

# Pylightnix

Sergey Mironov

February 2021

## Abstract

Pylightnix is a minimalistic Python library for managing filesystem-based immutable data objects, inspired by Purely Functional Software Deployment Model thesis by Eelco Dolstra and the Nix package manager.

The library may be thought as of low-level API for creating caching wrappers for a subset of Python functions. In particular, Pylightnix allows user to

- Prepare (or, in our terms, **instantiate**) the computation plan aimed at creating a tree of linked immutable stage objects, stored in the filesystem.
- Implement (**realize**) the prepared computation plan, access the resulting artifacts. Pylightnix is able to handle **non-deterministic** results of the computation. As one example, it is possible to define a stage depending on best top-10 instances (in a user-defined sense) of prerequisite stages.
- Handle changes in the computation plan, re-use the existing artifacts whenever possible.
- Gain full control over all aspects of the cached data including the garbage-collection.

# Contents

<b>1</b>	<b>Quick start</b>	<b>2</b>
1.1	The problem . . . . .	2
1.2	Stages basics . . . . .	2
1.2.1	Parameter stage . . . . .	2
1.2.2	Annealing stage . . . . .	3
1.2.3	Plotting stage . . . . .	4
1.2.4	Running the experiment . . . . .	4
1.3	Nested stages . . . . .	4
1.4	Collaborative work . . . . .	5
1.5	Synchronization . . . . .	5

## 1 Quick start

We illustrate Pylightnix principles by showing how to use it for making a mock data science task. We'll follow the SciPy anneal example but imagine that we want to save some intermediate results and also that there are several people working on this task.

We start by importing the required Python modules.

---

```
1 import numpy as np
2 import scipy.optimize as o
3 import matplotlib.pyplot as plt
4
5 from pylightnix import *
```

---

### 1.1 The problem

The annealing example begins with defining a complex parametric function **f** for what we need to find a minimum point as close as possible.

---

```
1 def f(z, *params):
2     x, y = z
3     a, b, c, d, e, f, g, h, i, j, k, l, scale = params
4     return (a * x**2 + b * x * y + c * y**2 + d*x + e*y + f) + \
5           (-g*np.exp(-((x-h)**2 + (y-i)**2) / scale)) + \
6           (-j*np.exp(-((x-k)**2 + (y-l)**2) / scale))
```

---

### 1.2 Stages basics

Lets assume that we are to get the **f**'s parameters from elsewhere and that we want to get the following results: (a) the starting parameters themselves (b) the annealing result and statistics and (c) visual plot of the results.

#### 1.2.1 Parameter stage

The Pylightnix Stage is a Python function that registers a user-defined action in the Manager and returns a DRef reference stating that the action is accepted.

---

```

1 def stage_params(m:Manager)->DRef:
2     def _config():
3         name = 'params'
4         out = [selfref, "params.npy"]
5         return locals()
6     def _make(b:Build):
7         np.save(mklens(b).out.syspath, (2, 3, 7, 8, 9, 10, 44, -1, 2, 26, 1,
8             ↪ -2, 0.5))
9     return mkdrv(m, mkconfig(_config()), match_only(), build_wrapper(_make))

```

---

The actual registration is done by the `mkdrv` function which accepts the manager and the following three arguments (a) the action Config object; (b) the Matcher callback function; (c) the action callback accepting the Build task description. In this example we are pretending to receive input parameters from a thirdparty. We could just as well have downloaded them from the Internet using the `fetchurl` pre-defined stage.

A few notes on the above code:

- We define configuraiton dictionary by reading a local variables of a helper function named `_config()`. The resulting dictionary should match the Config requirements. As often in Pylighnix, we use a standalone `mkconfig` construcor rather than create the object directly.
- By specifying a `selfref` output path we promise to produce an artifact with this name. Pylighnix will throw an error if it does not find such an artifact after the build action completes.
- When inside the build action, we refer to the pylighnix config items by using the `mklens` swiss knife helper. `mklens` returns the Lens object which has an overloaded dot `"."` operator. The middle part of the lens expression encodes the path in the configuration and the last word `"syspath"` encodes the format of the result.
- We pass a `match_only()` matcher in the third argument of `mkdrv`. By this we expect our action to return a determenistic result. Pylighnix will throw an error if it is not the case.

### 1.2.2 Annealing stage

---

```

1 def stage_anneal(m:Manager, ref_params:DRef)->DRef:
2     def _config():
3         name = 'anneal'
4         nonlocal ref_params
5         trace_xs = [selfref, 'tracex.npy']
6         trace_fs = [selfref, 'tracef.npy']
7         out = [selfref, 'result.npy']
8         return locals()
9     def _make(b:Build):
10        params = np.load(mklens(b).ref_params.out.syspath)
11        xs = []; fs = []
12        def _trace(x,f,ctx):
13            nonlocal xs,fs
14            xs.append(x.tolist())

```

```

15     fs.append(f)
16     res = o.dual_annealing(f, [[-10,10],[-10,10]],
17                           x0=[2.,2.],args=params,
18                           maxiter=500, callback=_trace)
19     np.save(mklens(b).trace_xs.syspath, np.array(xs))
20     np.save(mklens(b).trace_fs.syspath, np.array(fs))
21     np.save(mklens(b).out.syspath, res['x'])
22     return mkdrv(m, mkconfig(_config()), match_only(), build_wrapper(_make))

```

---

### 1.2.3 Plotting stage

---

```

1  def stage_plot(m:Manager, ref_anneal:DRef)->DRef:
2      def _config():
3          name = 'plot'
4          nonlocal ref_anneal
5          out = [selfref, 'plot.png']
6          return locals()
7      def _make(b:Build):
8          xs=np.load(mklens(b).ref_anneal.trace_xs.syspath)
9          fs=np.load(mklens(b).ref_anneal.trace_fs.syspath)
10         res=np.load(mklens(b).ref_anneal.out.syspath)
11         plt.figure()
12         plt.title(f"Min {fs[-1]}, found at {res}")
13         plt.plot(range(len(fs)),fs)
14         plt.grid(True)
15         plt.savefig(mklens(b).out.syspath)
16     return mkdrv(m, mkconfig(_config()), match_latest(),
17                 ↪ build_wrapper(_make))

```

---

### 1.2.4 Running the experiment

---

```

1  ds=instantiate_inplace(stage_params)
2  cl=instantiate_inplace(stage_anneal,ds)
3  vis=instantiate_inplace(stage_plot,cl)
4  rref=realize_inplace(vis)
5  print(rref)

```

---

```
rref:cb252162ae8a0fe2c90fe48193f851f8-ff1dd5c155ee3faa445572e3246f0566-plot
```

## 1.3 Nested stages

In the previous sections we implicitly used a global Manager to register our stages. A more clean approach requires us to define a nested stage covering the whole experiment. This way we postpone the problem of choosing which Manager to use.

---

```

1 def stage_all(m:Manager):
2     ds=stage_params(m)
3     cl=stage_anneal(m,ds)
4     vis=stage_plot(m,cl)
5     return vis

```

---

The nested stage may be instantiated and realized in one go.

---

```

1 rref_all=realize(instantiate(stage_all))
2 assert rref_all==rref
3 print(rref_all)

```

---

```
rref:cb252162ae8a0fe2c90fe48193f851f8-ff1dd5c155ee3faa445572e3246f0566-plot
```

---

## 1.4 Collaborative work

Alice and Bob runs the same experiment using their machines. We model this situation by running Pylighnix with different storage settings.

---

```

1 Sa=mksettings('_storageA')
2 Sb=mksettings('_storageB')
3 fsinit(Sa,remove_existing=True)
4 fsinit(Sb,remove_existing=True)
5 rrefA=realize(instantiate(stage_all, S=Sa))
6 rrefB=realize(instantiate(stage_all, S=Sb))
7 print(rrefA)
8 print(rrefB)

```

---

Our annealing problem is stochastic by nature so the results naturally differ.

## 1.5 Synchronization

Alice sends her results to Bob, so he faces a problem of picking the best realizaion.

---

```

1 print("Bob's storage before the sync:")
2 for rref in allrrefs(S=Sb):
3     print(rref)
4 arch=Path('archive.zip')
5 pack([rrefA], arch, S=Sa)
6 # Alice transfers archive to Bob using her favorite pigeon post service.
7 unpack(arch, S=Sb)
8 print("Bob's storage after the sync:")
9 for rref in allrrefs(S=Sb):
10     print(rref)

```

---

```
Bob's storage before the sync:
rref:6f906ca9ad8d9b128c57ed0df95ab899-7c308bad0e7fa364790b71e6ffab924a-params
rref:71f3be199a7abe8024f3441d008d0f4c-aeeb4bea649e22497b49dfc1ec5ebedf-anneal
rref:3e38836b267a5e3f211d4d66c18c74c9-ff1dd5c155ee3faa445572e3246f0566-plot
Bob's storage after the sync:
rref:6f906ca9ad8d9b128c57ed0df95ab899-7c308bad0e7fa364790b71e6ffab924a-params
rref:b034af3d2708e41b9f24eb1d684f45d0-aeeb4bea649e22497b49dfc1ec5ebedf-anneal
rref:71f3be199a7abe8024f3441d008d0f4c-aeeb4bea649e22497b49dfc1ec5ebedf-anneal
rref:3e38836b267a5e3f211d4d66c18c74c9-ff1dd5c155ee3faa445572e3246f0566-plot
rref:9de94205d5ed12ed9aee73e555c0383-ff1dd5c155ee3faa445572e3246f0566-plot
```

Bob writes a custom matcher that selects the realization which has the best annealing result.

```
1 def match_min(S, rrefs:List[RRef])>List[RRef]:
2     avail=[np.load(mklens(rref,S=S).trace_fs.syspath)[-1] for rref in rrefs]
3     best=sorted(zip(avail,rrefs))[0]
4     if best[1] in allrrefs(Sa):
5         print(f"Picking Alice ({best[0]}) out of {avail}")
6     else:
7         print(f"Picking Bob ({best[0]}) out of {avail}")
8     return [best[1]]
9
10 def stage_all2(m:Manager):
11     ds=stage_params(m)
12     cl=redefine(stage_anneal, new_matcher=match_min)(m,ds)
13     vis=stage_plot(m,cl)
14     return vis
15
16 rref_best=realize(instantiate(stage_all2,S=Sb))
17 print(rref_best)
```

```
Picking Bob (-3.408582123541778) out of [-3.302737056408237, -3.408582123541778]
rref:3e38836b267a5e3f211d4d66c18c74c9-ff1dd5c155ee3faa445572e3246f0566-plot
```