# Unit 10: Content Processing Framework
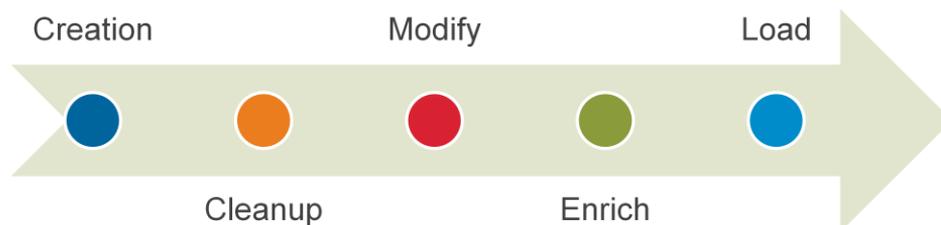
# Learning Objectives

- Define Content Processing

- Describe CPF use cases

- Define the following CPF concepts:

  - Domains

  - Pipelines

  - Functions

  - Pre/post commit triggers

- Debug CPF

- Build and deploy a CPF pipeline

## What is Content Processing?

- A pipeline guiding a document through the stages of its lifecycle in a series of automated steps or tasks
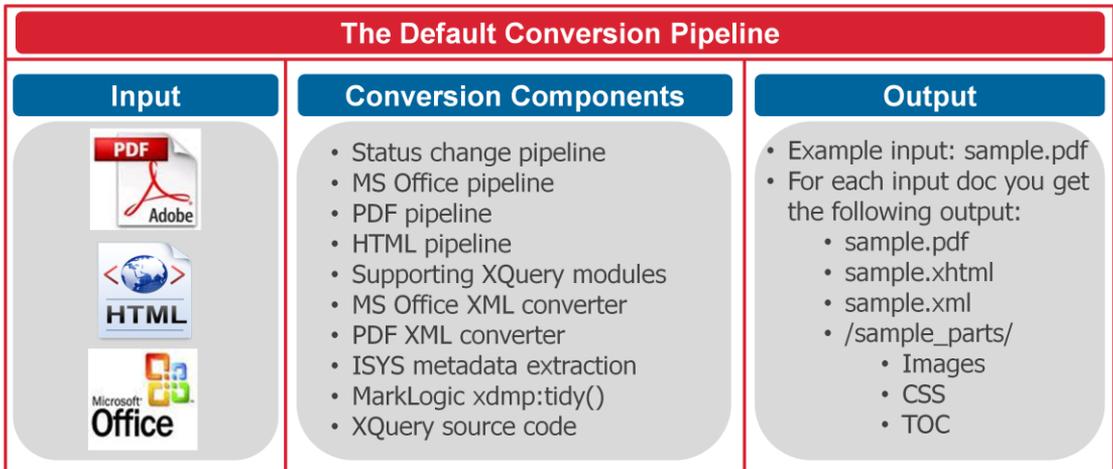- Goal: Make content more useful
  - Example:

| Creation | | Modify | | Load |
|---|---|---|---|---|
| | Cleanup | | Enrich | |

Content can go through many stages before it is ready to use in an application. Sometimes, these
stages might include making the content well-formed XML, some stages might transform one
XML structure to another, and some stages might add (or take away) value to the content by
examining the contents and combining it with other content or information. The process of
content going from one stage to another is called *content processing*.
Every document has a lifecycle. The lifecycle of different documents can vary significantly. A
lifecycle typically begins when a document is created, and then continues through various stages
of content processing; content processing enables a document to move through the phases of its
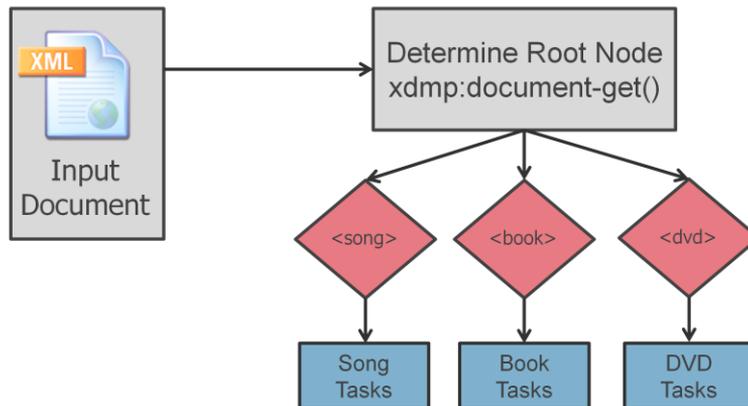lifecycle.

# CPF Uses Cases

- The ability to change content from one form to another
- Source code:
  - /MarkLogic Install Directory/Modules/MarkLogic/conversion

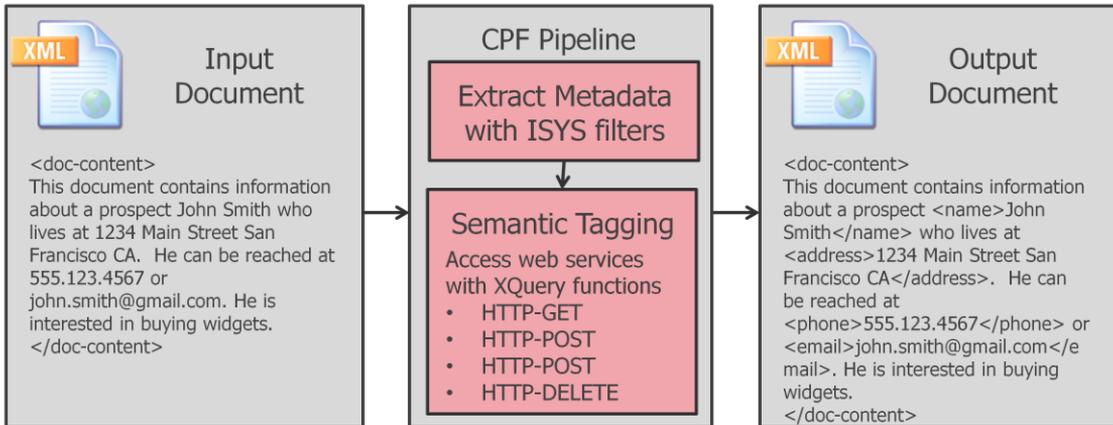| The Default Conversion Pipeline | | |
| --- | --- | --- |
| **Input** | **Conversion Components** | **Output** |
|  | • Status change pipeline<br>• MS Office pipeline<br>• PDF pipeline<br>• HTML pipeline<br>• Supporting XQuery modules<br>• MS Office XML converter<br>• PDF XML converter<br>• ISYS metadata extraction<br>• MarkLogic xdmp:tidy()<br>• XQuery source code | • Example input: sample.pdf<br>• For each input doc you get the following output:<br>  • sample.pdf<br>  • sample.xhtml<br>  • sample.xml<br>  • /sample_parts/<br>    • Images<br>    • CSS<br>    • TOC |

# CPF Uses Cases

- Ability to link together different units of processing as a sequence
- Automation of processing
- Separate documents for different types of processing

# CPF Uses Cases

- Enriching content
  - Semantic tagging | Metadata | Language Conversion
  - Accessing custom built or 3rd party web services

**Input Document**

```
<doc-content>
This document contains information
about a prospect John Smith who
lives at 1234 Main Street San
Francisco CA.  He can be reached at
555.123.4567 or
john.smith@gmail.com. He is
interested in buying widgets.
</doc-content>
```

**CPF Pipeline**

**Extract Metadata with ISYS filters**

**Semantic Tagging**

Access web services with XQuery functions
- HTTP-GET
- HTTP-POST
- HTTP-POST
- HTTP-DELETE

**Output Document**

```
<doc-content>
This document contains information
about a prospect <name>John
Smith</name> who lives at
<address>1234 Main Street San
Francisco CA</address>.  He can
be reached at
<phone>555.123.4567</phone> or
<email>john.smith@gmail.com</e
mail>. He is interested in buying
widgets.
</doc-content>
```

# CPF Concepts: Domains

- A domain defines the scope of documents to be processed.
  - Process some docs this way, other docs that way
- The domain defined in a CPF application answers two fundamental questions:

<table>
<tr>
<td>Which documents will be processed by this application?</td>
<td>Where is the code that makes up this application?</td>
</tr>
</table>

**Mark**Logic®
UNIVERSITY

# CPF Concepts:  Pipelines

- An XML document describing a set of content processing steps
- Defines the **steps** that occur while processing documents and the **actions** that occur under certain **conditions** at each step.
- After each step, the document being processed is committed to the database
- When the commit occurs, the CPF catches the change event with a trigger, which can in turn execute (trigger) more steps.

| Condition: An XQuery module that evaluates to TRUE or FALSE | Action: An XQuery module that is called when the condition linked to it evaluates TRUE |

**MarkLogic®**
UNIVERSITY

# CPF Concepts:  Functions

- MarkLogic provides many useful functions for the CPF

```
xdmp:pdf-convert
xdmp:word-convert
xdmp:excel-convert
xdmp:powerpoint-convert
xdmp:document-filter
xdmp:tidy
xdmp:document-load
xdmp:document-get
xdmp:http-get
xdmp:http-post
xdmp:http-put
xdmp:http-delete
```

**MarkLogic**
UNIVERSITY

10 - 9

# CPF Concepts: Triggers

- Triggers allow you to capture events and then perform some tasks after the event occurs
- Triggers capture document change events for documents under a domain
- Pre-commit triggers execute before the transaction commits
- Post-commit triggers execute after the transaction commits

Example Document Events:
- Create
- Update
- Delete
- Property Change

Example System Events:
- Database Online

# Debugging CPF

- Limit the scope of your testing
    - A single document is a good unit for initial testing
- Check the properties fragment for each document URI

| Document | Root Element | Properties |
|---|---|---|
| /songs/Bananarama.xml | E top-song | (properties) |

top-songs-admin (7000-top-songs-admin-HTTP)  1 Documents

```xml
<?xml version="1.0" encoding="UTF-8"?>
<prop:properties xmlns:prop="http://marklogic.com/xdmp/property">
  <cpf:processing-status xmlns:cpf="http://marklogic.com/cpf">done</cpf:processing-status>
  <cpf:property-hash xmlns:cpf="http://marklogic.com/cpf">d41d8cd98f00b204e9800998ecf8427e</cpf:property-hash>
  <cpf:last-updated xmlns:cpf="http://marklogic.com/cpf">2012-02-24T10:32:36.524-08:00</cpf:last-updated>
  <cpf:state xmlns:cpf="http://marklogic.com/cpf">http://marklogic.com/states/error</cpf:state>
  <cpf:error xmlns:cpf="http://marklogic.com/cpf">
    <error:error xsi:schemaLocation="http://marklogic.com/xdmp/error error.xsd" xmlns:error="http://marklogic.com/x
      <error:code>XDMP-UNEXPECTED</error:code>
      <error:name>err:XPST0003</error:name>
      <error:xquery-version>1.0-ml</error:xquery-version>
      <error:message>Unexpected token</error:message>
      <error:format-string>XDMP-UNEXPECTED: (err:XPST0003) Unexpected token syntax error, unexpected Let_, expectin
      <error:retryable>false</error:retryable>
      <error:expr> </error:expr>
      <error:data>
```

Our genre facet would be more useful if there was better consistency amongst our genre data in our documents.

For example, look at the facet and notice an entry for "pop rock" and also "pop / rock".

Another example, what is "synthpop"? And should it differ from "dance-pop"? Or should both simply be "pop"? We can build a CPF application to define those rules and normalize the <genre> data as the documents are loaded.

# Unit 10: Applying the Learning Objectives

- Define Content Processing
- Describe CPF use cases
- Define the following CPF concepts:
    - Domains
    - Pipelines
    - Functions
    - Pre/post commit triggers
- Debug CPF
- Build and deploy a CPF pipeline
    - Exercise 1 | Exercise 2 | Exercise 3 | Exercise 4 | Exercise 5

**MarkLogic** ®
UNIVERSITY