# Project:
# Introducing Deep Learning Pipelines for Apache Spark

## Prachi Sethi 19963

# Table of Contents

- Introduction
- Overview and Philosophy
- Steps
- Cluster Setup
- Compatibility
- Tools and Capabilities
- Working with Images in Spark DataFrames
- Transfer Learning
- Applying Models at Scale
- Conclusion
- References

# Introduction

Deep Learning Pipelines: A new library by Databricks.

Purpose: High-level APIs for scalable deep learning model application and transfer learning.

Integration: Combines popular deep learning libraries with MLlib Pipelines and Spark SQL.

# Overview and Philosophy

Overview: For detailed usage examples, refer to the Deep Learning Pipelines README.

Philosophy: Check out the Databricks blog post for the philosophy behind the library.

# Steps

Step 1-
https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c9
3eaaa8714f173bcfc/5669198905533692/3647723071348946/3983381308530
741/latest.html

Step 2- Create a data bricks account

Step 3- Import the notebook and install packages

Step 4 - Create the cluster

Step 5- Run the cells one by one

```
%pip install tensorflow==2.5.0
%pip install sparkdl
%pip install kafka
%pip install keras
%pip install optree
%pip install tensorframes
%pip install numpy==1.22.0 scipy tensorflow
%pip install --upgrade kafka-python
%pip install tensorflowonspark
%pip install jieba
```

```
ⓘ  › ImportError: cannot import name 'dtensor' from 'tensorflow.compat.v2.experimental' (/local_disk0/.ephemeral_nfs/env
s/pythonEnv-3405dc6b-54b6-4d55-9f9a-e80132b8230d/lib/python3.9/site-packages/tensorflow/_api/v2/compat/v2/experimental/__i
nit__.py)
```

# Cluster Setup

Availability: Deep Learning Pipelines is available as a Spark Package.

Setup Steps:

Create a new library with Source option "Maven Coordinate".

Search Spark Packages and Maven Central for "spark-deep-learning".

Attach the library to a cluster.

# Compatibility

Version: Works with spark-deep-learning release 0.1.0-spark2.1-s_2.11.

Future Releases: Check the project's GitHub page for the latest examples and docs.

Libraries: Also create and attach these libraries via PyPI: tensorflow, keras, h5py.

Spark Version: Compatible with Spark versions 2.0 or higher.

Instance Types: Works with any instance type (CPU or GPU).

# Tools and Capabilities

Image Processing: Tools for working with images using deep learning.

Categories:

Working with Images: Natively in Spark DataFrames.

Transfer Learning: Leverage deep learning quickly.

Model Application: Apply deep learning models at scale.

SQL Functions: Deploy models as SQL functions (coming soon).

Hyper-parameter Tuning: Distributed tuning via Spark MLlib Pipelines (coming soon).

# Working with Images in Spark DataFrames

Capability: Load, transform, and analyze images directly within Spark DataFrames.



```sh
%sh
curl -O http://download.tensorflow.org/example_images/flower_photos.tgz
tar xzf flower_photos.tgz
```

| % Total | | % Received | % Xferd | Average Speed Dload | Upload | Time Total | Time Spent | Time Left | Current Speed |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | --:--:-- --:--:-- --:--:-- | 0 |
| 42 | 218M | 42 92.0M | 0 | 0 | 102M | 0 | 0:00:02 | --:--:-- 0:00:02 | 102M |
| 100 | 218M | 100 218M | 0 | 0 | 126M | 0 | 0:00:01 | 0:00:01 --:--:-- | 126M |

# Transfer Learning

Quick Leverage: Utilize pre-trained models for various tasks with minimal effort.

Example: Apply a pre-trained model to classify images.

# Applying Models at Scale

Scalability: Apply your own or popular models to image data at scale.

Use Cases: Image classification, feature extraction, etc.

# Conclusion

- Deep Learning Pipelines: Facilitates scalable deep learning with Spark.
- Future Developments: Stay updated with the latest releases and features on GitHub.

# References

https://github.com/Prachi1615/MachineLearning

https://community.cloud.databricks.com/?o=4454248012642049#notebook/9142 29071432723/command/914229071432726