

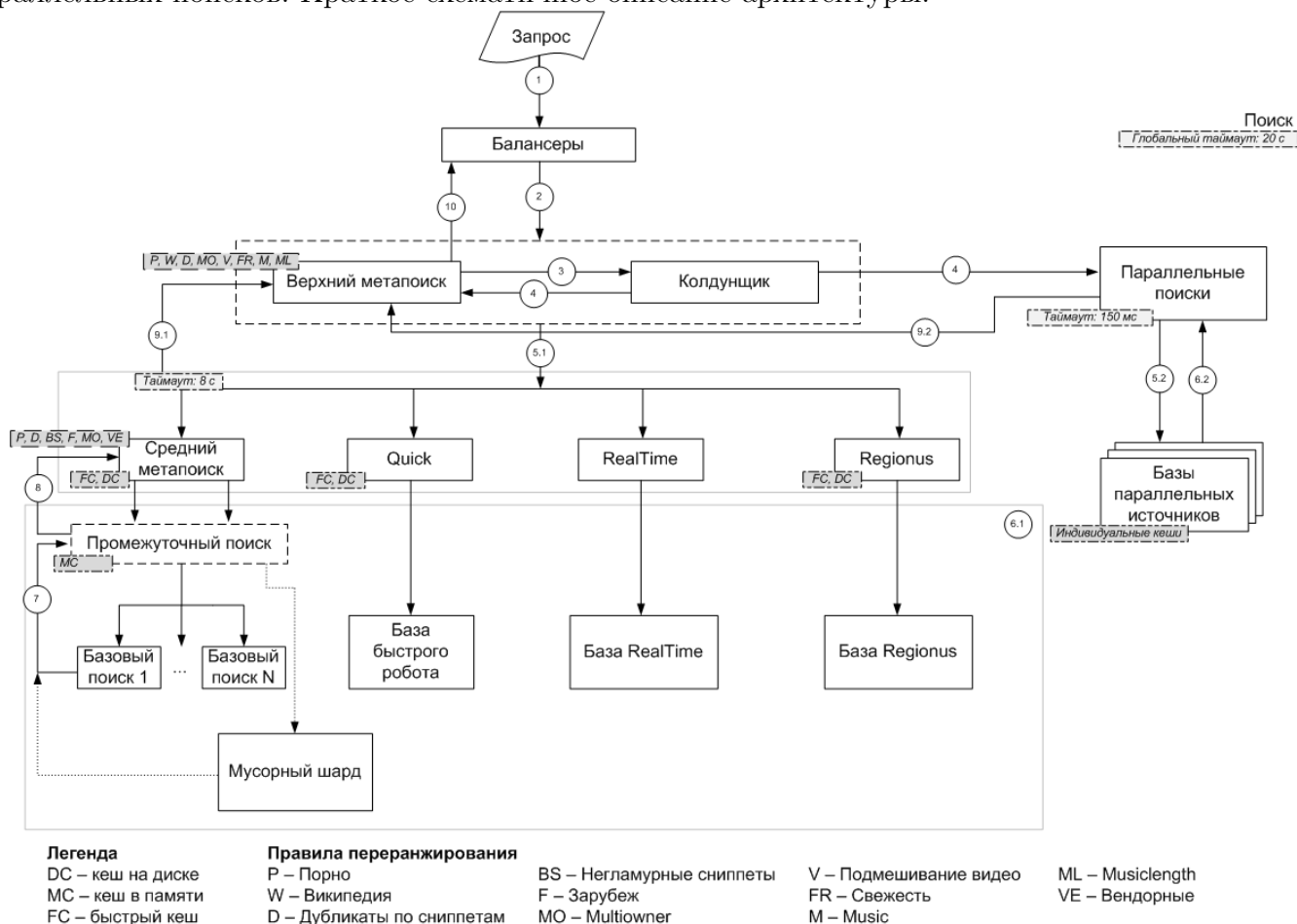
Разработка алгоритма для выделения частовстречающихся шаблонов ошибок из log-файлов программ

# Оглавление

<b>Введение</b>	<b>3</b>
<b>1 Анализ предметной области</b>	<b>5</b>
1.1 Представление задачи в терминах MapReduce . . . . .	5
1.2 Выбор языка программирования и средств разработки . . . . .	5
1.3 Структура данных . . . . .	5
1.4 Стадии выполнения задачи . . . . .	5
1.4.1 Написание программного кода . . . . .	5
1.4.2 Написание функциональных тестов . . . . .	5
<b>2 Имплементация задачи</b>	<b>6</b>
2.1 Диаграмма компонентов . . . . .	6
2.2 Описание способов запуска . . . . .	6
2.3 Анализ полученных результатов . . . . .	6
<b>3 Дальнейшее развитие</b>	<b>7</b>
<b>Заключение</b>	<b>8</b>

# Введение

Большинство программных систем, имеющих сложную структуру и состоящих из нескольких сотен различных компонент, обладают рядом схожих проблем. Например веб-поиск содержит следующие компоненты: балансеры, верхние, средние метапоиски, промежуточные и базовые поиски, колдунчики, антироботы, свежесть, региональные поиск, несколько десятков параллельных поисков. Краткое схематичное описание архитектуры:



Всего одновременно запущено несколько \*\*\*\*\* тысяч инстансов <sup>1</sup> приложений. Каждый инстанс генерирует множество ошибок и записывает каждую из них в log-файл.

Некоторые приложения, близкие по функционалу, пишут в один и тот же файл. Log-файлы ротируются согласно определённому алгоритму. Тем не менее объем log-файла для одного инстанса может достигать нескольких сотен мегабайт, что препятствует быстрому ручному анализу в случае инцидента и инженеры вынуждены тратить ценные секунды на просмотр сотен тысяч строк файла в поисках сообщения с описанием элемента, вызвавшего сбой работы системы.

Начальным требованием к системе для эффективного использования алгоритма является наличие сопоставимого с количеством поисковых приложений количества нод на которых может быть запущена программа, реализующая алгоритм.

В организации, в которой выполнялась учебно-исследовательская работа, развёрнута боль-

<sup>1</sup>инстанс — приложение, запущенное в контейнере и описываемое парой host:port

шая поисковая инфраструктура, которая не лишена недостатков и существует вероятность поломки некоторой её части. Существует множество средств мониторинга состояния веб-поиска и противодействия инцидентам, но в некоторых случаях инженерам их недостаточно и приходится вручную анализировать log-файлы отдельных экземпляров приложений на отдельных host'ах, что, в свою очередь, замедляет скорость реакции на непредвиденную ситуацию. Но даже автоматизация процесса анализа log-файла на одного экземпляра не решает проблему полностью, поэтому необходима возможность быстрого анализа log-файлов сразу множества экземпляров.

Таким образом, целью этой учебно-исследовательской работы является разработка алгоритма, позволяющего собирать статистику по ошибкам, встречающимся в log-файлах экземпляров поисковых приложений на основе существующих шаблонов и выделять новые шаблоны.

# Глава 1

## Анализ предметной области

### 1.1 Представление задачи в терминах MapReduce

Для распределённого запуска приложений была выбрана модель распределённых вычислений MapReduce <sup>1</sup> по ряду причин.

Как оказалось задача без особых сложностей выражается в терминах чистых функций flat map и flat reduce, так как основная структура, используемая в алгоритме — это пара(паттерн<sup>2</sup>, количество совпадений) и благодаря этому однопоточный код легко запускается на множестве нод MapReduce-кластера. Это обстоятельство освобождает от реализации сложного механизма сетевого взаимодействия.

В качестве реализации модели MapReduce была выбрана реализация Yet Another Map Reduce. Код программы, реализующий алгоритм легко переносится на другие реализации данной модели распределённых вычислений.

### 1.2 Выбор языка программирования и средств разработки

### 1.3 Структура данных

### 1.4 Стадии выполнения задачи

#### 1.4.1 Написание программного кода

#### 1.4.2 Написание функциональных тестов

---

<sup>1</sup>MapReduce — модель распределённых вычислений, представленная компанией Google, используемая для параллельных вычислений над очень большими, несколько петабайт, наборами данных в компьютерных кластерах.

<sup>2</sup>паттерн — регулярное выражение

## Глава 2

### Имплементация задачи

2.1 Диаграмма компонентов

2.2 Описание способов запуска

2.3 Анализ полученных результатов

## Глава 3

### Дальнейшее развитие

# Заключение