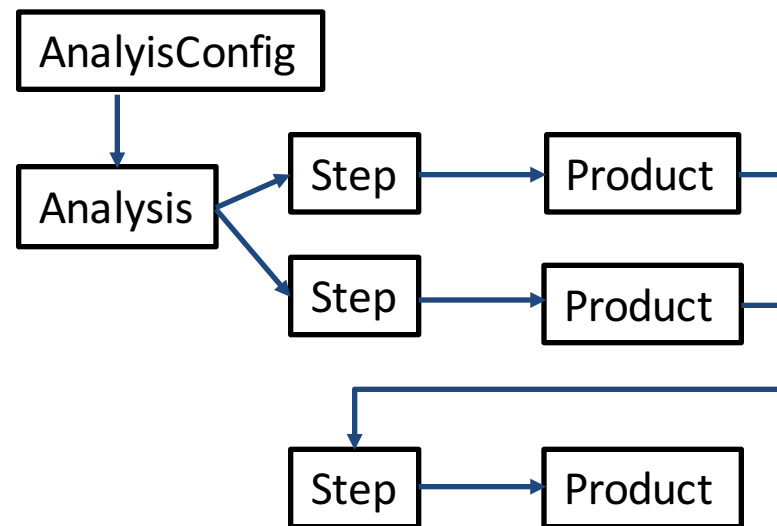


Run each analysis functions :

```
1 analysis = Analysis(config)
2 analysis.get_observations()
3 analysis.get_datasets()
4 analysis.get_excess_map()
```

Problems :

- Only one hard-coded workflow
- Everything into one class
- Config not fully validated so can fail late
- Not parallel



Run the full workflow :

```
1 analysis = Analysis(config)
2 analysis.run()
```

or from command line:

```
1 gammapy analysis run --filename config.yaml --out .
```

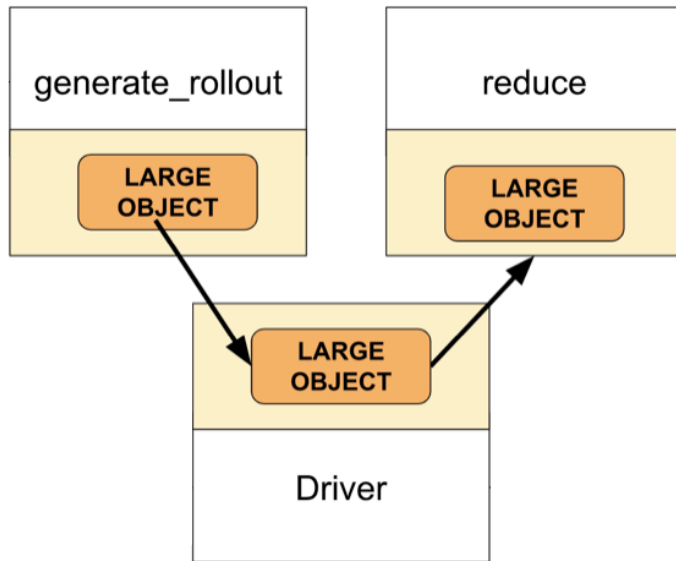
Advantages:

- More flexible and extensible
- Config not fully validated
- can check steps config on init
- Parallel

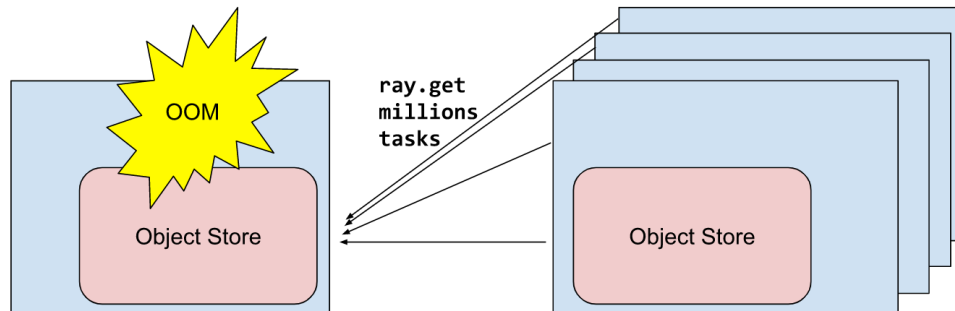
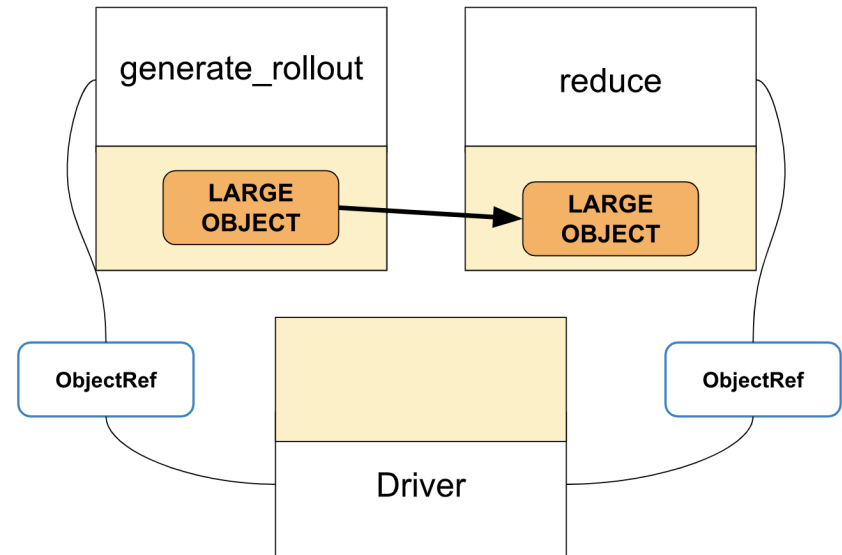
Exchange products between steps



Anti-pattern



Better Method



- Avoid to store all products in the Analysis class but pass only references
- Exchange required products directly between steps

Example



```
class Products(MutableSequence):
    def __init__(self, products=None):
        if products is None:
            products = []
        self._products = products

    @property
    def names(self):
        return [p.name for p in self._products]

    @property
    def data(self):
        return [p.data for p in self._products]

    def index(self, key):
        if isinstance(key, (int, slice)):
            return key
        elif isinstance(key, str):
            return self.names.index(key)

    def select(self, tag=None):
        #could select by step_name, tag, name)
        mask = np.ones(len(self), dtype=bool)
        for idx, p in enumerate(self._products):
            mask[idx] = p.tag==tag
        return self[mask]

    def get(self):
        ind = np.where([p.needs_update for p in self._products])[0]
        refs_to_run = {data for data in np.array(self.data)[ind]} #unique objects
        results = ray.get(list(refs_to_run))
        products = [item for sublist in results for item in sublist] #flatten
        for product in products:
            ind = np.where(np.array(self.names) == product.name)[0]
            for idx in ind:
                self[int(idx)] = product
```

```
class Product:
    def __init__(self, tag=None, step_name=None, data=None):
        self.tag = tag
        self.step_name = step_name
        self.name = make_name() #unique id from make_name
        self.data = data

    @property
    def needs_update(self):
        return isinstance(self.data, ray.ObjectRef)
```

```

class AnalysisStep:
    def __init__(self, name=None, parallel=True):
        self.parallel = parallel
        self.name = make_name(name)
        self._data = Products([])

    @property
    def data(self):
        return self._data

    def set_data(self, data, extend=True):
        if isinstance(data, Product):
            data = Products([data])
        selection = self._select_data(data)
        if extend:
            self._data.extend(selection)
        else:
            self._data = selection

    def _select_data(self, data):
        selection = Products([])
        for req in self.required_products:
            selection.extend(data.select(**req))
        return selection

    def run(self, data, parallel=None, extend=True):
        self.set_data(data, extend=extend)
        # copied or reference if data is on the object store ?
        if parallel is None:
            parallel = self.parallel
        if parallel:
            run_parallel = ray.remote(self.__class__._run)
            ref = run_parallel.remote(self)
            for p in self.products:
                p.data = ref
        else:
            self._run()
        return self.products

```

```

class SumStep(AnalysisStep):
    tag = "sum"
    required_products = [{"tag": "value"}]

    def __init__(self, **kwargs):
        super().__init__(**kwargs)
        self.products = Products([Product(tag="value",
                                            step_name=self.name,
                                            data=None),
                                  Product(tag="log",
                                            step_name=self.name,
                                            data=None)
                                  ])

    def _sum(self):
        return np.sum(self.data.data)

    def _run(self):
        self.data.get() #wait other remote and set results on self.data
        time.sleep(1)
        self.products[0].data = self._sum()
        self.products[1].data = "done"
        return self.products

```

```

step = SumStep()
step.run(data=Product(tag="value", data=1), parallel=True)
step.run(data=Product(tag="value", data=2), parallel=True)
step.run(data=step.products)
step.products.get()

```

Run step 1 & 2 in parallel, wait completion, run step 3

Only init no heavy computation

Config example



```
1 AnalysisConfig
2   general:
3     log: {level: info, filename: null, filemode: null, format: null, datefmt: null}
4     outdir: .
5     n_jobs: 1
6     datasets_file: null
7     models_file: null
8   observations:
9     datastore: $GAMMAPY_DATA/hess-dl3-dr1
10    obs_ids: []
11    obs_file: null
12    obs_cone: {frame: icrs, lon: 83.633 deg, lat: 22.014 deg, radius: 5.0 deg}
13    obs_time: {start: null, stop: null}
14    required_irf: [aeff, edisp, psf, bkg]
15  datasets:
16    type: 3d
17    stack: true
18    geom:
19      wcs:
20        skydir: {frame: icrs, lon: 83.633 deg, lat: 22.014 deg}
21        binsize: 0.02 deg
22        width: {width: 2.0 deg, height: 2.0 deg}
23        binsize_irf: 0.2 deg
24        selection: {offset_max: 2.5 deg}
25      axes:
26        energy: {min: 1.0 TeV, max: 10.0 TeV, nbins: 10}
27        energy_true: {min: 0.5 TeV, max: 20.0 TeV, nbins: 20}
28    map_selection: [counts, exposure, background, psf, edisp]
29    background:
30      method: fov_background
31      exclusion: $GAMMAPY_DATA/joint-crab/exclusion/exclusion_mask_crab.fits.gz
32      parameters: {method: scale}
33    safe_mask:
34      methods: [aeff-default]
35      parameters: {}
36    on_region: {frame: null, lon: null, lat: null, radius: null}
37    containment_correction: true
38  excess_map:
39    correlation_radius: 0.1 deg
40    parameters: {}
41    energy_edges: {min: null, max: null, nbins: null}
```

```
42 steps:
43 - name: hess1_data-reduction
44   tag: data-reduction
45   products: datasets_hess1
46   config:
47     observations:
48       datastore: data/hess1
49 - name: hesslu_data-reduction
50   tag: data-reduction
51   products: datasets_hesslu
52   config:
53     observations:
54       datastore: data/hesslu
55 - name: hess_excess_map
56   tag: excess_map
57   required_products: [datasets_hess1, datasets_hesslu]
58   products: hess_flux_maps
```

- PR 3852 : github.com/gammapy/gammapy/pull/3852
move analysis.get_...() methods to AnalysisStep classes
add analysis.run() method and cli equivalent
- Add AnalysisProduct and AnalysisProducts classes
- Add PARALLEL_BACKEND = ... (where ?)
and add a mechanism to switch from multiprocessing to ray.util.multiprocessing
(in TSmapEstimator and DatasetsMaker)
docs.ray.io/en/latest/ray-more-libs/multiprocessing.html
- Introduce parallel support in AnalysisStep and AnalysisProducts (optional by default)
- Refactor AnalysisStep to define all maker/Estimators on init (and run them on run)
- Provide new configs to run various analysis workflows
- Add tutorial notebook