## HIGH-PERFORMANCE DIMENSIONALITY REDUCTION IN LARGE LANGUAGE MODELS

Nahid.Emad@uvsq.fr Maison de la Simulation, Li-Parad University of Paris Saclay, France

First workshop on

Distributed and Parallel Programming for Extreme-scale AI June 16-17th, 2025 | Mines-Paris, PSL University, France

#### DIMENSIONALITY REDUCTION IN AI

Reducing the cost of LLMs through dimensionality reduction techniques that simplify data representation while retaining relevant information

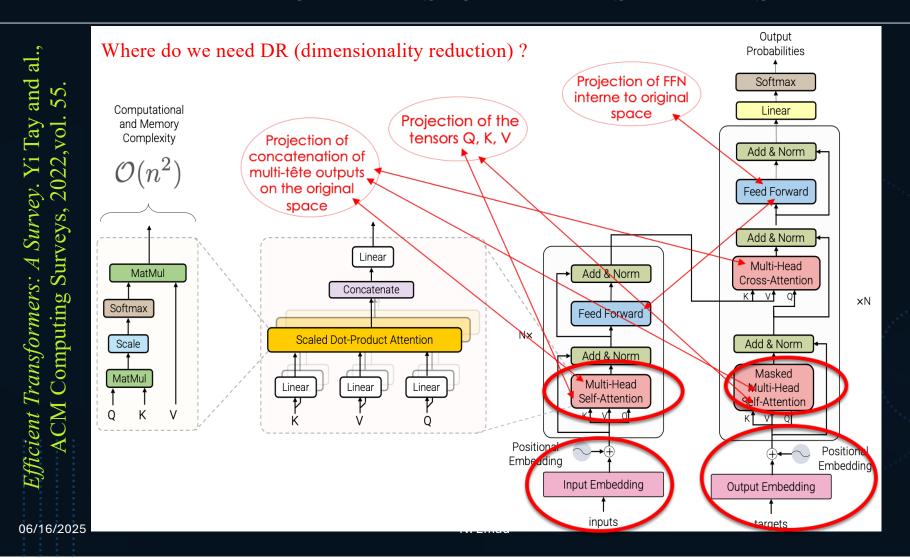
Distilling big/complex/raw data to produce unified/focused/usable ones

- Direct impact on applications as clustering, data visualization, anomaly detection, data compression, data preprocessing, ...
- Optimization of LLMs like Transformer by applying them to almost every layer

Processing reduced data provides accelerated & efficient LLMs

Challenges: Robust distributed dimensionality reduction

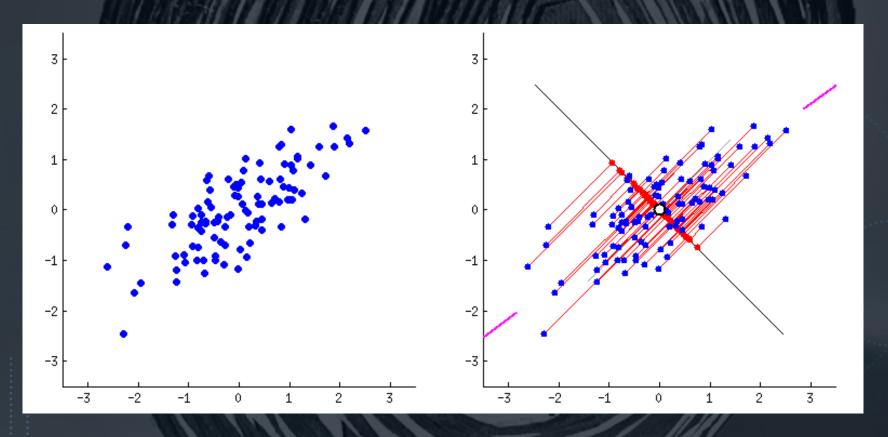
#### DR FOR TRANSFORMER-BASED LLMS



3

### PROJECTION ON SMALLER OR LARGER SPACES

## 1. Projection on smaller subspace



06/16/2025

N. Emad

# An Example application of PCA to a financial data set in <u>WEKA</u>, a Java data mining software <a href="https://youtu.be/OdlNM96sHio?si=9I3OO4dZgRE8pMxJ">https://youtu.be/OdlNM96sHio?si=9I3OO4dZgRE8pMxJ</a>

2. Projection on larger space

#### **O**UTLINE

- A brief overview of DR techniques
- Unite and Conquer methods in high-performance linear algebra
- Multiple implicit restarted Arnoldi method (MIRAMns)
- Concluding remarks

## A BRIEF OVERVIEW OF DR TECHNIQUES

Curse of dimensionality: Problems like increase of data size and parameters, points too far apart making it difficult to estimate data distributions, overfitting risk of the model, loss of the meaning of proximity notions, ...

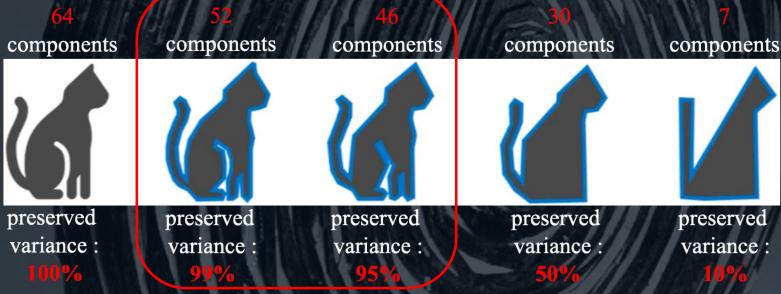
Major challenge: Learn from a simplified data representation of the original dataset to gain more insights from the original dataset.

Methodology: Apply DR techniques to map the raw information into a new feature space where data analysis methods can be used.

06/16/2025

#### DIMENSIONALITY REUCTION

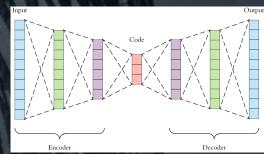
Reducing data by simplifying its representation while keeping relevant information



Guillaume Saint-Cirgue: https://www.youtube.com/c/MachineLearnia

#### DR TECHNIQUES BASED ON NON-LINEAR TRANSFORMATION

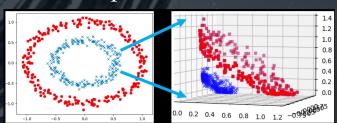
- Autoencoder: Encoding (DR) + decoding for reconstructing data. In case of linear activations or a single hidden layer of sigmoid, the ideal solution for an autoencoder is heavily linked to PCA.
- ✓ Pruning and quantization to compress models after training a mobile application (integration in on-ship SRAM)



- Locally Linear Embedding: Build a weight matrix representing the local reconstruction of each data point from its nearest neighbors, then, computes the eigenvectors of the global weight matrix (PCA)
- ✓ t-distributed Stochastic Neighbor Embedding (PCA + t-SNE),
- ✓ Kernel PCA: Transforms the dataset into a higher dimensional feature space + PCA

✓ ...

Some are often linked to spectral calculation



#### DR TECHNIQUES BASED ON LINEAR TRANSFORMATION

- ✓ Principal/Independent Component Analysis (PCA, ICA),
- ✓ Linear Discriminant Analysis (known #variables classification),
- ✓ Low-Rank approximations : Singular Value and CUR Decompositions (importance of dimensions)
- ✓ Random Projections (data.Random matrix),
- ✓ Partial Least Squares (PLS),
- ✓ t-distributed Stochastic Neighbor Embedding (PCA + t-SNE),

✓ ...

Methods closely related to the Spectral Calculation

## SPECTRAL COMPUTATION WITH LARGE SPARSE MATRICES

For  $A \in \mathbb{C}^{nxn}$ , compute a small number of (less)dominant eigenpairs  $\Lambda_k = (\lambda_1, ..., \lambda_k) \in \mathbb{C}^k$  and  $U_k = (u_1, ..., u_k)$ :  $Au_i = \lambda_i u_i$  for  $i \in [1, k]$   $(P_n)$ 

The vectors  $u_1, ..., u_k$  are the privileged axes, according to which the matrix application behaves like a dilation, multiplying the vectors by a constant.

This dilation ratio is called eigenvalue, the vectors to which it applies are

called eigenvectors, united in an eigenspace.

The calculation of the eigenelements of the matrices is commonly called, spectral calculation.

$$A \begin{bmatrix} u_1, u_2 \end{bmatrix} = \begin{pmatrix} 0.5 & 0 \\ 0 & 3 \end{pmatrix} \begin{bmatrix} u_1, u_2 \end{bmatrix}$$

#### ITERATIVE PROJECTION METHODS LARGE SPARSE EIGENPROBLEM

Random choice of an initial subspace  $\mathbb{K}_0$ ,

For i = 0, 1, ... until convergence do:

nalicitly Restarted Arnoldi Method

- 1. Projection of the problem onto  $\mathbb{K}_i$
- 2. Solve the projected problem in the subspace
- 3. Use 2 to compute approximated solutions in original space
- 4. If no convergence, with a better subspace go to 1 Improvement of  $\mathbb{K}_m(A, v)$  by that of v

A very commonly used subspace is Krylov subspace:

 $\mathbb{K}_m(A, v) = span(v, Av, ..., A^{m-1}v)$  with  $(k \le m \ll n)$ 

#### Unite and Conquer Methods in HP Linear Algebra

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

	C				
Unite and Conquer method	Bi-Lanczos	+	Bi-Lanczos	+	•••
	ERAM	+	ERAM	+	
	ERAM	+	GMRES	+	
	Projection method	+	Convergence accelerator	+	
	//	+		+	•••

#### Before with hybrid methods:

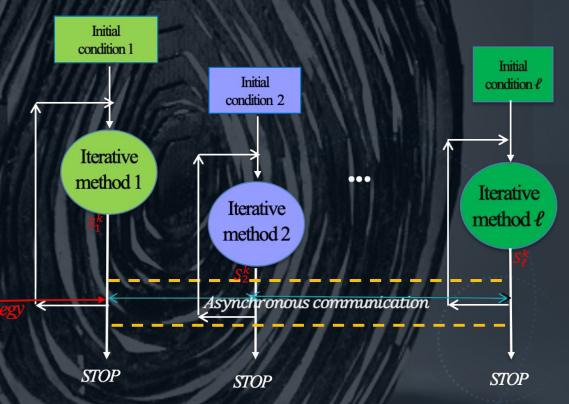
Y. Saad, *Chebyshev acceleration techniques for solving nonsymmetric eigenvalue; 1984*, C. Brezinski Hybrid procedures for solving linear systems; 1994), Code coupling (in simulation), ...

### Unite and Conquer Methods in HP Linear Algebra

#### Characteristics of UC methods

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

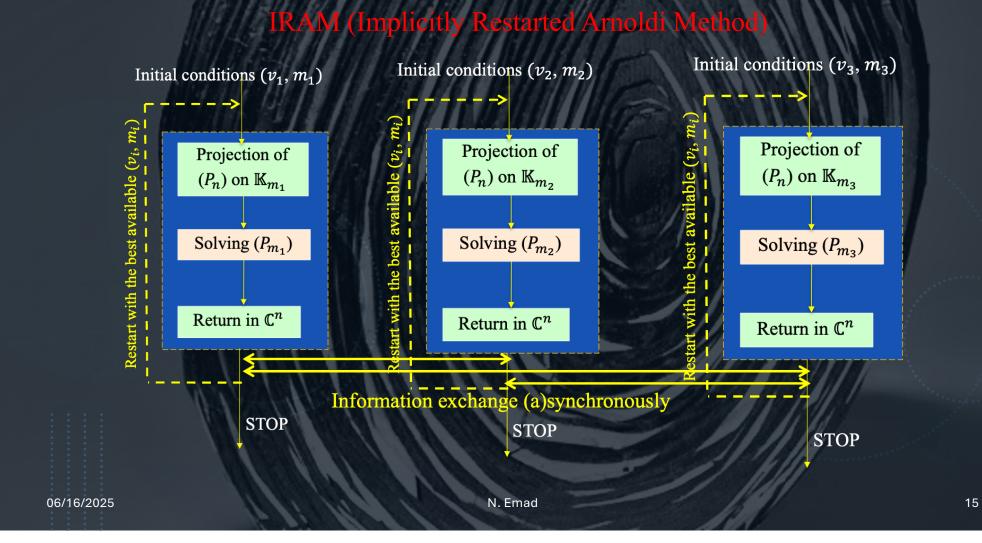




Well suited to high-scale computing systems

N. Emad, S. Petiton. Unite and Conquer Approach for High Scale Numerical Computing, Journal of Computational Science, ISSN-1877-7503, 2016

#### A SPECIFIC CASE OF UC METHODS: MULTIPLE-METHODS



## HOW TO FURTHER IMPROVE SUBSPACE QUALITY?

Improving  $\mathbb{K}_m(A, \mathbf{v}) = span(\mathbf{v}, Av, ..., A^{m-1}\mathbf{v})$  using both m and  $\mathbf{v}$ 

#### Different subspaces

```
\mathbb{K}_{m_i}(A, \mathbf{v_i}) = span \ (v_i, Av, \dots, A^{m_i-1}v_i), with m_i \neq m_j and v_i \neq v_j (for i, j \in [1, \dots, \ell] and i \neq j)
```

Nested subspaces  $(m_1 < m_2 < ... < m_\ell)$ 

$$\mathbb{K}_{m}(A, v) =$$

$$span(v, Av, ..., A^{m_{1}-1}v, A^{m_{1}}v, ..., A^{m_{2}-1}v, ..., A^{m_{\ell-1}}v, ..., A^{m_{\ell}}v),$$
with  $\mathbb{K}_{m_{1}} \subset \mathbb{K}_{m_{2}} \subset ... \subset \mathbb{K}_{m_{\ell}}.$ 

06/16/2025

#### PERFORMANCE OF UNITE AND CONQUER METHODS

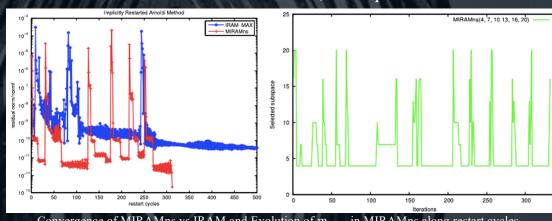
even for difficult input

the best subspace in which we search for principal axes is not always the largest.

Track for researching the number of clusters in unsupervised learning

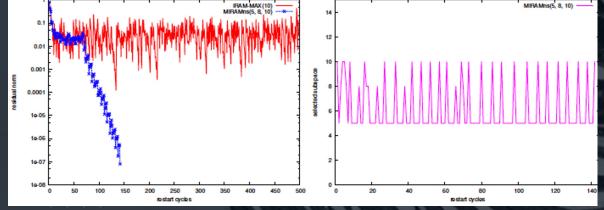
Iterations/collaboration concept in AI algorithms like EL or MixExpert



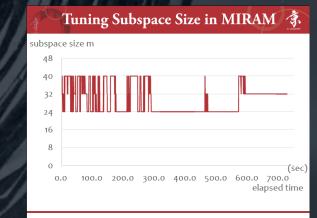


Convergence of MIRAMns vs IRAM and Evolution of  $m_{best}$  in MIRAMns along restart cycles

#### Autotuned IRAM and Evolution of $m_{best}$ within ARPACK



Bfw782a, k = 2,  $tol = 10^{-8}$  and a random initial guess



Schenk/nlpkkt240, n=27993600

06/16/2025

N. Emad

#### APPLICATION OF DDR WITH MIRAMNS IN NN

					ARCHA 18 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	
	Date Fruit		Fetal Health		Radar	
	Red Dim: $34 \rightarrow 5$		Red Dim: $21 \rightarrow 5$		Red Dim: $175 \rightarrow 12$	
	No Emb.	Emb.	No Emb.	Emb.	No Emb.	Emb.
Number of parameters	939,015	233,991	1,063,427	233,475	1,143,815	233,991
Training Exec. time (s)	2.085	1.354	3.7349	2.371	1112.81	716.668
Train accuracy (%)	97.63	93.04	95.94	91.88	98.13	97.95
Evaluation accuracy (%)	94.44	89.99	93.66	90.84	98.79	98.36

#### Performance with and without dimension reduction for different data

- A quarter reduction in network size (at least) while minimizing loss of precision.
- The performance improvement is an increasing function of the size of the input (implying reduction in large model training time).

Work in progress with Mines-Paris & Huawei (Chong Li, Quentin Petit, ...)

#### CLUSTERING AND ANOMALY DETECTION WITH UC METHODS

#### Clustering

NVIDIA (nvGraph library)
MIRAMns, Mlnzos ns

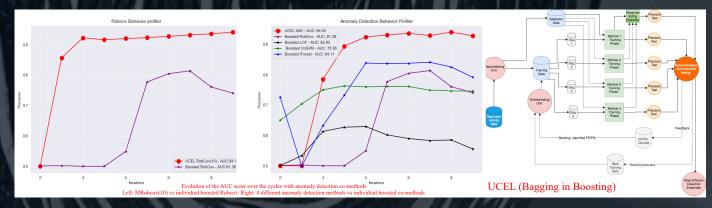
Profiling: modularity clustering

The eigensolver takes 90% of the time

The sparse matrix vector multiplication takes 90% of the time in the eigensolver

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

# Anomaly detection UCEL (ATOS) UC2B Generalization of UCEL(NUMERYX) Z. Ziani (ANOTHER BRAIN)



A. Diop, et al. A Unite and Conquer Based Ensemble Learning Method for User Behavior Modeling, IEEE IPCCC, 2020

#### CONCLUDING REMARKS

- Nonlinear approaches often use linear approaches
- Need to improve spectral computation techniques
- The extensibility of the UC approach opens multiple avenues for new techniques
  - Excellent convergence properties even for difficult input
  - Big finding: a key to choose the best search subspace to mine for dominant eeigenespace
    - This best subspace for searching principal axes is not always the largest.
- Using the concept of iterations/collaboration in AI algorithms like EL or MoE

06/16/2025

#### CONCLUDING REMARKS

Trade-off between reducing costs and maintaining the quality of results

Two important challenges in dimensionality reduction:

- A. What is the optimal number of principal axes to search?
- B. What is the optimal size of the search space for A?

06/16/2025