

W A h y n o t h e r w a l u i g i ?

# waluigi - Beyond luigi

Benjamin Fischer

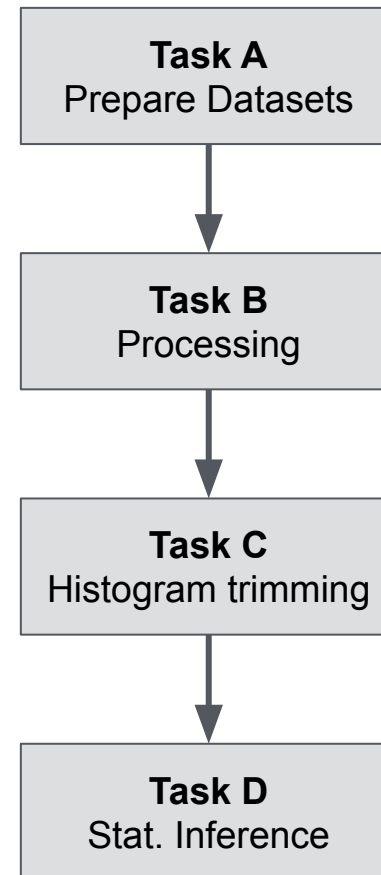
PyHEP.dev  
August 28<sup>th</sup> 2024

### Luigi's strengths:

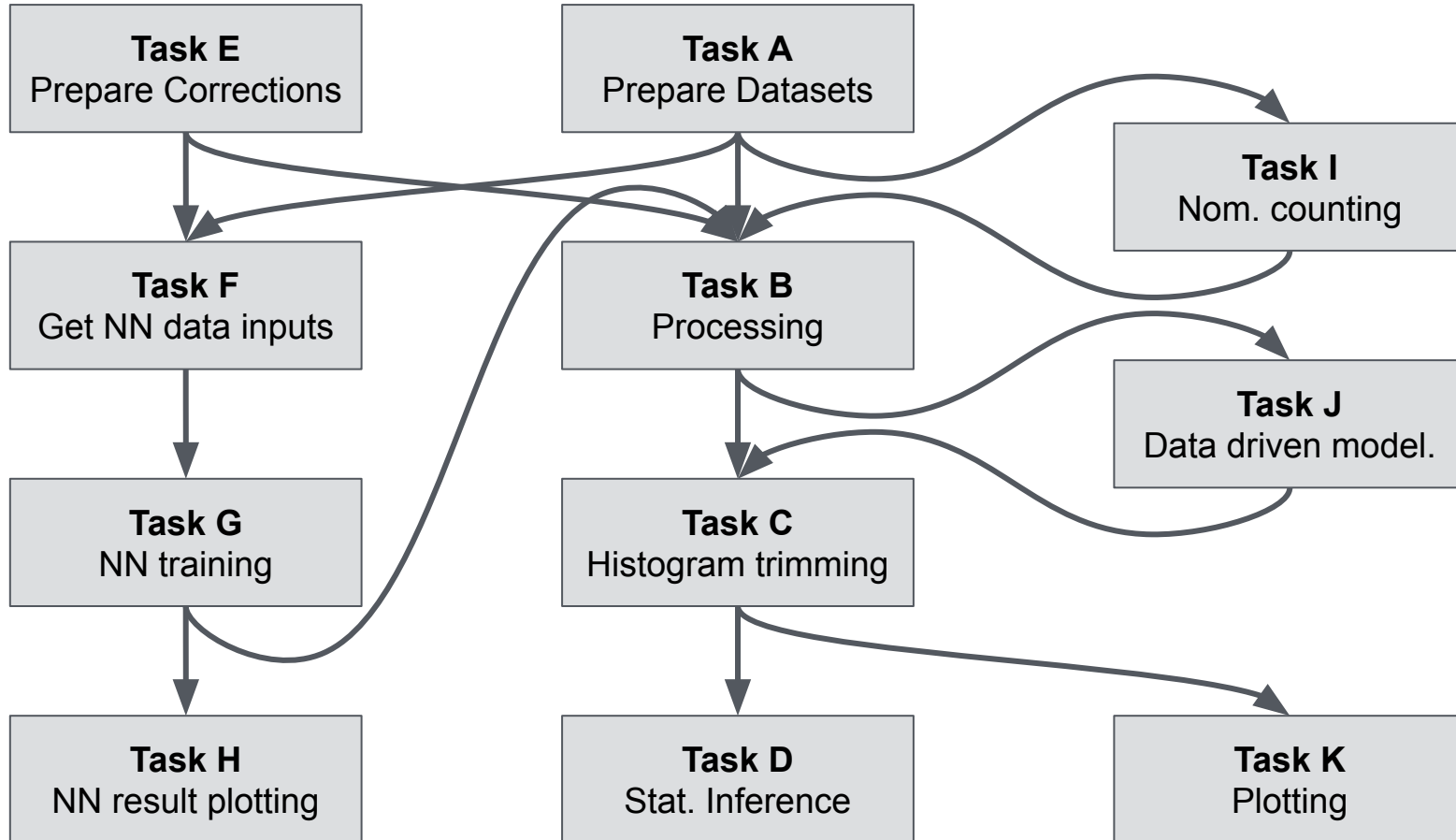
- nice web interface
  - well supported
  - powerful workflow description
    - multiple inputs & outputs
    - customizable “targets”  
i.e. local & remote files
    - **full task parametrization**
- ↙
- similar to dataclasses
  - effectively unique feature

Example Graph:

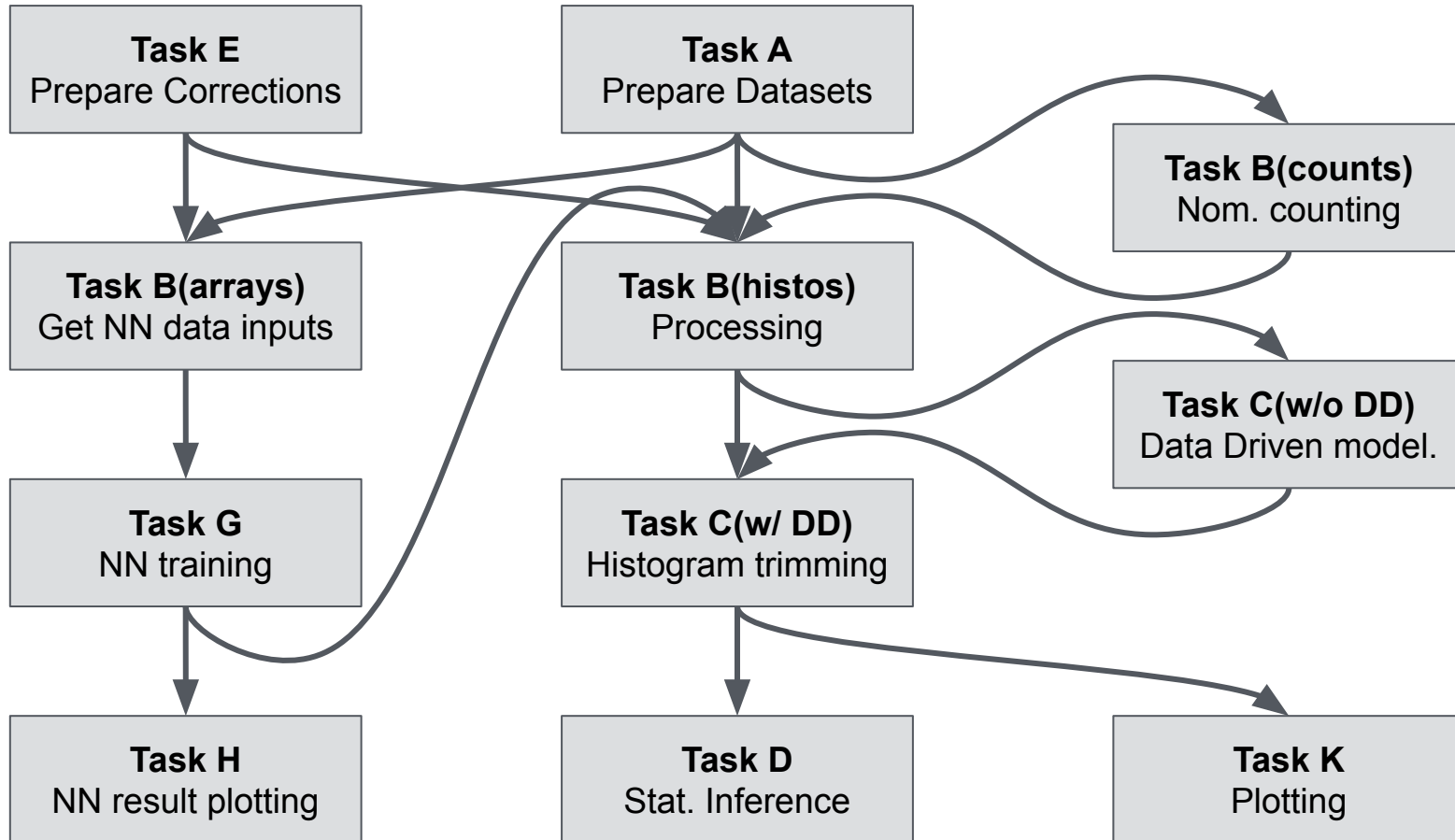
↓ Data flow



### 3 A more complex/realistic graph

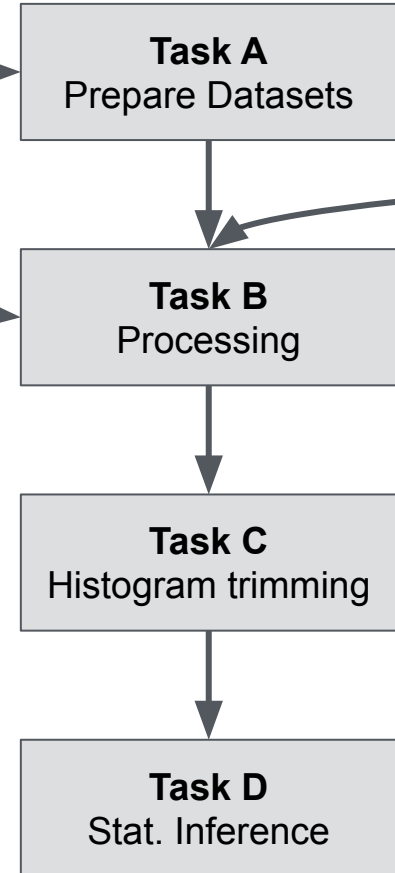


# 4 Now with parameters



# 5 An Example Task

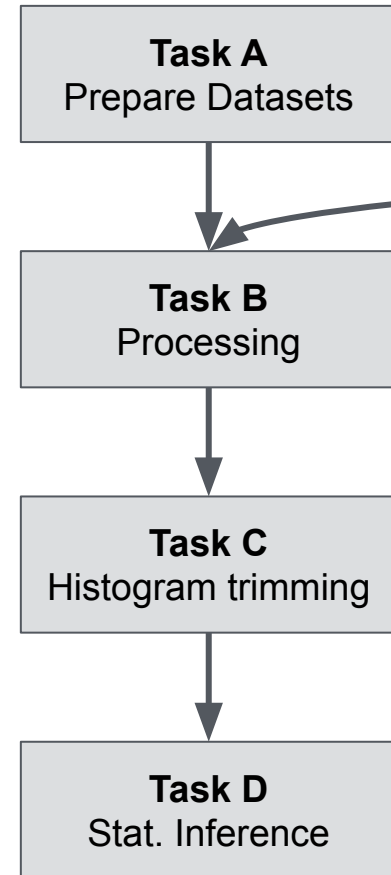
```
1 import luigi
2 from somewhere import TaskE
3
4 class TaskA(luigi.Task):
5     ...
6
7     def get_all_paths():
8         ...
9
10 class TaskB(luigi.Task):
11     what = luigi.ChoiceParameter(
12         choices=["histos", "array", "counts"],
13         description="what to produce when processing data files"
14     )
15
16     def requires(self):
17         return {
18             "datasets": TaskA(),
19             "corrections": TaskE(),
20         }
21
22     def output(self):
23         luigi.target.FileSystemTarget("/some/path/somewhere.pickle")
24
25     def run(self):
26         from coffea import magical_do_it_all_function, load, save
27
28         output = magical_do_it_all_function(
29             datasets = self.requires()["datasets"].get_all_paths(),
30             corrections = load(self.inputs()["corrections"].path)
31         )
32
33         save(self.output().path, output)
34
```



- extremely useful - but
  - cumbersome with large graphs
- somewhat addressed with: `luigi.util.inherits`, ...
  - not sufficient for dynamic/complex graphs

Example: Task A need a new parameter


```
5 class TaskA(luigi.Task):
6     year = luigi.IntParameter(default=2016, description="year of data taking")
7
8     ...
9
10 @luigi.util.inherits(TaskA)
11 class TaskB(luigi.Task):
12     def requires(self):
13         return {
14             "datasets": TaskA(year=self.year),
15             "corrections": TaskE(),
16         }
```



 is great!

But it could be better.

## Up next:

- a rough task concept
  - based on several ideas
- addressing pain points:
  - first & second hand
- no implementation (or name) yet
  - only mockup API usage
- no intent to reinvent the wheel:
  - make use of , etc
- feedback & discussion wanted!
  - what bugs you about luigi?

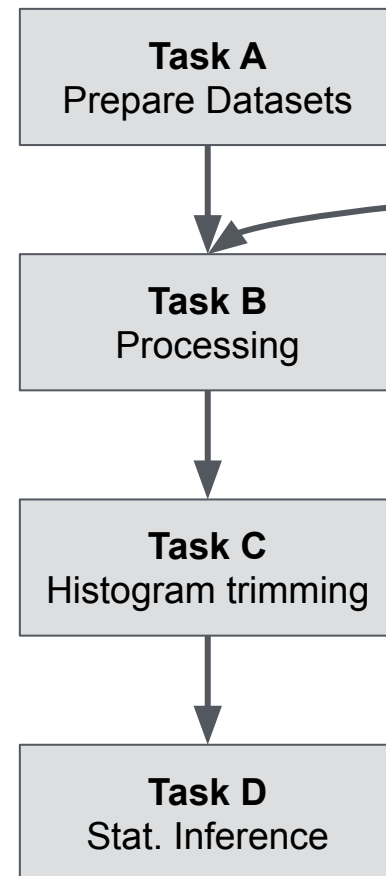
```
287
288 class FlatDatacardsTask(FlatDatacardsCommand, DHI, DatacardProvider):
289     pass
290
291
292 class MultiDatacardsTask(MultiDatacardsCommand, DHI, DatacardProvider):
293     pass
294
...
578 class FitDiagnosticsCombinedPostprocessed(FitDiagnosticsPostprocessed, FlatDatacardsTask):
579     def requires(self):
580         return FitDiagnosticsCombined.req(self)
581
582
583 class FitDiagnosticsFrequentistToysPostprocessed(
584     FitDiagnosticsFrequentistToysMixin, FitDiagnosticsPostprocessed, FlatDatacardsTask, PoolMap
585 ):
586
1182 class InferenceCombined(
1183     ModelMixin,
1184     RecipeMixin,
1185     CreateIssueMixin,
1186     CombinationTask,
1187     law WrapperTask,
1188 ):
1189     dhi_command = DHITask.dhi_command
1190     no_poll = DHI.no_poll
```

- it get worse! i.e. with:
  - dynamic dependencies
  - multiple dependencies
  - same Task w/ different parameters

- this needs addressing
- but, unlikely to work within luigi's design

## Idea:

1. dynamically “inherit” parameters of dependencies
2. dependencies are parameters
3. recurse through “parameterized objects”



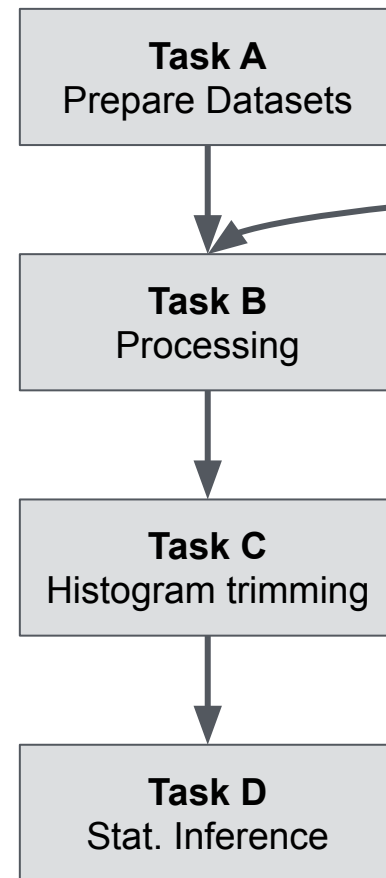


- it get worse! i.e. with:
  - dynamic dependencies
  - multiple dependencies
  - same Task w/ different parameters

- this needs addressing
- but, unlikely to work within luigi's design

## Idea:

1. dynamically “inherit” parameters of dependencies
2. dependencies are parameters
3. recurse through “parameterized objects”

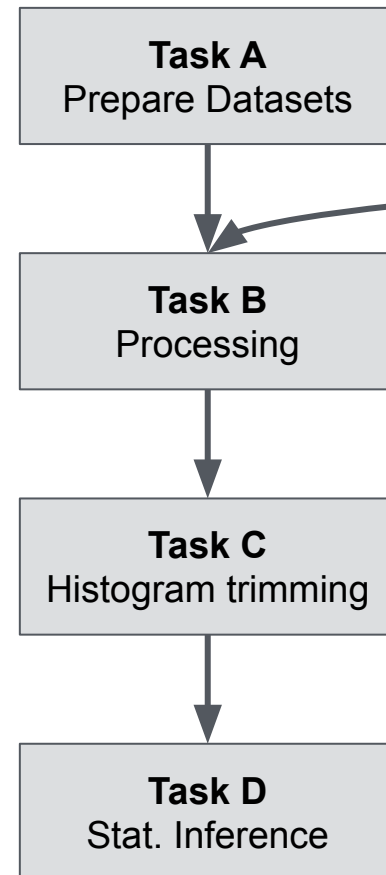


use path-like schema to address parameter, for both:

1. **value access**
2. defining values

Example:

```
1 from not_luigi import Task, Parameter
2
3 class TaskA(Task):
4     year: int = Parameter(default=2012)
5
6 class TaskB(Task):
7     datasets: TaskA = Parameter()
8
9     def run(self):
10         self.datasets.year
```

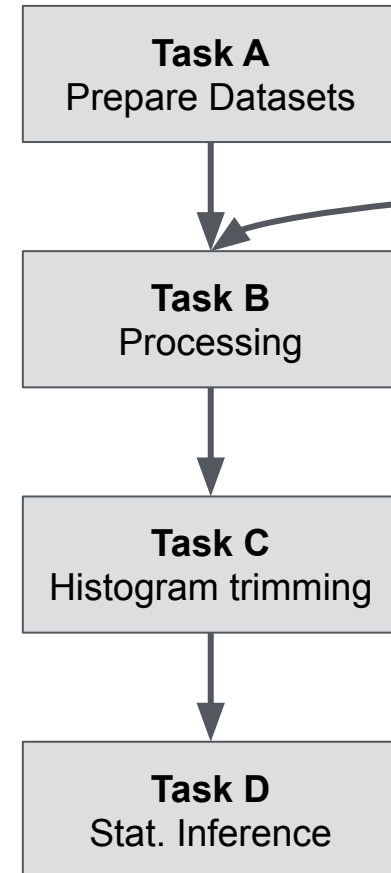


use path-like schema to address parameter, for both:

1. value access
2. **defining values (python)**

Example:

```
13 TaskC({
14     # increasing specificity
15     "year": 2024,
16     ("TaskA", "year"): 2024,
17     (TaskA, "year"): 2024,
18     TaskA.year: 2024,
19     # or even more explicit
20     (TaskB, TaskA.year): 2024,
21     (TaskB.datasets, TaskA.year): 2024,
22     "TaskA.datasets.year": 2024,
23 })
```

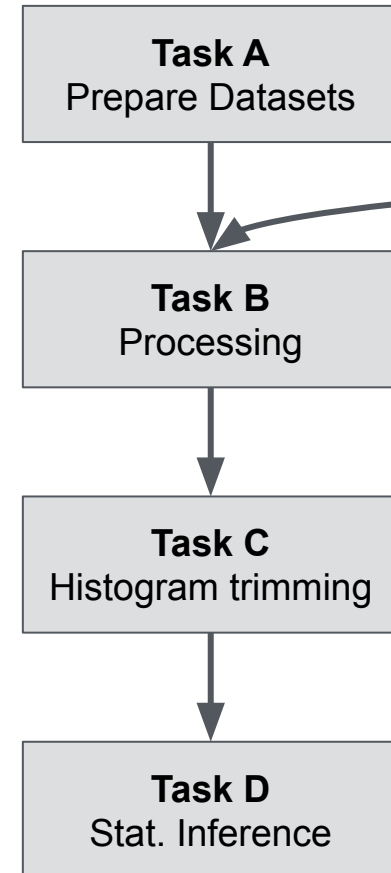


use path-like schema to address parameter, for both:

1. value access
2. **defining values (shell)**

Example:

```
1 not_luigi TaskC --year=2024
2 not_luigi TaskC --TaskA.year=2024
3 not_luigi TaskC --TaskA.year=2024
4 not_luigi TaskC --TaskB:TaskA.year=2024
5 not_luigi TaskC --TaskB.datasets:TaskA.year=2024
6 not_luigi TaskC --TaskB.datasets.year=2024
```



## Task Collection:

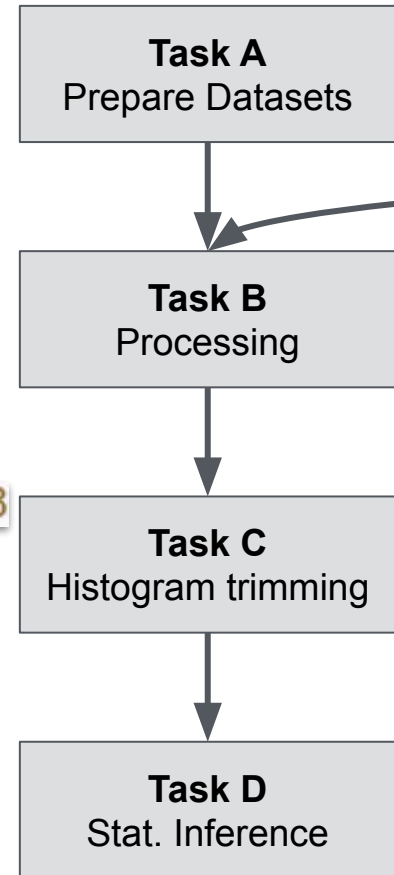
- instances of task with different parameters
- in luigi: need to write dedicated tasks “Wrappers”
- should be automatic feature

## Example:

- python: `TaskA(year=[2016, 2017, 2018])`
- shell: `not_luigi TaskA --year=2016,2017,2018`

## Details:

- can occur in deeply nested dependencies
- propagated via Exceptions (similar to StopIteration)
- caught & handled whenever desired (reduction Task)



- replace a task (implementation) with another one
- deeply within the dependency graph
- without, need to rewrite everything

Examples:

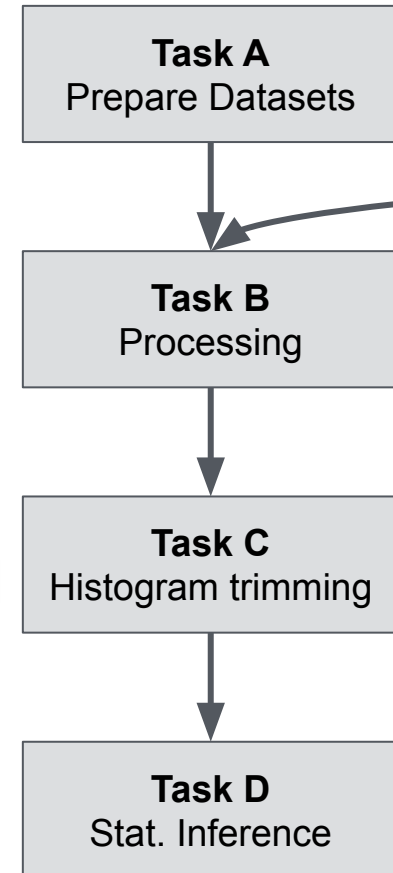
- python: 

```
21 class TaskBcustom(TaskB):  
22     |     ...  
23  
24 TaskZ({TaskB: TaskBcustom})
```
- shell: 

```
not_luigi TaskZ --TaskB=TaskBcustom
```

improves reusability, both way:

- reuse a graph, plugging in new tasks
- transplanting task/graph portions into own graph



- parameter carry arbitrary auxiliary data  
e.g. sampling information for hyperopt.

## special parameter types:

- const parameter:
  - e.g. indicate code changes
- dynamic parameter:
  - like a property, return anything
- actions:
  - methods to be called
  - in order given (even multiple times)

```
not_luigi TaskG --cleanup
```

```
5 class TaskG(Task):
6     learning_rate: float = Parameter(
7         default=1e-3,
8         range=[1e-2, 1e-5],
9         sample="loguniform",
10    )
11
12    code_version = ConstParameter(1)
13
14    @Parameter.dynamic
15    def maybe_important(self):
16        if self.other_parameter:
17            return Task0(self)
18
19    @Action
20    def cleanup(self):
21        ... # do someting
22
23    ...
```



is nice - should be be better - but not compatible with core design  
→ time to think about something new

## Design Ideas

- Parameterized Objects
- Dependencies via parameters
- Path-link parameter addressing
- Automatic Task Collections
- Task Substitution
- Special Parameters:  
Const~, Dynamic~, Action

```
1 from not_luigi import Task, Parameter
2
3 class TaskA(Task):
4     year: int = Parameter(default=2012)
5
6 class TaskB(Task):
7     datasets: TaskA = Parameter()
8
9     def run(self):
10        self.datasets.year
```

```
13 TaskC({
14     # increasing specificity
15     "year": 2024,
16     ("TaskA", "year"): 2024,
17     (TaskA, "year"): 2024,
18     TaskA.year: 2024,
19     # or even more explicit
20     (TaskB, TaskA.year): 2024,
21     (TaskB.datasets, TaskA.year): 2024,
22     "TaskA.datasets.year": 2024,
23 })
```