# Data Basics: Defining, Sourcing and Harmonizing Data

Masterclass Series (Class 1)
February 17th, 2021

*https://events.ethisphere.com/events/data-analytics-masterclass-series/*

Lextegrity     ETHISPHERE
GOOD. SMART. BUSINESS. PROFIT.

# ETHISPHERE®
## GOOD. SMART. BUSINESS. PROFIT.®

# Lextegrity

## Thank you for joining! Before we get started…

### Q&A
Please submit your questions using the Questions feature in your Zoom Experience.

### CHAT
Need assistance? The Chat feature will be open throughout the webcast.

### RECORDING & PPT
Today's presentation, recording and additional resources will be provided via email after the webcast.

### DON'T FORGET
Join us for the extended masterclass series!
**March 3rd** | Using Data: Generating Compliance
**March 17th** | Practical Data Analytics: Deep Dives Into Specific Use Cases

# Speaker



**Andy Miller**

Chief Analytics Officer @ Lextegrity

*Helping organizations mitigate risk across the enterprise, specifically through employing advanced analytics to expose fraud, bribery, corruption, conflicts of interest and sanctions violations within enterprise data.*

Contact me @ amiller@lextegrity.com

Learn more @ https://www.lextegrity.com/

ETHISPHERE
GOOD. SMART. BUSINESS. PROFIT.

Lextegrity

# The Lextegrity & Ethisphere Partnership

- Knowledge & Best Practices Sharing
  https://ethisphere.com/what-we-do/lextegrity/

- Data Analytics Working Group

- Integrity Analytics Collective
  https://www.lextegrity.com/collective

ETHISPHERE
GOOD. SMART. BUSINESS. PROFIT.

Lextegrity

# Lextegrity's Mission

**Data-Driven + Intuitive + Integrated**

## FIRST GENERATION



Third-Party Due Diligence — Gifts & Conflicts — ERP, T&E, P2P, HR — Internal Audit

Process, Data & Functional Silos

Lack of Integration

Periodic Audits

Weak Controls

Sampling

Manual

Basic Analytics

User Frustration

## INTEGRITY GATEWAY



Cross-Risk Pre-Approval Workflows

Connecting Approval Information to Spend

Integration with Financial Systems

CONTINUOUS RISK FEEDBACK

Spend Filtered Through Risk Algorithms

Intuitive, Multilingual, Desktop & Mobile

Continuous Monitoring on All Spend

ETHISPHERE
GOOD. SMART. BUSINESS. PROFIT.

Lextegrity

# Agenda

- Introductions
- **Raw Data to Information to Insights**
- Q&A

Lextegrity

ETHISPHERE
GOOD. SMART. BUSINESS. PROFIT.

# Raw Data to Information to Insights

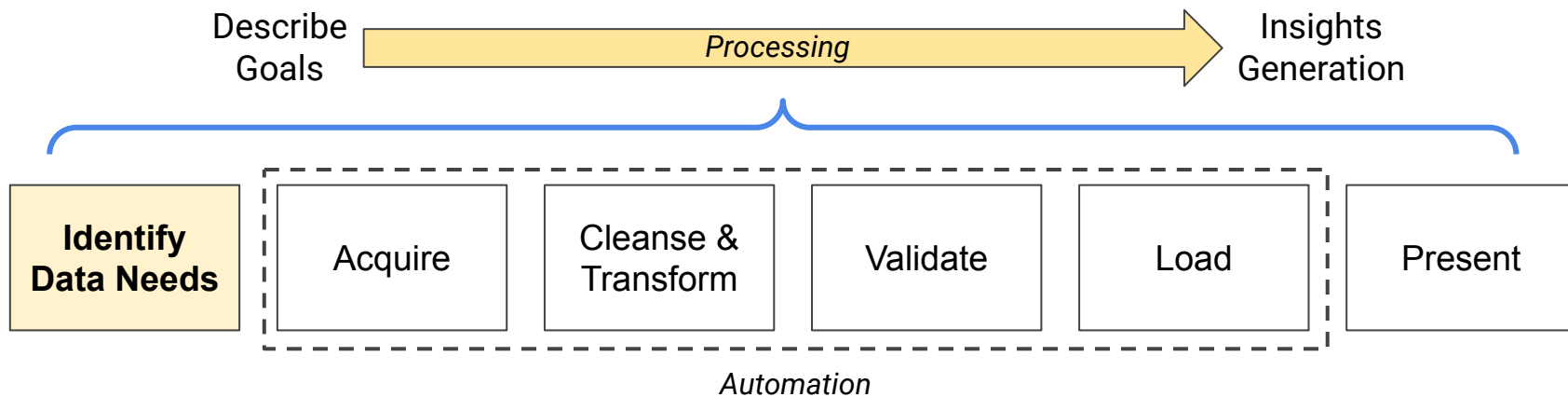**Describe Goals** → *Processing* → **Insights Generation**

- "Having a pulse on organizational spend"

- "Preventing & detecting fraud, bribery, corruption, conflicts of interest, and sanctions violations"

- "Monitoring higher-risk spend / vendors / employees"

- "Prioritizing spend items for additional review when exhibiting higher-risk attributes"

# Raw Data to Information to Insights

**Describe Goals** — Processing → **Insights Generation**

| Identify Data Needs | Acquire | Cleanse & Transform | Validate | Load | Present |
|---|---|---|---|---|---|

*Automation*

# The 5 V's of Data

**Volume**
Scale of Data

**Variety**
Different Forms
of Data

**Value**
Potential of Data

**Velocity**
Fast Moving Data

**Veracity**
Uncertainty of Data

# Relevant Data for Compliance

**Transactional**
- ERP (e.g. Oracle, SAP)
- T&E Systems (e.g. Concur, Expensify)
- Third Party Due Diligence

**Master Data**
- HR Systems (e.g. Workday, SuccessFactors)
- Master Data Sources (e.g. vendor/customer lists, HCPs, etc.)

**Supportive**
- Risk Program Data (e.g. investigations, training statistics, hotline calls, audit exceptions)
- Aggregated Metrics (e.g., budget to actual financial variances, HR turnover metrics)

**Open Source / Publicly Available**
- Industry Competitive Sources / Benchmarks
- Transparency Data (Life Sciences)
- Corruption Perception Index (CPI)

| | Volume | Velocity | Variety | Veracity | Value |
|---|---|---|---|---|---|
| Transactional | High | High | Low | Medium | Medium |
| Master Data | Medium | Low | Low | Medium | Low |
| Supportive | Low | Low | High | High | High |
| Open Source / Publicly Available | High | Low | High | Low | Medium |

*Likeness*  ● High  ● Medium  ○ Low

# Data Accessibility



| Example Data | | | | |
|---|---|---|---|---|
| Name | Date Modified | Size | Kind | |
| compresed_data.tdsx | Jan 30, 2021 at 5:53 AM | 6.8 MB | Tableau...Source | |
| concur_expenses.xlsx | Jan 29, 2021 at 9:23 AM | 47 KB | Spreadsheet | |
| CPI2020_GlobalTablesTS_210125.xlsx | Jan 25, 2021 at 3:41 PM | 82 KB | Spreadsheet | |
| CPI2020_SignificantChanges_210125.xlsx | Jan 25, 2021 at 3:41 PM | 42 KB | Spreadsheet | |
| critical process details.log | Jan 31, 2021 at 3:07 PM | 8 KB | Log File | |
| Entity Hierarchies.json | Today at 8:46 AM | 262 KB | JSON | |
| Extremely Important Receipt.png | Jan 19, 2021 at 2:03 PM | 40 KB | PNG image | |
| Important Contract.docx | Yesterday at 4:28 PM | 81 KB | Microso...(.docx) | |
| Important x2 Contract.pdf | Feb 8, 2021 at 2:48 PM | 148 KB | PDF Document | |
| sap_payments.xlsx | Jan 29, 2021 at 9:23 AM | 199 KB | Spreadsheet | |
| transactional_data.csv | Jan 29, 2021 at 9:23 AM | 327 KB | comma...d values | |
| very_important_database.sql | Nov 25, 2020 at 1:54 PM | 110.58 GB | SQL File | |

*Flat Files*



*Application Programming Interface*



*Databases*



*Data Lake*

# Data Formats: Structured

| | id | key | value | code2 | region | cpi_score ▾ 1 |
|---|---|---|---|---|---|---|
| 1 | 6d6b6bb7-0883-4444-a80b-b8e6875d5256 | DK | Denmark | DK | Europe | 88 |
| 2 | 7396158d-dcff-4c1c-98e9-2ea7fbb4cd4c | NZ | New Zealand | NZ | Western Pacific | 87 |
| 3 | b82b74fe-1a7f-4a28-92c0-3e9d42af7359 | CH | Switzerland | CH | Europe | 85 |
| 4 | 2d807e73-ba03-49cb-a4b9-e4e2e446c87c | SG | Singapore | SG | Western Pacific | 85 |
| 5 | 347bf5cb-22e7-4d7c-a02e-d3b9e7ad84d9 | SE | Sweden | SE | Europe | 85 |
| 6 | 099ac6ab-742f-4d0f-8abe-38f382a517ec | FI | Finland | FI | Europe | 85 |
| 7 | 551d41f2-15b1-4c3a-97e8-0241d1047d1d | NO | Norway | NO | Europe | 84 |
| 8 | 1c7c0205-86ed-4d8e-ab70-eacf5229e573 | NL | Netherlands | NL | Europe | 82 |
| 9 | 43d9d6fa-3501-47a4-b6fe-819cbd6bb116 | CA | Canada | CA | Americas | 81 |
| 10 | 83ab8d9e-8356-4913-9304-28e6591095a4 | LU | Luxembourg | LU | Europe | 81 |
| 11 | ab81ea9e-3619-4f04-a9e9-7c96b1d49c4f | DE | Germany | DE | Europe | 80 |
| 12 | 0e2a4fc4-9f1f-4843-a184-5278364a3cfc | GB | United Kingdom | GB | Europe | 80 |
| 13 | d9833297-1155-44df-80e0-d58c17418aaa | AU | Australia | AU | Western Pacific | 77 |
| 14 | 2dd87ef5-b836-46a0-aa9d-bae689480e60 | IS | Iceland | IS | Europe | 76 |
| 15 | bae72bf1-7dba-41cd-b8eb-35130eda90fa | HK | Hong Kong, SAR China | HK | No Region Specified | 76 |
| 16 | 0ce4b142-428e-4b7d-8615-65c6539c6d16 | AT | Austria | AT | Europe | 76 |
| 17 | 213a714a-0179-4f42-baab-e2a6d5c7e9bf | BE | Belgium | BE | Europe | 75 |
| 18 | 0d5809ac-af61-4ebe-8b76-95139413d89a | EE | Estonia | EE | Europe | 73 |
| 19 | 1b6ed6c9-9da4-44c7-9741-1eb9432a8ab6 | JP | Japan | JP | Western Pacific | 73 |
| 20 | 25036070-e159-4f29-8437-e610053b97df | IE | Ireland | IE | Europe | 73 |
| 21 | 10d5257c-b574-4da3-ab09-ed8c21c17c95 | FR | France | FR | Europe | 72 |
| 22 | 804f14f3-1fbd-4afe-b46c-5ef2c631897b | US | United States of America | US | Americas | 71 |
| 23 | 1529cb16-70ff-423b-8a2f-601eb772e207 | AE | United Arab Emirates | AE | Eastern Mediterranean | 70 |
| 24 | 8fc07b50-65ba-466e-adbb-ff8b2bcff2ad | UY | Uruguay | UY | Americas | 70 |
| 25 | 4cf7a44a-a521-41de-a311-960d1814d2f3 | BT | Bhutan | BT | South-East Asia | 68 |
| 26 | df0011d7-3764-4f97-a1a1-c9ce241740bf | BB | Barbados | BB | Americas | 68 |
| 27 | 812f382a-dbfb-4c97-845c-c436e4ee5c34 | CL | Chile | CL | Americas | 67 |
| 28 | 4fae67b2-596f-4ce1-864f-cc07bbac1d09 | SC | Seychelles | SC | Africa | 66 |
| 29 | 9409e504-ea69-41fa-a694-da95c2a24e64 | BS | Bahamas | BS | Americas | 65 |
| 30 | b509dc2e-6d55-486d-9c0c-3b65705c5d52 | PT | Portugal | PT | Europe | 64 |
| 31 | bffec4f2-099c-4d1b-b41c-3dcea60af682 | TW | Taiwan, Republic of China | TW | No Region Specified | 63 |
| 32 | f48fd5b5-4a4d-4195-9af1-82e91ca8607a | BN | Brunei Darussalam | BN | Western Pacific | 63 |
| 33 | 0e22de58-15f3-41fb-9199-70c0c34b9286 | QA | Qatar | QA | Eastern Mediterranean | 62 |

# Data Formats: Semi-structured

*JSON*

```
[
    {
        "id": "ab7131e9-b328-443a-ad4a-7d566dacf251",
        "key": "AD",
        "value": "Andorra",
        "code2": "AD",
        "region": "Europe"
    },
    {
        "id": "1529cb16-70ff-423b-8a2f-601eb772e207",
        "key": "AE",
        "value": "United Arab Emirates",
        "code2": "AE",
        "cpi_score": 70
    },
    {
        "id": "1e7d02b3-69fd-4111-8d9a-cbda88792097",
        "key": "AF",
        "value": "Afghanistan",
        "code2": "AF",
        "languages": ["Pashto", "Dari Persian"],
        "president": "Ashraf Ghani",
        "population": 32225560,
        "gdp": {
            "amount": 72911000000,
            "currency": "USD"
        },
        "cpi_score": 16
    },
    {
        "id": "10b4a9f7-1e7f-4b74-959e-7b2f211dc0e6",
        "key": "AG",
        "value": "Antigua and Barbuda",
        "code2": "AG",
        "region": "Americas"
    },
    {
        "id": "29d14a7b-1203-4d06-854e-83da0655de3a",
        "key": "AI",
        "value": "Anguilla",
        "code2": "AI",
        "region": "Americas"
    },
    {
        "id": "cfc79600-400f-4d3f-8aef-28a36764a187",
        "key": "AL",
        "value": "Albania",
        "code2": "AL",
        "region": "Europe",
        "cpi_score": 36
    }
]
```

*XML*

```
<?xml version="1.0" encoding="UTF-8"?>
<rdData>
    <dtOrders ShipCountry="France" ShipPostalCode="51100" ShipCity="Reims" ShipAddress="59 rue de
        l'Abbaye" ShipName="Vins et alcools Chevalier" Freight="32.38" ShipVia="3" ShippedDate="2010-07-
        16T00:00:00.000" RequiredDate="2010-08-01T00:00:00.000" OrderDate="2010-07-04T00:00:00.000"
        EmployeeID="5" CustomerID="VINET" OrderID="10248"/>
    <dtOrders ShipCountry="Germany" ShipPostalCode="44087" ShipCity="Münster" ShipAddress="Luisenstr.
        48" ShipName="Toms Spezialitäten" Freight="11.61" ShipVia="1" ShippedDate="2010-07-
        10T00:00:00.000" RequiredDate="2010-08-16T00:00:00.000" OrderDate="2010-07-05T00:00:00.000"
        EmployeeID="6" CustomerID="TOMSP" OrderID="10249"/>
    <dtOrders ShipCountry="Brazil" ShipPostalCode="05454-876" ShipCity="Rio de Janeiro" ShipAddress="Rua
        do Paço #67" ShipName="Hanari Carnes" Freight="65.83" ShipVia="2" ShippedDate="2010-07-
        12T00:00:00.000" RequiredDate="2010-08-05T00:00:00.000" OrderDate="2010-07-08T00:00:00.000"
        EmployeeID="4" CustomerID="HANAR" OrderID="10250" ShipRegion="RJ"/>
    <dtOrders ShipCountry="France" ShipPostalCode="69004" ShipCity="Lyon" ShipAddress="2 rue du
        Commerce" ShipName="Victuailles en stock" Freight="41.34" ShipVia="1" ShippedDate="2010-07-
        15T00:00:00.000" RequiredDate="2010-08-05T00:00:00.000" OrderDate="2010-07-08T00:00:00.000"
        EmployeeID="3" CustomerID="VICTE" OrderID="10251"/>
    <dtOrders ShipCountry="Belgium" ShipPostalCode="B-6000" ShipCity="Charleroi" ShipAddress="Boulevard
        Tirou #255" ShipName="Suprêmes délices" Freight="51.3" ShipVia="2" ShippedDate="2010-07-
        11T00:00:00.000" RequiredDate="2010-08-06T00:00:00.000" OrderDate="2010-07-09T00:00:00.000"
        EmployeeID="4" CustomerID="SUPRD" OrderID="10252"/>
</rdData>
```

# Data Formats: Unstructured

Doctor's Clinical Notes



Figure 2. Example of a Clinical Note

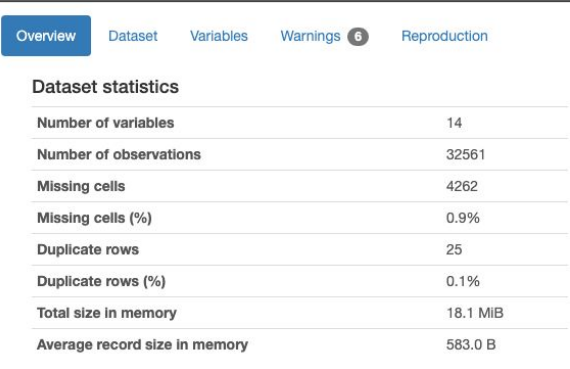## Example of Profiling Report



### Dataset Analysis



**Dataset statistics**

| | |
|---|---|
| Number of variables | 14 |
| Number of observations | 32561 |
| Missing cells | 4262 |
| Missing cells (%) | 0.9% |
| Duplicate rows | 25 |
| Duplicate rows (%) | 0.1% |
| Total size in memory | 18.1 MiB |
| Average record size in memory | 583.0 B |

### Variable Analysis

**age**
Real number ($\mathbb{R}_{\geq 0}$)

| | |
|---|---|
| Distinct | 73 |
| Distinct (%) | 0.2% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Infinite | 0 |
| Infinite (%) | 0.0% |

| | |
|---|---|
| Mean | 38.58164676 |
| Minimum | 17 |
| Maximum | 90 |
| Zeros | 0 |
| Zeros (%) | 0.0% |
| Memory size | 254.5 KiB |

Toggle details

**Statistics** | Histogram | Common values | Extreme values

**Quantile statistics**

| | |
|---|---|
| Minimum | 17 |
| 5-th percentile | 19 |
| Q1 | 28 |
| median | 37 |
| Q3 | 48 |
| 95-th percentile | 63 |
| Maximum | 90 |
| Range | 73 |
| Interquartile range (IQR) | 20 |

**Descriptive statistics**

| | |
|---|---|
| Standard deviation | 13.64043255 |
| Coefficient of variation (CV) | 0.3535471837 |
| Kurtosis | -0.1661274596 |
| Mean | 38.58164676 |
| Median Absolute Deviation (MAD) | 10 |
| Skewness | 0.5587433694 |
| Sum | 1256257 |
| Variance | 186.0614002 |
| Monotocity | Not monotonic |

Reference:
https://github.com/pandas-profiling/pandas-profiling

## Data Transformation

Analyze → Augment & Harmonize (cycle)

### *Augmentation*

High-Risk Subject Logic
- Vendor Type is "Government Related"
- Vendor Title/Profession contains "Government, Govt, Minister, Official, GO, etc."
- Vendor ID matches Third Party Due Diligence review "Risk Level" = "High"
- Vendor pre-approval request government official question is answered yes ("Is the vendor/individual associated with the the government?")
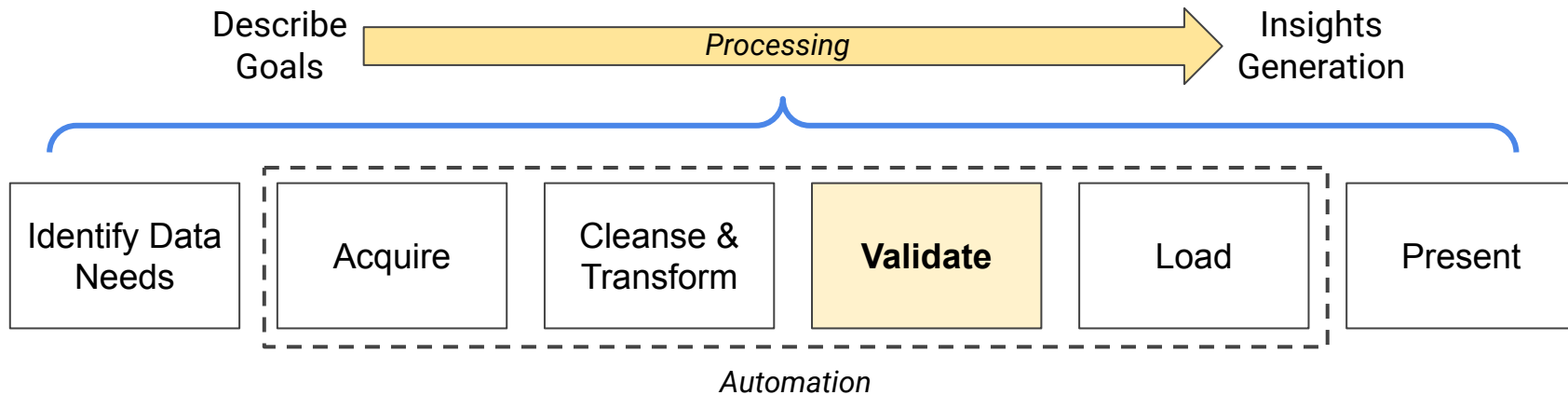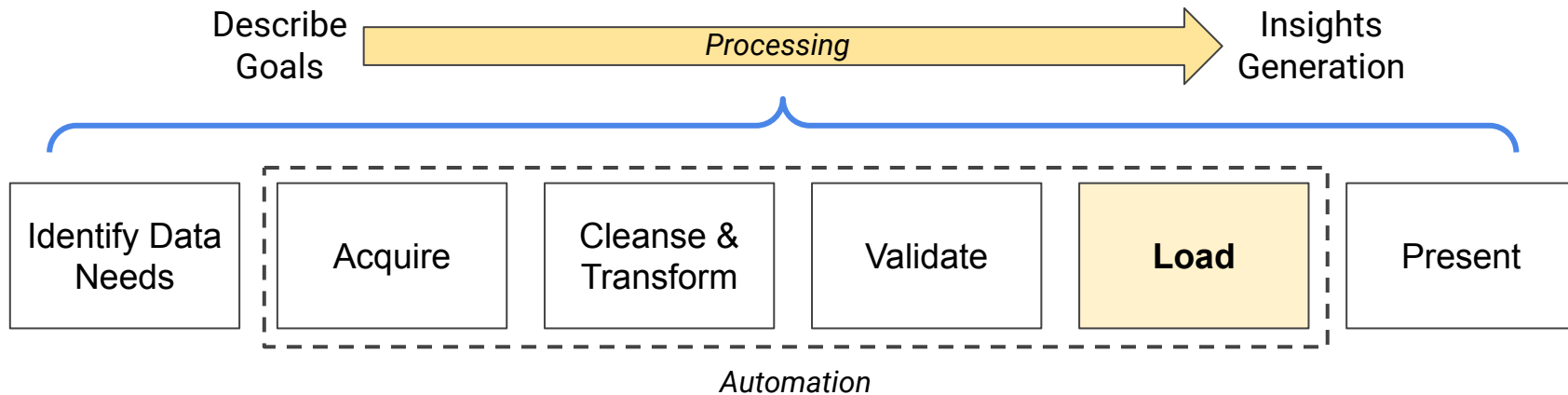
### *Elimination*

**Components**

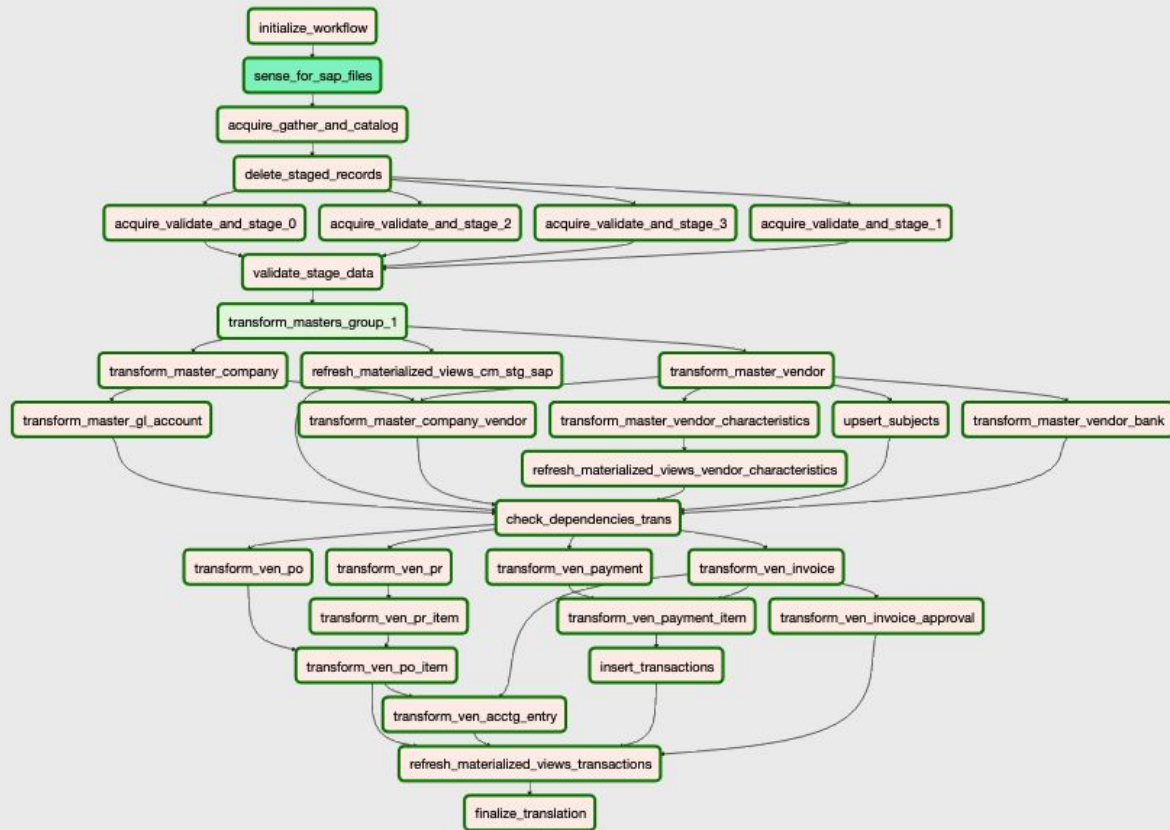| 123 | Field | Key | Data Element | Domain | Data Type | Length | Decimal Places | Short Description |
|---|---|---|---|---|---|---|---|---|
| 1 | MANDT | ☑ | MANDT | MANDT | CLNT | 3 | 0 | Client |
| 2 | LIFNR | ☑ | LIFNR | LIFNR | CHAR | 10 | 0 | Account Number of Vendor or Creditor |
| 3 | ANRED | ☐ | ANRED | TEXT15 | CHAR | 15 | 0 | Title |
| 4 | BAHNS | ☐ | BAHNS | BAHNH | CHAR | 25 | 0 | Train station |
| 5 | BBBNR | ☐ | BBBNR | NUM07 | NUMC | 7 | 0 | International location number (part 1) |
| 6 | BBSNR | ☐ | BBSNR | NUM05 | NUMC | 5 | 0 | International location number (Part 2) |
| 7 | BEGRU | ☐ | BRGRU | BRGRU | CHAR | 4 | 0 | Authorization Group |
| 8 | BRSCH | ☐ | BRSCH | BRSCH | CHAR | 4 | 0 | Industry key |
| 9 | BUBKZ | ☐ | BUBKZ | NUM01 | NUMC | 1 | 0 | Check digit for the international location number |
| 10 | DATLT | ☐ | DATLT | TEXT14 | CHAR | 14 | 0 | Data communication line no. |
| 11 | DTAMS | ☐ | DTAMS | DTAMS | CHAR | 1 | 0 | Report key for data medium exchange |
| 12 | DTAWS | ☐ | DTAWS | DTAWS | CHAR | 2 | 0 | Instruction key for data medium exchange |
| 13 | ERDAT | ☐ | ERDAT_RF | DATUM | DATS | 8 | 0 | Date on which the Record Was Created |

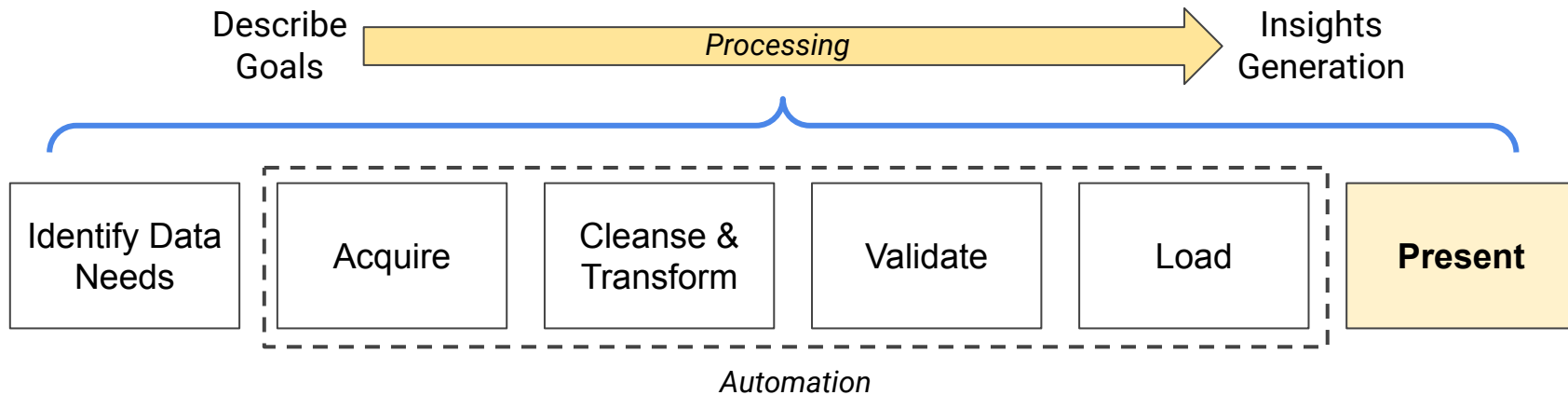*SAP LFA1 table components (Vendor Details)*

### *Harmonization*

- Airfare - Overseas
- Airfare - Domestic
- Airfare - First Class (EU)
- Airfare - Business Class (EU)
- Non-Employee Airfare
- Travel - Airfare Costs
- Airfare > $5K

} Airfare

Describe Goals → Processing → Insights Generation

| Identify Data Needs | Acquire | Cleanse & Transform | **Validate** | Load | Present |

Automation

Describe Goals → *Processing* → Insights Generation

| Identify Data Needs | Acquire | Cleanse & Transform | Validate | **Load** | Present |

*Automation*

Describe Goals → *Processing* → Insights Generation

Identify Data Needs | Acquire | Cleanse & Transform | Validate | Load | **Present**

*Automation*

*Summary Data*

*Detailed Data*

*Aggregation*

# Questions?