*Kickoff ExaDost, Paris 2023.09.13*

# Scaling up Scikit-Learn for Exascale

**Scikit-learn:**

The reference library for ML on tabular data

- Many reference algorithms (no deep learning)

- Tools designed to create ML Pipelines (preprocessing)

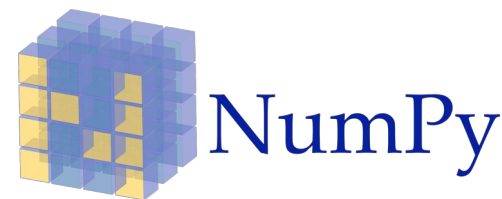- Utilities to perform hyperparameter optimization

- Tools to evaluate the results

**Scikit-learn:**

The reference library for ML on tabular data

**Current contraints**

- Input are pandas or numpy arrays

- Parallelism is done using joblib

- No GPU acceleration (for now)

# A zoom in the future development of Scikit-Learn

**I - Parallelism in scikit-learn**

**II – The Array API**

**III - Exploring a plugin system in scikit-learn**

## A high level API with joblib:

```
>>> from math import sqrt
>>> [sqrt(i ** 2) for i in range(10)]
[0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0]
```
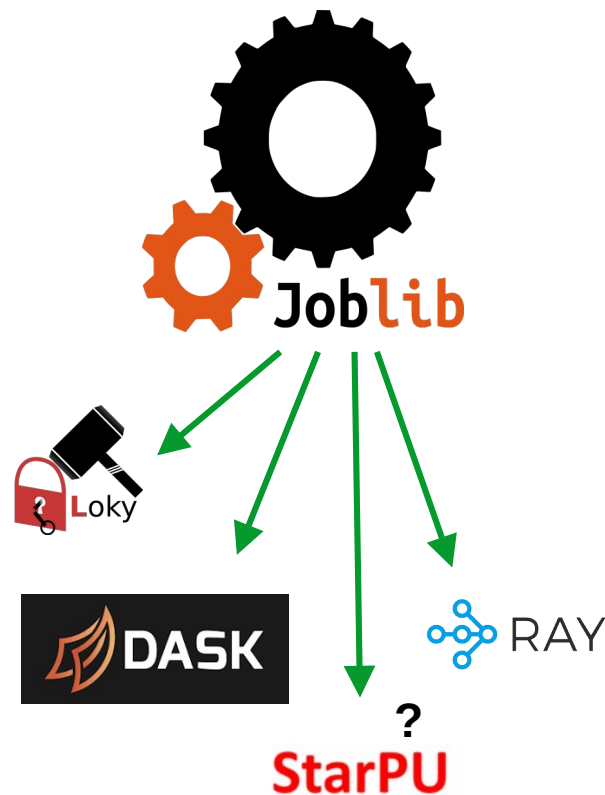
```
>>> from math import sqrt
>>> from joblib import Parallel, delayed
>>> Parallel(n_jobs=2)(delayed(sqrt)(i ** 2) for i in range(10))
[0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0]
```

```
>>> from joblib import parallel_config
>>> with parallel_config(backend='threading', n_jobs=2):
...     Parallel()(delayed(sqrt)(i ** 2) for i in range(10))
[0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0]
```

**> Embarassingly Parallel tasks, with pluggable backends**



Joblib

Loky

DASK

RAY

?

StarPU

**Low-level parallelism:**

Calls to BLAS (through numpy) are parallelised

Low level parallelism using openMP



**> Intrinsic parallelism, dependent on the data structure and the algorithm.**

**Challenge:** sometime hard to make the different parallelism levels cohexist!

# A zoom in the future development of Scikit-Learn

### I - Parallelism in scikit-learn

### II – The Array API

### III - Exploring a plugin system in scikit-learn

## Introduction to the array API standard:

```python
1  import numpy as np
2  data = np.arange(9).reshape(3, 3)
3  # with numpy, only get one output
4  # uses keyword `axis`
5  max_over_cols = np.max(data, axis=1)
```

# Introduction to the array API standard:

```python
import numpy as np
data = np.arange(9).reshape(3, 3)
# with numpy, only get one output
# uses keyword `axis`
max_over_cols = np.max(data, axis=1)
```

```python
import torch
data = torch.arange(9).reshape(3, 3)
# with pytorch, `torch.max` on nd tensors
# returns a length 2 tuple
# `max` expects parameter `dim`
max_over_cols, _ = torch.max(data, dim=1)
```

# Introduction to the array API standard:

```python
1  import numpy as np
2  data = np.arange(9).reshape(3, 3)
3  # with numpy, only get one output
4  # uses keyword `axis`
5  max_over_cols = np.max(data, axis=1)
```

```python
1  import torch
2  data = torch.arange(9).reshape(3, 3)
3  # with pytorch, `torch.max` on nd tensors
4  # returns a length 2 tuple
5  # `max` expects parameter `dim`
6  max_over_cols, _ = torch.max(data, dim=1)
```

**Different array libraries can have different syntax for the same operation !**

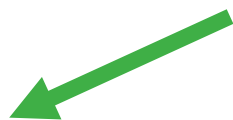## Introduction to the array API standard:

The Array API enables writing **agnostic single-source code** and interface it with all arrays !

```python
1  from array_api_compat import get_namespace
2  def max_over_cols_fn(data):
3      xp = get_namespace(data)
4      return xp.max(data, axis=1)
```

## Introduction to the array API standard:

The Array API enables writing **agnostic single-source code** and interface it with all arrays !

```
1  from array_api_compat import get_namespace
2  def max_over_cols_fn(data):
3      xp = get_namespace(data)
4      return xp.max(data, axis=1)
```

```
1  import numpy as np
2  data = np.arange(9).reshape(3, 3)
3  max_over_cols = max_over_cols_fn(data)
```

```
1  import torch
2  data = torch.arange(9).reshape(3, 3)
3  max_over_cols  = max_over_cols_fn(data)
```

# Array API spec excerpt:





**Assumptions**
**API**
Design topics & constraints
Future API standard evolution
API specification
　Array object
　Broadcasting
　Constants
　Creation Functions
　Data Type Functions
　Data Types
　Element-wise Functions
　Function and method signatures
　Indexing
　Indexing Functions
　Linear Algebra Functions
　Manipulation Functions
　Searching Functions

A conforming implementation of the array API standard must provide and support the

## Objects in API

| | |
|---|---|
| `broadcast_arrays`(*arrays) | Broadcasts one or m |
| `broadcast_to`(x, /, shape) | Broadcasts an array |
| `concat`(arrays, /, *, axis=0) | Joins a sequence of |
| `expand_dims`(x, /, *, axis=0) | Expands the shape (dimension) of size |
| `flip`(x, /, *, axis=None) | Reverses the order |
| `permute_dims`(x, /, axes) | Permutes the axes ( |
| `reshape`(x, /, shape, *, copy=None) | Reshapes an array w |
| `roll`(x, /, shift, *, axis=None) | Rolls array elements |
| `squeeze`(x, /, axis) | Removes singleton |
| `stack`(arrays, /, *, axis=0) | Joins a sequence of |

https://data-apis.org/array-api/latest/API_specification/index.html

## ExaDosT Kick-off meeting

# NumPEx
*Exascale computing*

## Compatible array libraries:
(done or ongoing)
(non-exhaustive)

**NumPy**

**CPU**

**PyTorch**

**CPU**
**GPU AMD/INTEL/NVIDIA**
**OTHERS**

**dask**

**CPU**
**DISTRIBUTED**

**Dpctl + DPNP**
**CPU**
**GPU INTEL**

**CuPy**

**GPU NVIDIA**

**Toward Array API compliancy in scikit-learn:**

*When possible, use the Array Api rather than numpy calls !*

```
636  -          np.square(X, out=X)                          659  +          X **= 2
637  -          np.sum(X, axis=0, out=X[0])                  660  +          total_var = xp.sum(xp.sum(X, axis=0) / N)
638  -          total_var = (X[0] / N).sum()

639                                                          661

640          self.explained_variance_ratio_ =               662          self.explained_variance_ratio_ =
        self.explained_variance_ / total_var                        self.explained_variance_ / total_var
641  -      self.singular_values_ = S.copy()  # Store the singular   663  +      self.singular_values_ = xp.asarray(S, copy=True)  #
        values.                                                     Store the singular values.
642                                                          664

643          if self.n_components_ < min(n_features, n_samples):   665          if self.n_components_ < min(n_features, n_samples):
644  -          self.noise_variance_ = total_var -           666  +          self.noise_variance_ = total_var -
        self.explained_variance_.sum()                              xp.sum(self.explained_variance_)
```

Excerpt of  https://github.com/scikit-learn/scikit-learn/pull/26315
*ENH Array API support for PCA*

**ExaDosT Kick-off meeting**

**Towards Array API compliancy in scikit-learn**:

✓ **sklearn.lda.LDA**

✓ **sklearn.decomposition.PCA**    (next release)

⚙ **sklearn.preprocessing.MinMaxScaler**    (ongoing)

Discussions and roadmaps #22352 #26024 #26083

## Towards Array API compliancy in scikit-learn:

🔥 *How to use it ?* 🔥

```
1  from sklearn import config_context
2  from sklearn.decomposition import PCA
3  pca_estimator = PCA(n_components=10)
4  with config_context(array_api_dispatch=True):
5      pca_estimator.fit(X)
```
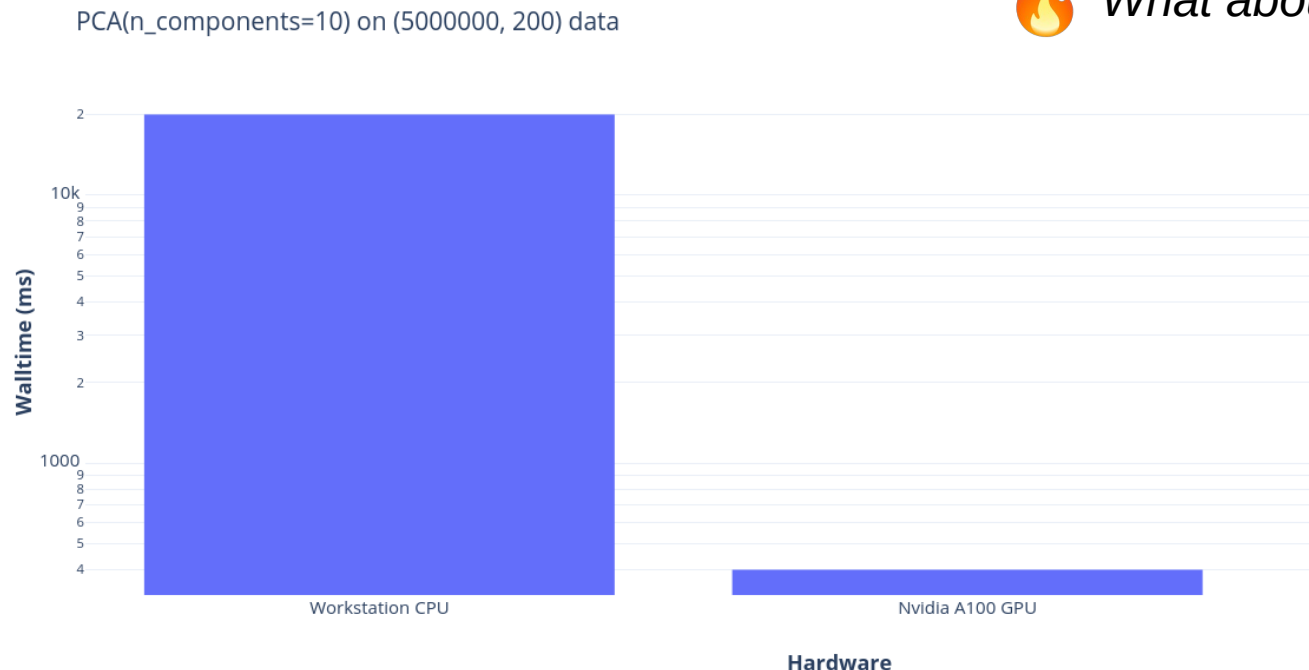
⚠️ Experimental features require explicit activation in scikit-learn config class

```
1  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
2  lda = LinearDiscriminantAnalysis()
3  with config_context(array_api_dispatch=True):
4      lda.fit(X)
```

# Towards Array API compliancy in scikit-learn:

PCA(n_components=10) on (5000000, 200) data



🔥 *What about performance ?* 🔥

YES when CPU takes too long
    40x speedup with pytorch
    on GPU

NB: pytorch backend for intel
GPUs currently lacks optimization
for QR factorization, so Intel GPUs
are not included in the figure

NumPEx
*Exascale computing*

**Towards Array API compliancy in scikit-learn**:

🔥 *Will all estimators support all compatible array libraries ?* 🔥

**Likely not :-(**

Current limitations of the Array API approach:

❌ Does not cover all functions for all libraries (torch.topk, np.argpartition, RNG…)

❌ Lacks synchronization functions for libraries that evaluate lazily (dask, jax...)

❌ Not all algorithms can be implemented (efficiently) with high level array operations !

(random forest, histogram gradient boosting trees,…)

**ExaDosT Kick-off meeting**

**Possible drawbacks for the Array API:**

Each call to a high-level array function writes the result into memory.

Example of pairwise_distance + argmin common pattern:

```
1   import numpy as np
2   from scipy.spatial.distance import cdist
3
4   # create some data
5   rng = np.random.default_rng(123)
6   data = rng.random((500000, 200))
7   query = rng.random((1, 200))
8
9   # writes all pairwise distances to memory
10  pairwise_distances = cdist(query, data)
11
12  closest = np.argmin(pairwise_distances)
```

# A zoom in the future development of Scikit-Learn

### I - Parallelism in scikit-learn

### II – The Array API

### III - Exploring a plugin system in scikit-learn

**Exploring a plugin system in scikit-learn**:

> Preview of user workflow:

```
$
$ pip install scikit-learn some-compute-plugin
```

**Exploring a plugin system in scikit-learn:**

> Preview of user workflow:

```
$
$ pip install scikit-learn some-compute-plugin
```

```python
1    from sklearn import config_context
2    from sklearn.cluster import KMeans
3    kmeans_estimator = KMeans()
4    with config_context(engine_provider="some_compute_plugin"):
5        y_pred = kmeans_estimator.fit_predict(X)
6        y_transform = kmeans_estimator.transform(X)
```

**Exploring a plugin system in scikit-learn**:

> Example for plugin developers: register a custom engine class

```python
class MyKMeansEngine:

    def __init__(self, estimator):
        self.estimator = estimator

    def prepare_fit(self, X, y=None, sample_weight=None):
        "TODO: insert custom engine implementation"

    def kmeans_single(self, X, sample_weight, centers_init):
        "TODO: insert custom engine implementation"

    def prepare_prediction(self, X, sample_weight):
        "TODO: insert custom engine implementation"

    def get_labels(self, X, sample_weight):
        "TODO: insert custom engine implementation"
```

> Test against native scikit-learn unit tests

```
~# pytest --sklearn-engine-provider some-compute-plugin --pyargs sklearn.cluster.tests.test_k_means
```

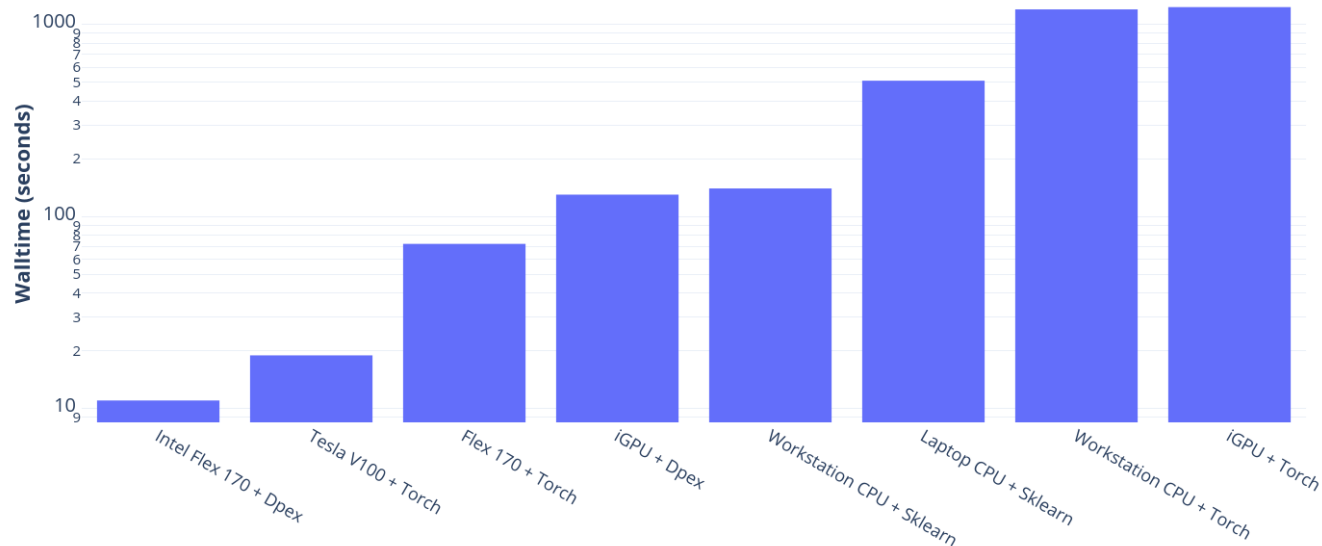**Exploring a plugin system in scikit-learn**:

Foundational issue     https://github.com/scikit-learn/scikit-learn/issues/22438

Pull request     https://github.com/scikit-learn/scikit-learn/pull/25535

Development branch     https://github.com/scikit-learn/scikit-learn/tree/feature/engine-api

Example plugin for KMeans on intel GPUs
     https://github.com/soda-inria/sklearn-numba-dpex

# A plugin for k-means with the ▲ oneAPI toolchain:

**Walltime for 100 k-means lloyd iterations on 50_000_000 * 14 data**



**sklearn_numba_dpex**
https://github.com/soda-inria/sklearn-numba-dpex

k-means is <u>one order of magnitude faster</u> on GPU, more significantly so when optimized with a low-level implementation

NB: cuml implementation of k-means for Nvidia GPUs could not be included in the figure because it could not scale to this amount of data memory-wise.

## ROADMAP in ExaDost:

- 1 Engineer position to work on scikit-learn and joblib (MIND)

  - Work on the lazy Array API
  - Work on online computations with scikit-learn (partial_fit),
  - Work on improving nested parallelism handeling.

- 1 PhD position with H. Hendrikx, (Thot)

  - Improve the distributed algorithms in scikit-learn