# FAST AND FURIOUS THINGS IN AGRUM/PYAGRUM

## CHRISTOPHE GONZALES

### LIS – AIX-MARSEILLE UNIVERSITÉ

Parallelization in Learning algorithms

Parallelization in Learning algorithms

Parallelization in Inferences

Parallelization in Learning algorithms

Parallelization in Inferences

aGrUM's multithreading facility

## BN structure learning : Greedy Hill Climbing

```
 1  𝒢_best ← initial graph (empty)
 2  sc_best ← Score(𝒢_best)
 3
 4  repeat
 5  │   found ← false
 6  │   foreach 𝒢' ∈ neighborhood of 𝒢_best do
 7  │   │   sc' ← Score(𝒢')
 8  │   │   if sc' > sc_best then
 9  │   │   │   𝒢_best ← 𝒢', sc_best ← sc'
10  │   │   │   found ← true
11  until found = false;
12
13  return 𝒢_best
```

**1** $\mathcal{G}_{best} \leftarrow$ initial graph (empty)
**2** $sc_{best} \leftarrow Score(\mathcal{G}_{best})$
**3**
**4** **repeat**
**5**      found $\leftarrow$ false
**6**      **foreach** $\mathcal{G}' \in$ *neighborhood of* $\mathcal{G}_{best}$ **do**
**7**          $sc' \leftarrow Score(\mathcal{G}')$
**8**          **if** $sc' > sc_{best}$ **then**
**9**              $\mathcal{G}_{best} \leftarrow \mathcal{G}'$, $sc_{best} \leftarrow sc'$
**10**              found $\leftarrow$ true
**11** **until** *found* = *false*;       ▶ 2 parallelization opportunities :
**12**
**13** **return** $\mathcal{G}_{best}$

**1** $\mathcal{G}_{best} \leftarrow$ initial graph (empty)
**2** $sc_{best} \leftarrow Score(\mathcal{G}_{best})$
**3**
**4 repeat**
**5**   found $\leftarrow$ false
**6**   **foreach** $\mathcal{G}' \in$ *neighborhood of* $\mathcal{G}_{best}$ **do**
**7**     $sc' \leftarrow Score(\mathcal{G}')$
**8**     **if** $sc' > sc_{best}$ **then**
**9**       $\mathcal{G}_{best} \leftarrow \mathcal{G}'$, $sc_{best} \leftarrow sc'$
**10**       found $\leftarrow$ true
**11 until** *found = false*;
**12**
**13 return** $\mathcal{G}_{best}$

▶ 2 parallelization opportunities :

   ▶ One thread per graph $\mathcal{G}'$

**1** $\mathcal{G}_{best} \leftarrow$ initial graph (empty)
**2** $\text{sc}_{best} \leftarrow \text{Score}(\mathcal{G}_{best})$
**3**
**4 repeat**
**5** | found $\leftarrow$ false
**6** | **foreach** $\mathcal{G}' \in$ *neighborhood of* $\mathcal{G}_{best}$ **do**
**7** | | $\text{sc}' \leftarrow \text{Score}(\mathcal{G}')$
**8** | | **if** $sc' > sc_{best}$ **then**
**9** | | | $\mathcal{G}_{best} \leftarrow \mathcal{G}'$, $\text{sc}_{best} \leftarrow \text{sc}'$
**10** | | | found $\leftarrow$ true
**11 until** *found = false*;
**12**
**13 return** $\mathcal{G}_{best}$

▶ 2 parallelization opportunities :

   ▶ One thread per graph $\mathcal{G}'$

   ▶ Several threads for each score

**1** $\mathcal{G}_{best} \leftarrow$ initial graph (empty)
**2** $sc_{best} \leftarrow Score(\mathcal{G}_{best})$
**3**
**4** **repeat**
**5**      found $\leftarrow$ false
**6**      **foreach** $\mathcal{G}' \in$ *neighborhood of* $\mathcal{G}_{best}$ **do**
**7**          $sc' \leftarrow Score(\mathcal{G}')$
**8**          **if** $sc' > sc_{best}$ **then**
**9**              $\mathcal{G}_{best} \leftarrow \mathcal{G}'$, $sc_{best} \leftarrow sc'$
**10**              found $\leftarrow$ true
**11** **until** *found = false*;

▶ 2 parallelization opportunities :

**12**
**13** **return** $\mathcal{G}_{best}$

     ▶ One thread per graph $\mathcal{G}'$

     ▶ Several threads for each score

*The BD score*

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

*The BD score*

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

-

**The BD score**

$$\text{Score}_{BD}(X_i|\mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$
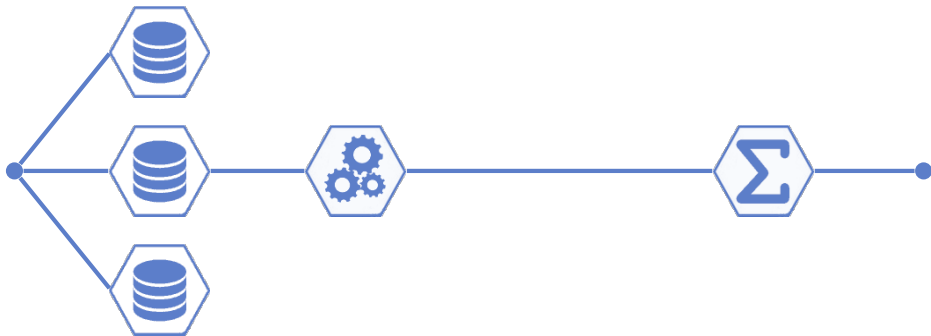
**The BD score**

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

**The BD score**

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$
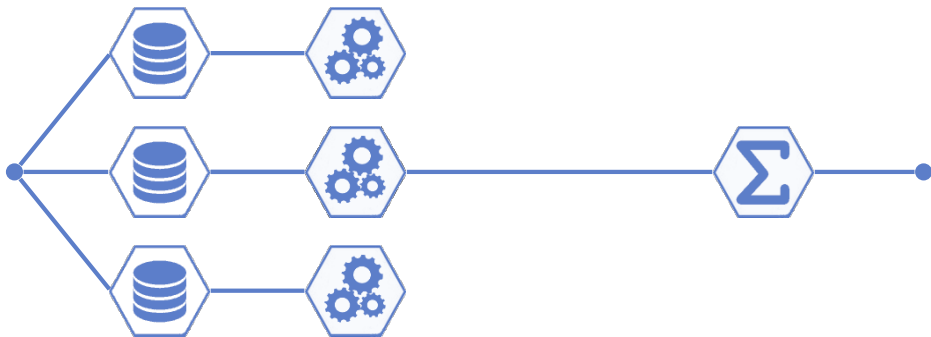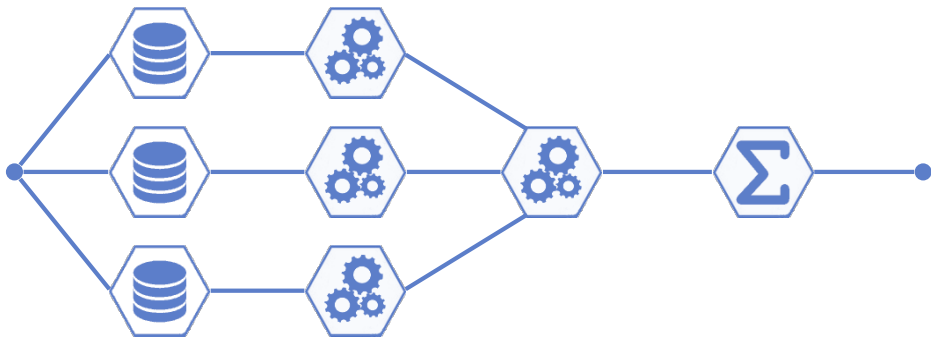
# Parallelizing the scores

## The BD score

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

**The BD score**

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

# Parallelizing the scores

## The BD score

$$\text{Score}_{BD}(X_i|\textbf{Pa}(X_i), \textbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$
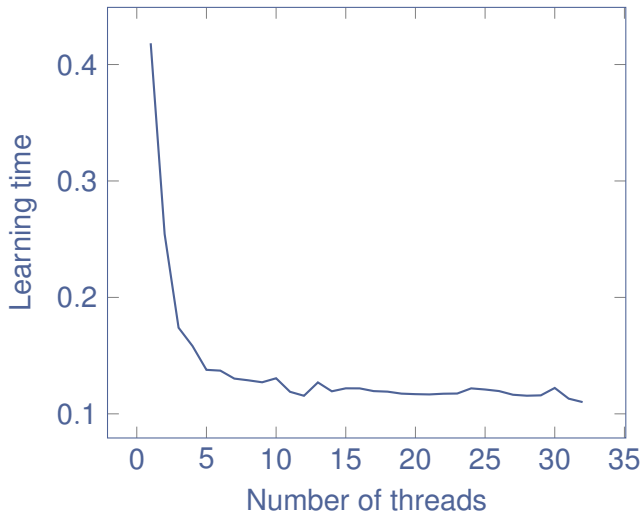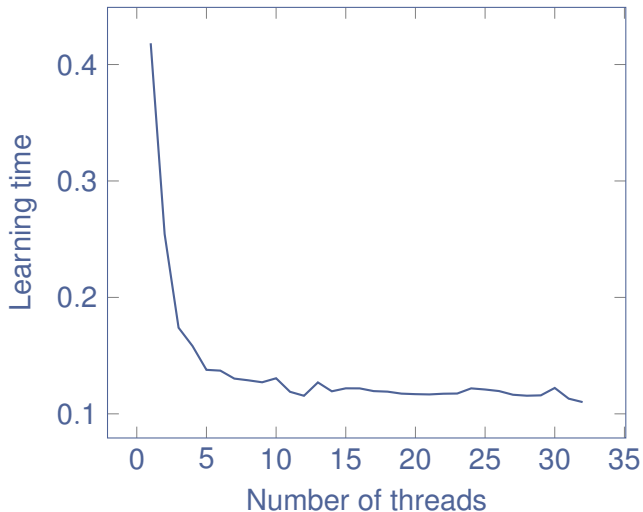
## The BD score

$$\text{Score}_{BD}(X_i | \mathbf{Pa}(X_i), \mathbf{D}) = \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(N_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(N_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})}$$

► **aGrUM's rule :** Number of records per thread ≥ 512

**aGrUM's rules :**

1. pyAgrum manages a « package-wide » number of threads

   $\implies$ by default, same number for all multithreaded objects

**aGrUM's rules :**

1. pyAgrum manages a « package-wide » number of threads
   $\implies$ by default, same number for all multithreaded objects

2. Objects can override this number

```
import pyAgrum as gum
```

▶ **Package-wide functions :**

| function | meaning |
|----------|---------|
| gum.getMaxNumberOfThreads | The number of processors of the computer |

```
import pyAgrum as gum
```

▶ **Package-wide functions :**

| function | meaning |
|---|---|
| gum.getMaxNumberOfThreads | The number of processors of the computer |

```
import pyAgrum as gum
```

▶ **Package-wide functions :**

| function | meaning |
|---|---|
| gum.getMaxNumberOfThreads | The number of processors of the computer |
| gum.getNumberOfLogicalCores | The number of processors |

# pyAgrum's Package-wide API

```
import pyAgrum as gum
```

▶ **Package-wide functions :**

| function | meaning |
|----------|---------|
| gum.getMaxNumberOfThreads | The number of processors of the computer |
| gum.getNumberOfLogicalCores | The number of processors |
| gum.getNumberOfThreads | The number of threads used by default by pyAgrum objects |

```
import pyAgrum as gum
```

► **Package-wide functions :**

| function | meaning |
|---|---|
| gum.getMaxNumberOfThreads | The number of processors of the computer |
| gum.getNumberOfLogicalCores | The number of processors |
| gum.getNumberOfThreads | The number of threads used by default by pyAgrum objects |
| gum.setNumberOfThreads | Sets the number of threads used by default |

# Using the pyAgrum-wide API – an example

```
1  import pyAgrum as gum
2
3  print("Nb procs =", gum.getMaxNumberOfThreads())
```

# Using the pyAgrum-wide API – an example

```
1  import pyAgrum as gum
2
3  print("Nb procs =", gum.getMaxNumberOfThreads())
```

```
Nb procs = 64
```

# Using the pyAgrum-wide API – an example

```python
import pyAgrum as gum

print("Nb procs =", gum.getMaxNumberOfThreads())

# Default number of threads for all pyAgrum objects
print("Nb used =", gum.getNumberOfThreads())
```

```
Nb procs = 64
```

# Using the pyAgrum-wide API – an example

```python
import pyAgrum as gum

print("Nb procs =", gum.getMaxNumberOfThreads())

# Default number of threads for all pyAgrum objects
print("Nb used =", gum.getNumberOfThreads())
```

```
Nb procs = 64
Nb used = 64
```

```
1  import pyAgrum as gum
2
3  print("Nb procs =", gum.getMaxNumberOfThreads())
4
5  # Default number of threads for all pyAgrum objects
6  print("Nb used =", gum.getNumberOfThreads())
7
8  # Changing this default number
9  gum.setNumberOfThreads(10)
```

```
Nb procs = 64
Nb used = 64
```

# Using the pyAgrum-wide API – an example

```python
1  import pyAgrum as gum
2
3  print("Nb procs =", gum.getMaxNumberOfThreads())
4
5  # Default number of threads for all pyAgrum objects
6  print("Nb used =", gum.getNumberOfThreads())
7
8  # Changing this default number
9  gum.setNumberOfThreads(10)
10
11 # Default number of threads for all pyAgrum objects
12 print("Nb used =", gum.getNumberOfThreads())
```

```
Nb procs = 64
Nb used = 64
Nb used = 10
```

▶ **Object overriding methods :**

▶ **Object overriding methods :**

| function | meaning |
|---|---|
| obj.getNumberOfThreads | returns the current number of threads used by obj |

▶ **Object overriding methods :**

| function | meaning |
|---|---|
| `obj.getNumberOfThreads` | returns the current number of threads used by `obj` |
| `obj.setNumberOfThreads` | changes the number of threads used by `obj` |

► **Object overriding methods :**

| function | meaning |
|---|---|
| obj.getNumberOfThreads | returns the current number of threads used by obj |
| obj.setNumberOfThreads | changes the number of threads used by obj |
| obj.isGumNumberOfThreadsOverriden | indicates whether obj uses its own number or that of pyAgrum |

# Using the object API – an example

```python
import pyAgrum as gum
learner = gum.BNLearner("data/alarm.csv")

print("pyAgrum threads =", gum.getNumberOfThreads())
print("Learner threads =", learner.getNumberOfThreads())
print("Learner override =",
      learner.isGumNumberOfThreadsOverriden())
```

# Using the object API – an example

```
1  import pyAgrum as gum
2  learner = gum.BNLearner("data/alarm.csv")
3
4  print("pyAgrum threads =", gum.getNumberOfThreads())
5  print("Learner threads =", learner.getNumberOfThreads())
6  print("Learner override =",
7        learner.isGumNumberOfThreadsOverriden())
```

```
pyAgrum threads = 64
Learner threads = 64
Learner override = False
```

## Using the object API – an example

```python
import pyAgrum as gum
learner = gum.BNLearner("data/alarm.csv")

print("pyAgrum threads =", gum.getNumberOfThreads())
print("Learner threads =", learner.getNumberOfThreads())
print("Learner override =",
      learner.isGumNumberOfThreadsOverriden())

# changing the number of threads only for the learner
learner.setNumberOfThreads(10)
print("pyAgrum threads =", gum.getNumberOfThreads())
print("Learner threads =", learner.getNumberOfThreads())
print("Learner override =",
      learner.isGumNumberOfThreadsOverriden())
```

```
pyAgrum threads = 64
Learner threads = 64
Learner override = False
```

## Using the object API – an example

```python
import pyAgrum as gum
learner = gum.BNLearner("data/alarm.csv")

print("pyAgrum threads =", gum.getNumberOfThreads())
print("Learner threads =", learner.getNumberOfThreads())
print("Learner override =",
      learner.isGumNumberOfThreadsOverriden())

# changing the number of threads only for the learner
learner.setNumberOfThreads(10)
print("pyAgrum threads =", gum.getNumberOfThreads())
print("Learner threads =", learner.getNumberOfThreads())
print("Learner override =",
      learner.isGumNumberOfThreadsOverriden())
```

```
pyAgrum threads = 64
Learner threads = 64
Learner override = False
```

```
pyAgrum threads = 64
Learner threads = 10
Learner override = True
```

## Using the object API – an example

```
1   import pyAgrum as gum
2   learner = gum.BNLearner("data/alarm.csv")
3
4   print("pyAgrum threads =", gum.getNumberOfThreads())
5   print("Learner threads =", learner.getNumberOfThreads())
6   print("Learner override =",
7         learner.isGumNumberOfThreadsOverriden())
8
9   # changing the number of threads only for the learner
10  learner.setNumberOfThreads(10)
11  print("pyAgrum threads =", gum.getNumberOfThreads())
12  print("Learner threads =", learner.getNumberOfThreads())
13  print("Learner override =",
14        learner.isGumNumberOfThreadsOverriden())
15
16  # making the learner use the package-wide number again
17  learner.setNumberOfThreads(0) # 0 = package-wide
18  print("Learner threads =", learner.getNumberOfThreads(),
19        " Learner override =",
20        learner.isGumNumberOfThreadsOverriden())
```

```
pyAgrum threads = 64
Learner threads = 64
Learner override = False
```

```
pyAgrum threads = 64
Learner threads = 10
Learner override = True
```

# Using the object API – an example

```python
1  import pyAgrum as gum
2  learner = gum.BNLearner("data/alarm.csv")
3
4  print("pyAgrum threads =", gum.getNumberOfThreads())
5  print("Learner threads =", learner.getNumberOfThreads())
6  print("Learner override =",
7        learner.isGumNumberOfThreadsOverriden())
8
9  # changing the number of threads only for the learner
10 learner.setNumberOfThreads(10)
11 print("pyAgrum threads =", gum.getNumberOfThreads())
12 print("Learner threads =", learner.getNumberOfThreads())
13 print("Learner override =",
14       learner.isGumNumberOfThreadsOverriden())
15
16 # making the learner use the package-wide number again
17 learner.setNumberOfThreads(0) # 0 = package-wide
18 print("Learner threads =", learner.getNumberOfThreads(),
19       " Learner override =",
20       learner.isGumNumberOfThreadsOverriden())
```

```
pyAgrum threads = 64
Learner threads = 64
Learner override = False
```

```
pyAgrum threads = 64
Learner threads = 10
Learner override = True
```

```
Learner threads = 64   Learner override = False
```

**Different situations :**

▶ Performing just one learning :
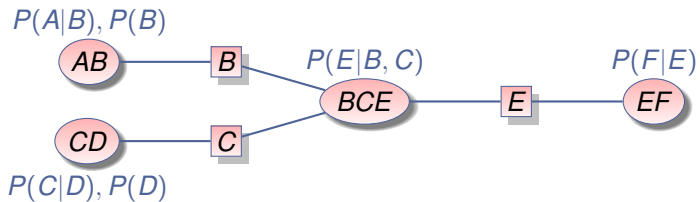  use pyAgrum number of threads

**Different situations :**

▶ Performing just one learning :
use pyAgrum number of threads

▶ Performing a large amount of learning experiments :
use 1 thread per BNlearner
perform experiments in parallel

**Different situations :**

▶ Performing just one learning :
  use pyAgrum number of threads

▶ Performing a large amount of learning experiments :
  use 1 thread per BNlearner
  perform experiments in parallel

▶ Performing a small amount $K$ of learnings :
  let nb processors $\approx A \times B$, with $K$ multiple of $A$
  perform $A$ experiments in parallel
  use $B$ thread per BNlearner

▶ Allow Databases to be column-wise instead of row-wise

$\implies$ improved cacheline use

Parallelization in Inferences

1. $P(A, B) = P(A|B) \times P(B)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$

$P(A|B), P(B)$

$AB$ — $B$

$P(E|B, C)$

$BCE$ — $E$ — $EF$

$P(F|E)$

$CD$ — $C$

$P(C|D), P(D)$

$P(A|B)$   $P(B)$

$\times$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$

$P(A|B), P(B)$

$AB$ — $B$

$P(E|B, C)$

$P(F|E)$

$BCE$ — $E$ — $EF$

$CD$ — $C$

$P(C|D), P(D)$

**1** $P(A, B) = P(A|B) \times P(B)$

**2** $P(B) = \sum_A P(A, B)$

$P(A|B)$     $P(B)$

$\times$

$+$

$P(A|B), P(B)$

$AB$   $B$   $P(E|B,C)$     $P(F|E)$

$BCE$   $E$   $EF$

$CD$   $C$

$P(C|D), P(D)$

$P(A|B)$   $P(B)$

$\times$

$+$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$

$P(A|B), P(B)$

$AB$    $B$    $P(E|B, C)$    $BCE$    $E$    $P(F|E)$    $EF$

$CD$    $C$

$P(C|D), P(D)$

$P(A|B)$    $P(B)$

$\times$

$+$
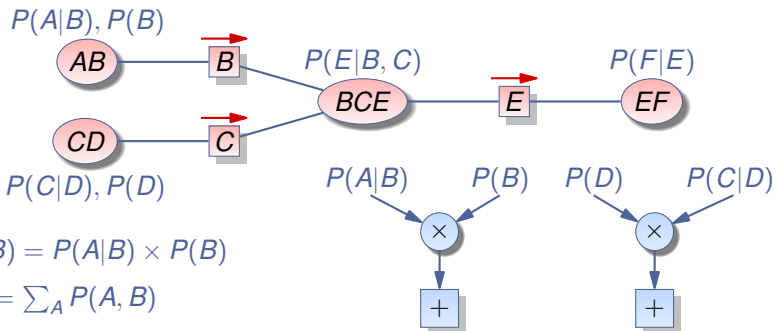
1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$

$P(A|B), P(B)$

$AB$ — $B$

$P(E|B, C)$

$P(F|E)$

$BCE$ — $E$ — $EF$

$CD$ — $C$

$P(C|D), P(D)$

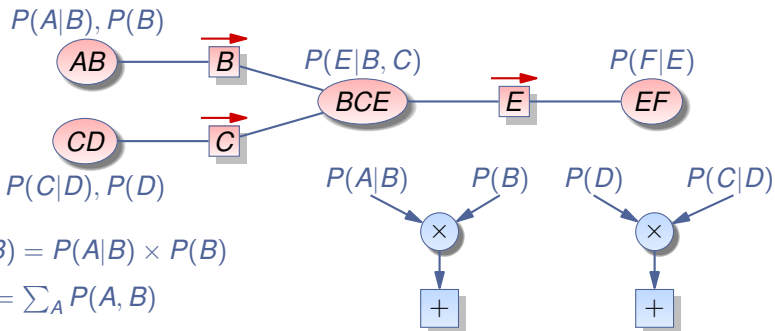$P(A|B)$   $P(B)$       $P(D)$   $P(C|D)$

$\times$              $\times$

$+$              $+$
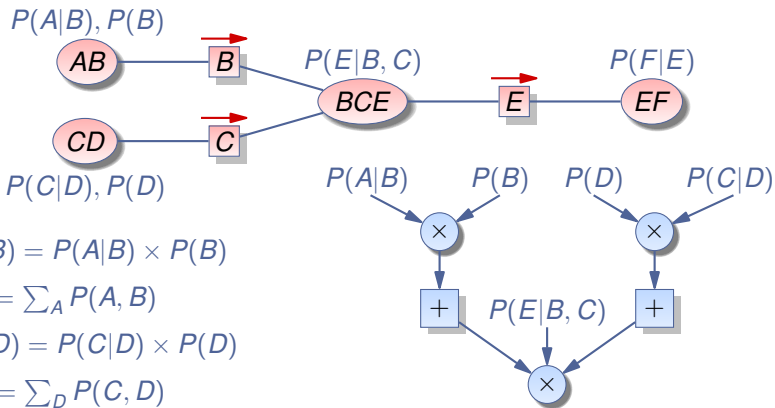
1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$

$P(A|B), P(B)$

$AB$ — $B$

$P(E|B, C)$

$P(F|E)$

$BCE$ — $E$ — $EF$

$CD$ — $C$

$P(C|D), P(D)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$

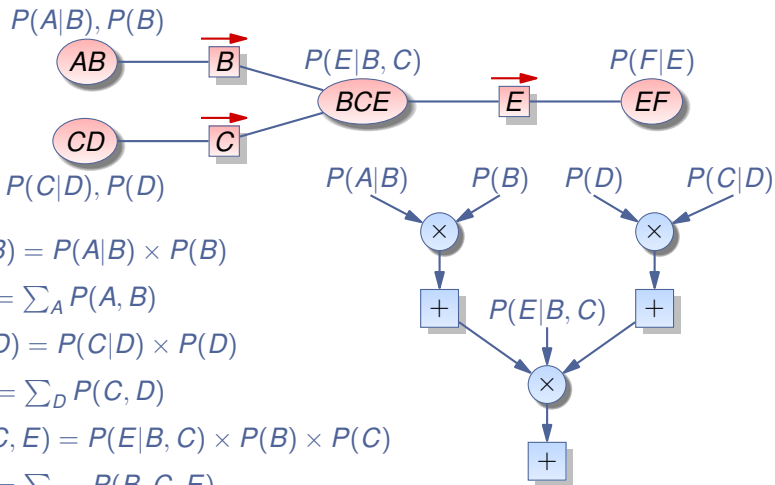$P(A|B)$   $P(B)$   $P(D)$   $P(C|D)$

$\times$   $\times$

$+$   $+$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$
5. $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$
6. $P(E) = \sum_{B,C} P(B, C, E)$

$P(A|B), P(B)$

$P(E|B, C)$

$P(F|E)$

$AB$  $B$

$BCE$  $E$  $EF$

$CD$  $C$

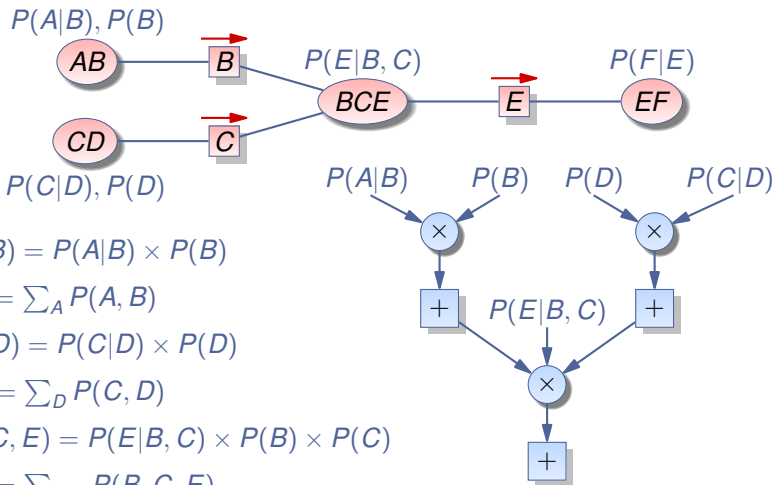$P(C|D), P(D)$

$P(A|B)$  $P(B)$  $P(D)$  $P(C|D)$

$P(E|B, C)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$
5. $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$
6. $P(E) = \sum_{B,C} P(B, C, E)$

$P(A|B), P(B)$

$AB$ — $B$

$P(E|B, C)$

$BCE$ — $E$

$P(F|E)$

$EF$

$CD$ — $C$

$P(C|D), P(D)$

$P(A|B)$   $P(B)$   $P(D)$   $P(C|D)$

$P(E|B, C)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$
5. $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$
6. $P(E) = \sum_{B,C} P(B, C, E)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$
5. $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$
6. $P(E) = \sum_{B,C} P(B, C, E)$
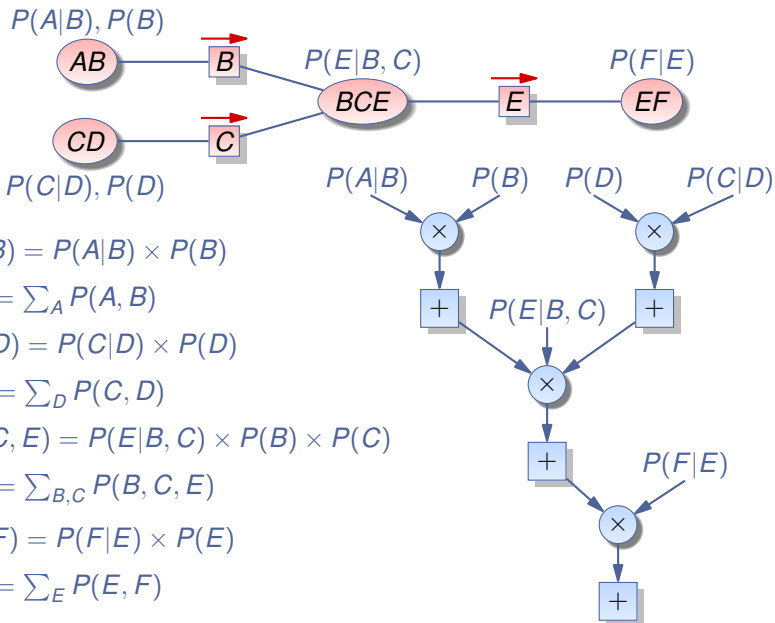7. $P(E, F) = P(F|E) \times P(E)$
8. $P(F) = \sum_E P(E, F)$

$P(A|B), P(B)$

AB — B

$P(E|B, C)$

BCE — E

$P(F|E)$

EF

CD — C

$P(C|D), P(D)$

$P(A|B)$ $P(B)$   $P(D)$ $P(C|D)$

① $P(A, B) = P(A|B) \times P(B)$

② $P(B) = \sum_A P(A, B)$

③ $P(C, D) = P(C|D) \times P(D)$

④ $P(C) = \sum_D P(C, D)$

⑤ $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$

⑥ $P(E) = \sum_{B,C} P(B, C, E)$

⑦ $P(E, F) = P(F|E) \times P(E)$

⑧ $P(F) = \sum_E P(E, F)$

$P(A|B), P(B)$

$P(E|B, C)$

$P(F|E)$

$P(C|D), P(D)$

$P(A|B)$   $P(B)$   $P(D)$   $P(C|D)$

$P(E|B, C)$

$P(F|E)$

1. $P(A, B) = P(A|B) \times P(B)$
2. $P(B) = \sum_A P(A, B)$
3. $P(C, D) = P(C|D) \times P(D)$
4. $P(C) = \sum_D P(C, D)$
5. $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$
6. $P(E) = \sum_{B,C} P(B, C, E)$
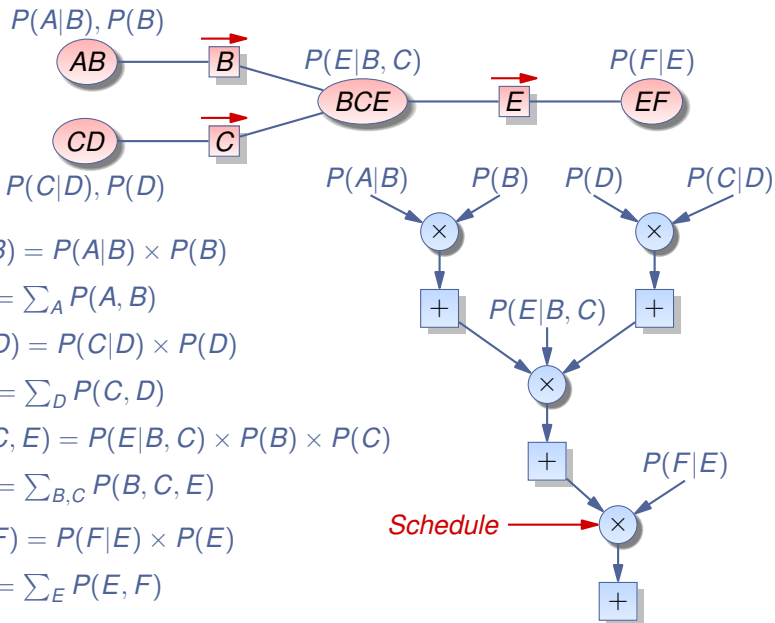7. $P(E, F) = P(F|E) \times P(E)$
8. $P(F) = \sum_E P(E, F)$

$P(A|B), P(B)$

$P(E|B, C)$

$P(F|E)$

$AB$   $B$   $BCE$   $E$   $EF$

$CD$   $C$

$P(C|D), P(D)$

① $P(A, B) = P(A|B) \times P(B)$

② $P(B) = \sum_A P(A, B)$

③ $P(C, D) = P(C|D) \times P(D)$

④ $P(C) = \sum_D P(C, D)$

⑤ $P(B, C, E) = P(E|B, C) \times P(B) \times P(C)$

⑥ $P(E) = \sum_{B,C} P(B, C, E)$

⑦ $P(E, F) = P(F|E) \times P(E)$

⑧ $P(F) = \sum_E P(E, F)$

*Schedule*

1     Create a junction tree

1 Create a junction tree

2 Create a schedule from the JT

1    Create a junction tree

2    Create a schedule from the JT

3    Execute the schedule

1   Create a junction tree

2   Create a schedule from the JT

3   Execute the schedule

▶ **Used by LazyPropagation and Shafer-Shenoy**

1    Create a junction tree

2    Create a schedule from the JT

3    Execute the schedule

▶ **Used by LazyPropagation and Shafer-Shenoy**

▶ **2 schedulers :**

  ▶ Sequential scheduler

  ▶ Parallel scheduler

1     Create a junction tree

2     Create a schedule from the JT

3     Execute the schedule

- ▶ **Used by LazyPropagation and Shafer-Shenoy**
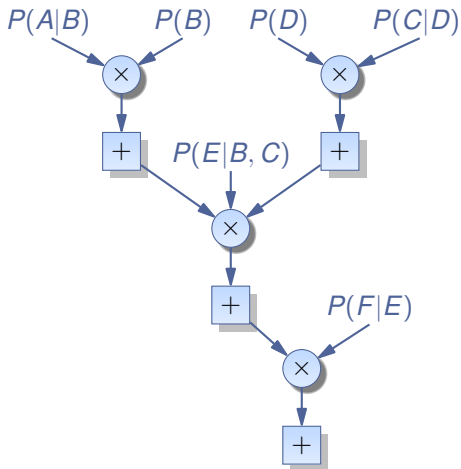
- ▶ **2 schedulers :**
  - ▶ Sequential scheduler
  - ▶ Parallel scheduler

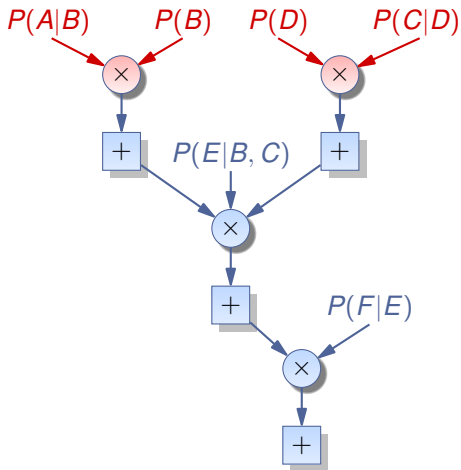- ▶ **1 Rule :** Use the sequential scheduler if and only if :
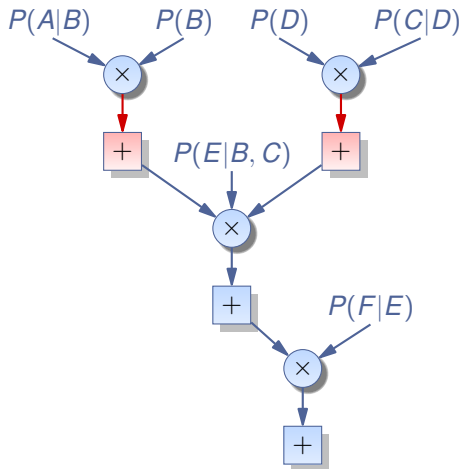
  1 thread or nb elementary operations $< 10^6$

▶ **Synchronization mechanisms :**

- ▶ Mutexes / locks
- ▶ Condition variables
- ▶ Atomics

▶ **Synchronization mechanisms :**

  ▶ Mutexes / locks

  ▶ Condition variables

  ▶ Atomics

► **Synchronization mechanisms :**

  ► Mutexes / locks

  ► Condition variables

  ► Atomics

  . . . but there is still an overhead.

▶ **Synchronization mechanisms :**

   ▶ Mutexes / locks

   ▶ Condition variables

   ▶ Atomics

   . . . but there is still an overhead.

▶ **Experiments on large BNs :**

   ▶ Major time reduction from 1 to 4-6 threads

   ▶ More limited gain above 6 threads

- **Synchronization mechanisms :**
  - Mutexes / locks
  - Condition variables
  - Atomics
  - ... but there is still an overhead.

- **Experiments on large BNs :**
  - Major time reduction from 1 to 4-6 threads
  - More limited gain above 6 threads
  - Explanation : clique sizes imbalanced

**Different situations for large BNs :**

▶ Performing just one inference :
    parallellize with pyAgrum number of threads

**Different situations for large BNs :**

► Performing just one inference :
   parallellize with pyAgrum number of threads

► Performing a large amount of inference experiments :
   use 1 thread per LazyPropagation instance
   perform experiments in parallel

**Different situations for large BNs :**

▶ Performing just one inference :
   parallellize with pyAgrum number of threads

▶ Performing a large amount of inference experiments :
   use 1 thread per LazyPropagation instance
   perform experiments in parallel

▶ Performing a small amount $K$ of inferences :
   perform $K/4$ experiments in parallel
   use 4 threads per LazyPropagation instance

# Schedules and their operators

Objects contained into schedules :

▶ ScheduleMultiDim : abstraction of potentials

# Schedules and their operators

Objects contained into schedules :

- ► ScheduleMultiDim : abstraction of potentials

- ► ScheduleBinaryCombination : $T \otimes T \mapsto T$

- ► ScheduleProjection : $T \Downarrow v \mapsto T$

## Schedules and their operators

Objects contained into schedules :

- ▶ ScheduleMultiDim : abstraction of potentials
- ▶ ScheduleBinaryCombination : $T \otimes T \mapsto T$
- ▶ ScheduleProjection : $T \Downarrow v \mapsto T$
- ▶ ScheduleDeletion : remove a ScheduleMultiDim from memory
- ▶ ScheduleStorage : store a ScheduleMultiDim into a container

# Schedules and their operators

Objects contained into schedules :

- ▶ ScheduleMultiDim : abstraction of potentials

- ▶ ScheduleBinaryCombination : $T \otimes T \mapsto T$

- ▶ ScheduleProjection : $T \Downarrow v \mapsto T$

- ▶ ScheduleDeletion : remove a ScheduleMultiDim from memory

- ▶ ScheduleStorage : store a ScheduleMultiDim into a container

$\implies$ very general-purpose

▶ Expose `numberOfOperations` and `setMaxMemory` to pyAgrum

▶ Expose `numberOfOperations` and `setMaxMemory` to pyAgrum

▶ Reduce the overhead of using schedulers

- Expose `numberOfOperations` and `setMaxMemory` to pyAgrum

- Reduce the overhead of using schedulers

- Add schedulers parallelizing both operators and operations
  $\implies$ requires splitting Potential operators computations

- ▶ Expose `numberOfOperations` and `setMaxMemory` to pyAgrum

- ▶ Reduce the overhead of using schedulers

- ▶ Add schedulers parallelizing both operators and operations
  - $\implies$ requires splitting Potential operators computations

- ▶ Add a scheduler exploiting GPU
  - $\implies$ requires ScheduleOperator for changing the order of variables in potentials

aGrUM's multithreading facility

► Multithreaded objects support both openMP and STL threads

## openMP vs. STL threads

- ► Multithreaded objects support both openMP and STL threads

- ► By default, openMP is used

## openMP vs. STL threads

▶ Multithreaded objects support both openMP and STL threads

▶ By default, openMP is used except if :

    ▶ either the user compiled aGrUM with `--threads=stl` option

    ▶ or the compiler does not support openMP

## openMP vs. STL threads

- ▶ Multithreaded objects support both openMP and STL threads

- ▶ By default, openMP is used except if :

  - ▶ either the user compiled aGrUM with `--threads=stl` option

  - ▶ or the compiler does not support openMP

- ▶ Parallelism achieved by using `ThreadExecutor` instances

## openMP vs. STL threads

- ▶ Multithreaded objects support both openMP and STL threads

- ▶ By default, openMP is used except if :

  - ▶ either the user compiled aGrUM with `--threads=stl` option

  - ▶ or the compiler does not support openMP

- ▶ Parallelism achieved by using `ThreadExecutor` instances

- ▶ Advantages :

  - ▶ Multithreaded objects are agnostic

## openMP vs. STL threads

- ▶ Multithreaded objects support both openMP and STL threads

- ▶ By default, openMP is used except if :

  - ▶ either the user compiled aGrUM with `--threads=stl` option

  - ▶ or the compiler does not support openMP

- ▶ Parallelism achieved by using `ThreadExecutor` instances

- ▶ Advantages :

  - ▶ Multithreaded objects are agnostic

  - ▶ Exceptions can be caught

## openMP vs. STL threads

- ▶ Multithreaded objects support both openMP and STL threads

- ▶ By default, openMP is used except if :

  - ▶ either the user compiled aGrUM with `--threads=stl` option

  - ▶ or the compiler does not support openMP

- ▶ Parallelism achieved by using `ThreadExecutor` instances

- ▶ Advantages :

  - ▶ Multithreaded objects are agnostic

  - ▶ Exceptions can be caught

  - ▶ When one thread : no overhead

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads) -> void {
  std::cout << "thread #" << this_thread << std::endl;
};

try {
  gum::ThreadExecutor::execute(5, func);
} catch(...) {
  std::cout << "Exception catched" << std::endl;
}
```

# ThreadExecutors – an example

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads) -> void {
  std::cout << "thread #" << this_thread << std::endl;
};

try {
  gum::ThreadExecutor::execute(5, func);
} catch(...) {
  std::cout << "Exception catched" << std::endl;
}
```

```
thread #0
thread #4
thread #3
thread #thread #2
1
```

# ThreadExecutors – an example

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads) -> void {
  std::cout << "thread #" << this_thread << std::endl;
};

try {
  gum::ThreadExecutor::execute(5, func);
} catch(...) {
  std::cout << "Exception catched" << std::endl;
}
```

```
thread #0
thread #4
thread #3
thread #thread #2
1
```

$\implies$ Exceptions can be catched in Python !

# ThreadExecutors – another example

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads,
               int nb,
               const std::string& str) -> void {
  std::cout << str << nb << " #"
            << this_thread << std::endl;
};

gum::ThreadExecutor::execute(5, func, 8, "thread ");
```

# ThreadExecutors – another example

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads,
               int nb,
               const std::string& str) -> void {
  std::cout << str << nb << " #"
            << this_thread << std::endl;
};

gum::ThreadExecutor::execute(5, func, 8, "thread ");
```

```
thread 8 #thread 0thread 8 #4thread 8 #2
8
thread 8 #3
 #1
```

# ThreadExecutors – another example

```cpp
auto func = [](const std::size_t this_thread,
               const std::size_t nb_threads,
               int nb,
               const std::string& str) -> void {
  std::cout << str << nb << " #"
            << this_thread << std::endl;
};

gum::ThreadExecutor::execute(5, func, 8, "thread ");
```

```
thread 8 #thread 0thread 8 #4thread 8 #2
8
thread 8 #3
 #1
```

$\implies$ Functions can have as many parameters as wished

Only constraint : first 2 params : this_thread and nb_threads

## Conclusion

▶ Parallelism speeds-up learning and inference computations

▶ Parallelism speeds-up learning and inference computations

▶ Many things to do yet for inferences

  ▶ in particular, check `const` objects. . .

## Conclusion

- Parallelism speeds-up learning and inference computations

- Many things to do yet for inferences

  - in particular, check `const` objects. . .

  - Reduce schedules' creations overhead

## Conclusion

▶ Parallelism speeds-up learning and inference computations

▶ Many things to do yet for inferences

  ▶ in particular, check `const` objects...

  ▶ Reduce schedules' creations overhead

  ▶ inferences over GPU