

# APACHE NIFI

## THE COMPLETE **GUIDE**

# WHAT IS DATA FLOW, DATA PIPELINE & ETL?

Let's start with the basic terms in the data world.



## WHAT IS DATA FLOW?

- X Data Flow → Moving data/content from Source to Destination
- X Data can be CSV, JSON, XML, Logs, HTTP data, Image, Videos, Telemetry data, etc..





## WHAT IS DATA PIPELINE?

- X Data Pipeline → Movement and Transformation of data/content from Source to Destination





## WHAT IS ETL?

X ETL → E - Extract, T - Transform, L - Load



## ETL VS DATA PIPELINE

- X ETL → Batch
- X Data Pipeline → Batch / Stream

# WHY WE SHOULD USE A FRAMEWORK FOR DATA FLOW

Considerations while building your own Data Flow / Data Pipeline.



SIMPLE PROBLEM TO SOLVE !





## FOUR V's

- X Volume -> refers to the vast amounts of data generated every second.
- X Velocity -> refers to the speed at which new data is generated and the speed at which data moves around.
- X Variety -> refers to the different types of data we can now use.
- X Veracity - refers to the messiness or trustworthiness of the data.



## CONSIDERATIONS

- X Support for multiple data format – CSV, JSON, Plaintext, Images, Videos, etc.
- X Support for various types of sources and destinations – FTP, HTTP, SQL Databases, NoSQL Databases, Search Engines, Cache Server, etc.
- X Scalable and Reliable for large volume and high-velocity data.
- X You should also consider Data Cleansing and Data Validation logics.

Welcome to



A robust open-source Data Ingestion and  
Distribution framework.



WHAT IS  
APACHE NIFI?



Apache NiFi supports powerful and scalable directed graphs of data routing, transformation, and system mediation logic.



NiFi was built to automate the flow of data between systems. It can propagate any data content from any source to any destination.

# INSTALLING APACHE NIFI IN A MAC / LINUX

# INSTALLING APACHE NIFI IN A WINDOWS MACHINE

# NiFi USER INTERFACE

# CORE NiFi TERMINOLOGIES



## FLOW-BASED PROGRAMMING (FBP)

- X NiFi is based on Flow Based Programming (FBP).
- X Flow Based Programming (FBP) is a programming paradigm that defines applications as networks of "black box" processes, which exchange data across predefined connections by message passing, where the connections are specified externally to the processes. These black box processes can be reconnected endlessly to form different applications without having to be changed internally. FBP is thus naturally component-oriented.



## How NiFi WORKS?

- X NiFi consists of atomic elements which can be combined to groups to build a data flow.
- X NiFi consists of Processors and Process Groups.



## WHAT IS A PROCESSOR?

- X Processors are atomic elements in NiFi which can do some simple tasks.
- X At the time of recording this video, NiFi has 280+ Processors.
- X Each Processor in NiFi is unique in its own way.
- X NiFi have processors for almost anything.



## WHAT A PROCESSOR CAN DO?

- X We have tons of Data Source and Data Sink Processors.
- X The Source and Sink can be anything - SQL, NoSQL, Search Engine, Cache Server, Messaging Queue, AWS Entities, etc.
- X NiFi have processors for all your data needs.
- X NiFi supports custom processors as well.



## WHAT IS A FLOWFILE?

- X FlowFile = Actual Data -> CSV, JSON, XML, Plaintext, SQL, Binary, etc.
- X FlowFile - Abstraction of data in NiFi.
- X A Processor can generate new Flow File by processing an existing Flow File or ingesting new data from any source.



## WHAT IS A CONNECTION?

- X In NiFi all processors can be connected with each other to create a data processing flow.
- X Processors are linked via connections.
- X Each connection will act as a Queue for FlowFiles.



## WHAT IS A PROCESS GROUP?

- X Set of Processors can be combined in NiFi to form a Process Group.
- X Process Groups helps to maintain large and complex data flow.
- X Input and Output ports are used to move data between Process Groups.



## WHAT IS A CONTROLLER SERVICE?

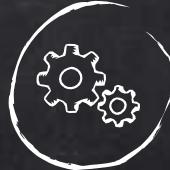
- X A Controller Service is a shared service that can be used by a Processor.
- X Controller Service can hold DB Connection details.
- X We can create Controller Service for CSV Reader, JSON Writer and many more.

# MORE ON FLOWFILES



## COMPONENTS OF FlowFILE

- X A FlowFile contains data.
- X A FlowFile is composed two components
  - Content
  - Attributes



## CONTENT & ATTRIBUTES OF FlowFILE

### X Content:

- That's the actual content of the FlowFile.
  - It's the actual content of a file you would read using GetFile, GetHTTP, etc.

### X Attributes:

- These are the metadata from the FlowFile.
- Contains information about the content:
  - When was it created, it's name, where is it from, what does it represent.



## CONSTRUCT USING NiFi FLOWFILES

- X A processor can (either or both)
  - Add, update, or remove attributes of a FlowFile.
  - Change content of a FlowFile.
  
- X Update the Attributes or Content or both using various processors available in NiFi to design your Dataflow.



## LIFE CYCLE OF FLOWFILE

- ✗ FlowFiles are persisted in the disk.
- ✗ FlowFiles are passed-by-reference.
- ✗ A new FlowFile will be created if the content of the existing FlowFile is modified or new data is ingested from source.
- ✗ New FlowFile will not be created if the attributes of the existing FlowFile is modified.

# WHAT PROCESSORS ARE AVAILABLE IN NiFi?

The Swiss Army Knife of Data Flow.



## WHY IT'S IMPORTANT?

- X In order to create an effective dataflow, the users must understand the various types of Processors available for them to use.
- X NiFi out of the box have different types of processors for all your data needs.
- X At the time of recording this video, NiFi has 280+ Processors.



## DATA INGESTION PROCESSORS

X GenerateFlowFiles  
X GetFile  
X GetFTP  
X GetSFTP  
X GetJMSQueue

X GetJMSTopic  
X GetHTTP  
X ListenHTTP  
X ListenUDP  
X GetHDFS

X GetKafka  
X QueryDatabaseTable  
X GetMongo  
X GetTwitter  
X ListHDFS / FetchHDFS



## DATA TRANSFORMATION PROCESSORS

- |                    |                       |                     |
|--------------------|-----------------------|---------------------|
| X ConvertRecord    | X ReplaceText         | X EncryptContent    |
| X UpdateRecord     | X CompressContent     | X TransformXml      |
| X ConvertJSONToSQL | X ConvertCharacterSet | X JoltTransformJSON |



## DATA EGRESS / SENDING DATA PROCESSORS

X PutEmail  
X PutFile  
X PutFTP

X PutSFTP  
X PutJMS  
X PutSQL

X PutKafka  
X PutMongo  
X PutHDFS



## ROUTING AND MEDIATION PROCESSORS

ControlRate

DetectDuplicate

DistributeLoad

RouteOnAttribute

RouteOnContent

ScanAttribute

ScanContent

ValidateXml

ValidateCSV



## DATABASE ACCESS PROCESSORS

ConvertJSONToSQL  
 ExecuteSQL

PutSQL  
 SelectHiveQL

PutHiveQL  
 ListDatabaseTables



## ATTRIBUTE EXTRACTION PROCESSORS

EvaluateJsonPath  
 EvaluateXPath  
 EvaluateXQuery

ExtractText  
 HashAttribute  
 HashContent

IdentifyMimeType  
 UpdateAttribute  
 LogAttribute



## SYSTEM INTERACTION PROCESSORS

- X ExecuteProcess
- X ExecuteStreamCommand



## SPLITTING AND AGGREGATION PROCESSORS

SplitText  
 SplitJson  
 SplitXml

SplitRecord  
 SplitContent  
 UnpackContent

SegmentContent  
 MergeContent  
 QueryRecord



## HTTP AND UDP PROCESSORS

- |              |                      |                   |
|--------------|----------------------|-------------------|
| X GetHTTP    | X PostHTTP           | X Listen UDP      |
| X ListenHTTP | X HandleHttpRequest  | X PutUDP          |
| X InvokeHTTP | X HandleHttpResponse | X ListenUDPRecord |



## AMAZON WEB SERVICES PROCESSORS

- |                 |             |               |
|-----------------|-------------|---------------|
| X FetchS3Object | X GetSQS    | X GetDynamoDB |
| X PutS3Object   | X PutSQS    | X PutDynamoDB |
| X PutSNS        | X DeleteSQS | X PutLambda   |

# PROCESSOR CONFIGURATIONS, CONNECTIONS & RELATIONSHIPS IN NiFi

Configuration over coding !!



## TYPE OF NiFi PROCESSOR CONFIGURATIONS

- X Standard Configurations – Configuration common across all processors
- X Unique Configurations – Configuration specific to a processor



## WHAT IS A RELATIONSHIP?

- X Each Processor has zero or more Relationships defined for it.
- X Relationships are named to indicate the result of processing a FlowFile.
- X After a Processor has finished processing a FlowFile, it will transfer the FlowFile to one of the Relationships.
- X The FlowFile creator needs to handle all the relationship of a processor or terminate unhandled relationship.

# CONNECTION QUEUE & BACK PRESSURE IN NiFi



## WHAT IS BACK PRESSURE?

- X Each Processor has zero or more Relationships defined for it.
- X Relationships are named to indicate the result of processing a FlowFile.
- X After a Processor has finished processing a FlowFile, it will transfer the FlowFile to one of the Relationships.
- X The FlowFile creator needs to handle all the relationship of a processor or terminate unhandled relationship.

# WORKING WITH ATTRIBUTES & CONTENT IN NiFi



## KEY STEPS IN DATA FLOW



# WORKING WITH EXPRESSION LANGUAGE IN NiFi



## KEY TAKEAWAY

- X NiFi is Swiss Army Knife of Data Flow with Various Processors and Components.
- X There is no single right approach.
- X But there can be one efficient approach compared to the others.
- X Always choose the efficient approach.



## KEY FACTORS WHICH DETERMINES THE EFFICIENCY OF A DATA FLOW

- X Number of IO Operations
  
- X Memory Utilization
  - CPU Utilization
  - For RAM Utilization

# MORE ON EXPRESSION LANGUAGE FUNCTIONS IN NiFi

WORKING WITH  
PROCESS GROUP,  
INPUT PORT & OUTPUT PORT  
IN NiFi



## PROBLEM STATEMENT

- X How will you share a data flow with your friend or colleague to help in debugging?
- X How will you move your data flow from one machine to another machine and continue your flow design from the new machine?
- X How will you move your completed data flow from a development environment to a testing environment?

# WORKING WITH TEMPLATES IN NiFi

Import and Export data flow using Templates



## VERSION CONTROL IN NiFi

- X Though templates are good to share data flow it's not ideal for Version Control.
- X “NiFi Registry” is the right tool for Version Control of Data Flow in NiFi.

# WORKING WITH FUNNEL IN NiFi

Combine data from several connections to single connection !!

# WORKING WITH CONTROLLER SERVICES IN NiFi



## WHAT IS CONTROLLER SERVICE?

- X A Controller Service is a shared service which can be used across Processors or other Controller Services.
- X "Controller Service" is not limited to Database configuration.
- X You can do various abstraction of functionality using it.

# WORKING WITH VARIABLE REGISTRY IN NiFi



## WHY TO USE VARIABLE REGISTRY?

- X To hold environmental or system specific properties.
- X Helps to simplify the configuration management and migration of Data Flow across environments.
- X Helps in CI/CD of Data Flow.



## HOW TO DEFINE VARIABLE REGISTRY?

X Variable Registry can be defined using two ways

- Using “**Variables**” window available in NiFi UI.
- Using “**Custom Configuration File**” and referring them in “**nifi.properties**” file.



## ADVANTAGES OF UI BASED VARIABLES

- X Used to add variables without restarting the NiFi server.
- X Helps to override an incorrect variable defined using property file based approach.



## TEMPLATE & VARIABLE REGISTRY

- X Templates created will have all the variables along with its value.
- X The variables pointing to sensitive field like password will be ignored during template creation.

# FLOWFILE PRIORITIZATION IN NiFi



## WHAT IS FLOWFILE PRIORITIZATION?

- X Prioritizers comes handy when you have data coming from multiple sources, and you would like to process some data immediately after it arrives compared to other data.



## TYPES OF FLOWFILE PRIORITIZATION

### X **FirstInFirstOutPrioritizer:**

- Given two FlowFiles, the one that reached the connection first will be processed first.

### X **NewestFlowFileFirstPrioritizer:**

- Given two FlowFiles, the one that is newest in the dataflow will be processed first.

### X **OldestFlowFileFirstPrioritizer:**

- Given two FlowFiles, the one that is oldest in the dataflow will be processed first. 'This is the default scheme that is used if no prioritizers are selected.'



## TYPES OF FLOWFILE PRIORITIZATION (CONT.)

### X PriorityAttributePrioritizer:

- Given two FlowFiles that both have a "priority" attribute, the one that has the highest priority value will be processed first. Note that an UpdateAttribute processor should be used to add the "priority" attribute to the FlowFiles before they reach a connection that has this prioritizer set. Values for the "priority" attribute may be alphanumeric, where "a" is a higher priority than "z", and "1" is a higher priority than "9", for example.

# FlowFile Expiration IN NiFi



## WHAT IS FlowFILE EXPIRATION?

- X FlowFile expiration is a concept by which data that cannot be processed within a particular timeframe can be automatically removed from the flow.
- X This will come in handy when the volume of data is expected to exceed the amount that can be processed by NiFi.



## How FLOWFILE EXPIRATION WORKS?

- X The expiration period is based on the time that the data entered the NiFi instance.
- X We can use the expiration in conjunction with Prioritizers to ensure that the highest priority data is processed first and then anything that cannot be processed within a specified period can be dropped.

# MONITORING NiFi

With great power comes great responsibility !!

# MONITORING NiFi USING REPORTING TASK



## WHAT IS A REPORTING TASK?

- X A reporting task in NiFi runs in the background and provides various statistics of the NiFi instance.

# REMOTE MONITORING NiFi USING REPORTING TASK



## REMOTE MONITORING USING REPORTING TASK?

- X It's practically impossible for someone to sit in front of a computer and monitor for any errors using the UI.
- X NiFi provides various Reporting Tasks to monitor your NiFi instance remotely.
- X We will explore "SiteToSiteBulletinReportingTask" and "SiteToSiteMetricsReportingTask."



## WHAT IS SITE-TO-SITE?

- X Site-to-Site is a protocol used to send data from one NiFi instance to another NiFi instance, in a smooth, effective, and secure way.
- X You can also use the Site-to-Site protocol to transmit data from any application which produces data to a NiFi instance.



# DATA PROVENANCE IN NiFi



## WHAT IS DATA PROVENANCE?

- X NiFi keeps a comprehensive track of the data by recording all the events applied on a flowfile starting from its ingestion point till it's removed from the data flow.
- X As the data is processed through the system, NiFi captures all the details like when the data got transformed, split, routed, aggregated, or dispersed to other endpoints.
- X All these information's are stored and indexed in a separate repository called the Provenance Repository.



## WHY DATA PROVENANCE?

- X To Debug your Data Flow at Production level.
- X To Identify the Origin, Destination & Transformation happened to your data.

# APACHE NiFi ARCHITECTURE

# HIGH LEVEL OVERVIEW OF KEY APACHE NiFi FEATURES

# OVERVIEW OF NiFi REGISTRY

Version Control for your Data Flow !!



## WHAT IS NiFi REGISTRY?

- X NiFi Registry is a complementary project that provides a central location for storing and managing shared resources across one or more NiFi instances.
- X It is a separate sub-project of Apache NiFi.
- X This means we must download it separately and it will follow a different release cycle and version.



## WHY NiFi REGISTRY?

- X Assume you are working in a team, and more than one person is working on a Data Flow, then version control of the flow becomes complex.
- X For a very long-time, people used NiFi templates to enable version control and it's such a pain.
- X NiFi template is not created or optimized for doing version control in the first place.



## WHY NIIFI REGISTRY? (CONT.)

- X We have to manually download the template and commit your changes to any other version control tools like TFS or SVN or GIT.
- X Merge is going to be a nightmare.
- X More complexity is involved when we need to take the latest version of the template from the repository and merge it with your un-committed local version.

RESCUE



IMPLEMENTATION OF A FLOW REGISTRY FOR  
STORING AND MANAGING VERSIONED FLOWS !!

# INSTALLATION OF NiFi REGISTRY

# CONFIGURING NiFi & NiFi Registry TO ENABLE VERSION CONTROL



## WHAT IS A BUCKET?

- X A bucket is nothing but logical segregation of the Data Flow, and one bucket can have more than one Flow associated with it.
- X You can consider a Bucket like a Folder available in the File System where we will keep related files inside it.
- X In our case, the files are nothing but the related flows.

# CONFIGURING NiFi REGISTRY WITH MULTIPLE NiFi INSTANCES



## CURRENT LIMITATION OF NiFi REGISTRY

- X You can't take the latest version of the flow without reverting your local changes.
- X To put it in simpler words, you can't merge your local changes with a newer version available in the registry.
- X You can only revert your local changes and take the latest version and redo your changes again and commit the same.

# CONFIGURING NiFi REGISTRY TO ENABLE GIT PERSISTENCE

# OVERVIEW ON NiFi CLUSTERING



## WHY NiFi CLUSTER?

- X Sometimes you may find it difficult to use one NiFi instance to process huge amount of data.
- X Instead you can use multiple NiFi servers to process large data sets.
- X This can be done by segregating the large data set into multiple smaller data sets and sending these small data sets to different servers to process them separately.



## COMMON PROBLEMS IN CLUSTERING

- X Each time when you want to change or update the dataflow, you must make those changes on each server and then monitor each server separately.
- X Segregating big data set into multiple smaller datasets and processing it separately using different servers will have its own complexity.



## FEATURES OF NiFi CLUSTER

- X In NiFi we can cluster multiple NiFi servers and each server in the cluster can perform the same set of tasks on the data, but each can operate on a different set of data.
- X In NiFi cluster, if you make change in one node it automatically gets replicated to all the nodes of the cluster.
- X Using a single interface, the Data Flow Manager can monitor the health and status of all the nodes.



## How NiFi CLUSTER WORKS?

- X NiFi follows a Zero-Master Clustering paradigm.
- X In NiFi, we can use Apache ZooKeeper for Cluster Management, and any failover is handled by ZooKeeper.



## WHAT IS ZooKEEPER?

- X ZooKeeper is an open-source server which enables highly reliable distributed coordination.
  
- X Many open source software like SOLR uses ZooKeeper for its Cluster Management.



## How NiFi CLUSTER WORKS? (CONT.)

- ✗ ZooKeeper elects one of the NiFi nodes as the Cluster Coordinator and all nodes in the cluster will send status information to this node.
- ✗ The Cluster Coordinator is responsible to disconnect nodes that do not emit any heartbeat status for some amount of time.
- ✗ When a new node opts to join the cluster, that node must first connect to the currently elected Cluster Coordinator in order to obtain the latest flow.



## How NiFi CLUSTER WORKS? (CONT.)

- X If the Cluster Coordinator decides to allow the node to join the cluster, the current flow is provided to that node, and that node will then able to join the cluster.
  
- X If the version of the flow configuration available in the new node differs from the version available in the Cluster Coordinator, the node will not be able to join the cluster.

# LIMITATION IN NiFi CLUSTERING



## LIMITATION IN NiFi CLUSTERING

- X In a NiFi cluster, data is “Distributed” for processing, but it’s not “Replicated.”
- X In a NiFi Cluster when we split large data set into multiple smaller data sets, each node will be given the smaller data set it needs to process, and the same node keeps this data in its disk storage.
- X A copy of this data is not maintained or replicated in any other nodes.



## SIDE-EFFECT OF NOT HAVING DATA REPLICATION

- X In a NiFi cluster, if one node goes down, we have a provision to offload that node and the data in that node will be distributed to other active nodes slowly, provided that node is still connected to the network.
  
- X But, if the node completely goes out of the network for some reason, says someone pulled the network cable of that machine. The data inside that node will not be processed or distributed to other active nodes until it comes back to the network.

# NiFi CLUSTER CONFIGURATION USING EMBEDDED ZOOKEEPER



## EMBEDDED ZOOKEEPER

- X Each NiFi instance is bundled with a ZooKeeper instance inside it.
- X We will be using three NiFi instances and the three embedded Zookeeper instances available within it.
- X The only key prerequisites here is you should have more than one machine available to follow along and connectivity between these machines must be enabled.



## DIGITAL OCEAN DROPLETS

- X I will not be using my local machine like I used to do previously. Instead I will be using three “Digital Ocean Droplets”.
- X If you are coming across Digital Ocean Droplets for the first time, just remember it’s something similar to “AWS EC2”.
- X To put it simpler, I have loaned commodity machines according to my required specification and I will be billed based on the number of hours I have used these machines.

# NiFi CLUSTER CONFIGURATION USING EXTERNAL ZOOKEEPER



## PITFALL USING EMBEDDED ZOOKEEPER

- X Using Embedded Zookeeper is perfectly fine and it can be used in production.
- X Drawback in this approach is, both NiFi and Zookeeper are running in the same server and if the server goes down, we will lose the NiFi and Zookeeper instance at the same time.
- X Using an external Zookeeper, we can run the NiFi and the Zookeeper instances in different machines.

# OVERVIEW ON NiFi CUSTOM PROCESSOR

Hope for the Best !! Plan for the Worst !!



## NiFi PROCESSOR RECAP

- X We have tons of Data Source and Data Sink Processors.
- X At the time of recording this video, NiFi has 280+ Processors.
- X Each Processor in NiFi is unique in its own way.
- X NiFi have processors for all your data needs.



## WHY CUSTOM PROCESSORS?

- X There will be situations where you will not be able to use any of these in-built processors which comes bundled with NiFi to cater your requirements.
  
- X This is where the Custom Processors of NiFi comes in handy.



## THE MILLION DOLLAR QUESTION

- X NiFi provides maven archetypes to create our own processor or controller service which is compatible and easy to include as part of our data flow.
  
- X Whenever someone tells, you can extend their tool. You will get tons of questions regarding how easy it's going to be and how can we can migrate this custom code when a new version of the tool arrives.



## WHY ITS SIMPLE?

- X A processor in NiFi takes an input flow file and do some processing on top of it and produces an output flow file. It can also have some properties using which you can configure the way the processor processes the input data.
  
- X The FlowFile abstraction of NiFi makes it so simple and easy to create custom processors.

RESCUE



# MAVEN ARCHETYPE

AUTO GENERATES EVERYTHING WE NEED !!



## ADVANTAGE OF USING MAVEN ARCHETYPE

- X The NiFi custom processor maven archetype makes your life easy by auto-generating the code required for us to get started.
  
- X The auto-generated code has everything we need, and all we have to do is to write your custom Java code inside it to process the input flow file to produce the output flow file.

OUR FIRST  
CUSTOM PROCESSOR



## WHAT IS A NAR?

- X NAR is NiFi Archive file.
- X It's similar to a WAR or JAR file.
- X We can deploy the NAR files in NiFi, similar to deployment of war files in tomcat or any other application servers.