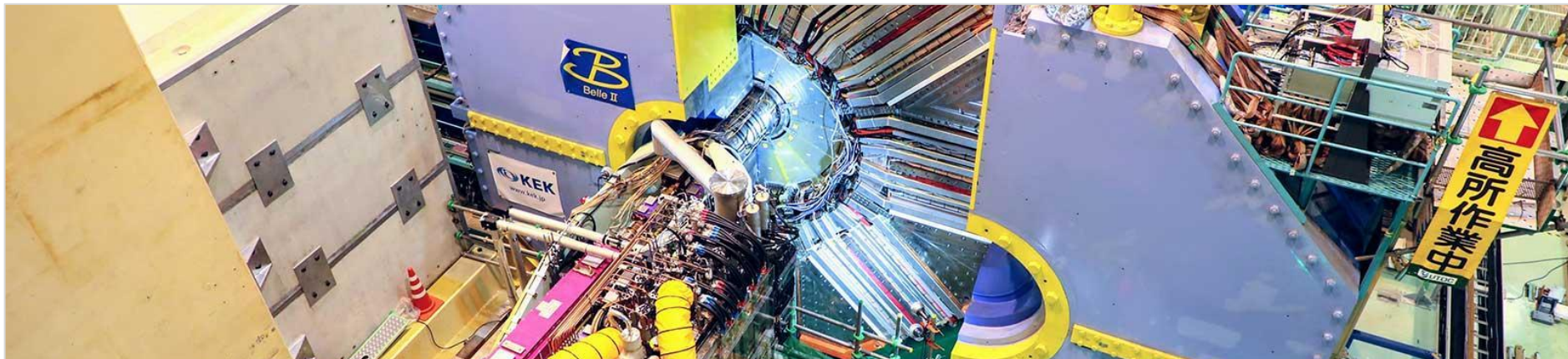# b2luigi - Bringing Batch 2 luigi!

PyHEP.dev 2024 - "Python in HEP" Developer Workshop
Alexander Heidelbach, Jonas Eppelt
alexander.heidelbach@kit.edu

# luigi in a Nutshell

- Building a Pipeline
  - Dependency Resolution
  - Workflow Management
  - Visualisation
  - Handling Failures
  - Command Line Integration
  - …
- https://github.com/spotify/luigi

"Hello World" in luigi:

```python
class MyTask(luigi.Task):
    parameter = luigi.Parameter()

    def run(self):
        do_smth(self.parameter)

    def output(self):
        return Target("some/file")

    def requires(self):
        yield OtherTask()
```
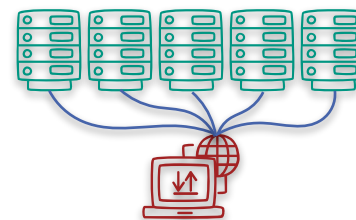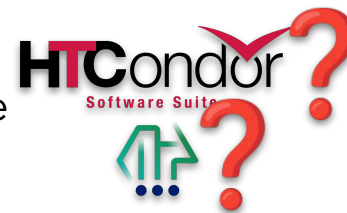
Belle II Group,
ETP, Karlsruhe Institute of Technology

# Why b2luigi?

- Many many many jobs!
  - Running batch job = running process
  - **Problem**: Limitation on the number of processes per user
  - **Solution**: Single process on submission machine
- Many many many tasks!
  - **Problem**: Tasks in luigi need to adjust for batch execution specifically
  - **Solution**: Abstract batch submission away from the task
- Which batch system?
  - **Problem**: Batch system usage defined by task instance
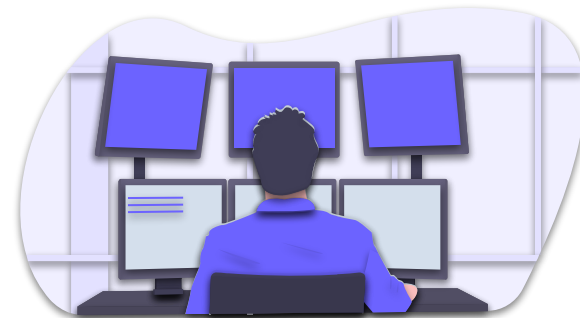  - **Solution**: Batch system usage <u>only</u> defined by config variable

b2luigi.process(...,batch=True)

28 August 2024     b2luigi - Bringing Batch 2 luigi!

Belle II Group,
ETP, Karlsruhe Institute of Technology

# Before we dive into it…

- b2luigi was written by a group of **PhD students** for their analyses
  - **Goal:** Make everyday work a little bit easier
  - **Not a goal:** Invent the next workflow management system
    - A lot of helper functions on top of luigi
    - Easy transition between luigi and b2luigi



- Since this year Belle II has been the official maintainer
  Completely new team of developers
  Targeting collaboration (and beyond…) wide use

# b2luigi vs luigi

```python
import luigi

class MyTask(luigi.Task):
    parameter = luigi.Parameter()

    def run(self):
        with open(f"/my/target/dir/my_output{self.parameter}.file", "w") as f:
            f.write(f"{self.parameter}")

    def output(self):
        return luigi.Target(f"/my/target/dir/my_output{self.parameter}.file")

    def requires(self):
        yield OtherTask()

if __name__ == "__main__":
    luigi.build(
        [MyTask(some_parameter=i) for i in range(100)],
        workers=20
    )
```

```python
import b2luigi as luigi

class MyTask(luigi.Task):
    parameter = luigi.Parameter()

    def run(self):
        with open(self.get_output_file_name("my_output.file"), "w") as f:
            f.write(f"{self.parameter}")

    def output(self):
        yield self.add_to_output("my_output.file")

    def requires(self):
        yield OtherTask()

if __name__ == "__main__":
    luigi.set_setting("result_dir", "/my/target/dir/")
    luigi.process(
        [MyTask(some_parameter=i) for i in range(100)],
        workers=20
    )
```

# Automated Output Bookkeeping

```
===== Luigi Execution Summary =====

Scheduled 100 tasks of which:
* 100 ran successfully:
    - 100 MyTask(some_parameter=0,1,10,11,12,13,14,15,16,17,18,...)

This progress looks :) because there were no failed tasks or missing dependencies

===== Luigi Execution Summary =====
```
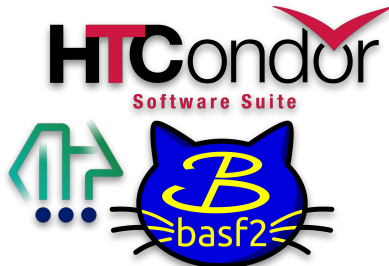
- The self.add_to_output function automatically ensures each output file is unique
- The task's parameter are used to construct a path structure:
  - result_dir
    - some_parameter=1
      - my_output.file
    - some_parameter=2
      - my_output.file
    - ....

- Use luigi.Parameter's significant and hash_function, as well as the declaration order to control this behaviour

Belle II Group,
ETP, Karlsruhe Institute of Technology

# Batch Processing

- Initial motivation: Make batch submission easy!
- Currently fully supported:
  - HTCondor
  - LSF
  - Gbasf2 (BelleII@WLCG)
- Challenges:
  - Using your environment
    - Settings: env_script, env,..
  - Ensuring consistent locations
    - Settings: working_dir, result_dir,…

```python
class MyTask(b2luigi.Task):
    batch_system = "htcondor"
```
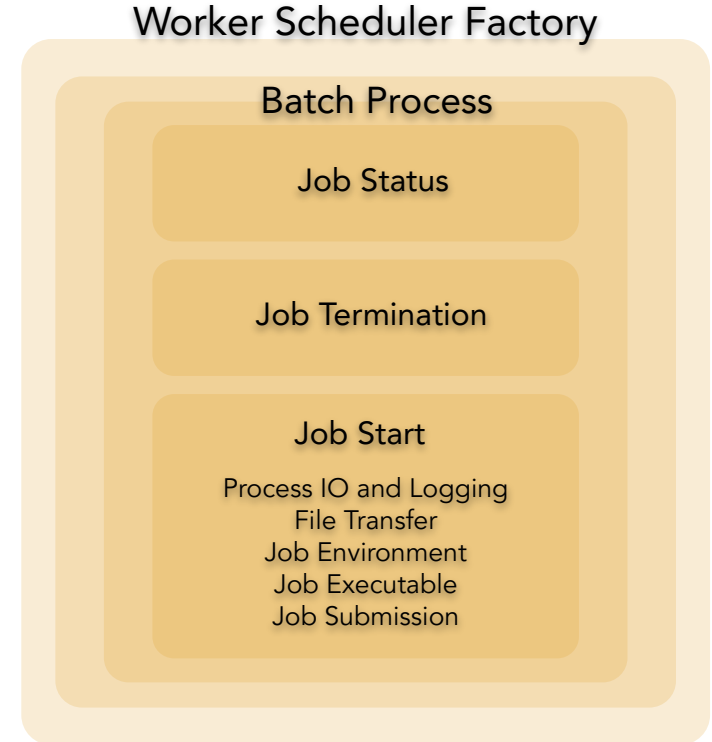
```python
class MyTask(b2luigi.Task):
    @property
    def htcondor_settings(self):
        return {"request_memory": 4096}
```

```python
b2luigi.set_setting("result_dir", "path/to/result")
```

```
settings.json
{
    "env_script": "setup.sh"
}
```

# Batch Processing in b2luigi

- Use of luigi's Worker Scheduler Factory
- Different processes for each batch system
- Batch process handles:
  - Job status
  - Job start
  - Job termination
- Job start consists (mostly) of two steps:
  - Creation of **executable wrapper**
  - **Submission** via batch system mechanism

**Worker Scheduler Factory**

**Batch Process**

Job Status

Job Termination

Job Start

Process IO and Logging
File Transfer
Job Environment
Job Executable
Job Submission

Belle II Group,
ETP, Karlsruhe Institute of Technology

# Task Settings

```python
class MyTask(b2luigi.Task):
    some_setting = "some value"
```
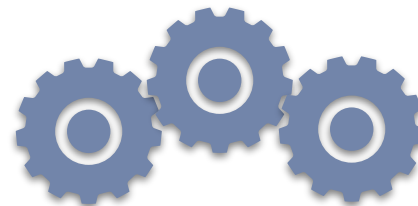
```python
class MyTask(b2luigi.Task):
    @property
    def some_setting(self):
        return "some value"
```

```python
b2luigi.set_setting("some_setting", "some value")
```

```
settings.json
{
    "some_setting": "some value"
}
```

Priority

- Settings are handled by b2luigi!
- Control:
  - Batch system choice
  - Output and log path
  - Environment
  - Workflow specific settings

# Example: HTCondor

```python
class MyTask(b2luigi.Task):
    parameter = b2luigi.IntParameter()
    batch_system = "htcondor"

    @property
    def executable(self):
        return ["$MY_PYTHON"]

    def output(self):
        yield self.add_to_output("test.txt")

    def run(self):
        with open(self.get_output_file_name("test.txt"), "w") as f:
            f.write(f"Test {self.parameter}")

class Wrapper(b2luigi.WrapperTask):
    def requires(self):
        for i in range(100):
            yield MyTask(parameter=i)


if __name__ == "__main__":
    b2luigi.set_setting("env_script", "setup.sh")

    b2luigi.set_setting("result_dir", "results")


    b2luigi.process(Wrapper(), batch=True, workers=100)
```

Script executed on the batch job side

Location needs to be accessible on the batch job side

Wrapper tasks need no output

Definition of the batch system

Executable for this specific task

- E.g. environment variable set in setup.sh

# Summary & Outlook

- b2luigi provides a simple and flexible implementation to run your workflow on batch systems!
- The abstraction of the batch processing to global settings allows for:
  - Quick change in the submission strategy
  - Simple code and execution
- Outlook
  - Code Maintainability
  - XRootDTarget
  - Grid Tool API

## b2luigi

`sphinx` `latest`  `license` `GPL-3.0`  `pypi` `v1.0.1`  `DOI` `10.5281/zenodo.11207742`

`b2luigi` is a helper package constructed around `luigi` that helps you schedule working packages (so-called tasks) locally or on a batch system. Apart from the very powerful dependency management system by `luigi`, `b2luigi` extends the user interface and has a built-in support for the queue systems, e.g. LSF and HTCondor.

You can find more information in the documentation. Please note that most of the core features are handled by `luigi`, which is described in the separate luigi documentation, where you can find a lot of useful information.

If you find any bugs or want to add a feature or improve the documentation, please send me a pull request! Check the development documentation on information how to contribute.

Contributors are listed here.

This project is in still beta. Please be extra cautious when using in production mode.

To get notified about new features, (potentially breaking) changes, bugs and their fixes, I recommend using the `watch` button on GitHub to get notifications for new releases and/or issues or to subscribe the releases feed (requires no GitHub account, just a feed reader).
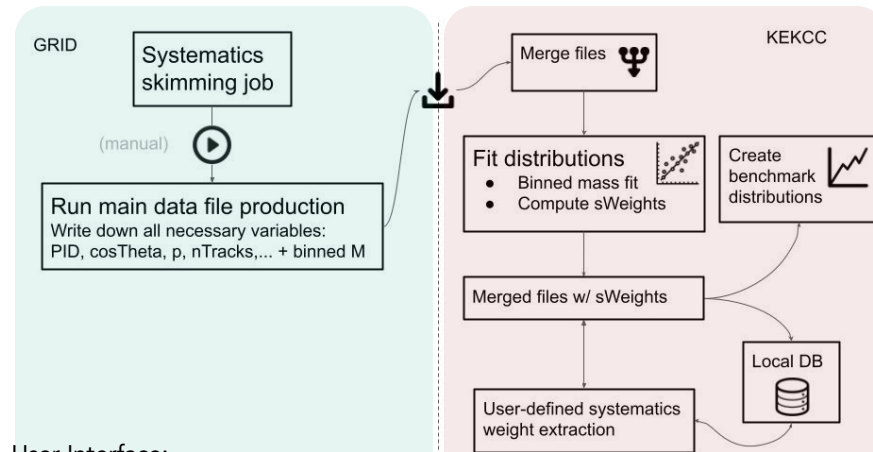
[Github](#)  [Zenodo](#)

[Documentation](#)  [PyPi](#)

Belle II Group,
ETP, Karlsruhe Institute of Technology

# Systematic Corrections Framework



- 2 Stage Algorithm
  - Ntuple Production
    - Centrally run for every campaign
    - Running: Gbasf2
  - Data/MC Corrections
    - User runs their specific selection
    - Running: Locally, HTCondor, LSF
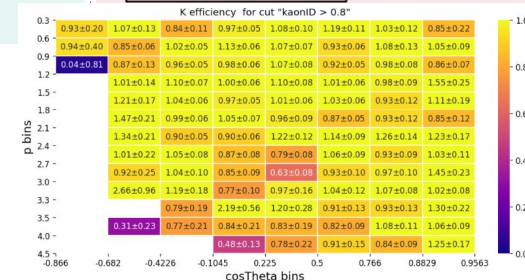- Also: Validation of different datasets

[Documentation](#)

Belle II Group,
ETP, Karlsruhe Institute of Technology

# Validation Interface for the Belle II Experiment

Belle II Group,
ETP, Karlsruhe Institute of Technology