# What is YAML?

References :

## 1. What is YAML

From the article (FreeCodeCamp: *What is YAML? The YML File Format*):

- YAML = *YAML Ain't Markup Language* (earlier "Yet Another Markup Language")

- It is a data-serialization language used commonly for configuration files.

- It's human-readable, using indentation instead of lots of brackets/tags.

- Supports collections (sequences, mappings), scalar values, explicit types.

- Supports multi-document files (multiple "documents" inside one .yaml or .yml file separated by ---).

## 2. Advantages of YAML vs JSON or XML

| Feature | YAML | JSON | XML |
|---|---|---|---|
| Readability for Humans | Very good: uses indentation, less punctuation. Easier to glance through. | Good, but more braces/brackets. | Verbose with tags — harder to read. |
| Conciseness | More concise — less boilerplate (no closing tags, fewer delimiters). | More compact than XML, but more punctuation than YAML. | Least concise. |
| Support for Complex Structures | Supports mappings, sequences, multi-documents, explicit data types, anchors & references, etc. | Supports objects, arrays; but no built-in anchors/references. | Very expressive, schema support. Good for validating complex configuration. |
| Schema / Validation | Less formal by default (though schema/validators exist). Some risk of indentation errors. | Many tools, schema definitions (JSON Schema). | Very strong schema support (XSD etc.). |
| Ease of writing & editing configs | Easier for humans to write, maintain. | Relatively easy. | More error prone, more verbose to maintain. |
| Whitespace / Formatting sensitivity | Indentation critical; spaces vs tabs matter. Can make mis-formatting errors. | Less reliant on indentation (structure via braces). | Tags define structure; indentation is purely for readability in many parsers. |

So in many settings, YAML is preferred for configuration files or places where humans will be editing configs frequently.

## 3. How Matillion (or similar ETL tools) could use/generate YAML files

While I couldn't find exactly from that article how Matillion generates YAMLs, in general ETL / data-orchestration tools use YAML for:

- **Configuration files**: defining environments, connections, credentials, parameters (e.g., define Snowflake connection, S3 bucket paths, schema mappings).

- **Job definitions / workflows**: some tools allow exporting job/workflow definitions in YAML (or similar) so that you can version control them, share them, or re-deploy.

- **Parameterization**: You might have files specifying variables, thresholds, alerting rules etc. in YAML.

- **Metadata and schema definitions**: describing table schema, column metadata, types, mappings etc. might be stored as YAML so that jobs can dynamically load schema.

If Matillion supports exporting of workflows into YAML (or some representation in YAML), it probably includes:

- Component names (e.g., "Fixed Flow", "Excel Query", "Join")

- Metadata: input source (S3, files, tables), destination (Snowflake table)

- Parameter definitions (e.g. variable names, sheet names, iterators)

- Transformations and dependencies

- Possibly credentials or references to stored environment objects

---

**4. Business Use Case: ETL + YAML + Matillion (or similar)**

Here's a detailed business scenario that shows when and how YAML files might be used in an ETL pipeline (in Matillion or similar), and how they give value.

---

**Use Case: Retail Analytics Platform**

**Scenario:**

A retail company wants to build a nightly pipeline that:

1. Picks up daily sales files (Excel workbooks) from various regional stores. Each workbook has multiple sheets: *Orders*, *Order_Items*, *Customers*.

2. Some customer information arrives in nested JSON format (e.g., preferences, addresses).

3. The pipeline should dynamically process any workbook regardless of how many sheets (within a known set), generate appropriate staging tables, flatten nested JSON, convert data types, and then join all data to produce KPIs like *Total Revenue*, *Total Profit*, *Revenue by Month & Store*, *Customer Segmentation*.

4. They want to version control the pipeline's configuration (e.g., which sheets to expect, schema of each sheet, thresholds, etc.), so changes in source files or schema can be handled without manually editing DAGs each time.

5. They also want to deploy the pipeline across multiple environments (Dev / QA / Prod), so connection details and parameters should be externalized.

## How YAML Can Be Used Here

| Purpose | What YAML Contains | Benefit |
|---|---|---|
| **Configuration of source files** | e.g. which S3 path, naming conventions of Excel, sheet names, file-patterns etc. YAML file: `sources.yaml` with mapping region → file location, expected sheets, schema of each sheet. | Can update when a new region is added or a sheet is removed, without changing the Matillion job logic. Better maintainability. |
| **Schema mappings** | Define what columns each sheet has, data types, nested JSON schema (e.g. in `customer_info`), which columns to flatten or drop. YAML file: `schema_customers.yaml`, `schema_orders.yaml`. | If a column is added or removed, update schema file; can feed into job to auto-generate staging tables. Avoids manual component editing. |
| **Workflow definitions / Dependencies** | Define in YAML the sequence of components (Excel Query → Fixed Flow → Iterator → Create Table → etc.), variable names, dependencies, retry policies. | Reuse across environments, version control the workflow, potentially auto-generate the orchestration DAG from YAML. |
| **Parameter files for environments** | Connection strings for S3, Snowflake, Matillion project variables etc. YAML per environment. | Security (don't hardcode), easier dev/test/prod promotion, consistent setups. |
| **Alerts / validation rules** | Define conditions e.g. if total revenue drop > X% vs previous day, or missing sheets, the pipeline fails. YAML file for rules. | Easier to manage business rules, update them as business changes. |

## How it would integrate in a Matillion ETL Pipeline

Here's a step-by-step flow of how YAML could be integrated:

1. **YAML Files are stored in a repository** (e.g. Git):
   o environments/prod.yaml
   o schemas/orders.yaml, schemas/order_items.yaml, schemas/customers.yaml
   o sources/regions.yaml
2. **Matillion Job/Orchestration** reads configuration from the YAML files:
   o For example, a component or a Python/Script component in Matillion loads sources/regions.yaml to see which S3 path and sheet names to process.
   o Schema YAML tells what columns to expect in each sheet, data type conversions needed.
3. **Dynamic Component Construction**:

- Use "Fixed Flow" or "Table Iterator" to dynamically create staging tables in Snowflake per sheet as per schema.
- Use schema definition to drive "Create Table", "Convert Types" components automatically.

4. **Transformation Steps**:
   - Extract nested JSON from customer_info, flatten using definitions in schema YAML.
   - Perform joins, aggregations as per business-defined KPIs.

5. **Validation / Checkpoints**:
   - After processing, perhaps a YAML file of rules defines acceptable thresholds (e.g. "Revenue cannot be null", "Number of records per sheet >= 1", etc.). If validation fails, alert or rollback.

6. **Deployment / Promotion**:
   - Same Matillion job logic, but environment YAML (credentials, bucket paths etc.) differ between Dev / QA / Prod.

7. **Export / Version Control**:
   - Matillion or other tools may allow exporting the pipeline's metadata (component graph, parameter settings) into YAML so that the job definition is version controlled, compared, diffed.

---

**Advantages in this Use Case**

- **Flexibility / Maintainability**: If source schema changes (new column, removed sheet), update YAML rather than editing many components manually.
- **Consistency across Environments**: Use same logic, change parameters via YAML.
- **Auditability / Change Tracking**: YAML in Git gives history of config changes.
- **Speed of Onboarding / Repeatability**: New region or file type can be onboarded by defining new YAML, minimal change to pipeline logic.
- **Separation of Logic and Data/Schema**: Business rules, schema definitions, source locations are external to pipeline logic; logic focuses only on execution.

---

**5. Potential Challenges / Considerations**

- YAML's strictness: indentation or syntax errors can cause failures. Needs validation.
- Need tooling or scripting inside Matillion (or via external process) to read YAML and feed pipeline components. Not all ETL tools have built-in YAML schema ingestion, so may need custom scripting or wrapper.
- Security: environment credentials in YAML must be protected (secrets management).

**What is YAML?**



- Complexity: over-abstracting everything into YAML may make troubleshooting harder if not well documented.

**How Matillion Uses / Exports YAML**

From Matillion docs:

1. **Pipeline Export in YAML format**

   o In Matillion Data Productivity Cloud (DPC), when you export pipelines (single pipeline, folder, or entire project), it creates a **zipped file** containing all pipeline definitions in **.yaml** format. docs.matillion.com

   o The export contains component definitions, component properties, connections between components (the workflow graph), canvas notes, pipeline variables, etc. docs.matillion.com

   o It does *not* include project-level settings such as schedules, secrets, OAuth credentials, or project variables. Those you need to configure separately in the target environment. docs.matillion.com

2. **Shared Pipelines config file (shared-pipelines.yaml)**

   o When pipelines are shared, a file called shared-pipelines.yaml is managed in the .matillion folder in the Git repo. docs.matillion.com

   o That file tracks which pipelines are shared, their IDs, their public variables (parameters), variable defaults, etc. docs.matillion.com

3. **Importing Pipelines**

   o The YAML files exported (pipeline definitions) can be imported into other projects (as.zip which contains the .yaml files) in Matillion DPC. exchange.matillion.com

   o This allows pipeline reuse, version control, sharing between projects, etc. docs.matillion.com+2docs.matillion.com+2

---

**Sample Structure / What a Pipeline YAML Might Contain**

Here's a typical structure of what Matillion might have in a pipeline YAML export (based on docs and examples):

# What is YAML?

```yaml
# Example: MyPipeline.orch.yaml

pipeline: MyPipeline            # name
type: orchestration            # or transformation
version: "1.0"

variables:
  - name: var_sheet_names
    type: grid / list / scalar
    default: ...
    description: "Names of sheets to process"

components:
  - id: FixedFlow_1
    type: FixedFlow
    properties:
      property_a: ...
      property_b: ...
      expression: "..."
    outputs:
      - id: TableIterator_1

  - id: TableIterator_1
    type: TableIterator
    properties:
      iterator_list: "${var_sheet_names}"
      child_component: ExcelQuery_1
```

```yaml
# ... more components: S3 Load, Extract Nested Data, Convert Type, Joins, Aggregate, Rewrite etc

connections:
  - from: FixedFlow_1
    to: TableIterator_1
  - from: TableIterator_1
    to: ExcelQuery_1
  - from: ExcelQuery_1
    to: CreateTable_Items
  # etc

notes:
  canvas_notes: |
    This pipeline dynamically loads sheets from Excel,
    flattens JSON, runs joins, aggregates profit/revenue per account.
```

Also, the shared-pipelines.yaml might look something like:

```yaml
version: "1.0"
type: shared-pipelines-config
pipelines:
  - pipeline: folder1/SharedPipeline.orch
    displayName: Shared Pipeline
    id: shared_pipeline_id
    enabled: true
    parameters:
      var1:
        type: TEXT
        description: "..."
        scope: COPIED
      var2:
        ...
```

**Practical Business Example with Matillion + YAML**

Here's a real-life scenario:

**Company:** A retail chain that has many stores across regions.

**Need:** Standardize nightly ETL pipelines, enable reuse, maintain versioning, allow quick changes when schemas/source files change.

**Solution using Matillion + YAML:**

- Build your pipeline in Matillion DPC (orchestration + transformation).

- Export the pipeline → YAML definition. Commit this into your organization's Git repository.

- Maintain schema definitions per region or file type in YAML, e.g., schemas/orders.yaml, schemas/customers.yaml.

- Define shared pipelines for repeated logic (data cleaning, flatten nested JSON, etc.), declare them in shared-pipelines.yaml so multiple other pipelines can call them.

- Use variables in your pipelines (referenced in YAML) so sheet names, region names, target table names can be parameterized.

- When a new sheet is added in the Excel sources, just update the YAML (sheet name list) and redeploy/import. The pipeline logic handles it because it uses Table Iterator and Fixed Flow etc., defined in the YAML.

- For deployment across environments (Dev / QA / Prod), each environment has its own YAML config for connection strings, S3 bucket, schema mapping etc.

**Benefits:**

- Consistency: shared logic is in central YAML, pipelines referencing it will behave similarly.

- Faster changes: schema or source changes handled by editing configuration rather than editing many pipelines.

- Version control: YAML in Git, diffing changes, rollback.

- Reusability: shared pipelines.

- Easier collaboration across teams.