**DEUTSCHE
NATIONAL
BIBLIOTHEK**

Christoph Böhme

# Analysis of library metadata with Metafacture

---

**DEUTSCHE
NATIONAL
BIBLIOTHEK**

## Agenda

**13:00** — a short introduction to Metafacture

**13:30** — warm-up exercises

**15:00** — triples and counting

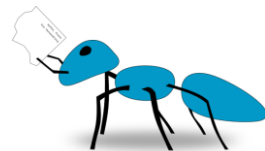**15:30** — exercises on counting data
        (incl. coffee break)

**16:00** — joining data sets and analysing them
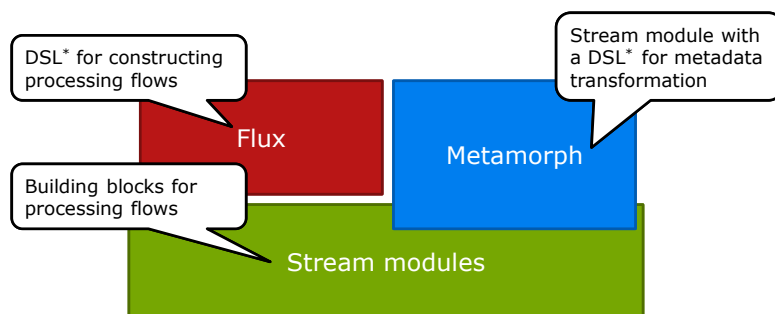
**16:30** — exercises on joining data

**18:40** — wrapping up

# Part 1
## A short introduction to Metafacture
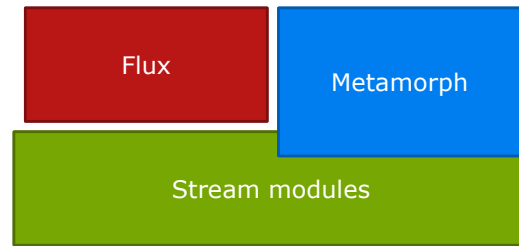
---

## Overview of Metafacture

DSL* for constructing processing flows

Stream module with a DSL* for metadata transformation

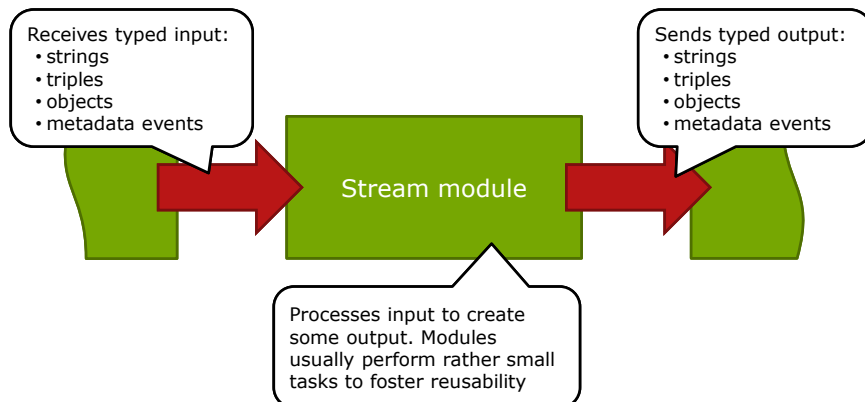Flux

Metamorph

Building blocks for processing flows

Stream modules

*DSL: Domain specific Language

# Overview of Metafacture

# The basic building block of Metafacture



Receives typed input:
• strings
• triples
• objects
• metadata events

Stream module

Sends typed output:
• strings
• triples
• objects
• metadata events

Processes input to create some output. Modules usually perform rather small tasks to foster reusability

# A simple processing flow

Read and print a file containing pica records:



- → file name
  file handle →
- → string
  metadata events →
- → a string
  nothing →

String → `open-file` → `as-lines` → `decode-pica` → `encode-formeta` → `write("stdout")`

- → file handle
  a string for each line →
- → metadata events
  string →

---

# Module configuration



**Module configuration**
- either a single mandatory value
- or optional key-value pairs

Stream module

# Overview of Metafacture

Flux

Metamorph

Stream modules

# Describing flows with Flux

A string as the initial input

Key-value based configuration

Modules are connected with a pipe character

Mandatory parameter

Flow ends with a semi-colon

```
"file.name"
|open-file
|as-lines
|decode-pica
|encode-formeta(style="multiline")
|write("stdout");
```

# Variables and comments in Flux

Define default values for the variables **in** and **out**

```
default in = "file.name";
default out = "stdout";

in
|open-file
// ...
|write(out);
```

Comments start with two slashes

Use variables instead of directly entering a string

---

# Neat features for strings in Flux

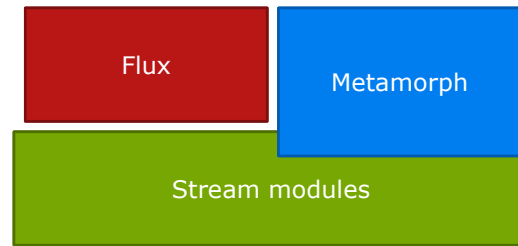**FLUX_DIR** always contains the path of the parent folder of the flux file

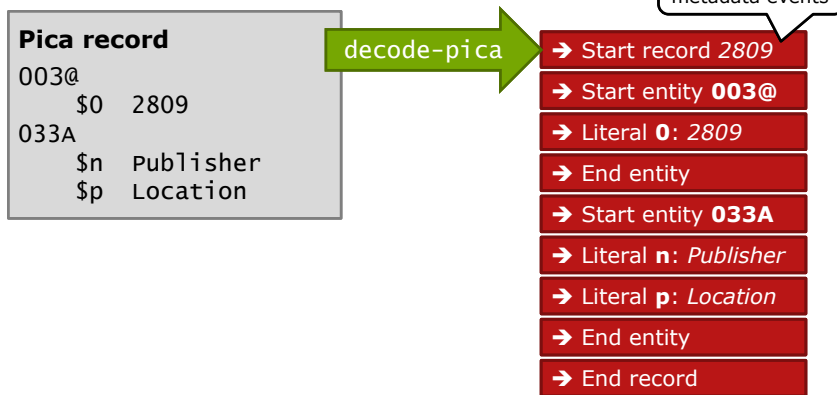Strings can be concatenated

```
default in = FLUX_DIR + "file.name";

default haiku =
"A dropped ice cream cone
One red ant, then two, then four
Suddenly – hundreds!";
```

Strings can be multi-line

# Overview of Metafacture

Flux

Metamorph

Stream modules

---

# Representation of metadata in Meta-facture: a stream of events

**Pica record**
```
003@
    $0   2809
033A
    $n   Publisher
    $p   Location
```

decode-pica

Sequence of metadata events

➜ Start record *2809*

➜ Start entity **003@**

➜ Literal **0**: *2809*

➜ End entity

➜ Start entity **033A**

➜ Literal **n**: *Publisher*

➜ Literal **p**: *Location*

➜ End entity

➜ End record

# Processing metadata events with Metamorph

**morph**

➔ Start record *id*
➔ Start entity **021A**
➔ Literal **a**: *The Trial*
➔ End entity
➔ End record

Listen for 021A.a
⌐ Output as Title

➔ Start record *id*
➔ Literal **Title**:
    *The Trial*
➔ End record

---

# Metamorph: data statements

```xml
<?xml version="1.0" encoding="UTF-8"?>

<metamorph xmlns="http://www.culturegraph.org/metamorph"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    version="1" entityMarker=".">

    <rules>
        <data source="021A.a" name="Title" />
    </rules>

</metamorph>
```

Separator for entities and literal names

Name of the literal to listen for

Name of the literal that is output

# Metamorph: modifying data

```
...

<rules>

   <data source="021A.a" name="Title">
       <regexp match="^(The) (.*)$" format="${2}, ${1}" />
   </data>

</rules>

...
```

Process the data value before outputting it. You can specify multiple functions here

# Metamorph: combining data

Name of the generated literal. It can include variables, too

Literal value constructed from the variables from the data statements below

```
...

<rules>

   <combine name="Publisher" value="${Pub}: ${Loc}">
       <data source="033A.n" name="Pub" />
       <data source="033A.p" name="Loc" />
   </combine>

</rules>

...
```
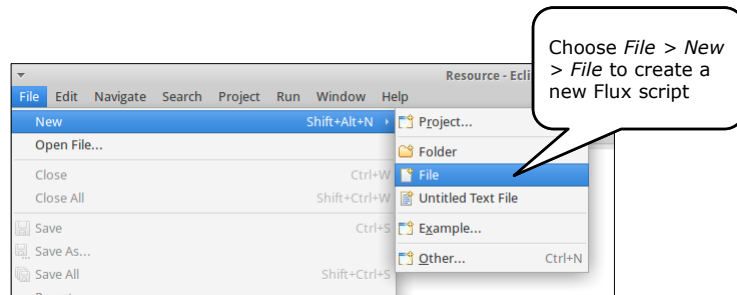
The data statements do not generate output but create variables instead

# Overview of Metafacture

---

# Creating Flux scripts in Eclipse



Choose *File > New > File* to create a new Flux script

# Creating Flux scripts in Eclipse

Select the project folder where you want to save your script

Enter a file name with the extension ".flux"

---

# Running Flux scripts in Eclipse

! The Flux script must be saved and selected

Choose "Run with Flux" to execute the selected Flux script

Choose "Flux Help" to get a list of all supported modules

1 Run with Flux
2 Run with Flux (No Parameters)
3 Flux Help
Run As
External Tools Configurations...
Organize Favorites...

# Running Flux scripts in Eclipse



Pass variables to the Flux script

---

# Running Flux scripts from the prompt

```
C:\> flux.bat myscript.flux
```

```
C:\> flux.bat myscript.flux in=in.file out=out.file
```

Use variables to pass user-defined values into a script

**Running Flux scripts from the prompt**

```
$ ./flux.sh myscript.flux
```

```
$ ./flux.sh myscript.flux in=in.file out=out.file
```

Use variables to pass user-defined values into a script

---

**Exercises: Documentation**

– The hand outs equip you with everything you need for the exercises

– Running "Flux Help" in Eclipse (or flux.sh/flux.bat without additional arguments from the prompt)

– Use auto-completion in Eclipse (press CTRL+Space)

– The Metafacture-Wiki:
https://github.com/culturegraph/metafacture-core/wiki

Exercises part 1
**Warm-up**

Part 2
**Triples and counting**

# The triple

Triple: | *Subject* | **Predicate** | *Object* |

Inspired by RDF triples but subject und predicate do not need to be URIs

---

# Generating triples

Metadata events ➜ stream-to-triples ➜ Triples

Literals on top level

➜ Start record *id*
➜ Literal **name**: *Klaus*
➜ Start entity **died**
➜ Literal **when**: *1401*
➜ Literal **where**: *HH*
➜ End entity
➜ End record

Entities on top level

| *record-id* | **name** | *Klaus* |

| *record-id* | **died** | *…* |

Serialised with Formeta

# Counting triples



count-triples
(countBy="object")

count 4

count 2

count 3

# Outputting triples



*red* **count** *4*

template
("${o} times ${s}")

"4 times red"

Use ${s}, ${p} and
${o} as placeholders for
subject, predicate and
object

# Counting data values

Converts the literals to triples

Optionally: converts triples into formatted text

```
→ morph → stream-to-triples → count-triples → template →
```

Extracts data items that should be counted and outputs them as top level literals

Counts the different object values of the triples

---

# Counting data values: flow of data

➔ Start record *id*
➔ Start entity **033A**
➔ Literal **p**: *Hamburg*
➔ End entity
➔ End record

**Morph**

➔ Start record *id*
➔ Literal **loc**: *Hamburg*
➔ End record

| *id* | **loc** | *Hamburg* |
|------|---------|-----------|

**Count**

| *Hamburg* | **count** | *1* |
|-----------|-----------|-----|

## Metamorph: choosing data

Only the value of the topmost data-statement that generates output is returned by the choose-statement

```
...
<rules>

   <choose name="Location">
     <data source="033A.p">
        <regexp match="^Ffm$" format="Frankfurt a. M." />
     </data>
     <data source="033A.p" />
   </choose>

</rules>
...
```

## Metamorph: generating constant values

```
...
<rules>

   <data source="021A.a" name="Title">
       <constant value="All books have the same name" />
   </data>

</rules>
...
```

No matter what the value of literal 021A.a is, always output the defined value

Exercises part 2
**Triples and counting**

Part 3
**Joining data sets and analysing them**

# Joining streams of data



**How is this done?**

---

# Converting triples into records

Triples ➡ `collect-triples` ➡ Metadata events

| 2 | **PA** | *OA* |
| 2 | **PB** | *OB* |

Sequences of triples with the same subject are merged

| 1 | **PC** | *OC* |
| 1 | **PD** | *OD* |

| 2 | **PE** | *OE* |
| 2 | **PF** | ... |

**record *2***
 **PA**: *OA*
 **PB**: *OB*

**record *1***
 **PC**: *OC*
 **PD**: *OD*

**record *2***
 **PE**: *OE*
 **PF** {
  ...
 }

Serialised entities are deserialised into entities

# Sorting triples

---

# Linking streams in Flux with wormholes

Receives triples from a "wormhole"

```
"file1"              "file2"              @x
|open-file           |open-file           |wait-for-inputs("2")
// ...               // ...               |sort-triples
|stream-to-triples   |stream-to-triples   |collect-triples
|@x;                 |@x;                  |encode-formeta
                                          |write("stdout");
```

Sends the triples into a "wormhole"

These three flows must be defined in the same Flux script

# Advanced triplification: ID redirection

Metadata events → **stream-to-triples (redirect="true")** → Triples

**Replaces record id in subjects**

➜ Start record *id*    ✗

➜ Literal **_id**: new *id*

➜ Literal **name**: *Klaus*

| *new id* | **name** | *Klaus* |
|---|---|---|

...

**Sets subject for individual triples**

➜ Literal **{to:*id*}n**: *v*

➜ End record

| *id* | **n** | *v* |
|---|---|---|

---

# Using _id-redirection

✗

| **id: gnd1** ... | **id: gnd2** ... | **id: gnd3** ... |
|---|---|---|

| **id: HH** gnd: 3 ... | **id: Ffm** gnd: 1 ... | **id: Ks** gnd: 2 ... |
|---|---|---|

➜ Start record *HH*

➜ Literal **gnd**: *3*

➜ Literal **country**: *D*

➜ End record

**Morph** →

➜ Start record *HH*

➜ Literal **_id**: *gnd3*

➜ Literal **country**: *D*

➜ Literal **wiki-id**: *HH*

➜ End record

| *gnd3* | | |
|---|---|---|

| *gnd3* | **country** | *D* |
|---|---|---|

| *gnd3* | **wiki-id** | *HH* |
|---|---|---|

# Using {to:ID}-redirection

Finding the backlinks:
Who links to this record?

| id: Ffm | id: Ks | id: HH |
|---------|--------|--------|
| … | … | … |

| id: gnd1 | id: gnd2 | id: gnd3 |
|----------|----------|----------|
| loc: HH | loc: Ffm | loc: HH |
| … | … | … |

| HH | | |
|----|---|---|

➜ Start record *gnd1*
➜ Literal **loc**: *HH*
➜ End record

Morph

➜ Start record *gnd1*
➜ Literal
   **{to:*HH*}ref**: *gnd1*
➜ End record

| HH | **ref** | *gnd1* |
|----|---------|--------|

---

# Putting the pieces together

Conversion with
id redirection

Joins the two
streams of triples
using a "wormhole"

Ensures that
triples with
the same
subject form
a sequence

morph → stream-to-triples →

@ → wait-for-inputs("2") → sort-triples

morph → stream-to-triples →

Waits for all flows
writing triples into
the "wormhole"

collect-triples

Prepares records to
generate the right
subject values

Turns the triples back
into metadata events

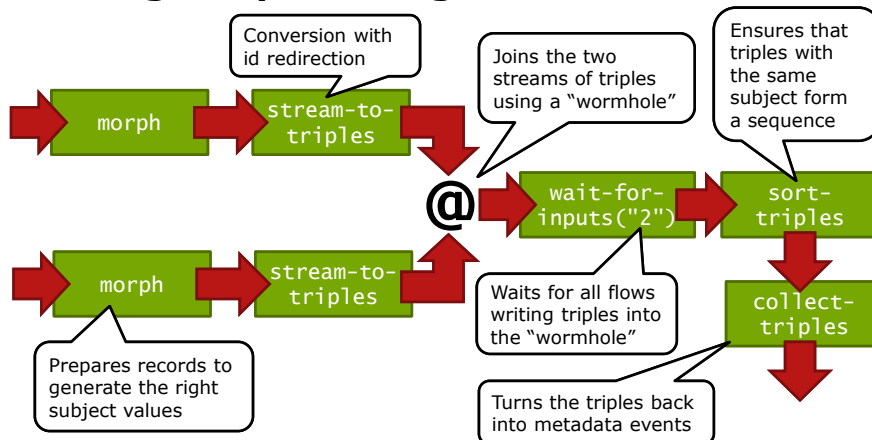# Metamorph: what else?

```
...

<rules>

    <data source="_id" name="recordId" />
    <data source="033A.n" name="Publisher" />
    <data source="_else" />

</rules>

...
```

> Metamorph makes the record id of the current record as **_id** available

> Any literal not handled by any other data-statement is passed to this statement. It can be used to pass data through Metamorph

---

Exercises part 3
## Joining data sets and analysing them

# Wrapping up

# What did we learn today?

– Foundations of processing metadata with Flux and Metamorph

– Exploring data sets by quantifying data values

– Joining data sets and analysing their relations

– Typical patterns for analysing data with Metafacture

  These patterns are similar to the way Hadoop operates: This makes migration from your desktop to a Hadoop cluster easy

## Metafacture

– Not only designed for data analysis but for metadata
   processing in general

– Software tool and library: It can easily be integrated
   into other applications

– Flux and Metamorph are extendable

– It is open source at http://culturegraph.github.io/

## Thank you very much!

**Further questions?**

Contact me at c.boehme@dnb.de
or join the mailing list:
http://lists.dnb.de/mailman/
listinfo/metafacture

**Workshop materials**

Download from http://b3e.net/
metafacture-swib14/workshop.zip