# Modern Analytics Platform
## Lambda Architecture in Azure

Simon Whiteley | Adatis

10/03/2017

Microsoft Partner
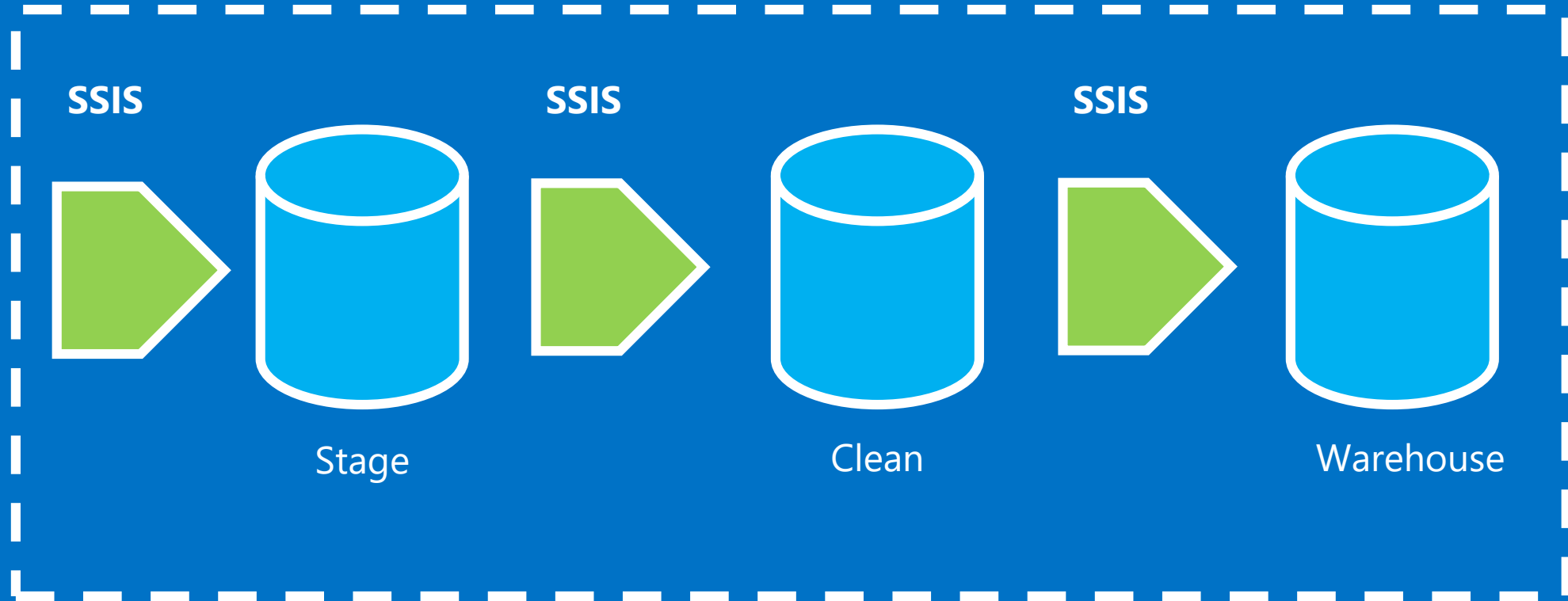Gold Data Analytics
Gold Data Platform
Gold Cloud Platform
Microsoft

MVP Microsoft® Most Valuable Professional

adatis

# Overview

Cloud BI

What is Lambda?

The Native Azure Approach

Alternative Models

# History

# One-Box SQL BI Architecture



**SSIS** → Stage → **SSIS** → Clean → **SSIS** → Warehouse
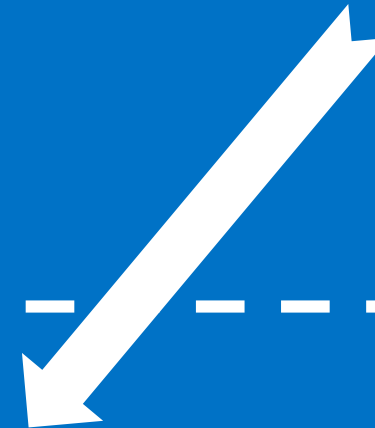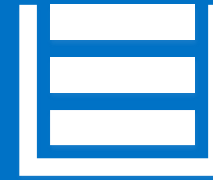
On-Prem SQL Server

"Big Data" Solutions

Modern Analytics Platform

My Life Goal:
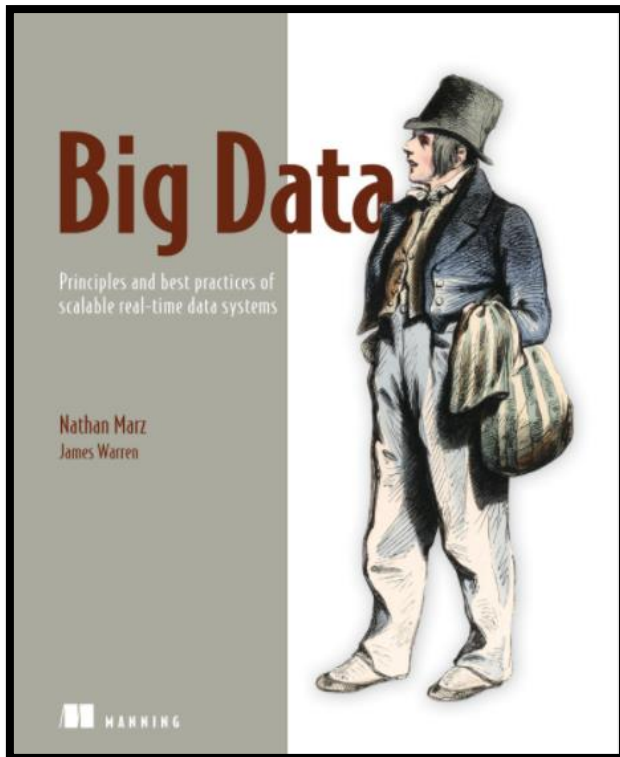Never to manage another Server

# MAP Wish List

- Can Handle Massive Datasets

- Linearly Scalable

- Near Real-Time

- Fault Tolerant

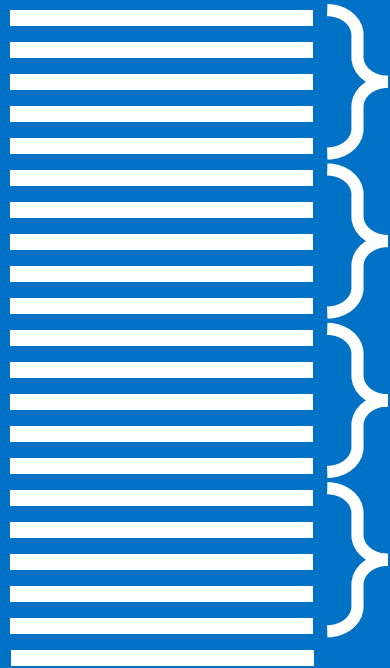- Low Tech Barrier for SQL Devs

LAMBDA

# Lambda Architecture

Use Batch and Stream technologies together to balance latency, throughput and fault-tolerance
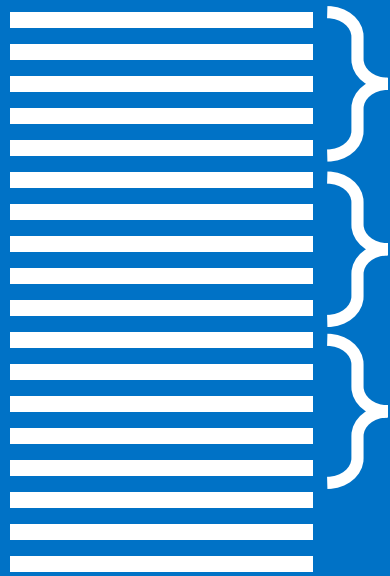


Nathan Marz & James Warren

adatis

# PROBLEM

# SOLUTION

Speed Layer

Message

Publish / Subscriber

Spout / Bolt
Storm Topology

Serving Layer
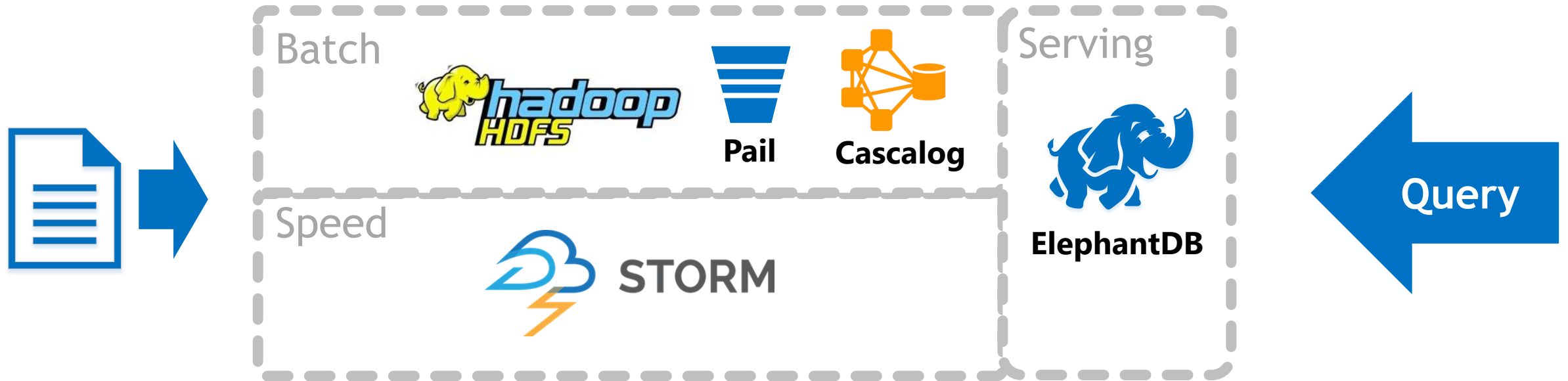
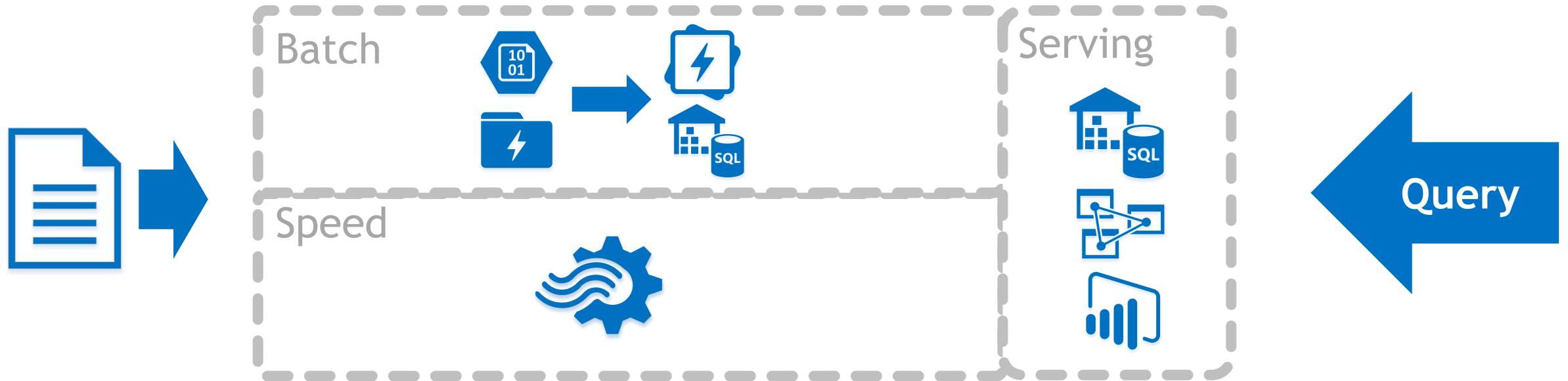Batch Views

Sharded
Database Layer

# The Marz Lambda Architecture

# The Azure Approach
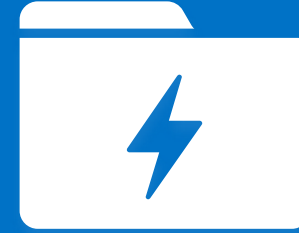
# Applying Lambda to Azure
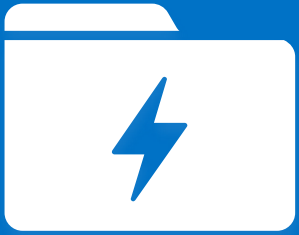
# Batch - Storage

## Blob Storage

- HDFS

- Hot/Cold Storage Tiers

- Limited Security

- File Size Limitations
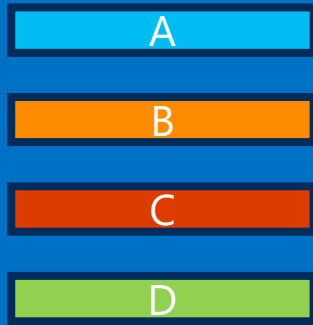
- Widely Compatible / Available

## Azure Data Lake Store

- WHDFS

- Single Pricing Model

- AAD-Integrated Security
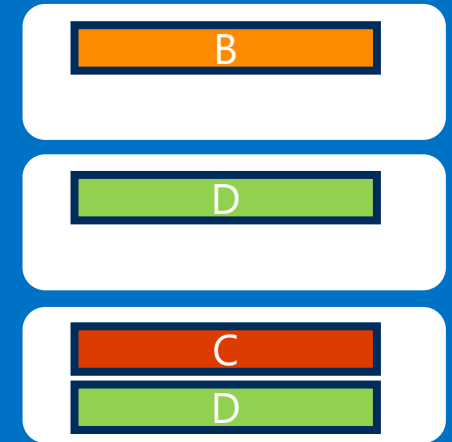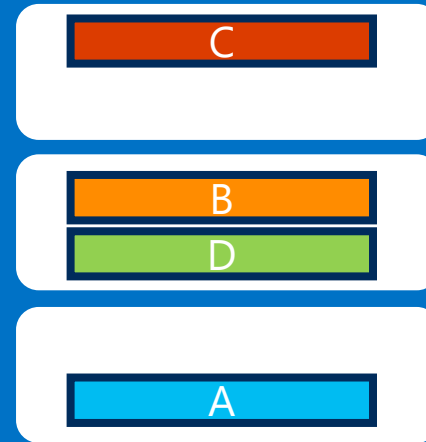
- No Limitations

- Still Immature

# Batch - Compute

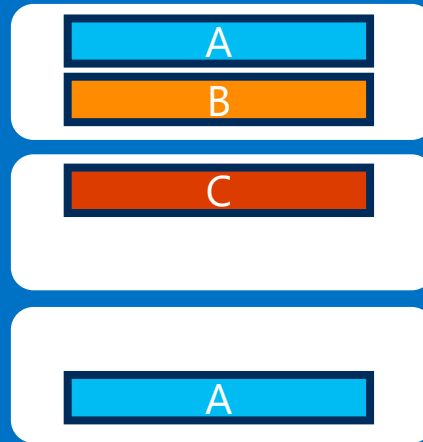## Azure Data Lake Analytics

- Pay Per Query / Unit

- U-SQL

- Outputs Structured/Unstructured

- Uses MapReduce-style processing

- Batch Mode

## Azure SQL DataWarehouse

- Pay Per Hour / Node

- T-SQL

- Fully Structured

- Can use MapReduce via Polybase

- Batch or Live Query

Azure Data Lake Analytics

# Azure Streaming Analytics

- Only PaaS Native Offering

- Uses SQL Language

- Built-in Azure Integrations

- Can Vertically Partition Files

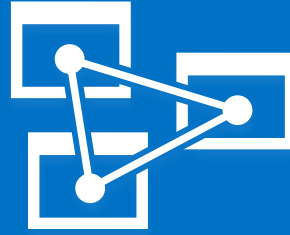- Can Write to Multiple Outputs

/Input/2017/06/19/0900.csv
/Input/2017/06/19/1000.csv
/Input/2017/06/19/1100.csv
/Input/2017/06/19/1200.csv

## Azure SQL DataWarehouse

- Low Concurrency (32!)
- Direct Query via Polybase
- Huge data capacity

## Azure Analysis Services

- High Concurrency
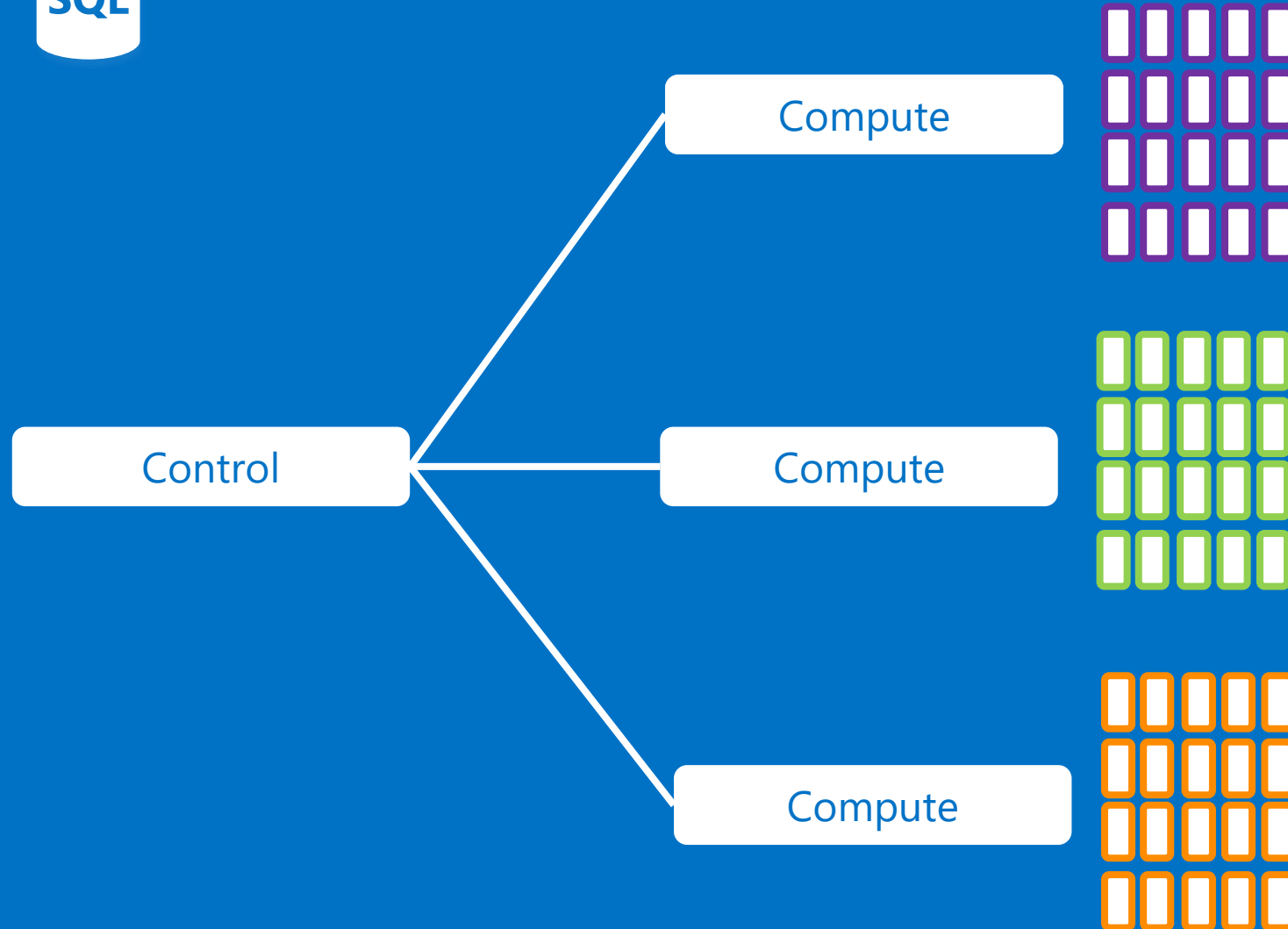- Scheduled Refresh / Direct over DBs
- Model Size Limits

## PowerBI

- High Concurrency
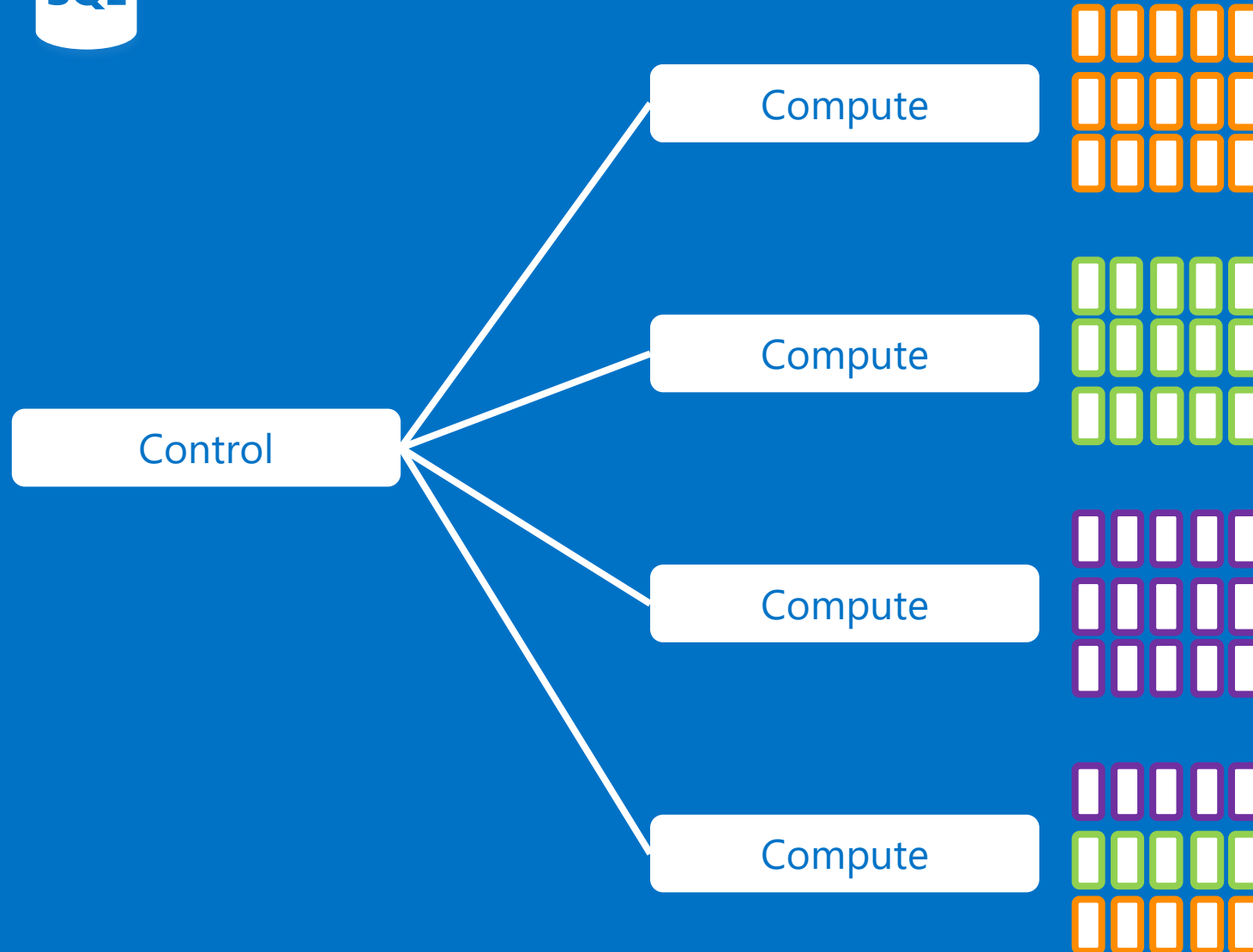- Scheduled Refresh / Direct over DBs
- Model Size Limits
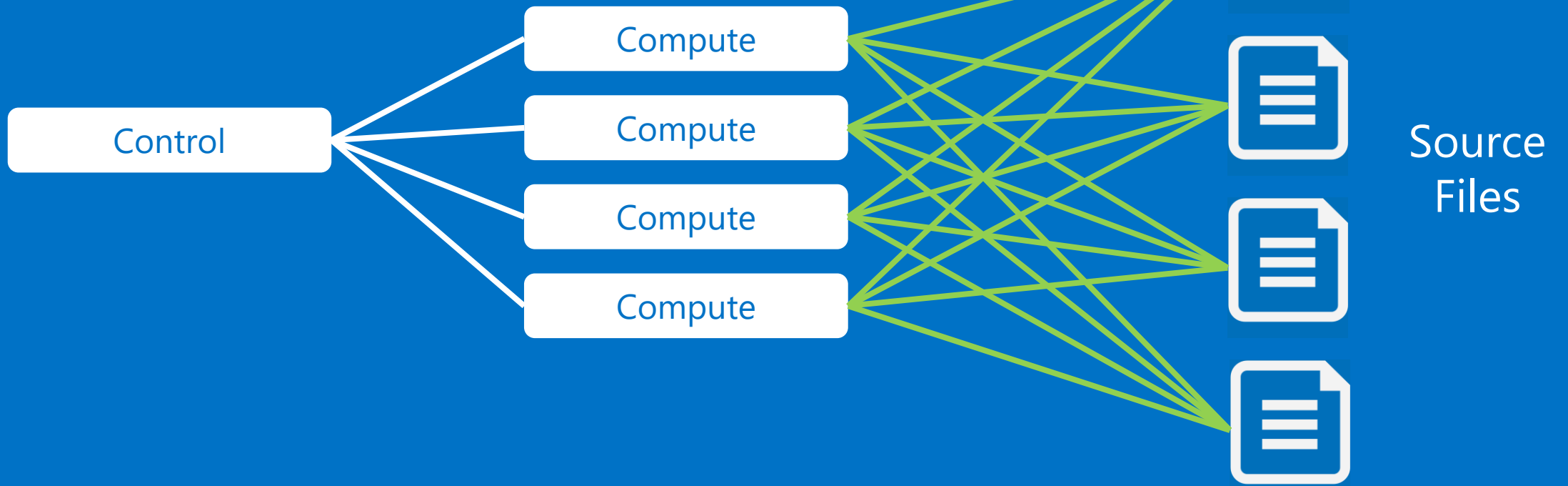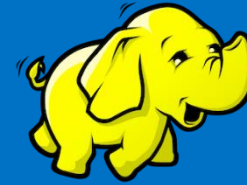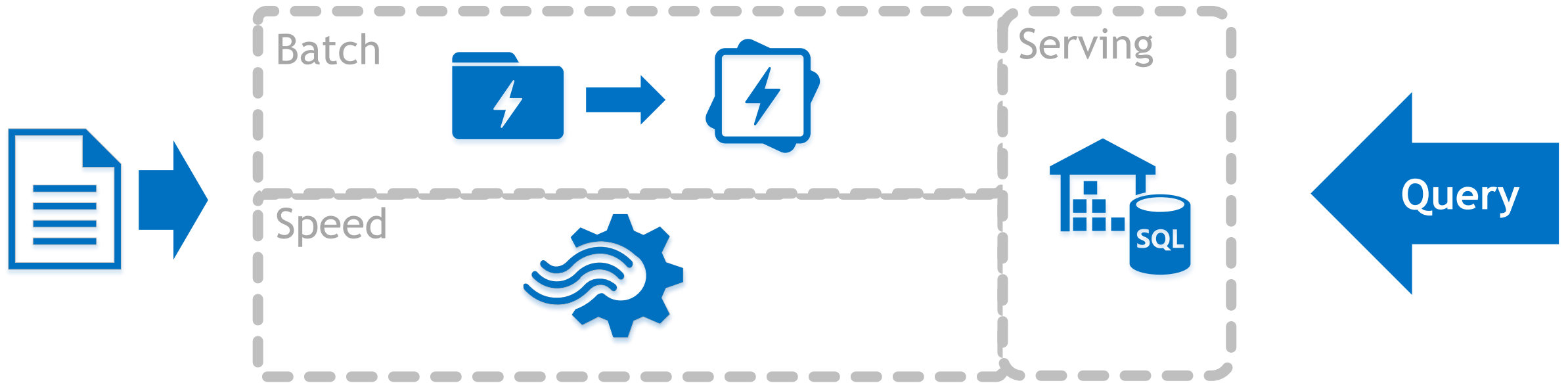
Azure SQL DataWarehouse

Azure SQL DataWarehouse

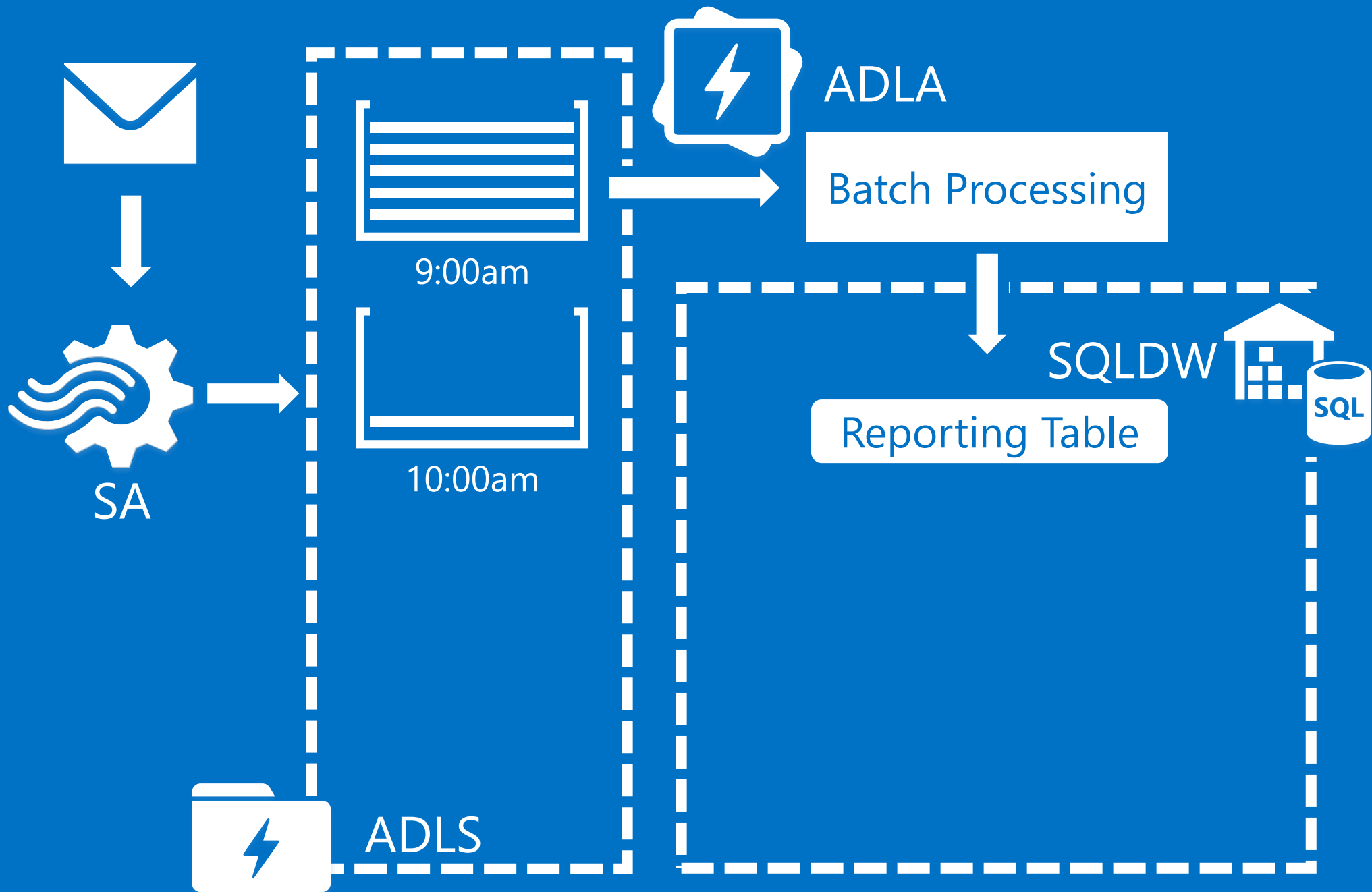# Applying Lambda to Azure

Start U-SQL Job

Call SQLDW Loading Procedure(s)

Every Hour

Azure Data Factory

Demo

Reset External Tables · Start U-SQL Job · Call SQLDW Loading Procedure(s) · Reset External Tables · Every Hour · Azure Data Factory
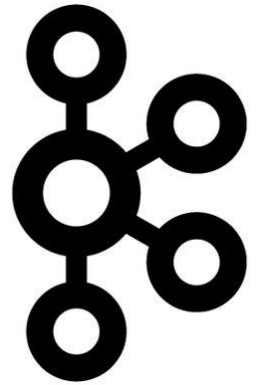
# Conclusions & Variations

Who Needs a Database?

Parallel Streaming

# Thanks for Listening

Simon Whiteley

@MrSiWhiteley

adatis

http://blogs.adatis.co.uk