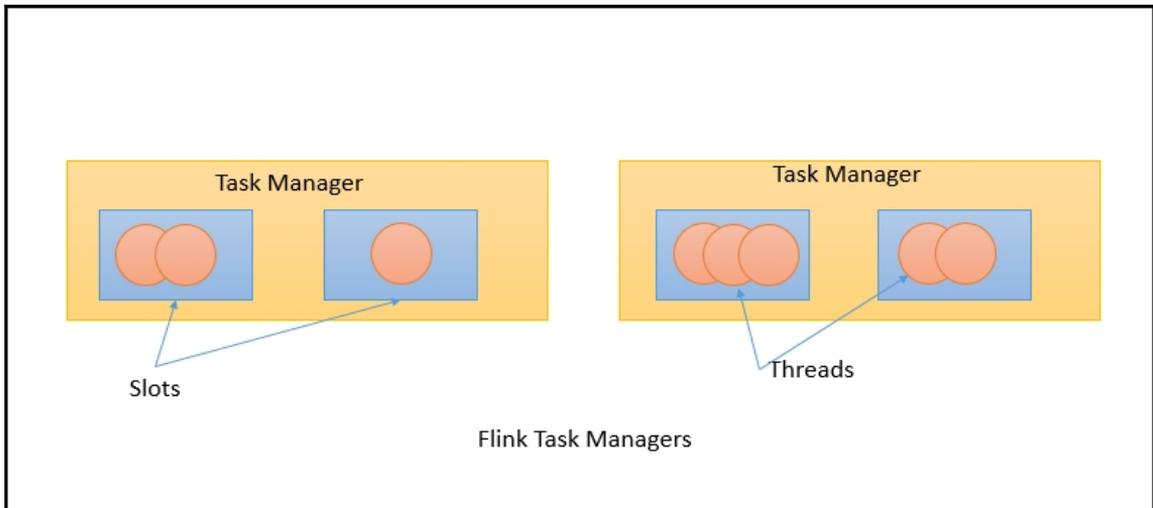
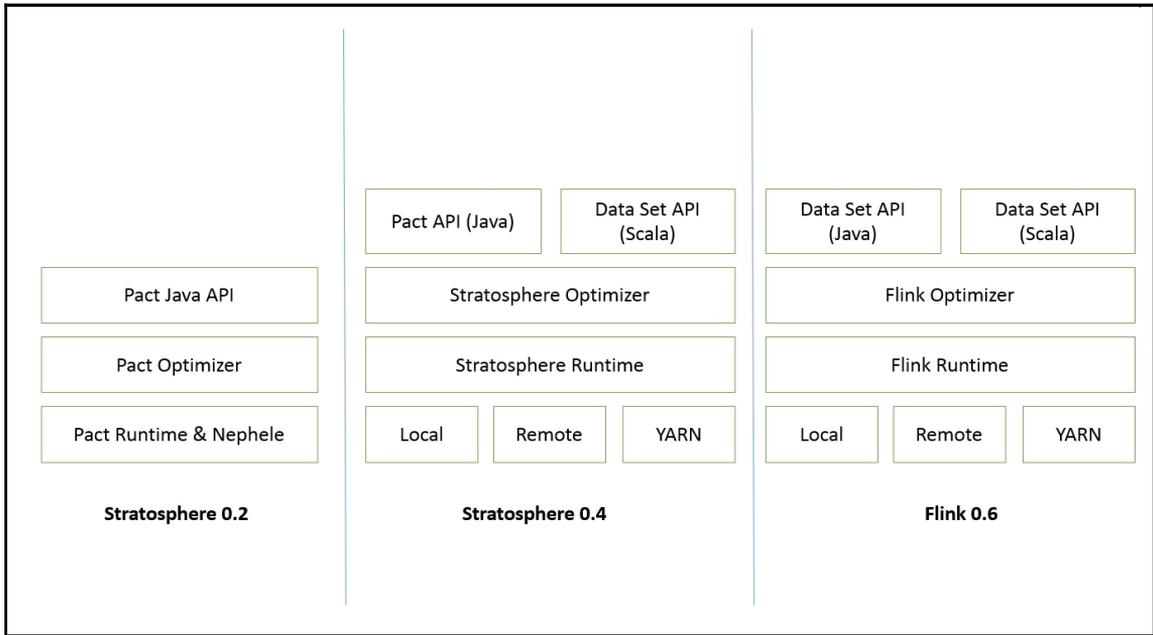


Chapter 1: Introduction to Apache Flink



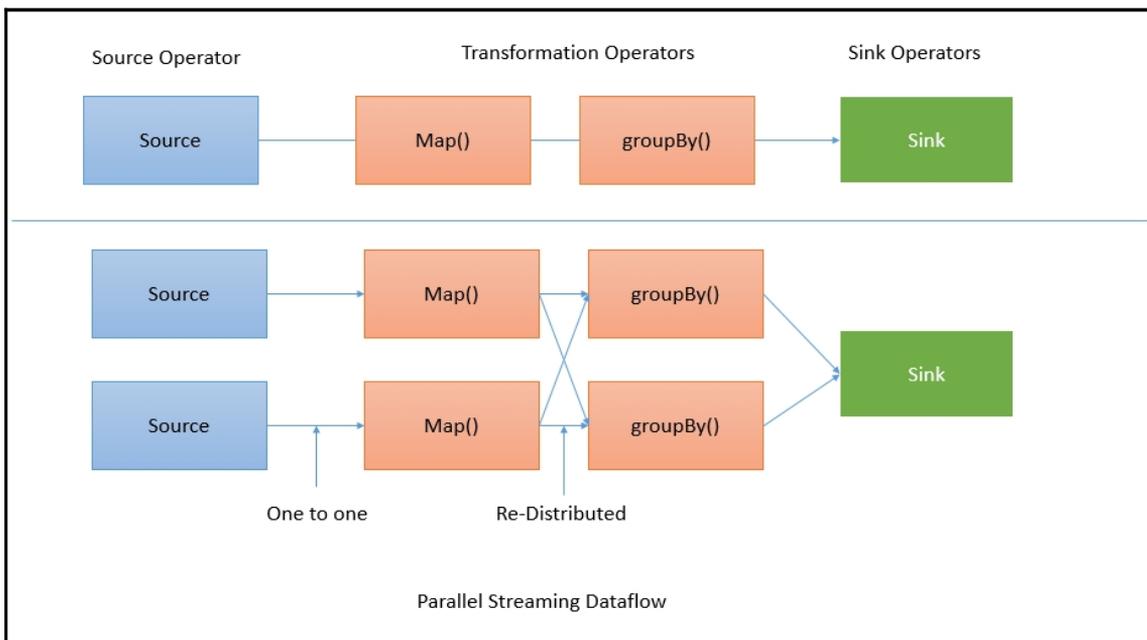
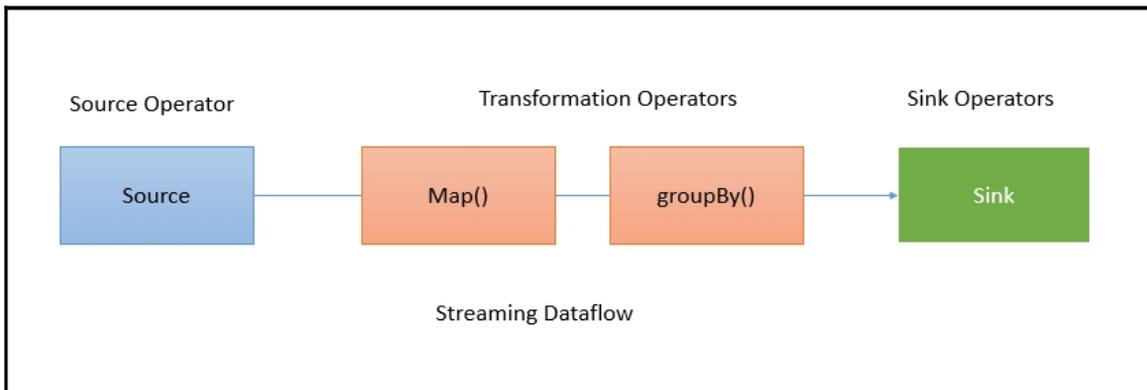
```

val text = env.readTextFile("input.txt") // Source

val counts = text.flatMap { _.toLowerCase.split("\\W+") filter { _.nonEmpty } }
    .map { (_, 1) }
    .groupBy(0)
    .sum(1) // Transformation

counts.writeAsCsv("output.txt", "\n", " ") // Sink

```



Binaries

	Scala 2.10	Scala 2.11
Hadoop 1.2.1	Download	
Hadoop 2.3.0	Download	Download
Hadoop 2.4.1	Download	Download
Hadoop 2.6.0	Download	Download
Hadoop 2.7.0	Download	Download

```
D:\>java -version
java version "1.8.0_92"
Java(TM) SE Runtime Environment (build 1.8.0_92-b14)
Java HotSpot(TM) 64-Bit Server VM (build 25.92-b14, mixed mode)
```

Apache Flink Dashboard Overview

Version: 1.0.3 Commit: f3a6b5f

	1
Task Managers	
	1
Task Slots	
	1
Available Task Slots	

Total Jobs

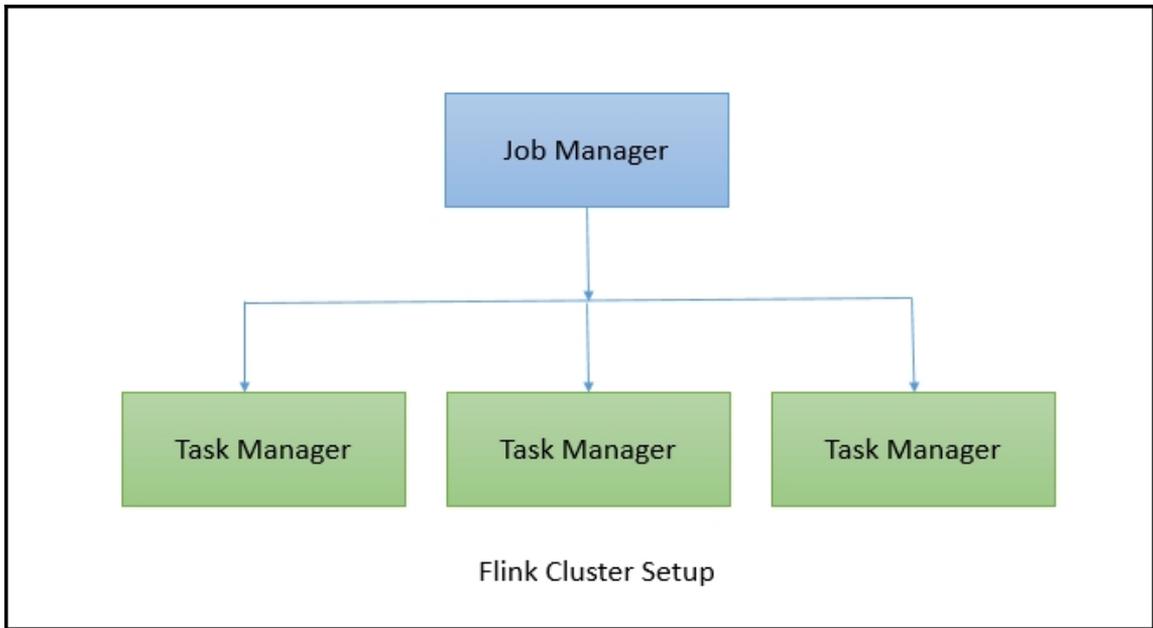
Running	0
Finished	0
Canceled	0
Failed	0

Running Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status
------------	----------	----------	----------	--------	-------	--------

Completed Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status
------------	----------	----------	----------	--------	-------	--------



Apache Flink Dashboard

Overview | Version: 1.1.4 | Commit: 8fb0fc8

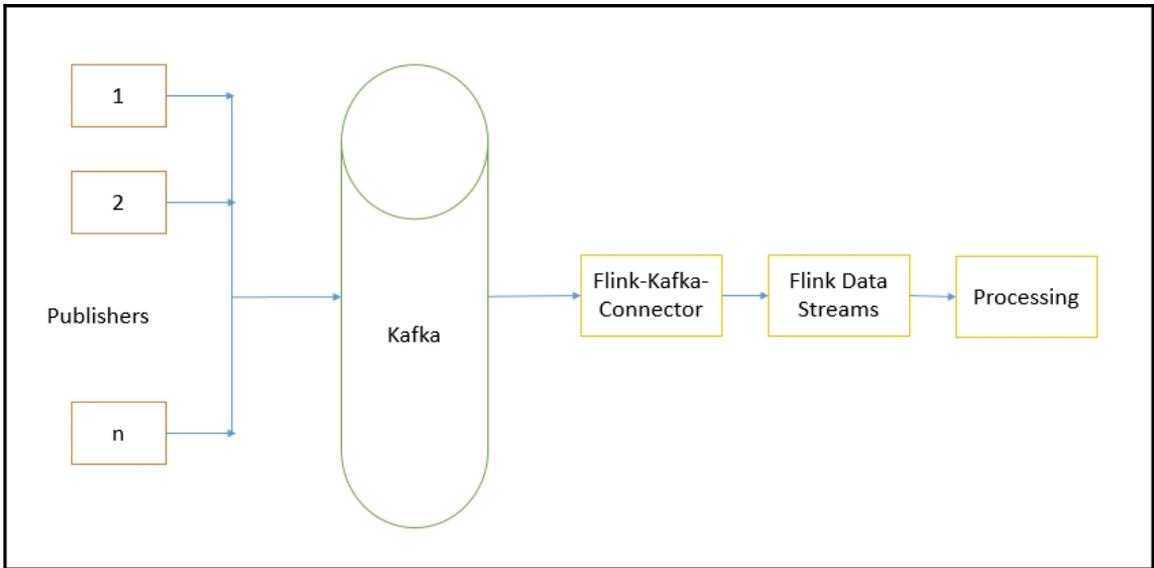
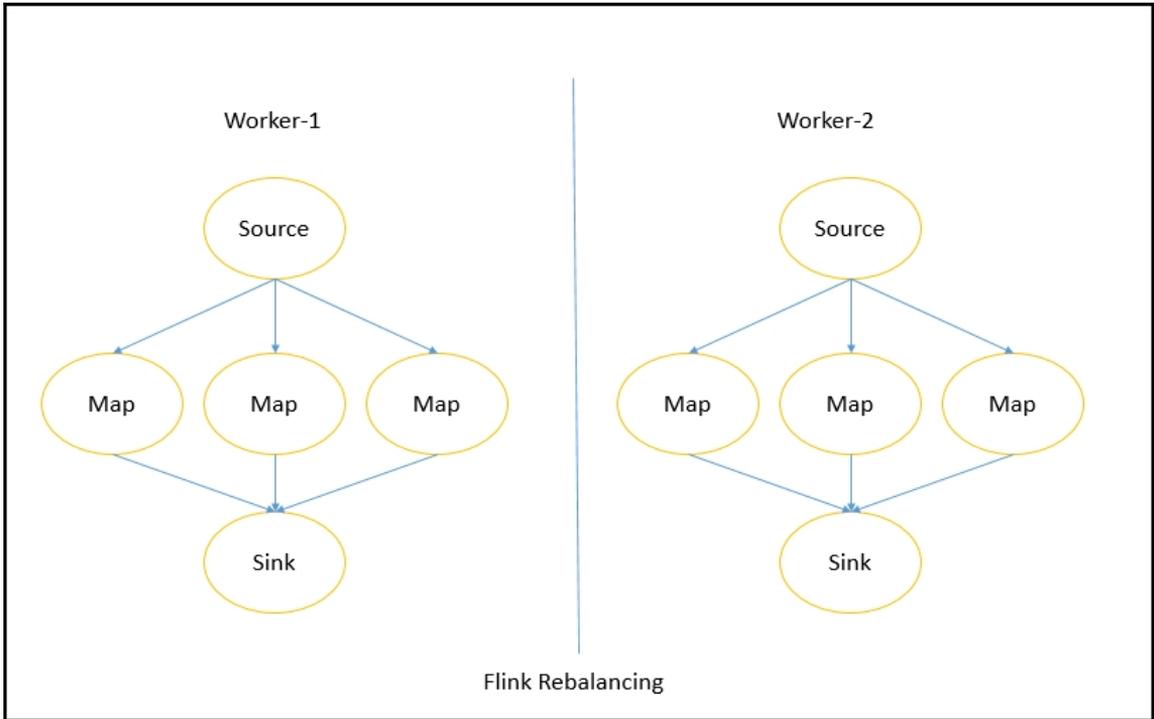
Icon	Count	Label
	1	Task Managers
	1	Task Slots
	1	Available Task Slots

Total Jobs	
Running	0
Finished	0
Canceled	0
Failed	0

Running Jobs							
Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status	

Completed Jobs							
Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status	

Chapter 2: Data Processing Using the DataStream API



HadoopTrendingTopics

Test OAuth

Details Settings **Keys and Access Tokens** Permissions

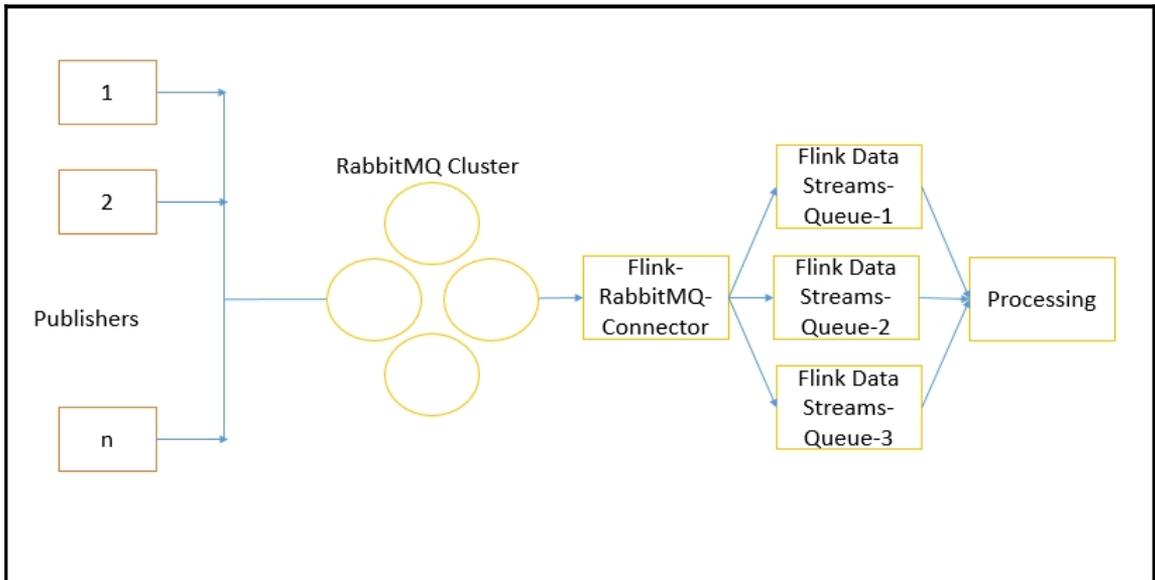
Application Settings

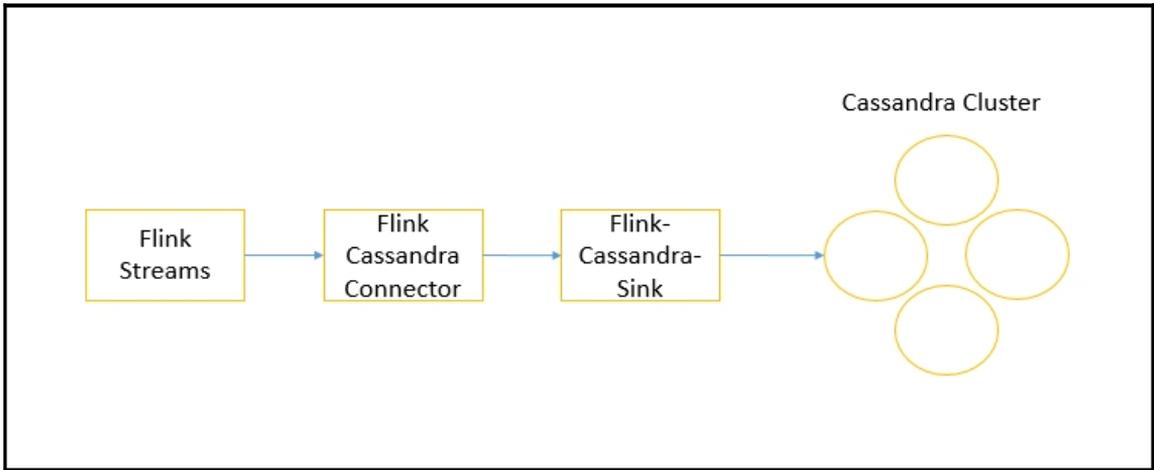
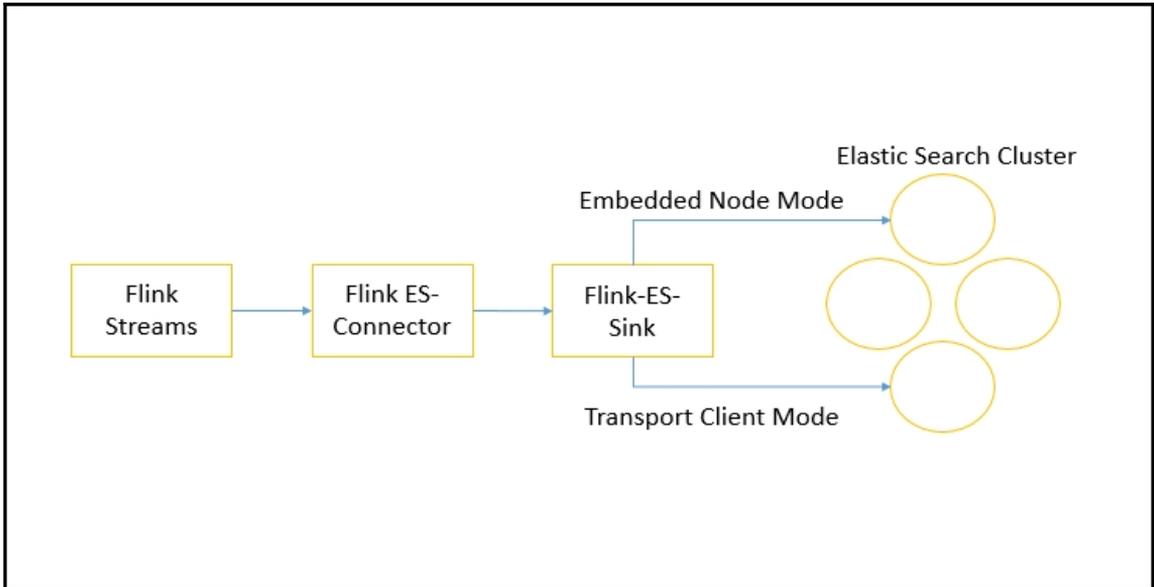
Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

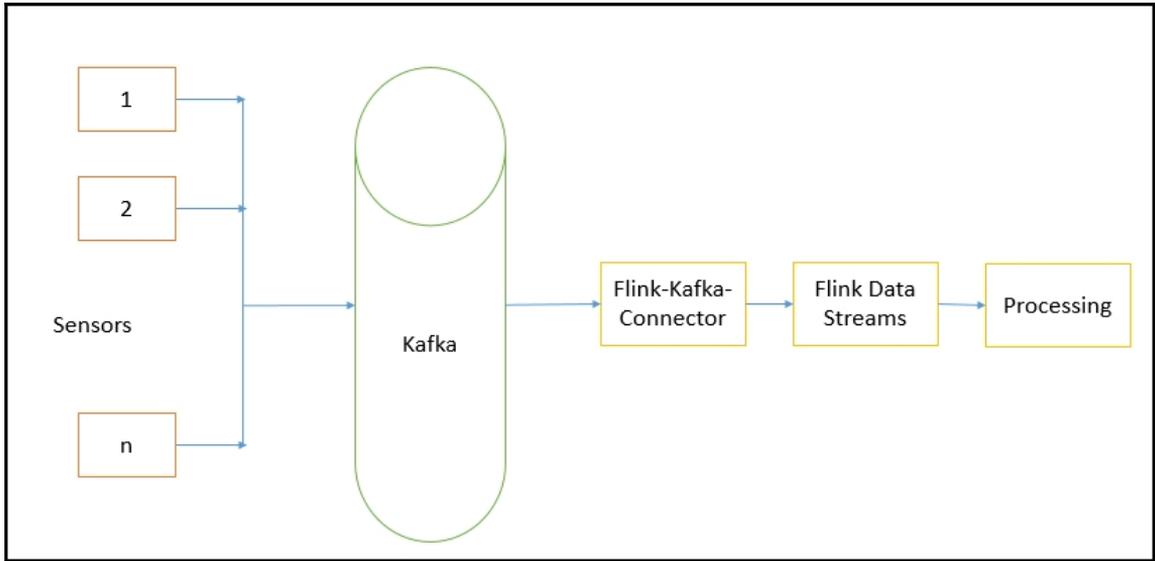
Consumer Key (API Key)	tP686cWPcJAdCsuF4Alv
Consumer Secret (API Secret)	HKP6yxcByJqHAypaGGII75Npcu5GkbrGMgOQRqIT8
Access Level	Read-only (modify app permissions)
Owner	HadoopTutorials
Owner ID	2825680861

Application Actions

Regenerate Consumer Key and Secret Change App Permissions







Obtain Execution Environment



Load Data from Source



Transform Data



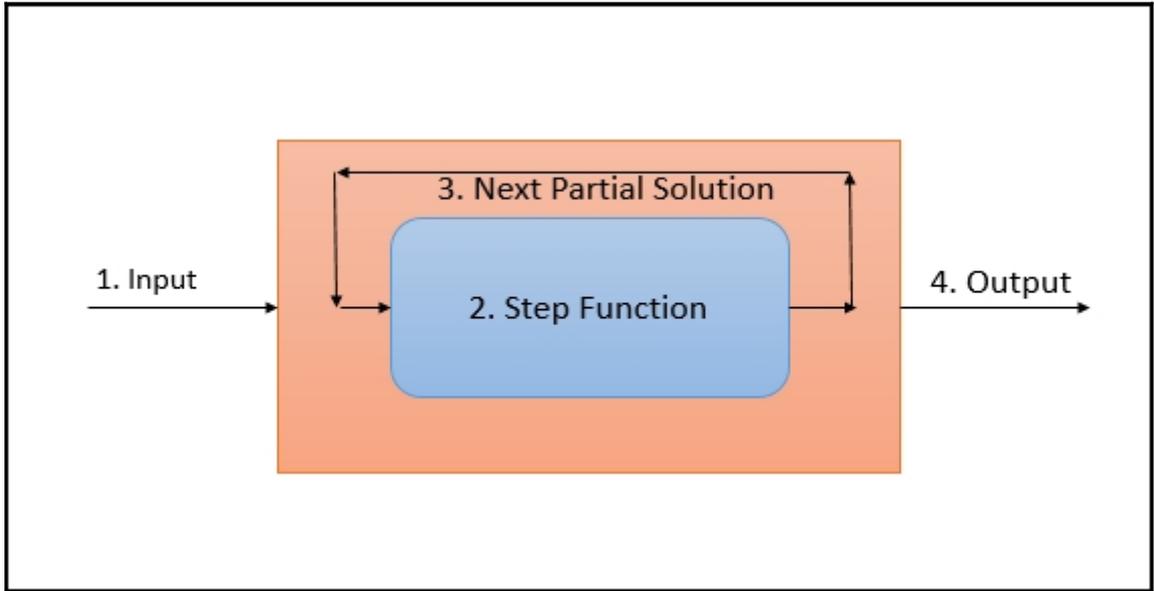
Transform Data

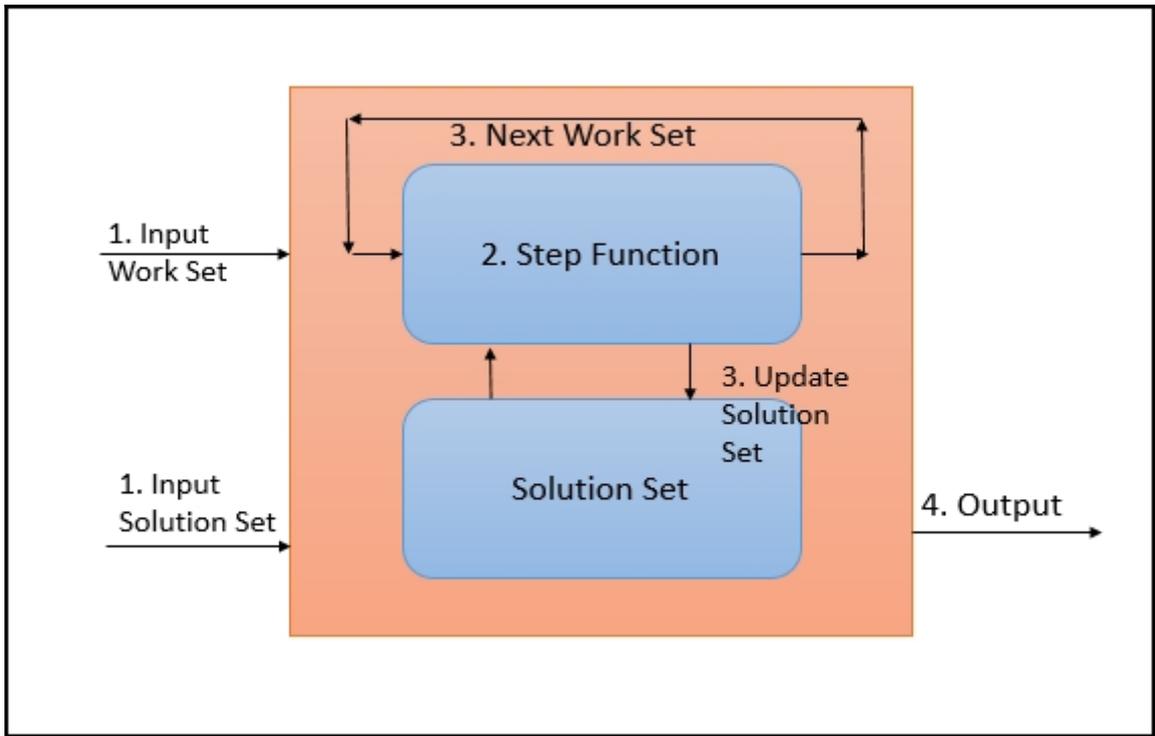


Store the Processed Data

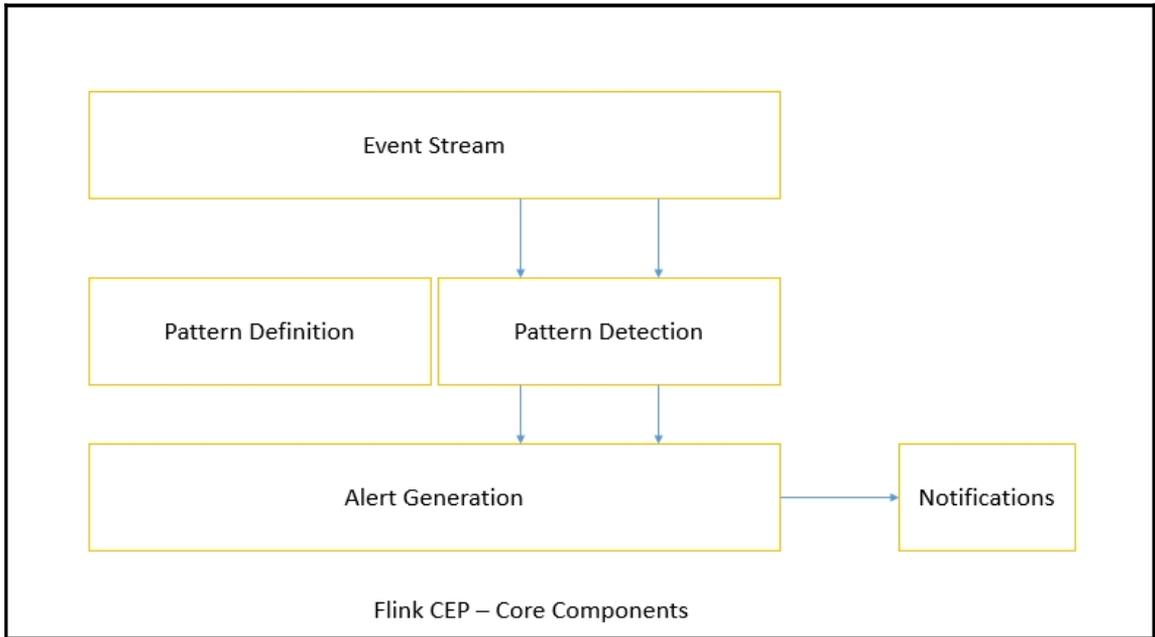
Anatomy of a Flink Program

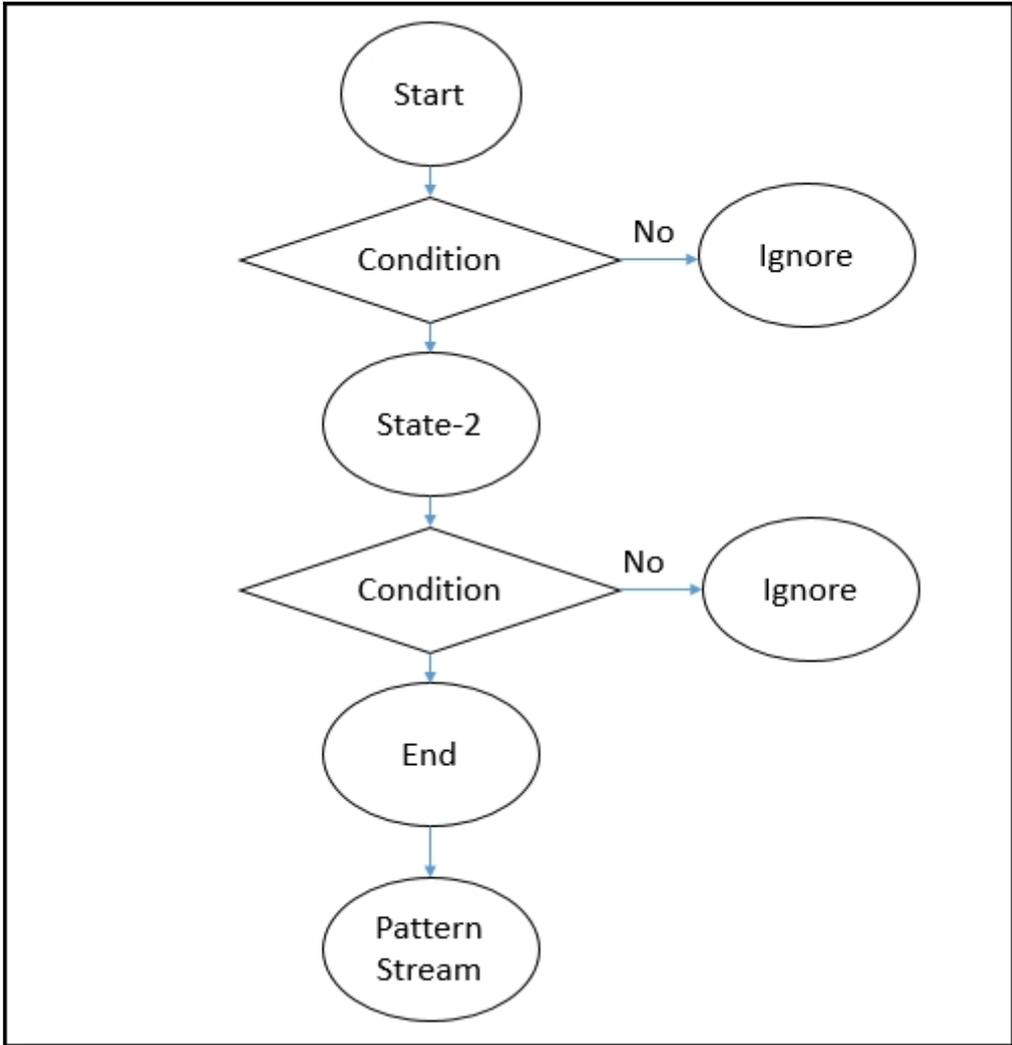
Chapter 3: Data Processing Using the Batch Processing API

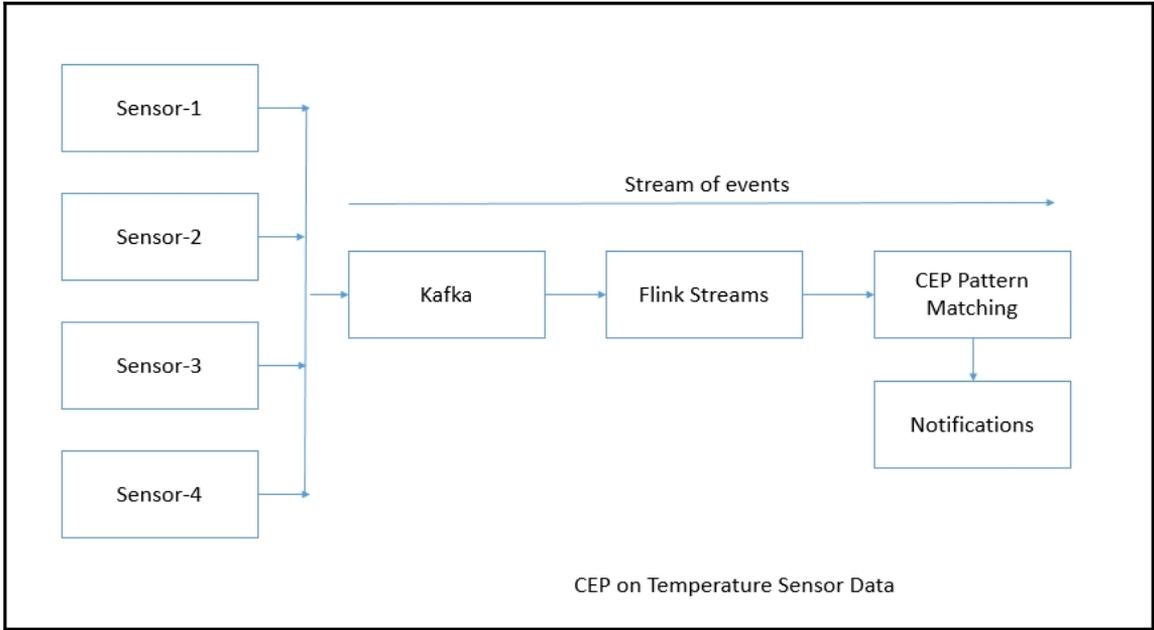




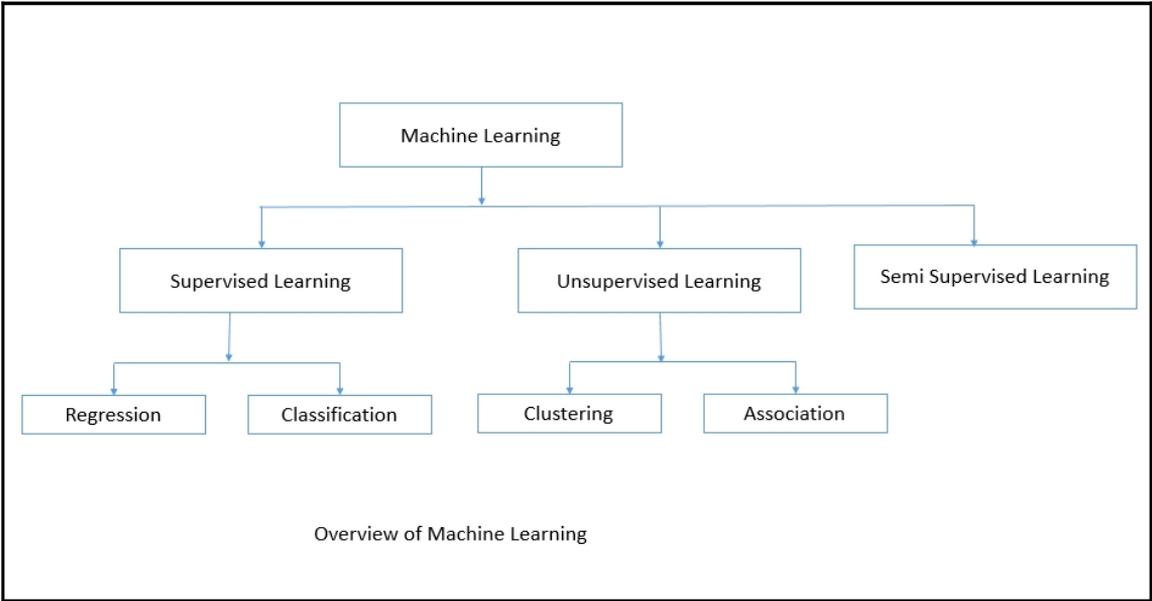
Chapter 5: Complex Event Processing



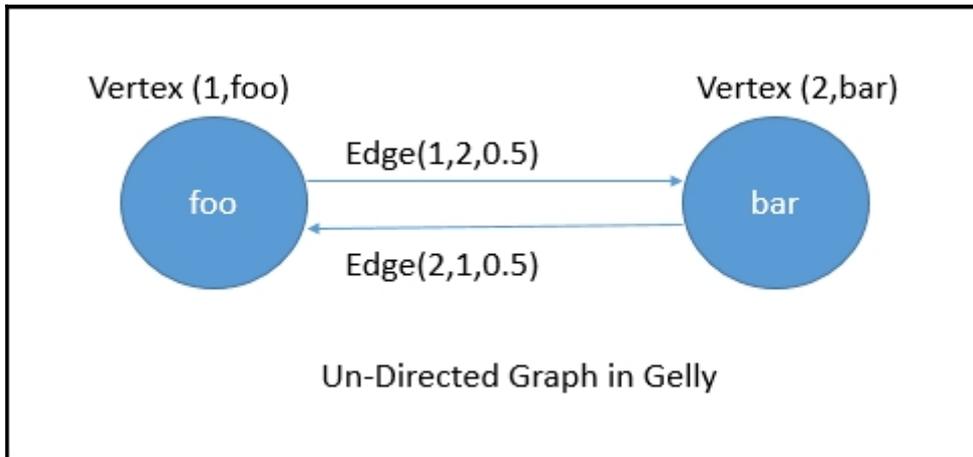
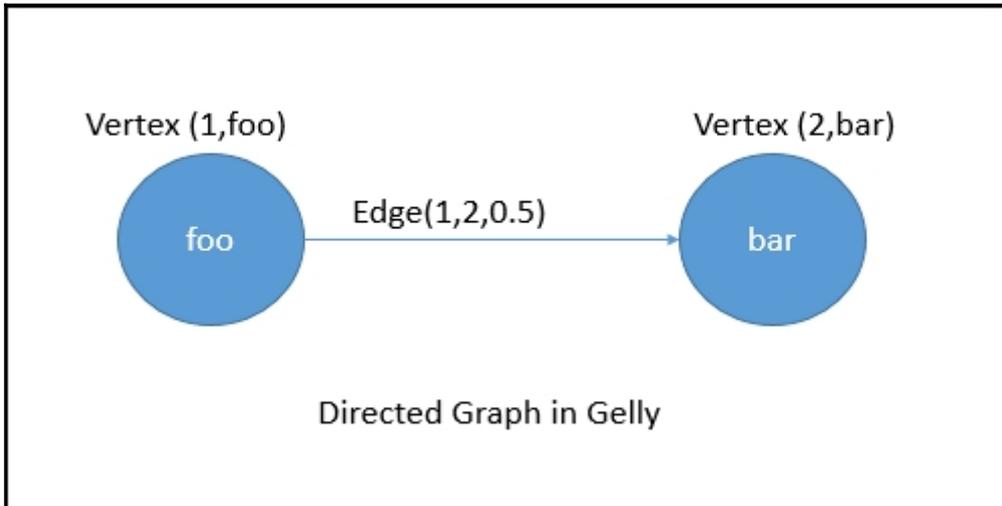


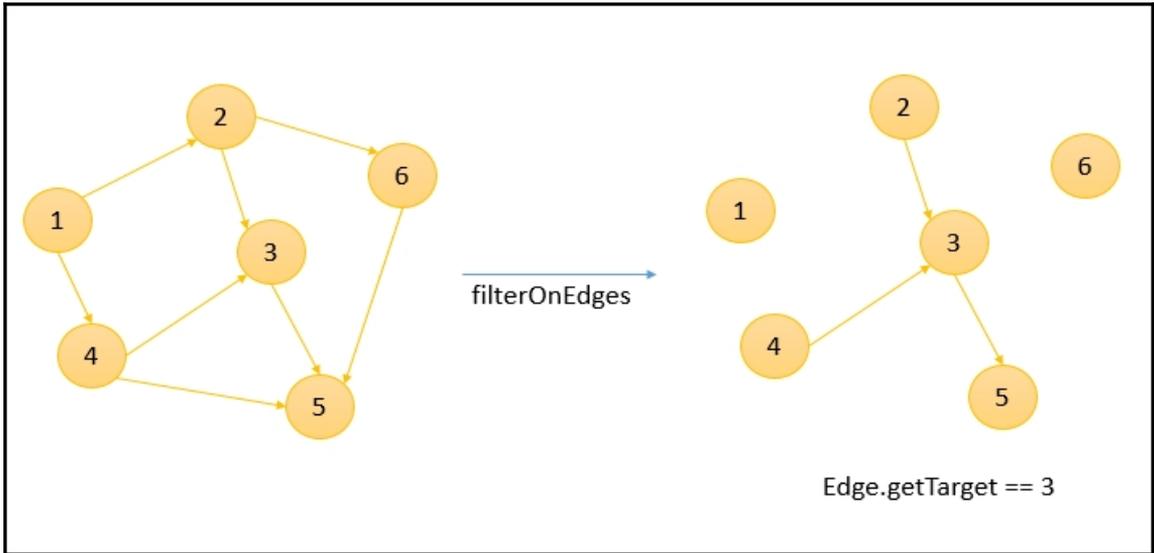
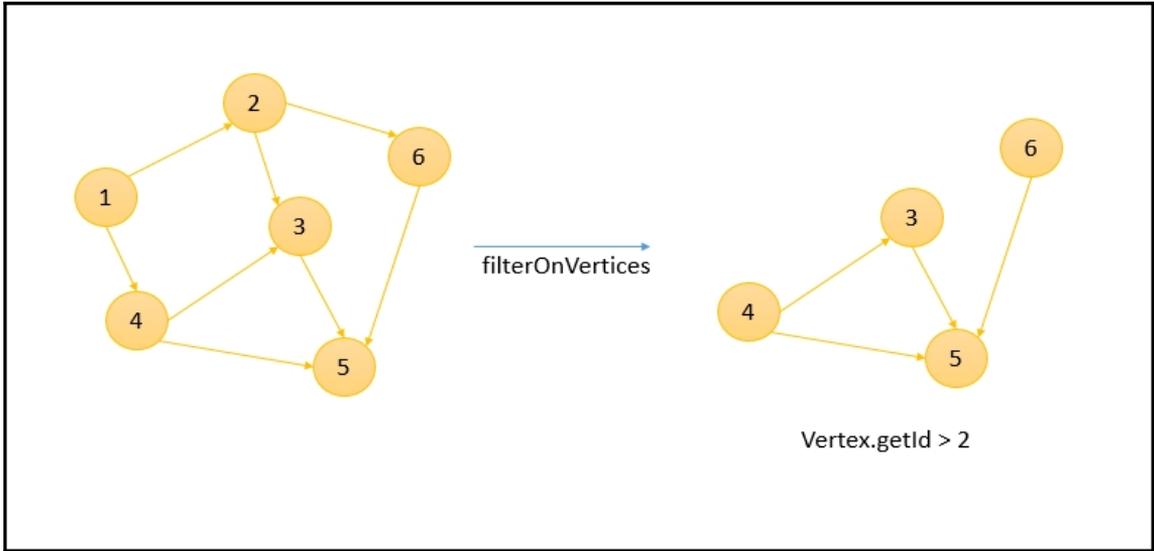


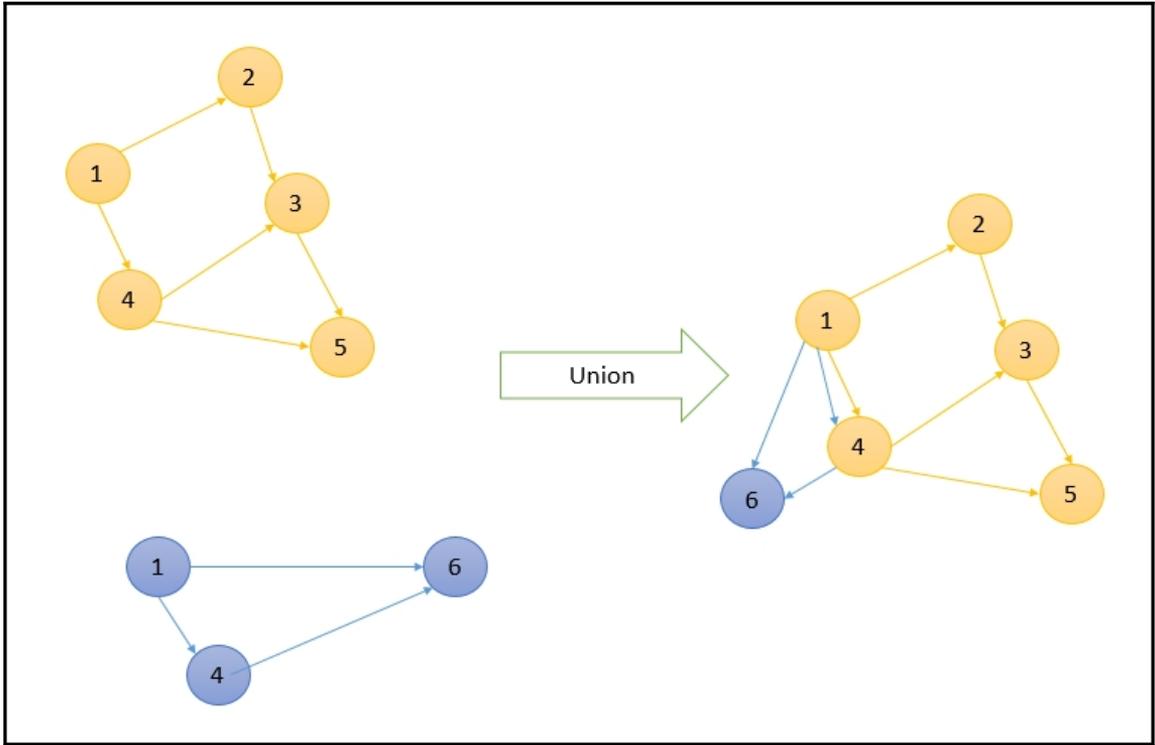
Chapter 6: Machine Learning Using FlinkML

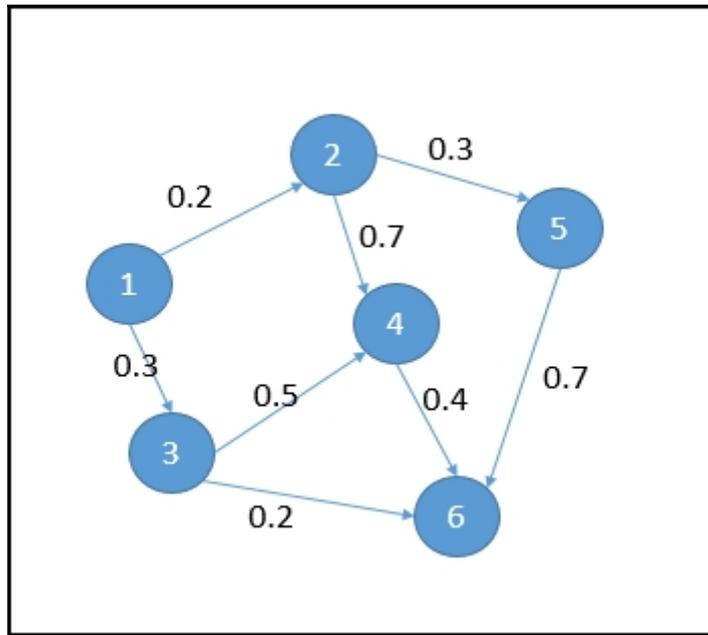


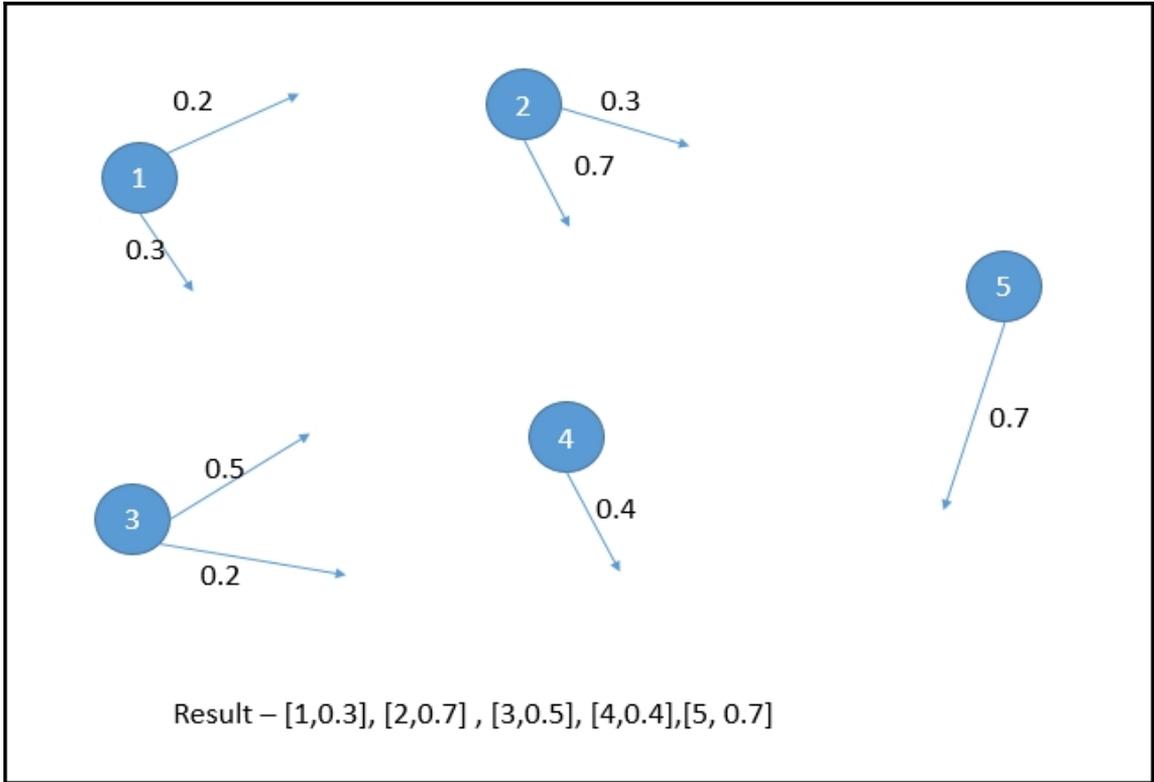
Chapter 7: Flink Graph API – Gelly











Chapter 8: Distributed Data Processing with Flink and Hadoop

Apache Flink Dashboard Overview | Version: 1.1.3 | Commit: 8e8d454

Task Managers: 2
Task Slots: 20
Available Task Slots: 20

Total Jobs:

Running	0
Finished	0
Canceled	0
Failed	0

Running Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status

Completed Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status

hadoop Application `application_1478079131011_0107` | Logged in as: dr.who

Kill Application | Application Overview

User: root
Name: Flink session with 2 TaskManagers
Application Type: Apache Flink
Application Tags:
Application Priority: 0 (Higher Integer value indicates higher priority)
YarnApplicationState: RUNNING: AM has registered with RM and started running.
Queue: default
FinalStatus Reported by AM: Application has not completed yet.
Started: Mon Nov 14 10:46:03 +0530 2016
Elapsed: 58mins, 10sec
Tracking URL: ApplicationMaster
Log Aggregation Status: NOT_START
Diagnostics:
Unmanaged Application: false
Application Node Label expression: <Not set>
AM container Node Label expression: <DEFAULT_PARTITION>

Application Metrics

Total Resource Preempted: <memory:0, vCores:0>
Total Number of Non-AM Containers Preempted: 0
Total Number of AM Containers Preempted: 0
Resource Preempted from Current Attempt: <memory:0, vCores:0>
Number of Non-AM Containers Preempted from Current Attempt: 0
Aggregate Resource Allocation: 26746984 MB-seconds, 10446 vcore-seconds

Application Overview

User: [root](#)
Name: Flink Application: org.apache.flink.examples.java.wordcount.WordCount
Application Type: Apache Flink
Application Tags:
Application Priority: 0 (Higher Integer value indicates higher priority)
YarnApplicationState: FINISHED
Queue: [default](#)
FinalStatus Reported by AM: SUCCEEDED
Started: Mon Nov 14 12:00:16 +0530 2016
Elapsed: 9sec
Tracking URL: [History](#)
Log Aggregation Status: SUCCEEDED
Diagnostics: Flink YARN Client requested shutdown
Unmanaged Application: false
Application Node Label expression: <Not set>
AM container Node Label expression: <DEFAULT_PARTITION>

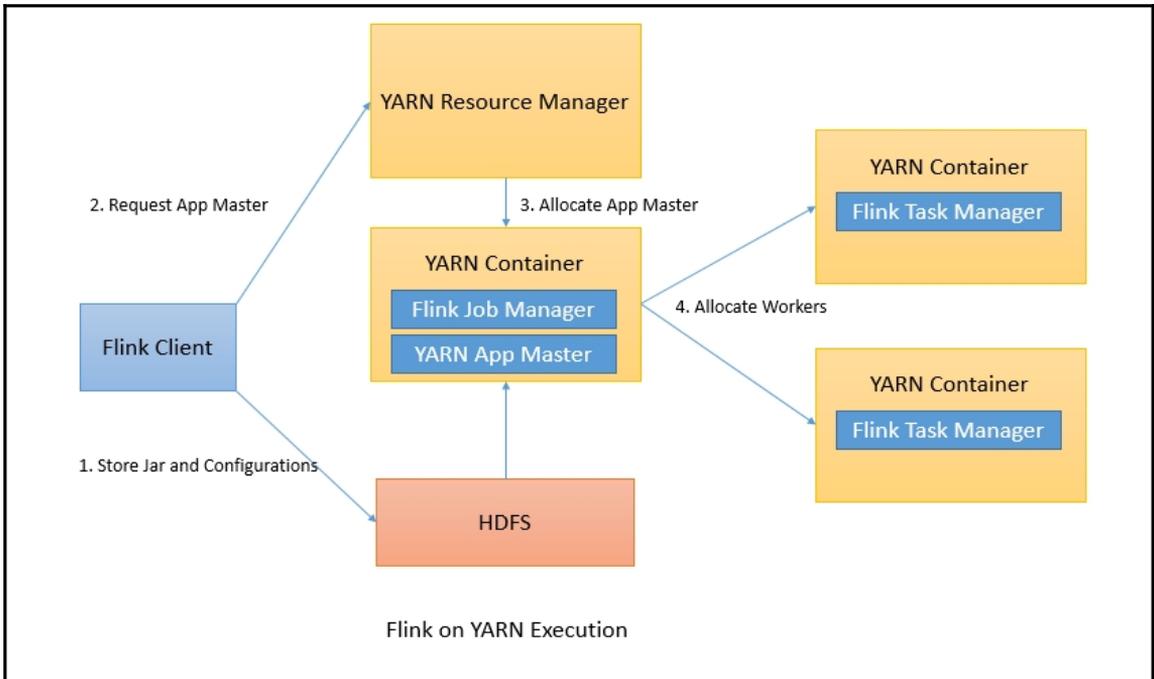
Application Metrics

Total Resource Preempted: <memory:0, vCores:0>
Total Number of Non-AM Containers Preempted: 0
Total Number of AM Containers Preempted: 0
Resource Preempted from Current Attempt: <memory:0, vCores:0>
Number of Non-AM Containers Preempted from Current Attempt: 0
Aggregate Resource Allocation: 48130 MB-seconds, 17 vcore-seconds

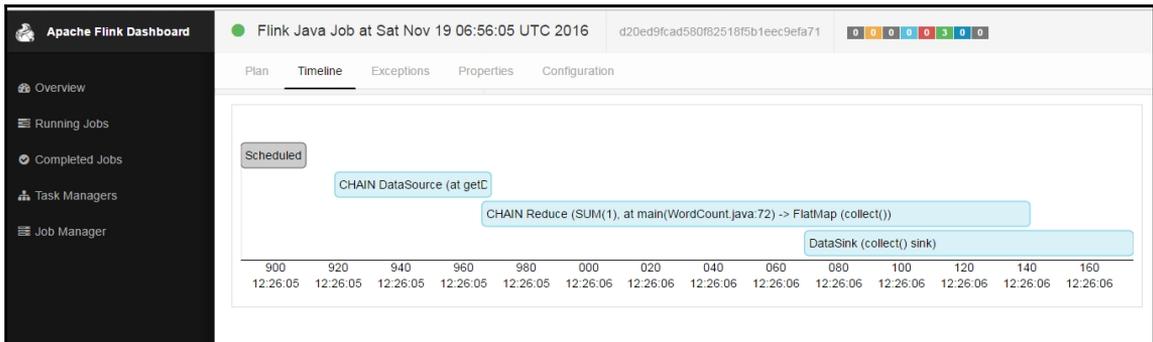
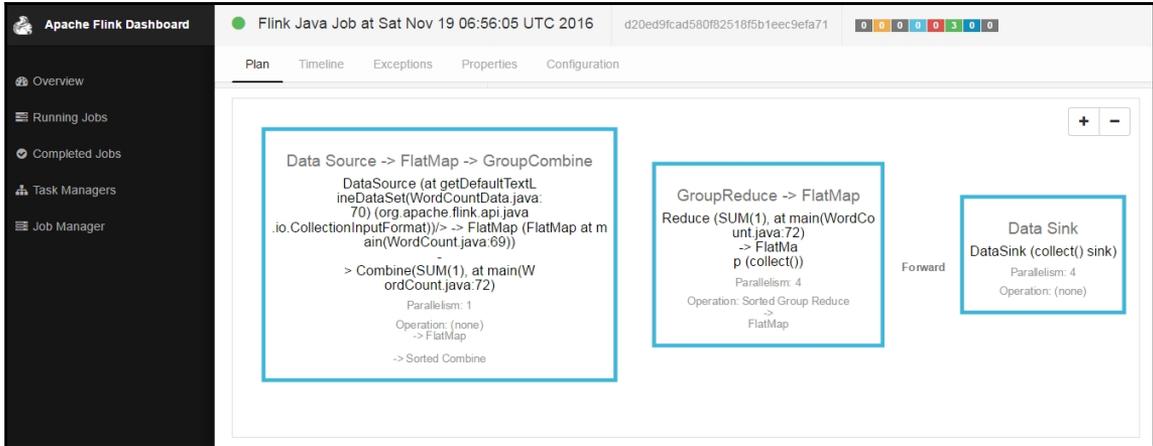
Show 20 entries Search:

Attempt ID	Started	Node	Logs	Blacklisted Nodes
appattempt_1478079131011_0108_000001	Mon Nov 14 12:00:16 +0550 2016	http://hdpdev005.pune-in0145.slb.com:8042	Logs	N/A

Showing 1 to 1 of 1 entries First Previous 1 Next Last



Chapter 9: Deploying Flink on Cloud



```

Sat Nov 19 08:34:05 UTC 2016: Using local tmp dir for staging files: /tmp/bdutil-20161119-083405-GdK
Sat Nov 19 08:34:05 UTC 2016: Using custom environment-variable file(s): bdutil_env.sh extensions/flink/flink_env.sh
Sat Nov 19 08:34:05 UTC 2016: Reading environment-variable file: ./bdutil_env.sh
Sat Nov 19 08:34:05 UTC 2016: Reading environment-variable file: extensions/flink/flink_env.sh
Sat Nov 19 08:34:05 UTC 2016: No explicit GCE_MASTER_MACHINE_TYPE provided; defaulting to value of GCE_MACHINE_TYPE: n1-standard-2
Delete cluster with following settings?
CONFIGBUCKET='bdutil-flink-bucket'
PROJECT='...'
GCE_IMAGE='https://www.googleapis.com/compute/v1/projects/debian-cloud/global/images/backports-debian-7-wheezy-v20160531'
GCE_ZONE='europe-west1-c'
GCE_NETWORK='default'
GCE_TAGS='bdutil'
PREEMPTIBLE_FRACTION=0.0
PREFIX='hadoop'
NUM_WORKERS=2
MASTER_HOSTNAME='hadoop-m'
WORKERS='hadoop-w-0 hadoop-w-1'
BDUTIL_GCS_STAGING_DIR='gs://bdutil-flink-bucket/bdutil-staging/hadoop-m'
(y/n) █

```

The screenshot shows the AWS Management Console interface. At the top, there are navigation tabs for 'Services' and 'Resource Groups'. Below this, the 'AWS services' section is active, with a search bar containing 'EMR'. A dropdown menu shows 'EMR' with the description 'Managed Hadoop Framework'. Below the search results, there are icons for 'EMR' and 'EC2', and a link to '> All services'.

The 'Build a solution' section is also visible, featuring several quick-start options:

- Launch a virtual machine**: With EC2, ~1 minutes
- Build a web app**: With Elastic Beanstalk, ~6 minutes
- Deploy a serverless microservice**: With Lambda, API Gateway, ~2 minutes
- Host a static website**: With S3, CloudFront, Route 53, ~5 minutes
- Create a backend for your mobile app**: With Mobile Hub, ~5 minutes
- Register a domain**: With Route 53, ~3 minutes

Services Resource Groups Tanmay Deshpande Mumbai Support

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

Logging [?]

S3 folder

Launch mode Cluster [?] Step execution [?]

Software configuration

Vendor Amazon

Release [?]

Applications

- Core Hadoop: Hadoop 2.7.3 with Ganglia 3.7.2, Hive 2.1.0, Hue 3.10.0, Mahout 0.12.2, Pig 0.16.0, and Tez 0.8.4
- HBase: HBase 1.2.3 with Ganglia 3.7.2, Hadoop 2.7.3, Hive 2.1.0, Hue 3.10.0, Phoenix 4.7.0, and ZooKeeper 3.4.8
- Presto: Presto 0.152.3 with Hadoop 2.7.3 HDFS and Hive 2.1.0 Metastore
- Spark: Spark 2.0.1 on Hadoop 2.7.3 YARN with Ganglia 3.7.2 and Zeppelin 0.6.2

Services Resource Groups Tanmay Deshpande Mumbai Support

EC2 Dashboard

Events

Tags

Reports

Limits

INSTANCES

Instances

Spot Requests

Reserved Instances

Dedicated Hosts

IMAGES

AMIs

Bundle Tasks

ELASTIC BLOCK STORE

Volumes

Snapshots

NETWORK & SECURITY

Security Groups

Elastic IPs

Placement Groups

Key Pairs

Create Security Group Actions

Filter by tags and attributes or search by keyword

Name	Group ID	Group Name	VPC ID	Description
<input type="checkbox"/>	sg-10dd5079	ElasticMapReduce-slave	vpc-31a80e58	Slave group for Elastic MapReduce created on 2016-11-1...
<input checked="" type="checkbox"/>	sg-11dd5078	ElasticMapReduce-master	vpc-31a80e58	Master group for Elastic MapReduce created on 2016-11-1...
<input type="checkbox"/>	sg-aac24fc3	default	vpc-31a80e58	default VPC security group

Security Group: sg-11dd5078

Description Inbound Outbound Tags

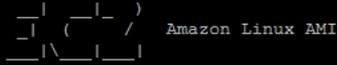
Edit

Type	Protocol	Port Range	Source
All TCP	TCP	0 - 65535	sg-10dd5079 (ElasticMapReduce-slave)
All TCP	TCP	0 - 65535	sg-11dd5078 (ElasticMapReduce-master)

```

Using username "hadoop".
Authenticating with public key "imported-openssh-key"
Last login: Sun Nov 20 06:24:42 2016

```



```

https://aws.amazon.com/amazon-linux-ami/2016.09-release-notes/
8 package(s) needed for security, out of 13 available
Run "sudo yum update" to apply all updates.

```

```

EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::M          M:::::M R:::::R
EE:::::EEEEEEEEEEEEE M:::::M          M:::::M R:::::RRRRRRR:::::R
E:::E          EEEEE M:::::M          M:::::M RR:::R          R:::R
E:::E          M:::::M:::M M:::M:::::M R:::R          R:::R
E:::::EEEEEEEEEE M:::::M M:::M M:::::M M:::::M R:::::RRRRRR:::::R
E::::::::::::::::::E M:::::M M:::M:::M M:::::M R:::::RR
E:::::EEEEEEEEEE M:::::M M:::::M M:::::M R:::::RRRRRR:::::R
E:::E          M:::::M M:::::M M:::::M R:::R          R:::R
E:::E          EEEEE M:::::M MMM M:::::M R:::R          R:::R
EE:::::EEEEEEEE:::E M:::::M          M:::::M R:::R          R:::R
E::::::::::::::::::E M:::::M          M:::::M RR:::R          R:::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRR          RRRRRR

```

```
[hadoop@ip-172-31-2-68 ~]$
```



All Applications

Cluster

- About
- Nodes
- Node Labels
- Applications
- NEW
- NEW SAVING
- SUBMITTED
- ACCEPTED
- RUNNING
- FINISHED
- FAILED
- KILLED

Scheduler

Tools

Cluster Metrics													
Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	VCores Reserved	Active Nodes	Decommissioning Nodes	Decommission Nodes
2	0	0	2	0	0 B	12 GB	0 B	0	8	0	2	0	0

Scheduler Metrics													
Scheduler Type		Scheduling Resource Type				Minimum Allocation				M			
Capacity Scheduler		[MEMORY]				<memory:32, vCores:1>				<memory:6144, vCo			

Show 20 entries													
ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus					
application_1479621657204_0002	hadoop	Flink Application: org.apache.flink.examples.java.wordcount.WordCount	Apache Flink	default	Sun Nov 20 12:11:34 +0550 2016	Sun Nov 20 12:11:47 +0550 2016	FINISHED	SUCCEEDED					
application_1479621657204_0001	hadoop	Flink session with 2 TaskManagers	Apache Flink	default	Sun Nov 20 12:10:03 +0550 2016	Sun Nov 20 12:10:11 +0550 2016	FAILED	FAILED					

Showing 1 to 2 of 2 entries

Apache Flink Dashboard Overview Version: 1.1.3 Commit: 8e8d454

Task Managers 2

Task Slots 8

Available Task Slots 8

Total Jobs	
Running	0
Finished	0
Canceled	0
Failed	0

Running Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status

Completed Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status

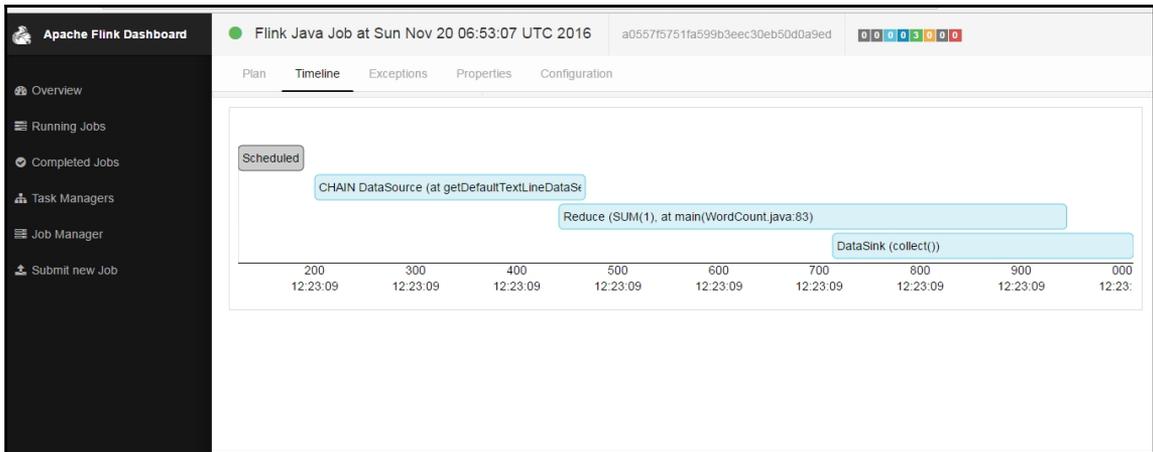
ec2-35-154-40-129.ap-south-1.compute.amazonaws.com:20888/proxy/application_1479621657204_0004/#/jobs/a0557f5751fa599b3eec30eb50d0a9ec

Apache Flink Dashboard Flink Java Job at Sun Nov 20 06:53:07 UTC 2016 a0557f5751fa599b3eec30eb50d0a9ed 0 0 0 0 3 0 0 0

Plan Timeline Exceptions Properties Configuration

```

graph LR
    A["Data Source -> FlatMap -> GroupCombine  
DataSource (at getDefaultTextL  
inDataSet(WordCountData.java:  
70) (org.apache.flink.api.java  
io.CollectionInputFormat)) -> FlatMap (FlatMap at m  
ain(WordCount.java:80))  
-> Combine(SUM(1), at main(W  
ordCount.java:83))  
Parallelism: 1  
Operation: none  
-> FlatMap  
-> Sorted Combine"]
    A --> B["GroupReduce  
Reduce(SUM(1), at main(WordCo  
unt.java:83))  
Parallelism: 2  
Operation: Sorted Group Reduce"]
    B -- Forward --> C["Data Sink:  
DataSink (collect())  
Parallelism: 1  
Operation: none"]
  
```



Services Resource Groups

Tanmay Deshpande Mumbai

Create Cluster - Quick Options

[Go to advanced options](#) [Click Here](#)

General Configuration

Cluster name

Logging [?]

S3 folder

Launch mode Cluster [?] Step execution [?]

Software configuration

Vendor Amazon

Release [?]

Applications

- Core Hadoop: Hadoop 2.7.3 with Ganglia 3.7.2, Hive 2.1.1, Hue 3.11.0, Mahout 0.12.2, Pig 0.16.0, and Tez 0.8.4
- HBase: HBase 1.2.3 with Ganglia 3.7.2, Hadoop 2.7.3, Hive 2.1.1, Hue 3.11.0, Phoenix 4.7.0, and ZooKeeper 3.4.9
- Presto: Presto 0.152.3 with Hadoop 2.7.3 HDFS and Hive 2.1.1 Metastore
- Spark: Spark 2.1.0 on Hadoop 2.7.3 YARN with Ganglia 3.7.2 and Zeppelin 0.6.2

Services Resource Groups Tanmay Deshpande Mumbai

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps

- Step 2: Hardware
- Step 3: General Cluster Settings
- Step 4: Security

Software Configuration

Vendor Amazon

Release

<input checked="" type="checkbox"/> Hadoop 2.7.3	<input type="checkbox"/> Zeppelin 0.6.2	<input type="checkbox"/> Tez 0.8.4
<input checked="" type="checkbox"/> Flink 1.1.4	<input type="checkbox"/> Ganglia 3.7.2	<input type="checkbox"/> HBase 1.2.3
<input checked="" type="checkbox"/> Pig 0.16.0	<input checked="" type="checkbox"/> Hive 2.1.1	<input type="checkbox"/> Presto 0.157.1
<input type="checkbox"/> ZooKeeper 3.4.9	<input type="checkbox"/> Sqoop 1.4.6	<input type="checkbox"/> Mahout 0.12.2
<input checked="" type="checkbox"/> Hue 3.11.0	<input type="checkbox"/> Phoenix 4.7.0	<input type="checkbox"/> Oozie 4.3.0
<input type="checkbox"/> Spark 2.1.0	<input type="checkbox"/> HCatalog 2.1.1	

Edit software settings (optional)

Enter configuration Load JSON from S3

```
classification=config-file-name,properties=[myKey1=myValue1,myKey2=myValue2]
```

Add steps (optional)

Step type

Welcome to the Google Cloud SDK!

To help improve the quality of this product, we collect anonymized usage data and anonymized stacktraces when crashes are encountered.. You may choose to opt out of this collection now (by choosing 'N' at the below prompt), or at any time in the future by running the following command:
gcloud config set disable_usage_reporting true

Do you want to help improve the Google Cloud SDK (Y/n)? Y

Your current Cloud SDK version is: 135.0.0
The latest available version is: 135.0.0

Components			
Status	Name	ID	Size
Not Installed	App Engine Go Extensions	app-engine-go	47.3 MiB
Not Installed	Cloud Datastore Emulator	cloud-datastore-emulator	15.4 MiB
Not Installed	Cloud Datastore Emulator (Legacy)	gcd-emulator	38.1 MiB
Not Installed	Cloud Pub/Sub Emulator	pubsub-emulator	16.3 MiB
Not Installed	Google Container Registry's Docker credential helper	docker-credential-gcr	2.2 MiB
Not Installed	gcloud Alpha Commands	alpha	< 1 MiB
Not Installed	gcloud Beta Commands	beta	< 1 MiB
Not Installed	gcloud app Java Extensions	app-engine-java	124.4 MiB
Not Installed	gcloud app Python Extensions	app-engine-python	7.2 MiB
Not Installed	kubectl	kubectl	15.9 MiB
Installed	BigQuery Command Line Tool	bq	< 1 MiB
Installed	Cloud SDK Core Libraries	core	5.1 MiB
Installed	Cloud Storage Command Line Tool	gsutil	2.8 MiB
Installed	Default set of gcloud commands	gcloud	

To install or remove components at your current SDK version [135.0.0], run:

```
$ gcloud components install COMPONENT_ID  
$ gcloud components remove COMPONENT_ID
```

To update your SDK installation to the latest version [135.0.0], run:

```
$ gcloud components update
```

Modify profile to update your \$PATH and enable shell command completion? (Y/n)?

← Create a bucket

Name ?

Must be unique across Cloud Storage. **Privacy:** Do not include sensitive information in your bucket name. Others can discover your bucket name if it matches a name they're trying to use.

bdutil-flink-bucket

Default storage class ?

[Learn about pricing](#)

- Multi-Regional**
Use to stream videos and host hot web content.
Best for data accessed frequently around the world.
- Regional**
Use to store data and run data analytics.
Best for data accessed frequently in one part of the world.
- Nearline**
Use to store rarely accessed documents.
Best for data accessed less than once per month.
- Coldline**
Use to store very rarely accessed documents.
Best for data accessed less than once per year.

Multi-Regional location

Redundant across 2+ regions within your selected location.

United States ▼

Create

Cancel

```

[tdeshpande2@dev-instance-1 bduutil-master]$ sudo ./bduutil -e extensions/flink/flink_env.sh deploy
Sat Nov 19 05:01:18 UTC 2016: Using local tmp dir for staging files: /tmp/bduutil-20161119-050118-FpC
Sat Nov 19 05:01:18 UTC 2016: Using custom environment-variable file(s): bduutil_env.sh extensions/flink/flink_env.sh
Sat Nov 19 05:01:18 UTC 2016: Reading environment-variable file: ./bduutil_env.sh
Sat Nov 19 05:01:18 UTC 2016: Reading environment-variable file: extensions/flink/flink_env.sh
Sat Nov 19 05:01:18 UTC 2016: No explicit GCE_MASTER_MACHINE_TYPE provided; defaulting to value of GCE_MACHINE_TYPE: n1-standard
Deploy cluster with following settings?
CONFIGBUCKET='bduutil-flink-bucket'
PROJECT='gcp-projects-123456789'
GCE_IMAGE='https://www.googleapis.com/compute/v1/projects/debian-cloud/global/images/backports-debian-7-wheezy-v20160531'
GCE_ZONE='europe-west1-d'
GCE_NETWORK='default'
GCE_TAGS='bduutil'
PREEMPTIBLE_FRACTION=0.0
PREFIX='hadoop'
NUM_WORKERS=2
MASTER_HOSTNAME='hadoop-m'
WORKERS='hadoop-w-0 hadoop-w-1'
BDUTIL_GCS_STAGING_DIR='gs://bduutil-flink-bucket/bduutil-staging/hadoop-m'
(y/n) █

```

Apache Flink Dashboard

Overview

- [Overview](#)
- [Running Jobs](#)
- [Completed Jobs](#)
- [Task Managers](#)
- [Job Manager](#)

	2
	Task Managers
	4
	Task Slots
	4
	Available Task Slots

Total Jobs

Running	0
Finished	1
Canceled	0
Failed	1

Running Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status