



Data Lake Agile Factory (DLAF) – PHD Data Collection Service **Technical Design Document**

Owner: UDA Program – Data Lake Team Maintained by: Dexter Obiso Effective Date: March 16, 2020

Change Summary

Version	Date	Change/ Updated by	Summary
1.0	03/16/2020	Dexter Obiso	Adding Azure Service Design/Implementation

Pending Issues

Date	Summary

Referenced Document(s)/Position Paper(s)

Version	Name	Link to NGS Document Repository (e.g. SharePoint)

Review/Approve Summary

Date	Role (Reviewer/ Approver)	By	Comments

Contents

<i>Service Background and Overview</i>	4
<i>Service Architecture Diagram</i>	5
Logical Architecture	5
Physical Architecture (DEV Environment)	7
Physical Architecture (PROD Environment)	8
<i>Data Model Diagram</i>	9
<i>PHD Adhoc Collection Service</i>	10
Installation and Setup	10
Create Installer	10
Configuration Settings	10
Install Instructions	11
System Design Considerations	11
Software Design Pattern	11
Multithreading	12
Error Handling & Logging	12
Service Layer	13
Program.cs	13
Logic Layer	13
PHDCollectionLogic.cs	13
TaskHelper.cs	13
ScheduleHelper.cs	13
PHDDataHelper.cs	13
ParseHelper.cs	13
CSVHelper.cs	14
AzureStorageHelper.cs	14
AzureLoggingHelper.cs	14
LocalLoggingHelper.cs	14
Data Layer	14
Domain Layer	14
Framework Layer	14
PHDAFService.cs	14
System Design Considerations	15
Software Design Pattern	15
How To Collect Tags Adhoc Runs	16

Service Background and Overview

The purpose PHD Adhoc Collection Service is to extract PHD Tag (and metadata) data from the on premise historian and push the data points in Azure. The service should be able to perform various types of interpolation calculations, and supports hourly scheduled collection along with adhoc collection. Collection can also be configured to occur on a tag-by-tag basis.

This service should satisfy the following system design requirements for PHD data collection:

- Focus on Scalability and Ease of Deployment
- High Degree of Concurrency
- Fault-Tolerant and Verbose Logging/Traceability
- Resiliency and Recoverability
- Atomic Transactions
- Decoupled Architecture

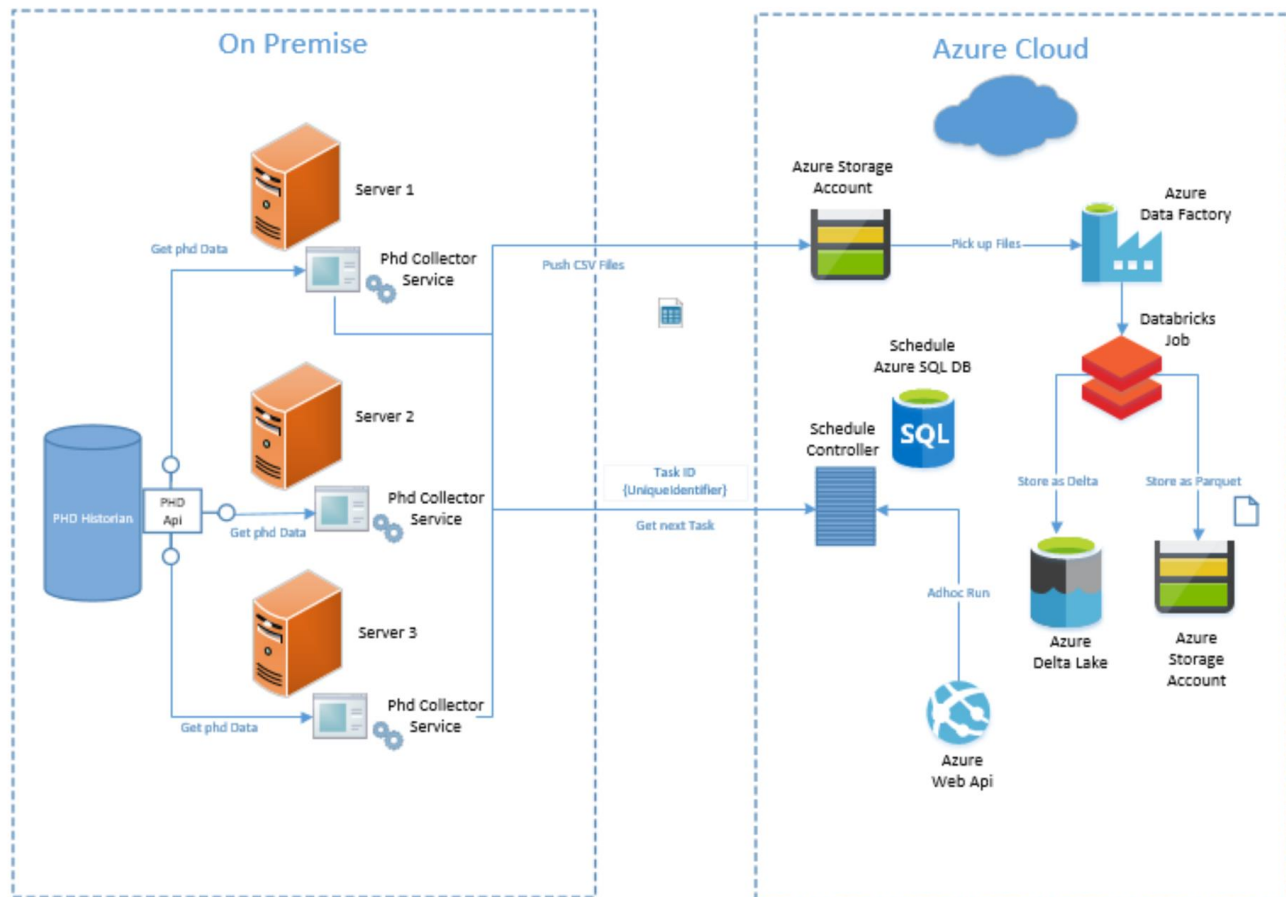
The PHD Adhoc Collection Service will be composed of the following systems:

- Windows Service deployed on On-premise VMs
- Web APHD deployed on Azure
- SQL Server Database deployed on Azure
- Azure Storage Accounts
- Azure Data Factory PHDpelines
- Azure Databricks notebooks
- Azure Delta Lake

Service Architecture Diagram

Logical Architecture

PHD Adhoc Collection – Logical Architecture



Azure Web APHD ADHOC Run – an aPHD to trigger an adhoc schedule record in the Azure SQL Database for any given time frame. It has two methods:

- Create Schedule – creates an adhoc schedule based on the timeframe
- Generate Tasks - generates a Task per PHD Tag for the timeframe

Azure SQL Schedule Controller – a group of stored procedures that act as the Business Layer to manage the schedules and tasks. The following stored procedures are:

- phd.sp_generate_tasks – generates all tasks for all schedules with a status of READY for all PHD tags in the PHD_HEADER table.
- phd.sp_generate_tasks_foraschedule – creates the tasks for a schedule.
- phd.sp_getnexttask – assigns a task for a caller (server/thread) based on a priority: Hourly before Adhoc then by status in the order of READY, FAILED with less than 3x of retries, IN PROGRESS/HOLD longer than 5 mins, ASSIGNED but not started over 5 mins. It handles concurrency by ensuring that no

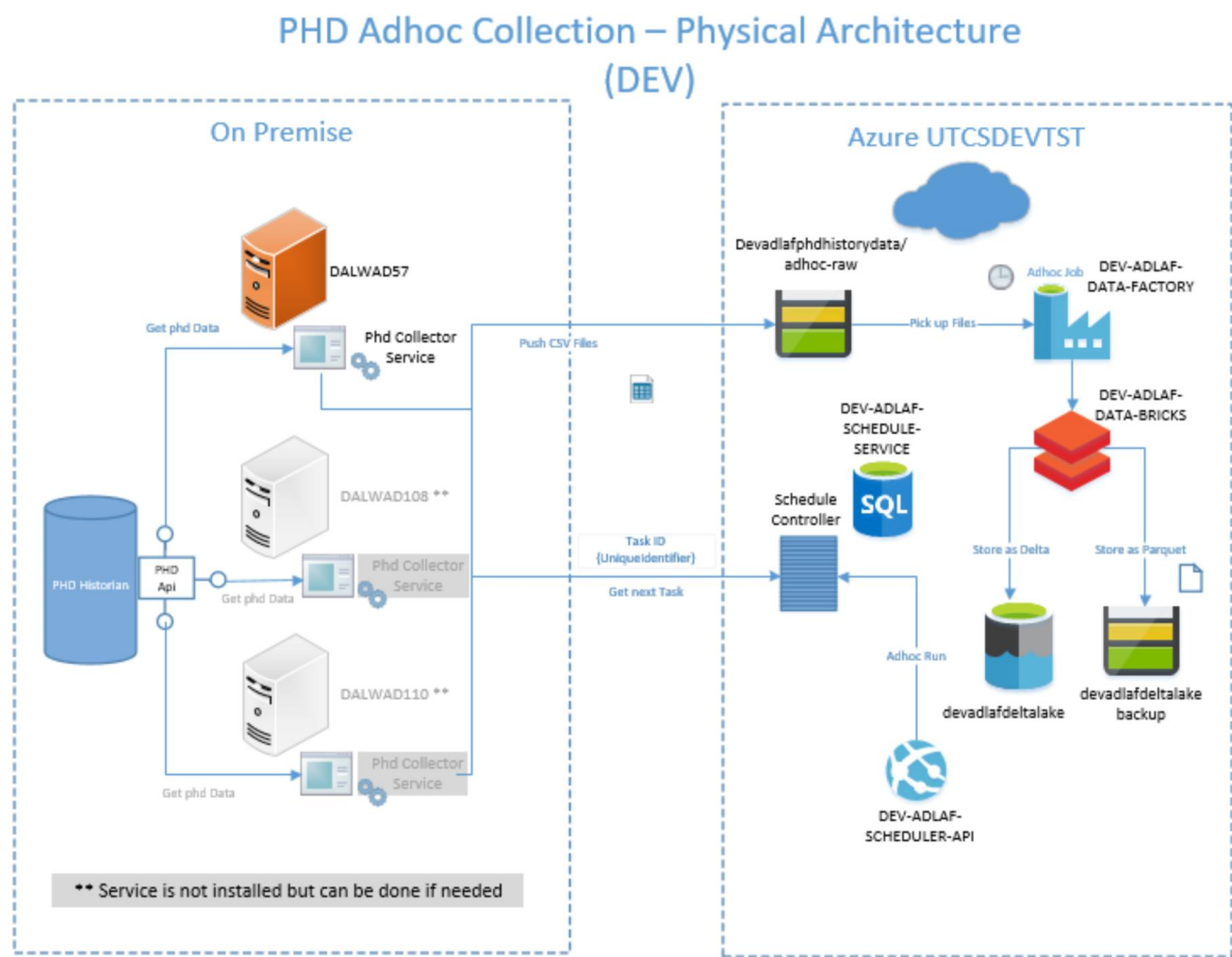
task is assigned to multiple callers. It assigns a unique sequence number for that task and locks it for the calling server/thread.

PHD Collector Service – a windows service that has multiple functions:

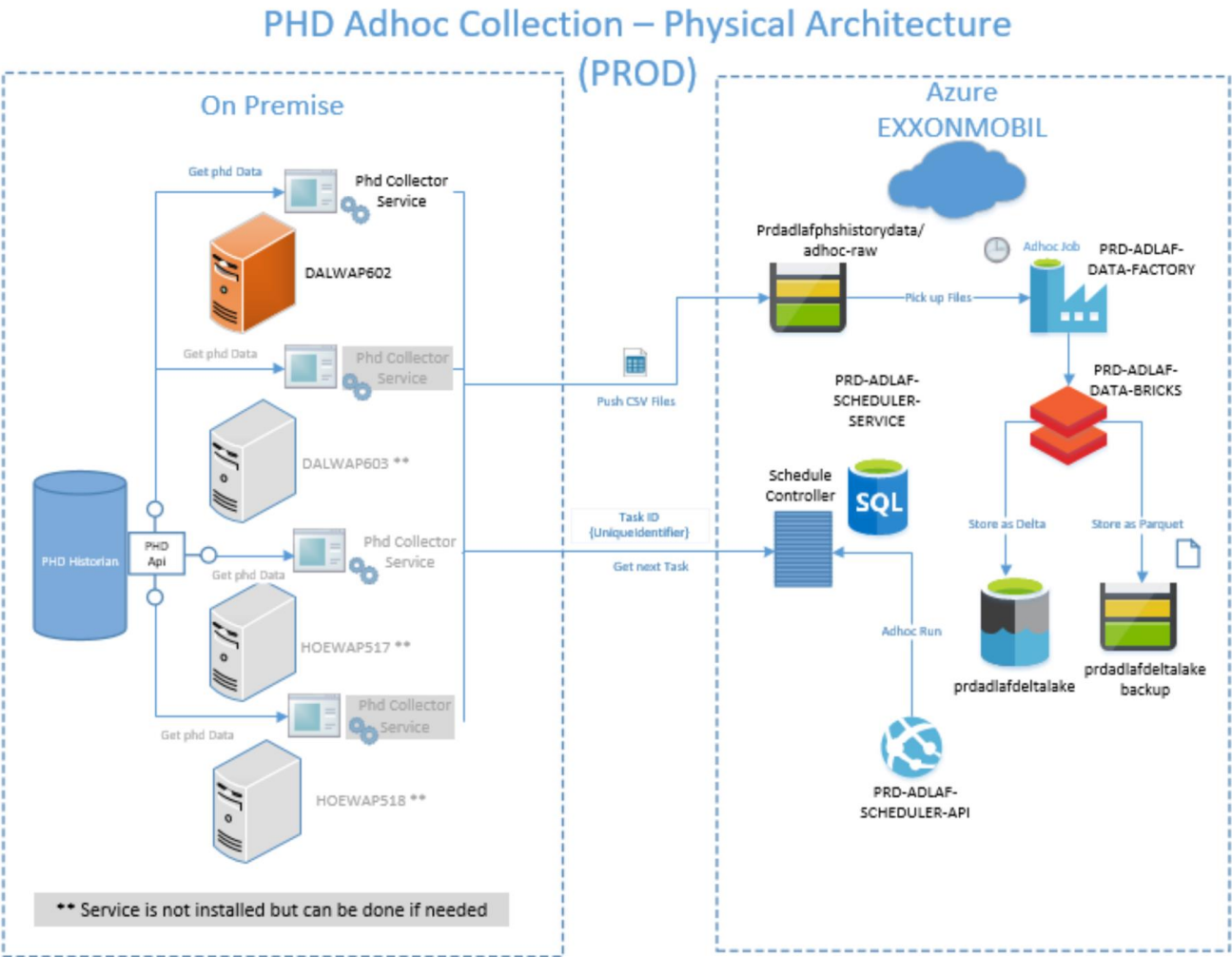
- Gets new tasks from Azure SQL Server
- Calls PHD APHD to get the data for a single tag for a given time frame
- Pushes the PHD Data into Azure Storage as a CSV file
- Can scale up multiple threads (upto 20) to perform each task concurrently
- Can scale out to as many servers as needed (currently 3 in DEV, 4 in PROD)
- Logs the task and schedule surveillance data

ADF Adhoc Job pick up CSV Files – an adhoc job that runs to pickup any CSV files from the storage account. It then calls Azure Databricks notebook to process the CSV and save it into an Azure Delta Lake (delta table) and Parquet file format in a storage account for backup.

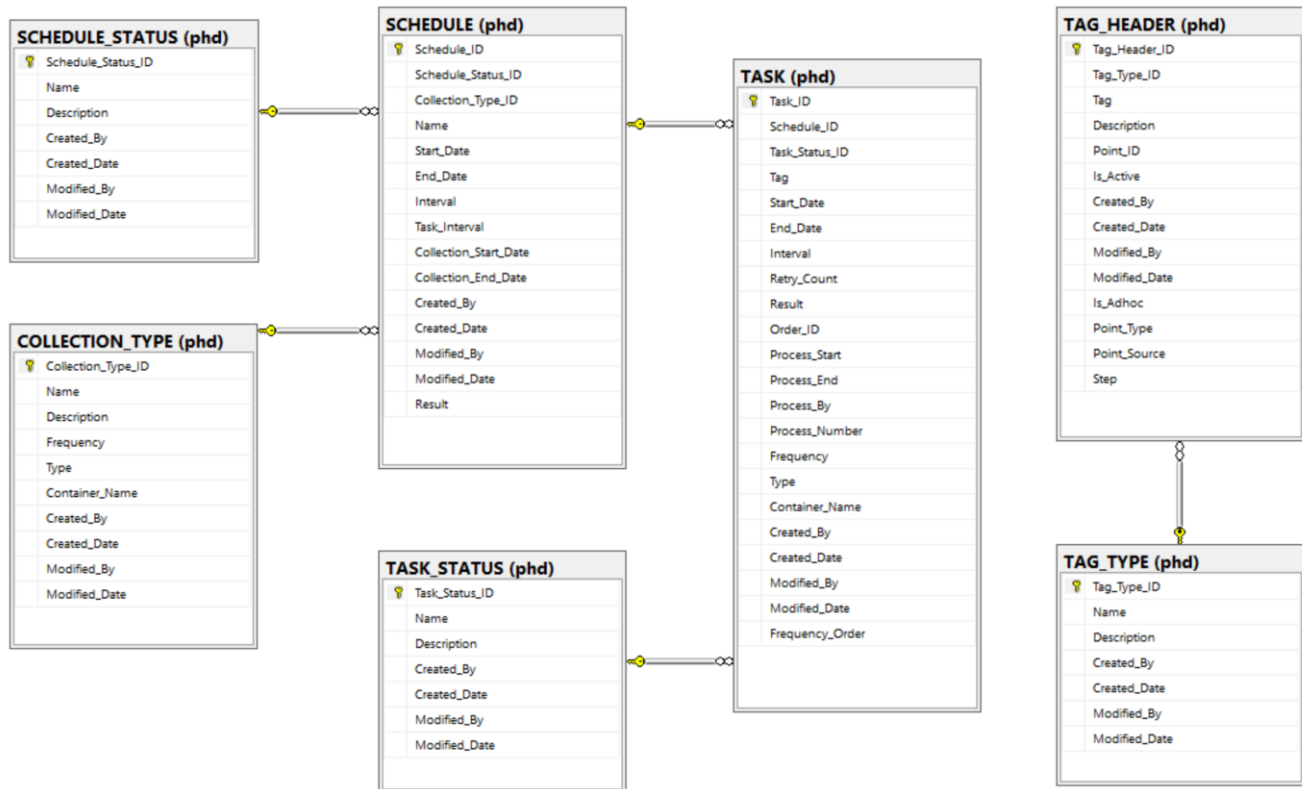
Physical Architecture (DEV Environment)



Physical Architecture (PROD Environment)



Data Model Diagram



TAG_HEADER – contains a full list of PHD Tags. An Is_Active flag indicates if it will be included in the hourly job. A Is_Adhoc flag indicates if it will included with the adhoc job run.

SCHEDULE – contains all the schedules and it's run details. Run details contain information on how the tasks should be broken out (HOUR/DAY), PHD parameters and schedule surveillance data (status, process_start, process_end).

TASK – contains all the tasks for a schedule along with PHD parameters and surveillance data (status, process_start, process_end, retry_count). Any errors or info regarding execution is also logged in the Result field.

SCHEDULE_STATUS – a look up table of status for schedules

- DRAFT – Not ready to be tasked
- READY – Ready to be tasked
- TASKED – schedule has been tasked (tasks created)
- IN PROGRESS – one or more task has started processing
- COMPLETED –all the tasks finished successfully
- FAILED – one or more task failed

TASK_STATUS – a look up table of status for tasks

- READY – Task is ready to be PHDcked up for processing

- ASSIGNED – Task has been PHDcked up by a thread to be processed
- IN PROGRESS – Task has started processing
- COMPLETED – Task finished successfully
- FAILED – Task finished with exception

COLLECTION_TYPE – a look up table of types for a schedule

- HOURLY-AVG - hourly job for average interpolation
- HOURLY-RAW – hourly job for raw data extraction
- ADHOC-AVG - adhoc run for average interpolation
- ADHOC-RAW - adhoc run for raw data extraction

PHD Adhoc Collection Service

The PHD Adhoc Collection Service is a windows service developed using .NET Framework 4.6.1 to execute PHD data collection tasks created by the Scheduling Service.

Installation and Setup

The following section goes into detail on how to make changes to the solution, deploy, and install on the required servers.

Create Installer

1. Build the solution using the “Release” configuration
2. Copy the contents of ... \DLAF_SERVICE_PHD\DLAF.SERVICE.PHD.ADHOC\bin\Release into a new directory
3. Copy the C:\Windows\Microsoft.NET\Framework64\vx.xxxx\InstallUtil.exe into the same directory
4. These are all the required files to perform the installation on the server

Configuration Settings

These are settings that can be configured for each particular installation
(DLAF.SERVICE.PHD.ADHOC.exe.config) **All config settings are required**

numberOfThreads: Number of concurrent executions per process (see multi-threading section for more details)

localExportDirectory: Directory for where local PHD data files should be written (ensure that the account executing the service has write access to that location)

localLogDirectory: Directory for where local log files should be written (ensure that the account executing the service has write access to that location)

aPHDmSubscriptionKey: Subscription key to APHDM

tenantEndpoint: OAuth 2.0 token endpoint (Get from Azure AD)

clientId: Client ID registration for this service in Azure

scope: Azure AD scope for APHDM

clientSecret: Client Secret for the Azure AD App registration for this service

loggingUrl: URL for Azure APHDM logging service

storageAccountName: Azure Storage account for exported PHD Data

storageAccountKey: Azure Storage account key for exported PHD Data

databaseServer: Database server containing the scheduling DB

databaseName: Database Name containing the scheduling DB

databaseId: Database ID (Read/Write Permissions)

databasePassword: Database Password

Install Instructions

1. Copy the installer folder to the server you would like to install in
2. Inside the coPHDed folder launch CMD
3. Run the following: installutil DLAF.SERVICE.PHD.ADHOC.exe
4. Supply the user credentials for the service account with access to the PHD server (saUIT-UDA-DLAF01-DEV@NA.xom.com / Password in ITPA)
5. Launch Windows Services and start the service
6. Set the windows service for delayed start & restart the service after 10 minutes in the event of failure

System Design Considerations

The following section goes over the thought process and design considerations for the choices of features that have been included or excluded from the system.

Software Design Pattern

This service has been developed using the repository pattern with nested projects with in a single solution:

- *DLAF-SERVICE-PHD-ADHOC (Service Layer)*
 - Implements multithreading on business logic methods
- *DLAF-SERVICE-PHD-ADHOC.Business (Logic Layer)*
 - Contains core business logic
 - Utilizes data access layer to extract data objects
- *DLAF-SERVICE-PHD-ADHOC.Data (Data Access Layer)*
 - Connects to SQL server and acts as a CRUD repository for data models

- *DLAF-SERVICE-PHD-ADHOC.Domain (Domain Layer)*
 - Contains a collection of entity and data transfer objects
- *DLAF-SERVICE-PHD-ADHOC.Framework (Framework Layer)*
 - Abstracts methods required to invoke the OSIssoft PHD SDK

The choice to use .NET 4.6.1 was due to the OSIssoft PHD SDK only supporting that version of .NET Framework.

The data access *DLAF-SERVICE-PHD-ADHOC.Data* uses .NET Standard 2.0 so that it's compatible with the Entity Framework Core (EFCore) library and .NET 4.6.1 projects.

Language version is set to C# 7.2 in order to support asynchronous methods on the Main method in Program.cs

The solution has been set up to execute as a console app when running locally, and a windows service when deployed. This is due to the the lack of troubleshooting capabilities when executing windows services using Visual Studio.

Multithreading

From initial benchmarking of collection times, we found that the average time to execute a single collection (single tag for one hour) averaged to be between 1-3 seconds. For a given hour, the minimum requirement would be to able process up to 14000 collection tasks (for ~7000 tags) within 60 minutes.

Using a single threaded design, this would mean that if this service was executing on a single VM it may take up to :

$$(14,000 \text{ tasks} \times 3 \text{ seconds}) / 60 = \text{Upwards of 700 minutes}$$

If executed across three VMs, this would go down to 233 minutes; however, with the limitations of an on-premise system we are limited in the degree to which we can horizontally scale.

When considering a multithreaded system design, there were a few things to take into consideration:

- Collection speeds and parallelism is further bottlenecked by the PHD Historian itself (only hosted on two servers for Kearn)
- At what degree of task parallelism do we experience no further benefit in collection time, either limited by the PHD Historian or the hosted server
- How do we ensure tasks are executed independently and execeptions will not affect other threads and the service overall

After testing, we discovered that 20 theads per process allowed us deliver sufficient performance while maintaining acceptable CPU and memory usage on the VMs.

With three VMs, this would mean 60 concurrently executing threads:

$$(14,000 \text{ tasks} \times 3 \text{ seconds}) / 60 / 60 \text{ threads} = \text{Upwards of 11.6 minutes}$$

Error Handling & Logging

Error handling within the service is primarily contained in the Service and Logic layers through a series of Try-Catches.

In the event of an exception, the error message is sent to the logging APHD in Azure, logged in the database on Task table, and also recorded in a log file on the server (in the event of no internet connectivity).

Log files are on the server are stored in two different ways:

- Service level exceptions are stored in the log file prefixed with the server name
- Task level exceptions are stored in the log folders prefixed with the Schedule GUID and the files are separated based on tag name

When an exception is caught the thread should immediately terminate, restart and start executing the next available task.

In the event the entire service crashes, the windows service is scheduled to continuously retry to restart after a 10 minute delay.

Service Layer

Program.cs

Start()

This method executes the collection and is determines how many task/threads to kick off based off the app configuration. It also contains the delay duration before the service checks to see if there is another active task.

Logic Layer

PHDCollectionLogic.cs

CollectData()

Main logic method that executes collection. Is invoked by the service layer to initialize collection for a given task.

TaskHelper.cs

Helper class that contains methods to perform business logic associated with PHD Scheduled Tasks

ScheduleHelper.cs

Helper class that contains methods to perform business logic associated with PHD Schedules

PHDDataHelper.cs

Helper class that contains methods to perform data cleansing on PHD data retrieved from the PHD SDK

ParseHelper.cs

Helper class that contains methods to perform generic data cleansing

CSVHelper.cs

Helper class that contains methods to write PHD data to CSV (if required)

AzureStorageHelper.cs

Helper class that contains methods to write PHD data to Azure blob storage

AzureLoggingHelper.cs

Helper class that contains methods to write logging data into the logging microservice in Azure

LocalLoggingHelper.cs

Helper class that contains methods to write logging data into the local on-premise file share

Data Layer

Implements Entity Framework Core through .NET Standard 2.0 (scaffolded off existing database)

Domain Layer

Contains classes for entities and data transfer objects.

Framework Layer

Contains classes associated with working with the PHD SDK

PHDAFService.cs

This class initializes all methods required to use the PHD SDK for data collection.

Constructor

The constructor calls the Init() method

Init()

This method initializes the connection to the PHD Server, using the default server connections specified on the VM. This can be viewed using the PHD SDK tool installed.

CollectDataAvg(string PHDTagName, DateTime StartTime, DateTime EndTime, double Interval)

This method collects PHD data by taking the average for a given time range and interpolation interval.

CollectData(string PHDTagName, DateTime StartTime, DateTime EndTime)

This method collects raw PHD data for a given time range.

System Design Considerations

The following section goes over the thought process and design considerations for the choices of features that have been included or excluded from the system.

Software Design Pattern

This service has been developed using the repository pattern with nested projects with in a single solution:

- *DLAF-SERVICE-PHD-ADHOC (Service Layer)*
 - Implements multithreading on business logic methods
- *DLAF-SERVICE-PHD-ADHOC.Business (Logic Layer)*
 - Contains core business logic
 - Utilizes data access layer to extract data objects
- *DLAF-SERVICE-PHD-ADHOC.Data (Data Access Layer)*
 - Connects to SQL server and acts as a CRUD repository for data models
- *DLAF-SERVICE-PHD-ADHOC.Domain (Domain Layer)*
 - Contains a collection of entity and data transfer objects
- *DLAF-SERVICE-PHD-ADHOC.Framework (Framework Layer)*
 - Abstracts methods required to invoke the OSIsoft PHD SDK

The choice to use .NET 4.6.1 was due to the OSIsoft PHD SDK only supporting that version of .NET Framework.

The data access *DLAF-SERVICE-PHD-ADHOC.Data* uses .NET Standard 2.0 so that it's compatible with the Entity Framework Core (EFCore) library and .NET 4.6.1 projects.

Language version is set to C# 7.2 in order to support asynchronous methods on the Main method in Program.cs

The solution has been set up to execute as a console app when running locally, and a windows service when deployed. This is due to the the lack of troubleshooting capabilities when executing windows services using Visual Studio.

How To Collect Tags Adhoc Runs

1. Identify the tags you want to collect for. Call the API in APIM, replace xxx with environment:

POST <https://xxx-adlaf.azure-api.net/scheduler/api/PHDSchedule/SetAdhocTags>

Body: {

```
"tags": ["tag1", "tag2", "tag3"...],  
"user": "testuser"  
}
```

2. Identify the date/time range you want to collect for. HOURLY is preferred than daily.
Call the API in APIM, replace xxx with environment:

POST <https://xxx-adlaf.azure-api.net/scheduler/api/PHDSchedule/CreateScheduleAdhoc>

Body:

```
{  
  "name": "Adhoc test run 1",  
  "user": "testuser",  
  "taskInterval": "HOURL", -- DAY or HOUR  
  "startDate": "2020-03-11 13:00:00",  
  "endDate": "2020-03-11 14:00:00"  
}
```