The following is a capstone presentation that showcases some of my work in creating a data validation process for automated data loads in SQL Server. Some information has been redacted to protect client confidentiality.

# INTRODUCTION

## Overview and Goals

# Sample Procedure

- Current procedure had three validation measures
    - Invoice Date
    - Price Changes
    - Revenue

```
23   DECLARE    @RESULT                NUMERIC
24   SET @RESULT =
25        (SELECT DATEDIFF(DD, GETDATE(), MAX(InvoiceDate))
26        FROM etl.Invoices
27        WHERE InvoiceType = 'INV'
28        --AND InvoiceDate <= '2018-01-01' --test screnario
29        )
30
31
32
33   IF (@RESULT < -2)
34        BEGIN
35             SET @MessageSubject = 'Validation Alert Notification 1: ' ;
36             SET @MessageBody =  'We have not loaded any new Invoices in ' + CAST(ABS(@RESULT) AS VARCHAR(10)) + ' days.'
37             EXEC msdb.dbo.sp_send_dbmail
38                  @recipients         = @DeliveryRecipients,
39                  @copy_recipients    = @ETLRecipients,
40                  @body_format        = 'HTML',
41                  @body               = @MessageBody,
42                  @subject            = @MessageSubject ;
43             --PRINT @MessageSubject + '        ' + @MessageBody
44        END
```
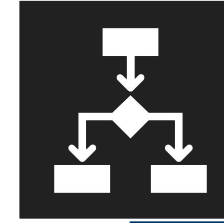
# Research

**Client Team Meetings**

- Figure out the metrics that are important to each client
- Establish thresholds and determine action upon failure

**Data Lead Meeting**

- Gain deeper understanding of the structure of each database
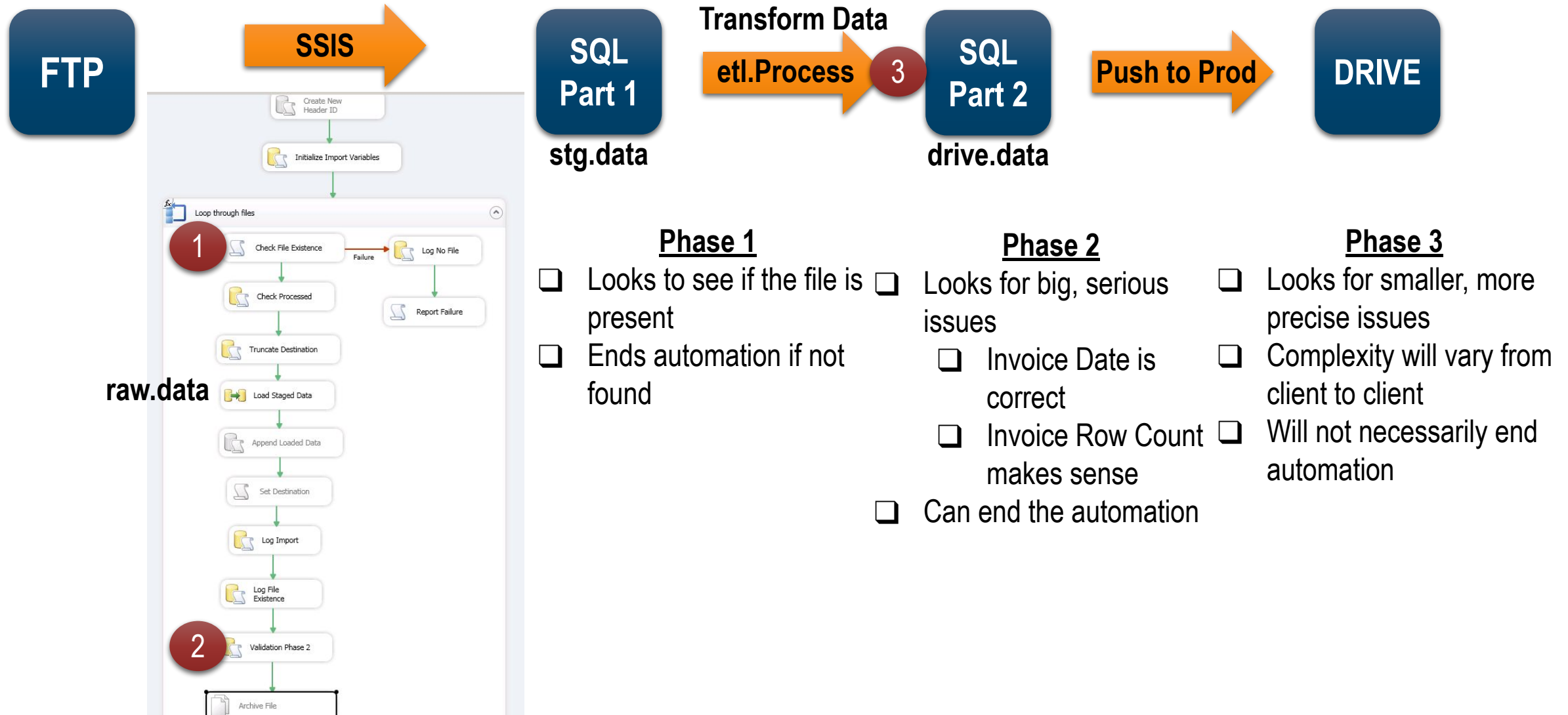- Go over clarifications specific to each database
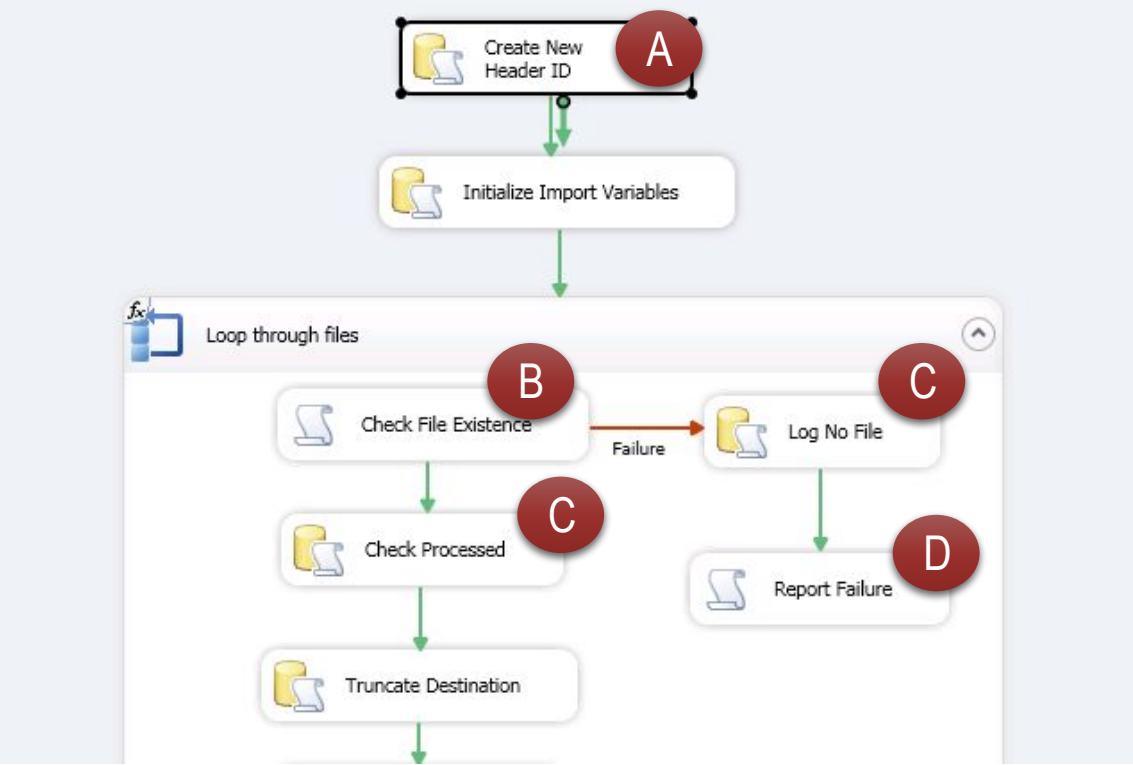
**Preparation of Structure**

- Individual creation of automated validation process
- Reach agreement on joint procedures

# THE PROCESS

# General Data Automation Process with Validation Phases

**FTP** → **SSIS** → **SQL Part 1** → Transform Data **etl.Process** ③ → **SQL Part 2** → **Push to Prod** → **DRIVE**

**stg.data**

**drive.data**



**raw.data**

### Phase 1
- ❑ Looks to see if the file is present
- ❑ Ends automation if not found

### Phase 2
- ❑ Looks for big, serious issues
  - ❑ Invoice Date is correct
  - ❑ Invoice Row Count makes sense
- ❑ Can end the automation

### Phase 3
- ❑ Looks for smaller, more precise issues
- ❑ Complexity will vary from client to client
- ❑ Will not necessarily end automation

# Phase One



A. Create new entry in the Load ID table

B. Check to see if file exists in the specified location

C. If success – continue with the load.     If failure – log failure

D. On failure – report failure, exit job, send email reporting that the file does not exist

**Load ID Table**

| | HeaderID | LoadName | LoadDate |
|---|---|---|---|
| 1 | 1 | Load | 2019-08-01 10:29:22.893 |
| 2 | 2 | Load | 2019-08-01 10:29:24.113 |
| 3 | 3 | Load | 2019-08-01 10:29:24.990 |
| 4 | 1 | ad | 2019-08-01 10:31:38.103 |
| 5 | 1 | Load | 2019-08-02 09:24:47.370 |
| 6 | 1 | ekly Load | 2019-08-02 09:50:24.167 |
| 7 | 2 | ekly Load | 2019-08-02 13:30:29.543 |
| 8 | 3 | ekly Load | 2019-08-02 13:37:39.293 |
| 9 | 4 | ekly Load | 2019-08-02 14:39:49.533 |
| 10 | 5 | ekly Load | 2019-08-02 14:43:55.117 |
| 11 | 6 | ekly Load | 2019-08-05 12:18:17.860 |

**Validation Log Table**

| | LogID | HeaderID | LoadDate | ImportID | DatabaseName | RuleName | RuleResult | Pass/Fail/Warning | Reason | LoadName |
|---|---|---|---|---|---|---|---|---|---|---|
| 16 | 16 | 1 | 2019-08-02 09:50:24.167 | 1 | | Check File Existence | 2 | Pass | The specified file was found | |
| 17 | 17 | 1 | 2019-08-02 09:50:24.167 | 1 | | Row Count | 3 | Pass | There were 3 records found | |
| 18 | 18 | 1 | 2019-08-02 09:50:24.167 | 455 | | Check File Existence | 0 | Fail | The specified file was not found | |

# Phase Two of Validation

- Found in the SSIS Package
- Looks for big issues where the data load would end if these issues are found
  - Invoice Date is incorrect
  - Invoice Row Count is severely off
- If one procedure fails, it will not go on to the next
- If you are using a SQL Agent Job, this procedure will cause the SSIS package to error out, thus causing the Agent Job to fail too
- Email is sent out if any of the validation procedures fail

```
9  ALTER PROC [etl].[ValidationPhase2]
10 (
11 @ImportID INT
12 )
13 AS
14
15 --Declare Variables that will be used in every Procedure
16 DECLARE @DatabaseName VARCHAR(255)
17 DECLARE @LoadName VARCHAR(255) = '
18 DECLARE @LoadDate DATETIME = (SELECT MAX(LoadDate) FROM etl.LoadIDHeader WHERE LoadName = @LoadName)
19 DECLARE @HeaderID INT = (SELECT MAX(HeaderID) FROM etl.LoadIDHeader WHERE LoadName = @LoadName)
20 DECLARE @RULERESULT VARCHAR(255)
21 DECLARE @PFW VARCHAR(255)
22 DECLARE @REASON VARCHAR(255)
23
24 BEGIN TRY
25 --List of Validation Procedures
26 EXEC etl.              tionInvoiceDate @HeaderID, @LoadDate, @ImportID, @DatabaseName, @RULERESULT, @PFW, @REASON, @LoadName
27 EXEC etl.PhaseTwoErrorTest @HeaderID, @LoadDate, @ImportID, @DatabaseName, @RULERESULT, @PFW, @REASON, @LoadName
28 END TRY
29
30 BEGIN CATCH
31     SELECT 1/0
32 END CATCH
```

# Phase Three of Validation

- Found at the end of or after etl.Process

- Looks for smaller issues in the data load
  - Revenue is 10% less or greater than the previous week's average
  - Number of customers is 10% less or greater than the previous week's average

- If one of the procedures finds an issue, it will log that result as a warning and move on to the next procedure

- Will not necessarily end the data load

- Email is sent out at the very end of this procedure

```
 8  ALTER PROC [etl].[ValidationPhase3]
 9  AS
10
11
12  --Declare variables that will be used in every procedure
13  DECLARE @IMPORTID VARCHAR(10) = 'N/A'
14  DECLARE @DatabaseName VARCHAR(50)
15  DECLARE @LoadName VARCHAR(50) =
16  DECLARE @LoadDate DATETIME = (SEL                    FROM etl.LoadIDHeader WHERE LoadName = @LoadName)
17  DECLARE @HeaderID INT = (SELECT MAX(HeaderID) FROM etl.LoadIDHeader WHERE LoadName = @LoadName)
18  DECLARE @RULERESULT VARCHAR(255)
19  DECLARE @PFW VARCHAR(50)
20  DECLARE @REASON VARCHAR(255)
21
22
23  --List of validation procedures
24  EXEC etl.AverageRevenue @HeaderID, @LoadDate, @ImportID, @DatabaseName, @RULERESULT, @PFW, @REASON, @LoadName
25  EXEC etl.AverageCustomers @HeaderID, @LoadDate, @ImportID, @DatabaseName, @RULERESULT, @PFW, @REASON, @LoadName
26  EXEC etl.AverageProducts @HeaderID, @LoadDate, @ImportID, @DatabaseName, @RULERESULT, @PFW, @REASON, @LoadName
27
28
29  --Procedures for the email
30  EXEC etl.EmailValidationResults @DatabaseName, @LoadDate, @HeaderID
```

# Email

```
9  ☐ALTER PROC [etl].[EmailValidationResults]
10 (     @DatabaseName VARCHAR(255)
11 ,     @HeaderID INT
12 ,     @LoadName VARCHAR(255)
13 )
14 AS
15
16
17 ☐CREATE TABLE #Temp
18 (     [LogID]              INT
19 ,     [HeaderID]           INT
20 ,     [LoadDate]           DATETIME
21 ,     [ImportID]           INT
22 ,     [DatabaseName]       VARCHAR(255)
23 ,     [RuleName]           VARCHAR(255)
24 ,     [RuleResult]         VARCHAR(255)
25 ,     [Pass/Fail/Warning]  VARCHAR(255)
26 ,     [Reason]             VARCHAR(255)
27 ,     [LoadName]           VARCHAR(255)
28 )
29
30
31 ☐INSERT INTO #Temp
32 SELECT *
33 FROM etl.ValidationLog
34 WHERE 1=1
35     AND LoadName = @LoadName
36     AND HeaderID = @HeaderID
37
```

Selects all entries in the Validation Log Table where the Load Name and Header ID are a match

Inserts those entries into a temporary table

Passes the temporary table through HTML and allows user to set the subject, body, and recipients

# Sample email sent after load fails.

Reply    Reply All    Forward    IM

Emily Hostetler

**Capstone Testing Validation Results**

## Load Validation Results

| HeaderID | LoadDate | DatabaseName | RuleName | Pass/Fail/Warning | Reason |
|----------|----------|--------------|----------|-------------------|--------|
| 8 | 2019-08-06T08:46:00.163 | | Load Date | Pass | The most recent file was loaded on 2019-07-26 |
| 8 | 2019-08-06T08:46:00.163 | | Records Count | Pass | There were 2685 records. This value is acceptable |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Other | Pass | Drive Quantity was 1025987. This value is above the monthly low. |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Orlando | Pass | Drive Quantity was 65431. This value is above the monthly low. |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Atlanta | Pass | Drive Quantity was 699397. This value is above the monthly low. |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Cleveland | Pass | Drive Quantity was 1242586. This value is above the monthly low. |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Detroit | Pass | Drive Quantity was 667525. This value is above the monthly low. |
| 8 | 2019-08-06T08:46:00.163 | | Drive Quantity - Miami | Warning | Drive Quantity was only 227094. The previous low was 251668 |

# CUSTOM PROCEDURES

# Custom Procedures – [_____] Row Count

- Goal: Check that the row count is above the minimum for the last 30 days
  - On Pass – record pass and continue
  - On Warn – record warning and continue

```
12  DECLARE @Minimum INT;
13  DECLARE @Result INT;
14  SET @Result = (SELECT COUNT(*) FROM stg.[_____]RE (CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
15                                           ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
16                                           END) = (SELECT MAX(CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
17                                           ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
18                                           END) FROM st[_____]--DATEADD(DAY, -2, CAST(CURRENT_TIMESTAMP AS DATE)))
19
20
21  SET @Minimum = (SELECT MinRecords FROM
22              (
23                  SELECT MIN(NumRecords) OVER (ORDER BY LoadDate ASC ROWS BETWEEN 24 PRECEDING AND 1 PRECEDING) AS MinRecords, LoadDate
24                  FROM (
25                      SELECT DISTINCT CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
26                                      ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
27                                      END AS LoadDate
28                                      COUNT(*) AS NumRecords
29                  FROM st[____]
30                  WHERE DATENAME(weekday, CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
31                                      ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
32                                      END
33                              ) <> 'Monday'
34                  AND dbo.fn_IsHoliday(DATEADD(DAY,-1, CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
35                                      ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
36                                      END)) <> 1
37
38                  GROUP BY SOURCEFILE
39                  ) b
40          ) a
41      WHERE LoadDate = (SELECT MAX(CASE WHEN TRY_CAST( LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE) IS NULL THEN TRY_CAST(REPLACE(LEFT(RIGHT(SOURCEFILE, 19), 8), '\', '') AS DATE)
42                                      ELSE TRY_CAST(LEFT(RIGHT(SOURCEFILE, 16), 8) AS DATE)
43                                      END) FROM stg.[_____]--CAST(DATEADD(DAY, -1, CURRENT_TIMESTAMP) AS DATE))
44      )
```
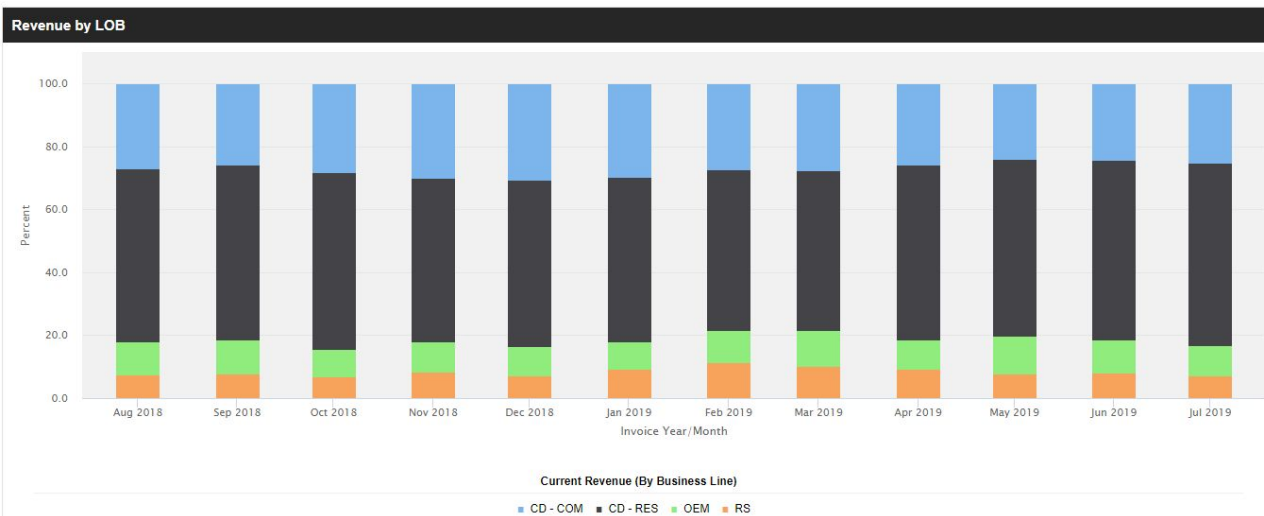
# Custom Procedures – [          ]Quantity

- Goal: Check that the Quantity for each branch or warehouse is above the minimum for the last 30 days
- If the given day is a holiday it doesn't count that as being part of the count.
  - On Pass – Log pass and continue
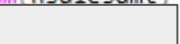  - On Warning – Log warning and continue

```sql
10  ALTER PROCEDURE [etl].[          ] @DatabaseName VARCHAR(255), @LoadName VARCHAR(255), @LoadDate DATETIME, @HeaderID INT, @ImportID VARCHAR(10)
11  AS
12
13  DECLARE @Minimum INT;
14  DECLARE @Result INT;
15  SET @Result = (SELECT SUM([DRIVE QUANTITY]) FROM stg[          ] WHERE [INVOICE DATE] = CAST('07-15-2019' AS DATE) AND [WAREHOUSE NAME] = 'OTHER')
16
17  SET @Minimum = (SELECT MinQuantity FROM
18      (
19          SELECT [INVOICE DATE], [WAREHOUSE NAME], MIN(Quantity) OVER (PARTITION BY [WAREHOUSE NAME] ORDER BY [INVOICE DATE] ASC ROWS BETWEEN 31 PRECEDING AND 1 PRECEDING) AS MinQuantity
20          FROM (
21              SELECT DISTINCT [INVOICE DATE]
22                          ,[WAREHOUSE NAME]
23                          ,SUM([DRIVE QUANTITY]) as Quantity
24          FROM stg[          ]
25          WHERE DATENAME(weekday, [INVOICE DATE]) <> 'Sunday'
26              AND dbo.fn_IsHoliday([INVOICE DATE]) <> 1
27          GROUP BY [INVOICE DATE], [WAREHOUSE NAME]
28          ) b
29      ) a
30  WHERE [INVOICE DATE] = CAST('07-15-2019' AS DATE) AND [WAREHOUSE NAME] = 'OTHER')
31
```

# Custom Procedures -

- Goal: Check Percent of Revenue by Business Line
- The percent of revenue should be pretty consistent across all months for the last year
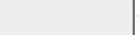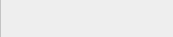
- Goal: Check if the source AR is present
    - If more than 5% is AR – Fail and exit load
    - Else warn and continue



Revenue by LOB

Current Revenue (By Business Line)

■ CD - COM    ■ CD - RES    ■ OEM    ■ RS

```sql
 8  ALTER PROCEDURE [etl].          @DatabaseName VARCHAR(255), @LoadName
 9      AS
10
11      DECLARE @RESULT NUMERIC(18,8)
12      SET @RESULT = NULLIF((
13          SELECT sum(nsalesamt)
14          FROM etl.
15          WHERE MappedDate >= DATEADD(WEEK, -1, CURRENT_TIMESTAMP)
16          AND SOURCE = 'AR'
17          ), 0)
18
19      DECLARE @REVENUE NUMERIC(18,8) = NULLIF((
20          SELECT sum(nsalesamt)
21          FROM etl.
22          WHERE MappedDate >= DATEADD(WEEK, -1, CURRENT_TIMESTAMP)
23          ), 0)
24
```

# Custom Procedures –

- Goal: Check that the Unit of Measure recorded in the invoice for a customer/product combination is correct

```sql
11
12  SELECT A.*
13      , B.UnitPriceBase
14      ,(B.UnitPriceBase - A.UnitPriceCurrent)/NULLIF(B.UnitPriceBase, 0) * 100 AS PctDifference
15      , CASE WHEN (A.UnitPriceCurrent - B.UnitPriceBase)/NULLIF(A.UnitPriceCurrent, 0) > .25 THEN 'FLAG' ELSE NULL END AS FLAG
16  INTO #Temp
17  FROM
18          (
19              select [InvoiceDate]
20                  ,[Customer Number]
21                  , [Item Number SKU_New]
22                  , [Unit of Measure]
23                  , (sum([Price Measurement Revenue])/nullif(sum([Transaction Quantity]), 0)) as UnitPriceCurrent
24          from raw.          where InvoiceDate = DATEADD(day, -1, '2019-07-24')
25          group by [Customer Number], [Item Number SKU_New], [Unit of Measure], InvoiceDate
26          ) A
27  JOIN (
28          select
29                  [Customer Number]
30                  , [Item Number SKU_New]
31                  , [Unit of Measure]
32                  , (sum([Price Measurement Revenue])/nullif(sum([Transaction Quantity]), 0)) as UnitPriceBase
33          from ra          where InvoiceDate BETWEEN DATEADD(YEAR, -1, '2019-07-24') AND DATEADD(day, -2, '2019-07-24')
34          group by [Customer Number], [Item Number SKU_New], [Unit of Measure]
35          ) B
36      ON A.[Customer Number] = B.[Customer Number] AND A.[Item Number SKU_New] = B.[Item Number SKU_New] AND A.[Unit of Measure] = B.[Unit of Measure]
37  WHERE CASE WHEN ABS((A.UnitPriceCurrent - B.UnitPriceBase)/NULLIF(A.UnitPriceCurrent,0)) > .25 THEN 'FLAG' ELSE NULL END IS NOT NULL
38
```