# Hybrid AI/HPC Approaches and Linear Algebra

HPC challenges for new extreme scale applications

Paris, March 6-7, 2023

*Nahid Emad*

*University of Paris Saclay / UVSQ*
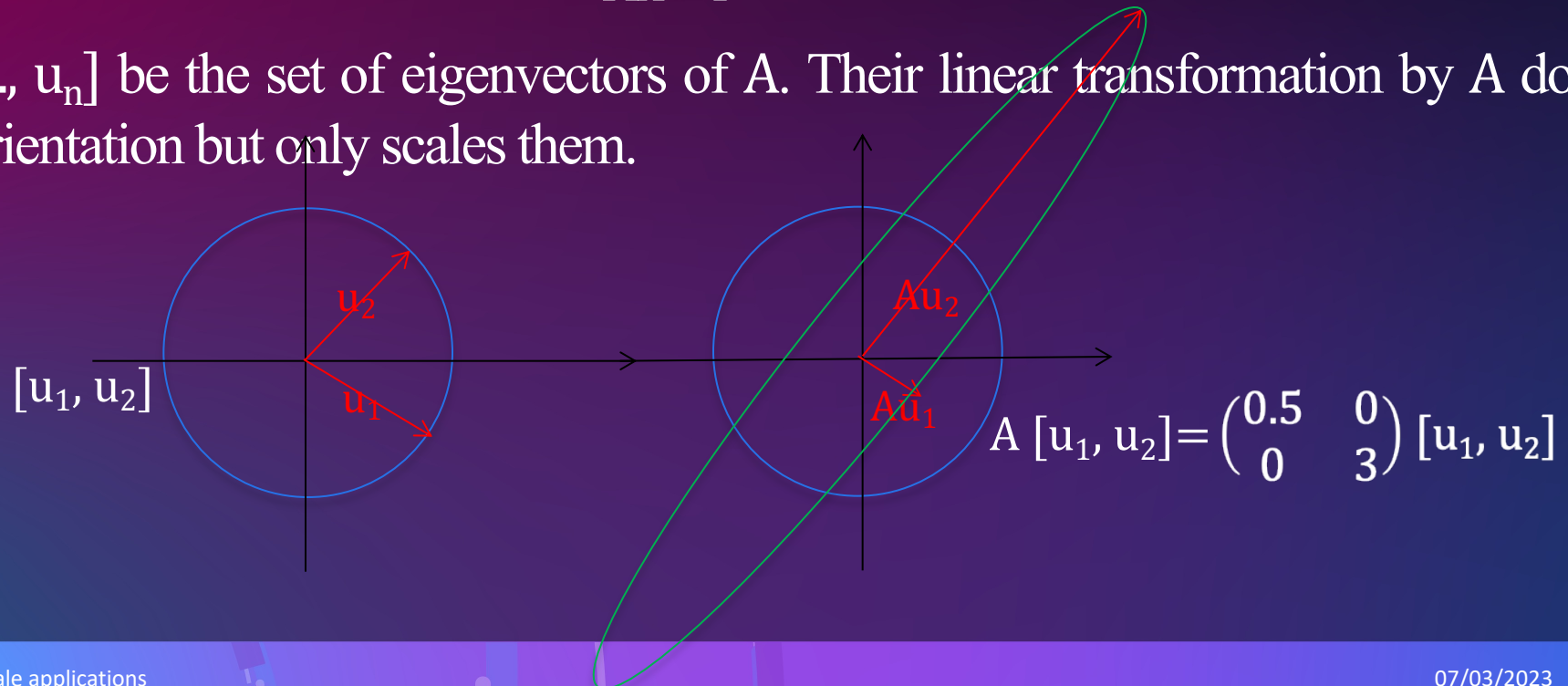*Maison de la Simulation & LI-PaRAD*

# Linear algebra main problem in ML/DL

- In machine learning, many problems can be solved by **linear transformations** and systems of linear equations.

- Let A and Y be n-size matrix representing a set of n observations and the vector of their labels. The search of a function f(A)=Y can be expressed as a linear system:

$$Ax=Y$$

- Let $U =[u_1, ..., u_n]$ be the set of eigenvectors of A. Their linear transformation by A does not change their orientation but only scales them.

$[u_1, u_2]$       $u_2$     $u_1$

$Au_2$     $Au_1$    $A\,[u_1, u_2] = \begin{pmatrix} 0.5 & 0 \\ 0 & 3 \end{pmatrix} [u_1, u_2]$

# Dominant eigenspace in ML

**Principal Component Analysis:** The goal is to find an orthonormal basis of the space of a dataset such that the variance of the dataset (degree of dispersion) in this basis is maximized. PCA helps reduce redundancies in datasets and extract important features while preserving accuracy.

- Let $X \in \mathbb{R}^{n \times p}$ be a centered matrix of $n$ observations of $p$ features. The PC of $X$ are the dominant eigenvectors of its covariance matrix $A = \frac{1}{n} X^T X$.

- The PC of $X$ are its dominant right singular vectors: $X = U \Sigma V^T$ with $U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{p \times p}$ unitary and $\Sigma \in \mathbb{R}^{n \times p}$ diagonal matrices of singular values. $A = X^T X = V \Sigma^2 V^T$. The columns of $V$ are the right singular vectors of $X$.

**PageRank algorithm example:** The Markov matrix leads to the equation which the steady state depends on one dominant component: $\lambda_1^k \boldsymbol{u_1} + \alpha_1 \lambda_2^k u_2 + \ldots + \alpha_n \lambda_1^k u_n$.

# ML methods and linear algebra

**Goal**: Build smarter machines thinking and acting on their own (needs of training –still- and more and more data)

- Supervised machine learning methods

  - Linear regression, logistic regression, recommendation systems, ANN, etc.
  - Linear algebra problem as linear systems and **eigenproblems**

- Unsupervised machine learning methods

  - K-means for partitioning, **dimensionality reduction**, **CPA**, etc.
  - Essentially **eigenproblems** and **SVD**

- Reinforcement learning methods (exploration & exploitation)

  - Bandit, **Markovian decision problems**, game trees.

# High performance data analysis

- **Data production** is now faster than compute capabilities

- **Applications** are classical simulation, social network-based simulation, ML algorithms

- Emerging **Exascale supercomputers** : Multi-level architectures (processor, memory, …), mixed arithmetic (16, 32, 64 bits,...), ..., and convergence of distributed and parallel computing inside them.

- Need of new **programming paradigms** for this extreme computational and data sciences programming.

- **New methods** must be developed (involving applied math, graph theory, Bayesian network, statistic, linear algebra, game theory, ...) but also, the new approaches such as transformer used in NLP.

- **Big Data analysis and HPC convergence** is crucial to propose future machine learning algorithm for high-scale platforms and supercomputers

**New paradigms for new intelligent applications**

# Outline

1. Main problems in linear algebra (moderate size)

2. Large and sparse linear algebra problem

3. A brief overview of IA : computing viewpoint (ML /DL)

4. Some applications of High-performance LA and AI

5. Concluding remarks

# Main problems in linear algebra (moderate size)

**Linear system (LS) :**

$$\text{Let } A \in \mathbb{C}^{n \times n}, \text{b} \in \mathbb{C}^n, \text{find } \text{x} \in \mathbb{C}^n, \text{such that}: A.x = b$$

**Eigenproblem (EIG) :**

$$\text{Let } A \in \mathbb{C}^{n \times n}, \text{find } \lambda_i \in \mathbb{C} \text{ and } u_i \in \mathbb{C}^n \text{ such that}: A.u_i = \lambda_i.u_i \quad (i = 1, \dots, n)$$

- Solving LS (topic well mastered overall)
  - ➢ **Direct** methods as Gauss and Gauss-Jordan, Cholesky, Householder based on LU, Cholesky, QR decomposition.
  - ➢ **Iterative** methods as Jacobi, Gauss-Seidel, Relaxation.
- Solving EIG (**topic not so well mastered**)
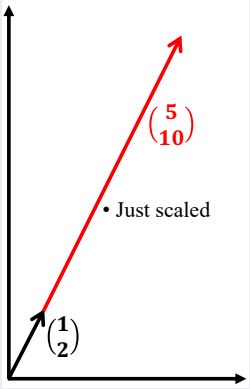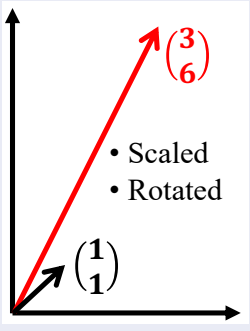  - ➢ **Only iterative** methods (Abel-Ruffini theorem) as Jacobi and QR

# Focus on Eigenproblem

**Eigenproblem (EIG) :**

$$\text{Let } A \in \mathbb{C}^{n \times n}, \text{ find } \lambda_i \in \mathbb{C} \text{ and } u_i \in \mathbb{C}^n \text{ such that} : A.u_i = \lambda_i.u_i \quad (i = 1, \dots, n)$$

**Power of eigenvectors :**

✓ A doesn't change the orientation of an eigenvector and/or eigenspace but just scales it.

✓ Principal components or axes of dataset:

| $A$ | $x$ | $Ax$ |
|---|---|---|
| $\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$ | $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ |  $\begin{pmatrix} 5 \\ 10 \end{pmatrix}$ • Just scaled $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ |
| | $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ |  $\begin{pmatrix} 3 \\ 6 \end{pmatrix}$ • Scaled • Rotated $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ |

Eigen-elements of A : $\lambda_1$=0, $\lambda_2$=5 and $v_1 = \begin{pmatrix} -2 \\ 1 \end{pmatrix}$, $v_2 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$
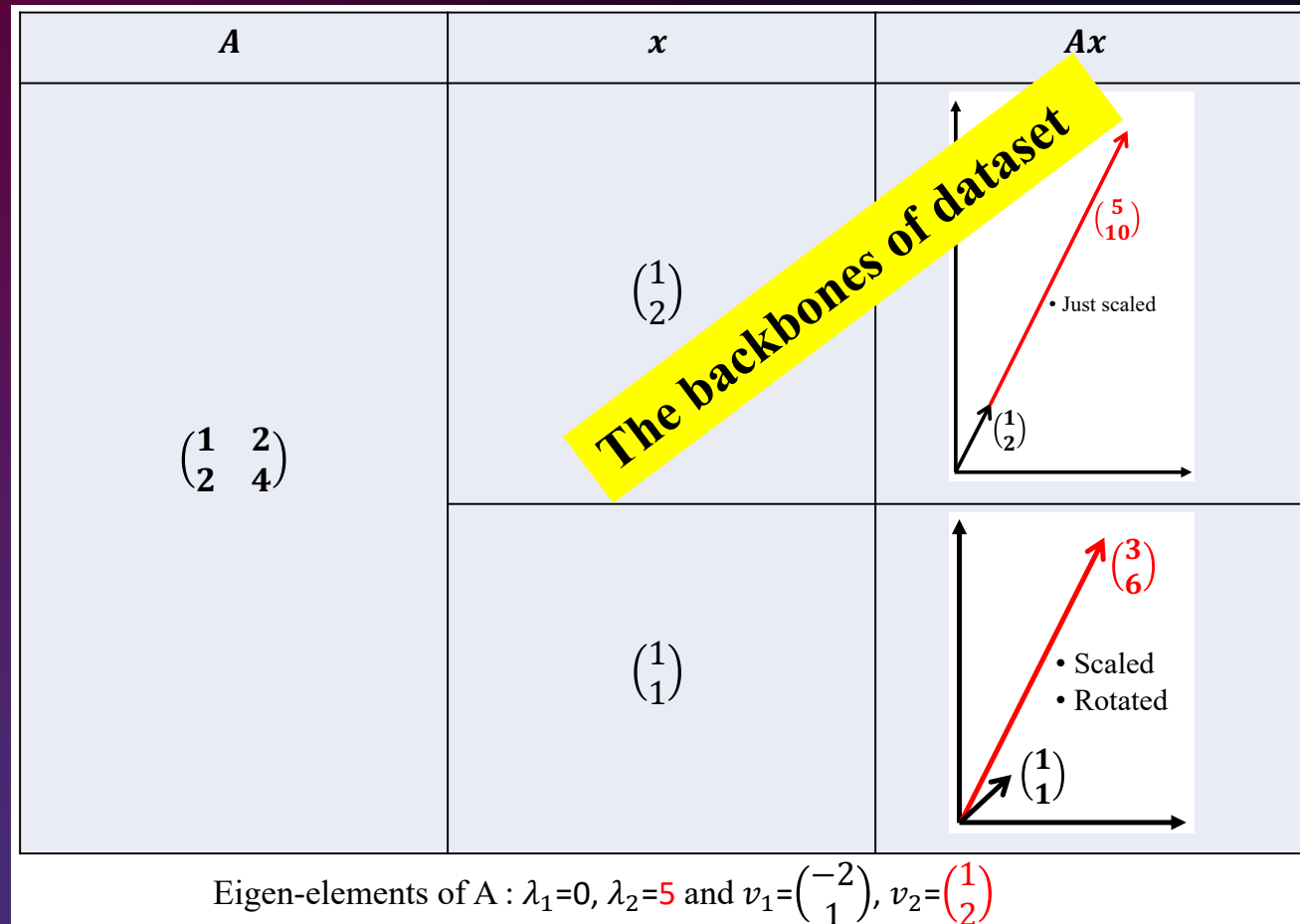
# Focus on Eigenproblem

**Eigenproblem (EIG) :**

Let $A \in \mathbb{C}^{n \times n}$, find $\lambda_i \in \mathbb{C}$ and $u_i \in \mathbb{C}^n$ such that : $A.u_i = \lambda_i.u_i$ $(i = 1, \ldots, n)$

**Power of eigenvectors :**

✓ A doesn't change the direction of an eigenvector and/or eigenspace but just scales it.
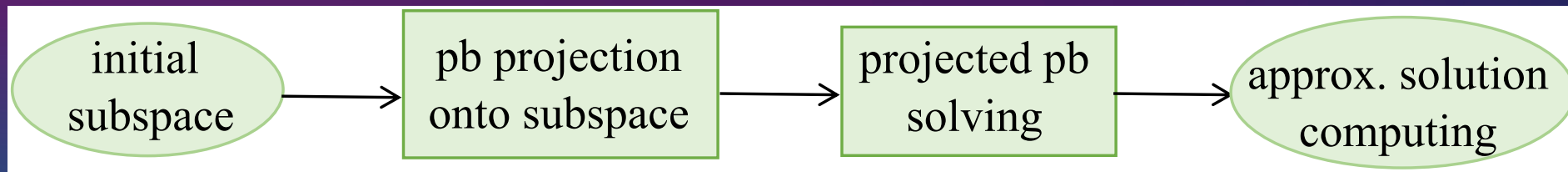
✓ Principal components or axes of dataset:

| $A$ | $x$ | $Ax$ |
|---|---|---|
| $\begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$ | $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ | The backbones of dataset  $\begin{pmatrix} 5 \\ 10 \end{pmatrix}$ • Just scaled $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ |
| | $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ | $\begin{pmatrix} 3 \\ 6 \end{pmatrix}$ • Scaled • Rotated $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ |

Eigen-elements of A : $\lambda_1 = 0$, $\lambda_2 = 5$ and $v_1 = \begin{pmatrix} -2 \\ 1 \end{pmatrix}$, $v_2 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$

# 2. Large and sparse linear algebra problems

- **Sparse dataset**

  ➢ Avoiding fill-in

  ➢ Problem : how to compress the dataset ? Use of ML methods
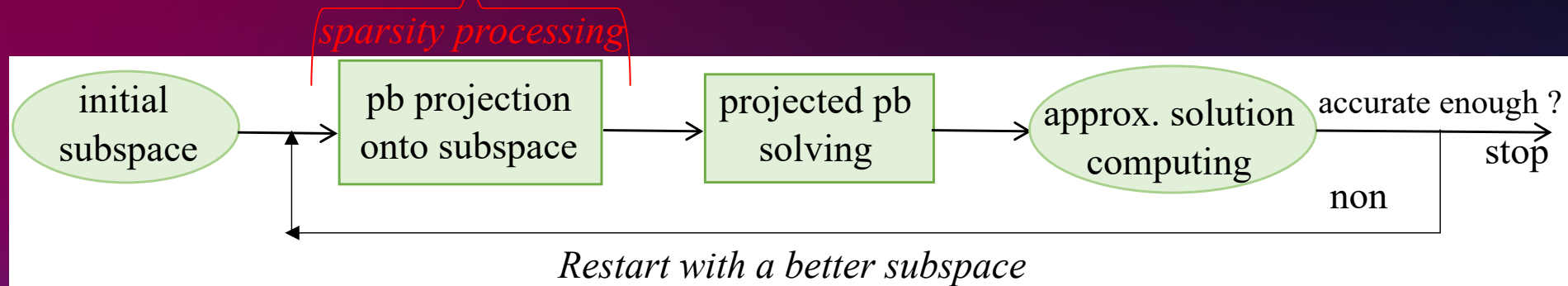
- **Large dataset**

  ➢ Dimensionality reduction - projection onto Krylov subspace

  ➢ Problem : how to choose the subspace-size ? Too large/small…

```
initial          pb projection       projected pb        approx. solution
subspace    →    onto subspace   →   solving        →    computing
```

# 2. Large and sparse linear algebra problems

**Iterative projection method**

➤ Preserve sparsity

➤ Reduce the problem size



*sparsity processing*

initial subspace → pb projection onto subspace → projected pb solving → approx. solution computing → accurate enough ? stop / non

*Restart with a better subspace*

**Main problems for these methods**

➤ Sparsity processing                    **What about m ?**

➤ Krylov subspace: optimal choice of $v$ for $\mathbb{K}_m(A, v) = span\,(v, Av, ..., A^{m-1}v)$

   optimal choice of m and $v$ ?

# Innovative Approach : Unite and Conquer methods

Suppose we have $\ell$ iterative methods to solve the same given problem. The unite and conquer approach consists of making collaborate these $\ell$ methods in order to accelerate the convergence of the whole system.

## Characteristics of UC methods

- Multi level parallelism (coarse grain and fine grain)
- Asynchronous communication
- Fault tolerance
- Great potential to dynamic load balancing
- Many parameters, many reuse software components
- Need well suited «standard» programming tools



*The key part of the method: restarting strategy*

$$f(S_1^k, ..., S_\ell^k)$$

*Well suited to petascale & emerging exascale computing systems*

# Innovative Approach : Unite and Conquer methods

Suppose we have $\ell$ iterative methods to solve the same given problem. The unite and conquer approach consists of making collaborate these $\ell$ methods in order to accelerate the convergence of the whole system.

**Characteristics of UC methods**

- Multi level parallelism (coarse
- Asy
- P
- G
- Ma
- Ne

The key part of the method: restarting strategy

$f(S_1^k, ..., S_\ell^k)$

Initial

Initial condition $\ell$

Iterative method $\ell$

method 2

$S_2^k$

$S_\ell^k$

Asynchronous communication

STOP          STOP          STOP

With S. Petiton, G. Edjlali, A. Fazeli, J. Dongarra, A. Drummond, T. Sakurai, M. Tsuji, C. Calvin, A. Fender, Z. Liu, and many others. N. Emad, S. Petiton. *Unite and Conquer Approach for High Scale Numerical Computing*, Journal of Computational Science, ISSN-1877-7503, 2016.

*Well suited to petascale & emerging exascale computing systems*

# Innovative Approach : Unite and Conquer methods

Due to the numerical and computational properties of a UC method, its overall convergence and computational performance are better than that of each of its co-methods individually.

➢ **Multiple-Method :** Case of UCM when the co-methods are the instances of the same iterative method. Example: MERAM, MIRAM, MLnczos, with different or nested subspaces.

➢ The asynchronism of communications implies better computational performance but introduces a certain *non-determinism*.

➢ The application of the UC approach to ML methods, which are inherently non-deterministic, does not suffer from this non-determinism.

# 3. A brief overview of AI : ML/DL computing viewpoint

1. HPC and AI convergence

2. High performance LA and AI ; some applications

   ➢ Sparse computation : a common topic of HPLA & ML/DL

   ➢ Focus on clustering (unsupervised ML) using UC methods

      o K-means and spectral computation

      o nvGraph of NVIDIA

   ➢ Focus on (semi)supervised (classification) using UC methods

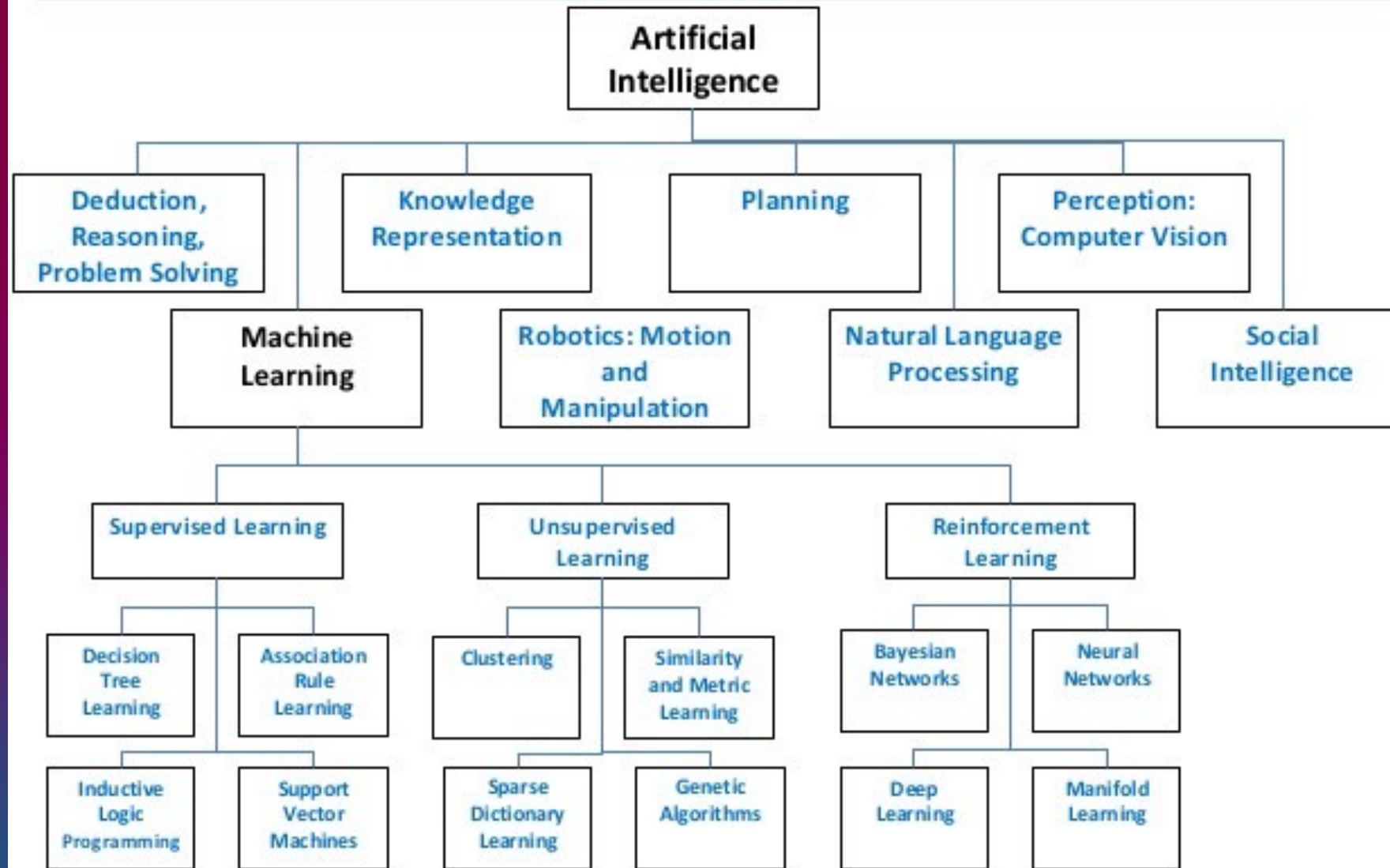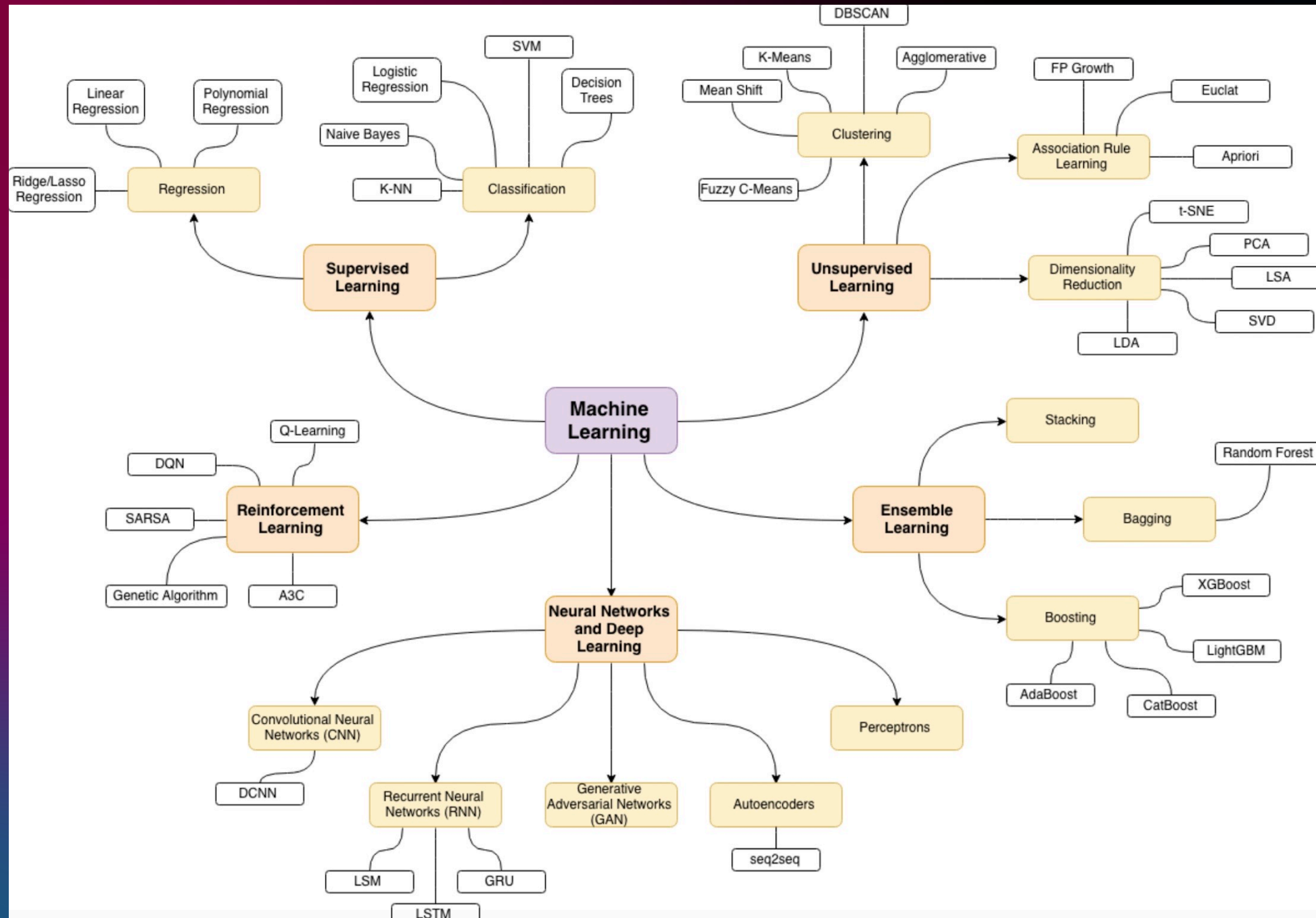      o UCM application to ensemble learning

      o UCEL framework

# HPC and AI convergence

- About ML/DL:

- **1943** : first NN

- **1957** : first NN with training

- **1974-1981** : "silence"

- **1981** : first perceptron multilayer

- **1990** : first CNN-LeCun

- **1997** : first RNN (LSTMs)

- **Circa 2012** :  The flight of ML/DL with **Big Data and computing power**

- **Circa 2017 :** Extension of NLP with transformers

Artificial Intelligence / Machine Learning Classification

http://image.slidesharecdn.com/deepdiveinaimlventurelandscape-150831132221-lva1-app6891/95/deepdive-in-aiml-venture-landscape-by-ajit-nazre-rahul-garg-3-638.jpg?cb=1441027412

https://raw.githubusercontent.com/trekhleb/homemade-machine-learning/master/images/machine-learning-map.png

# Sparse Computation : A common topic of HPLA & ML/DL

Automatic detection of the best sparse compression format as a function of the context (numerical method, parallel programming model, parallel/distributed architecture, etc.).

➢ Auto-tunig (PhD of Hamdi-Labri), **Expert System and then Machine Learning** (PhD of Mehrez).

➢ I. Mehrez, O. Hamdi-Larbi, T. Dufaud, N. Emad. Machine Learning for Optimal Compression Format Prediction on Multiprocessor Platform. HPCS 2018: 213-220.

➢ Ihrak Mehrez. Auto-tuning for automatic detection of the best compression format for sparse matrix (PhD, 2018). University of Paris Saclay, France.

➢ Olfa Hamdi-Larbi. Study of sparse matrix compression for numerical computation on high-scale distributed systems (PhD, 2010). UVSQ.

# Focus on clustering (unsupervised ML)



Example: detect relevant groups based on frequent co-purchasing on Amazon.com

Data: V. Krebs. 2004

Visualization: M. Bastian, S. Heymann, and M. Jacomy. "Gephi: An Open Source Software for exploring and manipulating networks" 2009
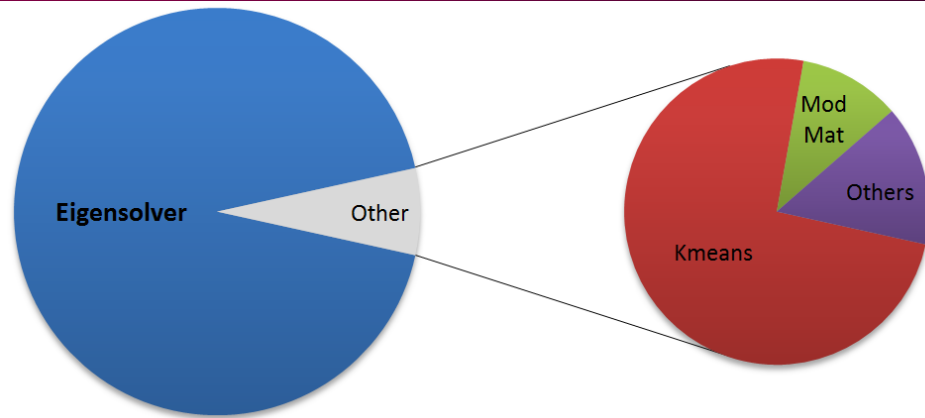
Two main methods allow partitioning vertices $V$ of a graph $G = (V, E)$ in a set of clusters $S_k \subseteq V$ such that $V = \cup_{k=1}^{p} S_k$ are **modularity maximization** and **minimum balanced cut**. This by computing the largest eigenpairs of the modularity matrix or the smallest of the Laplacian matrix.

The multiple implicitly restarted Arnoldi MIRAMns and multiple implicitly restarted Lanczos methods MIRLns with nested subspaces are used.



Pink     Liberal

Yellow   Neutral

Green    Conservative

Data: V. Krebs. 2004

Visualization: M. Bastian, S. Heymann, and M. Jacomy. "Gephi: An Open Source Software for exploring and manipulating networks" 2009

# Focus on Clustering (2)

1. Let $G = (V, E)$ be an input graph and $A$ be its weighted adjacency matrix.
2. Let $p$ be the number of desired clusters.
3. Set the modularity matrix $B = A - \dfrac{1}{2\,\omega} v\, v^T$
4. Find p largest eigenpairs $BU = U\Sigma$, where $\Sigma = diag\left(\lambda_1, \dots, \lambda_p\right)$
5. Scale eigenvectors $U$ by row or by column (optional).
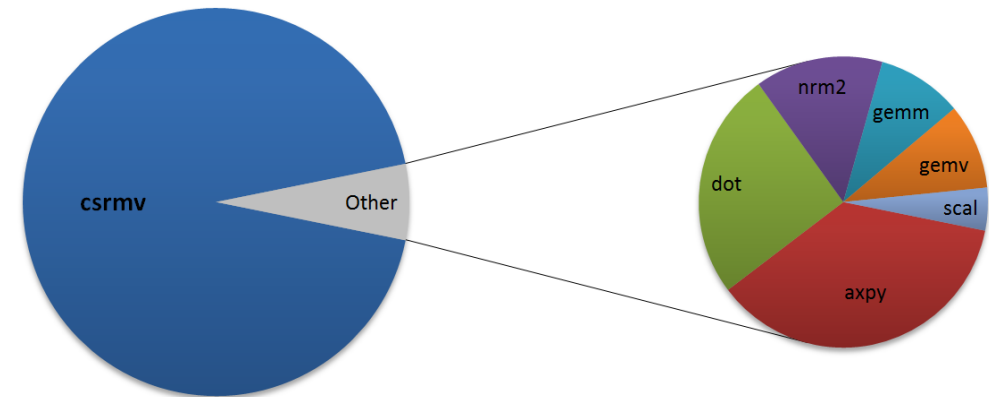6. Run clustering algorithm, such as k-means, on points defined by rows of $U$.

# Focus on Clustering (3)

Profiling: modularity clustering



The eigensolver takes **90%** of the time

The sparse matrix vector multiplication takes **90%** of the time in the eigensolver
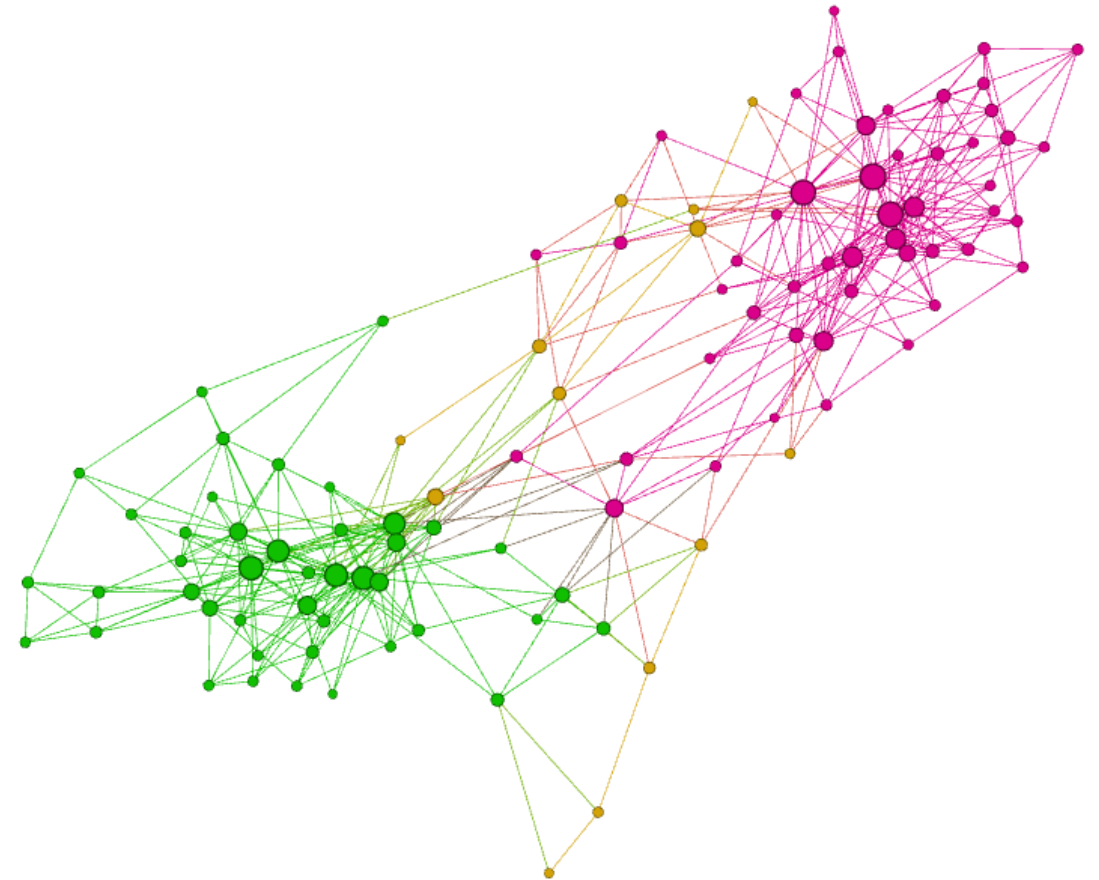
# Focus on clustering (4)



84% hit rate

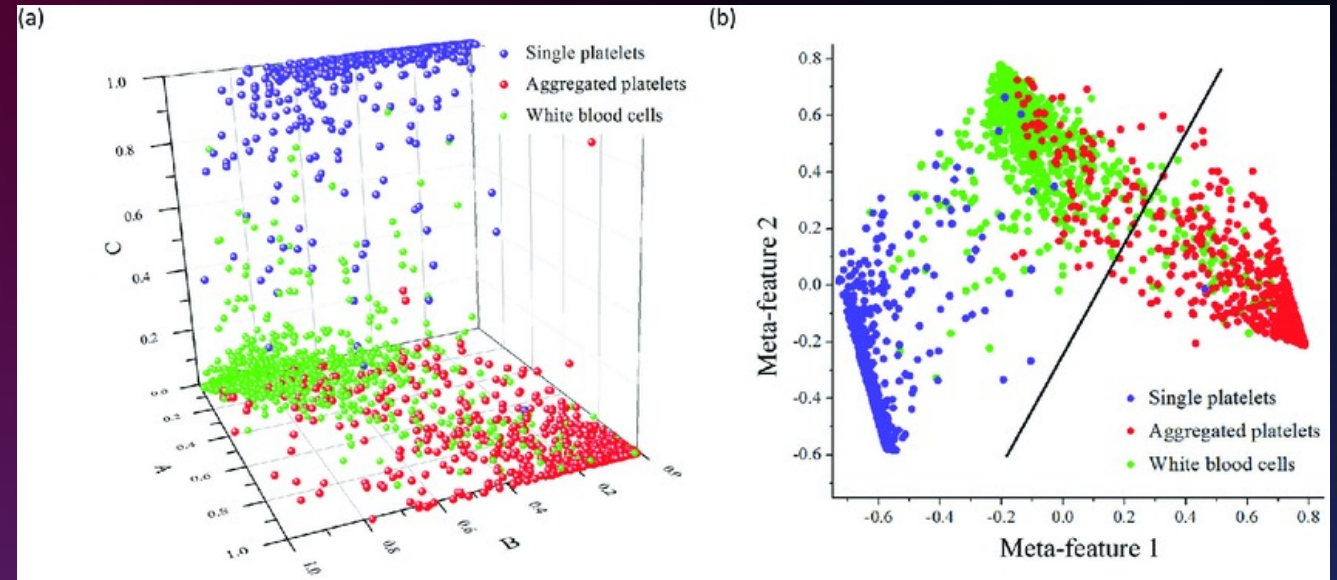Spectral Modularity maximization

Ground truth

A. Fender, N. Emad, S. Petiton, M. Naumov, *Parallel Modularity Clustering*, Procedia Computer Science, Volume 108, 2017, Pages 1793-1802

# Focus on (semi)supervised classification using UC methods

A process of categorizing a given dataset (structured or unstructured) into **predefined classes** (label or categories).

- Classification predictive modeling
- Binary classification
- Multi-class classification
- Multi-label classification
- Imbalanced classification



Jiang, Yiyue et al.*, Label-free detection of aggregated platelets in blood by machine-learning-aided optofluidic time-stretch microscop,*Lab Chip Journal, Vol. 17, 2426-2434, The Royal Society of Chemistry, 2017.
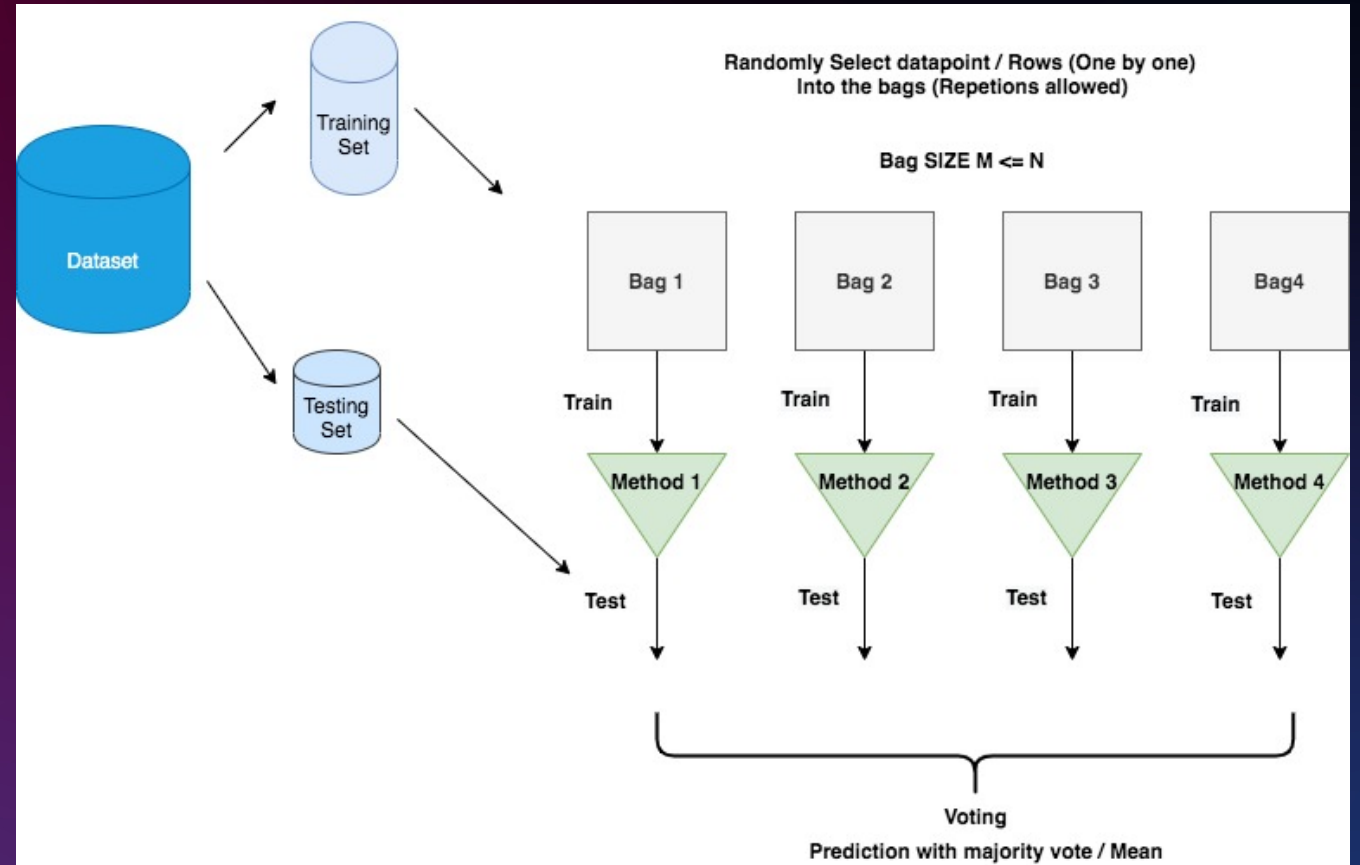
Classification methods: Logistic Regression, naive Bayes, stochastic Gradient Descent, KNN, Decision Tree, Random Forest, ANN, SVM, …

Application UC approach to ensemble learning for classification

# Ensemble Learning methods

**Bagging technique**

1. **Start.** Choose $\ell$ the number of the bags and $m$ ($\leq n$) the size of the bags.
2. **Iterate.** For $i = 1, \ldots, \ell$ do in parallel
   a) *Sampling.* Select the bag $B_i$ by a random sampling technique with replacement on LD,
   b) *Training and testing.* Train a model $L_i$ on the bag $B_i$ and test $L_i$ with TD dataset.
3. **Share.** On all the results of $\ell$ processes, use a selection system (like voting) to get the final prediction.
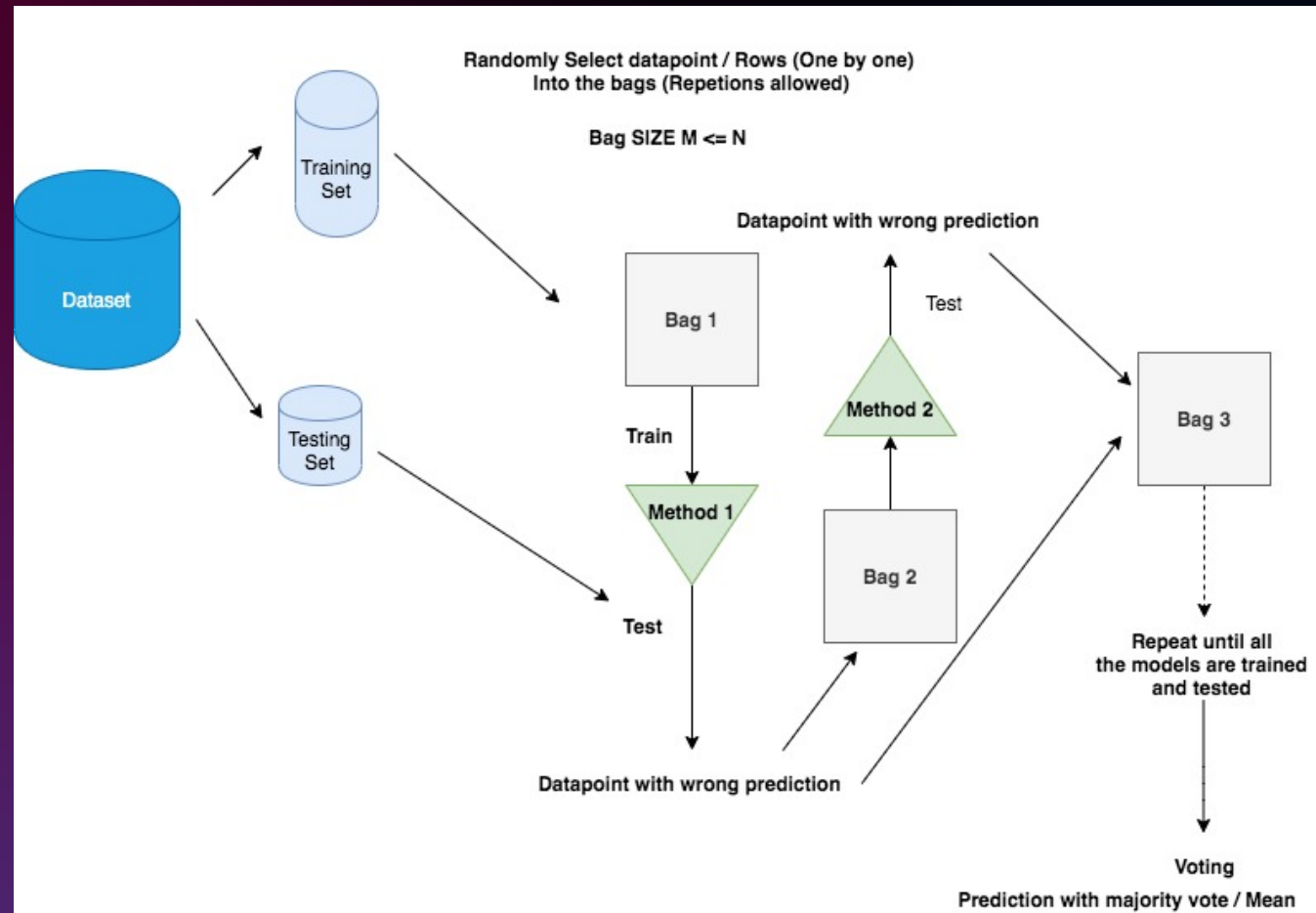


- The same weak learner
- Selection of the result by a voting system
- Intrinsic data parallelism

# Ensemble Learning methods

**Boosting technique**

1. **Start.** Choose $\ell, m$ the number and the size of the bags, the base week leaner $L_1$ and define the bag $B_1$ =LD.
2. **Iterate.** For $i = 1, \dots, \ell$ do
   a) *Training and testing.* Train $L_i$ leaner on dataset $B_i$, produce $L_i$ model, test it on $B_i$ and select $W_i$ the $k_i$-size miss-predicted sub-dataset of LD.
   If (P $(L_i) \geq \theta$) then put $best = i$ and stop.
   b) *Sampling.* Set the bag $B_{i+1}$= (1- $\alpha_i$) $R_i \cup \alpha_i W_i$, where $\alpha_i$ is the weight given to miss-predicted data and $R_i$ is the set of $(m - k_i)$ correctly predicted data in $B_i$ and go to 2.
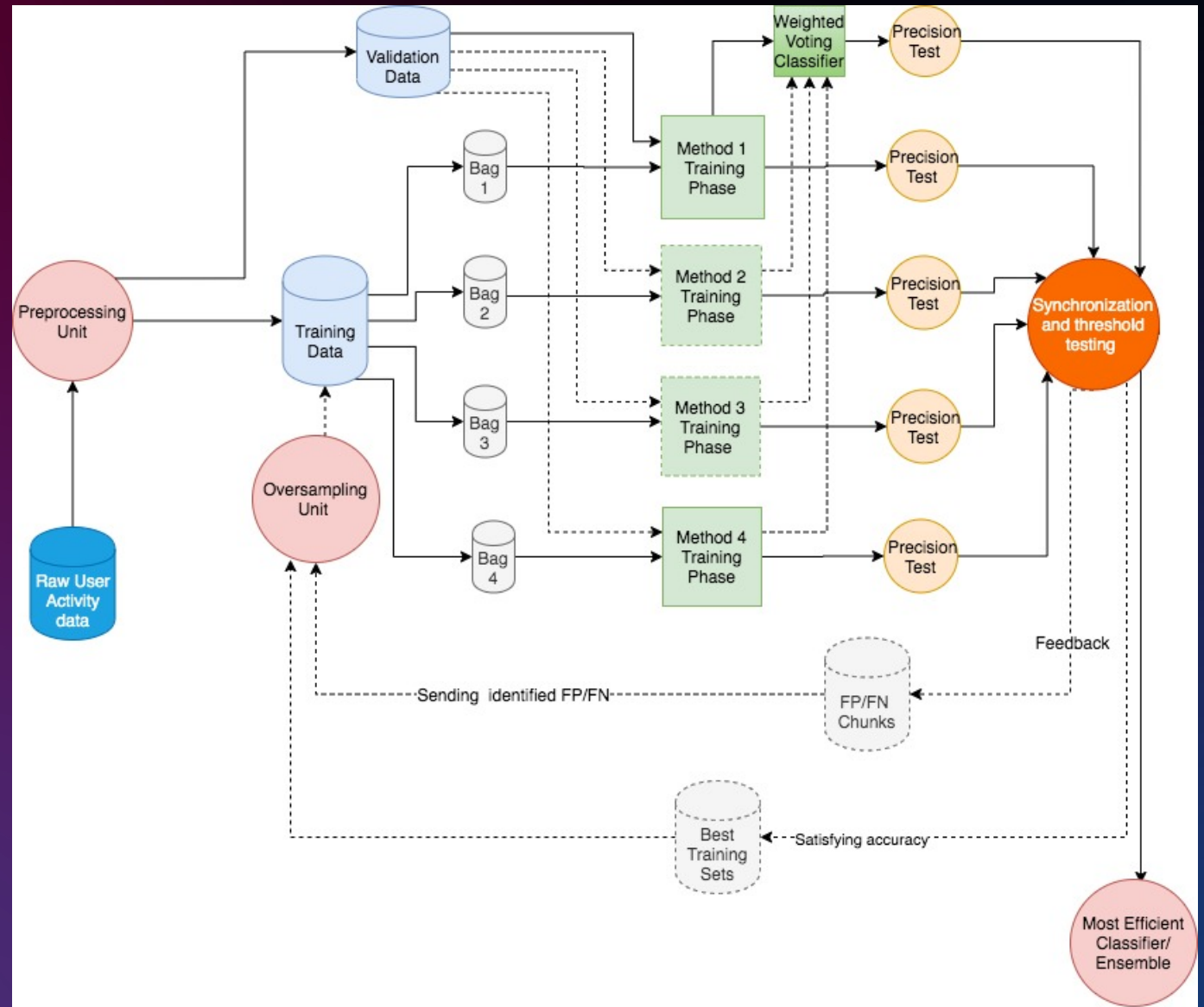3. **Result.** Set $L_{best}$ as a weighted combination of the previous $\ell$ leaners.

- Iterative process
- The miss-predicted data weighted more
- Selection of a weighted combination of the learners



Randomly Select datapoint / Rows (One by one)
Into the bags (Repetions allowed)

Bag SIZE M <= N

Dataset

Training Set

Testing Set

Datapoint with wrong prediction

Bag 1

Test

Method 2

Train

Method 1

Test

Bag 2

Datapoint with wrong prediction

Bag 3

Repeat until all the models are trained and tested
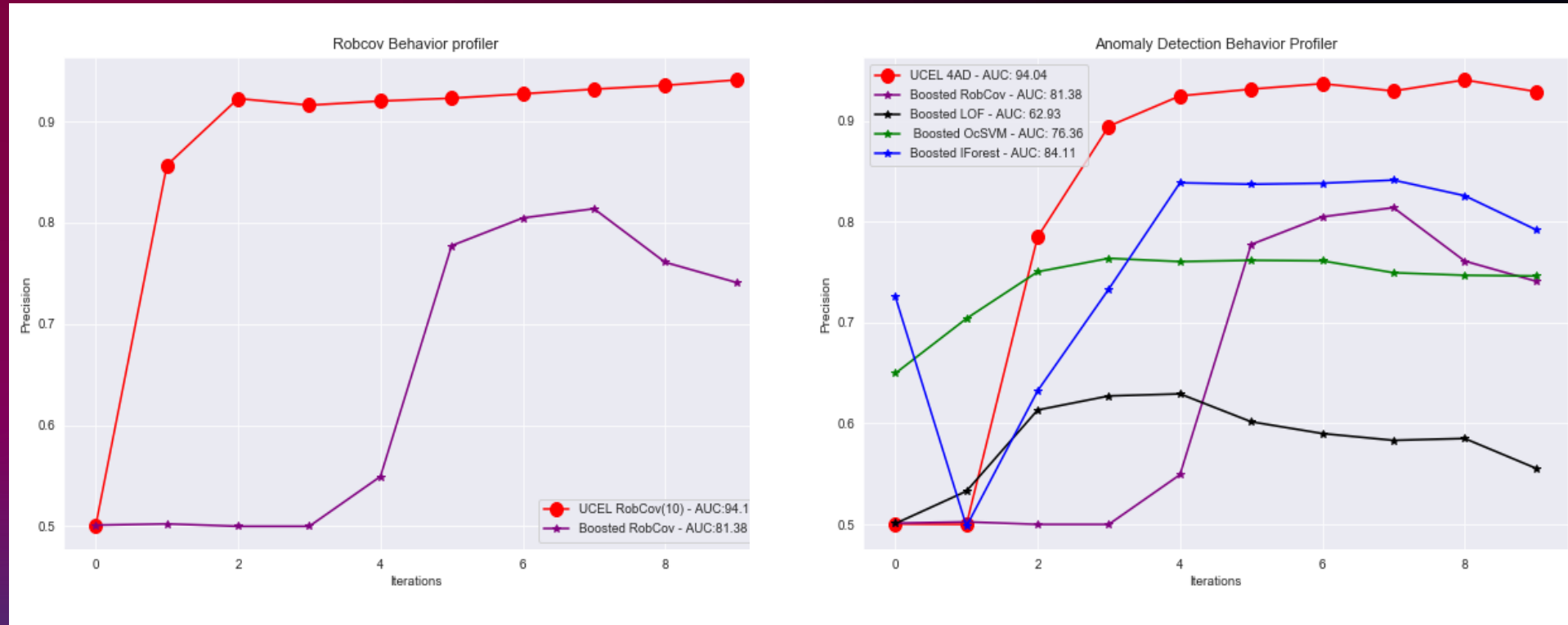
Voting
Prediction with majority vote / Mean

# UCEL: Unite and Conquer and Ensemble Learning

- **Co-methods:** Ensemble base-learners
- **Partial initial parameters:** Bags
- **Restarting strategy** is based on intermediate global result of the learners

- A. Diop, N. Emad, and T. Winter. *A Parallel and Scalable Framework for Insider Threat Detection*. In 27th IEEE International Conference HiPC , 16-19 Dec. 2020, Pune, India.
- A. Diop, N. Emad, and T. Winter. *A Unite and Conquer Based Ensemble Learning Method for User Behavior Modeling*. In 39th IEEE IPCCC Conference, Nov. 6th – 8th, 2020, Austin, Texas, USA.
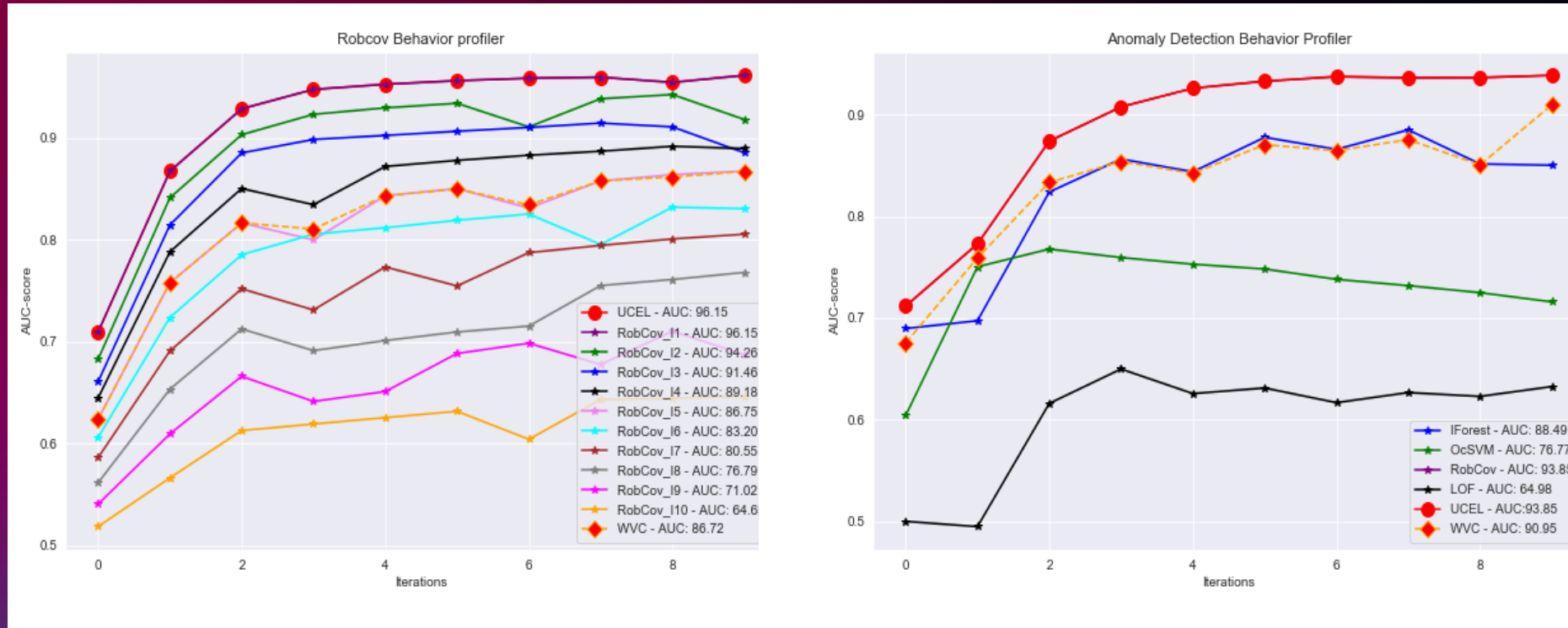
# UCEL performance



**Evolution of the AUC score over the cycles with anomaly detection co-methods**

Left: MRobcov(10) vs individual boosted Robcov,
Right: 4 different anomaly detection methods vs individual boosted co-methods
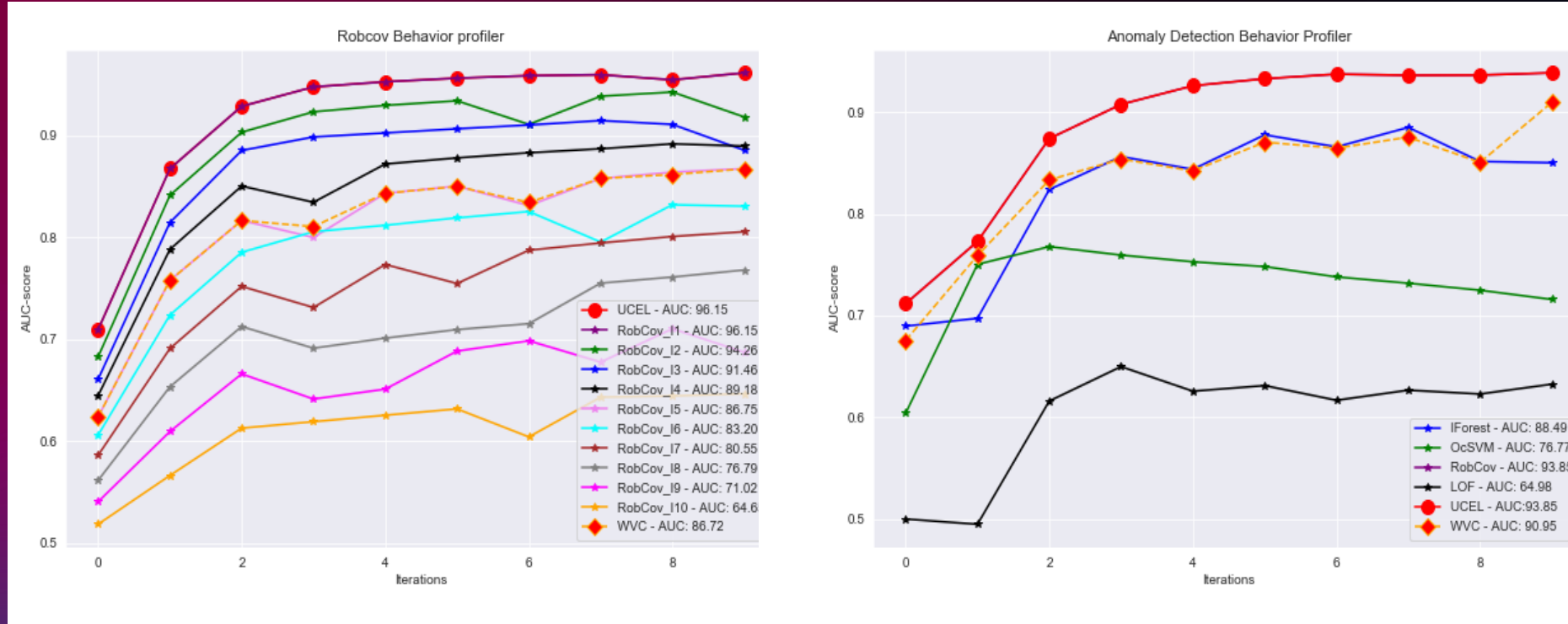
# UCEL performance



(a) multiple Robcov (10)   (b) BP (IForest, OcSVM, Robcov, LOF)

When an individual co-method suffers either from underfitting, overfitting, poor calibration of the hyperparameters or not having enough sample to establish a good decision boundary, UCEL boosts their low AUC-score over the cycles.

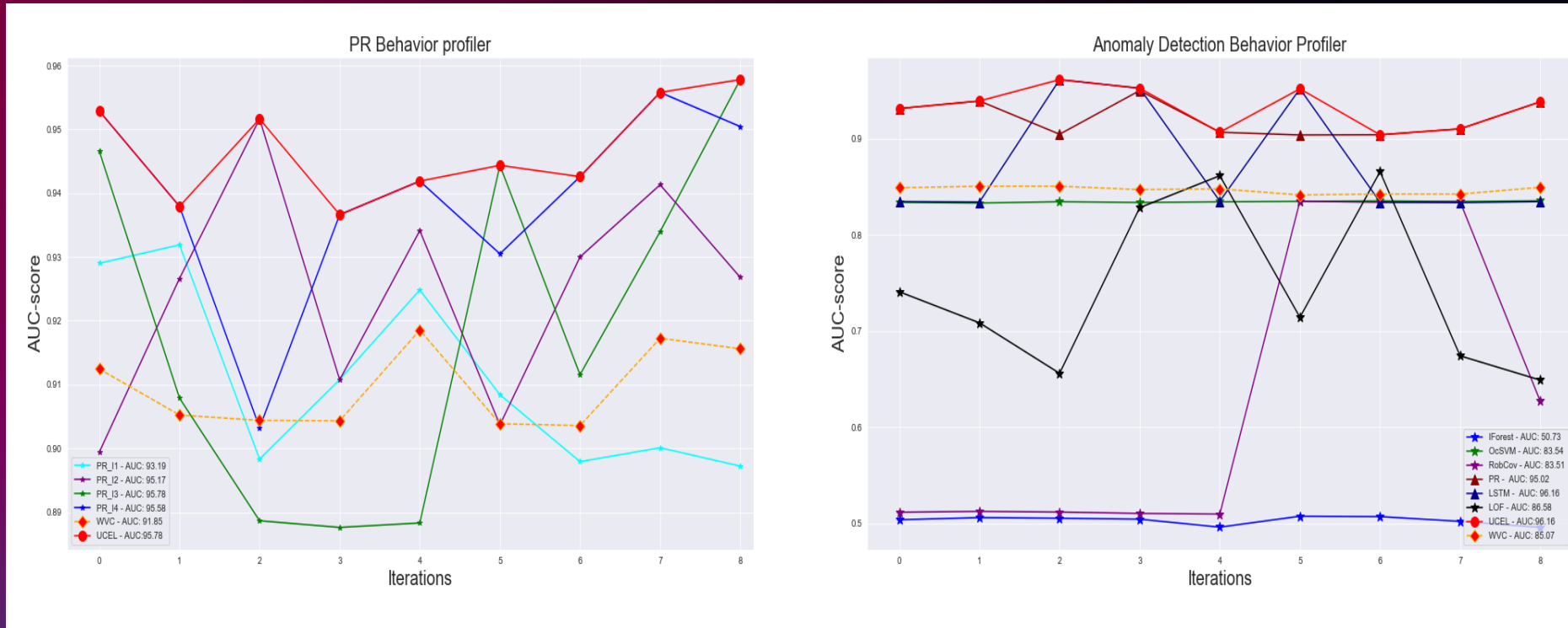# UCEL performance



(a) multiple Robcov (10)          (b) BP (IForest, OcSVM, Robcov, LOF)

Evolution of the AUC score over the cycles with Focus on graph-based anomaly detection using PageRank co-methods

Even a powerful co-method like PR is improved by UCEL.
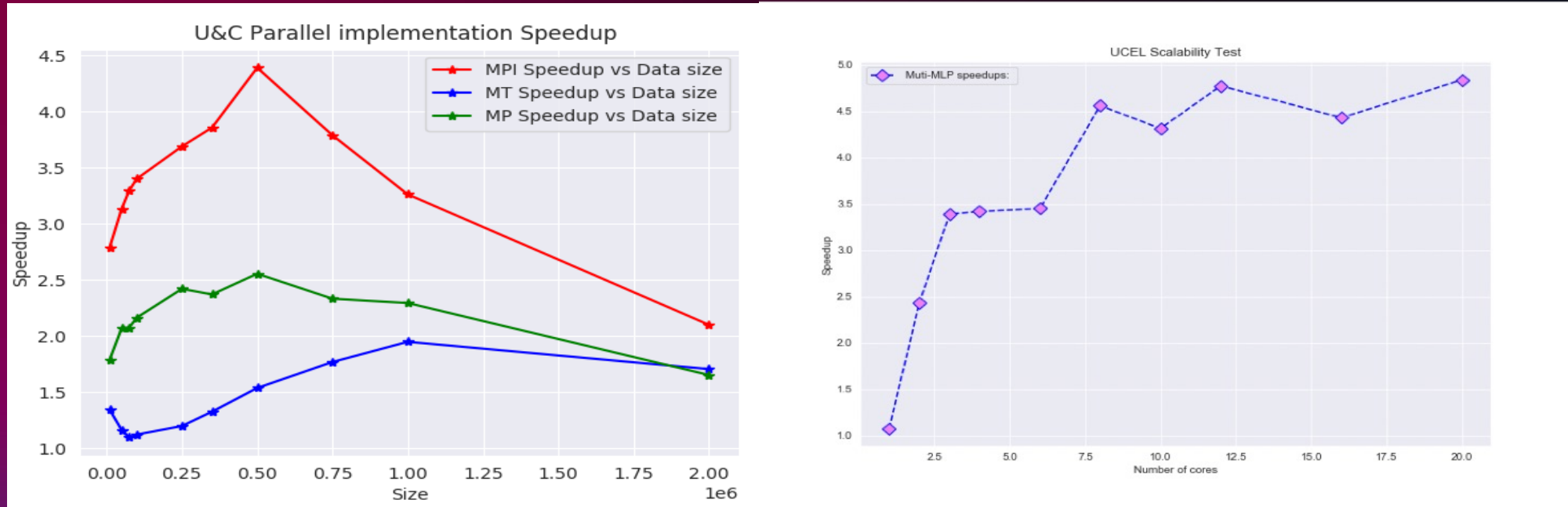
# UCEL performance



(a) MultiplePR(4)  (b) PB(4AD, PR, LSTM)

Evolution of the AUC score over the cycles with Focus on graph-based anomaly detection using PageRank co-methods

Even a powerful co-method like PR is improved by UCEL.

# UCEL performance



**Weak scalability:** due to the synchronization steps in the end of each cycle.

**Strong scalability**: the strong scalability of 10 MLP as co-methods, a dataset size of 500000 entries. The speedup rises from 1 to approximately 4.5 when we increase the number of processing cores.

# Concluding Remarks

- Important impact of HP eigenvalue computation in AI : Be aware of not always using the libraries.

- Interactions between machine learning and linear algebra approaches must be studied more.
  - ➢ UCEL  is a good example, and the approach is extensible.
  - ➢ Well adapted to new and emerging supercomputers

- ML presents a formidable tool bringing a tsunami of solutions to many problems
  - ➢ Experiment data (even when it does not seem interesting) must be saved...