

---

# **DKB Documentation**

**DKB team**

**Dec 06, 2018**



**CONTENTS:**

<b>1</b>	<b>pyDKB package</b>	<b>1</b>
1.1	Quickstart guide . . . . .	1
1.2	Subpackages . . . . .	3
<b>2</b>	<b>Stage 055</b>	<b>23</b>
<b>3</b>	<b>Indices and tables</b>	<b>29</b>
	<b>Python Module Index</b>	<b>31</b>
	<b>Index</b>	<b>33</b>



## PyDKB package

Common library for Data Knowledge Base Dataflow stages development.

**Dataflow** ETL process (extract-transform-load) for populating internal DKB storages and keeping them up to date

**Dataflow stage** Logical step of ETL process, implemented as standalone executable program (worker)

Dataflow stages are standalone programs, but can be combined into a pipeline by means of Kafka-based supervising program. For details about program compatibility with the supervisor please check documentation for the Metadata Integration Topology Management System (MInT MS) workers<sup>1</sup>. Worker program can be written in any language; pyDKB is intended to simplify this process for Python.

**Warning:** There are three types of stages corresponding three types of ETL operations: *source connector* (data extraction), *processor* (transformation) and *sink connector* (load to internal DKB storage). Currently pyDKB library can be used only for *processor* stages, but in future versions *connector* stages will also be supported.

### 1.1 Quickstart guide

To create simple processor stage application first decide input and output data format. In following examples we will work with data in JSON format (for the full list of supported formats check pyDKB.dataflow.messages section of this documentation).

Now let's start writing example processor `welcome.py` and implement message handler – functional part of the stage (operations to be performed on data flow units):

```
from pyDKB.dataflow.messages import JSONMessage

def my_process(stage, message):
    """ Single message processing. """
    input_data = message.content()
    name = input_data.get('name')
    if name:
        out_data = {'message': "Welcome, %s!" % name}
        out_message = JSONMessage(out_data)
        stage.output(out_message)
    return True
```

Function must take two arguments: `stage` (stage context object) and `message` (input message, which should be transformed by our stage). Message is a smallest data unit in the data flow running through the processor, and every message is to be processed independently of previous or following ones. `message.content()` and

---

<sup>1</sup> WIP

`JSONMessage(out_data)` statements are used to decode/encode message to/from Python `dict` object. Message, passed to the function, is taken from the input data flow; to write new message(s) to the output data flow, `stage.output(out_message)` is used. It can be used as many times as many output messages were generated (or once with the list of messages). In our example, messages without key 'name' will produce no output messages, so `stage.output()` will not be called at all. In terms of data flow it means that the input message is filtered out and will not reach the *sink connector*.

Boolean return value of `my_process()` indicates if the processing was successful or not. If processing failed (False is returned), produced output messages will be dropped to avoid loading sketchy information into the DKB storages.

Now as we have processing logic implemented, we need to turn it into fully functional application. Add following lines to `welcome.py`:

```
import sys
from pyDKB.dataflow.stage import JSONProcessorStage
from pyDKB.dataflow.messages import JSONMessage

def my_process(stage, message):
    <...function code...>

if __name__ == '__main__':
    stage = JSONProcessorStage()
    stage.process = my_process
    stage.parse_args(sys.argv[1:])
    stage.run()
```

First we create stage object and indicate that input and output message format is JSON: `stage = JSONProcessorStage()` (for full list of processors check [pyDKB.dataflow.stage package](#) section of this documentation); then set stage processing function to our function `my_process()`, parse command line arguments (`stage.parse_args(sys.argv[1:])`) and start the stage execution.

Easy, right?

It's time to run our application. Create input data sample `input.ndjson` with following lines:

```
{"name": "James", "city": "New York"}
{"user": "Jonathan", "role": "support"}
{"name": "John Smith"}
```

and type:

```
$ python welcome.py --dest s input.ndjson
{"message": "Welcome, James!"}
{"message": "Welcome, John Smith!"}
```

`--dest s` indicates that output destination is (s)tdout (default destination is file). For full information about modes in which the stage application can be used, run `python welcome.py -h`.

That's it, your first application is ready to be integrated into an ETL process as data processing node. For details about ETL process creation check *MInT Supervisor*<sup>2</sup> documentation.

---

<sup>2</sup> WIP

## 1.2 Subpackages

### 1.2.1 pyDKB.common package

Common modules.

#### Submodules

#### pyDKB.common.Type module

Abstract class for type definitions.

##### Example

```
>>> myType = Type("Orange", "Apple")
>>> myType.add("Plum")
>>> t = myType.Orange
>>> if t == myType.Orange:
...     print "Orange!"
... elif t == myType.member("Apple"):
...     print "Apple!"
...
Orange!
>>> if not myType.member("Walnut"):
...     print "Wrong type!"
...
Wrong type!
```

**class** pyDKB.common.Type.Type(\*args)

Bases: object

Abstract class for type definitions.

Member names (*str*) are passed to the constructor as positional arguments.

**add**(name)

Add member.

**Parameters** name (*str*) – name of the member to be added

**hasMember**(val)

Check if the member exists (by value).

**Parameters** val (*int*) – member to be checked

**Returns** True/False

**Return type** bool

**member**(name)

Check if the member exists (by name).

**Parameters** name (*str*) – name to be checked

**Returns** member value or False if member does not exist

**Return type** int, bool

**memberName**(val)

Return string name of the member.

**Parameters** `val` (*int*) – member to retrieve name for

**Returns** member name of False if member does not exist

**Return type** str, bool

### pyDKB.common.custom\_readline module

Implementation of “readline”-like functionality for custom separator.

---

**Todo:** make import of `fcntl` (or of this module) optional to avoid errors when library is used under Windows.

---

`pyDKB.common.custom_readline.custom_readline(f, newline)`

Read lines with custom line separator.

Construct generator with readline-like functionality: with every call of `next()` method it will read data from `f` until the `newline` separator is found; then yields what was read.

**Warning:** the last line can be incomplete, if the input data flow is interrupted in the middle of data writing.

#### Parameters

- `f` (*file*) – readable file object
- `newline` (*str*) – delimiter to be used instead of `\n`

**Returns** iterable object

**Return type** generator

---

#### Todo:

- make last “line” handling more strict: no `newline` == no line;
  - rethink function name (as “line” is actually a “message”);
  - move functionality to `pyDKB.dataflow.communication`<sup>1</sup> submodule)
- 

### pyDKB.common.exceptions module

Definition of common modules exceptions

**exception** `pyDKB.common.exceptions.HDFSException`

Bases: `exceptions.RuntimeError`

Base Exception for HDFS module.

### pyDKB.common.hdfs module

Utils to interact with HDFS.

---

<sup>1</sup> <https://github.com/PanDAWMS/dkb/pull/129>



`pyDKB.common.hdfs.File(fname)`

Get and open temporary local copy of HDFS file

Return value: open file object (TemporaryFile).

`pyDKB.common.hdfs.basename(path)`

Return file name without path.

`pyDKB.common.hdfs.check_stderr(proc, timeout=None, max_lines=1)`

Wait till the end of the subprocess and send its STDERR to STDERR.

Output only MAX\_LINES of the STDERR to the current STDERR; if MAX\_LINES == None, output all the STDERR.

Return value is the subprocess' return code.

`pyDKB.common.hdfs.dirname(path)`

Return dirname without filename.

`pyDKB.common.hdfs.getfile(fname)`

Download file from HDFS.

Return value: file name (without directory)

`pyDKB.common.hdfs.join(path, filename)`

Join path and filename.

`pyDKB.common.hdfs.listdir(dirname, mode='a')`

List files and/or subdirectories of HDFS directory.

**Parameters:** dirname – directory to list mode – 'a': list all objects

'f': list files 'd': list subdirectories

`pyDKB.common.hdfs.makedirs(dirname)`

Try to create directory (with parents).

`pyDKB.common.hdfs.movefile(fname, dest)`

Move local file to HDFS.

`pyDKB.common.hdfs.putfile(fname, dest)`

Upload file to HDFS.

## pyDKB.common.json\_utils module

Utils to work with JSON (dict) objects.

`pyDKB.common.json_utils.nestedKeys(key)`

Transform STRING with nested keys into LIST.

**Parameters:**

**STRING key – dot-separated list of nested keys.** If a key contains dot itself, the key must be put between quotation marks.

`pyDKB.common.json_utils.valueByKey(json_data, key)`

Return value by a chain (list) of nested keys.

**Parameters:** DICT json\_data – to search in STRING key – dot-separated list of nested keys

## 1.2.2 pyDKB.dataflow package

Dataflow organization utils.

### Subpackages

#### pyDKB.dataflow.communication package

pyDKB.dataflow.communication

pyDKB.dataflow.communication.**Message** (*msg\_type*)  
Return class XXXMessage, where XXX is the passed type.

### Subpackages

#### pyDKB.dataflow.communication.consumer package

Consumer submodule init file.

```
class pyDKB.dataflow.communication.consumer.ConsumerBuilder (config={})  
    Bases: object  
  
    Constructor for Consumer instance.  
  
    build (config={})  
        Return constructed consumer.  
  
    consumerClass = None  
  
    setSource (source)  
        Set data source for the consumer.  
  
    setType (Type)  
        Set message type for the consumer.
```

### Submodules

#### pyDKB.dataflow.communication.consumer.Consumer module

pyDKB.dataflow.communication.consumer.Consumer

```
class pyDKB.dataflow.communication.consumer.Consumer.Consumer (config={})  
    Bases: object  
  
    Data consumer implementation.  
  
    close ()  
        Close opened data stream and data source.  
  
    config = None  
  
    get_message ()  
        Get new message from current source.  
  
    Return values: Message object False (failed to parse message) None (all input sources are empty)
```

```

get_source_info()
    Return current source info.

get_stream()
    Get input stream linked to the current source.

    Return value: InputStream None (no sources left to read from)

init_stream()
    Init input stream.

log(message, level=3)
    Output log message with given log level.

message_class()
    Return message class.

message_type = None

next()
    Return new Message, read from input stream.

reconfigure(config={})
    (Re)initialize consumer with stage config arguments.

reset_stream()
    Reset input stream to the current source.

set_message_type(Type)
    Set input message type.

exception pyDKB.dataflow.communication.consumer.Consumer.ConsumerException
    Bases: pyDKB.dataflow.exceptions.DataflowException
    Dataflow Consumer exception.

```

### pyDKB.dataflow.communication.consumer.FileConsumer module

pyDKB.dataflow.communication.consumer.FileConsumer

Data consumer implementation for common (static) files.

**TODO: think about:**

- updatable files
- pipes (better, from the point of StreamConsumer)
- round-robin (for updatable sources)
- ...

```

class pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer(config={})
    Bases: pyDKB.dataflow.communication.consumer.Consumer.Consumer
    Data consumer implementation for HDFS data source.

    current_file = None

    get_source()
        Get nearest non-empty source (current or next).

    get_source_info()
        Return current source info.

```

**init\_sources ()**  
Initialize sources iterator if not initialized yet.

**next\_source ()**  
Reset \$current\_file to the next non-empty file.

**Return value:** File descriptor of the new \$current\_file None (no files left)

**reconfigure (config={})**  
(Re)initialize consumer with Stage configuration.

**source\_is\_empty ()**  
Check if current source is empty.

**Return value:** True (empty) False (not empty) None (no source)

### pyDKB.dataflow.communication.consumer.HDFSConsumer module

pyDKB.dataflow.communication.consumer.HDFSConsumer

**class** pyDKB.dataflow.communication.consumer.HDFSConsumer.**HDFSConsumer** (config={})  
Bases: *pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer*

Data consumer implementation for HDFS data source.

**reconfigure (config={})**  
Configure HDFS Consumer according to the config parameters.

### pyDKB.dataflow.communication.consumer.StreamConsumer module

pyDKB.dataflow.communication.consumer.StreamConsumer

Data consumer implementation for a single stream.

**TODO: think about multiple streams (like a number of named pipes, etc).** Perhaps, even merge this class with FileConsumer.

**class** pyDKB.dataflow.communication.consumer.StreamConsumer.**StreamConsumer** (config={})  
Bases: *pyDKB.dataflow.communication.consumer.Consumer.Consumer*

Data consumer implementation for Stream data source.

**fd = None**

**get\_source ()**  
Get Stream file descriptor.

**get\_source\_info ()**  
Return current source info.

**next\_source ()**  
Return None.

As currently we believe that there is only one input stream

**reconfigure (config={})**  
(Re)configure Stream consumer.

## pyDKB.dataflow.communication.producer package

Producer submodule init file.

```
class pyDKB.dataflow.communication.producer.ProducerBuilder (config={})
    Bases: object

    Constructor for Producer instance.

    build (config={})
        Return constructed producer.

    message_type = None

    producerClass = None

    setDest (dest)
        Set data destination for the producer.

    setSourceInfoMethod (src_info)
        Set method to get current source info.

    setType (Type)
        Set message type for the producer.

    src_info = None
```

## Submodules

### pyDKB.dataflow.communication.producer.FileProducer module

pyDKB.dataflow.communication.producer.FileProducer

Data producer implementation for common (static) files.

**TODO: think about:**

- pipes (better, from the point of StreamProducer)
- multiple parallel dests
- ...

```
class pyDKB.dataflow.communication.producer.FileProducer.FileProducer (config={})
    Bases: pyDKB.dataflow.communication.producer.Producer.Producer

    Data producer implementation for local file data dest.

    close ()
        Close opened files and remove temporary one.

    close_file ()
        Close current file.

    config_dir (config={})
        Configure output directory.

    current_file = None

    default_dir ()
        Get default directory name.
```

**dirname** (*dirname=None*)  
Set/get preferable directory name.

**ensure\_dir** ()  
Ensure that current directory for output files exists.

**file\_info** ()  
Return output file metadata (name, directory, ...).

**get\_dest** ()  
Get destination file descriptor.

**get\_dest\_info** ()  
Get current destination info.

**get\_dir** ()  
Get current directory for output files.

**get\_filename** ()  
Return filename, corresponding the source, or timestamp-based.

**get\_source\_info** ()  
Set current data source, if any.

**reconfigure** (*config={}*)  
(Re)configure producer according to the config hash.

**reset\_file** ()  
Resets current file according to the current source info.

**Metadata include:**

- fd – open file descriptor
- name – file name
- dir – directory name
- local\_path – local path to the file

**set\_default\_dir** ()  
Set default directory name.

**subdir** (*base\_dir, sub\_dir=""*)  
Construct full path for \$subdir of \$base\_dir.

## pyDKB.dataflow.communication.producer.HDFSProducer module

pyDKB.dataflow.communication.producer.HDFSProducer

Data producer implementation for common (static) files in HDFS.

### TODO: think about:

- pipes (better, from the point of StreamProducer)
- multiple parallel dests
- ...

**class** pyDKB.dataflow.communication.producer.HDFSProducer.**HDFSProducer** (*config={}*)  
Bases: *pyDKB.dataflow.communication.producer.FileProducer.FileProducer*

Data producer implementation for HDFS data dest.

```

close_file()
    Close current file and move it to HDFS.

config_dir (config={})
    Configure output directory.

ensure_dir()
    Ensure that current directory for output files exists.

file_info()
    Return output file metadata (name, directory, ...).

set_default_dir()
    Set default directory name.

subdir (base_dir, sub_dir="")
    Construct full path for $sub_dir of $base_dir.

```

## pyDKB.dataflow.communication.producer.Producer module

pyDKB.dataflow.communication.producer.Producer

```

class pyDKB.dataflow.communication.producer.Producer.Producer (config={})
    Bases: object

    Data producer implementation.

    close()
        Close opened data stream and data dest.

    config = None

    drop()
        Drop buffered messages.

    eop()
        Write EOP marker to the current dest.

    flush()
        Flush buffered messages to the current dest.

    get_dest()
        Return current destination.

    get_dest_info()
        Return current dest info.

    get_stream (actualize=True)
        Get output stream linked to the current dest.

        If $actualize parameter set to True, will try to reset current stream destination; else will use last known
        destination or None.

    init_stream()
        Init output stream (without real destination).

    log (message, level=3)
        Output log message with given log level.

    message_class()
        Return message class.

    message_type = None

```

**reconfigure** (*config*={})  
(Re)initialize producer with stage config arguments.

**reset\_stream** ()  
Reset input stream to the current dest.

**set\_message\_type** (*Type*)  
Set input message type.

**write** (*msg*)  
Put new message to the current dest (buffer).

**exception** `pyDKB.dataflow.communication.producer.Producer.ProducerException`  
Bases: `pyDKB.dataflow.exceptions.DataflowException`  
Dataflow Producer exception.

### **pyDKB.dataflow.communication.producer.StreamProducer module**

`pyDKB.dataflow.communication.producer.StreamProducer`

Data producer implementation for a single stream.

**TODO: think about multiple streams (like a number of named pipes, etc).** Perhaps, even merge this class with `FileProducer`.

**class** `pyDKB.dataflow.communication.producer.StreamProducer.StreamProducer` (*config*={})  
Bases: `pyDKB.dataflow.communication.producer.Producer.Producer`  
Data producer implementation for Stream data dest.  
  
**fd** = **None**  
  
**get\_dest** ()  
Get Stream file descriptor.  
  
**get\_dest\_info** ()  
Return current dest info.  
  
**reconfigure** (*config*={})  
(Re)configure Stream producer.

### **pyDKB.dataflow.communication.stream package**

`pyDKB.dataflow.communication.stream`

**class** `pyDKB.dataflow.communication.stream.StreamBuilder` (*fd*, *config*={})  
Bases: `object`  
Constructor for Stream object.  
  
**build** (*config*={})  
Create instance of Stream.  
  
**message\_type** = **None**  
  
**setStream** (*stream*)  
Set stream type: 'input' or 'output'.  
  
**setType** (*Type*)  
Set message type for the Stream.



```

streamClass = None

class pyDKB.dataflow.communication.stream.Stream (fd=None, config={})
    Bases: object

    Abstract class for input/output streams.

    close ()
        Close open file descriptors etc.

    configure (config)
        Stream configuration.

    get_fd ()
        Return open file descriptor or raise exception.

    log (message, level=3)
        Output log message with given log level.

    message_type ()
        Get type of the messages in the stream.

    reset (fd, close=True)
        Reset file descriptor in operation.

        Parameters fd – open file descriptor TODO: IOBase objects

    set_message_type (msg_type)
        Set type of the messages in the stream.

class pyDKB.dataflow.communication.stream.InputStream (fd=None, config={})
    Bases: pyDKB.dataflow.communication.stream.Stream.Stream

    Implementation of the input stream.

    get_message ()
        Get next message from the input stream.

        Return values: Message object False (failed to parse message) None (no messages left)

    next ()
        Get next message from the input stream.

    parse_message (message)
        Verify and parse input message.

        Retrun value: Message object False (failed to parse)

    reset (fd, close=True)
        Reset current stream with new file descriptor.

        Overrides parent method to reset __iterator property.

class pyDKB.dataflow.communication.stream.OutputStream (fd=None, config={})
    Bases: pyDKB.dataflow.communication.stream.Stream.Stream

    Implementation of the output stream.

    configure (config={})
        Configure instance.

    drop ()
        Drop buffer without sending messages anywhere.

    eop ()
        Signalize Supervisor about end of process.

```

**flush()**  
Flush buffer to the output stream.

**msg\_buffer = []**

**write(message)**  
Add message to the buffer.

## Submodules

### pyDKB.dataflow.communication.stream.InputStream module

pyDKB.dataflow.communication.stream.InputStream

**class** pyDKB.dataflow.communication.stream.InputStream.**InputStream** (*fd=None, config={}*)

Bases: *pyDKB.dataflow.communication.stream.Stream.Stream*

Implementation of the input stream.

**get\_message()**  
Get next message from the input stream.

**Return values:** Message object False (failed to parse message) None (no messages left)

**next()**  
Get next message from the input stream.

**parse\_message(message)**  
Verify and parse input message.

**Retrun value:** Message object False (failed to parse)

**reset(fd, close=True)**  
Reset current stream with new file descriptor.

Overrides parent method to reset `__iterator` property.

### pyDKB.dataflow.communication.stream.OutputStream module

pyDKB.dataflow.communication.stream.OutputStream

**class** pyDKB.dataflow.communication.stream.OutputStream.**OutputStream** (*fd=None, config={}*)

Bases: *pyDKB.dataflow.communication.stream.Stream.Stream*

Implementation of the output stream.

**configure(config={})**  
Configure instance.

**drop()**  
Drop buffer without sending messages anywhere.

**eop()**  
Signalize Supervisor about end of process.

**flush()**  
Flush buffer to the output stream.

```
msg_buffer = []

write(message)
    Add message to the buffer.
```

## pyDKB.dataflow.communication.stream.Stream module

pyDKB.dataflow.communication.stream.Stream

```
class pyDKB.dataflow.communication.stream.Stream.Stream (fd=None, config={})
    Bases: object

    Abstract class for input/output streams.

    close()
        Close open file descriptors etc.

    configure(config)
        Stream configuration.

    get_fd()
        Return open file descriptor or raise exception.

    log(message, level=3)
        Output log message with given log level.

    message_type()
        Get type of the messages in the stream.

    reset(fd, close=True)
        Reset file descriptor in operation.

        Parameters fd – open file descriptor TODO: IOBase objects

    set_message_type(msg_type)
        Set type of the messages in the stream.
```

## pyDKB.dataflow.communication.stream.exceptions module

pyDKB.dataflow.communication.stream.exceptions

```
exception pyDKB.dataflow.communication.stream.exceptions.StreamException
    Bases: pyDKB.dataflow.exceptions.DataflowException

    Exception for Stream operations.
```

## Submodules

### pyDKB.dataflow.communication.messages module

pyDKB.dataflow.communication.messages

Definition of abstract message class and specific message classes

```
class pyDKB.dataflow.communication.messages.AbstractMessage (message=None)
    Bases: object

    Abstract message
```

```
content ()
    Return message content.

decode (code)
    Decode original from CODE to TYPE-specific format.

    Raises ValueError

decoded = None

encode (code)
    Encode original message from TYPE-specific format to CODE.

    Raises ValueError

encoded = None

classmethod extension ()
    Return file extension corresponding this message type.

getOriginal ()
    Return original message.

msg_type = None

native_types = []

classmethod typeName ()
    Return message type name as string.

exception pyDKB.dataflow.communication.messages.DecodeUnknownType (code, cls)
    Bases: exceptions.NotImplementedError

    Exception to be thrown when message type is not decodable.

exception pyDKB.dataflow.communication.messages.EncodeUnknownType (code, cls)
    Bases: exceptions.NotImplementedError

    Exception to be thrown when message type is not encodable.

class pyDKB.dataflow.communication.messages.JSONMessage (message=None)
    Bases: pyDKB.dataflow.communication.messages.AbstractMessage

    Message in JSON format.

    decode (code=1)
        Decode original data as JSON.

    encode (code=1)
        Encode JSON as CODE.

    msg_type = 2

    native_types = [<type 'dict'>]

pyDKB.dataflow.communication.messages.Message (msg_type)
    Return class XXXMessage, where XXX is the passed type.

class pyDKB.dataflow.communication.messages.TTLMessage (message=None)
    Bases: pyDKB.dataflow.communication.messages.AbstractMessage

    Messages in TTL format

    Single message = single TTL statement
```

```

decode (code=1)
    Decode original data as TTL.

    Currently takes text as it is. TODO: check some formal matter to confirm the string is TTL.

encode (code=1)
    Encode TTL as CODE.

msg_type = 3

native_types = [<type 'str'>, <type 'unicode'>]

```

## pyDKB.dataflow.stage package

Stage submodule init file.

```

class pyDKB.dataflow.stage.ProcessorStage (description='DKB Dataflow data processing stage.')
    Bases: pyDKB.dataflow.stage.AbstractStage.AbstractStage
    Abstract class to implement Processor stages
    Processor stage – is a stage for data processing/transformation.
    Class/instance variable description:
        • communication.consumer.Consumer instance __input
        • Generator object for output file descriptor OR file descriptor (for (s)stream mode)
          __output
        • List of objects to be “stopped” __stoppable

    clear_buffer ()
        Drop buffered output messages.

    configure (args=None)
        Configure stage according to the config parameters.

        If $args specified, arguments will be parsed anew.

    defaultArguments ()
        Default parser configuration.

    flush_buffer ()
        Flush message buffer to the output.

    forward ()
        Send EOPMarker to the output stream.

    get_out_stream ()
        Get current output stream.

    get_source_info ()
        Get information about current source.

    input ()
        Generator for input messages.

        Returns iterable object. Every iteration returns single input message to be processed.

    input_message_class ()
        Get input message class.

```

**output** (*message*)

Put the (list of) message(s) to the output buffer.

**output\_message\_class** ()

Get output message class.

**static process** (*stage, input\_message*)

Transform input\_message -> output\_message.

To be implemented individually for every stage. Takes the stage as first argument to allow calling output() from inside the function.

**Return value:** True – processing successfully finished False – processing failed (skip the input message)

**run** ()

Run process() for every input() message.

**set\_input\_message\_type** (*Type=None*)

Set input message type.

**set\_output\_message\_type** (*Type=None*)

Set output message class.

**stop** ()

Finalize all the processes and prepare to exit.

## Submodules

### pyDKB.dataflow.stage.AbstractStage module

Definition of an abstract class for Dataflow Stages.

**class** pyDKB.dataflow.stage.AbstractStage.**AbstractStage** (*description='DKB Dataflow stage'*)

Bases: object

Class/instance variable description: \* Argument parser (argparse.ArgumentParser)

\_\_parser

- **Parsed arguments** (argparse.Namespace) ARGS
- **Stage config parser** (ConfigParser.SafeConfigParser) \_\_config
- **Stage custom config** (defaultdict(defaultdict(str))) CONFIG

**add\_argument** (*\*args, \*\*kwargs*)

Add specific (not common) arguments.

**args\_error** (*message*)

Output USAGE, error message and exit with code 2.

**config\_error** (*message='Failed to read config file:'*)

Output error message and exit with code 3.

**defaultArguments** ()

Config argument parser with parameters common for all stages.

**log** (*message, level=3*)

Output log message with given log level.

**output\_error** (*message=None, exc\_info=None*)  
Output traceback of the passed (or last) error with *message*.

**parse\_args** (*args*)  
Parse arguments and set dependant arguments if needed.

**Exits in case of error with code:** 2 – failed to parse arguments 3 – failed to read config file

**print\_usage** (*fd=<open file '<stderr>', mode 'w'>*)  
Print usage message.

**read\_config** ()  
Reads stage custom config file.

**Returns** (True|False)

**run** ()  
Run the stage.

**set\_error** (*err\_type, err\_val, err\_trace*)  
Set object *\_err* variable from the last error info.

**stop** ()  
Stop running processes and output error information.

## pyDKB.dataflow.stage.ProcessorStage module

Definition of an abstract class for Dataflow Data Processing Stages.

**USAGE:** ProcessorStage [<options>] [<input files>]

### OPTIONS:

<b>-s, --source</b>	{flsh} - where to get data from: local (f)iles, (s)tdin, (h)dfs
<b>-i, --input-dir</b>	DIR - base directory for relative input file names (for local and HDFS sources). If <input files> not specified, all files from the directory will be taken as the input.
<b>-d, --dest</b>	{flsh} - where to send data to: local (f)iles, (s)tdout, (h)dfs
<b>-o, --output-dir</b>	DIR - base directory for output files (for local and HDFS sources)
<b>--hdfs</b>	• equivalent to “--source h --dest h”
<b>-m, --mode</b>	MODE - MODE: (f)ile = --source f --dest f (can be rewritten with ‘s’ or ‘h’) (s)tream = --source s (can be rewritten with ‘h’) --dest s (m)apreduce = --source s (can be rewritten with ‘h’) --dest s

```
class pyDKB.dataflow.stage.ProcessorStage.ProcessorStage (description='DKB  
Dataflow data processing  
stage.')
```

Bases: `pyDKB.dataflow.stage.AbstractStage.AbstractStage`

Abstract class to implement Processor stages

Processor stage – is a stage for data processing/transformation.

Class/instance variable description:

- **communication.consumer.Consumer instance** `__input`
- Generator object for output file descriptor OR file descriptor (for (s)stream mode)  
`__output`
- **List of objects to be “stopped”** `__stoppable`

**clear\_buffer** ()  
Drop buffered output messages.

**configure** (*args=None*)  
Configure stage according to the config parameters.  
If \$args specified, arguments will be parsed anew.

**defaultArguments** ()  
Default parser configuration.

**flush\_buffer** ()  
Flush message buffer to the output.

**forward** ()  
Send EOPMarker to the output stream.

**get\_out\_stream** ()  
Get current output stream.

**get\_source\_info** ()  
Get information about current source.

**input** ()  
Generator for input messages.  
Returns iterable object. Every iteration returns single input message to be processed.

**input\_message\_class** ()  
Get input message class.

**output** (*message*)  
Put the (list of) message(s) to the output buffer.

**output\_message\_class** ()  
Get output message class.

**static process** (*stage, input\_message*)  
Transform input\_message -> output\_message.  
To be implemented individually for every stage. Takes the stage as first argument to allow calling output()  
from inside the function.

**Return value:** True – processing successfully finished False – processing failed (skip the input message)



```

run ()
    Run process() for every input() message.

set_input_message_type (Type=None)
    Set input message type.

set_output_message_type (Type=None)
    Set output message class.

stop ()
    Finalize all the processes and prepare to exit.

```

## Submodules

### pyDKB.dataflow.cds module

Extended CDSInvenioConnector allowing us to login via Kerberos

```

class pyDKB.dataflow.cds.CDSInvenioConnector (*args)
    Bases: invenio_client.contrib.cds.CDSInvenioConnector

    CDSInvenioConnector which closes the browser in most cases.

    delete (restore_handlers=True)

    handlers = False

    kill (signum, frame)
        Run del and propagate signal.

    orig_handlers = {}

class pyDKB.dataflow.cds.KerberizedCDSInvenioConnector (login='user',          pass-
                                                         word='password')
    Bases: pyDKB.dataflow.cds.CDSInvenioConnector

    Represents same CDSInvenioConnector, but this one is aware about SPNEGO: Simple and Protected GSSAPI
    Negotiation Mechanism

```

### pyDKB.dataflow.dkbID module

Utils to generate unique yet meaningful identifier for DKB objects.

```

pyDKB.dataflow.dkbID.dkbID (json_data, data_type)
    Return unique identifier for object of TYPE based on DATA.

```

### pyDKB.dataflow.exceptions module

Definition of DKB Dataflow exceptions

```

exception pyDKB.dataflow.exceptions.DataflowException
    Bases: exceptions.Exception

    Base Exception for Dataflow modules.

```

### **pyDKB.dataflow.types module**

Type definitions for library objects.

## Stage 055

Stage for converting JSON files(output of stage 015) into TTL files(input for stage 060).

JSON file should have the following structure:

```
{
  "GLANCE": {},
  "CDS" : {},
  "dkbID" : ...,
  "supporting_notes": [
    {
      "GLANCE": {},
      "CDS": {},
      "dkbID": ...,
    },
    {
      ...
    }
  ]
}
```

Resulting TTL file has the following structure:

```
PAPER a atlas:Paper .
PAPER atlas:hasGLANCE_ID __ .
PAPER atlas:hasShortTitle __ .
PAPER atlas:hasFullTitle __ .
PAPER atlas:hasRefCode __ .
PAPER atlas:hasCreationDate __ .
PAPER atlas:hasCDSReportNumber __ .
PAPER atlas:hasCDSInternal __ .
PAPER atlas:hasCDS_ID __ .
PAPER atlas:hasAbstract __ .
PAPER atlas:hasArXivCode __ .
PAPER atlas:hasFullTitle __ .
PAPER atlas:hasDOI __ .
PAPER atlas:hasKeyword __ .
JOURNAL_ISSUE a atlas:JournalIssue .
JOURNAL_ISSUE atlas:hasTitle __ .
JOURNAL_ISSUE atlas:hasVolume __ .
JOURNAL_ISSUE atlas:hasYear __ .
JOURNAL_ISSUE atlas:containsPublication> PAPER .
SUPPORTING_DOCUMENT a atlas:SupportingDocument .
SUPPORTING_DOCUMENT atlas:hasGLANCE_ID __ .
SUPPORTING_DOCUMENT atlas:hasLabel __ .
SUPPORTING_DOCUMENT atlas:hasURL __ .
```

```
SUPPORTING_DOCUMENT atlas:hasCreationDate __ .
SUPPORTING_DOCUMENT atlas:hasCDSInternal __ .
SUPPORTING_DOCUMENT atlas:hasCDS_ID __ .
SUPPORTING_DOCUMENT atlas:hasAbstract __ .
SUPPORTING_DOCUMENT atlas:hasKeyword __ .
PAPER atlas:isBasedOn SUPPORTING_DOCUMENT .
```

**TODO: This module doesn't convert authors metadata.** This task is still under consideration.

055\_documents2TTL.documents2ttl.**abstract\_extraction**(data)

Extract abstract from JSON.

**Parameters** data (dict) – JSON string

**Returns** abstract

**Return type** str

055\_documents2TTL.documents2ttl.**arxiv\_extraction**(data)

Extract arXiv code from JSON.

**Parameters** data (dict) – JSON string

**Returns** arXiv code

**Return type** str

055\_documents2TTL.documents2ttl.**cds\_id\_extraction**(data)

Extract CDS id from JSON.

**Parameters** data (dict) – JSON string

**Returns** CDS id

**Return type** int

055\_documents2TTL.documents2ttl.**cds\_internal\_extraction**(data)

Extract CDS internal report number parameter from JSON string.

**Parameters** data (dict) – JSON data from file or stream

**Returns** CDS internal report number

**Return type** unicode

055\_documents2TTL.documents2ttl.**cds\_parameter\_extraction**(param\_name,json\_data)

Extract CDS parameter value from CDS JSON.

**Parameters**

- **param\_name** (str) – name of the parameter, defined in \*\_CDS\_ATTRS dict
- **json\_data** – JSON with CDS parameters

**Returns** parameter value

**Return type** int, str

055\_documents2TTL.documents2ttl.**creation\_date\_extraction**(data)

Extract creation date from JSON.

**Parameters** data (dict) – JSON string

**Returns** creation date

**Return type** str

055\_documents2TTL.documents2ttl.**define\_globals** (*args*)

Define global variables for further usage in other functions.

**Parameters** *args* (*argparse.Namespace*) – stage arguments

055\_documents2TTL.documents2ttl.**document\_cds** (*data*, *doc\_iri*, *cds\_attrs*)

Convert CDS metadata from JSON to TTL.

**Parameters**

- **data** (*dict*) – JSON data from file or stream
- **doc\_iri** (*str*) – document IRI for current graph
- **cds\_attrs** (*list*) – PAPER\_CDS\_ATTRS | NOTE\_CDS\_ATTRS

**Returns** TTL string with CDS metadata

**Return type** *str*

055\_documents2TTL.documents2ttl.**document\_glance** (*data*, *doc\_iri*, *glance\_attrs*)

Convert GLANCE metadata from JSON to TTL.

**Parameters**

- **data** (*dict*) – JSON data from file or stream
- **doc\_iri** (*str*) – document IRI for current graph
- **glance\_attrs** (*list*) – PAPER\_GLANCE\_ATTRS | NOTE\_GLANCE\_ATTRS

**Returns** TTL string with GLANCE metadata

**Return type** *str*

055\_documents2TTL.documents2ttl.**documents\_links** (*data*)

Convert links from JSON to TTL.

The result looks as following: PAPER atlas:isBasedOn SUPPORTING\_DOCUMENT

**Parameters** *data* (*dict*) – JSON data from file or stream

**Returns** TTL string with links

**Return type** *str*

055\_documents2TTL.documents2ttl.**doi2ttl** (*doi*, *doc\_iri*)

Convert DOI parameter to TTL.

**Parameters**

- **doi** (*str*, *unicode* or *list*) – doi from JSON string
- **doc\_iri** (*str*) – document IRI for current graph

**Returns** TTL string with DOI

**Return type** *str*

055\_documents2TTL.documents2ttl.**fix\_list\_values** (*list\_vals*)

Apply fix\_string to each item in a list.

**Parameters** *list\_vals* (*list*) – list with strings to be fixed

**Returns** list with fixed strings

**Return type** *list*

055\_documents2TTL.documents2ttl.**fix\_string**(*wrong\_string*)

Fix escape sequences in a string.

**Parameters** **wrong\_string** (*str*, *unicode*) – string to be fixed

**Returns** fixed string

**Return type** *str*

055\_documents2TTL.documents2ttl.**generate\_journal\_id**(*journal\_dict*)

Generate a journal issue ID based on title, volume and year.

**Parameters** **journal\_dict** (*dict*) – journal parameters

**Returns** journal ID

**Return type** *str*

055\_documents2TTL.documents2ttl.**get\_document\_iri**(*doc\_id*)

Construct an IRI for a document.

**Parameters** **doc\_id** (*str*) – document id

**Returns** IRI

**Return type** *str*

055\_documents2TTL.documents2ttl.**glance\_parameter\_extraction**(*param\_name*,  
*json\_data*)

Extract a parameter value from GLANCE JSON.

**Parameters**

- **param\_name** (*str*) – name of the parameter
- **json\_data** (*dict*) – JSON with GLANCE metadata

**Returns** parameter value

**Return type** *str*, *unicode*

055\_documents2TTL.documents2ttl.**keywords2ttl**(*keywords*, *doc\_iri*)

Convert keywords from JSON string to TTL.

**Parameters**

- **keywords** (*dict or list of dicts*) – keywords parameters from JSON string
- **doc\_iri** (*str*) – document IRI for current graph

**Returns** TTL string with keywords

**Return type** *str*

055\_documents2TTL.documents2ttl.**main**(*argv*)

Parse command line arguments and run the stage.

**Parameters** **argv** (*list*) – arguments

055\_documents2TTL.documents2ttl.**process**(*stage*, *msg*)

Process a message. Convert the message's contents from JSON to TTL.

**Parameters**

- **stage** (*pyDKB.dataflow.stage.ProcessorStage*) – stage instance
- **msg** (*pyDKB.dataflow.Message*) – input message with JSON data

055\_documents2TTL.documents2ttl.**process\_journals**(*data*, *doc\_iri*)

Convert journal data from JSON to TTL.

**Parameters**

- **data** (*list*, *dict*) – JSON
- **doc\_iri** (*str*) – document IRI for current graph

**Returns** TTL string with journal issue with connection to paper

**Return type** str

055\_documents2TTL.documents2ttl.**report\_number\_extraction**(*data*)

Extract report number from JSON string.

**Parameters** **data** – JSON data from file or stream

**Returns** report number

**Return type** unicode

055\_documents2TTL.documents2ttl.**title\_extraction**(*data*)

Extracting title from JSON.

**Parameters** **data** (*dict*) – JSON string

**Returns** title

**Return type** str





## Indices and tables

- `genindex`
- `modindex`
- `search`



## PYTHON MODULE INDEX

### 0

055\_documents2TTL.documents2ttl, [23](#)

### p

pyDKB, [1](#)

pyDKB.common, [3](#)

pyDKB.common.custom\_readline, [4](#)

pyDKB.common.exceptions, [4](#)

pyDKB.common.hdfs, [4](#)

pyDKB.common.json\_utils, [5](#)

pyDKB.common.Type, [3](#)

pyDKB.dataflow, [6](#)

pyDKB.dataflow.cds, [21](#)

pyDKB.dataflow.communication, [6](#)

pyDKB.dataflow.communication.consumer,  
[6](#)

pyDKB.dataflow.communication.consumer.Consumer,  
[6](#)

pyDKB.dataflow.communication.consumer.FileConsumer,  
[7](#)

pyDKB.dataflow.communication.consumer.HDFSConsumer,  
[8](#)

pyDKB.dataflow.communication.consumer.StreamConsumer,  
[8](#)

pyDKB.dataflow.communication.messages,  
[15](#)

pyDKB.dataflow.communication.producer,  
[9](#)

pyDKB.dataflow.communication.producer.FileProducer,  
[9](#)

pyDKB.dataflow.communication.producer.HDFSProducer,  
[10](#)

pyDKB.dataflow.communication.producer.Producer,  
[11](#)

pyDKB.dataflow.communication.producer.StreamProducer,  
[12](#)

pyDKB.dataflow.communication.stream, [12](#)

pyDKB.dataflow.communication.stream.exceptions,  
[15](#)

pyDKB.dataflow.communication.stream.InputStream,  
[14](#)

pyDKB.dataflow.communication.stream.OutputStream,  
[14](#)

pyDKB.dataflow.communication.stream.Stream,  
[15](#)

pyDKB.dataflow.dkbID, [21](#)

pyDKB.dataflow.exceptions, [21](#)

pyDKB.dataflow.stage, [17](#)

pyDKB.dataflow.stage.AbstractStage, [18](#)

pyDKB.dataflow.stage.ProcessorStage, [19](#)

pyDKB.dataflow.types, [22](#)



## Symbols

055\_documents2TTL.documents2ttl (module), 23

## A

abstract\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

AbstractMessage (class in pyDKB.dataflow.communication.messages), 15

AbstractStage (class in pyDKB.dataflow.stage.AbstractStage), 18

add() (pyDKB.common.Type.Type method), 3

add\_argument() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18

args\_error() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18

arxiv\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

## B

basename() (in module pyDKB.common.hdfs), 5

build() (pyDKB.dataflow.communication.consumer.ConsumerBuilder method), 6

build() (pyDKB.dataflow.communication.producer.ProducerBuilder method), 9

build() (pyDKB.dataflow.communication.stream.StreamBuilder method), 12

## C

cds\_id\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

cds\_internal\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

cds\_parameter\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

CDSInvenioConnector (class in pyDKB.dataflow.cds), 21

check\_stderr() (in module pyDKB.common.hdfs), 5

clear\_buffer() (pyDKB.dataflow.stage.ProcessorStage method), 17

clear\_buffer() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

close() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 6

close() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

close() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

close() (pyDKB.dataflow.communication.stream.Stream method), 13

close() (pyDKB.dataflow.communication.stream.Stream.Stream method), 15

close\_file() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

close\_file() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 10

config (pyDKB.dataflow.communication.consumer.Consumer.Consumer attribute), 6

config (pyDKB.dataflow.communication.producer.Producer.Producer attribute), 11

config\_dir() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

config\_dir() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

config\_error() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18

configure() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 13

configure() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 14

configure() (pyDKB.dataflow.communication.stream.Stream method), 13

configure() (pyDKB.dataflow.communication.stream.Stream.Stream method), 15

configure() (pyDKB.dataflow.stage.ProcessorStage method), 17

configure() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

Consumer (class in pyDKB.dataflow.communication.consumer.Consumer), 6

ConsumerBuilder (class in pyDKB.dataflow.communication.consumer.ConsumerBuilder), 6

consumerClass (pyDKB.dataflow.communication.consumer.Consumer), 6

ConsumerException, 7

content() (pyDKB.dataflow.communication.messages.AbstractMessage method), 14

creation\_date\_extraction() (in module 055\_documents2TTL.documents2ttl), 24

current\_file (pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer attribute), 7

current\_file (pyDKB.dataflow.communication.producer.FileProducer.FileProducer attribute), 9

custom\_readline() (in module pyDKB.common.custom\_readline), 4

## D

DataflowException, 21

decode() (pyDKB.dataflow.communication.messages.AbstractMessage method), 10

decode() (pyDKB.dataflow.communication.messages.JSONMessage method), 11

decode() (pyDKB.dataflow.communication.messages.TTLMessage method), 11

decoded (pyDKB.dataflow.communication.messages.AbstractMessage attribute), 16

DecodeUnknownType, 16

default\_dir() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

defaultArguments() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18

defaultArguments() (pyDKB.dataflow.stage.ProcessorStage method), 17

defaultArguments() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

define\_globals() (in module 055\_documents2TTL.documents2ttl), 24

delete() (pyDKB.dataflow.cds.CDSInvenioConnector method), 21

dirname() (in module pyDKB.common.hdfs), 5

dirname() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

dkbID() (in module pyDKB.dataflow.dkbID), 21

document\_cds() (in module 055\_documents2TTL.documents2ttl), 25

document\_glance() (in module 055\_documents2TTL.documents2ttl), 25

documents\_links() (in module 055\_documents2TTL.documents2ttl), 25

doi2ttl() (in module 055\_documents2TTL.documents2ttl), 25

drop() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

drop() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 13

drop() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 14

## E

encode() (pyDKB.dataflow.communication.messages.AbstractMessage method), 10

encode() (pyDKB.dataflow.communication.messages.JSONMessage method), 11

encode() (pyDKB.dataflow.communication.messages.TTLMessage method), 17

encoded (pyDKB.dataflow.communication.messages.AbstractMessage attribute), 16

EncodeUnknownType, 16

ensure\_dir() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 9

ensure\_dir() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

eop() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

eop() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 13

eop() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 14

## F

fd (pyDKB.dataflow.communication.consumer.StreamConsumer.StreamConsumer attribute), 8

fd (pyDKB.dataflow.communication.producer.StreamProducer.StreamProducer attribute), 12

File() (in module pyDKB.common.hdfs), 4

file\_info() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

file\_info() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

FileConsumer (class in pyDKB.dataflow.communication.consumer.FileConsumer), 7

FileProducer (class in pyDKB.dataflow.communication.producer.FileProducer), 9

fix\_list\_values() (in module 055\_documents2TTL.documents2ttl), 25

fix\_string() (in module 055\_documents2TTL.documents2ttl), 25

flush() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

flush() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 13

flush() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 14

flush\_buffer() (pyDKB.dataflow.stage.ProcessorStage method), 17

flush\_buffer() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

forward() (pyDKB.dataflow.stage.ProcessorStage method), 17

forward() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

## G

generate\_journal\_id() (in module 055\_documents2TTL.documents2ttl), 26

get\_dest() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

get\_dest() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

get\_dest() (pyDKB.dataflow.communication.producer.StreamProducer.StreamProducer method), 12

get\_dest\_info() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

get\_dest\_info() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

get\_dest\_info() (pyDKB.dataflow.communication.producer.StreamProducer.StreamProducer method), 12

get\_dir() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

get\_document\_iri() (in module 055\_documents2TTL.documents2ttl), 26

get\_fd() (pyDKB.dataflow.communication.stream.Stream method), 13

get\_fd() (pyDKB.dataflow.communication.stream.Stream.Stream method), 15

get\_filename() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

get\_message() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 6

get\_message() (pyDKB.dataflow.communication.stream.InputStream.InputStream method), 13

get\_message() (pyDKB.dataflow.communication.stream.InputStream.InputStream method), 14

get\_out\_stream() (pyDKB.dataflow.stage.ProcessorStage method), 17

get\_out\_stream() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

get\_source() (pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer method), 7

get\_source() (pyDKB.dataflow.communication.consumer.StreamConsumer.StreamConsumer method), 8

get\_source\_info() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 6

get\_source\_info() (pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer method), 7

get\_source\_info() (pyDKB.dataflow.communication.consumer.StreamConsumer.StreamConsumer method), 8

get\_source\_info() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

get\_source\_info() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

get\_stream() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 7

get\_stream() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11

getfile() (in module pyDKB.common.hdfs), 5

getOriginal() (pyDKB.dataflow.communication.messages.AbstractMessage method), 16

glance\_parameter\_extraction() (in module 055\_documents2TTL.documents2ttl), 26

## H

handlers (pyDKB.dataflow.cds.CDSInvenioConnector at pyDKB.dataflow.cds), 21

hasMember() (pyDKB.common.Type.Type method), 3

HDFSConsumer (class in pyDKB.dataflow.communication.consumer.HDFSConsumer), 4

HDFSException, 4

HDFSProducer (class in pyDKB.dataflow.communication.producer.HDFSProducer), 10

## I

init\_sources() (pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer method), 7

init\_stream() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 7

init\_stream() (pyDKB.dataflow.communication.producer.Producer.Producer method), 1

input() (pyDKB.dataflow.stage.ProcessorStage method), 17

input() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

input\_message\_class() (pyDKB.dataflow.stage.ProcessorStage method), 17

InputStage\_class() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20

InputStream (class in pyDKB.dataflow.communication.stream), 13

InputStream (class in pyDKB.dataflow.communication.stream.InputStream), 14

## J

join() (in module pyDKB.common.hdfs), 5

JSONMessage (class in pyDKB.dataflow.communication.messages), 16

## K

KerberizedCDSInvenioConnector (class in pyDKB.dataflow.cds), 21  
 keywords2ttl() (in module 055\_documents2TTL.documents2ttl), 26  
 kill() (pyDKB.dataflow.cds.CDSInvenioConnector method), 21

## L

listdir() (in module pyDKB.common.hdfs), 5  
 log() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 7  
 log() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11  
 log() (pyDKB.dataflow.communication.stream.Stream method), 13  
 log() (pyDKB.dataflow.communication.stream.Stream.Stream method), 15  
 log() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18

## M

main() (in module 055\_documents2TTL.documents2ttl), 26  
 makedirs() (in module pyDKB.common.hdfs), 5  
 member() (pyDKB.common.Type.Type method), 3  
 memberName() (pyDKB.common.Type.Type method), 3  
 Message() (in module pyDKB.dataflow.communication), 6  
 Message() (in module pyDKB.dataflow.communication.messages), 16  
 message\_class() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 7  
 message\_class() (pyDKB.dataflow.communication.producer.Producer.Producer method), 11  
 message\_type (pyDKB.dataflow.communication.consumer.Consumer.Consumer attribute), 7  
 message\_type (pyDKB.dataflow.communication.producer.Producer.Producer attribute), 11  
 message\_type (pyDKB.dataflow.communication.producer.Producer.Producer attribute), 9  
 message\_type (pyDKB.dataflow.communication.stream.StreamBuilder attribute), 12  
 message\_type() (pyDKB.dataflow.communication.stream.Stream method), 13  
 message\_type() (pyDKB.dataflow.communication.stream.Stream.Stream method), 15  
 movefile() (in module pyDKB.common.hdfs), 5

msg\_buffer (pyDKB.dataflow.communication.stream.OutputStream attribute), 14  
 msg\_buffer (pyDKB.dataflow.communication.stream.OutputStream.OutputStream attribute), 14  
 msg\_type (pyDKB.dataflow.communication.messages.AbstractMessage attribute), 16  
 msg\_type (pyDKB.dataflow.communication.messages.JSONMessage attribute), 16  
 msg\_type (pyDKB.dataflow.communication.messages.TTLMessage attribute), 17

## N

native\_types (pyDKB.dataflow.communication.messages.AbstractMessage attribute), 16  
 native\_types (pyDKB.dataflow.communication.messages.JSONMessage attribute), 16  
 native\_types (pyDKB.dataflow.communication.messages.TTLMessage attribute), 17  
 nestedKeys() (in module pyDKB.common.json\_utils), 5  
 next() (pyDKB.dataflow.communication.consumer.Consumer.Consumer method), 7  
 next() (pyDKB.dataflow.communication.stream.InputStream method), 13  
 next() (pyDKB.dataflow.communication.stream.InputStream.InputStream method), 14  
 next\_source() (pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer method), 8  
 next\_source() (pyDKB.dataflow.communication.consumer.StreamConsumer.StreamConsumer method), 8

## O

orig\_handlers (pyDKB.dataflow.cds.CDSInvenioConnector attribute), 21  
 output() (pyDKB.dataflow.stage.ProcessorStage method), 17  
 output() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 20  
 output\_error() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 18  
 output\_message\_class() (pyDKB.dataflow.stage.ProcessorStage method), 18  
 OutputMessage class() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage attribute), 20  
 OutputStream (class in pyDKB.dataflow.communication.stream), 13  
 OutputStream (class in pyDKB.dataflow.communication.stream.OutputStream), 14  
 parse\_args() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 19



[parse\\_message\(\) \(pyDKB.dataflow.communication.stream.InputStream method\), 13](#)  
[parse\\_message\(\) \(pyDKB.dataflow.communication.stream.InputStream method\), 14](#)  
[print\\_usage\(\) \(pyDKB.dataflow.stage.AbstractStage.AbstractStage method\), 19](#)  
[process\(\) \(in module 055\\_documents2TTL.documents2ttl\), 26](#)  
[process\(\) \(pyDKB.dataflow.stage.ProcessorStage static method\), 18](#)  
[process\(\) \(pyDKB.dataflow.stage.ProcessorStage.ProcessorStage static method\), 20](#)  
[process\\_journals\(\) \(in module 055\\_documents2TTL.documents2ttl\), 26](#)  
[ProcessorStage \(class in pyDKB.dataflow.stage\), 17](#)  
[ProcessorStage \(class in pyDKB.dataflow.stage.ProcessorStage\), 19](#)  
[Producer \(class in pyDKB.dataflow.communication.producer.Producer\), 11](#)  
[ProducerBuilder \(class in pyDKB.dataflow.communication.producer\), 9](#)  
[producerClass \(pyDKB.dataflow.communication.producer.ProducerBuilder attribute\), 9](#)  
[ProducerException, 12](#)  
[putfile\(\) \(in module pyDKB.common.hdfs\), 5](#)  
[pyDKB \(module\), 1](#)  
[pyDKB.common \(module\), 3](#)  
[pyDKB.common.custom\\_readline \(module\), 4](#)  
[pyDKB.common.exceptions \(module\), 4](#)  
[pyDKB.common.hdfs \(module\), 4](#)  
[pyDKB.common.json\\_utils \(module\), 5](#)  
[pyDKB.common.Type \(module\), 3](#)  
[pyDKB.dataflow \(module\), 6](#)  
[pyDKB.dataflow.cds \(module\), 21](#)  
[pyDKB.dataflow.communication \(module\), 6](#)  
[pyDKB.dataflow.communication.consumer \(module\), 6](#)  
[pyDKB.dataflow.communication.consumer.Consumer \(module\), 6](#)  
[pyDKB.dataflow.communication.consumer.FileConsumer \(module\), 7](#)  
[pyDKB.dataflow.communication.consumer.HDFSConsumer \(module\), 8](#)  
[pyDKB.dataflow.communication.consumer.StreamConsumer \(module\), 8](#)  
[pyDKB.dataflow.communication.messages \(module\), 15](#)  
[pyDKB.dataflow.communication.producer \(module\), 9](#)  
[pyDKB.dataflow.communication.producer.FileProducer \(module\), 9](#)  
[pyDKB.dataflow.communication.producer.HDFSProducer \(module\), 10](#)  
[pyDKB.dataflow.communication.producer.Producer \(module\), 11](#)  
[pyDKB.dataflow.communication.producer.StreamProducer \(module\), 12](#)  
[pyDKB.dataflow.communication.stream \(module\), 12](#)  
[pyDKB.dataflow.communication.stream.exceptions \(module\), 15](#)  
[pyDKB.dataflow.communication.stream.InputStream \(module\), 14](#)  
[pyDKB.dataflow.communication.stream.OutputStream \(module\), 14](#)  
[pyDKB.dataflow.communication.stream.Stream \(module\), 15](#)  
[pyDKB.dataflow.dkbID \(module\), 21](#)  
[pyDKB.dataflow.exceptions \(module\), 21](#)  
[pyDKB.dataflow.stage \(module\), 17](#)  
[pyDKB.dataflow.stage.AbstractStage \(module\), 18](#)  
[pyDKB.dataflow.stage.ProcessorStage \(module\), 19](#)  
[pyDKB.dataflow.types \(module\), 22](#)

## R

[read\\_config\(\) \(pyDKB.dataflow.stage.AbstractStage.AbstractStage method\), 19](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.consumer.Consumer.Consumer method\), 7](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.consumer.FileConsumer.FileConsumer method\), 8](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.consumer.HDFSConsumer.HDFSConsumer method\), 8](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.consumer.StreamConsumer.StreamConsumer method\), 8](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.producer.FileProducer.FileProducer method\), 10](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.producer.Producer.Producer method\), 11](#)  
[reconfigure\(\) \(pyDKB.dataflow.communication.producer.StreamProducer.StreamProducer method\), 12](#)  
[report\\_number\\_extraction\(\) \(in module 055\\_documents2TTL.documents2ttl\), 27](#)  
[reset\(\) \(pyDKB.dataflow.communication.stream.InputStream.InputStream method\), 13](#)  
[reset\(\) \(pyDKB.dataflow.communication.stream.InputStream.InputStream method\), 14](#)  
[reset\(\) \(pyDKB.dataflow.communication.stream.Stream.Stream method\), 13](#)  
[reset\(\) \(pyDKB.dataflow.communication.stream.Stream.Stream method\), 15](#)  
[reset\\_file\(\) \(pyDKB.dataflow.communication.producer.FileProducer.FileProducer method\), 10](#)  
[reset\\_stream\(\) \(pyDKB.dataflow.communication.consumer.Consumer.Consumer method\), 7](#)  
[reset\\_stream\(\) \(pyDKB.dataflow.communication.producer.Producer.Producer method\), 12](#)  
[run\(\) \(pyDKB.dataflow.stage.AbstractStage.AbstractStage method\), 19](#)  
[run\(\) \(pyDKB.dataflow.stage.ProcessorStage method\), 18](#)

run() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage.stop() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 20 method), 19

**S**

stop() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 18

stop() (pyDKB.dataflow.stage.ProcessorStage.ProcessorStage method), 21

set\_default\_dir() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 10

set\_default\_dir() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

set\_error() (pyDKB.dataflow.stage.AbstractStage.AbstractStage method), 19

set\_input\_message\_type() (py-DKB.dataflow.stage.ProcessorStage method), 18

set\_input\_message\_type() (py-DKB.dataflow.stage.ProcessorStage.ProcessorStage method), 21

set\_message\_type() (py-DKB.dataflow.communication.consumer.Consumer.Consumer method), 7

set\_message\_type() (py-DKB.dataflow.communication.producer.Producer.Producer method), 12

set\_message\_type() (py-DKB.dataflow.communication.stream.Stream method), 13

set\_message\_type() (py-DKB.dataflow.communication.stream.Stream method), 15

set\_output\_message\_type() (py-DKB.dataflow.stage.ProcessorStage method), 18

set\_output\_message\_type() (py-DKB.dataflow.stage.ProcessorStage.ProcessorStage method), 21

setDest() (pyDKB.dataflow.communication.producer.ProducerBuilder class method), 16

setDest() (pyDKB.dataflow.communication.producer.ProducerBuilder class method), 9

setSource() (pyDKB.dataflow.communication.consumer.ConsumerBuilder method), 6

setSourceInfoMethod() (py-DKB.dataflow.communication.producer.ProducerBuilder method), 9

setStream() (pyDKB.dataflow.communication.stream.StreamBuilder method), 12

setType() (pyDKB.dataflow.communication.consumer.ConsumerBuilder method), 6

setType() (pyDKB.dataflow.communication.producer.ProducerBuilder method), 9

setType() (pyDKB.dataflow.communication.stream.StreamBuilder method), 12

source\_is\_empty() (py-DKB.dataflow.communication.consumer.FileConsumer.FileConsumer method), 8

src\_info (pyDKB.dataflow.communication.producer.ProducerBuilder attribute), 9

Stream (class in pyDKB.dataflow.communication.stream.Stream), 15

StreamBuilder (class in py-DKB.dataflow.communication.stream), 12

streamClass (pyDKB.dataflow.communication.stream.StreamBuilder attribute), 12

StreamConsumer (class in py-DKB.dataflow.communication.consumer.StreamConsumer), 8

StreamException (class in py-DKB.dataflow.communication.producer.StreamProducer), 15

StreamProducer (class in py-DKB.dataflow.communication.producer.StreamProducer), 12

subdir() (pyDKB.dataflow.communication.producer.FileProducer.FileProducer method), 10

subdir() (pyDKB.dataflow.communication.producer.HDFSProducer.HDFSProducer method), 11

**T**

title\_extraction() (in module 055\_documents2ttl.documents2ttl), 27

TTLMessage (class in py-DKB.dataflow.communication.messages), 16

Type (class in pyDKB.common.Type), 3

typeName() (pyDKB.dataflow.communication.messages.AbstractMessage method), 16

**V**

valueByKey() (in module pyDKB.common.json\_utils), 5

**W**

write() (pyDKB.dataflow.communication.producer.Producer.Producer method), 12

write() (pyDKB.dataflow.communication.stream.OutputStream method), 14

write() (pyDKB.dataflow.communication.stream.OutputStream.OutputStream method), 15