

The CERN ML Frameworks

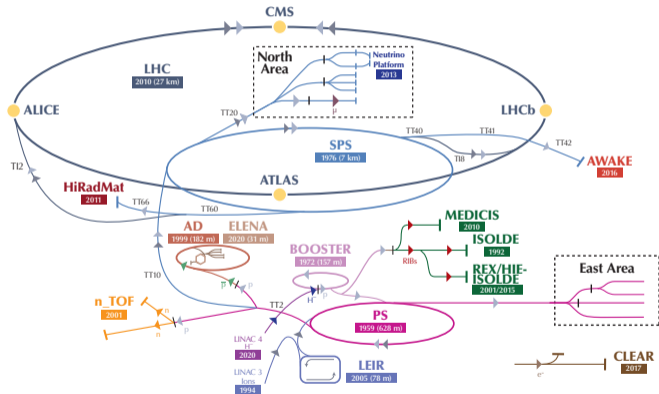
How to get ML into the Control Room

J.-B. de Martel ■ N. Madysa R. Gorbonosov V. Kain

3rd ICFA Beam Dynamics Mini-Workshop,
on Machine Learning Applications for Particle Accelerators,
3 November 2022

- many accelerators, extremely diverse
- lots of low-level problems already well automated
- **but:** many high-level problems still solved manually
- better turnaround time and beam quality **necessary** to reach target integrated luminosity

The CERN accelerator complex Complexe des accélérateurs du CERN



Lots of infrastructure already in place!

JAPC: uniform communication protocol with physical and virtual devices

LSA: access to the settings database and history

UCAP: *Function as a Service* development framework ([Link](#))

WRAP: web app to create browser-based monitoring dashboards ([Link](#))

Acc-Py: Python development and deployment framework

Two projects are more specialized on optimization and ML:

MLP: platform for ML model versioning, storage and deployment

GeOFF: Framework for operational use of RL and numerical optimization

Machine Learning Platform (MLP)

Volatile world of ML:

- code needs to run once
- bleeding-edge technology
- cloud services and custom tools
- maintainability not the first concern

Reliable world of accelerator controls:

- need to run reliably 24/7/365:
reproducibility, robustness, traceability
- highly reliable, battle-tested tools
- constraints of the technical network:
no internet access, restricted tooling,
security precautions
- minimize maintenance effort through
standardization and unification

Machine Learning Platform (MLP)

Objectives:

- common approach to storage, versioning, deployment and usage of models
- remove infrastructural concerns from model lifecycle
- fulfill needs specific to accelerator controls:
 - ▶ reliability
 - ▶ traceability
 - ▶ security
 - ▶ standardization
- stay out of the user's way:
 - ▶ minimize constraints on model developer's workflow
 - ▶ don't restrict choice of tools

Machine Learning Platform (MLP)

Concepts:

Model Type

- contains the model **code**
- fully fledged Python package
- lives as a Git repository on our servers
- different versions are different development snapshots

Model Parameters

- contains only the **training results**
- depends on the ML framework
- often an opaque binary file
- different versions are different training snapshots

Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package

```
1 $ tree
2 my-model-package/
3 |— my_model_package/
4 |   |— __init__.py
5 |   |— mlp-models.toml
6 |— .gitlab-ci.yaml
7 |— pyproject.toml
8 |— setup.cfg
```

Made easy with `acc-py init`
and `acc-py init-ci!`

Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package
- 2 Declare your model

```
1 # mlp-models.toml
2 [[model]]
3 name = "my_model_package:MyModel"
4 standalone = false
```



Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package
- 2 Declare your model
- 3 Implement a wrapper class

```
1 # __init__.py
2 from mlp_model_api import INPUTS, OUTPUTS
3 from mlp_tensorflow_model \
4     import TensorflowModel
5 class MyModel(TensorflowModel):
6     def __init__(self):
7         model = tf.keras.Model(...)
8         super().__init__(model)
9     def predict(self, inputs):
10        outputs = self.model(\
11            inputs[INPUTS])
12        return {OUTPUTS: outputs}
```



Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package
- 2 Declare your model
- 3 Implement a wrapper class
- 4 optional: other frameworks

```
1 # __init__.py
2 class MyModel(TensorflowModel):
3     ...
4     def export_parameters(
5         self, parameters_target: Path
6     ) -> None:
7         ...
8     def load_parameters(
9         self, parameters_src Path
10    ) -> None:
11        ...
```

Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package
- 2 Declare your model
- 3 Implement a wrapper class
- 4 optional: other frameworks
- 5 Push a Git tag to CI

```
1 $ git commit --all -m "New version."  
2 $ git push origin main  
3 $ git tag v1.0.0  
4 $ git push --tags origin
```

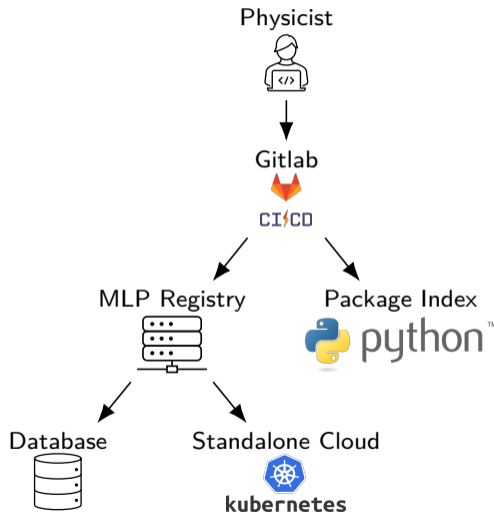


Machine Learning Platform (MLP)

Developing a model:

- 1 Set up your package
- 2 Declare your model
- 3 Implement a wrapper class
- 4 optional: other frameworks
- 5 **Push a Git tag to CI**

... and CI takes care of the rest!



Machine Learning Platform (MLP)

Publishing training results to the MLP:

```
1 # train_my_model.py
2 from my_model_package import MyModel
3 from mlp_client import AUTO, Client, Profile
4
5 mlp_model = MyModel()
6 train(mlp_model.model)
7 client = Client(Profile.PRO)
8 client.publish_model_parameters_version(
9     mlp_model, name="my_model.default", version=AUTO
10 )
```



Machine Learning Platform (MLP)

Embedding an MLP model in a Python application:

```
1 # operational_application.py
2 from my_model_package import MyModel
3 from mlp_client import AUTO, Client, Profile
4 from mlp_model_api import INPUTS, OUTPUTS
5
6 client = Client(Profile.PRO)
7 model = client.create_model(
8     MyModel, name="my_model.default", version=AUTO
9 )
10 response = model.predict({INPUTS: get_inputs()})
11 show_results(response[OUTPUTS])
```

Machine Learning Platform (MLP)

Connecting a Python application to a stand-alone model in the cloud:

```
1 # operational_application.py
2 from mlp_client import AUTO, Client, Profile
3 from mlp_model_api import INPUTS, OUTPUTS
4
5 client = Client(Profile.PRO)
6 model = client.create_standalone_model(
7     "my_model_package:MyModel",
8     name="my_model.default",
9     version=AUTO,
10 )
11 # Communicates via HTTP!
12 response = model.predict({INPUTS: get_inputs()})
13 show_results(response[OUTPUTS])
```

Machine Learning Platform (MLP)

Connecting a Java application to a stand-alone model in the cloud:

```
1 // OperationalApplication.java
2 public static void main() {
3     StandaloneModel<MyRequest, MyResponse> model =
4         new MlpClient(Environment.PRO)
5             .standaloneModelBuilder(MyRequest.class, MyResponse.class)
6             .modelType("my_model_package:MyModel")
7             .modelParameters("my_model.default")
8             .build();
9     MyRequest request = new MyRequest(getInputs());
10    MyResponse response = model.predict(request);
11    showResults(response.getOutput());
12 }
```

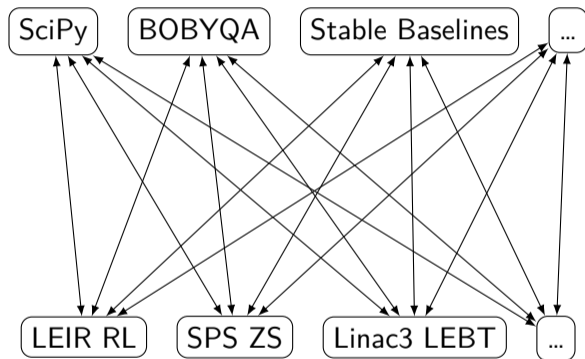

Machine Learning Platform (MLP)

Conclusion:

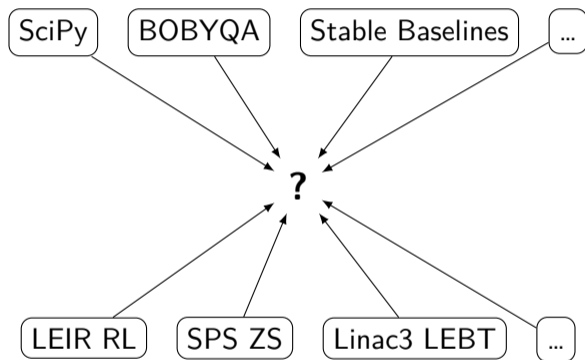
- MLP is ready for production*
- we separate model R&D from its operational use
- model developers can ignore infrastructural concerns
- we follow software engineering best practices
⇒ reliable and maintainable

**standalone model deployment in beta*

- many different optimizers/APIs
- many different optimization problems
- each problem involves complex machine communication
- operators don't want to juggle Python scripts!



- many different optimizers/APIs
- many different optimization problems
- each problem involves complex machine communication
- operators don't want to juggle Python scripts!

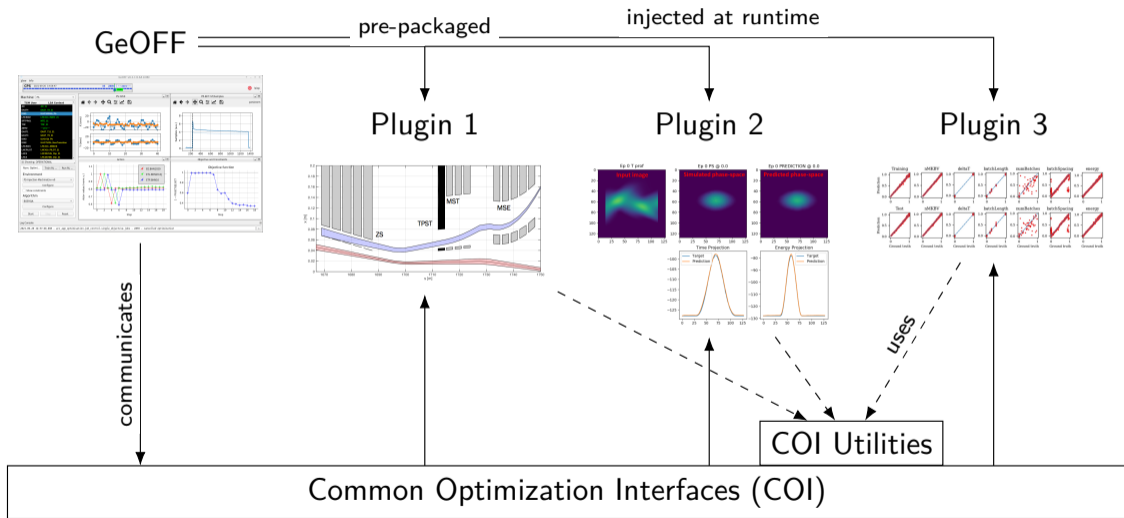


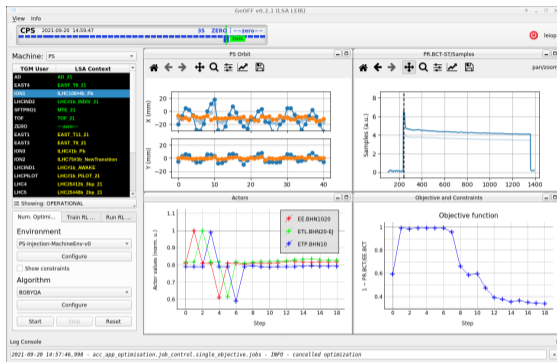
Goals:

- provide ecosystem for accelerator optimization and control
- provide compatibility with as many algorithms as possible
- facilitate the progression
manual tuning → *numerical optimization* → *machine learning*

Guiding principles:

- be agnostic over machine, communication protocol or devices
- minimize boilerplate code that does not solve the problem
- don't make people pay for features they don't use
- always leave an escape hatch open





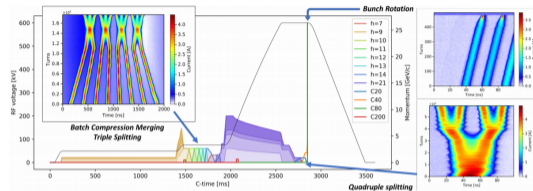
- lists, configures and runs optimization problems
- built-in list of optimizers (will be made user-extensible)
- optimization problems are loaded as plugins **pre-packaged or at runtime**

GeOFF Use Cases

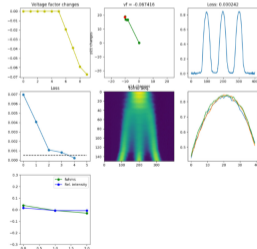
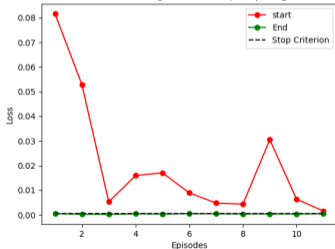
- Linac3: steering of beam transfer line
- Linac4: 2 expert tools
- PSB: operations (WIP) & commissioning
 - ▶ bunch recombination at PSB ejection
 - ▶ resonance compensation
 - ▶ RF optimization
 - ▶ injection to PS
- PS: used during commissioning
 - ▶ resonance compensation
 - ▶ transition gamma jump circuits
 - ▶ septa alignment for slow extraction
- LEIR: used during commissioning
 - ▶ transfer lines (from Linac3, to PS)
 - ▶ injection bumps
 - ▶ phase adjustment of RF cavities
- SPS: expert tool & operations
 - ▶ tune adjustments
 - ▶ septa alignment for slow extraction
 - ▶ spill noise reduction
 - ▶ splitter optimization
 - ▶ injection kicker optimization
 - ▶ crystal shadowing
- used at **almost all** accelerators
 - ▶ ISOLDE: lots of homogeneous devices
⇒ CPS Optimizer
 - ▶ LHC: fast acquisition, high safety req.
⇒ bespoke algorithms
- most often used as expert tool

Use Case: PS RF Manipulations via RL

- **LHC beam production** requires quadruple splitting at 26 GeV/c in PS
 - **RF phase errors** introduce spread in bunch-by-bunch intensity and emittance
- ⇒ RL agent corrects phases for uniform bunches



MD 6887: Start and end Agent Criterion (Comparing all bunches)



- faster than classical optimization due to reuse of experience
- trained on simulation, evaluated on real machine
- episode length $n \in [2, 18]$, $\bar{n} = 8.46$

presented yesterday by J. Wulff



Conclusion

- CERN fosters an ecosystem of independent efforts to improve operational ML usage
- solutions complement and integrate with each other
- **but:** still plenty to do!
- efforts to open-source MLP and GeOFF are underway
 - ▶ upcoming **EURO-LABS** project will drive this effort further

Links:

- <https://gitlab.cern.ch/nmadysa/icfa-tutorial>
- <https://cernml-coi.docs.cern.ch/>
- <https://gitlab.cern.ch/be-op-ml-optimization>
- <https://indico.cern.ch/event/1175862/> (seminar on MLP and GeOFF)

