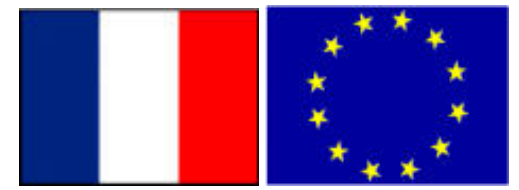




POLYTECH<sup>®</sup>  
PARIS-SACLAY



GP-GPU

# TD1: Dense matrix product on GPU using registers

Stéphane Vialle



université  
PARIS-SACLAY

ÉCOLE DOCTORALE

Sciences et technologies  
de l'information  
et de la communication (STIC)



Stephane.Vialle@centralesupelec.fr  
<http://www.metz.supelec.fr/~vialle>

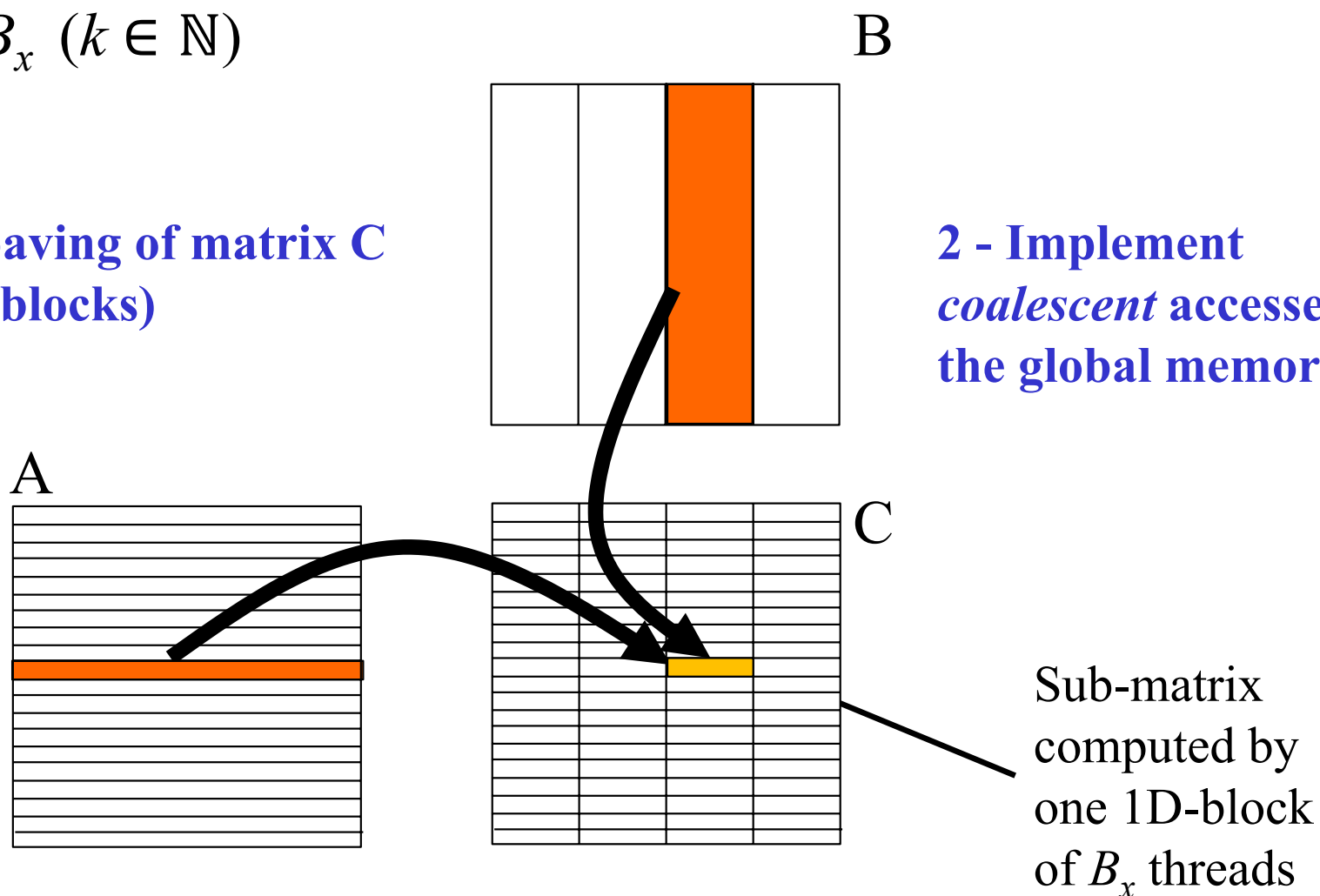
# Ex1: 2D-grid of 1D-blocks

## Product of matrices of $n \times n$ elements

- 1D-blocks of  $B_x$  threads (kernel  $k_0$ )
- With:  $n = k \cdot B_x$  ( $k \in \mathbb{N}$ )

1 - Design the paving of matrix C  
(2D-grid of 1D-blocks)

2 - Implement  
*coalescent* accesses to  
the global memory



# Ex1: 2D-grid of 1D-blocks

Q1.1 - Definition of the 2D-grid of 1D-blocks

Q1.2 - Definition of the computing kernel

Q1.3 - Analysis of the coalescence

Q1.4 - Number of global memory accesses (not including cache memory):

- Compute the total number of memory accesses requested by the threads

$$N_{RAM\ accesses}^{requested\ by\ all\ threads} = (n_{threads} \cdot n_{RAM\ accesses}^{requested\ by\ 1\ thread})$$

- Compute the total number of memory accesses achieved by the warps

When accesses are coalescent: 1 warp accesses 32 data in  $t_1$  RAM access

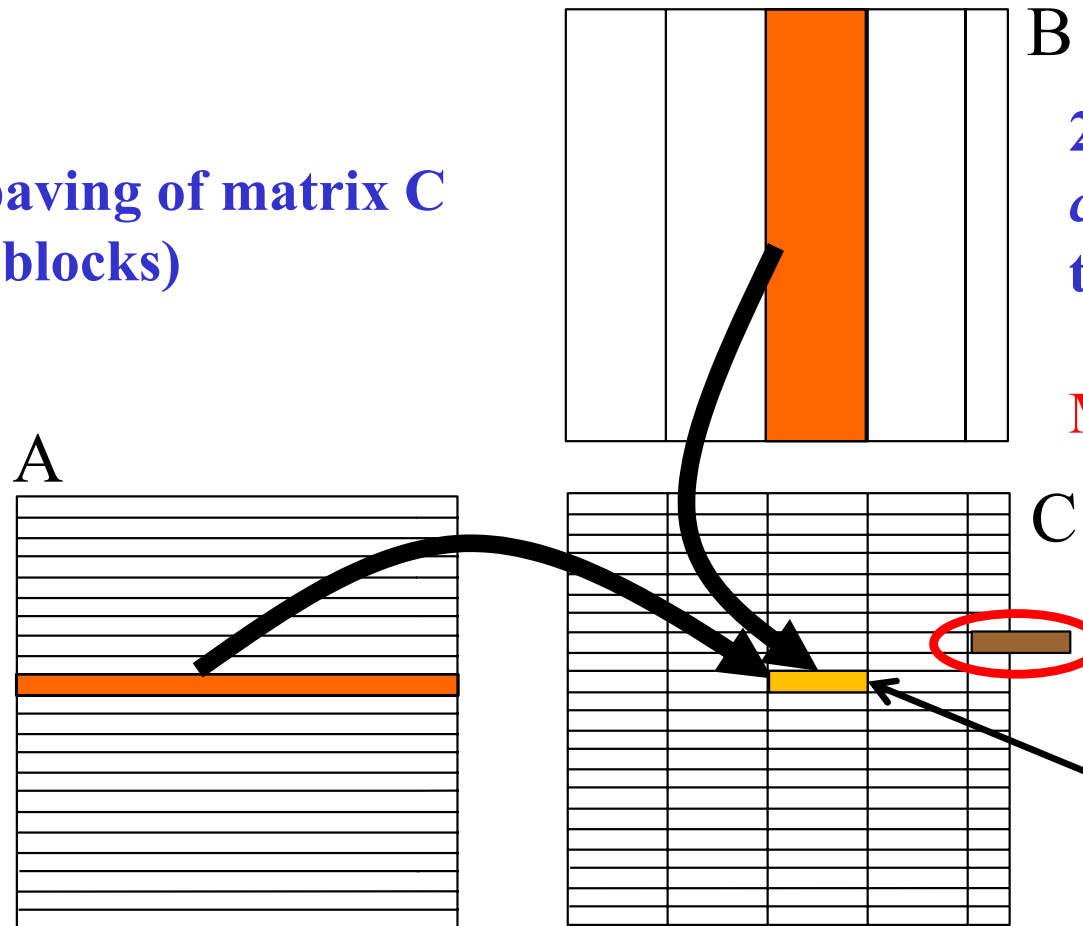
$$\text{Model : } T_{RAM\ access}^{total} = N_{RAM\ accesses}^{achieved\ by\ all\ warps} \cdot t_1\ RAM\ access$$

# Ex1: 2D-grid of 1D-blocks – Boundaries

## Product of matrices of $n \times n$ elements

- 1D-blocks of  $B_x$  threads (**upgraded** kernel  $k_0$ )
- With:  $n \neq k \cdot B_x$ ,  $k \in \mathbb{N}$

1 - Design the paving of matrix C  
(2D-grid of 1D-blocks)



2 - Implement  
*coalescent* accesses to  
the global memory  
&  
*Manage boundaries*

Sub-matrix  
computed by  
one 1D-block  
of  $B_x$  threads

# Ex1: 2D-grid of 1D-blocks – Boundaries

**Q1.5** - Upgrade of the 2D-grid of 1D-blocks

**Q1.6** - Upgrade of the computing kernel

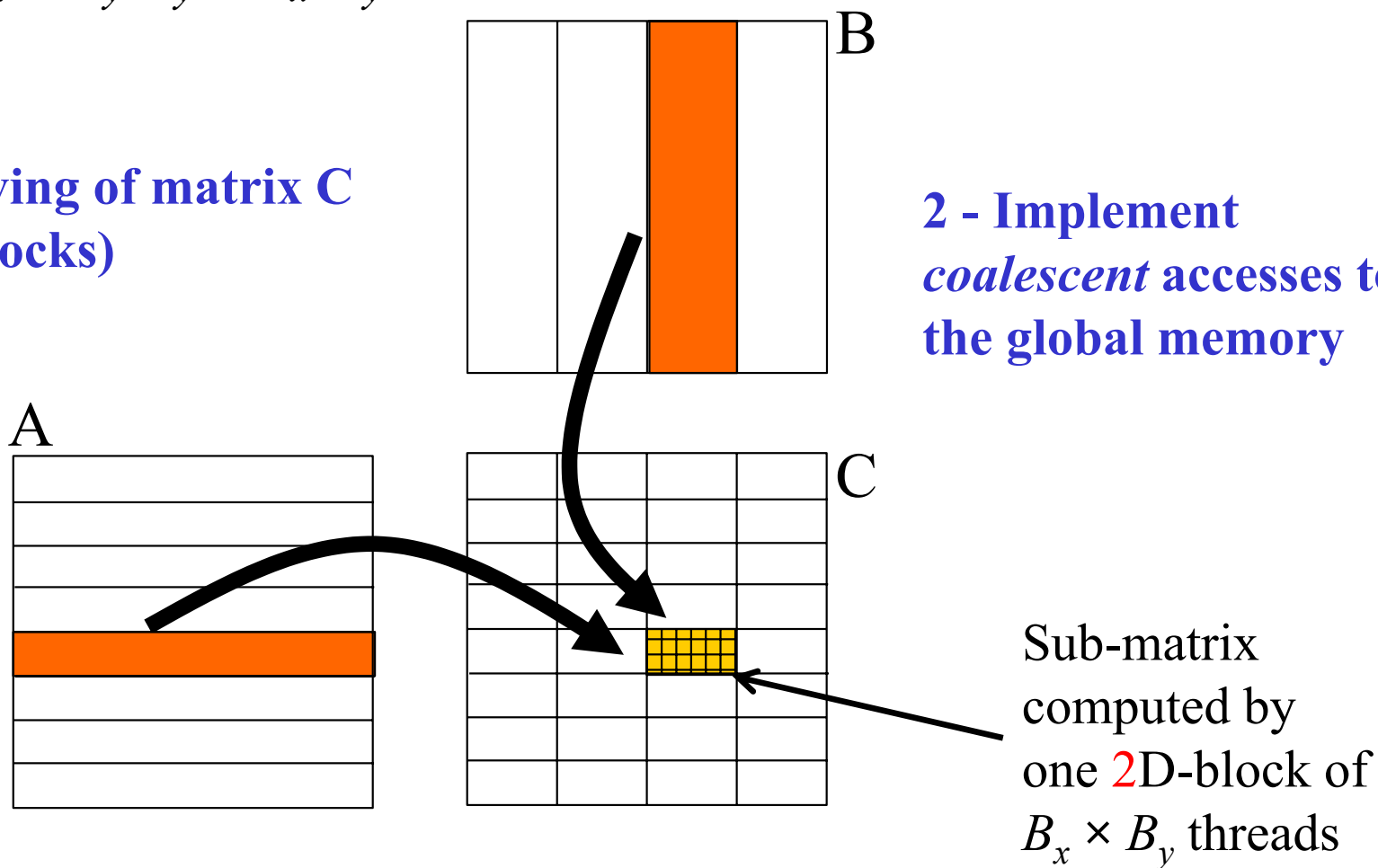
# Ex2: 2D-grid of 2D-blocks

## Product of matrices of $n \times n$ elements

- 2D-blocks of  $B_x \times B_y$  threads (kernel  $k_1$ )
- With:  $n = k_x \cdot B_x = k_y \cdot B_y$ ,  $(k_x, k_y) \in \mathbb{N}^2$

1 - Design the paving of matrix C  
(2D-grid of 2D-blocks)

2 - Implement  
*coalescent* accesses to  
the global memory



# Ex2: 2D-grid of 2D-blocks

Q2.1 - Definition of the 2D-grid of 2D-blocks

Q2.2 - Definition of the computing kernel

Q2.3 - Analysis of the coalescence

Q2.4 - Number of global memory accesses (not including cache memory):

- Compute the total number of memory accesses requested by the threads

$$N_{RAM\ access}^{requested\ by\ all\ threads} = (n_{threads} \cdot n_{RAM\ access}^{requested\ by\ 1\ thread})$$

- Compute the total number of memory accesses achieved by the warps

When accesses are coalescent: 1 warp accesses 32 data in  $t_1$  RAM access

$$\text{Model : } T_{RAM\ access}^{total} = N_{RAM\ access}^{achieved\ by\ all\ warps} \cdot t_1\ RAM\ access$$

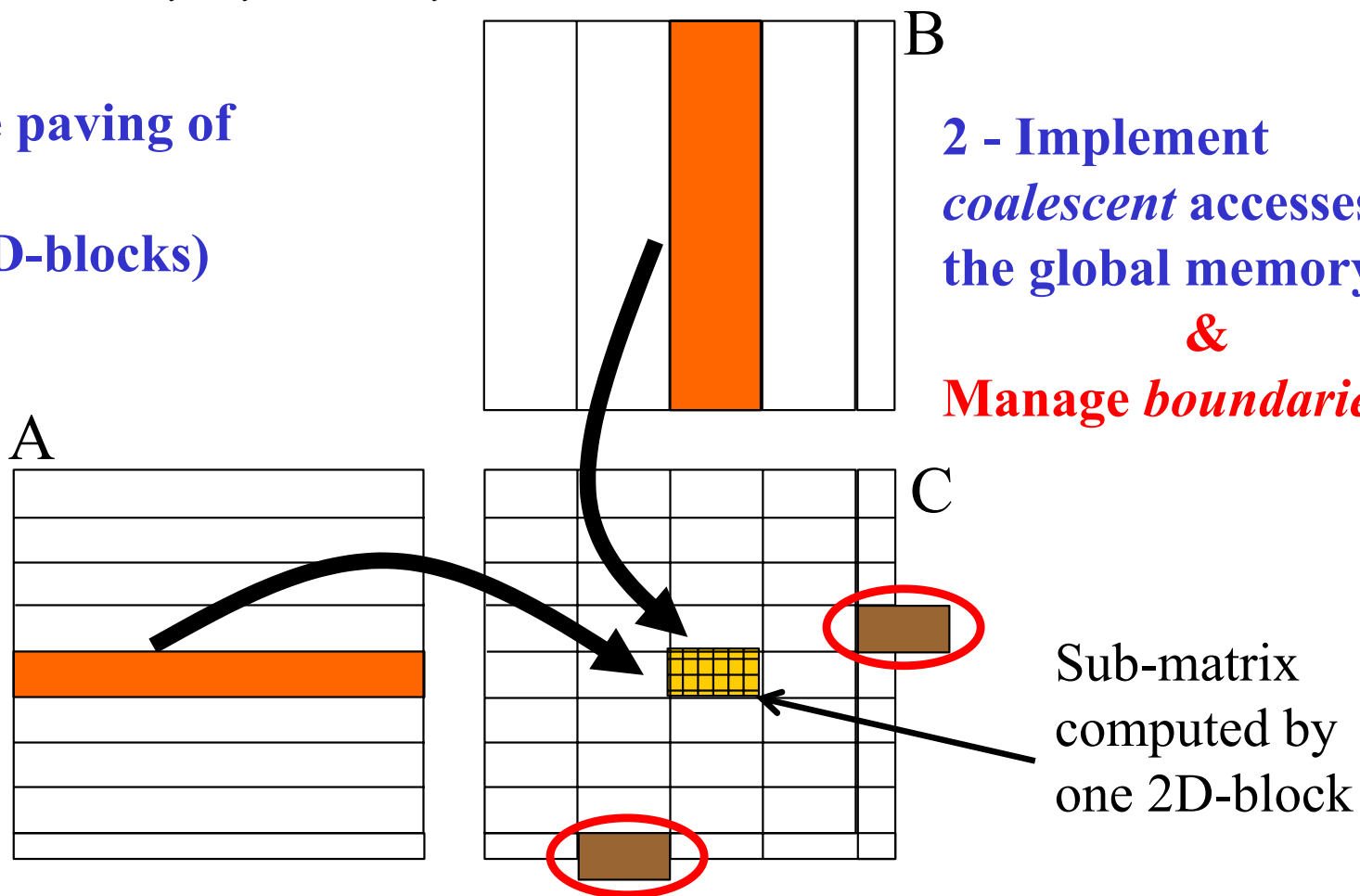
# Ex2: 2D-grid of 2D-blocks – Boundaries

## Product of matrices of $n \times n$ elements

- 2D-blocks of  $B_x \times B_y$  threads (**upgraded** kernel  $k_1$ )
- With:  $n \neq k_x \cdot B_x$ ,  $n \neq k_y \cdot B_y$ ,  $(k_x, k_y) \in \mathbb{N}^2$

1 - Design the paving of matrix C  
(2D-grid of 2D-blocks)

2 - Implement *coalescent* accesses to the global memory  
&  
*Manage boundaries*





# Ex2: 2D-grid of 2D-blocks – Boundaries

**Q2.5** - Upgrade of the 2D-grid of 2D-blocks

**Q2.6** - Upgrade of the computing kernel

# TD1 : Dense matrix product on GPU using registers

Fin