# Data Ingestor Auto Pilot – DIAP

## Easy Mass-Import of EXTERNAL Data (Files) into SAS

Stephan Weigandt

Sr Enterprise Data Engineer – Customer Advisory Support

FIRD – Field Innovation Research & Development

Stephan.Weigandt@sas.com

§sas

# Introduction

## Stephan Weigandt
Sr Enterprise Data Engineer
Field Innovation Research & Development, US Customer Advisory

➤ **Former SAS Partner (with various companies since 1998), joined SAS in 2016,**

  ➤ 24-years of architecting and creating automated solutions around CRM and closed loop predictive environments with focus on data processing and integrating.

  ➤ Co-created, designed and implemented a fully automated, closed loop predictive platform (building 350 models quarterly, which only needs to be supervised by 1.5 FTEs, including ETL, building model datasets, running QC on data and models, and creating automated PPTs for management consumption)

  ➤ Engaged across many industries and participated in customer references.

➤ **Everything is possible (…with SAS ☺).**

  ➤ And…if I have to do something twice, I will probably automate it.

  ➤ I enjoy streamlining and standardizing, "everything data" with the client needs always on top of the mind.

➤ **Love wearing hats and easily adjusted to the Californian lifestyle including doing yoga and more and more surfing.**

➤ **Certified in:**

  ➤ AWS Certified Solutions Architect – Associate

  ➤ SAS V6 and V8 Certified Professional Advanced

  ➤ Masters in Physics

➤ **SAS Computing Award winner 2000 (Europe, Dublin)**

§.sas

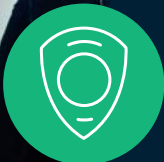# What is the situation?
## → COVID-19

### SAS helping and supporting local governments and researchers
- many urgent requests to make external data available for dashboards and/or modeling
- Data provided via websites or folder structures with missing and/or dirty metadata

### Time and Quality is of essence

Literally Life and Death could be on the line.

### Standard OOB tools quickly run out of options
- "difficult" files
- time consuming (one file at a time)
- → from taking hours to days to being a stumble block

§.sas

# Why DIAP?

Ingesting *(lots of) external data* into the SAS ecosystem *QUICKLY* by only *setting a few parameters* and *pushing a button* while *increasing the quality and confidence* into data insights.

# What are the Main Pain Points with External Data...

## Overwhelm with variety of data

- Files can be distributed over many and complex distributed directories
- XLSX, CSV, TXT, JMP, SHP, JSON,
- Fixed Width, XML
- variety of possible separators

## Challenging Metadata

- identify and deal with Metadata issues
  - file names and variable names can be too long and contain "weird" characters
- dealing with missing OR duplicate variable names

## Consistency and Standardization

- create reference tables
- no load date or load time tracking
- profiling option

§.sas

# How did we solve it?

Automation is key

DIAP allows to set a handful of parameters and let the machine do the job

**has power to traverse a folder structure**

**Import controlled into SAS Viya**

**consuming all the files in folder**

**Separator determination**

**review and create metadata**

**Keep track of Load Date/Time**

§.sas

# How did we solve it?

## Automation is key

DIAP allows to set a handful of parameters and let the machine do the job

**has power to traverse a folder structure**
- just provide top level directory name
- no limit in depth and complexity

**Import controlled into SAS Viya**
- Use data right away for dashboards/models
- Provide consistency for downstream referencing

**consuming all the files in folder**
- XLSX, CSV, TXT, JMP, SHP, JSON,
- Fixed Width
- XML

**Separator determination**
- Tab, CSV, Semicolon, Pipe

**review and create metadata**
- create reference tables
- identify and deal with Metadata issues
  - weird characters in names
- create profiles (integrates with "Autoprofiling on VA" solution)

**Keep track of Load Date/Time**
- Only upload modified files since last upload

§sas

# Why do I want to use DIAP?

Let the machine do the dirty work while getting a coffee or meditating ☺
- DIAP allows you to get your hands quicker on your data for better and deeper insights
- Consistency and Standardization with a higher quality

DIAP complements solutions like "SAS Information Catalog (IC)"
- IC will offer hooks to integrate easily DIAP

More confidence in 3rd party data

Avoid Data Swamp

§.sas

# "Heavier" version within VA –
# Interacting with DIAP via customized HTML within VA and
# Using Jobexecution

# Reference Table - Upload Log Table

# Reference Table - Variable Overview Table

# Integration with Profiling Solution (AutoProfiling from SAS Europe)

# Integration with Profiling Solution (AutoProfiling from SAS Europe)

# Demo…

§.sas

# Questions?

§sas.