**SanKir Technologies**

Accelerating Data Engineering Solution on Cloud

## Expertise in Data Engineering on Cloud

- Data Engineering pipeline – Architecture, Orchestration, Optimization and Monitoring
- End-to-End Automation
  - AWS CDK Toolkit
  - Terraform
- Databricks, Apache Airflow
- Big Data - Apache Spark, Hadoop, HDFS
- DWH - Google BigQuery, Snowflake and DBT

## Consultancy

## Engagement Partner for Business Opportunities

consultancy.sankir.com

aws        Google Cloud        Azure

## Solution Architecture

- Architecting the Data Pipeline on Cloud
  - Multiple Solution options for a data problem
- End-to-End DE Automation using AWS CDK Toolkit and Terraform
- Data Engineering Solution on Cloud
  - Optimizing the Data processing
  - Best Practices in data & network security
  - Spark Cluster - Sizing and Optimization for better performance of workloads
- Leverage SanKir AWS/Azure/GCP Cloud Infrastructure to quickly test the solution
- SanKir can Work with Client's CxO, Technical Managers or Engineering team to solve the organizational Data problems
- SanKir aims to align with Client's Core Business Objectives

## Leverage SanKir services in

- Cloud Services and Solution
- Data Engineering Tools and Solution
- CI/CD – DevOps
- Container Orchestration using Docker and Kubernetes
- Streaming using Spark streaming and Kafka
- API development using Golang

## PoC on Data Engineering pipeline powered by SanKir Framework and Assets

# Consultancy in following Cloud Services and Solution

## AWS

- Amazon S3
- Spark Cluster – AWS EMR
  - Cluster Sizing and Monitoring

- Cloud DWH - RedShift/Athena

- AWS CDK Toolkit

- AWS SDK using Boto3

- Orchestration
  - AWS Cloud Formation
  - AWS Lambda
  - CloudWatch

- IAM and VPC
- Secrets Manager, KMS
- RDS - PostgreSQL, MySQL

- Apache Airflow/MWAA
- Databricks on AWS
- Snowflake on AWS

- CodeBuild, CodeDeploy
- AWS Cost Management

## Azure

- ADLS
- Spark Cluster - Azure Databricks

- Cloud DWH
  - Azure Synapse Analytics
  - SQL Databases

- ABFS(S)/WASB(S) - API to access ADLS

- IAM
- Access Keys
- Shared Access Signature
- Azure Key Vaults, Keys & Secrets
- Service Principal

- Azure Data Factory

- Snowflake on Azure

- Azure Active Directory (AAD)

## GCP

- Google Cloud Storage
- Spark Cluster - Google Dataproc

- Cloud DWH - Google BigQuery

- Google Cloud SDK
  - Google Cloud CLI
  - Cloud Shell

- IAM
- Cloud Functions
- Secret Manager

- Cloud Run
- Cloud Build

- Terraform on GCP

# Consultancy in following Data Engineering Tools and Solution

## Databricks

- ELT Solution provided on Azure, Azure and GCP
- Databricks Spark Cluster
  - Cluster Types - All Purpose and Job Clusters
  - Cluster Configuration
  - Cluster Mode - Standard, High Concurrency, Single Node
  - Cluster Pools
- Databricks Runtime
- Auto Scaling & Auto Termination
- Notebook Workflow Utilities
- DBFS & Databricks mounts
- Databricks secrets
- File System Utilities

## Snowflake

- ELT Solution provided on AWS, Azure and GCP
- Loading Data from S3/ADLS/GCS
- DML for bulk data loading/unloading using COPY command
- Snow pipe - Load data fast, analyze even faster
- Snowflake Connector for Python & Spark
- Accelerating BI Queries with Caching
- Eliminating Concurrency Issues with Snowflake Virtual Warehouses
- Data Protection with Time Travel
- Zero copy cloning
- Standard and Extended SQL
- Advanced DML such as multi-table INSERT, MERGE, and multi-merge
- Statistical and Analytical aggregate functions
- Windowing functions

## Apache Airflow

- Data Pipeline Solution on AWS, Azure and GCP
- Dynamic DAG Authoring using task groups
- DAG Dependencies
- Airflow Variables
- Scheduling and Triggers
- Templating
- Task Flow API
- XCOMs

## Data Build Tool - DBT

- Types of Materializations, Incremental Models & Ephermal Models
- Hooks, Snapshots & Macros
- Integrate DBT with Snowflake, Databricks, BigQuery, RedShift & Apache Spark

# SanKir as Engagement Partner for Business Opportunities

- Architecting the Solution for Client's Potential Business opportunities
  - Recommendation on right choice of Cloud vendor, Data Pipeline Tools and services and Cloud Data Warehouse based on Client requirements

- Execute PoC to win a Project deal

- Customized technology upskill program for Corporate professionals
  - Enable Client to showcase their team's skills for resourcing opportunities
  - Hands-on sessions with PoC project

- Joint Go-To-Market Strategy
  - SanKir can partner with Client from initial bid stage to Architecting the Solution to Project Execution

# WHY SANKIR

- Differentiator in Pricing
  - Competitive pricing compared to Cloud vendor's professional services
  - Aim to deliver within Client's Budget

- Industry Experience in Multitude of Data Projects

Let SanKir own your Data Problems

# Proof of Work - PoC

- **PoCs done for Clients**
  - Retail data transformation using Spark, S3, AWS EMR, Athena
    - End-to-End DE Automation using AWS CDK Toolkit
      - EC2 creation, Airflow installation, EMR creation and Spark job submission
    - Data Profiling – Column profile detail & Data quality metrics
  - Orchestration using Airflow for Ed-Fi operational store used in K-12 Education
    - Dynamic DAGs
    - Schedule and monitoring of task-groups and tasks
  - Data Engineering using built-in and scalable Azure Databricks platform
    - Infra creation ( Spark Cluster ) – compute sizing
    - Data Transformation, Best practice using Key vaults
  - Loading retail data from AWS S3 file storage into Snowflake tables
  - Designing Data pipeline for ETH blockchain using Airflow and AWS EMR

*Cost of Cloud Infrastructure of PoC for the Organization will be very minimal*

# PoC – E2E Automated DE Pipeline using Airflow and AWS EMR for Retail data set

| Data Size | 1 GB |
|---|---|
| AWS Services | S3, Athena, EMR |

| AWS EMR – Spark Cluster | |
|---|---|
| Nodes | 1 Master and 2 Cores (3-Node Cluster) |
| Hardware | Instance type - m5.4xlarge 16 vCore, 64 GiB memory, EBS only storage EBS Storage:256 GiB |
| Release label | emr-5.35.0 |
| Hadoop distribution | Hadoop 2.10.1 |
| Applications | Spark 2.4.8 ; Ganglia 3.7.2 ; Airflow(MWAA) 2.0.2 |

- Spark Cluster sizing based on Hardware Configuration

- Spark Job monitoring using Spark UI/Ganglia

- Cost of Cloud Infrastructure of PoC for the Organization will be very minimal

- PoC setup can be customized to customer needs within a week

Full details of PoC will be shared upon request

# E2E DE Pipeline for Retail data set - Screenshots

# PoC - Data Engineering using built-in and scalable Azure Databricks platform

**Kiran Hiremath**, Director

IT Professional with 27+ years of Experience.

- Expert in Data Engineering, Cloud Services and Distributed Computing using Apache Spark
  - Data Engineering Pipeline – Architecture, Orchestration, Optimization and Monitoring
  - Spark Cluster sizing – AWS EMR, Google Cloud Dataproc
  - Cloud Storage and Datawarehouse – AWS S3, Athena, ADLS, GCS, BigQuery
  - Big Data Technologies: Hadoop, HDFS, Spark, Scala, pySpark
  - Databricks on Azure & AWS
- Experience in Pre-Sales, CoE, Alliance, Software Development & Management
- Worked for TCS, Wipro and has interfaced with MNCs like Nortel, Motorola and Alcatel-Lucent
  - Feature design in Network-Switches
  - Data-driven Contact Centre solution in Healthcare, Banking and Telcos
  - Data Migration, ETL to Salesforce.com application from Siebel CRM
  - IT Transformation Programs

https://github.com/kiranhm1972          https://www.linkedin.com/in/hiremath-kiran/

Experience in Consultancy, Conceptualization, Asset & Solution Development on Data Platforms

**Sanjay Bheemasenrao**, Director

27+ years of experience in building products and services

- Demonstrated leadership in building focused teams to achieve excellence

- Functioned in various capacities and has held many important positions in India and US. Has worked for reputed companies like Oracle, Tata Elxsi, GE and Prentice Hall.

- Expert in Data Engineering covering all the Data Management aspects

  - Data Engineering Pipeline – Architecture, Orchestration, Optimization and Monitoring

  - Big Data Technologies: Hadoop, HDFS, Spark, Scala, Python, Hive

  - Architecting Distributed Computing solutions using Apache Spark, Hive, BigQuery and AWS Athena, RedShift, DBT, AWS Lambda, Cloud Watch, Cloud Build, Apache Airflow.

- Worked in Oracle India for 16+ years

  - Oracle E-Business Release Management expert - 11i to R12.2

  - Developed SaaS solutions for Manufacturing industry on Oracle Cloud

GitHub https://github.com/sbheemas

LinkedIn https://www.linkedin.com/in/sanjaybheemasenarao/

Technology focus areas – AWS and GCP, Data Lake and Cloud Data warehouses

**Self-paced Online Courses**

For Experienced IT professionals

| Data Engineering on Cloud | Snowflake with pro-Spark | Databricks on Cloud |
|---|---|---|
| pro-Spark-aws | pro-Spark on aws | Learn Databricks on AWS |
| pro-Spark-gcp | pro-Spark on gcp | Learn Databricks on GCP |

- Quiz
- Assignment
- Bonus Lessons

Exclusive batches for corporate at competitive price

# About SanKir

**10,000+**

Lines Of Code

**1000+**

Coached Professionals

**20+**

Meetups

**55+**

IT Industry Experience

## SanKir History

**2022 – Consultancy**
Solution Architecture, PoC, Cloud Services & Solution, Data Engineering Tools & Solution

**2020 – Self-paced Online Courses**
Data Engineering on Cloud – AWS, Azure and GCP.

**2018 - Founded**
Classroom Courses on Big Data, Spark, Java with Capstone Projects.

SanKir Technologies

Thank you!

consultancy.sankir.com

info@sankir.com