# Working in the (Azure Data) Factory the right way

Riccardo Perico

rperico@lucient.com

# Riccardo Perico

BI & Power BI Engineer @ Lucient

PBIUG Italy Co-leader

Microsoft Data Platform MVP

Linkedin: https://linkedin.com/in/riccardo-perico

Email: rperico@lucient.com

Slideshare: https://www.slideshare.net/RiccardoPerico

GitHub: https://github.com/R1k91

# Partners

# Where are we? We're in the factory

*"**Big data** requires a service that can **orchestrate** and **operationalize** processes to refine these **enormous stores of raw data** into actionable business insights. Azure Data Factory is a **managed cloud service** that's built for these complex hybrid extract-transform-load (**ETL**), extract-load-transform (**ELT**), and **data integration** projects."*
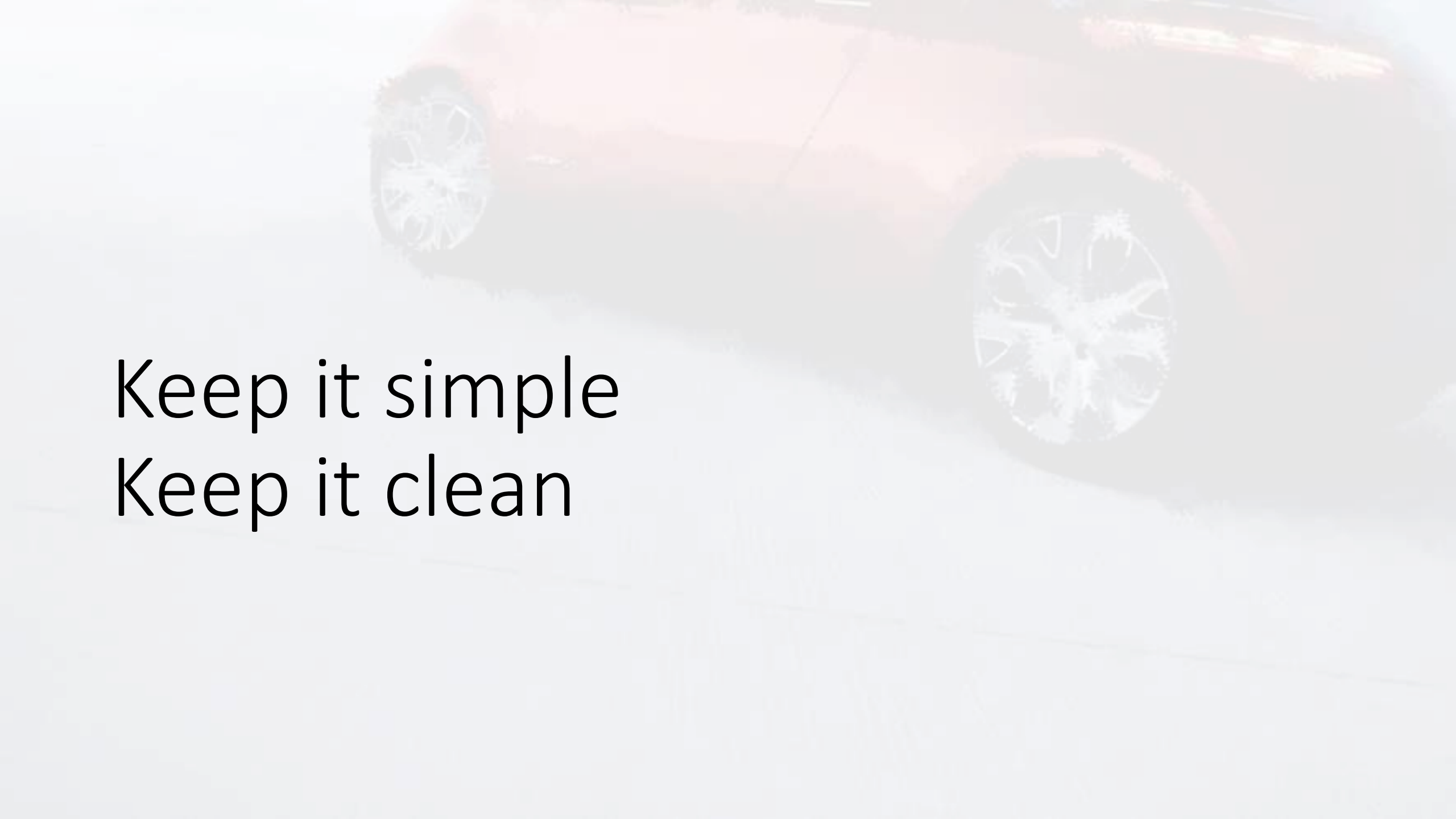
# What about this session?

*Starting to work in ADF it's easy but...*

*You'll learn by doing especially best practices and patterns*

*Let me share some lessons learnt from by background*
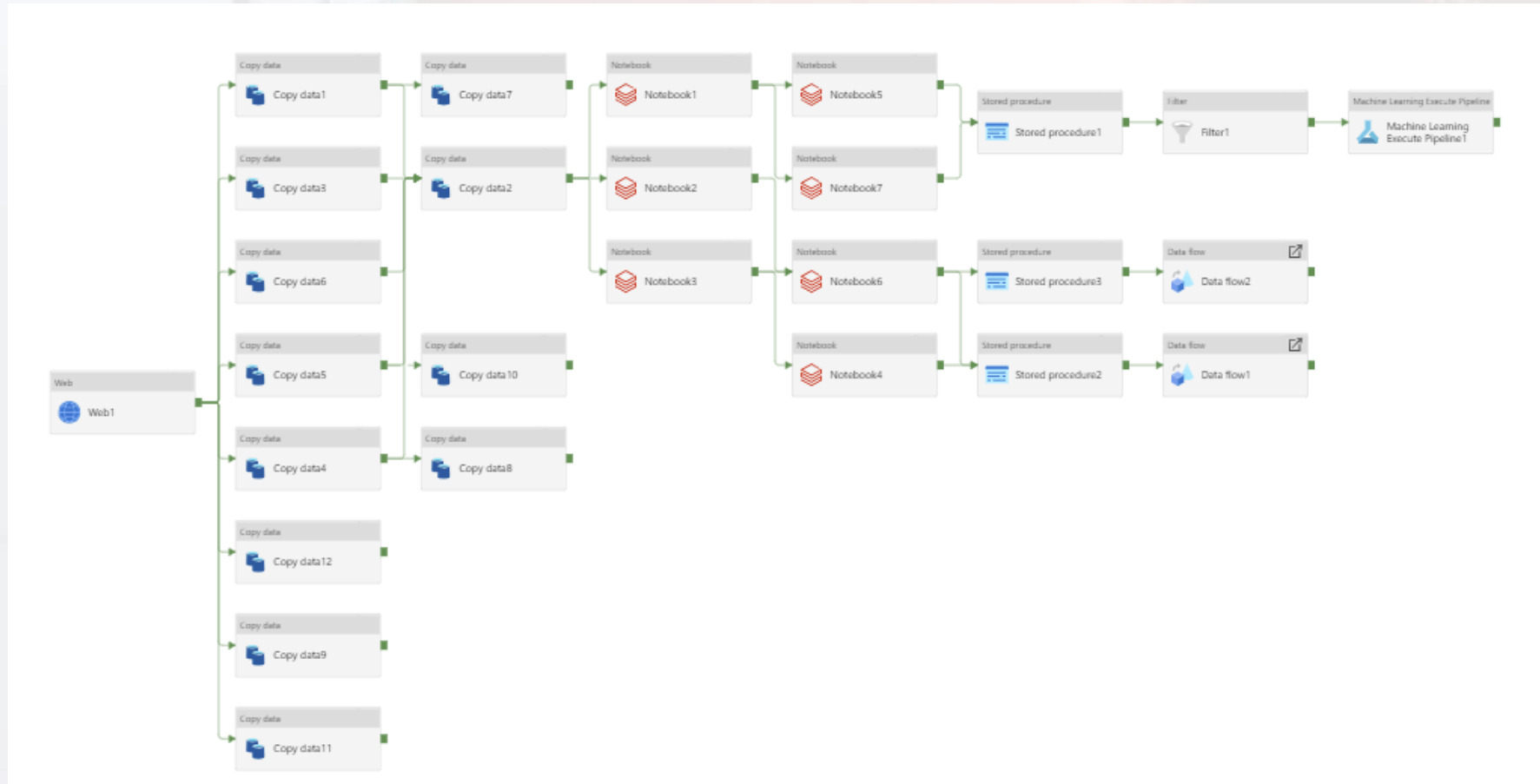
# Keep it simple
# Keep it clean

# General rules

- Give objects (all the objects!!!) a meaningful name
- Use a "pipeline approach"
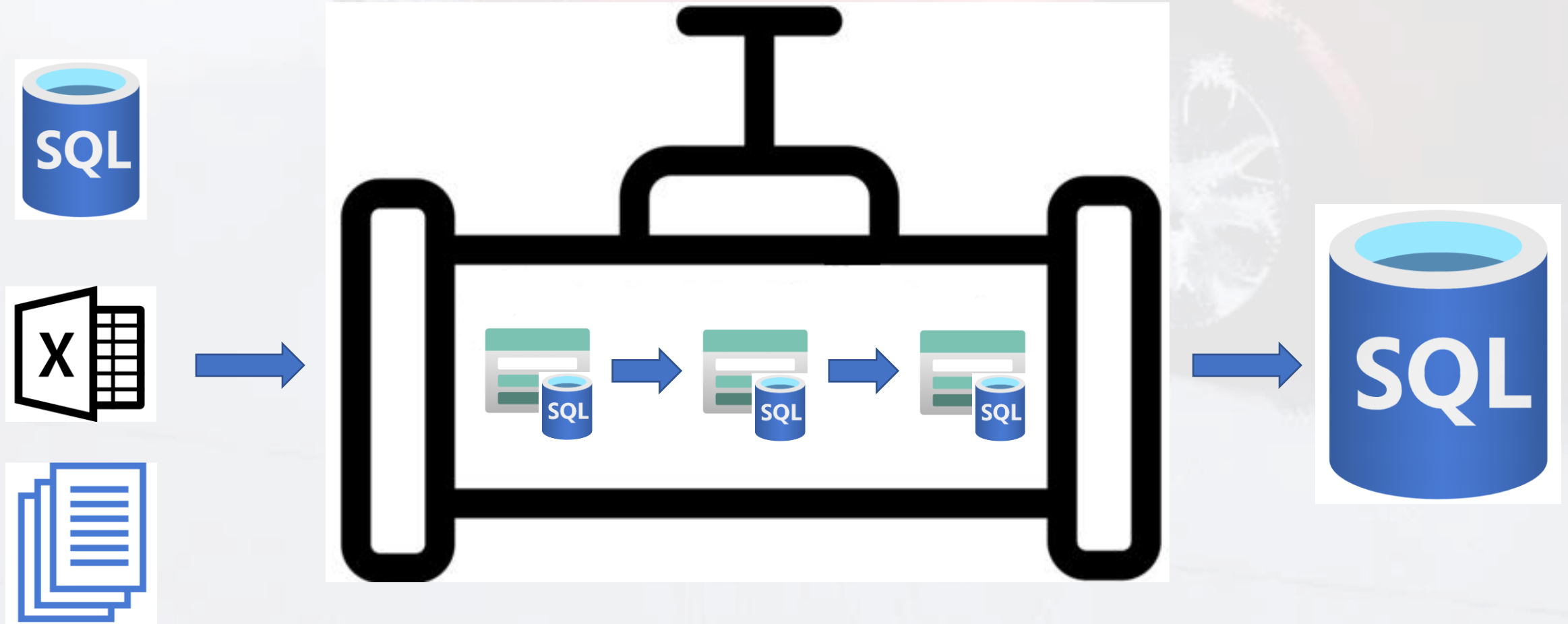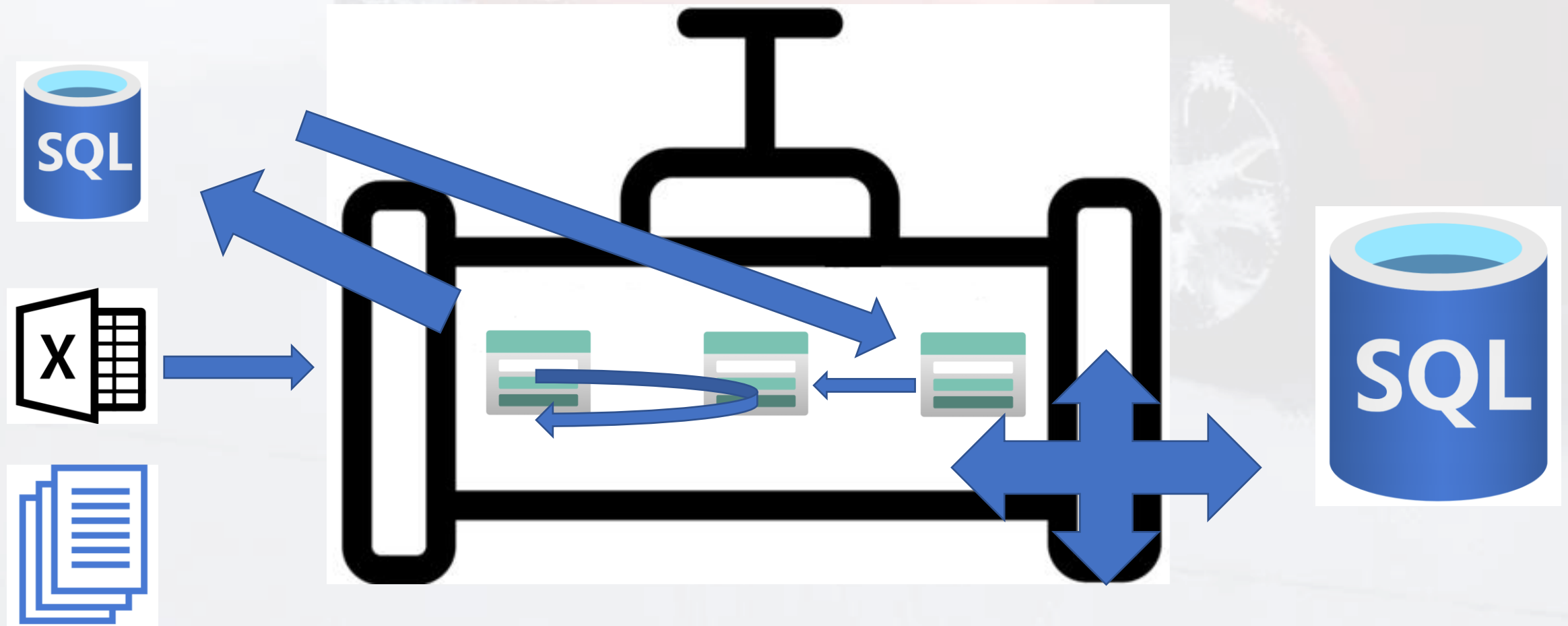- Organize objects accordingly
- …and many others

# Amarcord

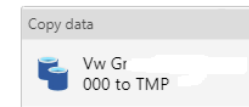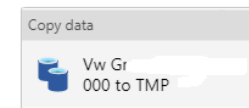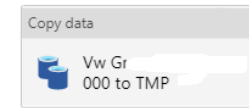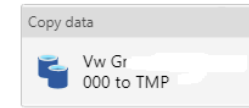# New tools same pitfalls

# "Pipeline approach": it flows in one way

# Pipeline approach: not that chaos… please

# Keep it clean

# Global Parameters

# Before and After

| Parameter Name | Parameter Value |
|---|---|
| pContainerName | myContainer |
| pArchiveFiles | True |

pContainerName
pArchiveFiles

pContainerName
pArchiveFiles

Can reference global directly

# One drawback

- Only pipelines can use them

> System variables
> Functions
∨ Global parameters
   gp_archive_container

Pipelines

> Functions
∨ Parameters
   dspContainer

Datasets

# Global Parameters

Demo

# Metadata Driven Approach

# What is a Metadata-driven approach?

- Don't do it manually
- Make it flexible and dynamic
- Make it configurable

# Is there anything out-of-the-box?

# Metadata-driven Approach

Demo

# Parametrizing Datasets

# Moreover…

# Parametrizing Datasets

Demo

# ADF loves AKV

# ADF loves AKV

- Best practice: always store **secrets outside ADF**
- Azure Key Vault securely stores and gives access to secrets
- PROS:
  - Developers can work **without knowing secrets**
  - Administrators can **setup and rotate secrets** without accessing ADF
- What should we store?
  - Passwords
  - Users Id / Identities
  - Servers' names / Services' names
  - …

# ADF loves AKV

Demo

# Managed Identities to get access

# What is this?

*"**Managed identities** provide an **identity for applications** to use when **connecting to resources** that support Azure Active Directory (Azure AD) authentication"*

# 2 types and both supported by ADF

| Property | System-assigned managed identity | User-assigned managed identity |
|---|---|---|
| Creation | Created as part of an Azure resource (for example, Azure Virtual Machines or Azure App Service). | Created as a stand-alone Azure resource. |
| Life cycle | Shared life cycle with the Azure resource that the managed identity is created with. When the parent resource is deleted, the managed identity is deleted as well. | Independent life cycle. Must be explicitly deleted. |
| Sharing across Azure resources | Can't be shared. It can only be associated with a single Azure resource. | Can be shared. The same user-assigned managed identity can be associated with more than one Azure resource. |
| Common use cases | Workloads that are contained within a single Azure resource. Workloads for which you need independent identities. For example, an application that runs on a single virtual machine. | Workloads that run on multiple resources and can share a single identity. Workloads that need pre-authorization to a secure resource, as part of a provisioning flow. Workloads where resources are recycled frequently, but permissions should stay consistent. For example, a workload where multiple virtual machines need to access the same resource. |

# Connect to ASQL via Managed Identity

CREATE USER [my-adf-name] FOR EXTERNAL PROVIDER;
GO


GRANT CONNECT TO [my-adf-name];
GO


ALTER ROLE [db_owner] ADD MEMBER [my-adf-name];
GO

# TTL integration runtime

# You definitely should use it

# TTL integration runtime

Demo

# My last advice for today…

No demo sorry…

# Useful Links

- ADF docs: http://tiny.cc/adfrw1

- Metadata driven out-of-the-box: http://tiny.cc/adfrw2

- Metadata driven by Paul Andrew: https://github.com/mrpaulandrew/procfwk

- Managed identities: http://tiny.cc/adfrw3

- Connect to ASQL using ADF's Managed Identity: http://tiny.cc/adfrw4

Q&A

Thank you !