



Tensor Core Computing: an example on Independent Component Analysis

May 2024

Zeyu Wang, Zoltan Juhasz

University of Pannonia
Veszprem, Hungary



Content outline

1. Tensor core

- 1.1 Abstract and hierarchy
- 1.2 Programming paradigm

2. Independent Component Analysis (ICA)

- 2.1 Algorithm flow
- 2.2 Parallel implementation
- 2.3 Variables analysis

3. Implementation strategy

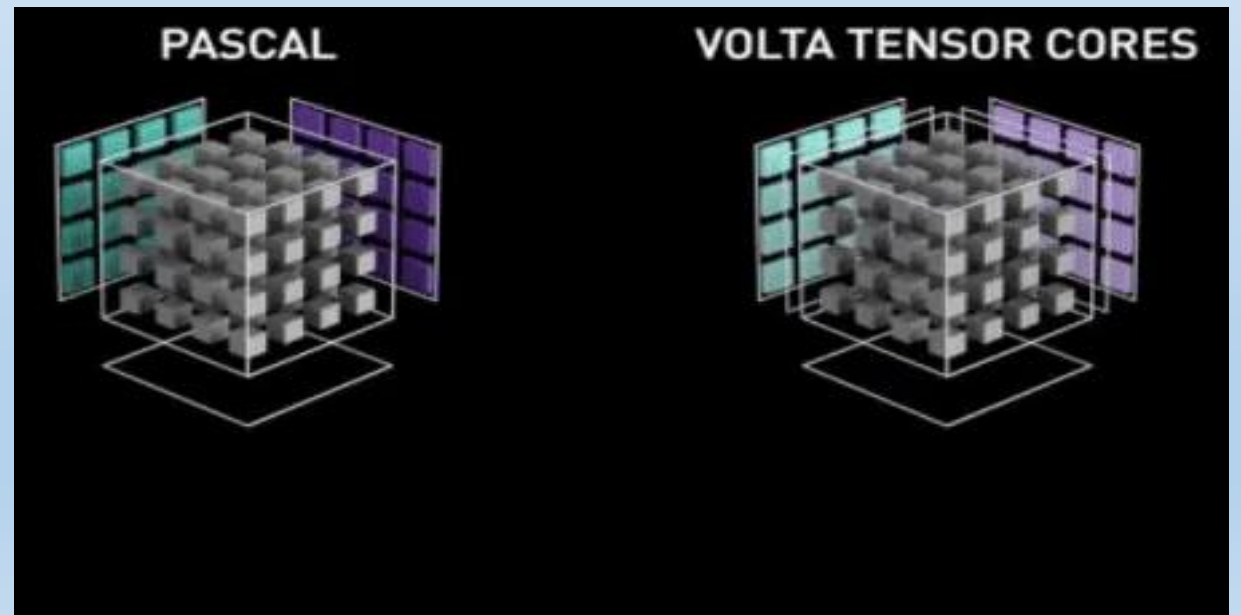
- 3.1 Locating matrix multiplication
- 3.2 Numerical correctness
- 3.3 Performance preview

4. Future works

1.1 Abstract and hierarchy of the tensor core



$$\begin{matrix}
 & & \text{4x4} & & & & \text{4x4} & & & & & \text{4x4} \\
 & & \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{matrix} & * & & \begin{matrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{matrix} & = & & \begin{matrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{matrix} \\
 & & \text{A} & & & & \text{B} & & & & & \text{C}
 \end{matrix}$$

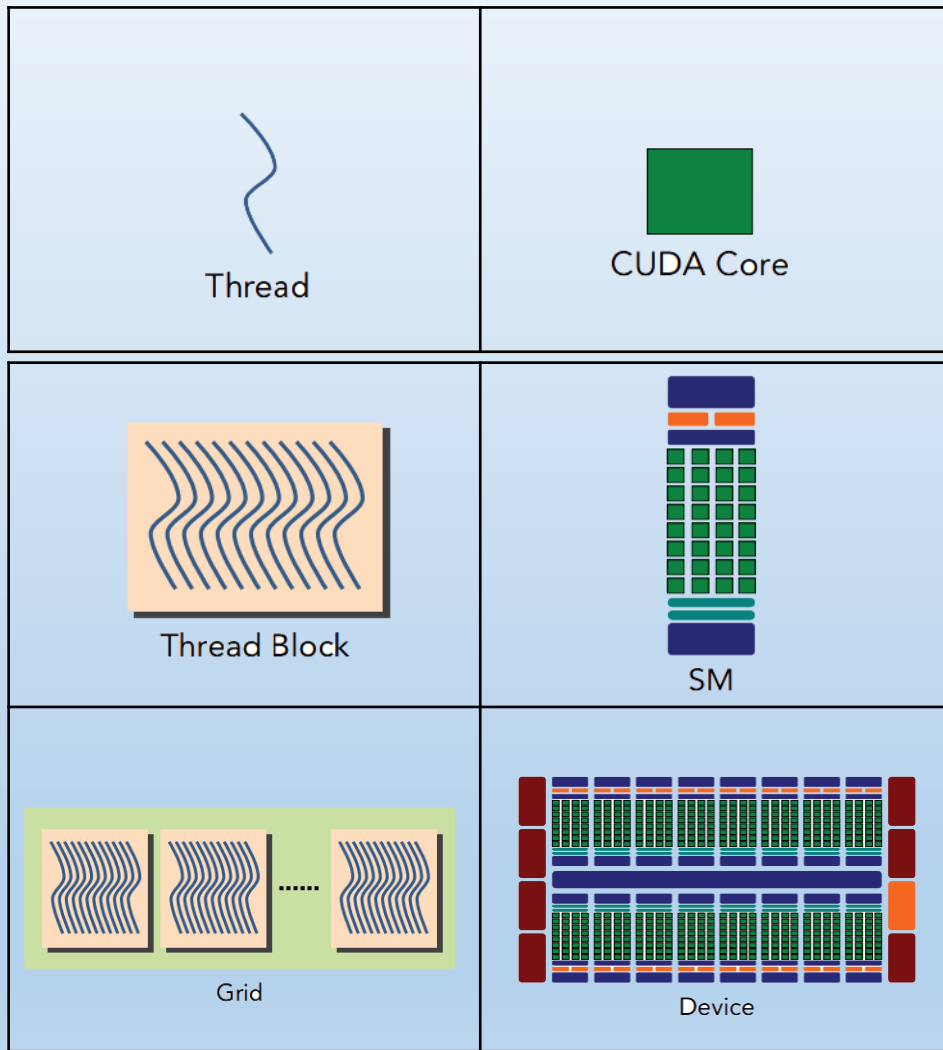


The structure of GA100 SM

Some GPUs equipped with tensor cores

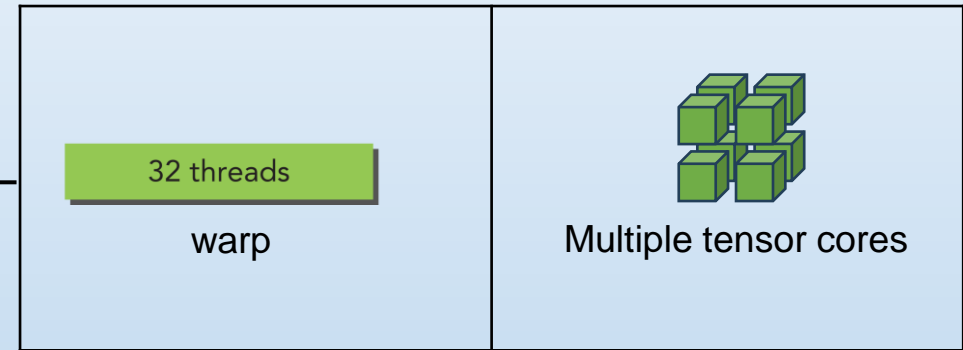
Model	Micro-architecture	Launch	CUDA cores	Memory size (GB)	Processing power (TFLOPS)		
					Tensor core FP32	Single precision	Double precision
V100	Volta	Jun 21, 2017	5120	16/32	112.2	14.0	7.0
T4	Turing	Sept 12, 2018	2560	16	64.8	8.1	0.2
A100	Ampere	May 14, 2020	6912	40/80	312.0	19.5	9.7
H100	Hopper	Mar 22, 2022	14592	80	756.4	51.2	25.6

1.1 Abstract and hierarchy of the tensor core



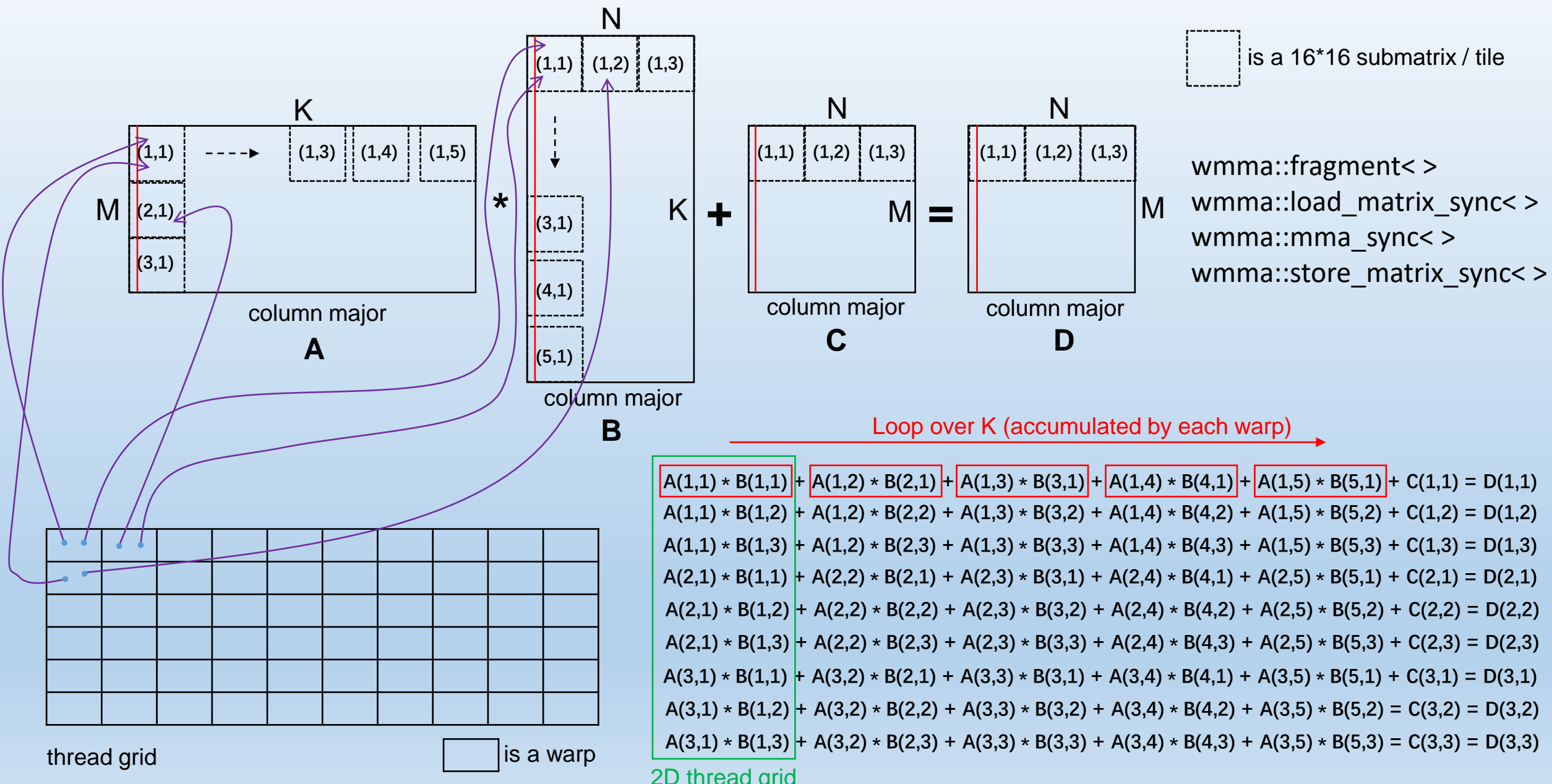
Software hierarchy

Hardware hierarchy

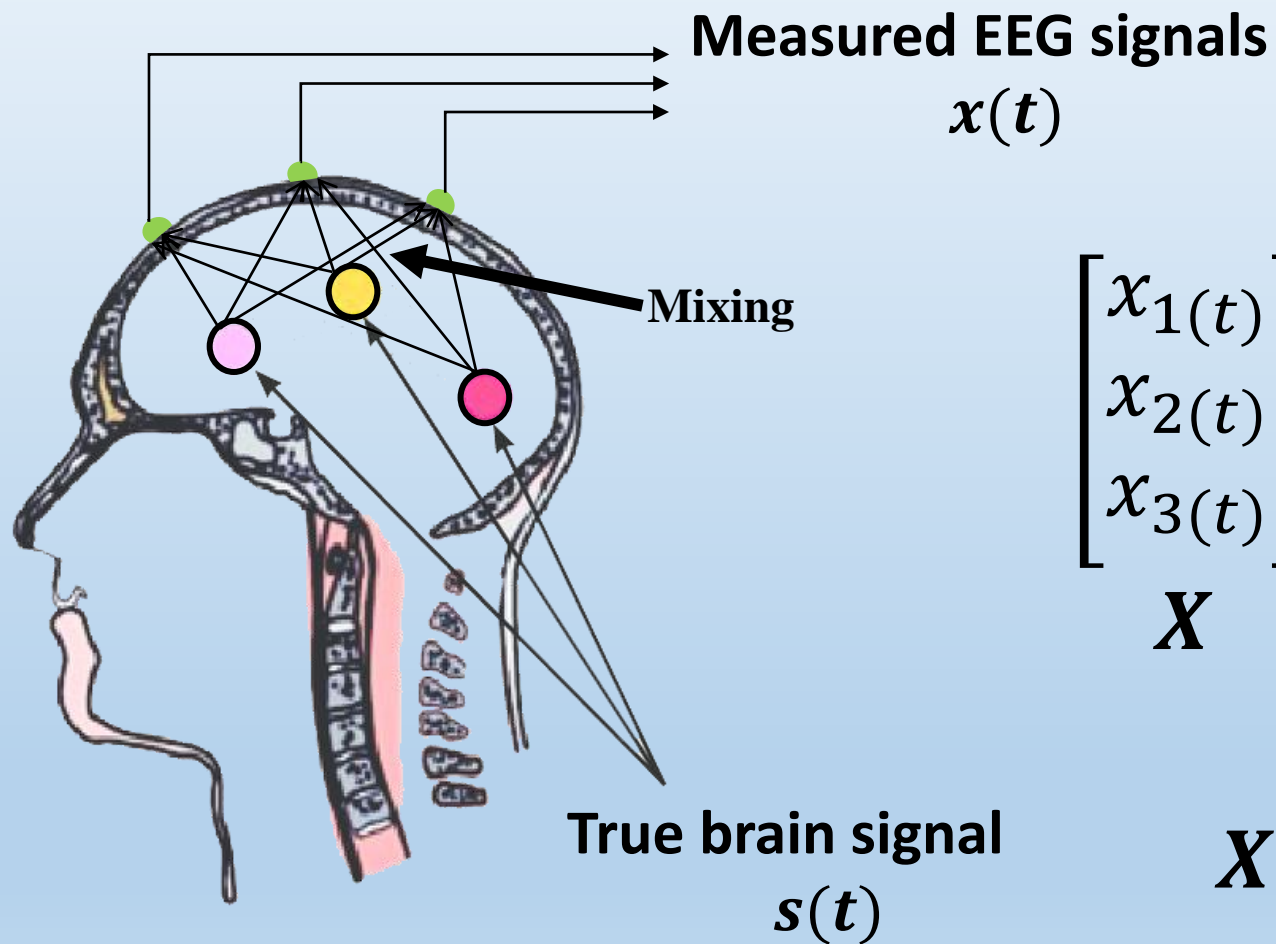


$$\begin{array}{c}
 \text{16x16} \\
 \begin{array}{|c|c|c|c|}
 \hline
 A_{0,0} & A_{0,1} & A_{0,\dots} & A_{0,15} \\
 \hline
 A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\
 \hline
 A_{\dots,0} & A_{2,1} & A_{2,2} & A_{2,3} \\
 \hline
 A_{15,0} & A_{3,1} & A_{3,2} & A_{3,3} \\
 \hline
 \end{array} \\
 \mathbf{A}
 \end{array}
 *
 \begin{array}{c}
 \text{16x16} \\
 \begin{array}{|c|c|c|c|}
 \hline
 B_{0,0} & B_{0,1} & B_{0,2} & B_{0,15} \\
 \hline
 B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\
 \hline
 B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\
 \hline
 B_{15,0} & B_{3,1} & B_{3,2} & B_{3,3} \\
 \hline
 \end{array} \\
 \mathbf{B}
 \end{array}
 =
 \begin{array}{c}
 \text{16x16} \\
 \begin{array}{|c|c|c|c|}
 \hline
 C_{0,0} & C_{0,1} & C_{0,2} & C_{0,15} \\
 \hline
 C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\
 \hline
 C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\
 \hline
 C_{15,0} & C_{3,1} & C_{3,2} & C_{3,3} \\
 \hline
 \end{array} \\
 \mathbf{C}
 \end{array}$$

1.2 Tensor core programming paradigm



2.1 Independent Component Analysis (ICA)



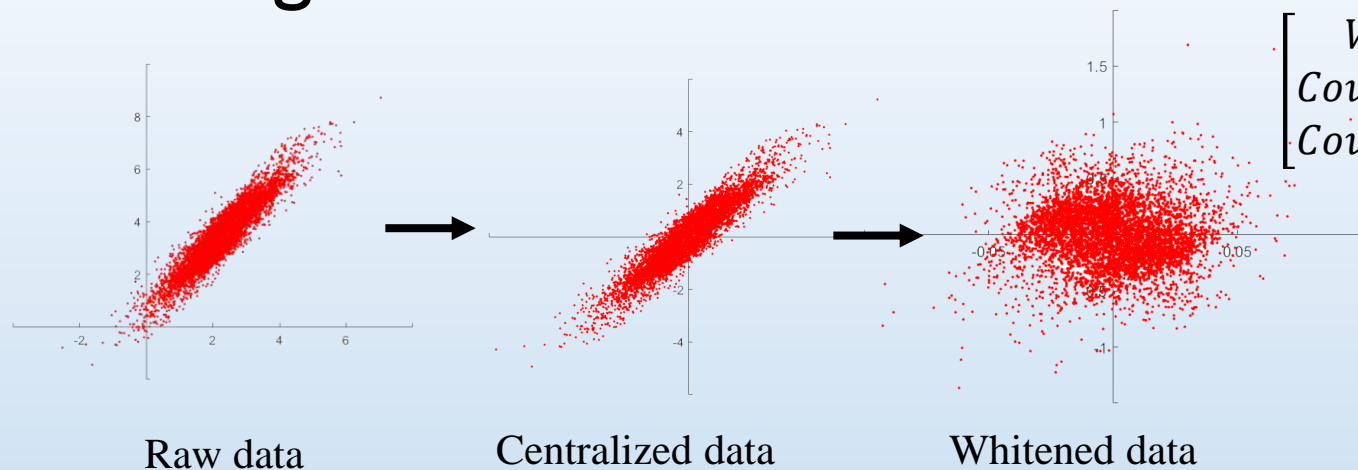
$$\begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{32} \end{bmatrix} \cdot \begin{bmatrix} s_1(t) \\ s_2(t) \\ s_3(t) \end{bmatrix}$$

$\mathbf{X} \qquad \qquad \mathbf{A} \qquad \qquad \mathbf{S}$

Unknown

$$\mathbf{X} = \mathbf{A} \cdot \mathbf{S} \implies \begin{cases} \mathbf{S} = \mathbf{W} \cdot \mathbf{X} \\ \mathbf{W} = \mathbf{A}^{-1} \end{cases}$$

2.1 Algorithm flow



$$Cov = XX^T$$

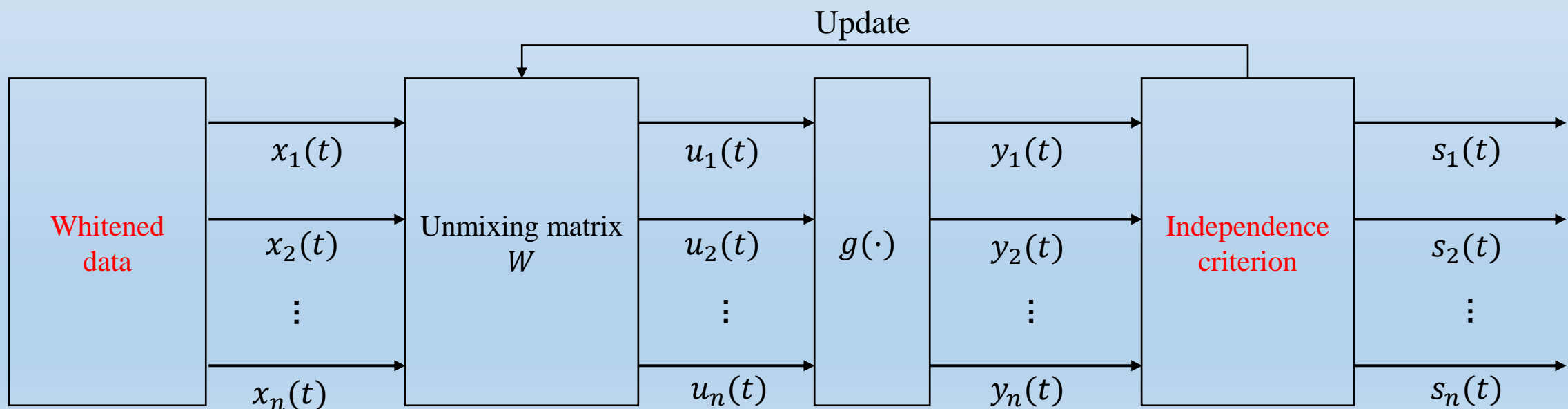
$$\begin{bmatrix} Var(CH1) & Cov(CH1, CH2) & Cov(CH1, CH3) \\ Cov(CH2, CH1) & Var(CH2) & Cov(CH2, CH3) \\ Cov(CH3, CH1) & Cov(CH3, CH2) & Var(CH3) \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Singular Value Decomposition

$$X = U\Sigma V^T$$

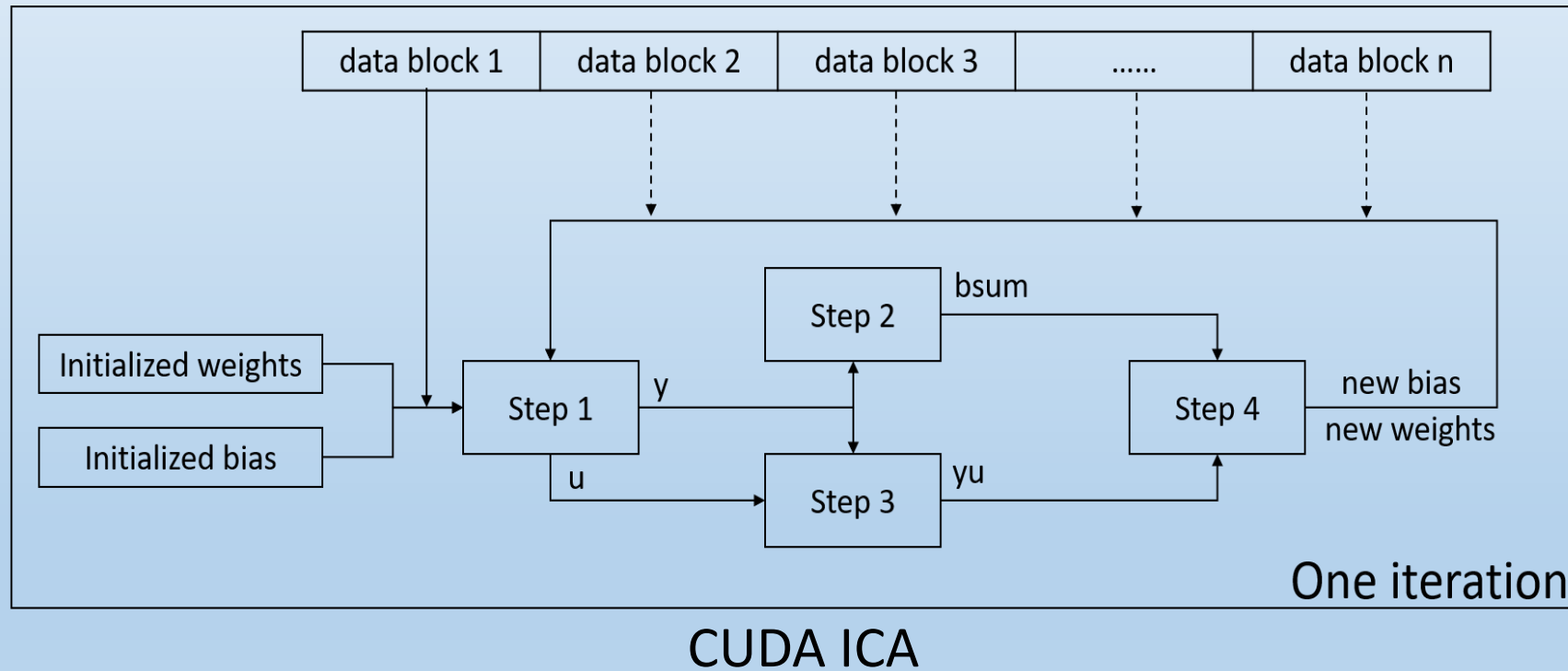
$$\Sigma = \text{diag}(\sigma_1, \sigma_2, \sigma_3, \dots)$$

$$X_{deco} = \sqrt{U^T X}$$



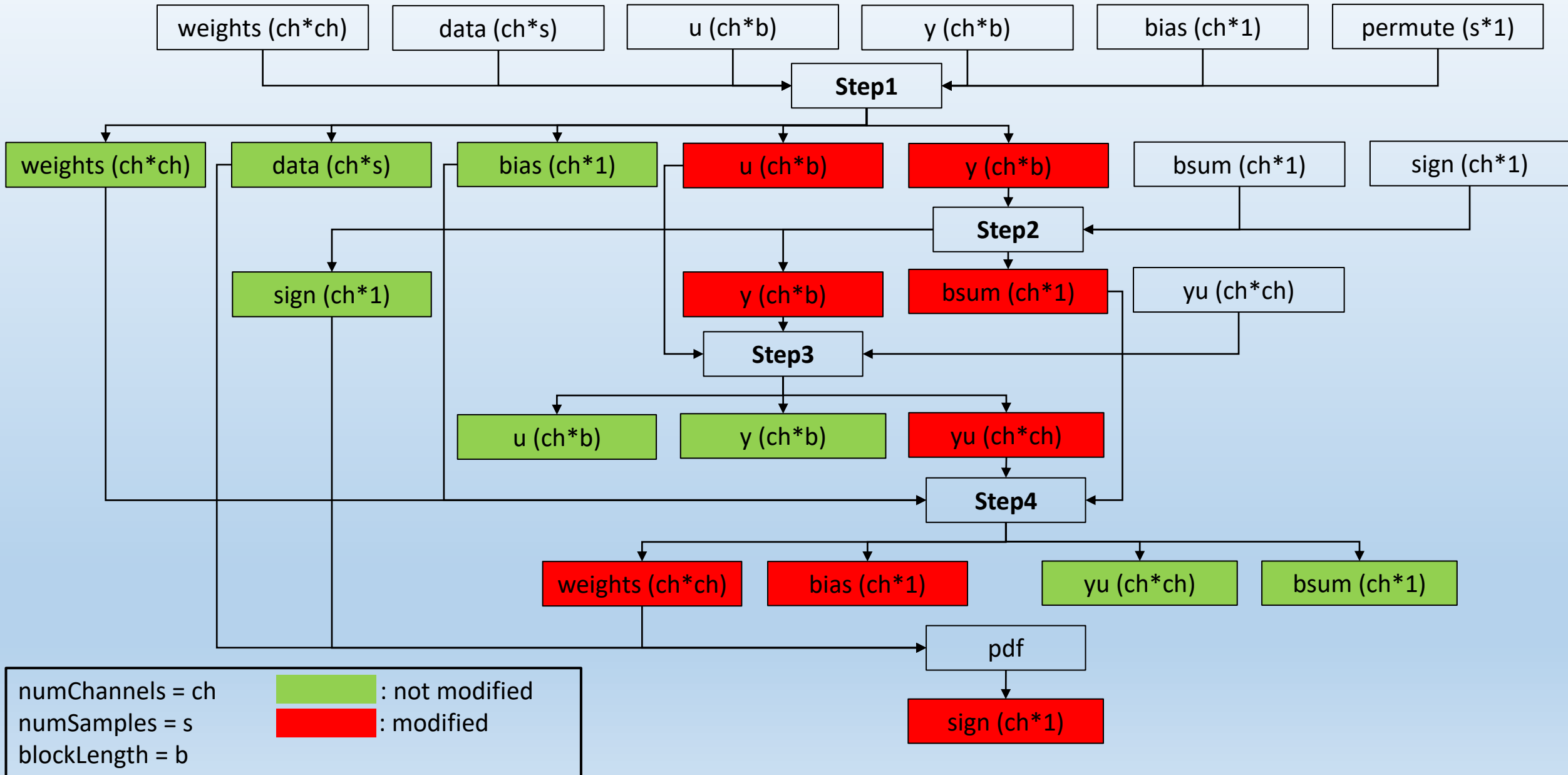
The process of solving the unmixing matrix

2.2 Parallel implementation

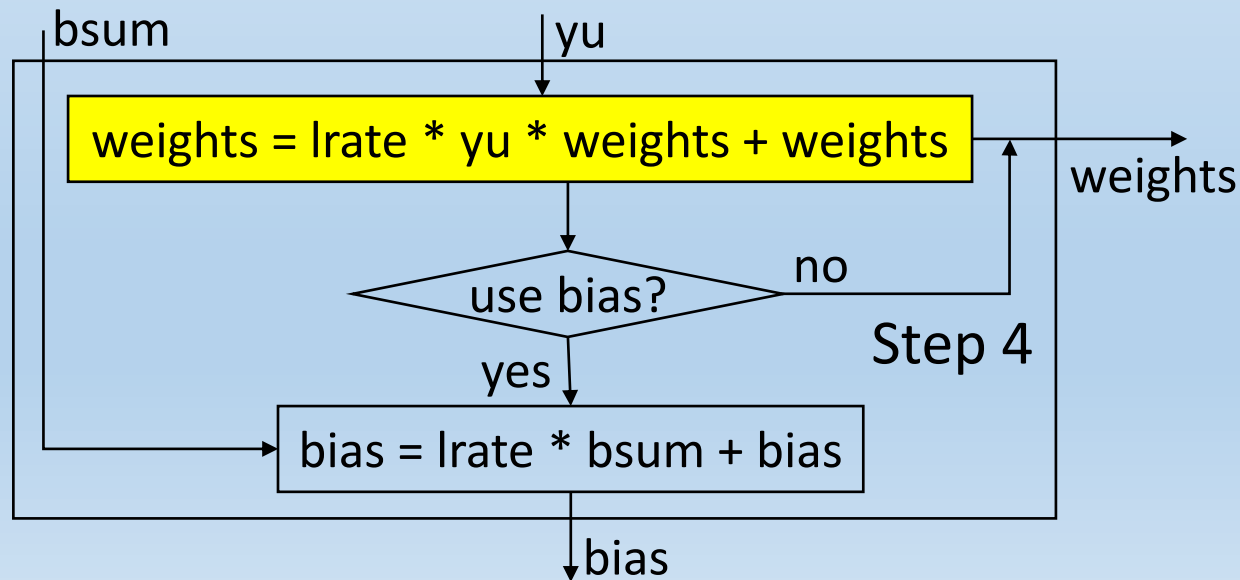
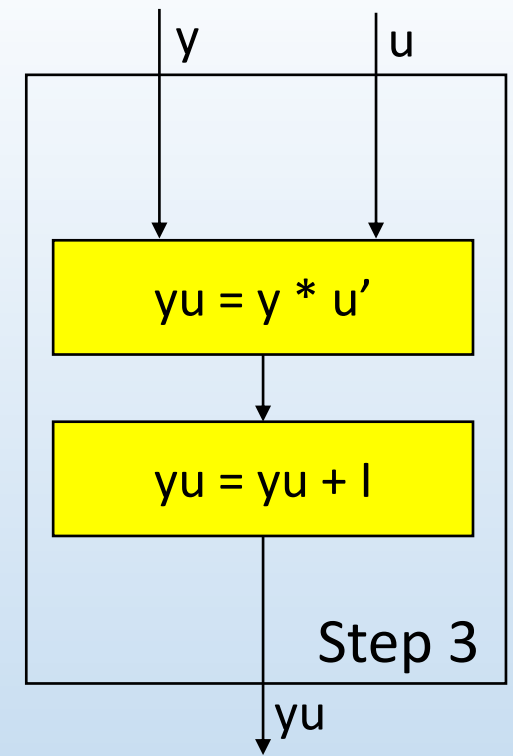
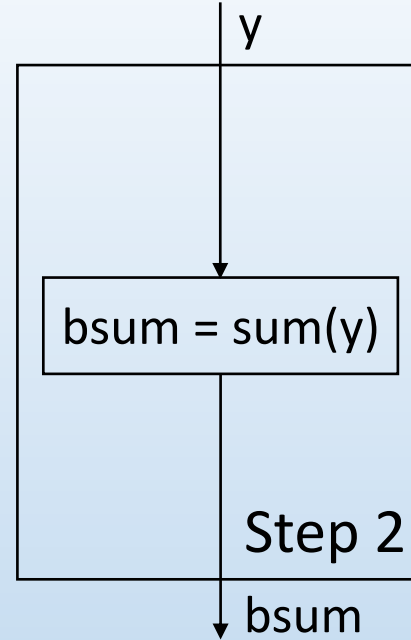
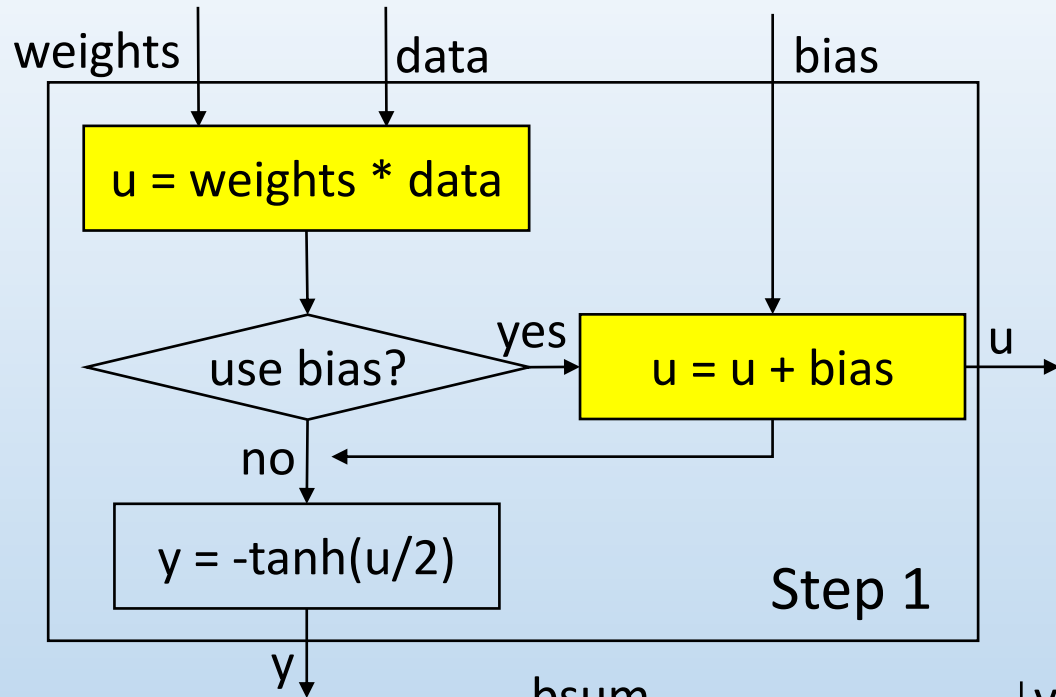


1. Initialize $W(0)$
2. **while** (not converged and $t < \text{limit}$)
 - 2.1 $\mathbf{Y} = \text{permutate}(\mathbf{X})$
 - 2.2 **for** each block \mathbf{B} in signal \mathbf{Y}
 - 2.2.1 $\mathbf{U} = \mathbf{W}\mathbf{Y} + \text{bias}$
 - 2.2.2 $\mathbf{Y} = \tanh(\mathbf{U})$
 - 2.2.3 $\mathbf{W}^* =$
 $\mathbf{W} + \eta(\mathbf{I} - \mathbf{Y}\mathbf{U}^T - \mathbf{U}\mathbf{U}^T)\mathbf{W}$
 - 2.2.4 update bias
 - end for**
- end while**

2.3 Variables analysis



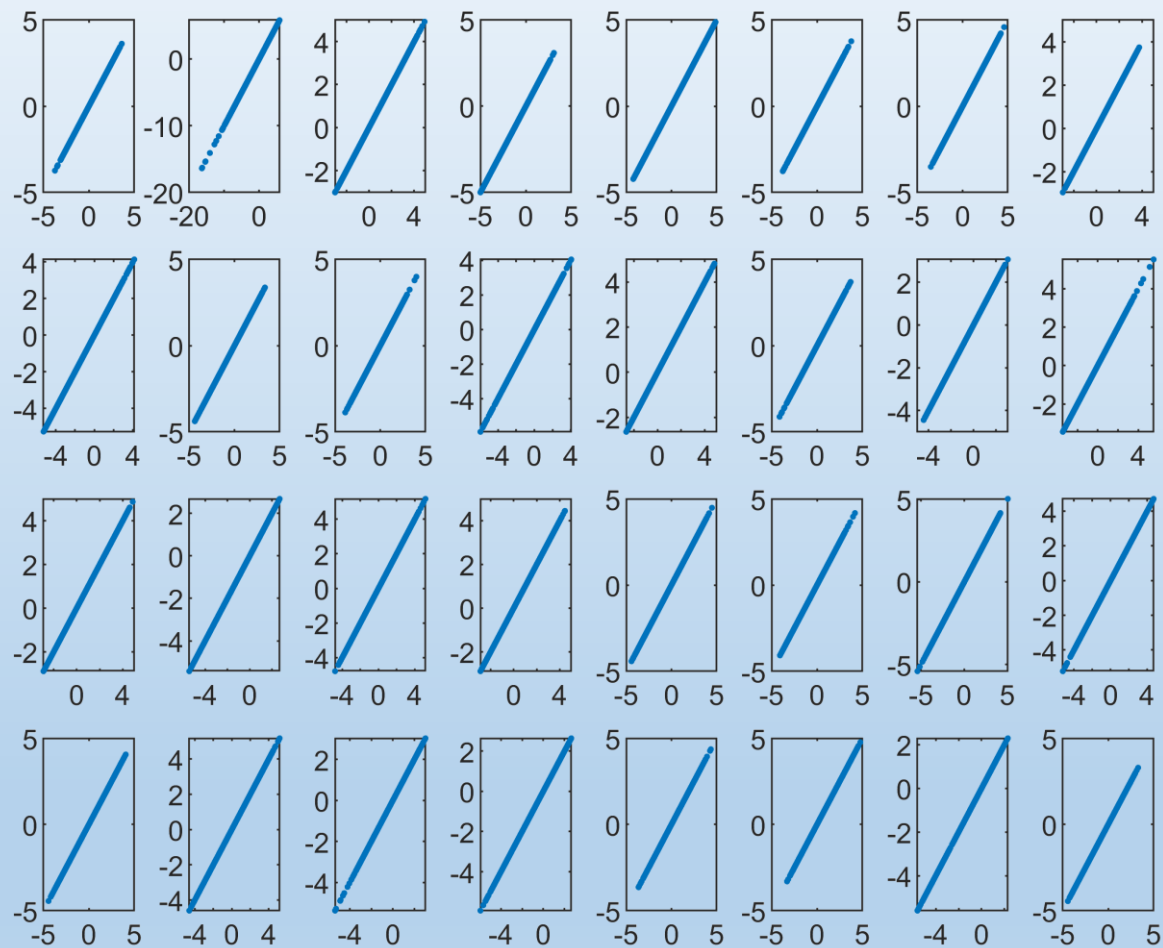
3.1 Locating matrix multiplication



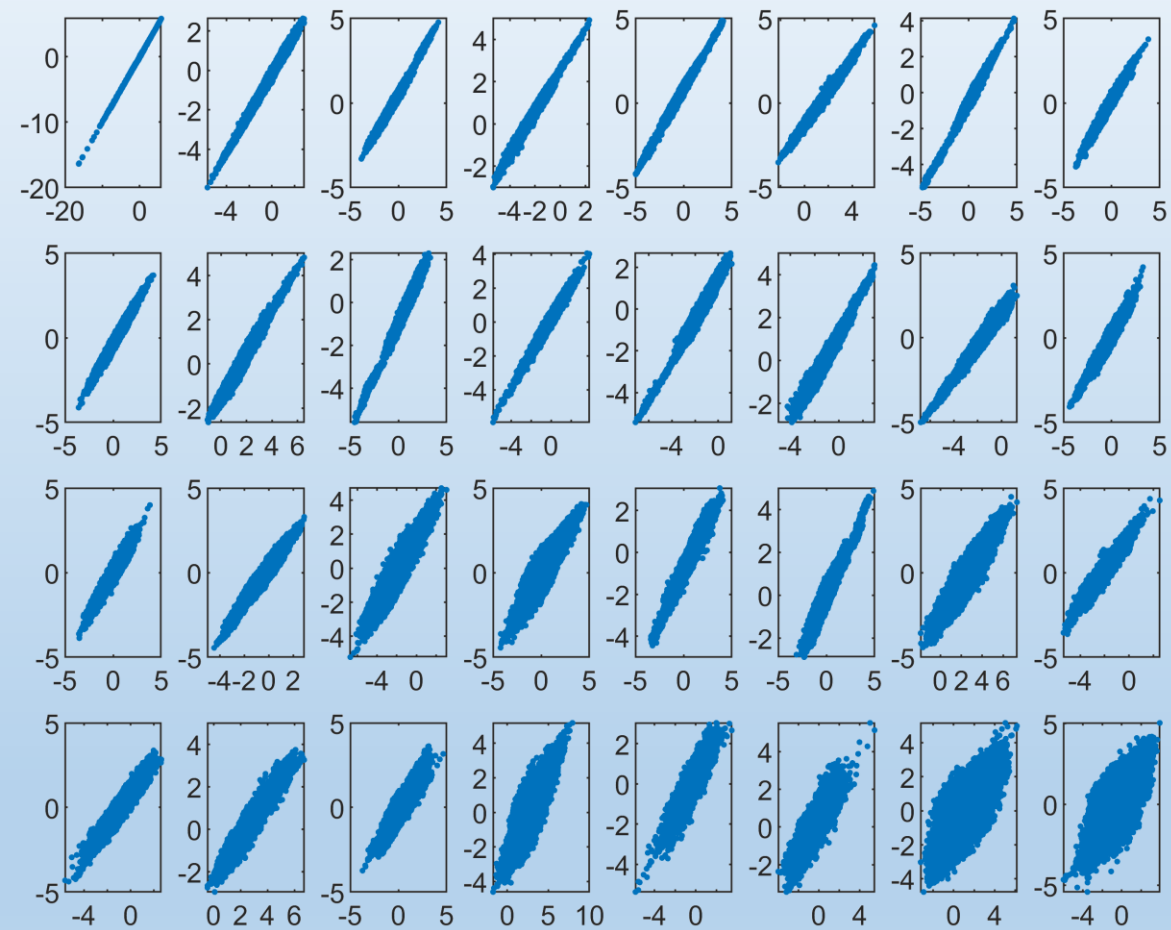
1. Initialize $W(0)$
2. **while** (not converged and $t < \text{limit}$)
 - 2.1 $\mathbf{Y} = \text{permute}(\mathbf{X})$
 - 2.2 **for** each block \mathbf{B} in signal \mathbf{Y}
 - 2.2.1 $\mathbf{U} = \mathbf{W}\mathbf{Y} + \text{bias}$
 - 2.2.2 $\mathbf{Y} = \text{tanh}(\mathbf{U})$
 - 2.2.3 $\mathbf{W}^* =$
 $\mathbf{W} + \eta(\mathbf{I} - \mathbf{Y}\mathbf{U}^T - \mathbf{U}\mathbf{U}^T)\mathbf{W}$
 - 2.2.4 update bias
- end for**
- end while**

: FMA operation

3.2 Numerical correctness



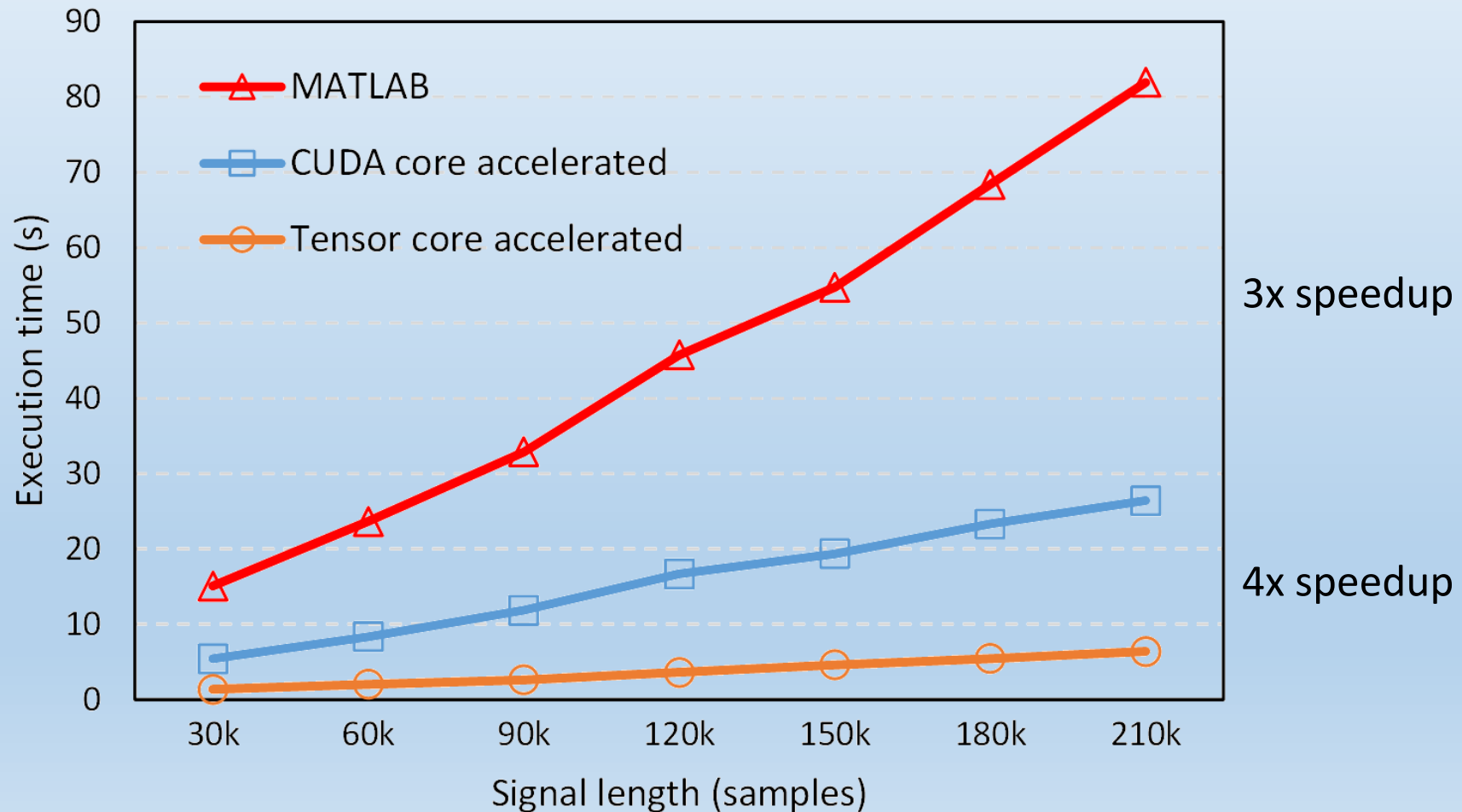
The diagonal lines produced by IC pairs from
MATLAB-ICA and CUDAICA



The diagonal lines produced by IC pairs from
tensor core ICA and CUDAICA

3.3 Performance preview

32-channel EEG data, up to 210k time points, CUDA core and tensor core execution time are produced by A100 GPU on Komondor supercomputer. The MATLAB execution time is produced by a 8-core CPU (i7-9700k)



4 Future works

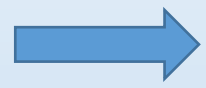
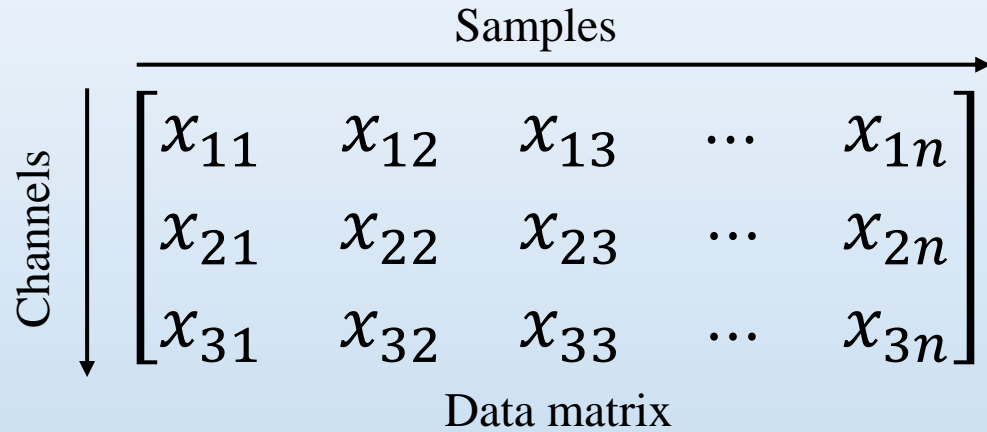
- Numerical correctness issue
- Tensor core + shared memory strategy
- Further performance testing
 - New strategy and larger dataset
 - EEG processing pipeline integration
- Impact of mixed precision computing
 - FP32, TF32, FP16, BF16
- **Other potential algorithms that tensor cores can optimize**
 - **Are there a lot of matrix multiplication operations in the algorithm? Why not consider using tensor cores to speed it up?**



UNKP Új Nemzeti
Kiválóság Program

This research is supported by
Új Nemzeti Kiválóság Program (UNKP)

Decorrelation



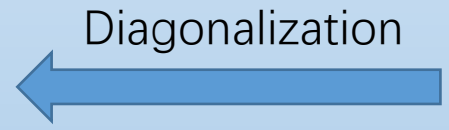
$$Cov(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - E(x))(y_i - E(Y))$$

Covariance



$$\begin{bmatrix} Var(CH1) & Cov(CH1, CH2) & Cov(CH1, CH3) \\ Cov(CH2, CH1) & Var(CH2) & Cov(CH2, CH3) \\ Cov(CH3, CH1) & Cov(CH3, CH2) & Var(CH3) \end{bmatrix}$$

Covariance matrix

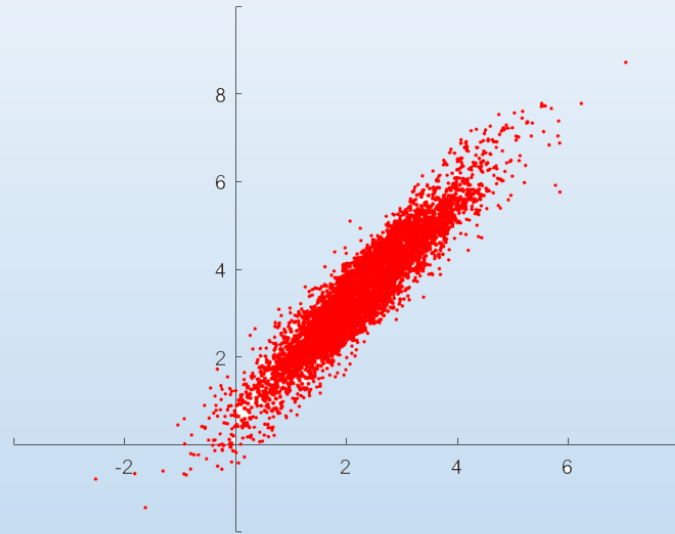


$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

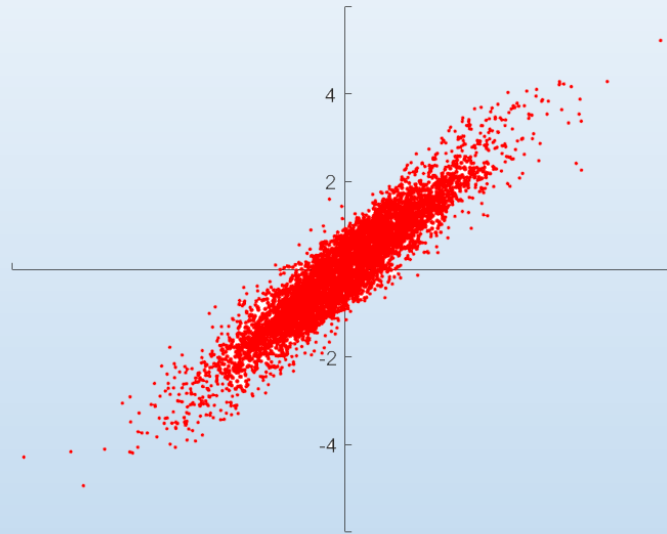
Diagonal matrix

$Cov(X, Y) \begin{cases} > 0: & \text{Positive correlation} \\ < 0: & \text{Negative correlation} \\ = 0: & \text{No correlation} \end{cases}$

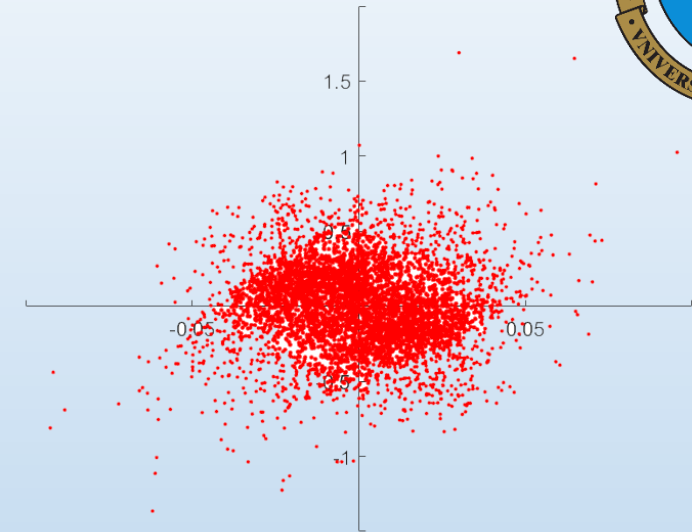
Decorrelation



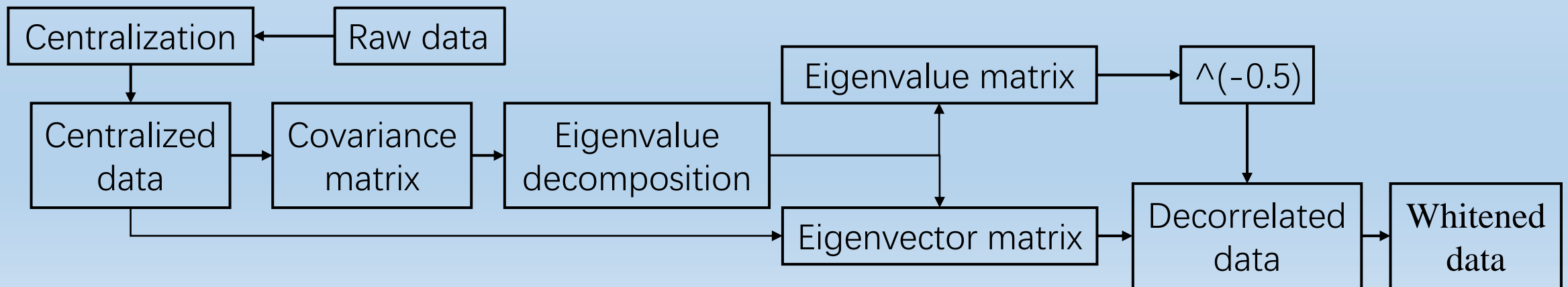
Raw data



Centralized data



Whitened data



$$X = U\Sigma V^T \quad \begin{array}{l} UU^T = I \\ V^T V = I \end{array}$$

$$XX^T = U\Sigma V^T * (U\Sigma V^T)^T$$

$$XX^T = U\Sigma V^T * V\Sigma U^T$$

$$XX^T = U\Sigma U^T$$

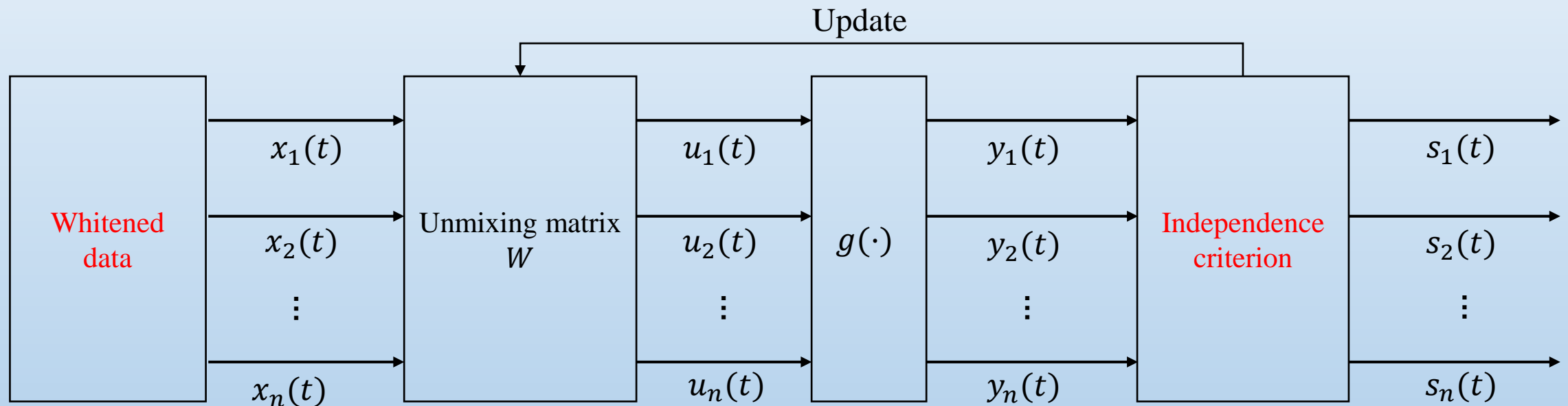
$$XX^T U = U\Sigma^2$$

$$U^T XX^T U = U^T U\Sigma^2$$

Set $Y = U^T X$

$$YY^T = \Sigma^2$$

Remove dependency



The process of solving the unmixing matrix