# Introduction to Apache NiFi

# What is Apache NiFi?

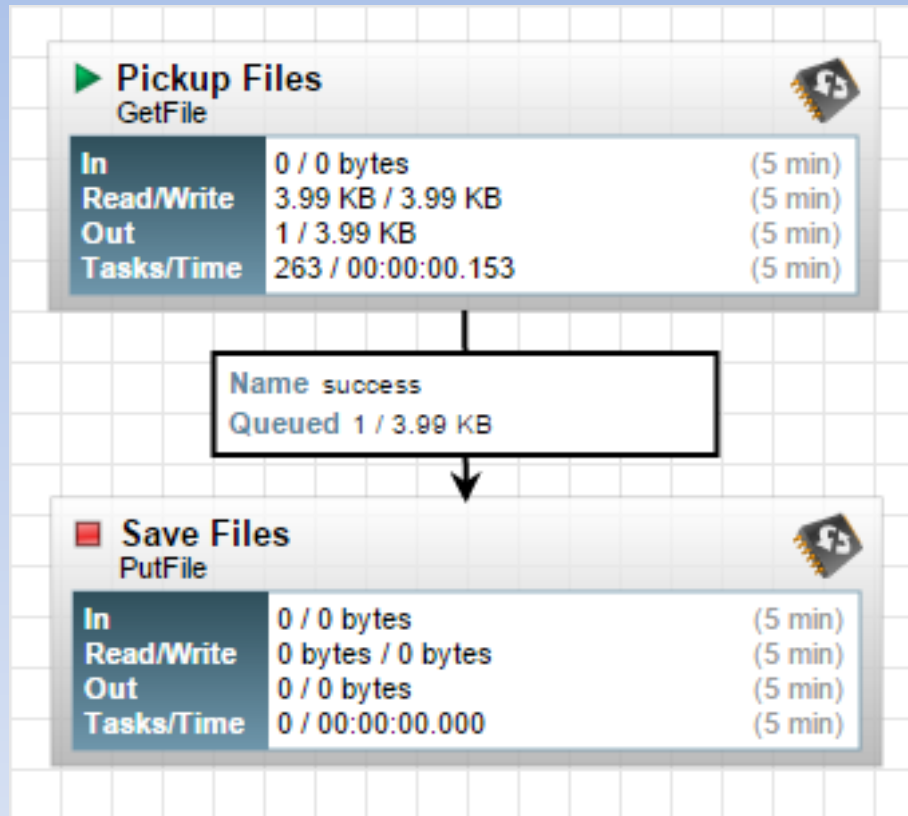An Open Source Data Distribution and Processing System

# What does that mean?

Apache NiFi provides a way to move data from one place to another, making routing decisions and transformations as necessary along the way
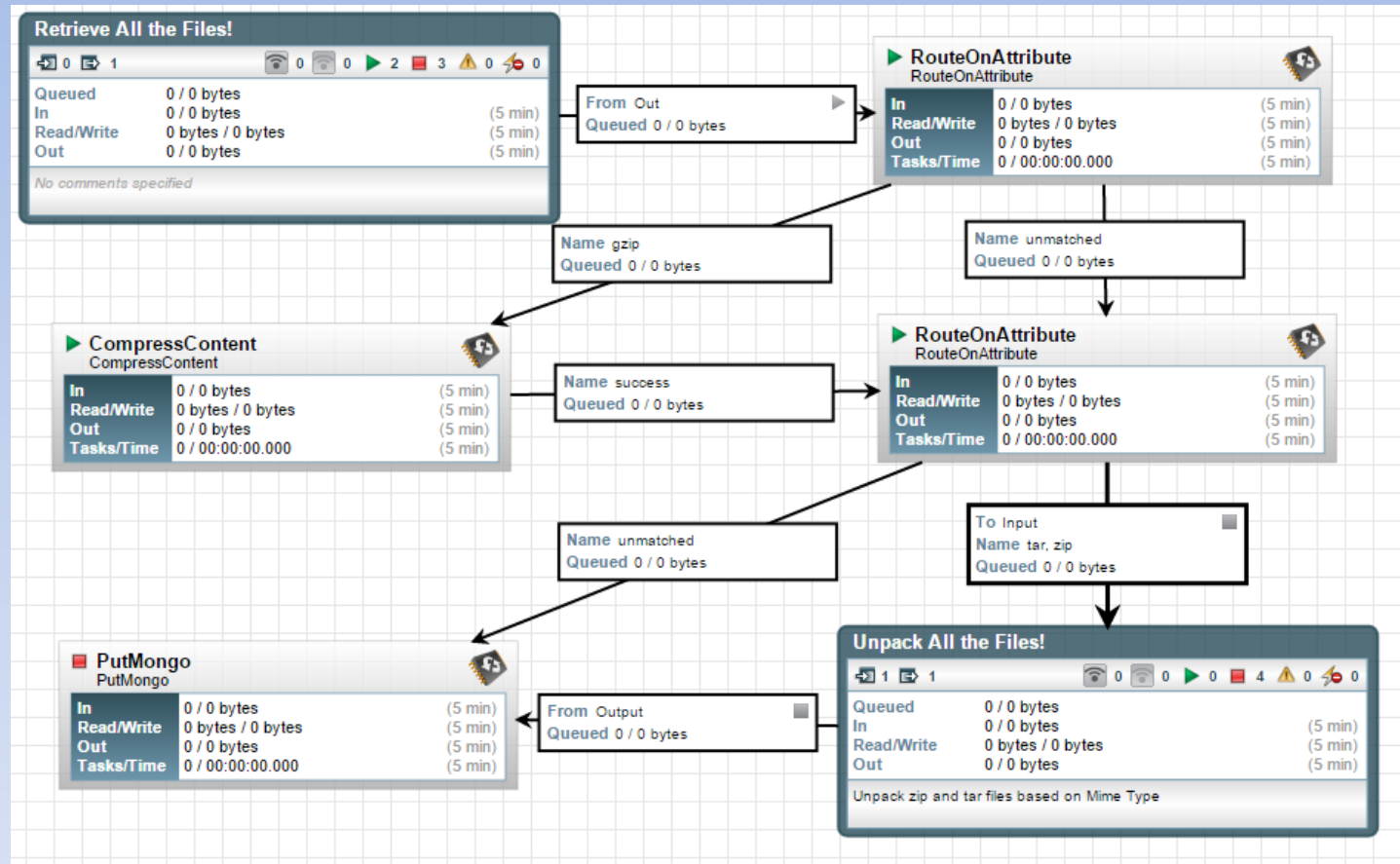
**REQUITEST**
*Define Your Success!*

# Why Use Apache NiFi?

- Easy to use

- Powerful

- Reliable

- Secure

- Scalable

**REQUITEST**
*Define Your Success!*

# Can handle Basic Flows...

# … To More Advanced Flows

# Features

- Web-based Interface
  - Flow construction, control, and monitoring all from a single easy to use interface

- Data Provenance
  - Track data throughout the entire flow
  - Information about FlowFiles as they traverse the flow are automatically indexed
  - Critical for supporting troubleshooting and flow optimization.

REQUITEST
Define Your Success!

# Features

- Data Recovery
  - Ages off content as space is needed
  - Allows for fine grained download, recovery, and replay of individual files.

- Secure
  - Provides content encryption, communication over secure protocols (SSL, SSH, HTTPS), etc.
  - Provides a pluggable role-based authentication/authorization mechanism for both data transfer and user management

**REQUITEST**
*Define Your Success!*

# Features

- Highly configurable
  - Fine grained Quality of Service control
  - Dataflow modifiable at runtime
  - Loss tolerant vs guaranteed delivery
  - Low latency vs high throughput
  - Back pressure

**REQUITEST**
*Define Your Success!*

# Features

- Extensible
  - Build your own processors, controller services, and more
  - Enables rapid development and effective testing
  - Allows for development of simple single function components that can be reused and combined to make more complex flows
  - Provides classloader isolation for easier management of dependencies

# Definitions

## FlowFile

- The data that moves through the flow
- Can be cloned, merged, split, modified, transferred, and deleted
- Consists of:
  - Map of key/value pair attribute strings
  - Content of zero or more bytes

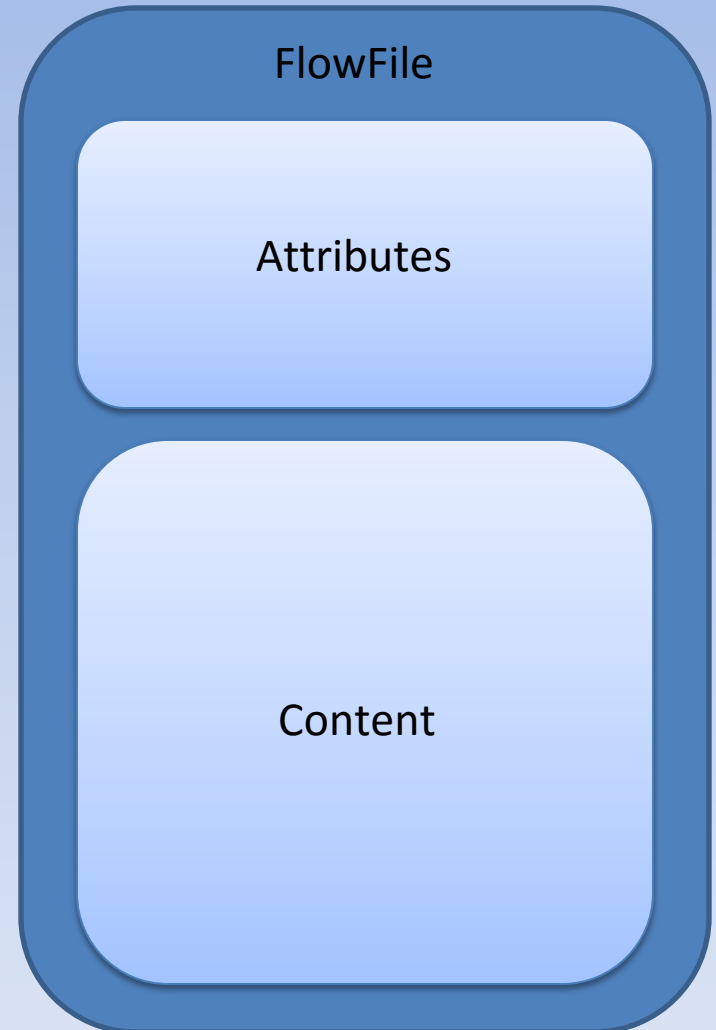**REQUITEST**
*Define Your Success!*

# FlowFile Breakdown

## Attributes

- Map of Key/Value pairs
- Heavily used to make routing decisions
- Values accessed using NiFi's Expression Language

## Content

- The actual data that is being routed through the dataflow
- May be manipulated multiple times throughout the course of a dataflow

FlowFile

Attributes

Content

# Common Attributes

**filename** – A filename that can used when storing data locally or on a remote system

**path** – the directory that can be used when storing data
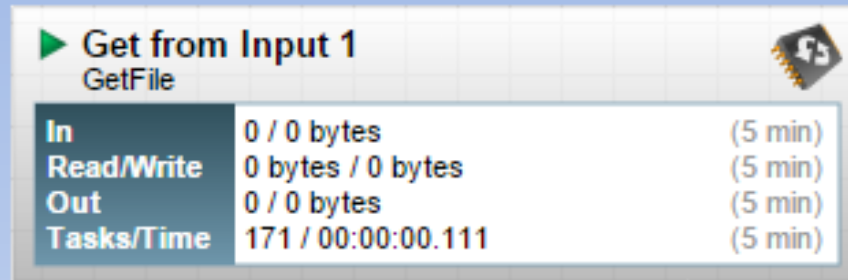
**uuid** – A Universally Unique Identifier that distinguishes FlowFile from other FlowFiles

**entryDate** – the date and time at which the FlowFile entered the system

**lineageStartDate** – The date and time at which the oldest ancestor of the FlowFile entered the system.

**fileSize** – Represents the number of bytes taken up by the FlowFile's Content

**REQUITEST**
*Define Your Success!*

# Definitions



## FlowFile Processor

- Single step in the flow
- Performs the work on the FlowFile
  - Routing
  - Data Transformation
  - Mediation between systems
- Has access to FlowFile content and attributes
- Can operate on zero or more FlowFiles in a single unit of work

# FlowFile Processor Examples

## Ingestion

- **GetFile** – Pull content from the local disk and delete the original file
- **GetSFTP** – Pull content from a remote system then delete the original file

## Routing

- **RouteOnAttribute** – Route FlowFiles based on the values of specific FlowFile attributes

## Data Transformation

- **CompressContent** – Compress or decompress content
- **ReplaceText** – Use Regular Expressions to modify textual content

**REQUITEST**
*Define Your Success!*

# FlowFile Processor Examples

## Data Egress

- **PutFile** – Writes the FlowFile contents to a directory on the local disk
- **PutSFTP** – Copies the contents of the FlowFile to a remote server

## Attribute Extraction

- **UpdateAttribute** – Adds or updates attributes using statically defined values or dynamically derived values using NiFi's Expression Language
- **ExtractText** – Creates attributes based on User defined Regular Expressions

## Splitting and Aggregation

- **UnpackContent** – Unpacks archive formats such as TAR and ZIP and sends each file within the archive as a separate FlowFile through the dataflow

**REQUITEST**
*Define Your Success!*

# Definitions

## Connection

- Provides linkage between processors
- Queues that allow rate control
- Dynamically prioritizable
- Enable back pressure via configurable upper bounds

**REQUITEST**
*Define Your Success!*

# Definitions

## Controller Service

- A single service that can be shared between multiple FlowFile processors

- Performs a specific task or maintains a common set of information

- Example: **StandardSSLContextService** provides a single configuration for a keystore and/or truststore that can be used throughout a dataflow

**REQUITEST**
*Define Your Success!*

# Definitions

## Flow Controller

- Scheduler
- Maintains processor and connection configuration
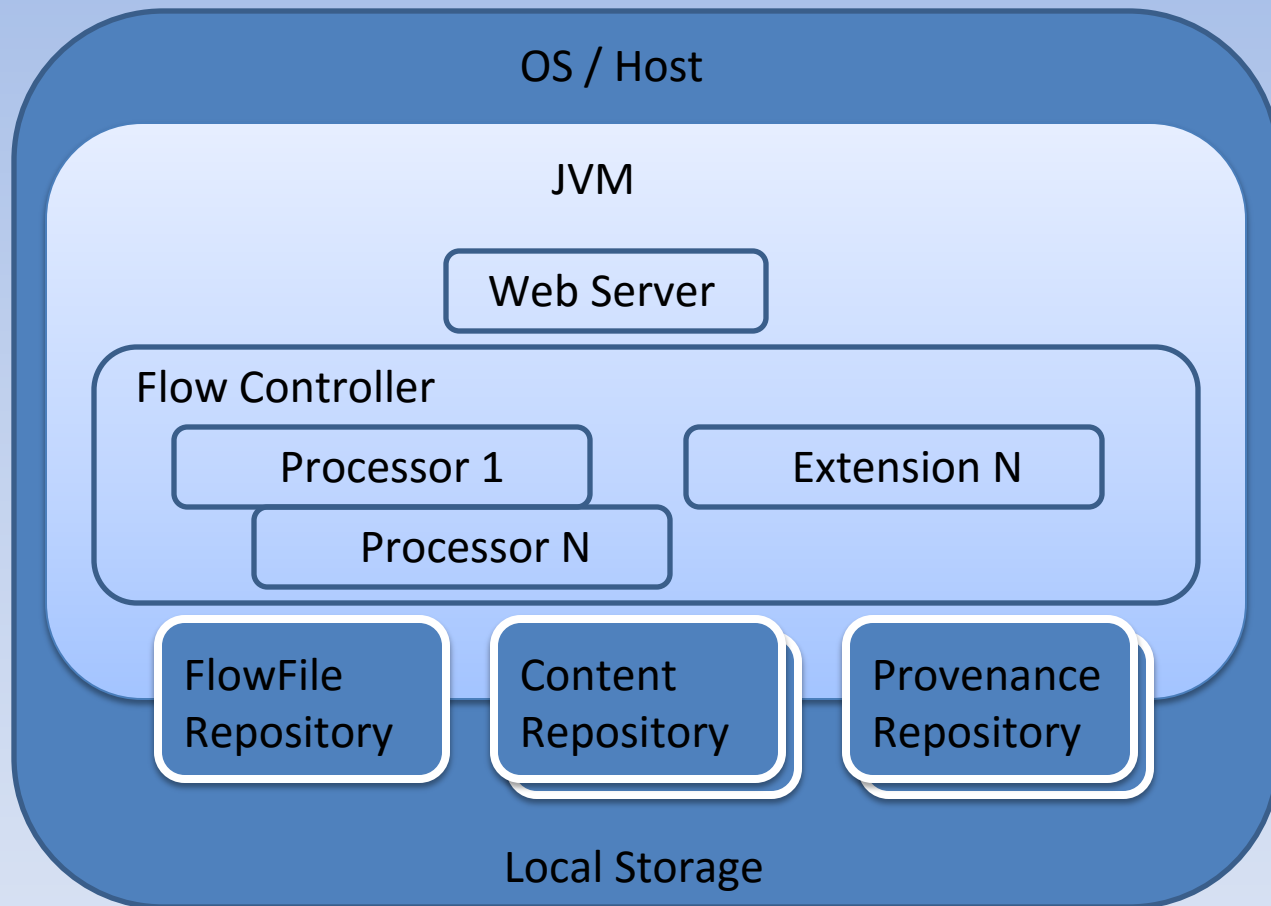- Handles scheduling of threads which processes use

## Process Group

- Set of processors and their connections
- Receives data via input port(s) and sends data via output port(s)

**REQUITEST**
*Define Your Success!*

# NiFi Architecture

- NiFi is a Java based system that executes within a JVM.
- Primary components are:
  - Web Server
    - Hosts NiFi's HTTP-based control API
  - Flow Controller
    - Provides and schedules threads for execution
  - Extensions
    - FlowFile Processors, Controller Services, etc.
  - Repositories
    - FlowFile
    - Content
    - Provenance

# NiFi Architecture

# Repositories

## FlowFile Repository

- Holds information pertaining to the FlowFile and its attributes

## Content Repository

- Holds all of the FlowFile content

## Provenance Repository

- Holds all information pertaining to the life of the FlowFile as it traverses the dataflow

# How Do I Get It?

http://nifi.apache.org/download.html

Two versions available:
- A "tarball" tailored for Linux
- A zip file tailored for Windows

Download the appropriate version and extract to the location from which you want to run NiFi.

Mac OSX Users may also use the tarball or can install via Homebrew by running:

```
brew install nifi
```

**REQUITEST**
*Define Your Success!*

# Running NiFi
# (Linux/Mac OSX)

Using a Terminal window, navigate to the directory where NiFi was installed.

To run NiFi in the foreground, run
```
bin/nifi.sh run
```
Use Ctrl-C to stop the application.

To run NiFi in the background, run
```
bin/nifi.sh start
```
To stop the application, use
```
bin/nifi.sh stop
```

**REQUITEST**
*Define Your Success!*

# Running NiFi
# (Windows)

Nagivate to the folder where NiFi was installed. Double-click the *bin/run-nifi.bat* file.

To stop the application, select the window that was launched and press Ctrl-C.

**REQUITEST**
*Define Your Success!*

# Congratulations!

To start using NiFi, open a web browser and navigate to [http://localhost:8080/nifi](http://localhost:8080/nifi)

Port 8080 is the default port and can be changed by editing the `nifi.properties` file in the NiFi `conf` directory.

# Further Resources

- RequiTest Website:
  http://requitest.com/

- Apache NiFi Website:
  http://nifi.apache.org/

- Apache NiFi Users Mailing List:
  http://mail-archives.apache.org/mod_mbox/nifi-users/

- Apache NiFi Developers Mailing List:
  http://mail-archives.apache.org/mod_mbox/nifi-dev/

**REQUITEST**
*Define Your Success!*