

KARPOV.COURSES >>> КОНСПЕКТ



>Конспект > 8 урок > Отладка, профилирование и мониторинг Spark Job

> Оглавление

> Оглавление

> Profiling

Profiling OS level

Profiling JVM level

Profiling Job level

Profiling cluster level

> Требования к идеальному профилированию

Что хочется?

Возможные решения

> Обычный flow анализа

> Profiling

Профилирование - сбор характеристик работы программы, таких как время выполнения отдельных фрагментов ([wikipedia](https://en.wikipedia.org/wiki/Profiling))

Spark написан на Scala -> все средства для профилирования JVM программ нам доступны

Сбор информации о выполнении обычно делают на уровне:

- **ОС** (top, vmstat, lsof, tcpdump, netstat, iostat и т.д.) - позволяют измерять нагрузку на жесткий диск, сеть.
 - **JVM** (обычно делается с помощью agent)
 - конкретной **Spark job** - Spark имеет UI, где каждая Spark job имеет свой интерфейс для отслеживания состояния.
 - **кластера** (YARN или k8s) - самый высокоуровневый мониторинг
-

Profiling OS level

Достоинства

- честные метрики, привязанные к конкретному серверу
- родные для Linux утилиты (наличие документации)

Недостатки

- требуется установка всех необходимых утилит на все сервера/контейнеры.
 - организация отправки метрик в единое хранилище
 - зачастую трудно изолировать конкретный процесс (можно промерить общую нагрузку, включая другие работающие процессы)
-

Profiling JVM level

Достоинства

- родные для JVM средства (jconsole, visualvm, Java Flight Recorder и т.д.)
- очень подробная детализация (вплоть до вызова native кода)
- можно использовать для любых JVM программ

Недостатки

- требуется подключение к конкретным серверам и процессам (что очень затруднительно в кластерах)

- предварительно желательно локализовать проблему (для подключения к конкретному executor)
 - организация сбора метрик и дампов памяти в единое хранилище (зачастую требуется много дискового пространства)
 - можно "утонуть" в обилии информации, требуется опыт для эффективного анализа
-

Profiling Job level

Достоинства

- базовые метрики можно посмотреть как online в процессе выполнения Job, так и по завершении на History Server
- есть графовое представление выполняемых задач
- работает для любых Job вне зависимости от используемого ЯП

Недостатки

- не всегда доступно (закрытый production контур)
 - гранулярность (job -> stage -> task)
 - требуются дополнительные настройки инфраструктуры (log aggregation, history server и т.д.)
-

Profiling cluster level

Достоинства

- удобно для высокоуровневого анализа и мониторинга утилизации ресурсов
- централизация

Недостатки

- зачастую трудно отделить проблемы job от проблем инфраструктуры
- при гонке за ресурсы картина может искажаться
- очень трудно оценивать производительность задач при коммуникации с внешними системами (БД, S3/HDFS, Kafka и т.д.)

> Требования к идеальному профилированию

Что хочется?

1. Сбор метрик в **едином** хранилище

Для избежания проблем, когда, к примеру, на уровне JVM мы можем получить информацию вплоть до нативных вызовов, а на уровне кластера - только информацию о потреблении ресурсов Job.

2. Добавлять **свои** метрики

Spark имеет собственные метрики, связанные с потреблением ресурсов, характеристиками шафла, но мы можем хотеть наблюдать за бизнес-метриками. Например, вести подсчет обработки физических лиц, исключая юридические.

3. Включать/выключать метрики **по необходимости**

Желательно иметь простой и гибкий механизм для такого функционала, без привлечения разработчиков.

4. **Подробная аналитика** (группировка, временной отрезок, скользящие окна и т.д.)

5. Удобные средства **визуализации**

Удобнее наблюдать за состоянием на графиках, а не таблицах.

6. Минимальные **изменения** в коде (в идеале - без изменений)

Возможные решения

1. Сбор метрик в **едином** хранилище

Использовать push/pull доставку в time-series database (InfluxDB, Prometheus)

Обычные реляционные БД и большинство нереляционных БД не позволяют эффективно хранить и обрабатывать информацию, связанную с метриками и мониторингом. Для такого рода задач популярны time-series БД.

2. Добавлять **свои** метрики

Поддержка с Spark 3.0 (требуется дополнительный код, что логично)

3. Включать/выключать метрики **по необходимости**

Метрики необходимо группировать по приемнику (sink). Приемники можно включать/отключать на момент запуска Job.

4. **Подробная аналитика** (группировка, временной отрезок, скользящие окна и т.д.)

Time-Series DB поддерживает широкие аналитические возможности.

5. Удобные средства **визуализации**

Существуют легковесные средства визуализации, например, Grafana.

6. Минимальные **изменения** в коде (в идеале - без изменений)

Настройка приемников (sink) доступна как параметр запуска, для отладки можно выводить в консоль, а на production - statsd, InfluxDB и т.д.

Подключение разного рода профайлеров производится через java agent (jar + параметры запуска)

> Обычный flow анализа

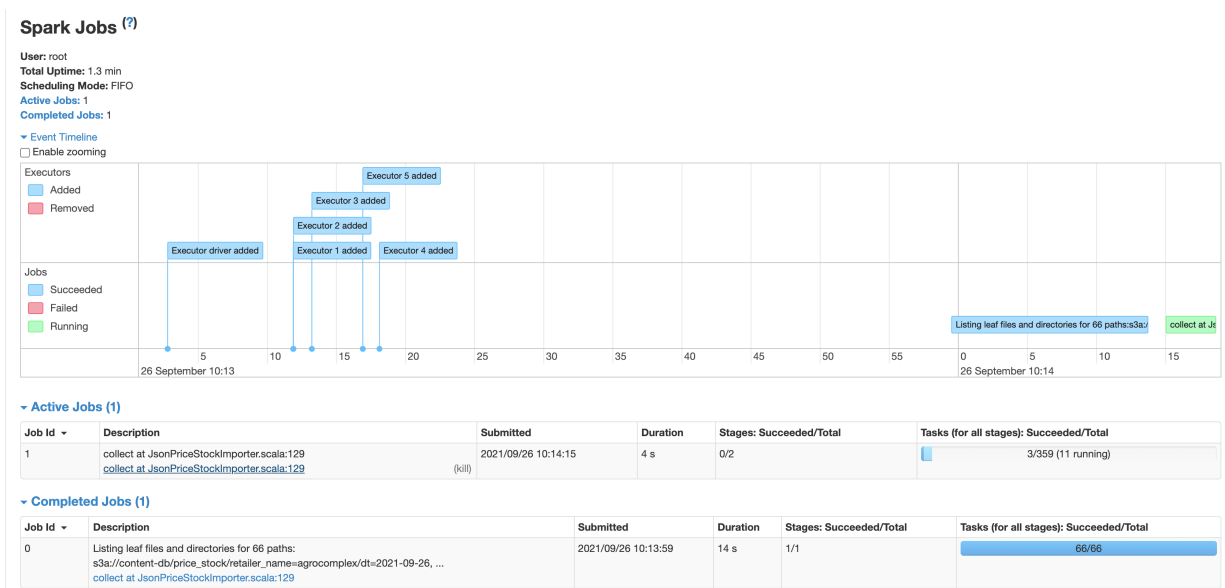
1. Смотрим на **YARN Resource Manager UI**

На уровне кластера можем посмотреть, какие Job и с каким статусом выполнялись.

application_1621415166805_0875	root	ru.sbermarket.MetricExporter	SPARK	default	0	Sun Sep 26 12:08:31 +0300 2021	Sun Sep 26 12:08:31 +0300 2021	Sun Sep 26 12:08:50 +0300 2021	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<input type="button" value="History"/>	0
--------------------------------	------	------------------------------	-------	---------	---	--------------------------------	--------------------------------	--------------------------------	----------	-----------	-----	-----	-----	-----	-----	-----	-----	--	---

2. Заходим в **History** (**Application Master**, если задание еще работает) - для более детального анализа упавшей или подозрительной Job.

3. Смотрим на **Jobs**, выбираем **подозрительную**.



Для начала стоит обратить внимание на выделение ресурсов под конкретную Job. Например, открыть **Event Timeline** и посмотреть, в какой момент добавлялись ехеситор. Могло случиться такое, что в момент запуска, ей не хватало ресурсов и она просто ждала их.

Details for Job 1

Status: RUNNING
Active Stages: 1
Pending Stages: 1

► Event Timeline
► DAG Visualization

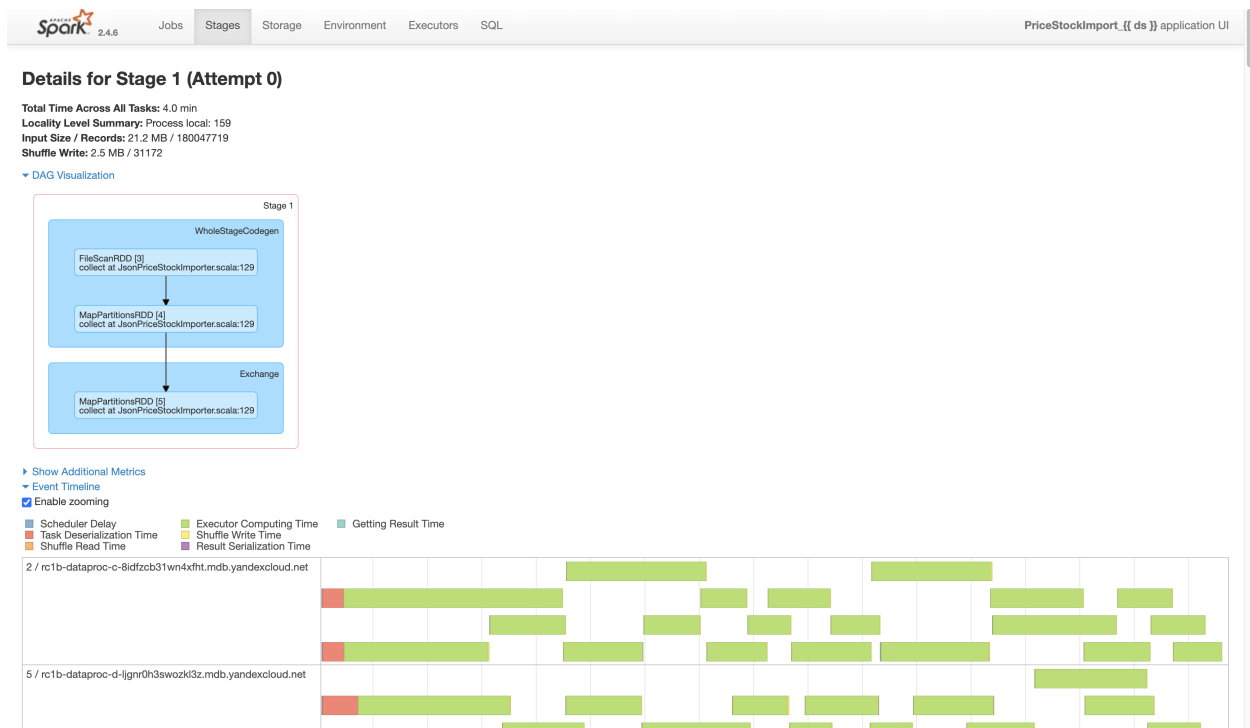
▼ Active Stages (1)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
1	collect at JsonPriceStockImporter.scala:129 +details (kill)	2021/09/26 10:14:15	16 s	94/159 (11 running)	6.5 MB			1070.9 KB

▼ Pending Stages (1)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
2	collect at JsonPriceStockImporter.scala:129 +details	Unknown	Unknown	0/200				

Перейдя на конкретную Job, мы можем смотреть из каких **Stage** (шагов) состоит выполнение терминального оператора.



По каждому **Stage** мы можем посмотреть из каких этапов он состоял, а также timeline (что происходит на каждом конкретном ехесутор в конкретный момент времени).

Summary Metrics for 159 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.3 s	1.0 s	1 s	2 s	5 s
Scheduler Delay	3 ms	5 ms	6 ms	7 ms	28 ms
Task Deserialization Time	0 ms	2 ms	2 ms	3 ms	0.7 s
GC Time	0 ms	0 ms	0 ms	8 ms	0.4 s
Result Serialization Time	0 ms	0 ms	0 ms	0 ms	2 ms
Getting Result Time	0 ms	0 ms	0 ms	0 ms	0 ms
Peak Execution Memory	64.3 MB	64.3 MB	64.3 MB	64.3 MB	64.3 MB
Input Size / Records	42.2 KB / 225	55.8 KB / 705925	94.1 KB / 942892	127.4 KB / 1128508	1034.9 KB / 4098613
Shuffle Write Size / Records	1284.0 B / 9	6.9 KB / 49	12.6 KB / 101	18.8 KB / 166	55.6 KB / 1317

▼ Aggregated Metrics by Executor

Executor ID ▲	Address	Task Time	Total Tasks	Failed Tasks	Killed Tasks	Succeeded Tasks	Input Size / Records	Shuffle Write Size / Records	Blacklisted
1 stdout stderr	rc1b-dataproc-c-97cq7ba0wj61569n.mdb.yandexcloud.net:37670	49 s	33	0	0	33	4.3 MB / 38500945	534.1 KB / 6222	false
2 stdout stderr	rc1b-dataproc-c-8idfzcb31wn4xfht.mdb.yandexcloud.net:35373	48 s	32	0	0	32	4.0 MB / 39353280	499.5 KB / 6150	false
3 stdout stderr	rc1b-dataproc-c-e60f4h554f686b7.mdb.yandexcloud.net:40540	50 s	33	0	0	33	4.4 MB / 34078689	518.8 KB / 6532	false
4 stdout stderr	rc1b-dataproc-c-e366rsg7pfnmjf0.mdb.yandexcloud.net:33036	54 s	29	0	0	29	4.7 MB / 32331460	490.4 KB / 5976	false
5 stdout stderr	rc1b-dataproc-d-ijgn0h3swozk3z.mdb.yandexcloud.net:39406	48 s	32	0	0	32	3.9 MB / 35783345	508.8 KB / 6292	false

▼ Tasks (159)

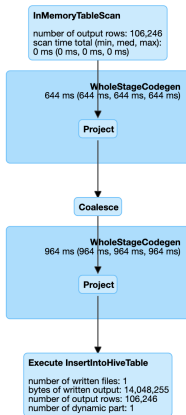
Page: 1 2 >

2 Pages. Jump to 1. Show 100 items in a page. Go

Index ▲	ID	Attempt	Status	Locality Level	Executor ID	Host	Launch Time	Duration	Scheduler Delay	Task Deserialization Time	GC Time	Result Serialization Time	Getting Result Time	Peak Execution Memory	Input Size / Records	Shuffle Write Size / Records	Write Time	Errors
0	66	0	SUCCESS	PROCESS_LOCAL	1	rc1b-dataproc-c-97cq7ba0wj61569n.mdb.yandexcloud.net stdout stderr	2021/09/26 10:14:15	4 s	18 ms	0.4 s		0 ms	0 ms	64.3 MB	78.3 KB / 1166024	19 ms	19.8 KB / 167	
1	67	0	SUCCESS	PROCESS_LOCAL	4	rc1b-dataproc-d-e366rsg7pfnmjf0.mdb.yandexcloud.net stdout stderr	2021/09/26 10:14:15	4 s	17 ms	0.6 s		0 ms	0 ms	64.3 MB	75.8 KB / 1089468	16 ms	18.0 KB / 161	

Details for Query 9

Submitted Time: 2021/09/26 10:15:54
Duration: 11 s
Succeeded Jobs: 10



▼ Details

```

== Parsed Logical Plan ==
InsertIntoTable 'UnresolvedRelation 'content.db'. 'price_stock', false, false
+- Project [external_offer_id#10394L, instamart_discount_price_per_item#12731, instamart_discount_price_per_kilo#12854, instamart_discount_price_per_pack#12977, instamart_product_type#9562,
  instamart_regular_price_per_item#12362, instamart_regular_price_per_kilo#12485, instamart_regular_price_per_pack#12608, instamart_sku#9561, instamart_sku_id#10763L, max_shop_stock#12116, min_shop_stock#12239,
  offer_items_per_pack#10517, offer_name#9556, offer_pack_type#9558, offer_vat#13346, retailer_discount_end#9573, retailer_discount_price_per_item#11624, retailer_discount_price_per_kilo#11747,
  retailer_discount_price_per_pack#11870, retailer_discount_start#9572, retailer_partner_price_per_item#11255, retailer_partner_price_per_kilo#11378, retailer_partner_price_per_pack#11501, ... 16 more fields]
  +- Repartition 1, false
     +- Project [external_offer_id#10394L, instamart_discount_price_per_item#12731, instamart_discount_price_per_kilo#12854, instamart_discount_price_per_pack#12977, instamart_product_type#9562,
       instamart_regular_price_per_item#12362, instamart_regular_price_per_kilo#12485, instamart_regular_price_per_pack#12608, instamart_sku#9561, instamart_sku_id#10763L, max_shop_stock#12116, min_shop_stock#12239,
       offer_items_per_pack#10517, offer_name#9556, offer_pack_type#9558, offer_vat#13346, retailer_discount_end#9573, retailer_discount_price_per_item#11624, retailer_discount_price_per_kilo#11747,
       retailer_discount_price_per_pack#11870, retailer_discount_start#9572, retailer_partner_price_per_item#11255, retailer_partner_price_per_kilo#11378, retailer_partner_price_per_pack#11501, ... 16 more fields]
  
```

Вкладка **SQL** показывает план **Catalyst** (оптимизатора).

Подводя итог, можно сказать, что самым простым и в 90% случаев достаточным средством для мониторинга и профилирования ситуации в Job, является **Spark UI**, который внутри себя собирает информацию, как на уровне конкретных JVM процессов, так и на уровне заданий.