# Introduction to Big Data

## ECAP456

Dr. Rajni Bhalla

Associate Professor

# Learning Outcomes

After this lecture, you will be able to

- Learn about data ingestion

- Learn about uploading data

# Introduction

Data ingestion in Splunk happens through the Add Data feature which is part of the search and reporting app. After logging in, the Splunk interface home screen shows the Add Data icon

# Introduction

## About uploading data

When you add data to your Splunk deployment, the data is processed and transformed into a series of individual events that you can view, search, and analyze

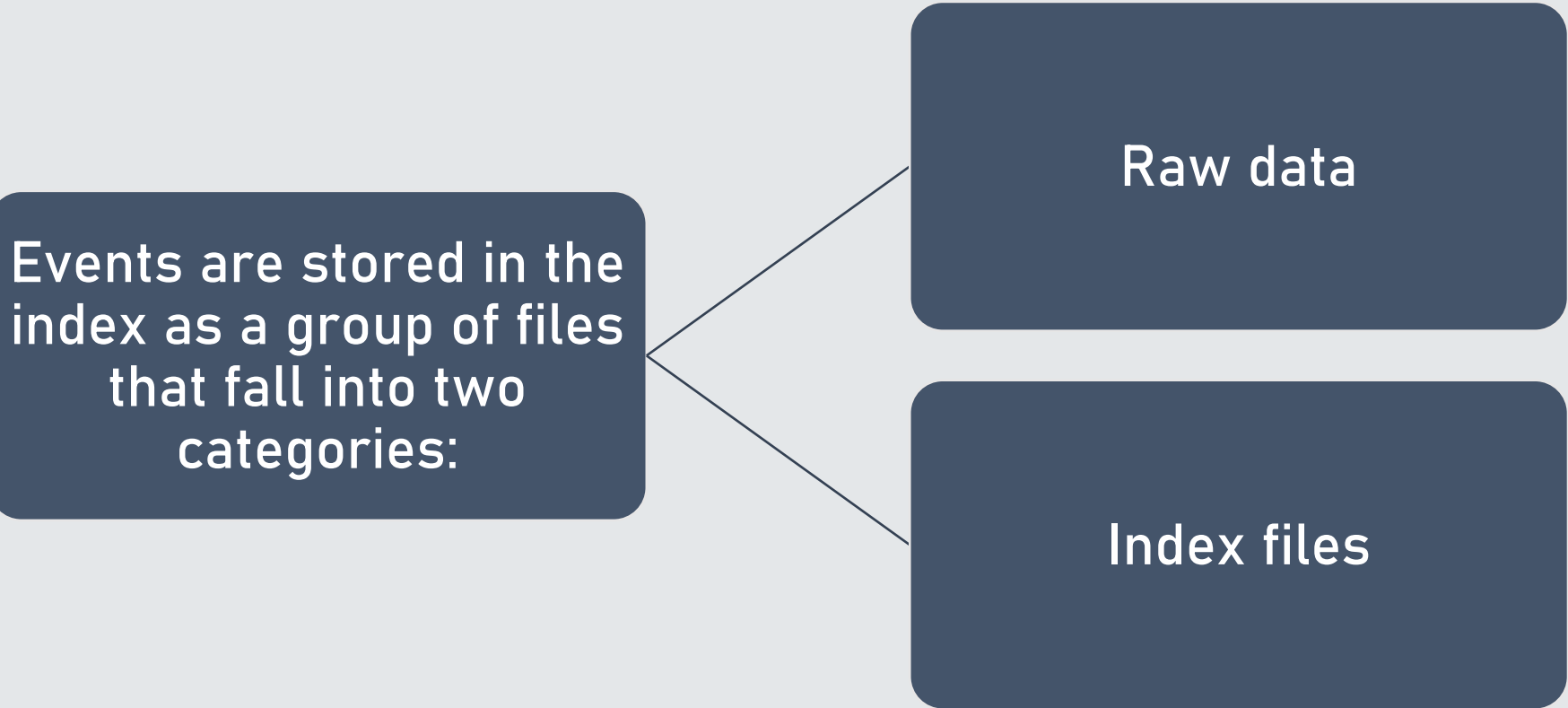# Introduction

## What kind of data?

| Data source |
| --- |
| Files and directories |
| Network events |
| IT Operations |
| Cloud services |
| Database services |
| Security services |
| Virtualization services |
| Application servers |
| Windows sources |
| Other sources |

# Where is the data stored?

- The process of transforming the data is called indexing.

- During indexing, the incoming data is processed to enable fast searching and analysis. The processed results are stored in the index as events.

- The index is a flat file repository for the data. The index resides on the computer where you access your Splunk deployment.

# Where is the data stored?

Events are stored in the index as a group of files that fall into two categories:

Raw data

Index files

# Where is the data stored?

Buckets

Main

Splunk Instance

# Use the Add Data wizard

- If there is a Welcome window displayed, close that window.

- Click Settings > Add Data.

# Use the Add Data wizard

At the bottom of the window, click Upload. There are other options for adding data, but for this tutorial you will upload the data files.

# Use the Add Data wizard

## Under Select Source, click Select File..

# Use the Add Data wizard

- In your download directory, select zip file and click Open.

- Click Next to continue to Input Settings.

- Under Input Settings, you can override the default settings for Host, Source type, and Index. Because this tutorial uses a ZIP file, you are going to modify the Host setting to assign the host values by using a portion of the path name for the files included in the ZIP file. The setting that you specify depends whether you are using Splunk Cloud or Splunk Enterprise, and on the operating system that you are using.

# Use the Add Data wizard

# Use the Add Data wizard

a. Select Regular expression on path.

b. Type \\(.*)\/ for the regex to extract the host

values from the path.

# Use the Add Data wizard

Click Review. The following screen appears where you can review your input settings.

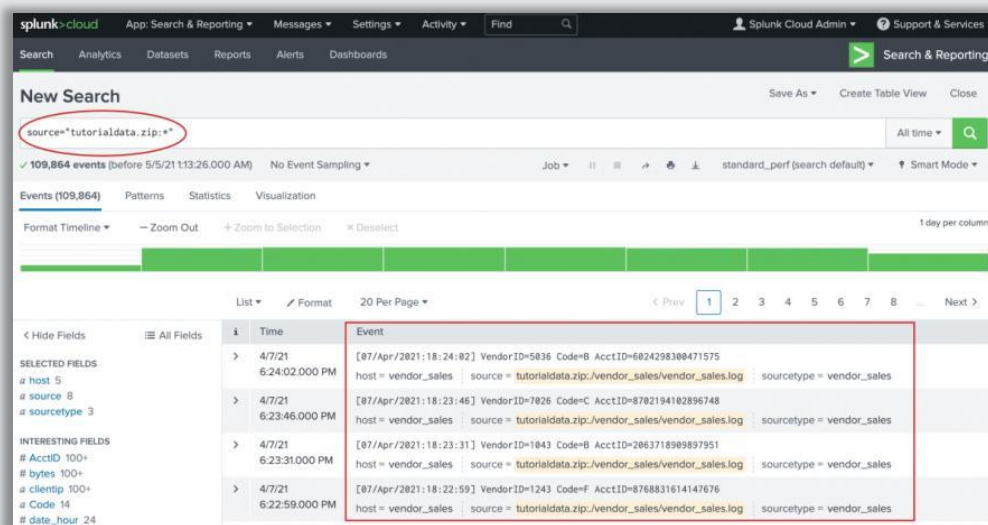# Use the Add Data wizard

Click Submit to add the data.

# Splunk cloud

- To see the data in the Search app, click Start Searching.

- You might see a screen asking if you want to take a tour. You can take the tour or click Skip. The Search app opens and a search is automatically run on the tutorial data source.

# Splunk cloud

- To see the data in the Search app, click Start Searching.

- You might see a screen asking if you want to take a tour. You can take the tour or click Skip.

- The Search app opens and a search is automatically run on the tutorial data source

# Splunk cloud

Success! The results confirm that the data in zip file was indexed and that events were created.

Click the Splunk logo to return to Splunk Home.

That's all for now...