

# RBatchGenerator()

# Uproot+Awkward

What are the problems?

- required more memory to test  
(maximum try 50 files for 15 GB RAM)

Next step to try!

- save each file to numpy file
- try to use torch dataloader and load from dir.
- compare speed of training and memory using in each approach.

# Load with torch loader

```
1 def load_dataset(data_path='/home/northnpk/Downloads/JetClass_dataset/',
2                 tree_name = "tree",
3                 batch_size = 128,
4                 max_num_particles=128,
5                 chunk_size = 5_000,
6                 vec_columns=['part_px', 'part_py', 'part_pz', 'part_energy'],
7                 columns=['jet_pt', 'jet_eta', 'jet_phi', 'jet_energy'],
8                 targets = ['label_QCD', 'label_Hbb', 'label_Hcc', 'label_Hgg', 'label_H4q',
9                           'label_Hqq', 'label_Zqq', 'label_Wqq', 'label_Tbqq', 'label_Tbl']) -> DataLoader:
10 # Passing vector columns
11 max_vec_sizes = dict(zip(vec_columns, [max_num_particles]*len(vec_columns)))
12
13 # get files path
14 train_path = [data_path + 'train/' + p for p in os.listdir(data_path + 'train')]
15
16 # load RDataFrame from path
17 r_df = RDataFrame(tree_name, train_path[:10])
18 # Getting RBatchGenerator from tmva RBatchGenerator
19 gen_train = ROOT.TMVA.Experimental.CreatePyTorchGenerators(r_df,
20                                                         batch_size,
21                                                         chunk_size,
22                                                         columns=columns+vec_columns+targets,
23                                                         max_vec_sizes=max_vec_sizes,
24                                                         target=targets,
25                                                         validation_split=0,
26                                                         drop_remainder=False)
27
28 # load Generator to torch dataloader
29 dataset = RBatchDataset(gen_train)
30 return DataLoader(dataset=dataset, collate_fn=dataset.collate_fn, batch_size=None)
```

# Load with torch loader

```
1 # Define load_dataset method to loading TMVA RBatchGenerator into Torch DataLoader
2 def load_dataset(data_path='/home/northnpk/Downloads/JetClass_dataset/',
3                 tree_name = "tree",
4                 batch_size = 128,
5                 max_num_particles=128,
6                 chunk_size = 5_000,
7                 vec_columns=['part_px', 'part_py', 'part_pz', 'part_energy'],
8                 columns=['jet_pt', 'jet_eta', 'jet_phi', 'jet_energy'],
9                 targets = ['label_QCD', 'label_Hbb', 'label_Hcc', 'label_Hgg', 'label_H4q',
10                          'label_Hqql', 'label_Zqq', 'label_Wqq', 'label_Tbqq', 'label_Tbl']) -> DataLoader:
```

# Load with torch loader

```
● ● ●  
1 # Passing vector columns  
2     max_vec_sizes = dict(zip(vec_columns, [max_num_particles]*len(vec_columns)))  
3  
4     # get files path  
5     train_path = [data_path + 'train/' + p for p in os.listdir(data_path + 'train')]  
6  
7     # load RDataFrame from path  
8     r_df = RDataFrame(tree_name, train_path[:10])
```

# Load with torch loader

```
1 # Getting RBatchGenerator from tmva RBatchGenerator
2 gen_train = ROOT.TMVA.Experimental.CreatePyTorchGenerators(r_df,
3                                                         batch_size,
4                                                         chunk_size,
5                                                         columns=columns+vec_columns+targets,
6                                                         max_vec_sizes=max_vec_sizes,
7                                                         target=targets,
8                                                         validation_split=0,
9                                                         drop_remainder=False)
10
11 # load Generator to torch dataloader
12 dataset = RBatchDataset(gen_train)
13 return DataLoader(dataset=dataset, collate_fn=dataset.collate_fn, batch_size=None)
```

# RBatchDataset(IterableDataset)

```
1 # Define Torch IterableDataset for TMVA RBatchGenerator
2 class RBatchDataset(IterableDataset):
3     def __init__(self, generator, vector_prep_fn=None):
4         self.generator = generator
5         self.batch_size = self.generator.base_generator.batch_size
6         self.vector_prep_fn = vector_prep_fn
7         self.apply_vector_prep_fn = False
8         if self.vector_prep_fn != None:
9             self.apply_vector_prep_fn = True
10        # print(f'Batch size: {self.batch_size}')
11        if self.generator.last_batch_no_of_rows != 0 and self.generator.number_of_batches > 1:
12            self.length = ((self.generator.number_of_batches-1) * self.batch_size) + self.generator.last_batch_no_of_rows
13        elif self.generator.number_of_batches == 1:
14            self.length = self.generator.last_batch_no_of_rows
15        else :
16            self.length = (self.generator.number_of_batches * self.batch_size)
17
18        self.vec_columns = [c not in self.generator.base_generator.given_columns for c in self.generator.base_generator.train_columns]
19        self._columns = [c in self.generator.base_generator.given_columns for c in self.generator.base_generator.train_columns]
20        self.label_columns = self.generator.base_generator.target_columns
21        self.vec_names = []
22        for n in self.generator.base_generator.given_columns:
23            if n not in self.label_columns and n not in self.generator.base_generator.train_columns:
24                self.vec_names.append(n)
25
26        def collate_fn(self, data):
27            tensors, targets = data
28            if self.apply_vector_prep_fn :
29                return self.vector_prep_fn(tensors[:, self.vec_columns].reshape(len(tensors), 4, 128), self.vec_names), tensors[:, self._columns], targets
30
31            else :
32                return tensors[:, self.vec_columns].reshape(len(tensors), 4, 128), tensors[:, self._columns], targets
33
34        def __iter__(self):
35            return self.generator.__iter__()
36
37        def __len__(self):
38            return self.generator.number_of_batches
39
```

# RBatchDataset(IterableDataset)

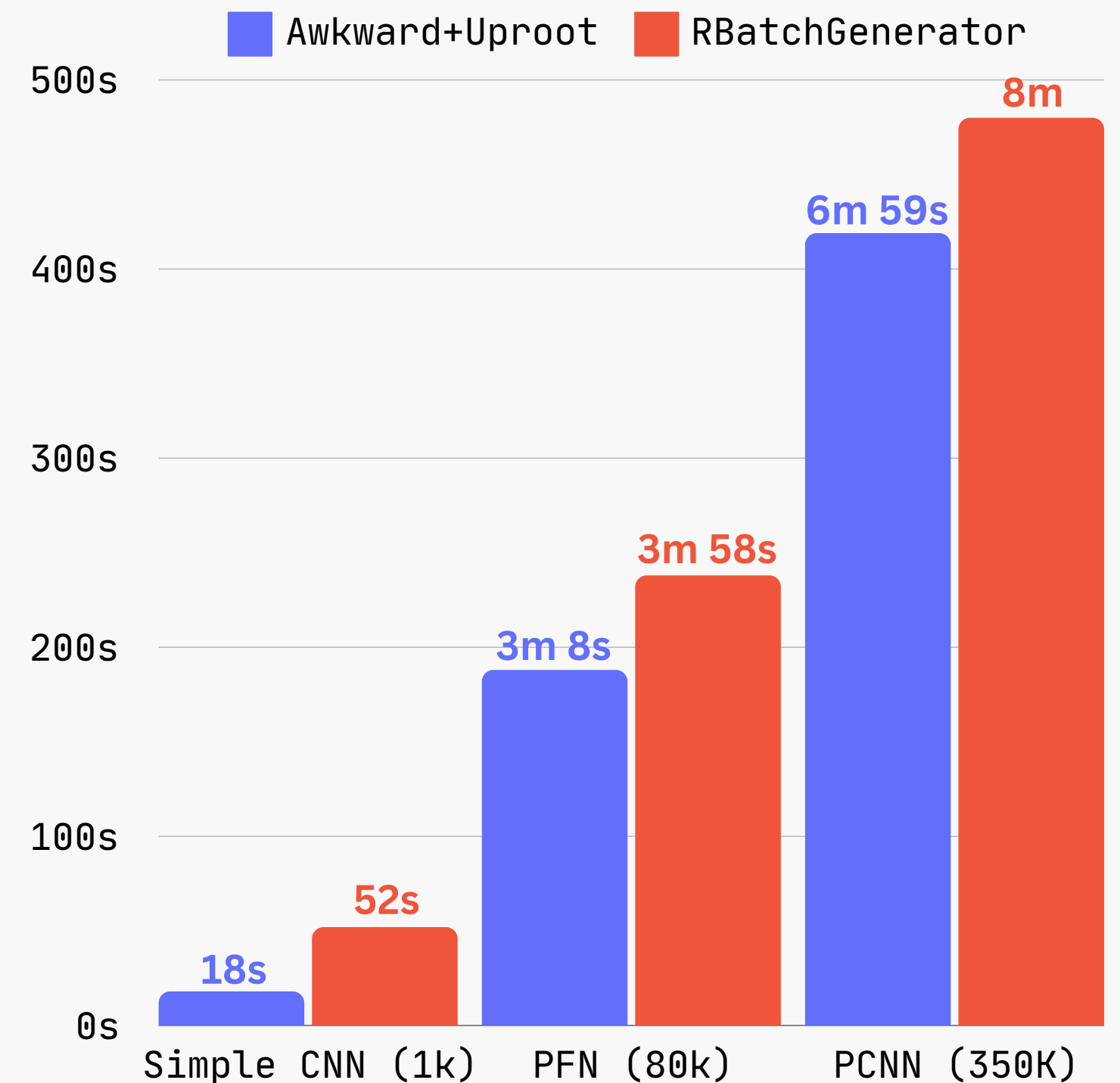
```
1 def collate_fn(self, data):  
2     tensors, targets = data  
3     if self.apply_vector_prep_fn :  
4         return self.vector_prep_fn(tensors[:, self.vec_columns].reshape(len(tensors), 4, 128), self.vec_names), tensors[:, self._columns], targets  
5  
6     else :  
7         return tensors[:, self.vec_columns].reshape(len(tensors), 4, 128), tensors[:, self._columns], targets
```



# RBatchGenerator()

training time with Pytorch  
CPU only for 1 epoch  
(Lower : Better)

- 10 files
- batch\_size : 128
- dataset : JetClass



# RBatchGenerator()

Next step to try!

- profiling with perf, memray to see What, Where, When, and How the memory use.
- Training Particle Transformers with RBatchGenerator

**Thank you**