

Machine Learning Platform: Deploying and Managing Models in the CERN Control System

Jean-Baptiste de Martel - Nico Madysa, Roman Gorbonosov, Verena Kain

01/07/2022 – BE Seminar

Summary

- **System overview**
 - Motivation
 - Objectives
 - Workflow
 - Continuous retraining
- **Standalone deployment prototype**
- **MLP in practice**

ML for accelerator controls - motivation

- **Some accelerators / processes difficult to model with means available in CCC**
-> **no online analytical model available**
 - Space charge dominated dynamics in LINACs
 - Transmission optimization during transition crossing
 - Collimator alignment, septa with many degrees of freedom
 - Multi-turn injection with accumulated intensity diagnostics only
- **Usually solved by trial & error, (semi-)manual scans**
 - Time consuming
 - Depends on experience of operator
- **ML techniques outperform previous techniques and open new doors to automation**
- **But the approach today is ad hoc and heterogeneous**

Finding a compromise

Volatile world of ML

- **Code needs to run once**
- **Bleeding edge technology**
- **Used to own tools and comfort, cloud services**
- **Maintainability is not always the main concern**

Reliable world of accelerator controls

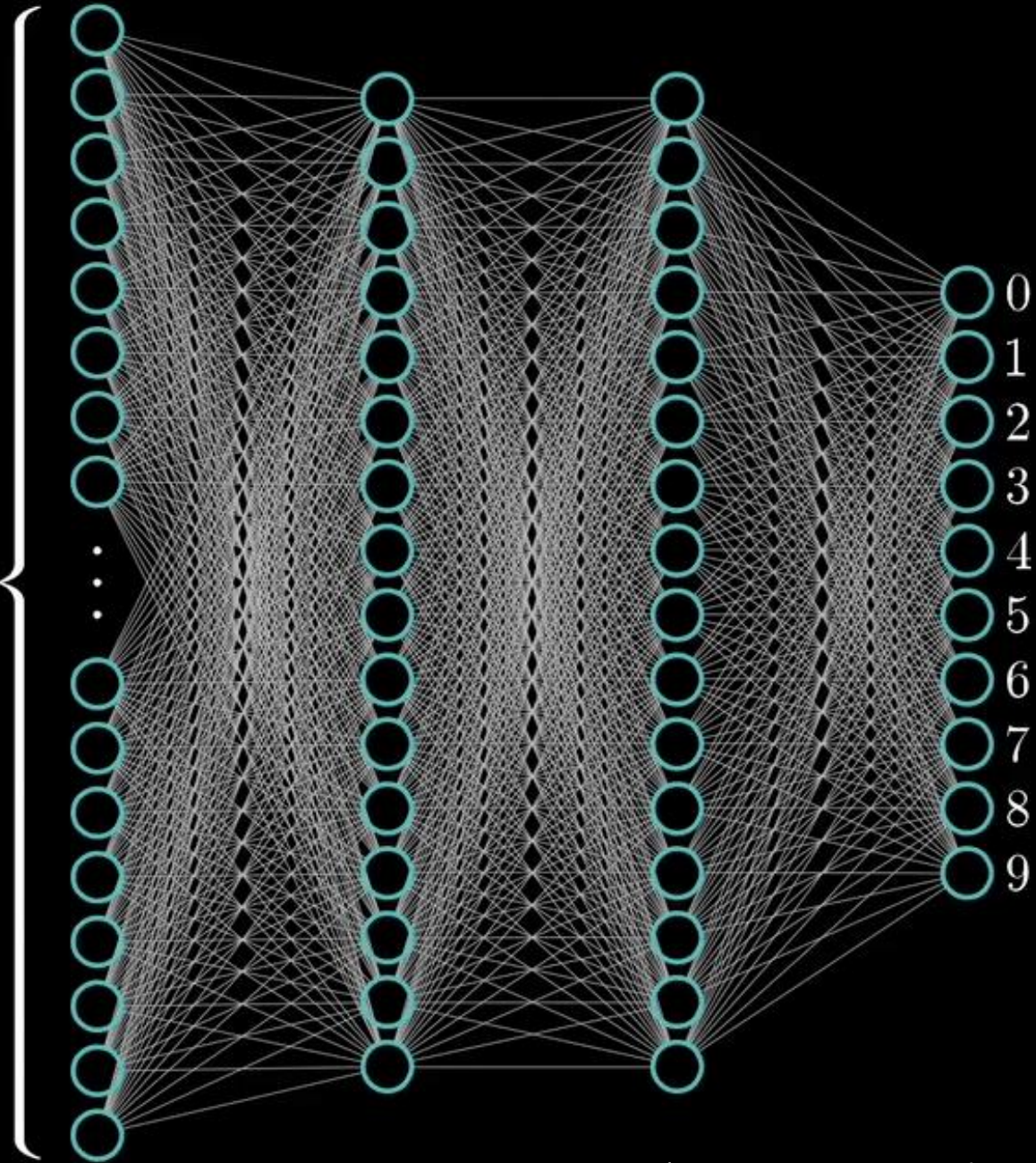
- **Need to run reliably 24/7/365: need reproducibility, robustness, traceability**
- **Use highly reliable, battle-tested tools**
- **Constraints of the technical network: no internet access, restricted tooling, security precautions**
- **Standardize and unify to minimize maintenance**

Objectives

- **Provide a common approach to storage, versioning, deployment and usage of models**
- **Accelerate and simplify the model lifecycle by abstracting infrastructural concerns**
- **Fulfill the specific needs of the accelerator control system**
 - Reliability
 - Traceability
 - Security
 - Standardization
- **Stay out of the user's way**
 - Minimize constraints on model developer's workflow
 - Avoid constraints on choice of tools



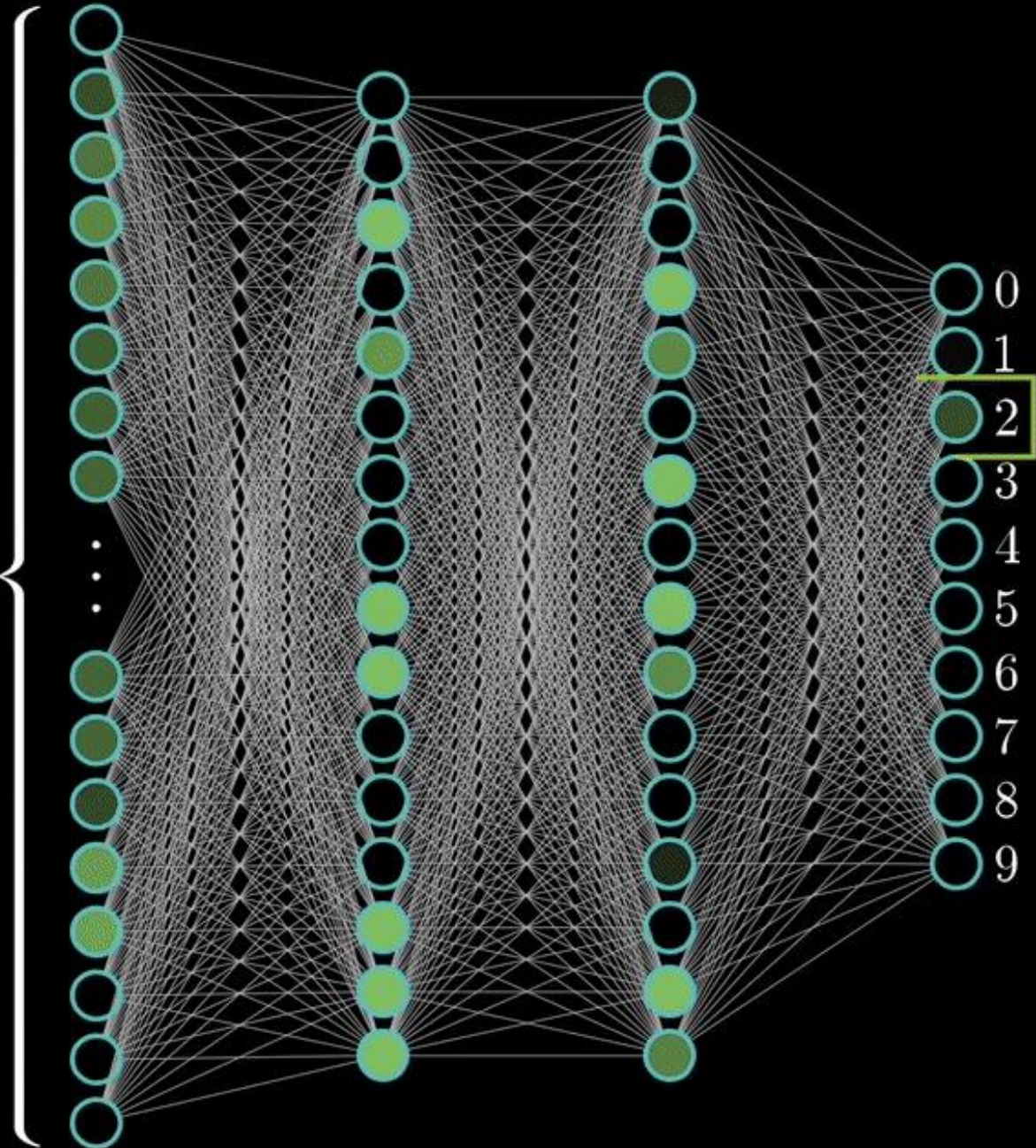
784





Model type:
Layout/architecture
of the neural
network –
i.e., number of
neurons, how they
are connected,
etc...

784



Model parameters:
“Trained weights”
- values assigned
to the neurons and
connections after
training

Model:
Combination of a
model type and
model parameters

ARTIFICIAL INTELLIGENCE

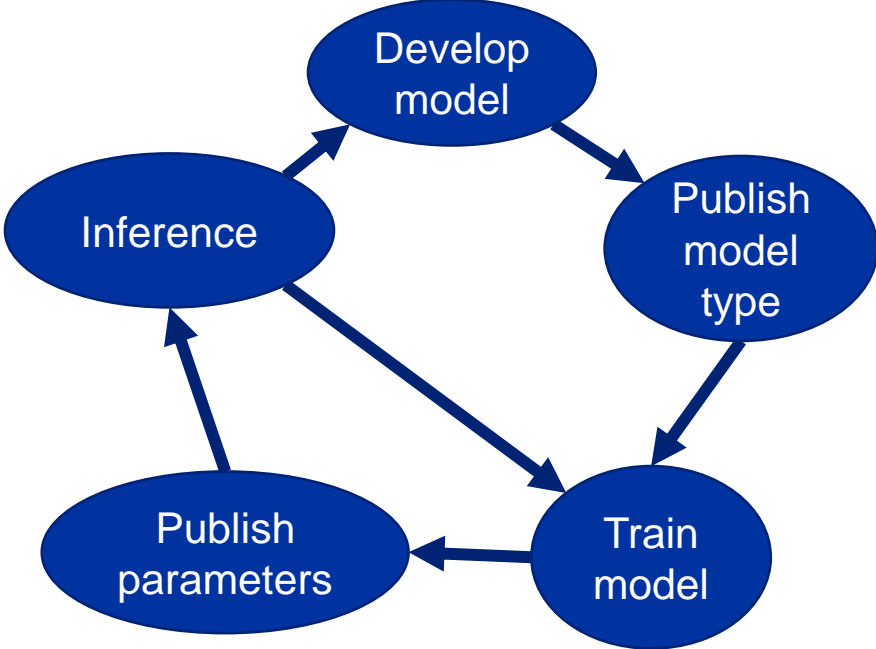
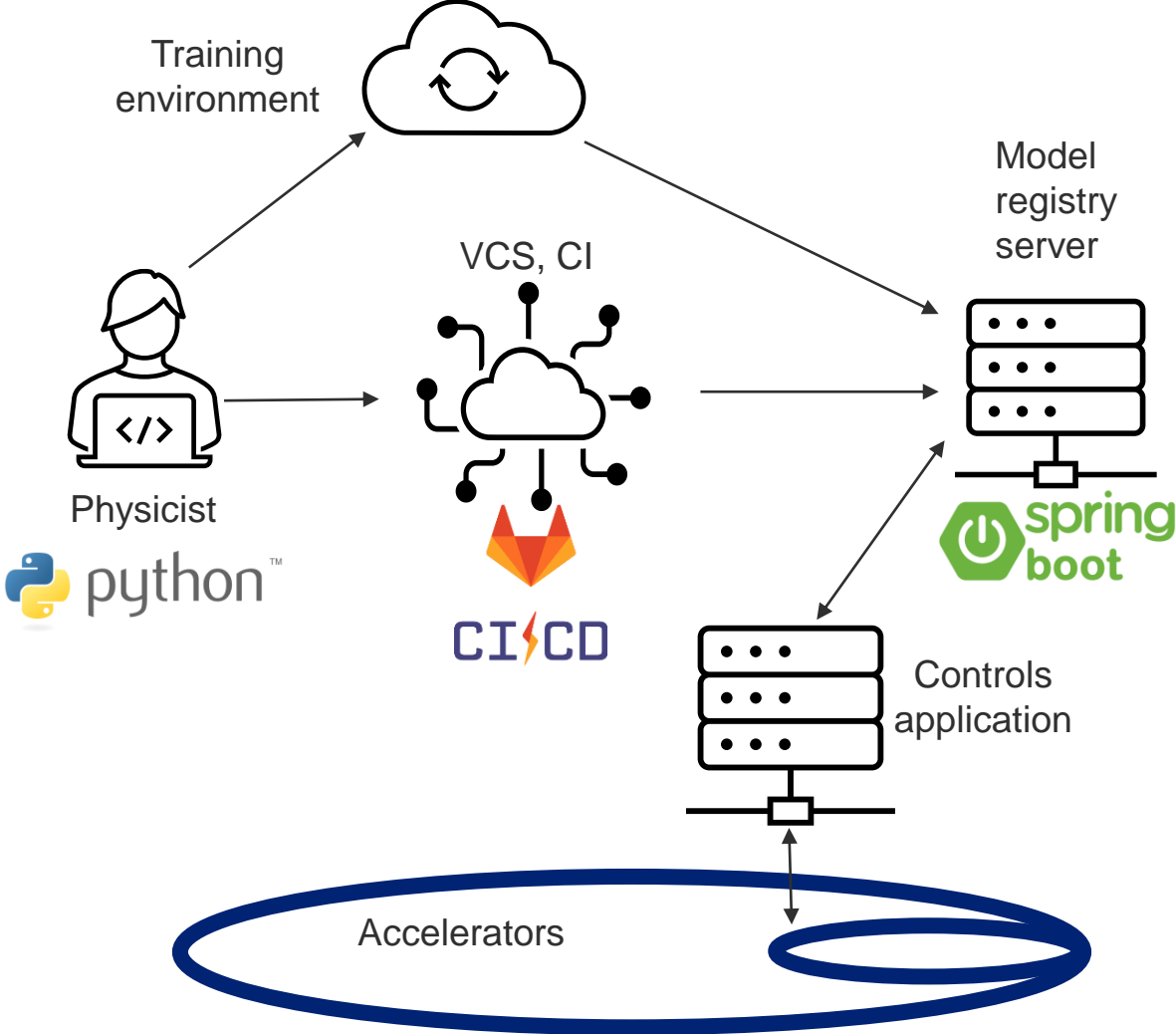
MACHINE LEARNING

NEURAL NETS

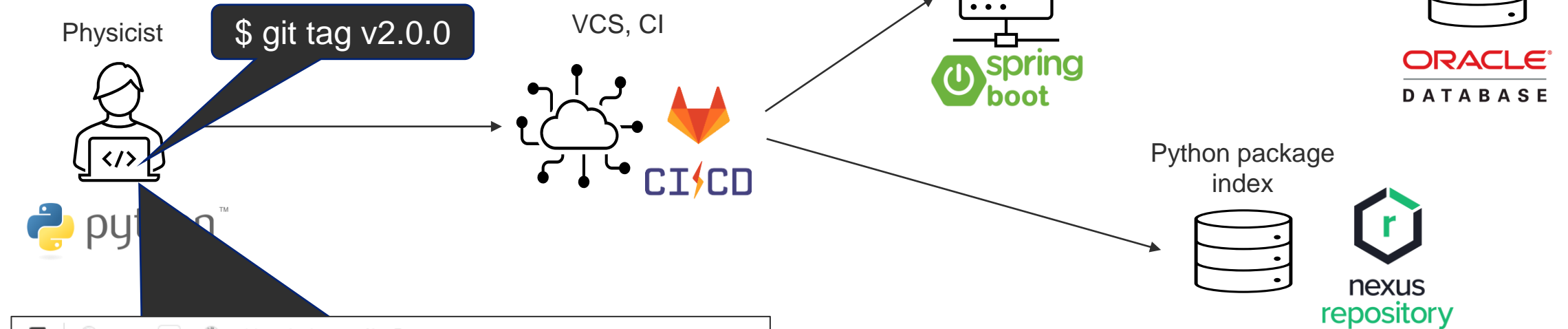
DEEP
LEARNING

dozens of
different ML
methods

Development workflow



Publishing model types



acc-co > models > simple-ann > New Tag

New Tag

Tag name

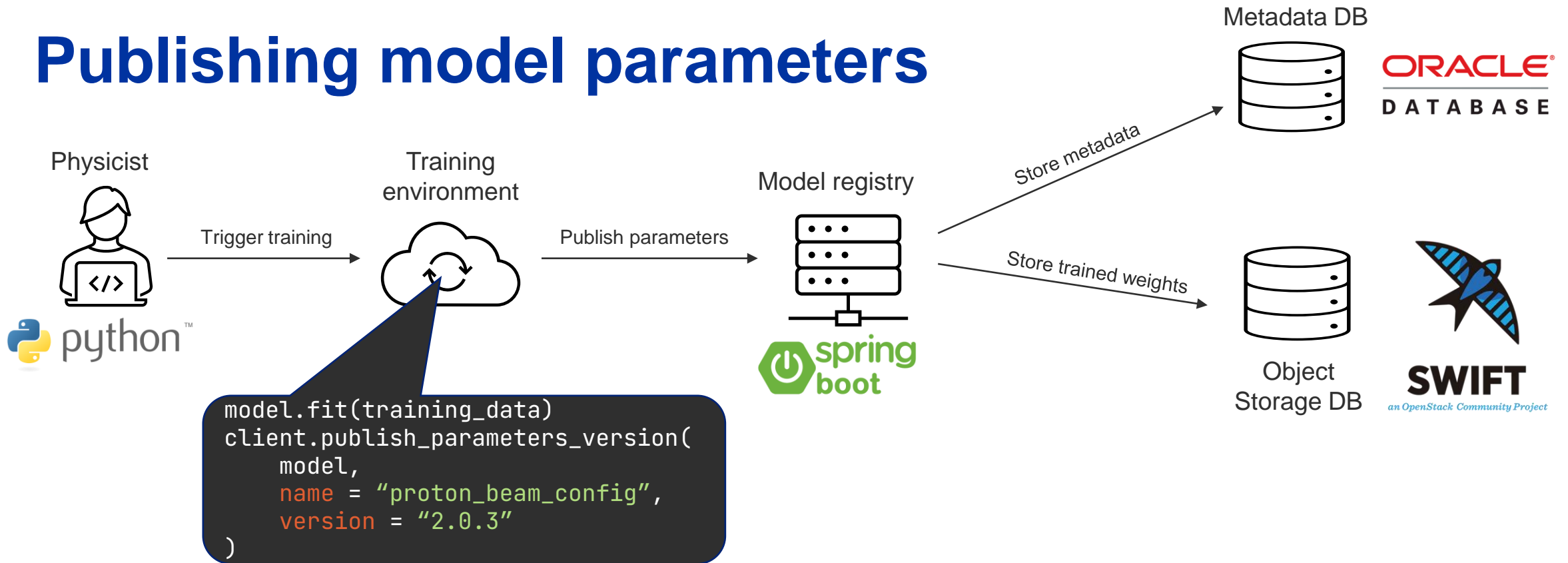
Create from Existing branch name, tag, or commit SHA

Message

Advantages

- **Access control and traceability for model types**
- **Quick & easy, no need to learn new tools, complexity is hidden**

Publishing model parameters



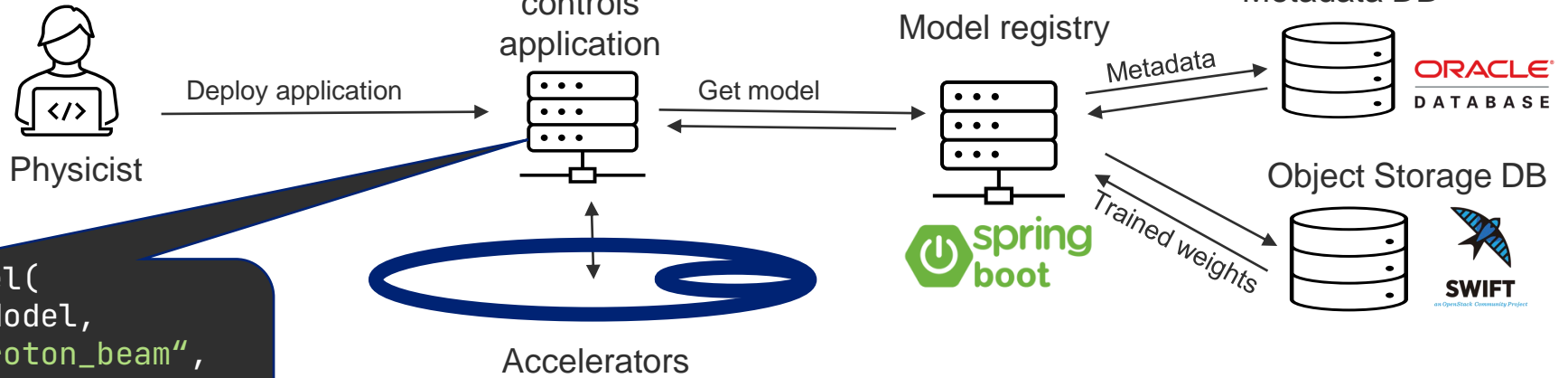
Usage

- Choose parameters name and version
- Use the client library to publish

Advantages

- All parameters stored centrally and reliably
- Compatibility is fully managed

Inference



```
model = client.create_model(  
    model_type = BeamLineModel,  
    model_parameters = "proton_beam",  
    params_version = "2.0.3"  
)  
result = model.predict(input)
```

Usage

- Use the MLP client library to instantiate the model
- Provide model type, parameters name and version

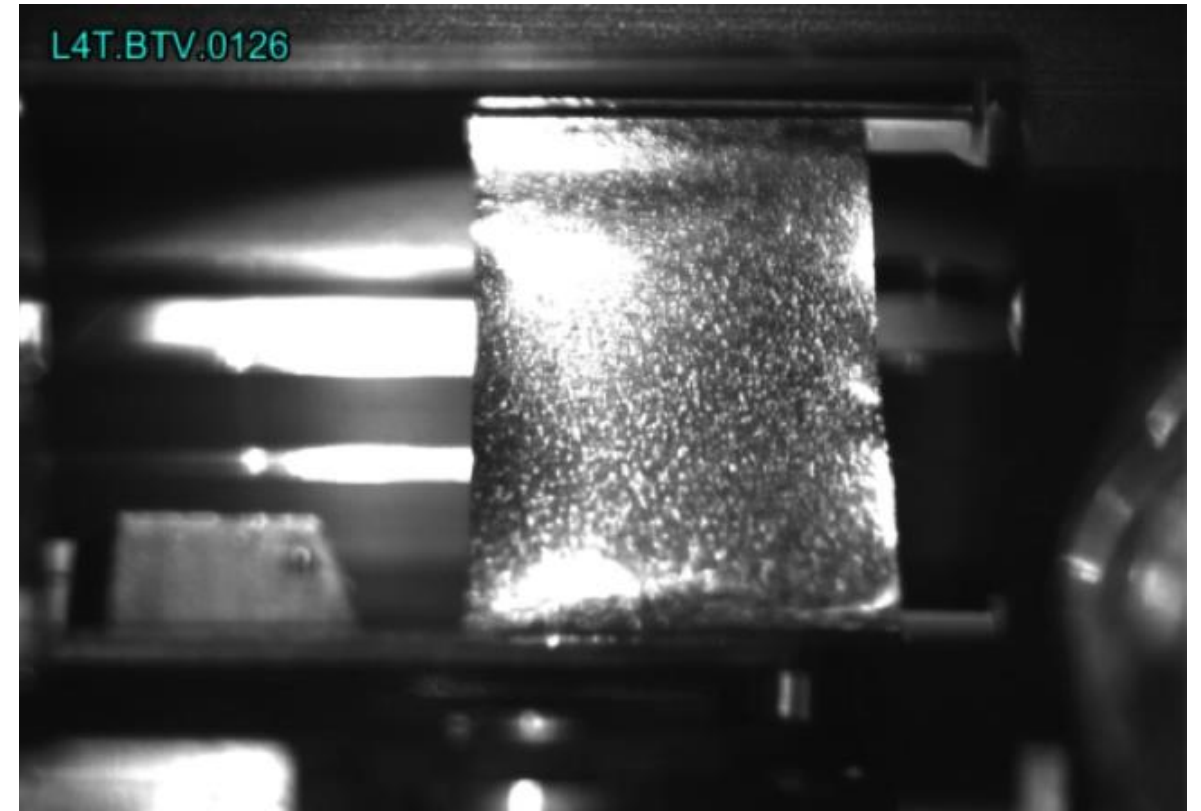
Advantages

- Parameters retrieved and loaded transparently
- Parameter traceability

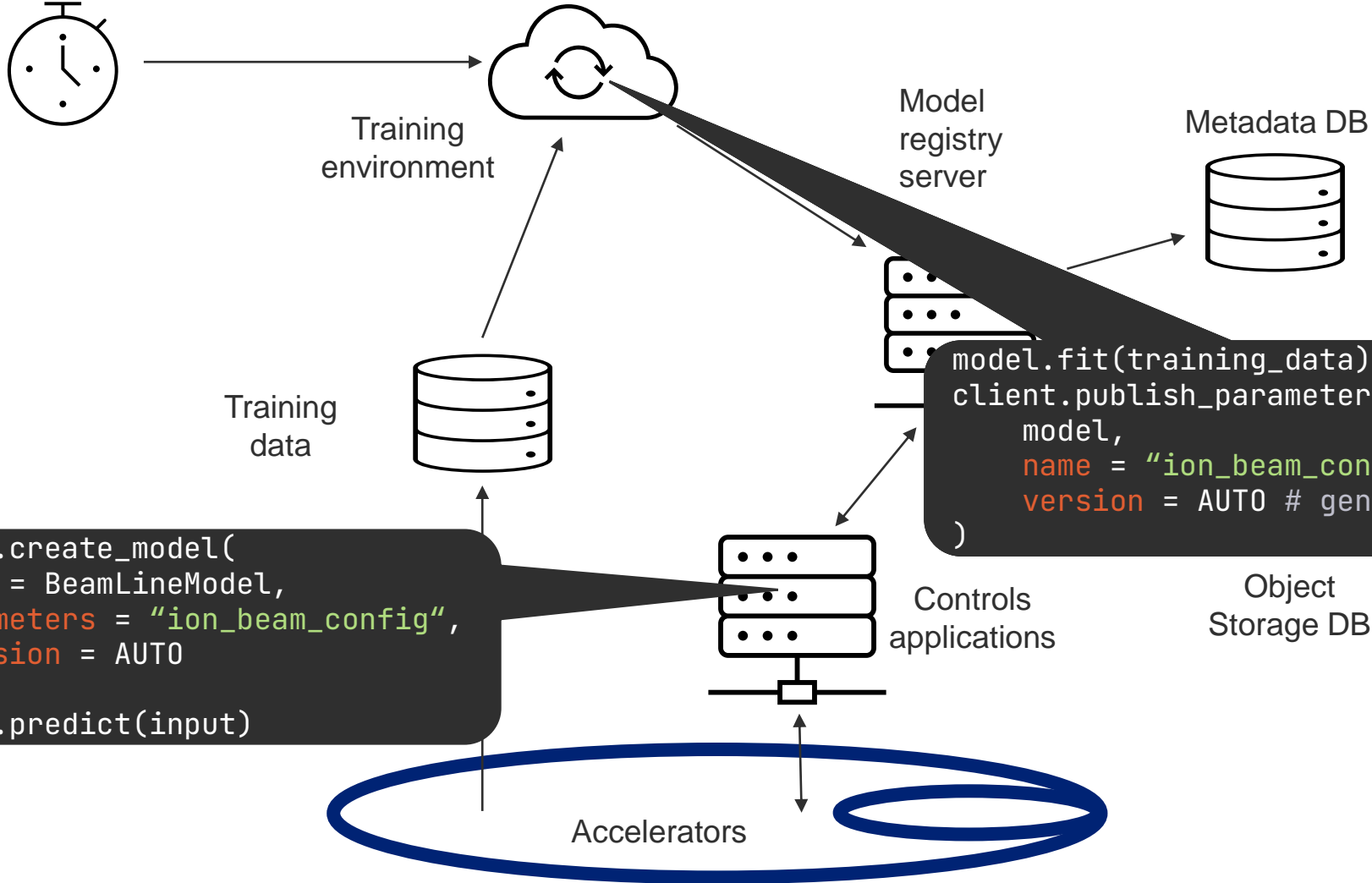
Continuous retraining - motivation

Example: stripper foil degradation

- The stripper foil is an essential component of our linacs
 - It degrades over time and is replaced regularly
 - Beam characteristics vary
 - Machine parameters need to adapt
- > need to re-train model continuously to keep it up to date



Continuous retraining - implementation



```
model = client.create_model(  
    model_type = BeamLineModel,  
    model_parameters = "ion_beam_config",  
    params_version = AUTO  
)  
result = model.predict(input)
```

```
model.fit(training_data)  
client.publish_parameters_version(  
    model,  
    name = "ion_beam_config",  
    version = AUTO # generated  
)
```

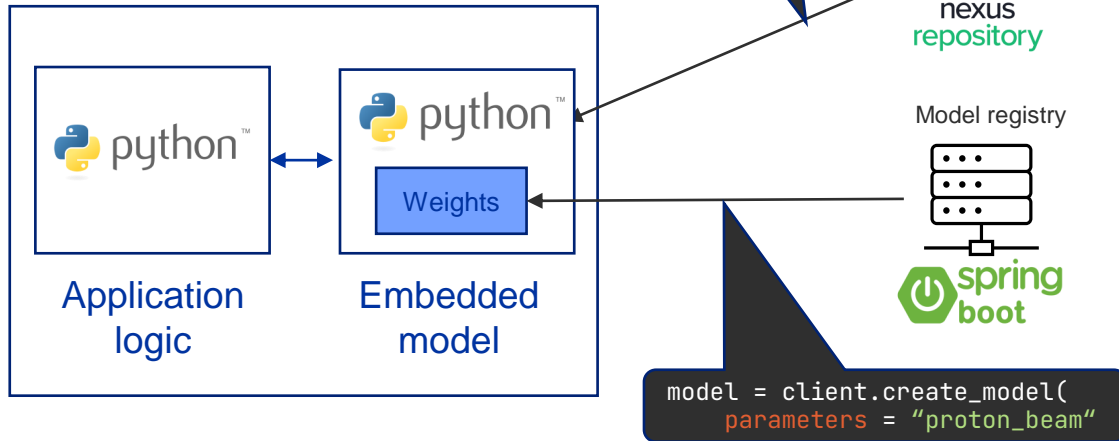
Standalone deployment prototype

Calling models remotely from any language

Inference - Embedded vs standalone

Embedded

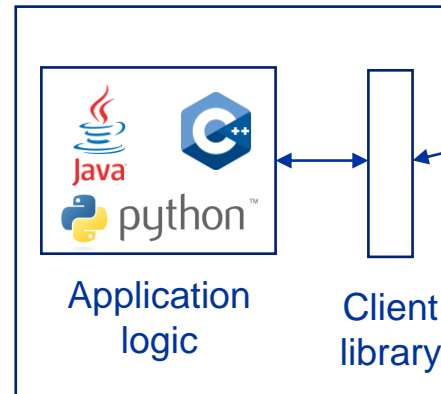
Controls application



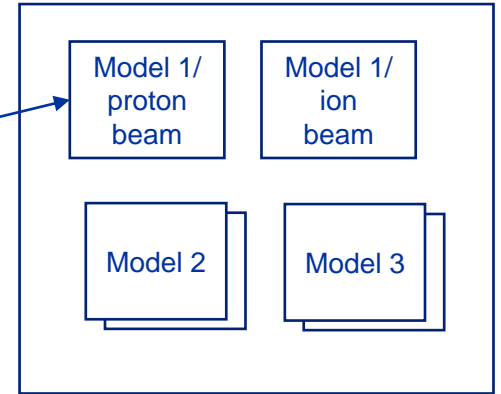
- Python only
- Model type must be installed
- Parameters retrieved then stored locally

Standalone

Controls application

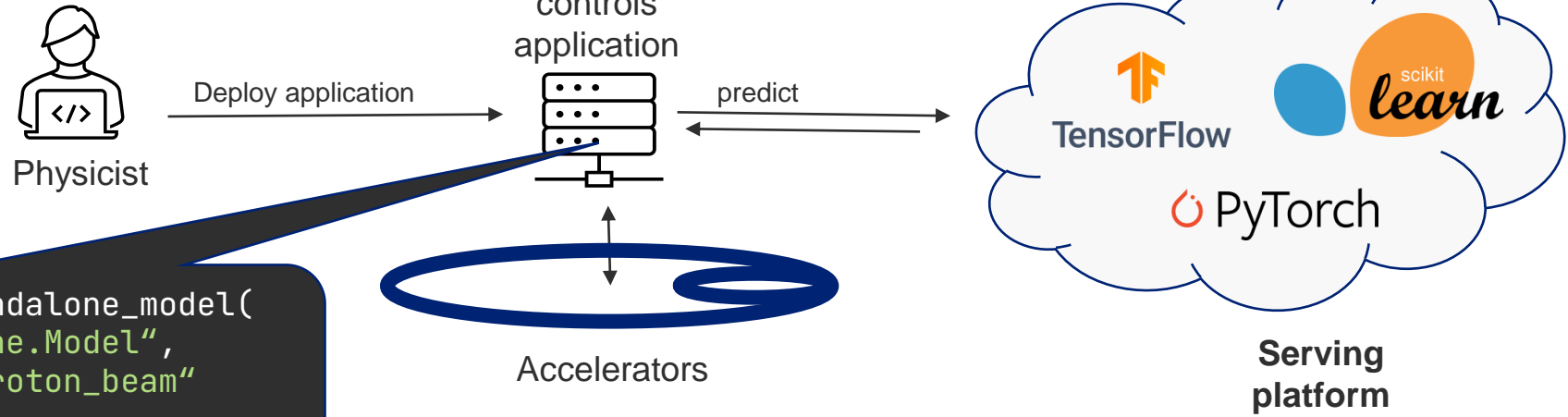


Standalone serving cluster



- Language-agnostic approach
- No local model installation needed
- Everything happens remotely

Inference (Standalone)



```
model = client.create_standalone_model(  
    model_type = "beam_line.Model",  
    model_parameters = "proton_beam"  
)  
result = model.predict(input)
```

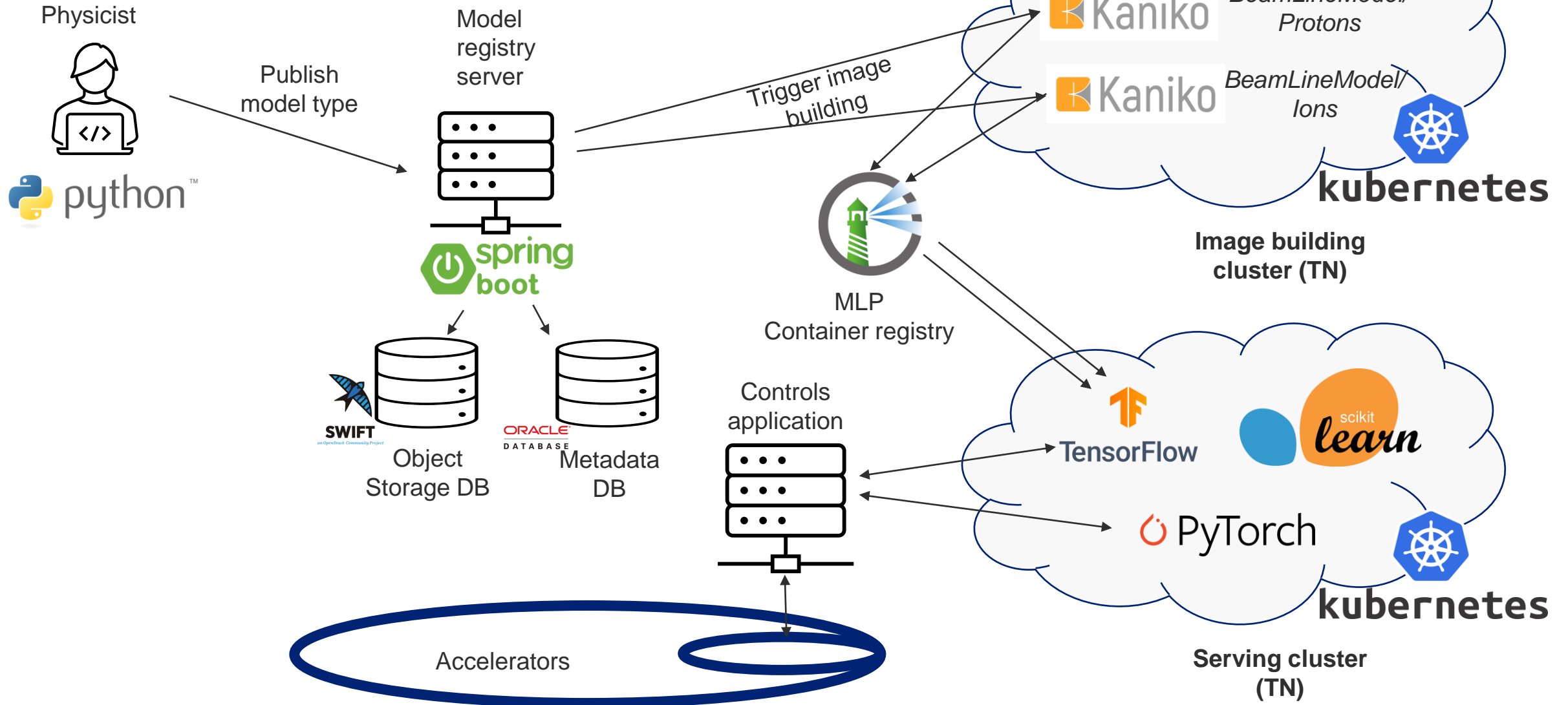
Usage

- Use the MLP client library to instantiate the model
- Provide model type and parameters name

Advantages

- Call models from any language
- Seamless model updates

Behind the scenes



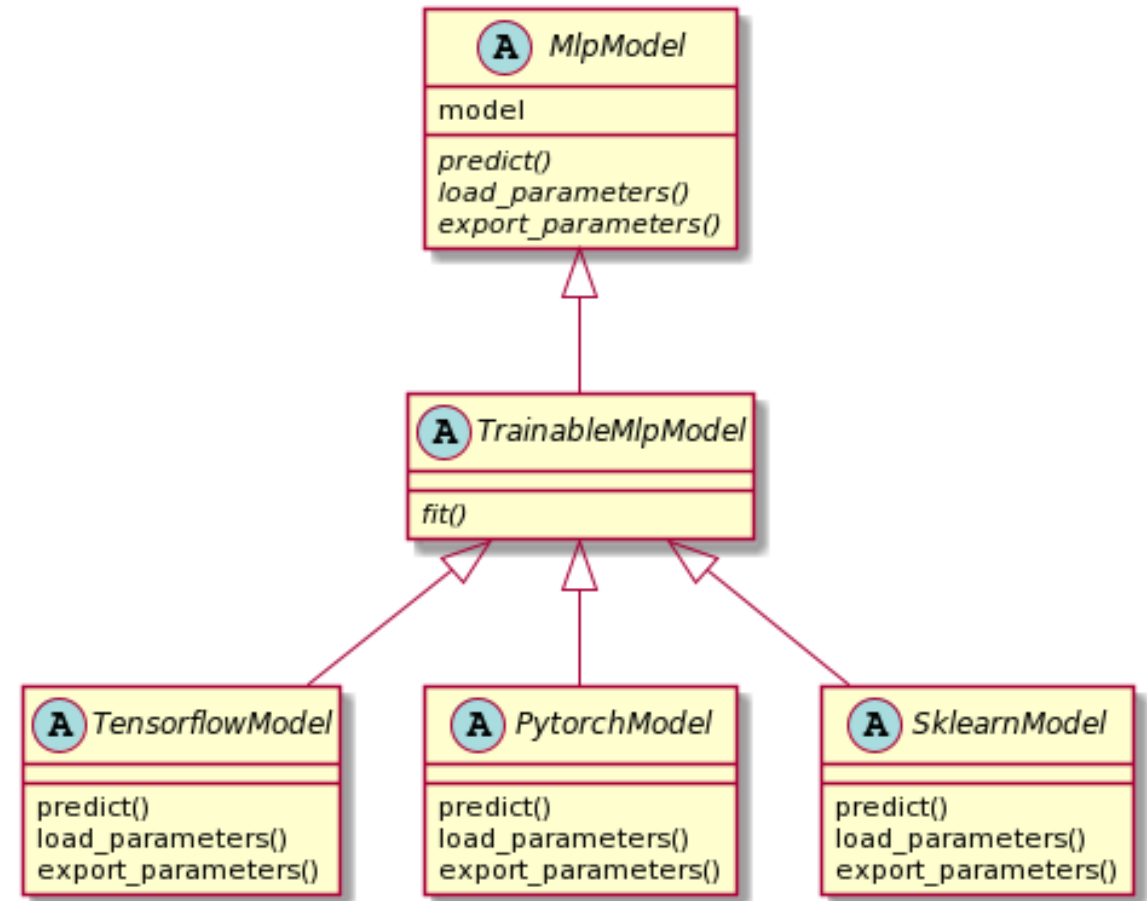
MLP in practice

MLP in Practice

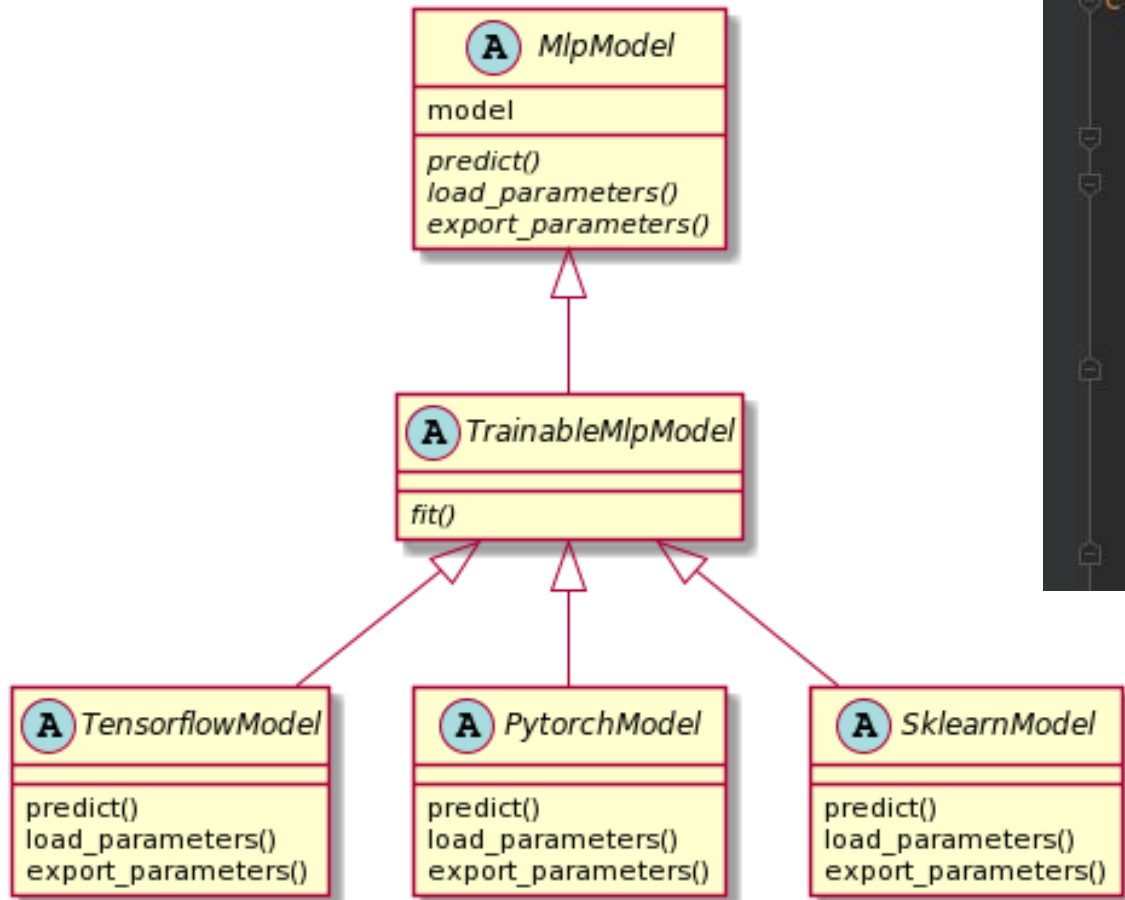
- **Adapting a model**
 - Implementing the model interface
 - Model declaration
 - CI template
- **Using a model**
 - Publishing trained parameters
 - Loading parameters
 - Standalone prediction

MLP in practice: Implementing a model

- **We define a common API for all controls models**
 - shared abstraction layer
- **Interface defines 4 methods:**
 - *Fit* – train the model on the provided data
 - *Export parameters* – extract current values of all model parameters
 - *Load parameters* – configure the model using the provided parameters
 - *Predict* – return a prediction from the input data
- **Default extensible implementations for common frameworks**



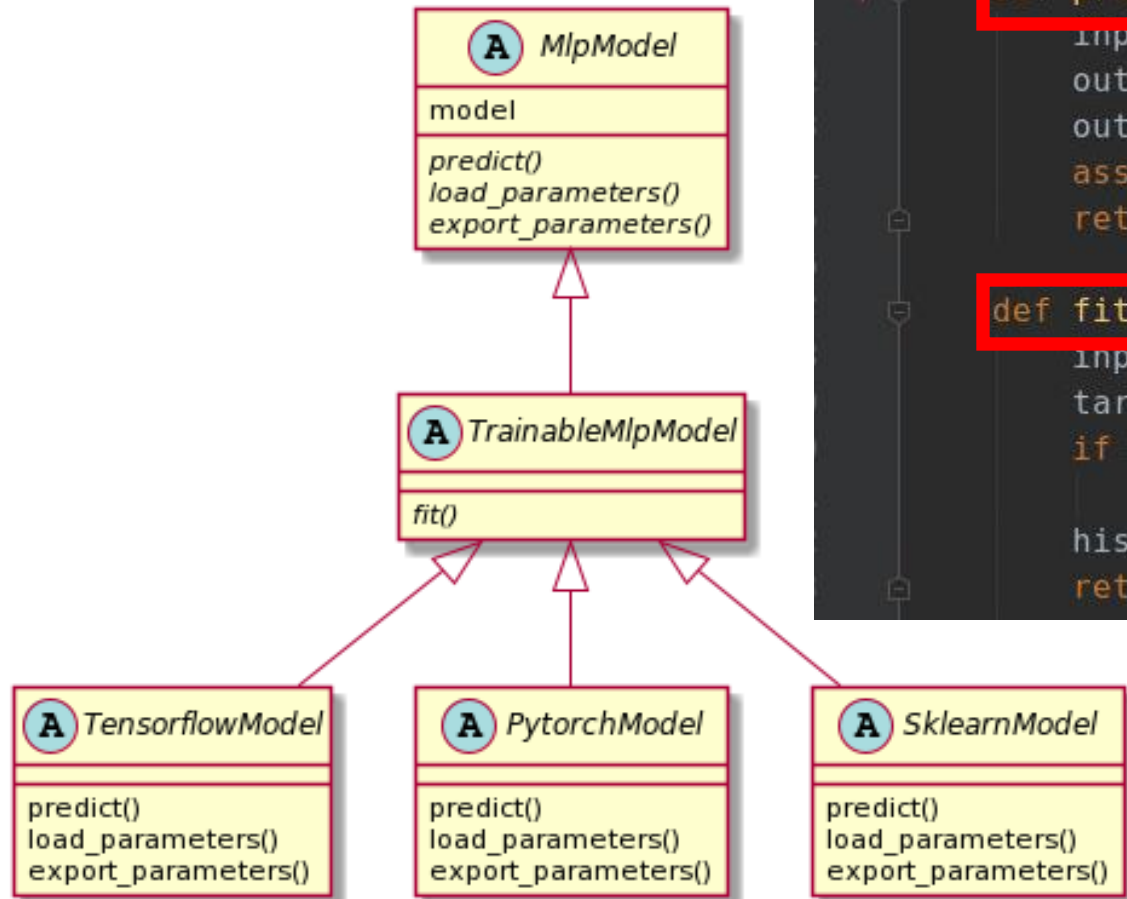
MLP in practice: Implementing a model



```
class Model(mlptf.TensorFlowModel):
    """ANN that fits a sine function."""

    def __init__(self):
        model: tf.keras.Model = tf.keras.models.Sequential([
            tf.keras.layers.Dense(16, activation="relu"),
            tf.keras.layers.Dense(16, activation="relu"),
            tf.keras.layers.Dense(1)
        ])
        model.compile(loss='mean_squared_error',
                      optimizer=tf.keras.optimizers.Adam(0.01))
        model.build(input_shape=self.input_shape)
        super().__init__(model)
```

MLP in practice: Implementing a model



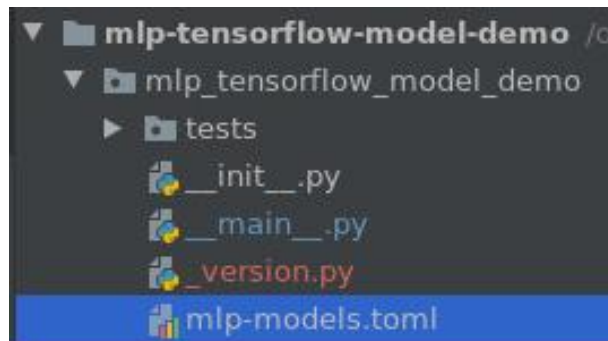
```
def predict(self, input_data):
    inputs = np.asarray(input_data[api.INPUTS])
    outputs = self.model.predict(inputs[:, np.newaxis])
    outputs = outputs[:, 0]
    assert np.ndim(outputs) == 1
    return {api.OUTPUTS: outputs}

def fit(self, input_data) -> api.FitHistory:
    inputs = np.asarray(input_data[api.INPUTS])
    targets = np.asarray(input_data[api.TARGETS])
    if np.ndim(targets) == 2:
        targets = np.squeeze(targets, 1)
    history = self.model.fit(inputs[:, np.newaxis], targets, epochs=500)
    return mlptf.convert_to_common_history(history)
```


MLP in practice: Implementing a model

Model declaration

- *mlp-models.toml*



```
1  [[model]]
2  name = "mlp_tensorflow_model_demo:Model"
3  standalone = false
4
5  [[model]]
6  name = "mlp_tensorflow_model_demo:SecondModel"
7  standalone = false
8
```

CI configuration

- *gitlab-ci.yml*

```
1  include:
2  - project: acc-co/machine-learning-platform/mlp-ci
3    file: pipeline-templates/mlp-model-gitlab-ci-template.yml
4
5  variables:
6  project_name: mlp_tensorflow_model_demo
7
8  #
9  # --- Custom CI jobs ---
10 #
11
12 # Full installation, tested with pytest.
13 test_install:
14   extends: .acc_py_full_test
15
16
17 # Development installation, tested with pytest.
18 test_dev:
19   extends: .acc_py_dev_test
```

MLP in practice: Using a model

Publish trained parameters

```
model = Model()
history = model.fit(dict(inputs=inputs, targets=targets))

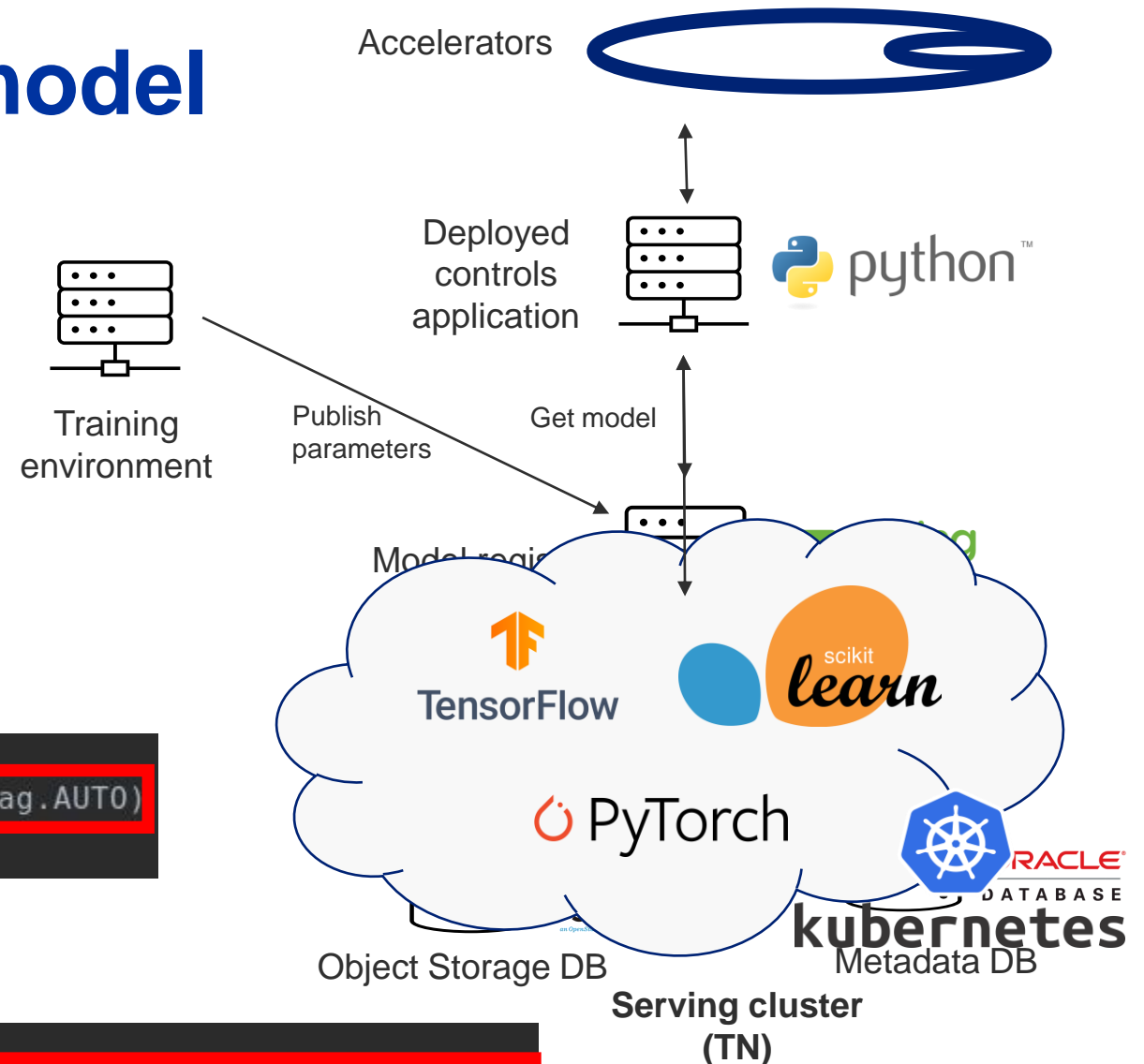
client = Client(Profile.DEV)
mpv = client.publish_model_parameters_version(
    model, name="tf_demo.tanh", version=AUTO)
```

Load remote parameters

```
client = Client(Profile.DEV)
model = client.create_model(Model, "tf_demo.sin", VersionFlag.AUTO)
preds = model.predict(dict(inputs=inputs))["outputs"]
```

Call standalone model (prototype)

```
client = Client(Profile.DEV)
model = client.create_standalone_model("mlp_tensorflow_model_demo:Model", "tf_demo.sin")
preds = model.predict(dict(inputs=inputs))["outputs"]
```



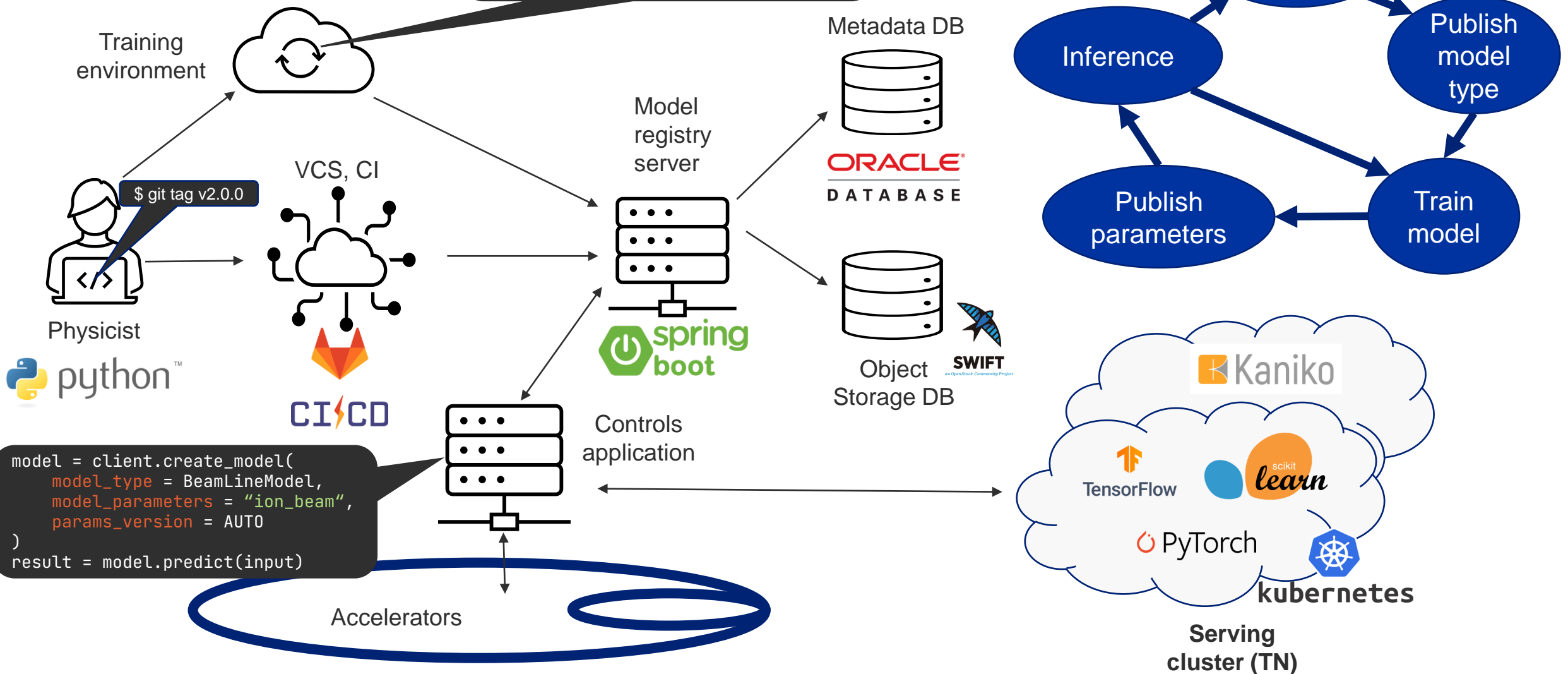
Conclusion

- **We want to help physicists develop models faster and unburden them from infrastructural concerns while minimizing constraints**
- **We also want to apply software engineering best practices to ensure reliability and maintainability of the control system**
- **MLP provides a basis to achieve these goals and is ready for production (standalone deployment in beta)**
- **Could not cover everything, simplified a lot – please contact us offline!**
 - [MLP Wikis](#) - [ICALEPCS paper](#) – machine-learning-platform-support@cern.ch
 - jean-baptiste.de.martel@cern.ch
 - nico.madysa@cern.ch
 - roman.gorbonosov@cern.ch

Thank you !

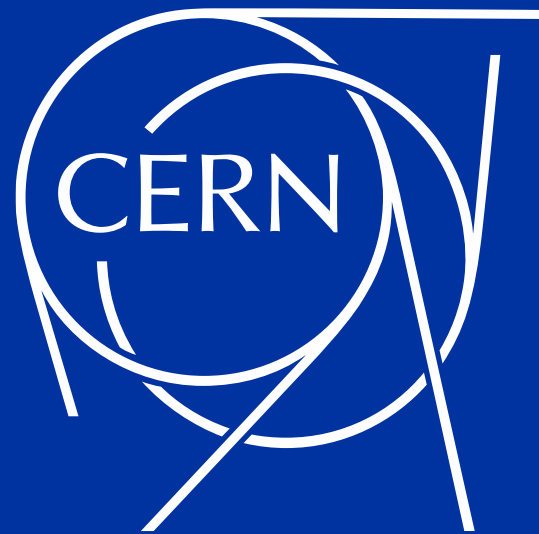
```

model.fit(training_data)
client.publish_parameters_version(
    model,
    name = "ion_beam",
    version = AUTO # generated
)
    
```



```

model = client.create_model(
    model_type = BeamLineModel,
    model_parameters = "ion_beam",
    params_version = AUTO
)
result = model.predict(input)
    
```

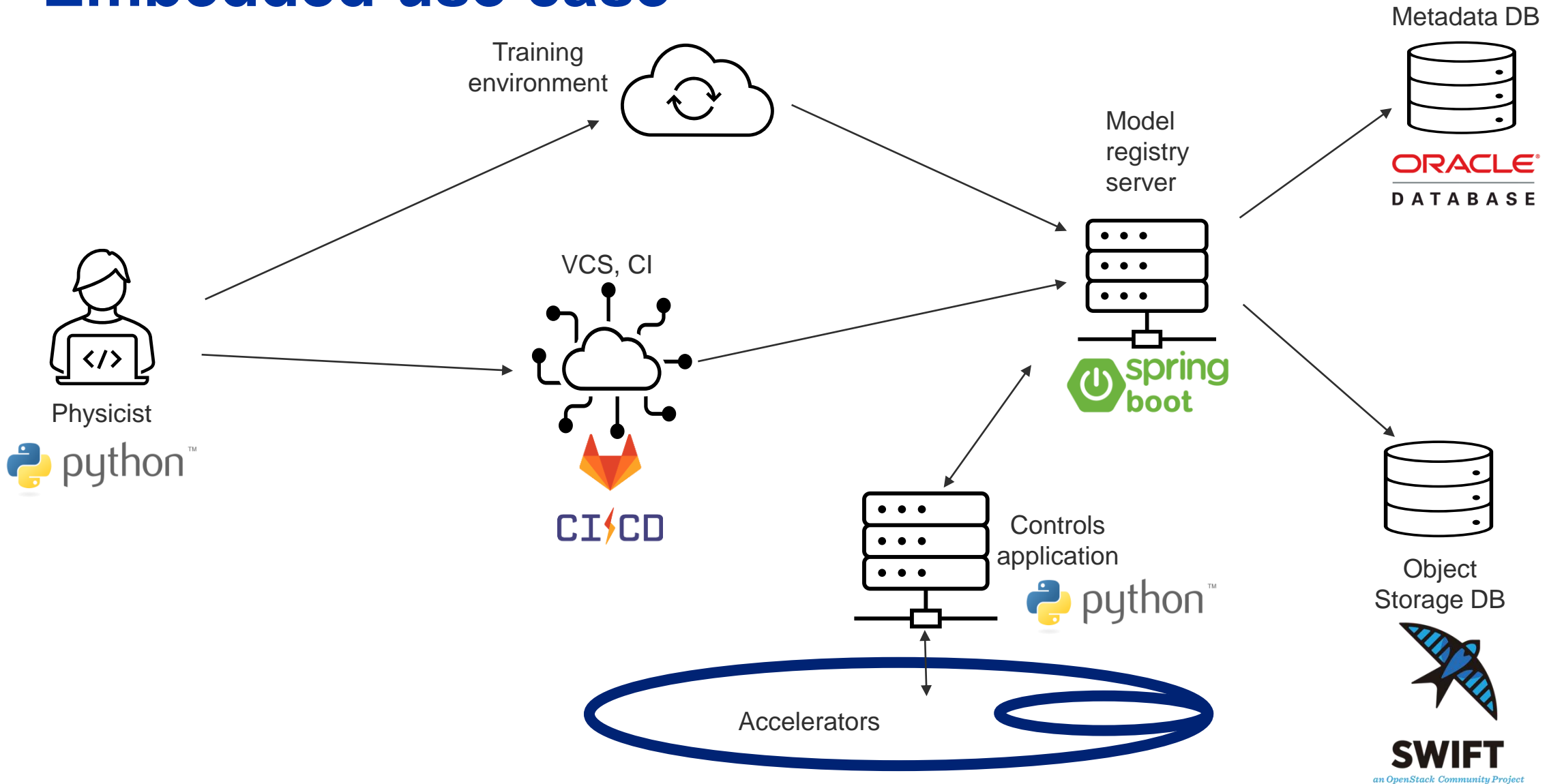


Reserve slides

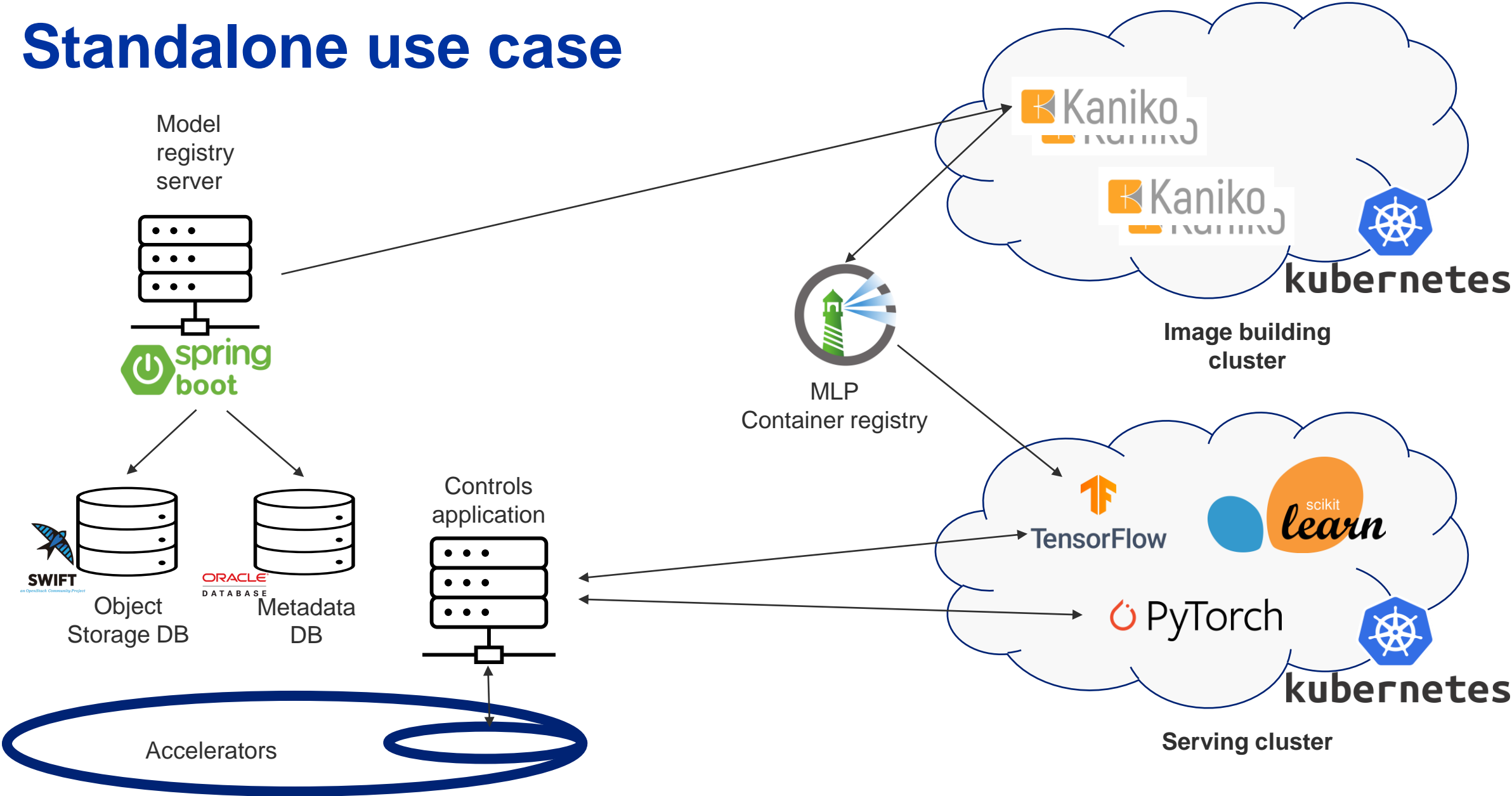
Use cases

- **SPS eddy current effect – V. Kain**
- **SPS deep hysteresis compensation – N. Madysa**
- **LEIR Schottky computer vision – N. Madysa**
- **SPS MKD dump pattern analysis – F. M. Velotti**
- **Awake auto-matching - F. M. Velotti**

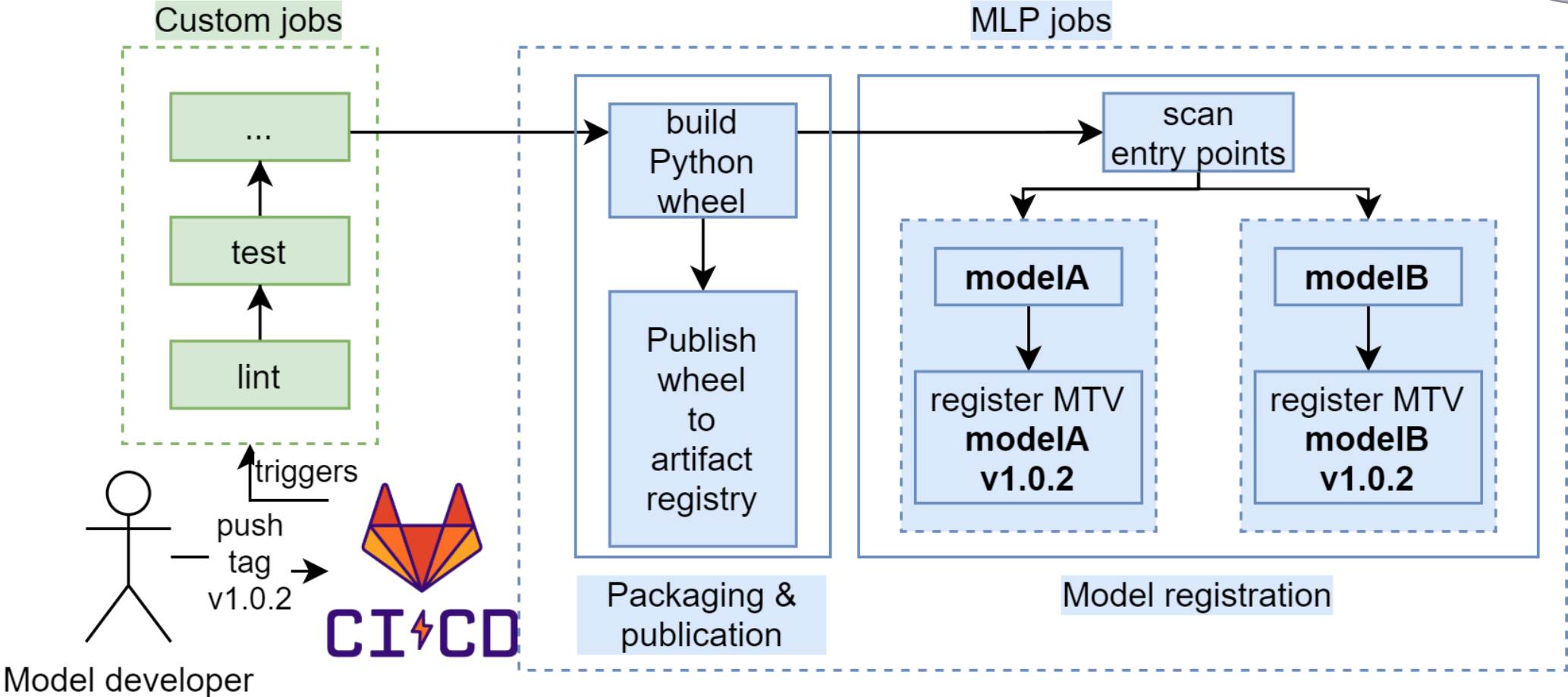
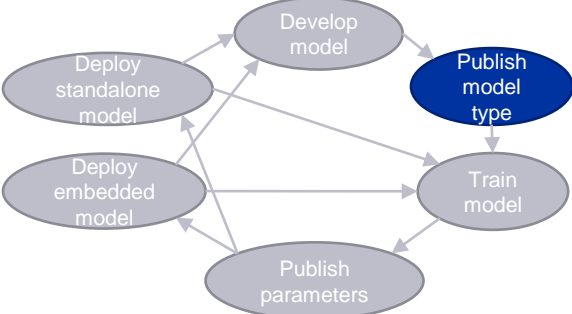
Embedded use case



Standalone use case



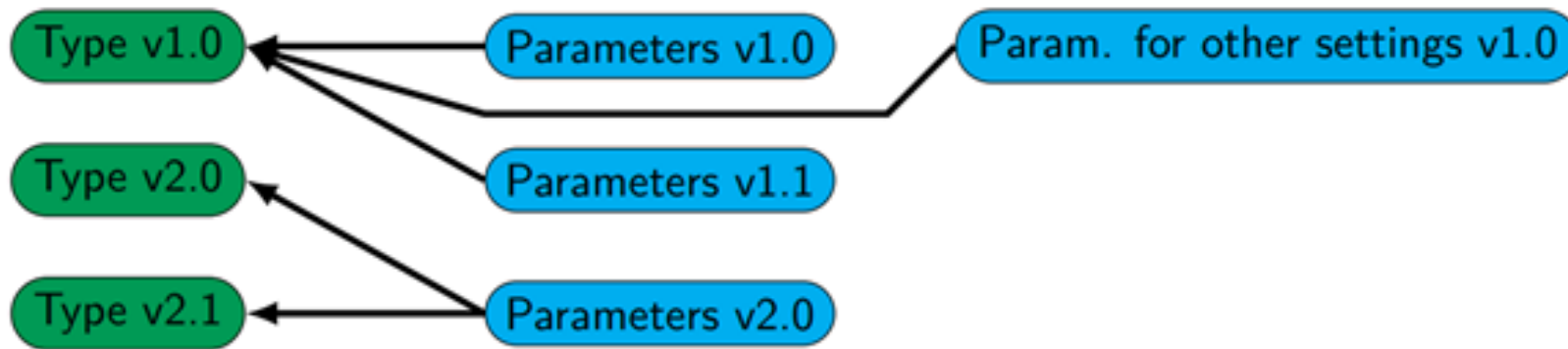
Publishing model types



Model parameters version number generation

Model type version	Highest existing parameters version	->	Generated parameters version
1.0.0	None exist yet	->	1.0
1.0.0	1.0	->	1.1
1.6.0	1.1	->	1.2
2.0.0	1.2	->	2.0
3.3.0	4.0 (no 3.x)	->	ambiguity
3.3.0	4.0 (3.3 exists)	->	3.4


Compatibility



Standalone deployment CI

Test

✓ test_dev 


✓ test_install 

Deploy


✓ acc_py_relea... 

✓ acc_py_relea... 


Register


✓ register mode... 


Standalone


✓ deploy stand... 

Downstream

✓ standalone-d...
#2758638 
Multi-project

✓ standalone-d...
#2758637 
Multi-project

✓ standalone-d...
#2758636 
Multi-project

✓ standalone-d...
#2758635 
Multi-project

Downstream


✓ standalone-d...
#2758638
Multi-project >

✓ standalone-d...
#2758637
Multi-project <

✓ standalone-d...
#2758636
Multi-project >

✓ standalone-d...
#2758635
Multi-project >

Build

✓ build contain... 

Replicate to acc registry

✓ replicate ima... 

Deploy

✓ deploy model 