

# Arenadata™ Hadoop

*Версия - v1.5.1*

**Руководство администратора по HDFS**

# Оглавление

<b>1</b>	<b>ACL на HDFS</b>	<b>3</b>
1.1	Настройка ACL на HDFS	3
1.2	Использование команд CLI для создания ACL	3
<b>2</b>	<b>Примеры ACL</b>	<b>5</b>
2.1	Введение: ACL в сравнении с разрешениями	5
2.2	Пример 1: Предоставление доступа к другой именованной группе	6
2.3	Пример 2: Использование ACL по умолчанию для автоматического применения к дочерним файлам и каталогам	6
2.4	Пример 3: Блокировка доступа конкретного пользователя к подкаталогу	8
<b>3</b>	<b>Особенности ACL для HDFS</b>	<b>9</b>
3.1	Реализация POSIX ACL	9
<b>4</b>	<b>Архивные хранилища</b>	<b>14</b>
4.1	Введение	14
4.2	Типы хранилищ HDFS	14
4.3	Политики хранения: “Hot”, “Warm”, и “Cold”	14
4.4	Настройка архивного хранилища	15
<b>5</b>	<b>Централизованное управление кэшем в HDFS</b>	<b>17</b>
5.1	Обзор	17
5.2	Кэширование	17
5.3	Архитектура кэширования	18
5.4	Терминология кэширования	19
5.5	Настройка централизованного кэширования	19
5.6	Ограничения ОС	21
5.7	Использование Cache Pools и Directives	21
5.8	Команды Cache Pools	22
5.9	Команды Cache Directives	23
<b>6</b>	<b>Настройка HDFS Compression</b>	<b>26</b>
<b>7</b>	<b>Настройка Rack Awareness на ADH</b>	<b>28</b>
7.1	I. Создание скрипта Rack Topology	28
7.2	II. Добавление свойства Script Topology в core-site.xml	29
7.3	III. Перезапуск HDFS и MapReduce	29
7.4	IV. Контроль Rack Awareness	30

---

<b>8</b>	<b>Архивы Hadoop</b>	<b>32</b>
8.1	Введение . . . . .	32
8.2	Компоненты архивов Hadoop . . . . .	33
8.3	Создание архива Hadoop . . . . .	33
8.4	Hadoop Archives и MapReduce . . . . .	35
<b>9</b>	<b>API-интерфейсы JMX Metrics для HDFS Daemons</b>	<b>36</b>
9.1	Использование веб-интерфейса HDFS Daemon . . . . .	36
9.2	Прямой доступ к удаленному агенту JMX . . . . .	36
<b>10</b>	<b>Память в качестве хранилища (техническое превью)</b>	<b>38</b>
10.1	Введение . . . . .	38
10.2	Типы хранилища HDFS . . . . .	39
10.3	Политика хранения LAZY_PERSIST . . . . .	39
10.4	Настройка памяти в качестве хранилища . . . . .	39
10.5	Использование “mover” для применения политик хранения . . . . .	41
<b>11</b>	<b>Запуск DataNodes от Non-root</b>	<b>42</b>
11.1	Введение . . . . .	42
11.2	Настройка DataNode SASL . . . . .	43
<b>12</b>	<b>Режим локального чтения данных на HDFS</b>	<b>45</b>
12.1	Необходимые компоненты . . . . .	45
12.2	Настройка локального чтения данных на HDFS . . . . .	45
<b>13</b>	<b>Руководство администратора по WebHDFS</b>	<b>49</b>

В руководстве приведены сведения по настройке Ambari и Hadoop для Kerberos, расширенные параметры безопасности для Ambari и аутентификация SPNEGO для Hadoop.

Документ может быть полезен администраторам, программистам, разработчикам и сотрудникам подразделений информационных технологий, осуществляющих внедрение и сопровождение кластера.

---

**Important:** Контактная информация службы поддержки – e-mail: [info@arenadata.io](mailto:info@arenadata.io)

---

# Глава 1

## ACL на HDFS

В данном руководстве описывается использование **списков контроля доступа (ACL)** в **распределенной файловой системе Hadoop (HDFS)**. **ACL** расширяет модель разрешения **HDFS** для поддержки более детального доступа к файлам на основе произвольных комбинаций пользователей и групп.

### 1.1 Настройка ACL на HDFS

По умолчанию **ACL** отключены, и при этом **NameNode** отклоняет все попытки установить **ACL**. Например:

```
<property>
  <name>dfs.namenode.acls.enabled</name>
  <value>true</value>
</property>
```

Для включения **ACL** в **HDFS** необходимо в файле *hdfs-site.xml* установить свойству *dfs.namenode.acls.enabled* значение *true*.

### 1.2 Использование команд CLI для создания ACL

В **FsShell** добавляется две новые под-команды: *setfacl* и *getfacl*. Они моделируются после одних и тех же команд **Linux**, но при этом реализуют меньше флагов («flags»). Поддержка дополнительных флагов может быть добавлена позже, если потребуется.

- **Setfacl**. Устанавливает ACL для файлов и каталогов. Применение:

```
-setfacl [-bkR] {-m|-x} <acl_spec> <path> -setfacl --set <acl_spec> <path>
```

Функции команды приведены в таблице.

Таблица 1.1.: Функции команды setfacl

Функция	Описание
-b	Удаление всех записей, но с сохранением записей ACL. Записи для пользователей и групп сохраняются для совместимости с разрешениями
-k	Удаление ACL по умолчанию
-R	Применение операции ко всем файлам и каталогам рекурсивно
-m	Изменение ACL. Новые записи добавляются в ACL, а существующие записи сохраняются
-x	Удаление указанных записей ACL. Все остальные записи ACL сохраняются
-set	Полная замена ACL и сброс всех существующих записей. «acl_spec» включает записи для пользователей и групп для совместимости с разрешениями
<acl_spec>	Список записей ACL, разделенных запятыми
<path>	Путь к файлу или директории для изменения

Например:

```
hdfs dfs -setfacl -m user:hadoop:rw- /file
hdfs dfs -setfacl -x user:hadoop /file
hdfs dfs -setfacl -b /file
hdfs dfs -setfacl -k /dir
hdfs dfs -setfacl --set user::rw-,user:hadoop:rw-,group::r--,other::r-- /file
hdfs dfs -setfacl -R -m user:hadoop:r-x /dir
hdfs dfs -setfacl -m default:user:hadoop:r-x /dir
```

Код выхода:

При успехе 0 и ненулевое значение при ошибке.

- **Getfacl.** Отображает ACL файлов и каталогов. Если каталог имеет ACL по умолчанию, *getfacl* также его отображает. Применение:

```
-getfacl [-R] <path>
```

Функции команды приведены в таблице .

Таблица 1.2.: Функции команды getfacl

Функция	Описание
-R	Список ACL всех рекурсивных файлов и каталогов
<path>	Путь к файлу или директории списка

Например:

```
hdfs dfs -getfacl /file
hdfs dfs -getfacl -R /dir
```

Код выхода:

При успехе 0 и ненулевое значение при ошибке.

## Глава 2

# Примеры ACL

### 2.1 Введение: ACL в сравнении с разрешениями

До реализации **списков контроля доступа (ACL)** модель разрешения **HDFS** была эквивалентна традиционным **UNIX-разрешениям**. В этой модели разрешения для каждого файла или каталога управляются набором из трех различных пользовательских классов: “Владелец”, “Группа” и “Другие”. Для каждого пользовательского класса существует три разрешения: чтение, запись и выполнение. Таким образом, для любого объекта файловой системы его разрешения могут быть закодированы в  $3 \times 3 = 9$  бит. Когда пользователь пытается получить доступ к объекту файловой системы, **HDFS** применяет разрешения в соответствии с конкретным классом пользователя, применимым к нему. Если пользователь является владельцем, **HDFS** проверяет разрешения класса “Владелец”. Если пользователь не является владельцем, но является членом группы объектов файловой системы, **HDFS** проверяет разрешения класса “Группа”. В противном случае **HDFS** проверяет разрешения класса “Другие”.

Эта модель может в достаточной степени удовлетворять большому количеству требований безопасности. Например, рассмотрим отдел продаж, который хотел бы, чтобы один пользователь – Брюс, менеджер отдела, – контролировал все изменения данных продаж. Другим сотрудникам отдела продаж необходимо просмотреть данные, но не следует изменять их. Все остальные компании (за пределами отдела продаж) не должны иметь возможность просматривать данные. Это требование может быть реализовано путем запуска `chmod 640` в файле со следующим результатом:

```
-rw-r-----1 brucesales22K Nov 18 10:55 sales-data
```

Только Брюс может изменить файл, только члены группы продаж могут прочитать файл, и никто другой не может получить доступ к файлу каким-либо образом.

Предположим, что возникают новые требования. Отдел продаж вырос, и Брюс не может контролировать все изменения в файле. Новое требование состоит в том, что Брюсу, Диане и Кларку разрешено вносить изменения. К сожалению, для реализации этого требования не существует возможности для разрешений, потому что может быть только один владелец и одна группа, и группа уже используется для реализации требования только для чтения для команды продаж. Типичным обходным решением является установка владельца файла на мнимую учетную запись пользователя, такую как `salesmgr`, и разрешение Брюсу, Диане и Кларку использовать учетную запись `salesmgr` с помощью `sudo` или аналогичных механизмов. Недостатком этого обходного пути является то, что он создает сложность для конечных пользователей, требуя от них использовать разные учетные записи для разных действий.

Предположим теперь, что помимо сотрудников по продажам все руководители компании должны иметь возможность читать данные о продажах. Это еще одно требование, которое не может быть выражено с помощью разрешенных битов, поскольку существует только одна группа, и она уже используется в продажах. Типичным обходным решением является установка группы файлов в новую мнимую группу, такую как `salesandexecs`, и добавление всех пользователей `sales` и всех пользователей `execs` к этой группе. Недостатком этого обходного

пути является то, что он требует от администраторов создание и управление дополнительными пользователями и группами.

На основании примера можно увидеть, что может быть неудобно использовать **Permission Bits** для удовлетворения требований к разрешению, которые отличаются от естественной организационной иерархии пользователей и групп.

Преимущество использования **ACL** заключается в том, что он позволяет более естественным образом решать эти требования, поскольку для любого объекта файловой системы несколько пользователей и несколько групп могут иметь разные наборы разрешений.

## 2.2 Пример 1: Предоставление доступа к другой именованной группе

Для решения одного из вопросов, затронутых в предыдущем разделе, установим **ACL**, который предоставляет доступ к данным о доступе к чтению членам группы *execs*.

- Установить **ACL**:  

```
> hdfs dfs -setfacl -m group:execs:r-- /sales-data
```
- Запустить *getfacl*, чтобы проверить результаты:

```
> hdfs dfs -getfacl /sales-data
# file: /sales-data
# owner: bruce
# group: sales
user::rw-
group::r--
group:execs:r--
mask::r--
other::---
```

- Если запустить команду *ls*, можно увидеть, что перечисленные разрешения были добавлены с символом “+” для обозначения наличия **ACL**. Символ “+” добавляется к разрешениям любого файла или каталога с **ACL**.

```
> hdfs dfs -ls /sales-data
Found 1 items
-rw-r-----+ 3 bruce sales          0 2014-03-04 16:31 /sales-data
```

Новая запись **ACL** добавляется к существующим разрешениям, определенным в разрешенных битах. Как владелец файла, Брюс имеет полный контроль. Члены группы *sales* или *execs* имеют доступ на чтение. У остальных нет доступа.

## 2.3 Пример 2: Использование **ACL** по умолчанию для автоматического применения к дочерним файлам и каталогам

В дополнение к **ACL**, выполняемому проверки во время разрешений, существует также отдельная концепция **ACL по умолчанию**. **ACL по умолчанию** может применяться только к каталогу, а не к файлу. **ACL по умолчанию** не имеют прямого влияния на проверки разрешений для существующих дочерних файлов и каталогов, но вместо этого определяют **ACL**, которые будут получать новые дочерние файлы и каталоги при их создании.

Предположим, есть каталог “monthly-sales-data”, который далее подразделяется на отдельные каталоги для каждого месяца. Установим **ACL по умолчанию**, чтобы гарантировать, что члены группы *execs* автоматически получают доступ к новым подкаталогам по мере их создания каждый месяц.



- Установить ACL по умолчанию в родительский каталог:  
> `hdfs dfs -setfacl -m default:group:execs:r-x /monthly-sales-data`
- Создать подкаталоги:  
> `hdfs dfs -mkdir /monthly-sales-data/JAN`  
> `hdfs dfs -mkdir /monthly-sales-data/FEB`
- Убедиться, что HDFS автоматически применил ACL по умолчанию в подкаталоги:

```
> hdfs dfs -getfacl -R /monthly-sales-data
# file: /monthly-sales-data
# owner: bruce
# group: sales
user::rwx
group::r-x
other:---
default:user::rwx
default:group::r-x
default:group:execs:r-x
default:mask::r-x
default:other:---

# file: /monthly-sales-data/FEB
# owner: bruce
# group: sales
user::rwx
group::r-x
group:execs:r-x
mask::r-x
other:---
default:user::rwx
default:group::r-x
default:group:execs:r-x
default:mask::r-x
default:other:---

# file: /monthly-sales-data/JAN
# owner: bruce
# group: sales
user::rwx
group::r-x
group:execs:r-x
mask::r-x
other:---
default:user::rwx
default:group::r-x
default:group:execs:r-x
default:mask::r-x
default:other:---
```

**ACL по умолчанию** копируется из родительского каталога в дочерний файл или каталог при его создании. Последующие изменения **ACL по умолчанию** в родительском каталоге не изменяют **ACL** существующих дочерних элементов.

## 2.4 Пример 3: Блокировка доступа конкретного пользователя к подкаталогу

Предположим, что необходимо немедленно заблокировать доступ ко всему подкаталогу для конкретного пользователя. Применение к данному пользователю **ACL** в корне этого подкаталога является самым быстрым способом без риска случайного отзыва разрешений у других пользователей.

1. Добавить запись **ACL** для блокировки всего доступа пользователя Диана к “monthly-sales-data”:

```
> hdfs dfs -setfacl -m user:diana:--- /monthly-sales-data
```

2. Запустить *getfacl* для проверки результатов:

```
> hdfs dfs -getfacl /monthly-sales-data
# file: /monthly-sales-data
# owner: bruce
# group: sales
user::rwx
user:diana:---
group::r-x
mask::r-x
other:---
default:user::rwx
default:group::r-x
default:group:execs:r-x
default:mask::r-x
default:other:---
```

Новая запись **ACL** добавляется к существующим разрешениям, определенным в **Permission Bits**. Брюс имеет полный контроль как владелец файла. Члены группы *sales* или *execs* имеют доступ на чтение. У остальных нет доступа.

Важно помнить о порядке оценки записей **ACL**, когда пользователь пытается получить доступ к объекту файловой системы:

- Если пользователь является владельцем файла, применяются разрешения “Владелец”;
- Если у пользователя есть запись **ACL**-пользователя, применяются соответствующие права;
- Если пользователь является членом группы файлов или любой именованной группы в **ACL**, то для всех соответствующих записей принудительно объединяются разрешения (пользователь может быть членом нескольких групп);
- Если ничто из вышеуказанного не применимо, назначаются разрешенные биты класса “Другие”.

В данном примере запись **ACL**-пользователя достигла установленной цели, поскольку пользователь не является владельцем файла, а именованная пользовательская запись имеет приоритет над всеми другими записями.

## Глава 3

# Особенности ACL для HDFS

### 3.1 Реализация POSIX ACL

Списки ACL на HDFS реализуются с помощью модели ACL POSIX. Если вы когда-либо использовали POSIX ACL в файловой системе Linux, ACL на HDFS работает также.

#### 3.1.1 Совместимость и применение

HDFS может связывать дополнительный ACL с любым файлом или каталогом. Все операции HDFS, которые обеспечивают соблюдение разрешений, выраженных с помощью Permission Bits, также должны обеспечить любой ACL, который определен для файла или каталога. Любая существующая логика, которая обходит Permission Bits, также обходит и ACL. Включая супер-пользователя HDFS и установку *false* в конфигурации *dfs.permissions*.

#### 3.1.2 Доступ через ориентированные на пользователя конечные точки

HDFS поддерживает операции по настройке и получению ACL, связанного с файлом или каталогом. Эти операции доступны через несколько ориентированных на пользователя конечных точек. Эти конечные точки включают FsShell CLI, программную манипуляцию через классы FileSystem и FileContext, WebHDFS и NFS.

#### 3.1.3 Обратная связь пользователя: индикатор CLI для ACL

К перечисленным разрешениям любого файла или каталога с соответствующим ACL добавляется символ “+” (результат команды *ls -l*).

#### 3.1.4 Обратная совместимость

Реализация ACL обратно совместима с существующим применением разрешенных битов. Изменения, применяемые посредством Permission Bits (то есть *chmod*), также отображаются как изменения в ACL. Аналогично, изменения, применяемые к записям ACL для базовых классов пользователей (“Владелец”, “Группа” и “Другие”), также отображаются в виде изменений в Permission Bits. Другими словами, операции Permission Bits и ACL управляют совместно используемой моделью, и операции Permission Bits можно рассматривать как подмножество операций ACL.

#### 3.1.5 Низкие накладные расходы

Добавление ACL не приводит к негативному влиянию на потребление системных ресурсов при развертывании, чтобы отказаться от использования ACL. Он включает в себя процессор, память, диск и

пропускную способность сети.

Использование **ACL** влияет на производительность **NameNode**. Поэтому рекомендуется использовать **Permission Bits**, если это возможно, прежде чем использовать **ACL**.

### 3.1.6 Ограничения ACL

Количество записей в одном **ACL** ограничено максимум до 32. Попытки добавить записи **ACL** сверх максимума выполняются с ошибкой, обращенной к пользователю. Это делается по двум причинам: упростить управление и ограничить потребление ресурсов.

Списки **ACL** с очень большим количеством записей, как правило, трудно понять и могут указывать на то, что требования лучше устраняются путем определения дополнительных групп или пользователей. **ACL** с очень большим количеством записей также требуют большей памяти и хранилища, и для каждой проверки разрешений требуется больше времени.

Число 32 соответствует максимальному количеству записей **ACL**, принудительно используемых семейством файловых систем *ext*.

### 3.1.7 Символические ссылки

У символических ссылок нет собственных списков **ACL**. **ACL** символической ссылки всегда рассматривается как разрешения по умолчанию (777 в **Permission Bits**). Операции, которые изменяют **ACL** символической ссылки, вместо этого изменяют саму символическую ссылку.

### 3.1.8 Снапшоты

При создании снапшота все списки **ACL** блокируются. Изменения в **ACL** в момент создания снапшота не фиксируются.

### 3.1.9 Инструментарий

Инструмент, который распространяет **Permission Bits**, не распространяет **ACL**. **ACL** включается командой `cp -p` и `distcp -p`.

### 3.1.10 Варианты использования ACL на HDFS

**ACL** на **HDFS** поддерживает следующие варианты использования:

#### 3.1.11 Доступ нескольким пользователям

В данном разделе описывается случай, когда нескольким пользователям требуется доступ для чтения к файлу. При этом ни один из пользователей не является владельцем файла. И пользователи не являются членами общей группы, поэтому невозможно использовать групповые разрешения.

В данном случае устанавливается **ACL** доступ, содержащий несколько именованных пользовательских записей:

```
ACLs on HDFS supports the following use cases:
```

#### 3.1.12 Доступ нескольким группам

В данном разделе описывается случай, когда нескольким группам требуется чтение и запись в файл. При этом нет группы, объединяющей всех необходимых пользователей, поэтому невозможно использовать групповые разрешения.

В данном случае устанавливается **ACL** доступ, содержащий несколько именованных групповых записей:

```
group:sales:rw-
group:execs:rw-
```

### 3.1.13 Hive Partitioned Tables

В данном случае **Hive** содержит секционированную таблицу данных о продажах. Ключ раздела – *country*. **Hive** сохраняет секционированные таблицы с помощью отдельного подкаталога для каждого определенного значения ключа раздела, поэтому структура файловой системы в **HDFS** выглядит так:

```
user
|-- hive
  |-- warehouse
    |-- sales
      |-- country=CN
      |-- country=GB
      |-- country=US
```

Группа *salesadmin* – это группа для всех этих файлов. Члены группы имеют доступ на чтение и запись ко всем файлам. Отдельные группы, зависящие от конкретной страны, могут запускать запросы на использование, которые только считывают данные для конкретной страны, например, *sales\_CN*, *sales\_GB* и *sales\_US*. У этих групп нет доступа на запись.

Этот вариант использования можно решить, установив **ACL** доступ в каждом подкаталоге, содержащем запись собственной группы и именованной группы:

```
country=CN
group::rwx
group:sales_CN:r-x

country=GB
group::rwx
group:sales_GB:r-x

country=US
group::rwx
group:sales_US:r-x
```

---

**Important:** Функциональность записи ACL группы владельца (запись группы без имени) эквивалентна установленным Permission Bits

---

Хранение в **Hive** в настоящее время не учитывает разрешения **ACL** в **HDFS**. Скорее, он проверяет доступ с использованием традиционной модели разрешений **POSIX**.

### 3.1.14 ACL по умолчанию

В данном случае администратор файловой системы или владелец поддерева хотел бы определить политику доступа, которая будет применяться ко всему поддереву. Эта политика доступа должна применяться не только к текущему набору файлов и каталогов, но также к любым новым файлам и каталогам, которые будут добавляться позже.

Этот вариант использования можно решить, установив **ACL по умолчанию** в каталог. **ACL по умолчанию** может содержать любую произвольную комбинацию записей. Например:

```
default:user::rwx
default:user:bruce:rw-
```

```
default:user:diana:r--
default:user:clark:rw-
default:group:r--
default:group:sales:rw-
default:group:execs:rw-
default:others:---
```

Важно отметить, что **ACL по умолчанию** копируется из каталога во вновь созданные дочерние файлы и каталоги во время их создания. Если изменить **ACL по умолчанию** в каталоге, это не влияет на **ACL** файлов и подкаталогов, которые уже существуют в каталоге. **ACL по умолчанию** никогда не рассматриваются во время исполнения разрешения. Они используются только для определения **ACL**, какие новые файлы и подкаталоги будут автоматически получены при их создании.

### 3.1.15 Минимальные ACL/Permissions

**ACL** на **HDFS** поддерживают развертывания, использующие **Permission Bits**, а не **ACL** с именованными пользователями и группами. Биты разрешения эквивалентны минимальному **ACL**, содержащему только 3 записи. Например:

```
user:rw-
group:r--
others:---
```

### 3.1.16 Блокировка доступа к поддереву для конкретного пользователя

В данном примере создано глубокое вложенное подэлемента файловой системы, читаемое во всем мире, к которому устанавливается требование блокировать доступ для конкретного пользователя ко всем файлам этого поддерева.

Этот пример можно решить, установив **ACL** в корневое дерево, с именованной пользовательской записью, которая удаляет весь доступ пользователя. Для этой файловой системы:

```
dir1
|-- dir2
    |-- dir3
        |-- file1
        |-- file2
        |-- file3

dir1|-- dir2|-- dir3|-- file1|-- file2|-- file3
```

Установка **ACL** на *dir2* блокирует доступ для Брюса к *dir3*, *file1*, *file2* и *file3*:

```
user:bruce:---
```

Удаление разрешений на выполнение на *dir2* означает, что Брюс не может получить доступ к *dir2* и, следовательно, не может видеть ни один из его дочерних элементов. Это также означает, что доступ блокируется автоматически для любых вновь добавленных файлов под *dir2*. То есть если *file4* создается под *dir3*, Брюс не сможет получить к нему доступ.

### 3.1.17 ACL с Sticky Bit

В данном случае нескольким именованным пользователям или группам требуется полный доступ к каталогу общего назначения, например, / *tmp*. Однако разрешения “Write” и “Execute” в каталоге также дают пользователям возможность удалить или переименовать любые файлы в каталоге, даже файлы, созданные другими пользователями. Разрешения пользователей необходимо ограничить, чтобы у них был допуск на удаление или переименование созданных только ими файлов.

Этот случай можно решить, объединив **ACL** с **Sticky bit**. **Sticky bit** – это функциональность, которая в настоящее время работает с **Permission Bits**. И она работает в сочетании с **ACL**.

## Глава 4

# Архивные хранилища

В данной главе описывается использование политик хранения файлов и каталогов, предназначенных для архивного хранения.

### 4.1 Введение

Архивные хранилища позволяют хранить данные на физических носителях с высокой плотностью хранения и низкими ресурсами обработки.

Для реализации архивного хранилища необходимо выполнить следующие действия:

- Выключить DataNode;
- Назначить тип хранилища ARCHIVE DataNodes, предназначенный для архивного хранения;
- Установить политики хранения “HOT”, “WARM” или “COLD” в файлах и каталогах HDFS;
- Перезапустить DataNode.

Для обновления параметра политики хранения в файле или каталоге, необходимо использовать инструмент переноса данных **HDFS** для перемещения блоков, как указано в новой политике хранения.

### 4.2 Типы хранилищ HDFS

Типы хранилищ **HDFS** могут использоваться для данных, предназначенных различным типам физических носителей.

Доступны следующие типы хранилищ:

- DISK – дисковое хранилище (тип хранилища по умолчанию);
- ARCHIVE – архивные хранилища (высокая плотность хранения, низкие ресурсы обработки);
- SSD – твердотельный накопитель;
- RAM\_DISK – память DataNode.

Если тип хранения не назначен, по умолчанию используется тип *DISK*.

### 4.3 Политики хранения: “Hot”, “Warm”, и “Cold”

Данные можно хранить на дисках типа *DISK* или *ARCHIVE*, используя следующие предварительно настроенные политики хранения:



- *HOT* – используется как для хранения, так и для вычисления. Данные, которые используются для обработки, остаются в этой политике. Все копии хранятся на *DISK*. Для создания резервной копии памяти нет, и для их хранения используется *ARCHIVE*;
- *WARM* – частично “HOT” и частично “COLD”. При “WARM” первая копия хранится на *DISK*, а остальные хранятся в *ARCHIVE*. Резервным хранилищем для создания и копирования является *DISK* или *ARCHIVE*, в случае если *DISK* недоступен;
- *COLD* – используется только для хранения, с ограниченным расчетом. Данные, которые больше не используются или которые необходимо заархивировать, переносятся из хранилища “HOT” в “COLD”. При “COLD” все копии хранятся в *ARCHIVE*, и нет резервного хранилища для создания или копирования.

В таблице приведена политика копирования.

Таблица 4.1.: Политика копирования

ID политики	Название политики хранения	Место размещения копии (для n копий)	Резервное хранилище для разработки	Резервное хранилище для копий
12	HOT	DISK: n	<нет>	ARCHIVE
8	WARM	DISK: 1, ARCHIVE: n-1	DISK, ARCHIVE	DISK, ARCHIVE
4	COLD	ARCHIVE: n	<нет>	<нет>

---

**Important:** В настоящее время политики хранения нельзя редактировать

---

## 4.4 Настройка архивного хранилища

Для настройки архивного хранилища необходимо выполнить следующие действия:

### 1. Выключить **DataNode**

Закрывать **DataNode** с помощью соответствующих команд.

### 2. Назначить тип хранения **ARCHIVE** в **DataNode**

Для назначения типа хранения *ARCHIVE* для **DataNode** можно использовать свойство *dfs.name.dir* в файле */etc/hadoop/conf/hdfs-site.xml*.

Свойство *dfs.name.dir* определяет, где в локальной файловой системе **DataNode** хранит свои блоки.

Чтобы назначить **DataNode** как хранилище *DISK*, необходимо использовать путь к локальной файловой системе, как обычно. Поскольку *DISK* является типом памяти по умолчанию, больше ничего не требуется. Например:

```
<property>
  <name>dfs.data.dir</name>
  <value>file:///grid/1/tmp/data_trunk</value>
</property>
```

Чтобы назначить **DataNode** как хранилище *ARCHIVE*, необходимо добавить *[ARCHIVE]* в начало пути локальной файловой системы. Например:

```
<property>
  <name>dfs.data.dir</name>
  <value>[ARCHIVE]file:///grid/1/tmp/data_trunk</value>
</property>
```

### 3. Установка и получение политики хранения

Необходимо установить политику хранения файла или каталога. Команда:

```
hdfs dfsadmin -setStoragePolicy <path> <policyName>
```

Аргументы:

- *<path>* – путь к каталогу или файлу;
- *<policyName>* – название политики хранения.

Пример:

```
hdfs dfsadmin -setStoragePolicy /cold1 COLD
```

Получение политики хранения файла или каталога осуществляется по команде:

```
hdfs dfsadmin -getStoragePolicy <path>
```

Аргументы:

- *<path>* – путь к каталогу или файлу.

Пример:

```
hdfs dfsadmin -getStoragePolicy /cold1
```

### 4. Запуск DataNode

Запустить **DataNode** с помощью соответствующих команд.

### 5. Использовать “mover” для применения политик хранения

При обновлении параметра политики хранения в файле или каталоге, новая политика не применяется автоматически. Необходимо использовать инструмент переноса данных **HDFS mover** для фактического перемещения блоков (как указано в новой политике хранения).

Средство миграции данных *mover* сканирует выбранные файлы в **HDFS** и проверяет, соответствует ли размещение блоков политике хранения. Копии блоков, нарушающих политику хранения, он перемещает в соответствующий тип хранилища для выполнения требований политики.

Команда:

```
hdfs mover [-p <files/dirs> | -f <local file name>]
```

Аргументы:

- *-p <files/dirs>* – список файлов / каталогов HDFS для переноса, разделенные пробелами;
- *-f <local file>* – локальный файл, содержащий список файлов / каталогов HDFS для миграции.

---

**Important:** Если оба параметра *-p* и *-f* опущены, путь по умолчанию является корневым каталогом

---

Пример:

```
hdfs mover /cold1/testfile
```

## Глава 5

# Централизованное управление кэшем в HDFS

В данной главе приведены инструкции по настройке и использованию централизованного управления кэшем в **HDFS**. Централизованное управление кэшем позволяет указать пути к каталогам или файлам, которые будут кэшироваться **HDFS**, тем самым повышая производительность приложений, которые неоднократно обращаются к одним и тем же данным.

### 5.1 Обзор

Централизованное управление кэшем в **HDFS** – это точный механизм кэширования, который позволяет указывать пути к каталогам или файлам, которые будут кэшироваться **HDFS**. **NameNode** связывается с **DataNodes**, которые имеют требуемые доступные на диске блоки, и поручает **DataNodes** кэшировать блоки в кэшах.

Централизованное управление кэшем в **HDFS** дает много существенных преимуществ:

- Точный кэш предотвращает вытеснение часто используемых данных из памяти. Это особенно важно, когда размер рабочего набора превышает размер оперативной памяти, который является общим для многих рабочих нагрузок **HDFS**;
- Поскольку кэши данных **DataNode** управляются **NameNode**, приложения могут запрашивать набор местоположений кэшированных блоков при принятии решений о размещении задач. Совмещение задачи с кэшированной блочной копией повышает производительность чтения;
- Когда блок кэшируется с помощью **DataNode**, клиенты могут использовать новый, более эффективный API-интерфейс с нулевой копией. Поскольку проверка контрольных сумм кэшированных данных выполняется **DataNode** один раз, при использовании этого нового API клиенты могут понести практически нулевые расходы;
- Централизованное кэширование может улучшить общую эффективность использования памяти кластера. Когда используется ОС буферного кэша на каждом **DataNode**, повторные чтения блока приводят к тому, что все  $\langle n \rangle$  копии блока перемещаются в буферный кэш. При централизованном управлении кэшем точно указывается только  $\langle m \rangle$  копий  $\langle n \rangle$ , тем самым сохраняя память  $\langle n - m \rangle$ .

### 5.2 Кэширование

Централизованное управление кэшем полезно:

- Для файлов, к которым неоднократно обращаются. Например, небольшая таблица фактов в Hive, которая часто используется, является хорошим кандидатом для кэширования. И наоборот, кэширование запроса годовой отчетности менее полезно, так как подобные данные, вероятно, могут быть прочитаны только один раз;
- Для смешанной рабочей нагрузки с SLA-производительностью. Кэширование рабочего набора высокоприоритетной рабочей нагрузки гарантирует, что она не конкурирует с низкоприоритетными рабочими нагрузками для дискового ввода-вывода.

### 5.3 Архитектура кэширования

На рисунке показана централизованная кэшированная архитектура управления (Рис.5.1.).

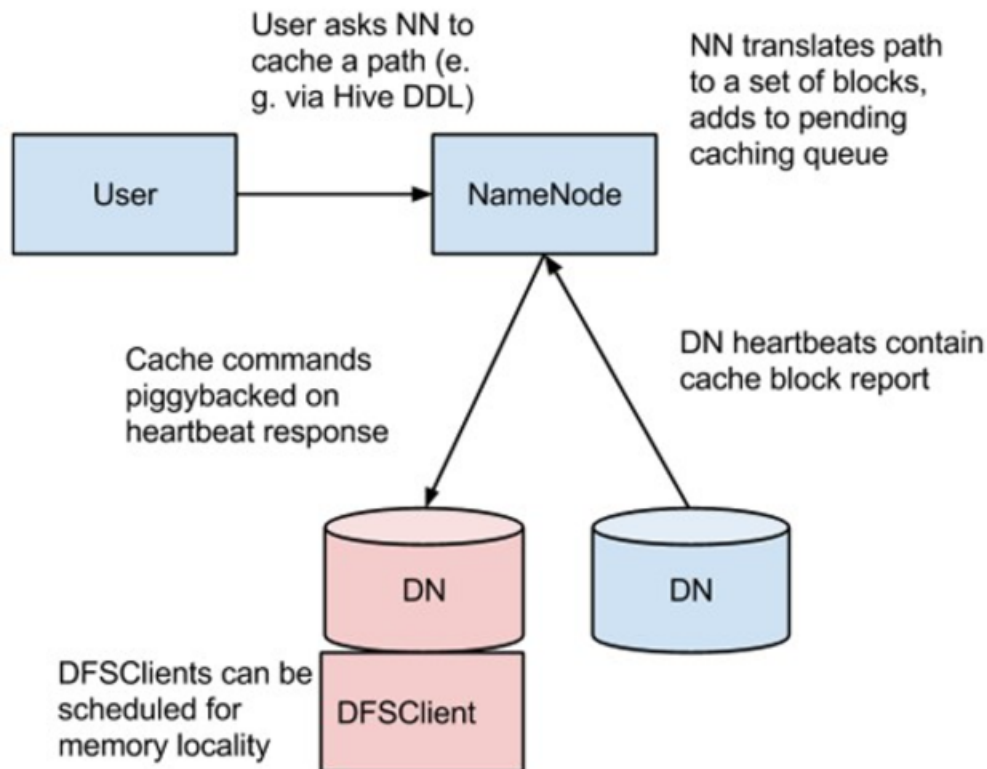


Рис.5.1.: Архитектура кэширования

В данной архитектуре **NameNode** отвечает за координацию всех кэш-файлов с неактивными данными в кластере. **NameNode** периодически получает кэш-отчет от каждого **DataNode**. Кэш-отчет описывает все блоки, кэшированные в **DataNode**. **NameNode** управляет кэшами **DataNode** с помощью кэш-копий и мгновенных команд *uncache*.

**NameNode** запрашивает набор **Cache Directives**, чтобы определить, какие контуры следует кэшировать. **Cache Directives** постоянно сохраняются в *fsimage* и журналах и могут быть добавлены, удалены и изменены с помощью **Java** и API-интерфейсов командной строки. В **NameNode** также хранится набор **Cache Pools**, являющийся административным объектом, который группирует **Cache Directives** для управления ресурсами, а также для обеспечения прав доступа.

**NameNode** периодически повторно сканирует пространство имен и актив **Cache Directives**, чтобы определить, какие блоки нужно кэшировать, а какие нет, и назначает задачи кэширования **DataNodes**.

Повторное сканирование также может быть вызвано действиями пользователя, такими как добавление или удаление **Cache Directives** или удаление **Cache Pools**.

Блоки кэша, находящиеся в стадии разработки, поврежденные или неполные, не кэшируются. Если **Cache Directives** содержит ссылку, адрес ссылки не кэшируется.

В настоящее время кэширование может применяться только к каталогам и файлам.

## 5.4 Терминология кэширования

### 5.4.1 Cache Directive

**Cache Directive** определяет контур для кэширования. Пути могут указывать либо каталоги, либо файлы. Каталоги кэшируются не рекурсивно, то есть кэшируются только файлы в листинге каталога первого уровня.

**Cache Directives** также указывают дополнительные параметры, такие как фактор репликации кэша и время окончания. Фактор репликации указывает количество блочных реплик в кэше. Если несколько **Cache Directives** относятся к одному файлу, применяется максимальный коэффициент репликации кэша.

Время окончания задается в командной строке как жизненный период (*time-to-live* - *TTL*), который представляет собой относительное время действия в будущем. После истечения срока действия **Cache Directive** больше не учитывается **NameNode** при принятии решений кэширования.

### 5.4.2 Cache Pool

**Cache Pool** - это административный объект, используемый для управления группами директив кэша. Кэш-пулы имеют UNIX-подобные разрешения, которые ограничивают доступ пользователей и групп к пулу. Разрешения на запись позволяют пользователям добавлять и удалять директивы кэша в пул. Разрешения на чтение позволяют пользователям просматривать директивы кэша в пуле и дополнительные метаданные. Execute-разрешение не используется.

**Cache Pools** также используются для управления ресурсами. Кэш-пулы могут обеспечить максимальный предел памяти, ограничивающий совокупное количество байтов, которые могут быть кэшированы директивами в пуле. Как правило, сумма лимитов пула приблизительно равна суммарной памяти, зарезервированной для кэширования **HDFS** в кластере. Кэш-пулы также мониторят ряд статистических данных, чтобы помочь пользователям кластера отслеживать, что в настоящее время кэшируется, и определить, что еще нужно кэшировать.

**Cache Pools** также могут обеспечить максимальный жизненный период, ограничив максимальное время истечения срока действия директив, добавляемых в пул.

## 5.5 Настройка централизованного кэширования

### 5.5.1 Собственные библиотеки

Для отгорождения блокировки файлов в памяти, **DataNode** использует собственный код *JNI* из *libhadoop.so*.

---

**Important:** Обязательно включите *JNI*, если используется централизованное управление кешем **HDFS**

---

### 5.5.2 Свойства конфигурации

Свойства конфигурации для централизованного кэширования указаны в файле *hdfs-site.xml*.

### 5.5.3 Требуемые свойства

В настоящее время требуется только одно свойство:

- *dfs.datanode.max.locked.memory*. Это свойство определяет максимальный объем памяти (в байтах), который будет использовать DataNode для кэширования. Также необходимо увеличить размер “заблокированного объема памяти” *ulimit (ulimit -l)* пользователя DataNode, чтобы превысить этот параметр (более подробно описано в следующем разделе “Дополнительные свойства”). При настройке данного значения необходимо помнить, что пространство в памяти также понадобится и для других целей, таких как JNM и DataNode, а также страниц кэша операционной системы.

Пример:

```
<property>
  <name>dfs.datanode.max.locked.memory</name>
  <value>268435456</value>
</property>
```

### 5.5.4 Дополнительные свойства

Следующие свойства не являются обязательными, но могут быть заданы в настройках:

- *dfs.namenode.path.based.cache.refresh.interval.ms* число миллисекунд, которое NameNode использует между последующими повторными сканированиями кэша. По умолчанию этот параметр установлен на *300000*, что составляет пять минут. Пример:

```
<property>
  <name>dfs.namenode.path.based.cache.refresh.interval.ms</name>
  <value>300000</value>
</property>
```

- *dfs.time.between.resending.caching.directives.ms* NameNode использует это значение как количество миллисекунд между повторным кэшированием директивов. Пример:

```
<property>
  <name>dfs.time.between.resending.caching.directives.ms</name>
  <value>300000</value>
</property>
```

- *dfs.datanode.fsdatasetcache.max.threads.per.volume* DataNode использует это значение как максимальное количество потоков на единицу объема для кэширования новых данных. По умолчанию этот параметр имеет значение *4*. Пример:

```
<property>
  <name>dfs.datanode.fsdatasetcache.max.threads.per.volume</name>
  <value>4</value>
</property>
```

- *dfs.cachereport.intervalMsec* DataNode использует это значение как число миллисекунд между отправкой отчета о состоянии кэша в NameNode. По умолчанию этот параметр установлен на *10000*, что составляет 10 секунд. Пример:

```
<property>
  <name>dfs.cachereport.intervalMsec</name>
  <value>10000</value>
</property>
```

- *dfs.namenode.path.based.cache.block.map.allocation.percent* Процент Java heap, распределенный по картам кэшированных блоков. Карта кэшированных блоков - это хеш-карта, которая использует связанное хэширование. Доступ к меньшим картам осуществляется медленнее, чем если количество кэшированных блоков велико; большие карты потребляют больше памяти. Значение по умолчанию равно 0,25%. Пример:

```
<property>
  <name>dfs.namenode.path.based.cache.block.map.allocation.percent</name>
  <value>0.25</value>
</property>
```

## 5.6 Ограничения ОС

Если выдается сообщение об ошибке “*Cannot start datanode because the configured max locked memory size... is more than the datanode’s available RLIMIT\_MEMLOCK ulimit*”, это означает, что операционная система накладывает более низкое ограничение на объем памяти, который можно заблокировать, чем настроено. Чтобы исправить это, необходимо настроить значение *ulimit -l*, с которым работает **DataNode**. Это значение обычно настраивается в файле */etc/security/limits.conf* (может варьироваться в зависимости от используемой ОС и дистрибутива).

Вы узнаете, что значение настроено правильно, когда сможете запустить *ulimit -l* из оболочки и получить либо более высокое значение, чем настроенное, либо строку “unlimited”, что указывает на отсутствие ограничения.

---

**Important:** Для *ulimit -l* характерно выводить ограничение блокировки памяти в килобайтах (КБ), но при этом *dfs.datanode.max.locked.memory* должно быть указано в байтах.

---

Например, значение *dfs.datanode.max.locked.memory* установлено в 128000 байт:

```
<property>
  <name>dfs.datanode.max.locked.memory</name>
  <value>128000</value>
</property>
```

Лучше установить *memlock* (максимальное адресное пространство с закрытой памятью) на несколько большее значение. Например, чтобы установить *memlock* на 130 KB (130 000 байт) для пользователя *hdfs*, необходимо добавить следующую строку в */etc/security/limits.conf*:

```
hdfs - memlock 130
```

---

**Important:** Информация в данном разделе не применяется к развертыванию в Windows. Windows не имеет прямого эквивалента *ulimit -l*.

---

## 5.7 Использование Cache Pools и Directives

Можно использовать **интерфейс командной строки (CLI)** для создания, изменения и перечисления **Cache Pools** и **Cache Directives** с помощью подкоманды *hdfs cacheadmin*.

**Cache Directives** идентифицируются уникальным не повторяющимся 64-битным ID. Идентификаторы не используются повторно, даже если **Cache Directive** удалена.

**Cache Pools** идентифицируются по уникальному имени строки.

Сначала следует создать **Cache Pools**, а затем добавить в него **Cache Directives**.

## 5.8 Команды Cache Pools

### 5.8.1 addPool

Команда добавления нового **Cache Pool**:

```
hdfs cacheadmin -addPool <name> [-owner <owner>] [-group <group>]
[-mode <mode>] [-limit <limit>] [-maxTtl <maxTtl>]
```

Функции команды *addPool* описаны в таблице.

Таблица 5.1.: Функции команды addPool

Функция	Описание
<name>	Имя нового Cache Pool
<owner>	Имя пользователя владельца Cache Pool. По умолчанию используется текущий пользователь
<group>	Группа, которой назначен Cache Pool. По умолчанию используется имя основной группы текущего пользователя
<mode>	Восьмеричные разрешения в стиле UNIX, назначенные Cache Pool. По умолчанию установлены <i>0755</i>
<limit>	Максимальное количество байтов, которые в совокупности могут быть кэшированы директивами в Cache Pool. По умолчанию ограничение не установлено
<maxTtl>	Максимальное допустимое время ожидания для директив, добавляемых в Cache Pool. Значение может быть указано в секундах, минутах, часах и днях, например, <i>120 s, 30 m, 4 h, 2 d</i> . Допустимыми единицами являются <i>[smhd]</i> . По умолчанию максимальное значение не задано. Значение <i>never</i> указывает, что предела нет

### 5.8.2 modifyPool

Команда изменения метаданных существующего **Cache Pool**:

```
hdfs cacheadmin -modifyPool <name> [-owner <owner>] [-group <group>]
[-mode <mode>] [-limit <limit>] [-maxTtl <maxTtl>]
```

Функции команды *modifyPool* описаны в таблице.

Таблица 5.2.: Функции команды removePool

Функция	Описание
<name>	Имя требующего изменения Cache Pool
<owner>	Имя пользователя владельца Cache Pool
<group>	Группа, которой назначен Cache Pool
<mode>	Восьмеричные разрешения в стиле UNIX, назначенные Cache Pool
<limit>	Максимальное количество байтов, которые в совокупности могут быть кэшированы директивами в Cache Pool
<maxTtl>	Максимальное допустимое время ожидания для директив, добавляемых в Cache Pool. Значение может быть указано в секундах, минутах, часах и днях, например, <i>120 s, 30 m, 4 h, 2 d</i> . Допустимыми единицами являются <i>[smhd]</i> . По умолчанию максимальное значение не задано. Значение <i>never</i> указывает, что предела нет

### 5.8.3 removePool

Команда удаления **Cache Pool**. Также удаляет пути, связанные с ним:



```
hdfs cacheadmin -removePool <name>
```

Функции команды *removePool* описаны в таблице.

Таблица 5.3.: Функции команды *removePool*

Функция	Описание
<name>	Имя удаляемого Cache Pool

### 5.8.4 listPools

Команда отображает информацию об одном или нескольких **Cache Pool**, например, имя, владельца, группу, разрешения и прочее:

```
hdfs cacheadmin -listPools [-stats] [<name>]
```

Функции команды *listPools* описаны в таблице.

Таблица 5.4.: Функции команды *listPools*

Функция	Описание
<-stats>	Отображение дополнительной статистики по Cache Pool
<name>	Если параметр задан, то выдается только упомянутый Cache Pool

### 5.8.5 help

Отображает подробную информацию о команде:

```
hdfs cacheadmin -help <command-name>
```

Функции команды *help* описаны в таблице.

Таблица 5.5.: Функции команды *help*

Функция	Описание
<command-name>	Отображение подробной информации по указанной команде. Если команда не указана, отображается подробная справка по всем командам

## 5.9 Команды Cache Directives

### 5.9.1 addDirective

Команда добавления нового **Cache Directive**:

```
hdfs cacheadmin -addDirective -path <path> -pool <pool-name> [-force]
[-replication <replication>] [-ttl <time-to-live>]
```

Функции команды *addDirective* описаны в таблице.

Таблица 5.6.: Функции команды `addDirective`

Функция	Описание
<path>	Путь к каталогу кэша или файлу
<pool-name>	Cache Pool, к которому добавляется Cache Directive. Необходимо разрешение для Cache Pool на запись, чтобы добавить новые директивы
<-force>	Пропуск проверки ограничений ресурсов Cache Pool
<-replication>	Восьмеричные разрешения в стиле UNIX, назначенные Cache Pool. По умолчанию установлены <code>0755</code>
<limit>	Используемый коэффициент репликации кэша. По умолчанию установлено значение <code>1</code>
<time-to-live>	Продолжительность действия директивы. Значение может быть указано в минутах, часах и днях, например, <code>30 m</code> , <code>4 h</code> , <code>2 d</code> . Допустимыми единицами являются [smhd]. Значение <code>never</code> означает, что директива никогда не истекает. Если параметр не установлен, директива никогда не истекает

### 5.9.2 removeDirective

Команда удаления **Cache Directive**:

```
hdfs cacheadmin -removeDirective <id>
```

Функции команды `removeDirective` описаны в таблице.

Таблица 5.7.: Функции команды `removeDirective`

Функция	Описание
<id>	Идентификатор Cache Directive для удаления. Необходимо разрешение <code>Write</code> Cache Pool, к которому принадлежит директива. Можно использовать команду <code>-listDirectives</code> для отображения списка идентификаторов Cache Directive

### 5.9.3 removeDirectives

Команда удаления всех **Cache Directives** по указанному пути:

```
hdfs cacheadmin -removeDirectives <path>
```

Функции команды `removeDirectives` описаны в таблице.

Таблица 5.8.: Функции команды `removeDirectives`

Функция	Описание
<path>	Путь Cache Directives для удаления. Необходимо разрешение <code>Write</code> Cache Pool, к которому относятся директивы. Можно использовать команду <code>-listDirectives</code> для отображения списка Cache Directives

### 5.9.4 listDirectives

Команда возврата списка **Cache Directives**:

```
hdfs cacheadmin -listDirectives [-stats] [-path <path>] [-pool <pool>]
```

Функции команды `listDirectives` описаны в таблице.

Таблица 5.9.: Функции команды listDirectives

Функция	Описание
<path>	Список Cache Directives данного пути. Если в <path>, принадлежащему Cache Pool, нет доступа <i>Read</i> , Cache Directive не указывается
<pool>	Список Cache Directives, относящихся только к данному Cache Pool
<-stats>	Статистика по Cache Directive указанного пути

## Глава 6

# Настройка HDFS Compression

В данном разделе описывается, как настроить **HDFS Compression** на **Linux**.

**Linux** поддерживает **GzipCodec**, **DefaultCodec**, **BZip2Codec**, **LzoCodec** и **SnappyCodec**. Как правило, для **HDFS Compression** используется **GzipCodec**.

Для **GzipCodec** необходимо выполнить следующие инструкции:

- Вариант I: использовать **GzipCodec** для одnorазовых заданий:

```
hadoop jar hadoop-examples-1.1.0-SNAPSHOT.jar sort sbr"-Dmapred.compress.map.output=true" sbr"-Dmapred.
↪map.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec"sbr "-Dmapred.output.
↪compress=true" sbr"-Dmapred.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec"sbr -
↪outKey org.apache.hadoop.io.Textsbr -outValue org.apache.hadoop.io.Text input output
```

- Вариант II: включить **GzipCodec** в качестве сжатия по умолчанию:

– Отредактировать файл *core-site.xml* на главной машине **NameNode**:

```
<property>
<name>io.compression.codecs</name>    <value>org.apache.hadoop.io.compress.GzipCodec,org.apache.
↪hadoop.io.compress.DefaultCodec,com.hadoop.compression.lzo.LzoCodec,org.apache.hadoop.io.
↪compress.SnappyCodec</value>
<description>A list of the compression codec classes that can be used for compression/
↪decompression.</description>
</property>
```

– Изменить файл *mapred-site.xml* на главной машине **JobTracker**:

```
<property>
  <name>mapred.compress.map.output</name>
  <value>true</value>
</property>
```

```
<property>
  <name>mapred.map.output.compression.codec</name>
  <value>org.apache.hadoop.io.compress.GzipCodec</value>
</property>
```

```
<property>
  <name>mapred.output.compression.type</name>
  <value>BLOCK</value>
</property>
```

- 
- (Опционально) Задать следующие два параметра конфигурации для включения сжатия задания. Изменить файл *mapred-site.xml* на главной машине Resource Manager:

```
<property>
  <name>mapred.output.compress</name>
  <value>true</value>
</property>
```

```
<property>
  <name>mapred.output.compression.codec</name>
  <value>org.apache.hadoop.io.compress.GzipCodec</value>
</property>
```

- Перезапустить кластер.

## Глава 7

# Настройка Rack Awareness на ADH

Для настройки **Rack Awareness** на кластере **ADH** необходимо выполнить следующие настройки:

### 7.1 I. Создание скрипта Rack Topology

**Hadoop** использует скрипты топологии для определения местоположения стойки узлов и применяет данную информацию для репликации данных блока в резервные стойки.

- Создать скрипт топологии и файл данных. Скрипт топологии должен быть исполняемым.

**Пример скрипта топологии.** Имя файла: *rack-topology.sh*

```
#!/bin/bash

# Adjust/Add the property "net.topology.script.file.name"
# to core-site.xml with the "absolute" path the this
# file. ENSURE the file is "executable".

# Supply appropriate rack prefix
RACK_PREFIX=default

# To test, supply a hostname as script input:
if [ $# -gt 0 ]; then

CTL_FILE=${CTL_FILE:-"rack_topology.data"}

HADOOP_CONF=${HADOOP_CONF:-"/etc/hadoop/conf"}

if [ ! -f ${HADOOP_CONF}/${CTL_FILE} ]; then
    echo -n "$RACK_PREFIX/rack "
    exit 0
fi

while [ $# -gt 0 ] ; do
    nodeArg=$1
    exec< ${HADOOP_CONF}/${CTL_FILE}
    result=""
    while read line ; do
        ar=( $line )
        if [ "${ar[0]}" = "$nodeArg" ] ; then
            result="${ar[1]}"
        fi
    done
done
```

```

    fi
done
shift
if [ -z "$result" ] ; then
    echo -n "$RACK_PREFIX/rack "
else
    echo -n "$RACK_PREFIX/rack_$result "
fi
done

else
    echo -n "$RACK_PREFIX/rack "
fi

```

**Пример файла данных топологии.** Имя файла: *rack\_topology.data*

```

# This file should be:
# - Placed in the /etc/hadoop/conf directory
#   - On the Namenode (and backups IE: HA, Failover, etc)
#   - On the Job Tracker OR Resource Manager (and any Failover JT's/RM's)
# This file should be placed in the /etc/hadoop/conf directory.

# Add Hostnames to this file. Format <host ip> <rack_location>
192.0.2.0 01
192.0.2.1 02
192.0.2.2 03

```

- Скопировать оба этих файла в каталог */etc/hadoop/conf* на всех узлах кластера;
- Запустить скрипт *rack-topology.sh*, чтобы убедиться, что он возвращает правильную информацию о стойке для каждого хоста.

## 7.2 II. Добавление свойства Script Topology в core-site.xml

- Остановить HDFS;
- Добавить в *core-site.xml* следующее свойство:

```

<property>
  <name>net.topology.script.file.name</name>
  <value>/etc/hadoop/conf/rack-topology.sh</value>
</property>

```

По умолчанию скрипт топологии обрабатывает до 100 заявок за запрос. Можно указать другое количество заявок в свойстве *net.topology.script.number.args*. Например:

```

<property>
  <name>net.topology.script.number.args</name>
  <value>75</value>
</property>

```

## 7.3 III. Перезапуск HDFS и MapReduce

Перезапустить **HDFS** и **MapReduce**.

## 7.4 IV. Контроль Rack Awareness

После запуска сервисов можно использовать следующие способы для контроля, что **Rack Awareness** активирована:

- Просмотреть журналы NameNode, расположенные в `/var/log/hadoop/hdfs/` (например: `hadoop-hdfs-namenode-sandbox.log`). Должна быть следующая запись:

```
014-01-13 15:58:08,495 INFO org.apache.hadoop.net.NetworkTopology: Adding a new node: /rack01/
-><ipaddress>
```

- Команда Hadoop `fsck` должна возвращать на подобии следующего (в случае двух стоек):

```
Status: HEALTHY
Total size: 123456789 B
Total dirs: 0
Total files: 1
Total blocks (validated): 1 (avg. block size 123456789 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)
Number of data-nodes: 40
Number of racks: 2
FSCK ended at Mon Jan 13 17:10:51 UTC 2014 in 1 milliseconds
```

- Команда Hadoop `dfsadmin -report` возвращает отчет, содержащий имя стойки рядом с каждой машиной. Отчет должен выглядеть примерно следующим образом (частично):

```
[bsmith@hadoop01 ~]$ sudo -u hdfs hadoop dfsadmin -report
Configured Capacity: 19010409390080 (17.29 TB)
Present Capacity: 18228294160384 (16.58 TB)
DFS Remaining: 5514620928000 (5.02 TB)
DFS Used: 12713673232384 (11.56 TB) DFS Used%: 69.75%
Under replicated blocks: 181
Blocks with corrupt replicas: 0
Missing blocks: 0

-----
Datanodes available: 5 (5 total, 0 dead)

Name: 192.0.2.0:50010 (h2d1.phd.local)
Hostname: h2d1.phd.local
Rack: /default/rack_02
Decommission Status : Normal
Configured Capacity: 15696052224 (14.62 GB)
DFS Used: 314380288 (299.82 MB)
Non DFS Used: 3238612992 (3.02 GB)
DFS Remaining: 12143058944 (11.31 GB)
DFS Used%: 2.00%
DFS Remaining%: 77.36%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
```



```
Cache Used%: 100.00%  
Cache Remaining%: 0.00%  
Last contact: Thu Jun 12 11:39:51 EDT 2014
```

## Глава 8

# Архивы Hadoop

**Распределенная файловая система Hadoop (HDFS)** предназначена для хранения и обработки больших наборов данных, но при этом **HDFS** может быть менее эффективна при хранении большого количества мелких файлов. Когда в **HDFS** хранится много мелких файлов, они занимают большую часть пространства имен. В результате место на диске недоиспользуется из-за ограничения пространства имен.

**Архивы Hadoop (HAR)** используются для устранения ограничений пространства имен, связанных с хранением большого количества мелких файлов. **Архив Hadoop** позволяет упаковывать небольшие файлы в блоки **HDFS** более эффективно, тем самым сокращая использование памяти **NameNode**, сохраняя прозрачный доступ к файлам. **Архивы Hadoop** также совместимы с **MapReduce**, обеспечивая прозрачный доступ к исходным файлам с помощью заданий **MapReduce**.

### 8.1 Введение

**HDFS** предназначена для хранения и обработки больших массивов данных (терабайт). Например, большой продуктивный кластер может иметь *14 ПБ* дискового пространства и хранить *60 миллионов* файлов.

Однако хранение большого количества мелких файлов в **HDFS** неэффективно. Обычно файл считается «маленьким», когда его размер существенно меньше размера блока **HDFS**, который по умолчанию равен *256 МБ* в **ADH**. Файлы и блоки являются объектами имен в **HDFS**, что означает, что они занимают пространство имен (место в **NameNode**). Таким образом, емкость пространства имен системы ограничена физической памятью **NameNode**.

Когда в системе хранится много мелких файлов, они занимают большую часть пространства имен. Как следствие, дисковое пространство недоиспользуется из-за ограничения пространства имен. В одном реальном примере кластер имел *57 миллионов* файлов размером менее *256 МБ*, при этом каждый из этих файлов занимал один блок в **NameNode**. Эти мелкие файлы использовали *95%* пространства имен, но занимали только *30%* дискового пространства кластера.

**HAR** могут использоваться для устранения ограничений пространства имен, связанных с хранением большого количества мелких файлов. **HAR** упаковывает несколько мелких файлов в большой, обеспечивая прозрачный доступ к исходным файлам (без расширения файлов).

**HAR** увеличивает масштабируемость системы за счет сокращения использования пространства имен и уменьшения нагрузки на работу в **NameNode**. Это улучшение оптимизирует память в **NameNode** и распределяет управление пространством имен по нескольким **NameNodes**.

**HAR** также совместим с **MapReduce**, он обеспечивает параллельный доступ к исходным файлам с помощью заданий **MapReduce**.

## 8.2 Компоненты архивов Hadoop

### 8.2.1 Формат модели данных HAR

Формат модели данных архивов **Hadoop** имеет следующий вид:

```
foo.har/_masterindex //stores hashes and offsets
foo.har/_index //stores file statuses
foo.har/part-[1..n] //stores actual file data
```

Файлы данных хранятся в нескольких файлах, которые индексируются для сохранения первоначального разделения данных. Кроме того, файлы доступны параллельно с помощью программ **MapReduce**. В индексах файлов также записываются исходные структуры дерева каталогов и статус файла.

### 8.2.2 Файловая система HAR

Большинство архивных систем, таких как **tar**, являются инструментами для архивирования и деархивации. Как правило, они не вписываются в фактический уровень файловой системы и, следовательно, не являются прозрачными для разработчика приложения, поскольку пользователь должен предварительно деархивировать (расширить) архив перед использованием.

**Архив Hadoop** интегрируется с интерфейсом **файловой системы Hadoop**. *HarFileSystem* реализует интерфейс *FileSystem* и предусматривает доступ через *har://*. Это обеспечивает прозрачность архивных файлов и структур дерева каталогов для пользователей. Доступ к файлам в **HAR** можно получить напрямую, без его расширения.

Например, следующая команда для копирования файла **HDFS** в локальный каталог:

```
hdfs dfs -get hdfs://namenode/foo/file-1 localdir
```

Предположим, что **Архив Hadoop** *bar.har* создан из каталога *foo*. С помощью **HAR** команда для копирования исходного файла становится следующей:

```
hdfs dfs -get har://namenode/bar.har/foo/file-1 localdir
```

Пользователям следует изменить пути URI. Но при этом пользователи могут создать символическую ссылку (из *hdfs://namenode/foo* для *har://namenode/bar.har/foo* в примере выше), и тогда изменять URI не будет необходимости. В любом случае, *HarFileSystem* вызывается автоматически для обеспечения доступа к файлам в **HAR**. Из-за этого прозрачного слоя **HAR** совместим с **API Hadoop**, **MapReduce**, интерфейсом командной строки **оболочки FS** и приложениями более высокого уровня, такими как **Pig**, **Zebra**, **Streaming**, **Pipes** и **DistCp**.

### 8.2.3 Инструмент архивации Hadoop

**Архивы Hadoop** могут быть созданы с использованием инструмента архивации **Hadoop**. Инструмент архивации использует **MapReduce** для эффективного параллельного создания **архивов Hadoop**. Инструмент вызывается с помощью команды:

```
hadoop archive -archiveName name -p <parent> <src>* <dest>
```

Список файлов генерируется путем рекурсивного перемещения исходных каталогов, а затем список разбивается на карту входящих задач. Каждая задача создает файл (около 2 ГБ, настраивается) из подмножества исходных файлов и выводит метаданные. Наконец, *reduce task* собирает метаданные и генерирует индексные файлы.

## 8.3 Создание архива Hadoop

Инструмент архивации **Hadoop** вызывается следующей командой:

```
hadoop archive -archiveName name -p <parent> <src>* <dest>
```

Где *-archiveName* - это имя создающегося архива. В имени архива должно быть указано расширение *.har*. Аргумент *<parent>* используется для указания относительного пути к папке, в которой файлы будут архивироваться в **HAR**. Например:

```
hadoop archive -archiveName foo.har -p /user/hadoop dir1 dir2 /user/zoo
```

В данном примере создается архив с использованием */user/hadoop* в качестве каталога архива. Каталоги */user/hadoop/dir1* и */user/hadoop/dir2* будут заархивированы в архиве */user/zoo/foo.har*.

---

**Important:** Архивирование не удаляет исходные файлы. При необходимости удаления входных файлов после создания архива (в целях сокращения пространства имен), исходные файлы удаляются вручную

---

Хотя команда архивации **Hadoop** может быть запущена из файловой системы хоста, файл архива создается в **HDFS** из существующих каталогов. Если сослаться на каталог в файловой системе хоста, а не на **HDFS**, выдается следующая ошибка:

```
The resolved paths set is empty. Please check whether the srcPaths exist, where srcPaths
= [</directory/path>]
```

Для создания каталогов **HDFS**, используемых в предыдущем примере, необходимо выполнить следующую команду:

```
hdfs dfs -mkdir /user/zoo
hdfs dfs -mkdir /user/hadoop
hdfs dfs -mkdir /user/hadoop/dir1
hdfs dfs -mkdir /user/hadoop/dir2
```

### 8.3.1 Просмотр файлов в архивах Hadoop

Команда *hdfs dfs -ls* может использоваться для поиска файлов в архивах **Hadoop**. Используя пример архива */user/zoo/foo.har*, созданный в предыдущем разделе, необходимо использовать следующую команду для вывода списка файлов в архиве:

```
hdfs dfs -ls har:///user/zoo/foo.har/
```

Результатом будет:

```
har:///user/zoo/foo.har/dir1
har:///user/zoo/foo.har/dir2
```

Данные архивы были созданы с помощью следующей команды:

```
hadoop archive -archiveName foo.har -p /user/hadoop dir1 dir2 /user/zoo
```

Если изменить данную команду на:

```
hadoop archive -archiveName foo.har -p /user/ hadoop/dir1 hadoop/dir2 /user/zoo
```

И затем выполнить следующую команду:

```
hdfs dfs -ls -R har:///user/zoo/foo.har
```

То результатом будет:

```
har:///user/zoo/foo.har/hadoop
har:///user/zoo/foo.har/hadoop/dir1
har:///user/zoo/foo.har/hadoop/dir2
```

Следует обратить внимание, что с измененным родительским аргументом файлы заархивированы относительно `/user/`, а не `/user/hadoop`.

## 8.4 Hadoop Archives и MapReduce

Чтобы использовать **архивы Hadoop** с **MapReduce**, необходимо ссылаться на файлы несколько иначе, чем на файловую систему по умолчанию. Если есть **архив Hadoop**, хранящийся в **HDFS** в `/user/zoo/foo.har`, следует указать каталог ввода как `har:///user/zoo/foo.har`, чтобы использовать его как **MapReduce**. Поскольку **архивы Hadoop** отображаются как файловая система, **MapReduce** может использовать все логические входные файлы в **архивы Hadoop** в качестве входных данных.

## Глава 9

# API-интерфейсы JMX Metrics для HDFS Daemons

Для доступа к показателям **HDFS** можно использовать методы с помощью API-интерфейсов **Java Management Extensions (JMX)**.

### 9.1 Использование веб-интерфейса HDFS Daemon

Доступ к метрикам **JMX** можно получить через веб-интерфейс **HDFS daemon**. Это рекомендуемый метод.

Например, для доступа к **NameNode JMX** необходимо использовать следующий формат команды:

```
curl -i http://localhost:50070/jmx
```

Для извлечения только определенного ключа можно использовать параметр *qry*:

```
curl -i http://localhost:50070/jmx?qry=Hadoop:service=NameNode,name=NameNodeInfo
```

### 9.2 Прямой доступ к удаленному агенту JMX

Данный метод требует, чтобы удаленный агент **JMX** был включен с опцией *JVM* при запуске сервисов **HDFS**.

Например, следующие параметры *JVM* в *hadoop-env.sh* используются для включения удаленного агента **JMX** для **NameNode**. Он работает на порту *8004* с отключенным **SSL**. Имя пользователя и пароль сохраняются в файле *mxremote.password*.

```
export HADOOP_NAMENODE_OPTS="-Dcom.sun.management.jmxremote
-Dcom.sun.management.jmxremote.password.file=$HADOOP_CONF_DIR/jmxremote.password
-Dcom.sun.management.jmxremote.ssl=false
-Dcom.sun.management.jmxremote.port=8004 $HADOOP_NAMENODE_OPTS"
```

Подробности о связанных настройках можно найти [здесь](#). Также можно использовать инструмент *jmxquery* для извлечения информации через **JMX**.

**Hadoop** также имеет встроенный инструмент запросов **JMX** *jmxget*. Например:

```
hdfs jmxget -server localhost -port 8004 -service NameNode
```

---

**Important:** Инструмент *jmxget* требует, чтобы аутентификация была отключена, так как она не принимает имя пользователя и пароль

---

Использование **JMX** может быть сложным для персонала, который не знаком с настройкой **JMX**, особенно **JMX** с **SSL** и **firewall tunnelling**. Поэтому обычно рекомендуется собирать информацию **JMX** через веб-интерфейс **HDFS daemon**, а не напрямую обращаться к удаленному агенту **JMX**.

## Глава 10

# Память в качестве хранилища (техническое превью)

В данной главе описывается как использовать память **DataNode** в качестве хранилища в **HDFS**.

---

**Important:** Данная возможность представлена как технический обзор и рассмотрена в рамках развития. Не следует использовать ее в продуктивной среде. При наличии вопросов относительно данной особенности необходимо обратиться в службу поддержки – [info@arenadata.io](mailto:info@arenadata.io)

---

### 10.1 Введение

**HDFS** поддерживает эффективную запись больших наборов данных в надежное хранилище, а также обеспечивает безотказный доступ к данным. Это хорошо работает для пакетных заданий, записывающих большое количество данных.

Новые классы приложений управляют вариантами использования для записи меньшего количества временных данных. При использовании памяти **DataNode** в качестве хранилища адресов вариантов использования приложений, записывается относительно небольшое количество промежуточных наборов данных с низкой задержкой.

Запись данных блока в память снижает надежность, поскольку данные могут быть потеряны из-за перезагрузки процесса до их сохранения на диск. При этом **HDFS** пытается своевременно сохранить копии данных на диск, чтобы уменьшить риск потери данных.

Память **DataNode** указывается при использовании типа хранения *RAM\_DISK* и политики хранения *LAZY\_PERSIST*.

Для использования памяти **DataNode** в качестве хранилища **HDFS** необходимо выполнить следующие шаги:

- Выключить **DataNode**;
- Установить часть памяти **DataNode** для использования **HDFS**;
- Назначить **DataNode** тип хранения *RAM\_DISK* и включить режим локального чтения данных;
- Установить политику хранения *LAZY\_PERSIST* в файлах и каталогах **HDFS**, которые будут использовать память в качестве хранилища;
- Перезапустить **DataNode**.



При обновлении параметра политики хранения в файле или каталоге, необходимо использовать инструмент переноса данных **HDFS mover** для фактического перемещения блоков (как указано в новой политике хранения).

Память как хранилище представляет собой один из аспектов возможностей управления ресурсами **YARN**, который включает **CPU scheduling**, **CGroups**, **node labels** и архивное хранилище.

## 10.2 Типы хранилища HDFS

Типы хранилища **HDFS** могут использоваться для назначения данных для различных типов физических носителей информации. Доступны следующие типы хранилища:

- **DISK** - дисковое хранилище (тип хранения по умолчанию);
- **ARCHIVE** - архивное хранилище (высокая плотность хранения, низкий ресурс обработки);
- **SSD** - твердотельный накопитель;
- **RAM\_DISK** - память **DataNode**.

Если тип хранилища не назначен, по умолчанию используется тип *DISK*.

## 10.3 Политика хранения LAZY\_PERSIST

С помощью политики хранения *LAZY\_PERSIST* можно хранить данные в сконфигурированной памяти **DataNode**. При этом первая копия данных хранится в *RAM\_DISK* (память **DataNode**), а остальные копии - на *DISK*. Резервным хранилищем для создания и копирования является *DISK*.

В таблице приведена политика хранения копий данных.

Таблица 10.1.: Политика хранения копий данных

ID политики	Название политики	Размещение блока (для n реплик)	Резервное хранилище для генерации	Резервное хранилище для копирования
15	LAZY_PERSIST	RAM_DISK: 1, DISK: n-1	DISK	DISK

---

**Important:** В настоящее время политики хранения нельзя редактировать

---

## 10.4 Настройка памяти в качестве хранилища

Для настройки памяти в качестве хранилища необходимо выполнить следующие действия:

### 1. Выключить **DataNode**

Закрывать **DataNode** с помощью соответствующих команд.

### 2. Установить часть памяти **DataNode** для **HDFS**

Для использования памяти **DataNode** в качестве хранилища, необходимо сначала установить часть памяти **DataNode** для использования **HDFS**.

Например, для выделения *2 ГБ* памяти для хранения **HDFS** необходимо использовать следующие команды:

```
sudo mkdir -p /mnt/hdfsramdisk
sudo mount -t tmpfs -o size=2048m tmpfs /mnt/hdfsramdisk
Sudo mkdir -p /usr/lib/hadoop-hdfs
```

### 3. Назначить тип памяти `RAM_DISK` и включить режим локального чтения данных

Чтобы присвоить `DataNodes` тип памяти `RAM_DISK` и включить режим локального чтения данных, необходимо изменить следующие свойства в файле `/etc/hadoop/conf/hdfs-site.xml`:

- Свойство `dfs.name.dir` определяет, где в локальной файловой системе DataNode хранит свои блоки. Чтобы указать DataNode в качестве хранилища `RAM_DISK`, необходимо добавить `[RAM_DISK]` в начало пути локальной файловой системы и в свойство `dfs.name.dir`;
- Установить для параметра `dfs.client.read.shortcircuit` значение `true`, чтобы включить режим локального чтения данных.

Например:

```
<property>
  <name>dfs.data.dir</name>
  <value>file:///grid/3/aa/hdfs/data/,[RAM_DISK]file:///mnt/hdfsramdisk</value>
</property>

<property>
  <name>dfs.client.read.shortcircuit</name>
  <value>true</value>
</property>

<property>
  <name>dfs.domain.socket.path</name>
  <value>/var/lib/hadoop-hdfs/dn_socket</value>
</property>

<property>
  <name>dfs.checksum.type</name>
  <value>NULL</value>
</property>
```

### 4. Установить политику хранения `LAZY_PERSIST` в файлах или каталогах

Для установки политики хранения `LAZY_PERSIST` в файлах или каталогах необходимо выполнить команду:

```
hdfs dfsadmin -setStoragePolicy <path> <policyName>
```

Аргументы:

- `<path>` - путь к каталогу или файлу;
- `<policyName>` - название политики хранения.

Пример:

```
hdfs dfsadmin -setStoragePolicy /memory1 LAZY_PERSIST
```

Для возврата политики хранения файла или каталога необходимо выполнить команду:

```
hdfs dfsadmin -getStoragePolicy <path>
```

Аргументы:

- `<path>` - путь к каталогу или файлу.

Пример:

```
hdfs dfsadmin -getStoragePolicy /memory1 LAZY_PERSIST
```

### 5. Запуск DataNode

Запустить **DataNode** с помощью соответствующих команд.

## 10.5 Использование “mover” для применения политик хранения

При обновлении параметра политики хранения в файле или каталоге, новая политика не применяется автоматически. Необходимо использовать инструмент переноса данных **HDFS mover** для фактического перемещения блоков (как указано в новой политике хранения).

Средство миграции данных *mover* сканирует указанные файлы в **HDFS** и проверяет, соответствует ли размещение блоков политике хранения. Копии блоков, нарушающих политику хранения, он перемещает в соответствующий тип хранилища для выполнения требований политики.

Команда:

```
hdfs mover [-p <files/dirs> | -f <local file name>]
```

Аргументы:

- *-p <files/dirs>* - список файлов/каталогов **HDFS** для переноса, разделенные пробелами;
- *-f <local file>* - локальный файл, содержащий список файлов/каталогов **HDFS** для переноса

---

**Important:** Если оба параметра *-p* и *-f* опущены, путь по умолчанию является корневым каталогом

---

Пример:

```
hdfs mover /memory1/testfile
```

## Глава 11

# Запуск DataNodes от Non-root

В данной главе описывается как запускать **DataNodes** от пользователя без прав *root*.

### 11.1 Введение

Исторически сложилось так, что часть конфигурации безопасности для **HDFS** задействовала запуск **DataNode** от пользователя *root* и привязала привилегированные порты для конечных точек сервера. Это было сделано для решения проблемы безопасности, то есть если задание **MapReduce** запущено, а **DataNode** остановился, задачу **MapReduce** можно привязать к порту **DataNode** и потенциально сделать что-то вредоносное. Решением подобных случаев стал запуск **DataNode** от пользователя *root* и использование привилегированных портов. При этом только пользователь *root* может получить доступ к привилегированным портам.

Теперь для безопасного запуска **DataNodes** от пользователя без прав *root* можно использовать **Simple Authentication and Security Layer (SASL)**. **SASL** используется для обеспечения безопасной связи на уровне протокола.

---

**Important:** Важно выполнить передачу данных с помощью *root* при запуске **DataNodes** с использованием **SASL** для его запуска в конкретной последовательности во всем кластере. В противном случае может возникнуть риск простоя приложения

---

Для переноса существующего кластера, использующего аутентификацию *root*, с целью использования **SASL**, сначала необходимо убедиться, что версия **2.6.0** (или более поздняя) развернута для всех узлов кластера, а также для любых внешних приложений, которые необходимо подключить к кластеру. Только версии **2.6.0** + из **HDFS**-клиента могут подключаться к **DataNode**, использующему **SASL** для аутентификации протокола передачи данных, поэтому очень важно, чтобы все абоненты имели необходимую версию перед переходом.

После развертывания версии **2.6.0** (или более поздней) необходимо обновить конфигурацию любых внешних приложений, чтобы включить **SASL**. Если для **SASL** включен клиент **HDFS**, он может успешно подключиться к **DataNode**, работающему с аутентификацией *root* или аутентификацией **SASL**. Изменение конфигурации для всех клиентов гарантирует, что последующие изменения конфигурации в **DataNodes** не нарушат работу приложений. Наконец, каждый отдельный **DataNode** может быть перенесен путем изменения его конфигурации и перезапуска. Допустимо временно сочетать некоторые **DataNodes** с аутентификацией *root* и некоторые **DataNodes**, работающие с аутентификацией **SASL**, в течение периода перехода, поскольку клиент **HDFS**, подключенный для **SASL**, может подключаться к обоим.

## 11.2 Настройка DataNode SASL

Для настройки **DataNode SASL** с безопасным запуском **DataNode** от *non-root* пользователя необходимо выполнить следующие действия:

### 1. Выключить DataNode

Выключить **DataNode** с помощью соответствующих команд.

### 2. Включить SASL

Чтобы включить **DataNode SASL** в файле `/etc/hadoop/conf/hdfs-site.xml` настроить свойства.

Свойство `dfs.data.transfer.protection` позволяет использовать **DataNode SASL**. Данному свойству можно установить одно из следующих значений:

- *authentication* - устанавливает взаимную аутентификацию между клиентом и сервером;
- *integrity* - в дополнение к аутентификации гарантирует, что man-in-the-middle не может вмешиваться в сообщения, обмен которыми осуществляется между клиентом и сервером;
- *privacy* - в дополнение к функциям *authentication* и *integrity* он также полностью шифрует сообщения, обмен которыми осуществляется между клиентом и сервером.

Также необходимо установить для свойства `dfs.http.policy` значение `HTTPS_ONLY`. При этом следует указать порты для **DataNode RPC** и HTTP-серверов.

Например:

```
<property>
  <name>dfs.data.transfer.protection</name>
  <value>integrity</value>
</property>

<property>
  <name>dfs.datanode.address</name>
  <value>0.0.0.0:10019</value>
</property>

<property>
  <name>dfs.datanode.http.address</name>
  <value>0.0.0.0:10022</value>
</property>

<property>
  <name>dfs.http.policy</name>
  <value>HTTPS_ONLY</value>
</property>
```

**Important:** Параметр шифрования `dfs.encrypt.data.transfer=true` похож на `dfs.data.transfer.protection=privacy`. Эти два параметра являются взаимоисключающими, поэтому они не должны устанавливаться вместе. В случае если оба параметра установлены, `*dfs.encrypt.data.transfer` не используется

### 3. Обновить настройки экосистемы

В файле `/etc/hadoop/conf/hadoop-env.xml` изменить параметр:

```
#On secure datanodes, user to run the datanode as after dropping privileges export HADOOP_SECURE_DN_USER=
```

Строка экспорта `HADOOP_SECURE_DN_USER=dfs` включает устаревшую конфигурацию безопасности и, чтобы включить **SASL**, должна быть установлена на пустое значение.

#### 4. Запустить **DataNode**

Запустить **DataNode** с помощью соответствующих команд.

## Глава 12

# Режим локального чтения данных на HDFS

В HDFS чтение обычно проходит через **DataNode**. Таким образом, когда клиент запрашивает **DataNode** для чтения файла, **DataNode** считывает этот файл с диска и отправляет данные клиенту через сокет TCP. Так называемое “локальное чтение” читает в обход **DataNode**, позволяя клиенту непосредственно прочитать файл. Очевидно, что это возможно только в тех случаях, когда клиент находится вместе с данными. Локальное чтение обеспечивают значительное повышение производительности для многих приложений.

### 12.1 Необходимые компоненты

Для настройки локального чтения данных необходимо включить *libhadoop.so*. Подробные сведения о включении библиотеки см. в “Native Libraries”.

### 12.2 Настройка локального чтения данных на HDFS

Для настройки локального чтения данных на HDFS, необходимо в файл *hdfs-site.xml* добавить свойства, приведенные в таблице. Локальное чтение данных должно быть настроено как для **DataNode**, так и для клиента.

Таблица12.1.: Свойства для файла hdfs-site.xml

Свойство	Значение	Описание
dfs.client.read.shortcircuit	true	При значении <i>true</i> включается режим локального чтения данных

Свойство	Значение	Описание
dfs.domain.socket.path	<code>/var/lib/hadoop-hdfs/dn_socket</code>	Путь к сокету домена. В сообщениях при локальном чтении данных используется сокет домена UNIX. Это особый путь в файловой системе, позволяющий связываться клиенту и DataNodes. Необходимо установить путь к этому сокету. DataNode должен иметь возможность создать этот путь. С другой стороны, создание этого пути не должно быть возможным для любого пользователя, кроме пользователя <i>hdfs</i> или <i>root</i> . По этой причине часто используются пути в <code>/var/run</code> или <code>/var/lib</code>
dfs.client.domain.socket.data.traffic	<code>false</code>	Контролирует, будет ли обычный трафик данных передаваться через сокет домена UNIX. Функция не была сертифицирована релизами ADH, поэтому рекомендуется установить значение <i>false</i>
dfs.client.use.legacy.blockreader.local	<code>false</code>	Установка значения <i>false</i> указывает, что используется новая версия локального чтения (на основе HDFS-347). Эта версия поддерживается и рекомендуется для использования с ADH. Значение <i>true</i> означает, что используется старый режим локального чтения
dfs.datanode.hdfs-blocks-metadata.enabled	<code>true</code>	Логический тип данных, который обеспечивает поддержку на стороне сервера DataNode для экспериментального <i>DistributedFileSystem#getFileVBlockStorageLocations</i> API
dfs.client.file-block-storage-locations.timeout	<code>60</code>	Таймаут для параллельных RPC, сделанных в <i>DistributedFileSystem#getFileBlockStorageLocations</i> (в секундах). Это свойство устарело, но по-прежнему поддерживается для обратной совместимости
dfs.client.file-block-storage-locations.timeout.millis	<code>60000</code>	Таймаут для параллельных RPC, сделанных в <i>DistributedFileSystem#getFileBlockStorageLocations</i> (в миллисекундах). Это свойство заменяет <i>dfs.client.file-block-storage-locations.timeout</i> и предлагает более точный уровень детализации
dfs.client.read.shortcircuit.skip.checksum	<code>false</code>	Если этот параметр конфигурации установлен локальное чтение будет пропускать контрольную сумму файлов. Обычно это не рекомендуется, но может быть полезно для специальных настроек. Может пригодиться, если есть собственные контрольные суммы файлов вне HDFS



Свойство	Значение	Описание
dfs.client.read.shortcircuit.streams.cache.size	256	DFSClient поддерживает кэш недавно открытых файловых дескрипторов. Параметр управляет размером кэша. При установке значения выше указанного будут использоваться дополнительные дескрипторы файлов, но они могут обеспечить лучшую производительность при рабочей нагрузке с большим количеством запросов
dfs.client.read.shortcircuit.streams.cache.expiry.ms	300000	Контролирует минимальный промежуток времени нахождения файловых дескрипторов в контексте кэша клиента, прежде чем они могут быть закрыты (в миллисекундах)

XML для вышеуказанных записей:

```
<configuration>
  <property>
    <name>dfs.client.read.shortcircuit</name>
    <value>true</value>
  </property>

  <property>
    <name>dfs.domain.socket.path</name>
    <value>/var/lib/hadoop-hdfs/dn_socket</value>
  </property>

  <property>
    <name>dfs.client.domain.socket.data.traffic</name>
    <value>false</value>
  </property>

  <property>
    <name>dfs.client.use.legacy.blockreader.local</name>
    <value>false</value>
  </property>

  <property>
    <name>dfs.datanode.hdfs-blocks-metadata.enabled</name>
    <value>true</value>
  </property>

  <property>
    <name>dfs.client.file-block-storage-locations.timeout.millis</name>
    <value>60000</value>
  </property>

  <property>
    <name>dfs.client.read.shortcircuit.skip.checksum</name>
    <value>false</value>
  </property>

  <property>
    <name>dfs.client.read.shortcircuit.streams.cache.size</name>
    <value>256</value>
  </property>
</configuration>
```

```
<property>
  <name>dfs.client.read.shortcircuit.streams.cache.expiry.ms</name>
  <value>300000</value>
</property>
</configuration>
```

## Глава 13

# Руководство администратора по WebHDFS

Для настройки **WebHDFS** необходимо использовать следующие настройки:

- Настроить WebHDFS. Добавить в файл *hdfs-site.xml* следующее свойство:

```
<property>
  <name>dfs.webhdfs.enabled</name>
  <value>true</value>
</property>
```

- (Опционально) При запуске защищенного кластера выполнить следующие действия:

- Создать пользователя-принципала сервиса HTTP, используя команду:

```
kadmin: addprinc -randkey HTTP/$(Fully_Qualified_Domain_Name)@$(Realm_Name).COM
```

где *Fully\_Qualified\_Domain\_Name* - хост, на котором разворачивается NameNode; *Realm\_Name* - название сферы Kerberos.

- Создать файлы *keytab* для принципалов HTTP:

```
kadmin: xst -norandkey -k /etc/security/spnego.service.keytab
HTTP/$(Fully_Qualified_Domain_Name)
```

- Убедиться, что файл *keytab* и принципал связаны с необходимым сервисом:

```
klist -k -t /etc/security/spnego.service.keytab
```

- Добавить в файл *hdfs-site.xml* следующие свойства:

```
<property>
  <name>dfs.web.authentication.kerberos.principal</name>
  <value>HTTP/$(Fully_Qualified_Domain_Name)@$(Realm_Name).COM</value>
</property>
<property>
  <name>dfs.web.authentication.kerberos.keytab</name>
  <value>/etc/security/spnego.service.keytab</value>
</property>
```

где *Fully\_Qualified\_Domain\_Name* - хост, на котором разворачивается NameNode; *Realm\_Name* - название сферы Kerberos.

- Перезапустить сервисы NameNode и DataNode с помощью соответствующих команд.