anchor link** ها را مستقیماً در صف لینک‌ها قرار دادیم و عدم مشاهده آن‌ها تاکنون را به یک صف دیگر که مدیریت آن بر عهده HBase است، سپردیم. دلیل عدم استفاده از **Kafka** ، عدم پیکربندی مناسب و از کار افتادن بی مورد آن بود. در اسرع وقت سعی می‌کنیم این مشکل را حل کنیم.

## ۵-۶ اجرای برنامه

برنامه توسط ابزار مدیریت پکیج **maven**<sup>۱</sup> ساخته شده است. نسخه‌ای که در کنار گزارش قرار دارد نسخه‌ای است که بر روی یک سرور اجرا می‌شود. برای مشاهده نسخه اصلی کد که در دو سرور اجرا می‌شود به آدرس زیر مراجعه کنید:

<https://github.com/Ahmad535353/NimboSearchEngine>

برای اجرای بخش خزنده برنامه کافی است روند زیر طی شود:

- اجرای **HBase** ، **Hadoop** و **ElasticSearch**

- ایجاد یک جدول به نام **CrawlerTable** در HBase با فامیلی ستون <sup>۲</sup> **data**

- **java -jar Crawler.jar**

<sup>۱</sup><https://maven.apache.org>

<sup>۲</sup>Column Family

برای ارسال **query** و دریافت پاسخ، کافی است فایل اجرایی **search.sh** را اجرا کنید. با اجرای فایل از شما درخواست می شود تا **query** را وارد کنید. پس از قرار دادن **query** و زدن دکمه **Enter**، پاسخ دریافت و نمایش داده می شود.

## ۶-۶ موارد باقی مانده

به دلیل ضیق وقت، فرصت نشد تا موارد زیر پیاده سازی شوند. سعی می شود در اسرع وقت این موارد پیاده سازی شوند:

- پشتیبانی گرفتن از صف های لینک، لینک-سند و صف کلاس **HBase** و جلوگیری از دست رفتن داده ها ( **data loss** ) هنگام قطع شدن برق.
- فیلتر کردن برخی پروتکل ها مانند **FTP** و ...
- برطرف کردن مشکل **Kafka**
- پیاده سازی و یا پیدا کردن یک ماژول تشخیص دهنده زبان یک متن
- جایگزین کردن کلاس های منسوخ شده ( **deprecated** ) برنامه مانند کلاس های استفاده شده در کلاس **Elastic** و ...

## ۶-۷ نتیجه گیری

با اتمام پیاده سازی، یک موتور جستجو در اختیار ماست که فقط «کار» می کند. نتیجه **query** ها اصلاً خوب نیستند. تیم شرکت به ما توصیه کردند که از الگوریتم **Page Rank**<sup>۳</sup> و مدل برنامه نویسی **Map/Reduce** برای امتیازدهی سایت ها و رتبه بندی آن ها استفاده کنیم و همچنین نواقص قبلی را برطرف کنیم. برای پیاده سازی **Page Rank** از فریم ورک **Spark**<sup>۴</sup> استفاده می کنیم. درباره ی **Spark** و **Page Rank** در فصل بعد توضیح داده می شود.

<sup>۳</sup><https://en.wikipedia.org/wiki/PageRank>

<sup>۴</sup><https://spark.apache.org/>

## فصل هفتم

### بهبود عملکرد موتور جستجو

## ۱-۷ مقدمه

همان‌طور که در قسمت نتیجه‌گیری فصل قبل ذکر شد، نتایج جستجو کیفیت خوبی ندارند. به همین دلیل مدیران تیم و اعضای شرکت پیشنهاداتی برای بهبود نتایج مطرح کردند که شامل موارد زیر می‌شود:

- بهبود جستجو با محاسبه تعداد لینک‌های ورودی

- بهبود جستجو با اعمال متن **anchor link** ها

- بهبود جستجو با **Page Rank**

به دلیل کمبود وقت و صرف وقت برای ارائه نهایی، تنها فرصت شد تا مورد اول را پیاده سازی کنیم. به همین دلیل در موارد دیگر تنها به توضیح بخشی از روند کاری که باید انجام شود می‌پردازیم.

## ۲-۷ بهبود جستجو با محاسبه تعداد لینک‌های ورودی

برای این‌که نتایج جستجو کیفیت بیشتری داشته باشند، باید صفحات مهم‌تر، شانس بیشتری برای ظاهر شدن در نتایج جستجو داشته باشند. ساده‌ترین راه، پیدا کردن تعداد لینک‌ها به یک صفحه است. اگر به یک صفحه، تعداد ارجاعات بیشتری باشد، طبعاً مهم‌تر است. تعداد ارجاعات را با **Map/Reduce** پیدا می‌کنیم، آن‌گاه آن‌ها را به یک امتیاز تبدیل کرده و در کنار امتیازی که **ElasticSearch** برای **query** ها در نظر می‌گیرد به کار برده و در نتیجه‌ی نهایی تاثیر می‌دهیم. برای پیاده‌سازی موارد بالا، دو برنامه نوشته شده است:

- برنامه اول برنامه‌ای است که تعداد لینک‌ها به یک صفحه را محاسبه کرده و در یک جدول در **HBase** قرار می‌دهد. این برنامه همراه گزارش ضمیمه شده است. برای اجرای آن کافی است مراحل زیر طی شوند:

یک **instance** از **HBase** با دستور **start-hbase.sh** اجرا شود.

یک جدول با نام **InnerlinksTable** با ستون فامیلی <sup>۱</sup> **NumOfLinks** ساخته شود.

**java -jar innerlinks-calculator.jar**

- برنامه‌ی دوم برنامه‌ای است که از کاربر **query** می‌گیرد و با در نظر گرفتن تعداد لینک‌های ورودی به هر صفحه و تطابق محتوای هر صفحه با کلمات **query** ورودی، بهترین نتایج را به صورت نزولی نمایش می‌دهد.

برای اجرای این برنامه باید مراحل زیر طی شود:

اجرای **HBase** با دستور **start-hbase.sh**

اجرای **ElasticSearch** با اجرای دستور **sudo systemctl start elasticsearch.service**

**java -jar query-processor-with-innerlinks.jar**

۳-۷ بهبود جستجو با اعمال متن **anchor link**

اگر صفحات عنوان خوبی داشتند، عالی بود ولی خیلی از صفحات یا عنوان خوبی ندارند یا در عنوان‌شان فقط به یکی از جنبه‌های اطلاعاتی که می‌توان در آن صفحه یافت اشاره می‌کنند. حال چگونه می‌توان عناوین خوب تولید کرد؟ پاسخ، متن‌هایی است که صفحات دیگر

<sup>۱</sup>Column Family

برای آن صفحه انتخاب کرده‌اند یا همان متن **anchor link** ها. حال باید با **Map/Reduce** ، **anchor text** های پر تکرار برای هر صفحه را بیابیم (البته نه عبارات بی اثری مانند «این‌جا» یا «لینک») و آن‌ها را در کوثری‌ها تاثیر دهیم تا جستجوهای بهتری داشته باشیم.

## ۴-۷ بهبود جستجو با Page Rank

برای محاسبه‌ی **Page Rank** ، به سراغ فریم‌ورک **Spark** می‌رویم. اما این دو چه هستند؟

• **Page Rank** الگوریتمی است که شرکت گوگل برای امتیازدهی به صفحات وب در موتور جستجویش از آن استفاده می‌کند. این الگوریتم با محاسبه‌ی تعداد لینک‌ها به یک صفحه و کیفیت آن‌ها میزان اهمیت آن صفحه را پیدا کرده و به آن یک عدد نسبت می‌دهد. فرض بر این است که سایت مهم‌تر تعداد ارجاعات بیشتری از طرف سایت‌های دیگر دارد. اطلاعات بیشتر در آدرس زیر:

<https://en.wikipedia.org/wiki/PageRank>

• **Spark**<sup>۲</sup> یک فریم‌ورک برای انجام محاسبات در یک خوشه از کامپیوترها (**cluster**) است. پردازش موازی داده‌ها در خوشه را امکان‌پذیر می‌سازد و به طور ضمنی **fault tolerant** است. کتابخانه‌های زیر از **Spark** استفاده می‌کنند:

MLIB □

GraphX □

Spark SQL □

Spark Streaming □

به دلیل این‌که **Spark** از عملیات‌های درون حافظه‌ای<sup>۳</sup> استفاده می‌کند، تا ۱۰ برابر سریع‌تر از **YARN** عمل می‌کند و این یکی از دلایلی است که به سراغ آن می‌رویم.

## ۵-۷ نتیجه‌گیری

علاوه بر ایده‌های مطرح شده برای بهبود عملکرد موتور جستجو، همچنان ایده‌ها و راه‌های زیادی برای بهتر کردن کیفیت **query** ها وجود دارد. موفق به پیاده سازی تنها یکی از این ایده‌ها شدیم و امید است بتوانیم خارج از تایم کارآموزی به بهبود عملکرد این موتور جستجو بپردازیم.

<sup>۲</sup>[https://en.wikipedia.org/wiki/Apache\\_Spark](https://en.wikipedia.org/wiki/Apache_Spark)

<sup>۳</sup>in-memory operations

## فصل هشتم

### نتیجه گیری

در این دوره با ابزارهای جدیدی مانند **Hadoop** ، **HBase** ، **ElasticSearch** ، **Kafka** و **Spark** آشنا شدیم. همچنین با ابزار مدیریت پکیج **maven** ، ابزار مدیریت نسخه **git** و سیستم عامل لینوکس پروژه را پیاده سازی کردیم و با کتابخانه‌های زیادی در زبان جاوا کار کردیم. از میان معماری‌های موجود برای طراحی نرم‌افزارها، با معماری **client-server** و **publisher-subscriber** آشنا شدیم و دید خوبی نسبت به این معماری‌ها پیدا کردیم. ترکیب مدل برنامه نویسی **Map/Reduce** و فریم‌ورک **Spark** و پردازش داده‌هایی که بر روی سرورهای مختلف قرار داشتند بخش هیجان انگیز کار بود که قبلاً با چنین مسئله‌ای مواجه نشده بودیم. حاصل این دوره کارآموزی، یک موتور جستجو بود که هر چند به قدرتمندی موتور جستجوهای امروزی نیست اما با توجه به میزان زمانی که برای ساخت آن صرف شد، سطح دانش توسعه دهندگان آن، قدرت سخت افزاری سرورها و موارد دیگر، نتایج مطلوبی با جستجوی **query** ها حاصل شد.

همچنین در کنار یادگیری نکات فنی، کار گروهی خوبی را تجربه کردیم و با نمونه کار واقعی آشنا شدیم. همچنان می‌توان بهبودهای زیادی در موتور جستجوی ساخته شده ایجاد کرد که چند مورد در فصل قبل بیان شدند. امیدوارم این گزارش مثمر ثمر واقع شود.



# واژه‌نامه‌ی انگلیسی به فارسی

## A

**agent** ..... عامل نرم‌افزاری یا انسانی

**architecture** ..... معماری

## B

**big data** ..... داده‌ی حجیم

نوع خاصی از صف که همزمان چندین نفر می‌توانند در آن داده بگذارند و چندین نفر داده‌ها را دریافت کنند بدون این

که مشکلی پیش آید. **blocking queue** .....

**block** ..... مسدود کردن

**body** ..... بدنه

**build** ..... ساختن. در این‌جا منظور ساخت فایل **jar** است.

**bulk** ..... توده، فله

## C

**channel** ..... کانال

**cluster** ..... خوشه

**clustering** ..... خوشه بندی

**compression** ..... فشرده سازی

**consistency** ..... سازگاری

**crawl** ..... خزیدن

## D

**data loss** ..... از دست دادن داده

**data science** ..... علم داده

**dataset** ..... مجموعه‌ای از داده‌ها

**deprecated** ..... منسوخ شده

**distributed** ..... توزیع شده

**distribution** ..... توزیع

**document-based** ..... سندگرا

**document** ..... سند

## F

**failure** ..... نقصان، خطا

**fault tolerant** ..... مقاوم در برابر خطا

**fetcher** ..... دریافت کننده

## G

موجودیتی در جاوا است که حافظه اشیایی که دسترسی

به آنها از بین رفته است را آزاد می‌کند **garbage collector**

## H

**hard disk** ..... دیسک سخت

**history** ..... تاریخچه

<b>host</b> . . . . . میزبان	<b>real time</b> . . . . . بلادرنگ
<b>I</b>	
<b>index</b> . . . . . فهرست کردن	<b>reduce</b> . . . . . کاهش دادن
<b>instance</b> . . . . . نمونه	<b>replication</b> . . . . . تکثیر
<b>L</b>	<b>request</b> . . . . . درخواست
<b>latency</b> . . . . . تاخیر	<b>response</b> . . . . . واکنش
مجموعه‌ای از کلاس‌ها و توابع (کدها) که کاربرد خاصی دارند و در کنار هم جمع شده‌اند.	<b>rtm (real time messaging)</b> . . . . . پیام رسانی لحظه ای
<b>library</b> . . . . .	<b>S</b>
<b>load</b> . . . . . بارگیری	<b>slave</b> . . . . . برده
<b>M</b>	<b>snapshot</b> . . . . . کپی
<b>map</b> . . . . . نگاشتن	<b>sparse</b> . . . . . تنک
<b>master</b> . . . . . ارباب	<b>stream</b> . . . . . جریانی از داده‌ها که با سرعت بالا ارسال می‌شوند
<b>O</b>	<b>structured</b> . . . . . ساخت‌یافته
<b>object</b> . . . . . شیء	در اینجا به معنای مشتری یک کانال که به تمام پیام‌های ارسال شده به کانال دسترسی دارد . . . . .
<b>operation</b> . . . . . عملیات	<b>subscriber</b> . . . . .
<b>P</b>	<b>T</b>
<b>pair</b> . . . . . جفت، دوتایی	<b>tag</b> . . . . . در اینجا منظور برچسب‌های فایل HTML است
<b>parser</b> . . . . . تجزیه کننده	<b>thread</b> . . . . . نخ
<b>platform</b> . . . . . سکو، در اینجا منظور سیستم عامل است.	<b>throughput</b> . . . . . گذردهی
<b>private</b> . . . . . خصوصی	<b>title</b> . . . . . عنوان
<b>public</b> . . . . . عمومی	<b>transaction</b> . . . . . معامله
<b>publisher</b> . . . . . انتشاردهنده داده (پیام) در کانال	<b>U</b>
<b>Q</b>	<b>unstructured</b> . . . . . بدون ساختار
<b>queue</b> . . . . . صف	<b>U</b>
<b>R</b>	<b>web page</b> . . . . . صفحه وب