# Exploiting System Level Heterogeneity to improve the performance of a GeoStatistics Multi-phase Task-based Application

50th International Conference on Parallel Processing

Lucas Leandro Nesi, Arnaud Legrand, Lucas Mello Schnorr August 11, 2021

Institute of Informatics, Federal University of Rio Grande do Sul (UFRGS), Brazil University Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, France













System-Level Heterogeneity: Nodes with Different Computational Power

System-Level Heterogeneity: Nodes with Different Computational Power

#### Santos Dummont



Copyright LNCC

Four Partitions/Systems:

- GPU: Xeon 2695 + K40
- Hybrid: Xeon 2695 + Xeon Phi
- CPU: Xeon 2695
- SDumont: Xeon 6252 + V100

### Jean Zay



Copyright Photothèque CNRS/Cyril Frésillon

Four Partitions/Systems:

- GPU 1: Xeon 6248 + V100
- GPU 2: Xeon 6226 + V100
- GPU 3: Xeon 6240R + A100
- · CPU Only: Xeon 6248

System-Level Heterogeneity: Nodes with Different Computational Power

#### Santos Dummont



Copyright LNCC

Four Partitions/Systems:

- GPU: Xeon 2695 + K40
- Hybrid: Xeon 2695 + Xeon Phi
- CPU: Xeon 2695
- SDumont: Xeon 6252 + V100

### Jean Zay



Copyright Photothèque CNRS/Cyril Frésillon

Four Partitions/Systems:

- GPU 1: Xeon 6248 + V100
- GPU 2: Xeon 6226 + V100
- GPU 3: Xeon 6240R + A100
- CPU Only: Xeon 6248

#### **Cloud in HPC**



Many Different Instances

System-Level Heterogeneity: Nodes with Different Computational Power

#### Santos Dummont



Copyright LNCC

Four Partitions/Systems:

- GPU: Xeon 2695 + K40
- Hvbrid: Xeon 2695 + Xeon Phi .
- CPU: Xeon 2695
- SDumont: Xeon 625

## Jean Zay



Copyright Photothèque CNRS/Cyril Frésillon

Four Partitions/Systems:

GPU 1: Xeon 6248 + V100

#### **Cloud in HPC**



## ExaGeoStat – Geostatistics on Manycore Systems

- Design at Kaust (Hatem Ltaief et al., 2018)
- · Predict missing observations in climate/weather forecasting applications
- · A Task-Based Application that uses Chameleon and StarPU
- Efficiency with Mixed-Precision shown with 4096 nodes



Source: https://github.com/ecrc/exageostat



4096 Shaheen-II Cray XC40 Nodes

Source: https://github.com/ecrc/exageostat

## ExaGeoStat – Fits the Gaussian Process

#### Their Gaussian Process uses:

• Matérn function, hyper-parameter  $\theta$  (scale)



• Optimize  $\theta$  to maximize Gaussian log-likelihood function:

$$I(\theta) = -\frac{n}{2}\log(2\pi) - \frac{1}{2}\log|\Sigma_{\theta}| - \frac{1}{2}\mathsf{Z}^{\mathsf{T}}\Sigma_{\theta}^{-1}\mathsf{Z}$$

## ExaGeoStat – Application DAG Structure

### **Multiple Phases**



## ExaGeoStat – Application DAG Structure



## **ExaGeoStat - Synchronous Version**

#### ExaGeoStat - Sync Version

#### dlacpy dpotrf dsvrk dzcov dcma daemm dmdet dsconv dtrsm sdmat reg A teration 40 В 60 80 C 100 CPU 0 CUDA 0 Node Occupation CPU 1 D CUDA 1 101423 C A B CPU 2 CUDA 2 CPU 3 D CUDA 3 Used (MB) 30000 D 20000 10000 0 50000 Time [ms] 100000 0

#### **Three Major phases**

- A: Generation Matérn
- B: Cholesky
- C: Solve, Determinant, Dot Product

## **Asynchronization Optimizations Results**



## **Asynchronization Optimizations Results**



6/10

- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- *F<sub>s</sub>* Factorization step *s* end time
- Q<sub>s,t</sub> Number of tasks t of step s
- $w_{t,r}$  Duration of task t on resource r
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)



- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- F<sub>s</sub> Factorization step s end time
- Q<sub>s,t</sub> Number of tasks t of step s
- $w_{t,r}$  Duration of task *t* on resource *r*
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

Minimize 
$$\sum_{s \in S} (G_s + F_s) s.t.$$
:



- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- F<sub>s</sub> Factorization step s end time
- Q<sub>s,t</sub> Number of tasks t of step s
- $w_{t,r}$  Duration of task t on resource r
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

$$\begin{split} & \text{Minimize} \sum_{s \in \mathcal{S}} (G_s + F_s) \ s.t. : \\ & \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = \mathcal{Q}_{s,t} \end{split}$$



- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- *F<sub>s</sub>* Factorization step *s* end time
- Q<sub>s,t</sub> Number of tasks t of step s
- $w_{t,r}$  Duration of task t on resource r
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

$$\begin{split} \text{Minimize} & \sum_{s \in \mathcal{S}} (G_s + F_s) \ s.t. : \\ & \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = \textit{Q}_{s,t} \end{split}$$

$$orall s > 1, orall r \in \mathcal{R}: \mathit{G}_{s-1} + lpha_{s, \mathsf{dcmg}, r} \mathit{w}_{\mathsf{dcmg}, r} \leq \mathit{G}_{s}$$



∀s

- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- *F<sub>s</sub>* Factorization step *s* end time
- Q<sub>s,t</sub> Number of tasks t of step s
- $w_{t,r}$  Duration of task t on resource r
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

$$\begin{split} \text{Minimize} & \sum_{s \in S} (G_s + F_s) \ s.t. : \\ & \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = Q_{s,t} \\ & > 1, \forall r \in \mathcal{R} : G_{s-1} + \alpha_{s,\text{dcmg},r} \text{W}_{\text{dcmg},r} \leq G_s \\ & \forall s \in \mathcal{S}, \forall r \in \mathcal{R} : G_s + \sum_{t \neq \text{dcmg}} \alpha_{s,t,r} \text{W}_{t,r} \leq F_s \end{split}$$



- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- *F<sub>s</sub>* Factorization step *s* end time
- $Q_{s,t}$  Number of tasks t of step s
- $w_{t,r}$  Duration of task *t* on resource *r*
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

$$\begin{split} \text{Minimize} & \sum_{s \in \mathcal{S}} (G_s + F_s) \ s.t. : \\ & \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = Q_{s,t} \\ & \forall s > 1, \forall r \in \mathcal{R} : G_{s-1} + \alpha_{s,\text{demg},r} \text{W}_{\text{demg},r} \leq G_s \\ & \forall s \in \mathcal{S}, \forall r \in \mathcal{R} : G_s + \sum_{t \neq \text{demg}} \alpha_{s,t,r} \text{W}_{t,r} \leq F_s \\ & \forall s > 1, \forall r \in \mathcal{R} : F_{s-1} + \sum_{t \neq \text{demg}} \alpha_{s,t,r} \text{W}_{t,r} \leq F_s \end{split}$$



- t a task from the T set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- *F<sub>s</sub>* Factorization step *s* end time
- $Q_{s,t}$  Number of tasks t of step s
- $w_{t,r}$  Duration of task *t* on resource *r*
- *α*<sub>s,t,r</sub> The tasks *t* of step *s* to be placement on *r* (Optimized Variable)

$$\begin{split} \text{Minimize} & \sum_{s \in \mathcal{S}} (G_s + F_s) \ s.t. : \\ & \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = Q_{s,t} \\ & \forall s > 1, \forall r \in \mathcal{R} : G_{s-1} + \alpha_{s,\text{dcmg},r} \textit{w}_{\text{dcmg},r} \leq G_s \\ & \forall s \in \mathcal{S}, \forall r \in \mathcal{R} : G_s + \sum_{t \neq \text{dcmg}} \alpha_{s,t,r} \textit{w}_{t,r} \leq F_s \\ & \forall s > 1, \forall r \in \mathcal{R} : F_{s-1} + \sum_{t \neq \text{dcmg}} \alpha_{s,t,r} \textit{w}_{t,r} \leq F_s \\ & \forall r \in \mathcal{R}, \forall s \in \mathcal{S} : \sum_{z \leq s,t \in \mathcal{T}} \alpha_{z,t,r} \textit{w}_{t,r} \leq F_s \end{split}$$



∀s

- t a task from the  $\mathcal{T}$  set
- r a resource from the  $\mathcal{R}$  set
- s a step from the S set
- G<sub>s</sub> Generation step s end time
- F<sub>s</sub> Factorization step s end time
- $Q_{s,t}$  Number of tasks t of step s
- $w_{t,r}$  Duration of task t on resource r
- $\alpha_{s,t,t}$  The tasks t of step s to be placement on *r* (Optimized Variable)

$$\begin{split} \text{Minimize} \sum_{s \in S} (G_s + F_s) \ s.t. : \\ \forall t \in \mathcal{T}, \forall s \in \mathcal{S} : \sum_{r \in \mathcal{R}} \alpha_{s,t,r} = Q_{s,t} \\ \forall s > 1, \forall r \in \mathcal{R} : G_{s-1} + \alpha_{s,\text{dcmg},r} \text{W}_{\text{dcmg},r} \leq G_s \\ \forall s \in \mathcal{S}, \forall r \in \mathcal{R} : G_s + \sum_{t \neq \text{dcmg}} \alpha_{s,t,r} \text{W}_{t,r} \leq F_s \\ \forall s > 1, \forall r \in \mathcal{R} : F_{s-1} + \sum_{t \neq \text{dcmg}} \alpha_{s,t,r} \text{W}_{t,r} \leq F_s \\ \forall r \in \mathcal{R}, \forall s \in \mathcal{S} : \sum_{z \leq s, t \in \mathcal{T}} \alpha_{z,t,r} \text{W}_{t,r} \leq F_s \\ \min_{r \in \mathcal{R}} (\text{W}_{\text{dcmg},r}) \leq G_1 \end{split}$$



#### Nodes Relative Power



## Multi Distributions: Communication Issue



#### Total Communication Phases: 890

#### Nodes Relative Power



## Multi Distributions: Communication Issue



Improved Generation distribution

1D-1D Factorization distribution



#### Total Communication Phases: 890

Total Communication Phases: 517 (-41.9%)

4 Chetemi + 4 Chifflet



4 Chetemi + 4 Chifflet







4 + 4 + 1 Chifflot









## Conclusion

How to distribute an HPC application on these heterogeneous resources?

## Conclusion

## How to distribute an HPC application on these heterogeneous resources?

- · Asynchronous Phase Execution with correct Overlapping
- · Multi-Phase Distributions is not a individual Task







## Conclusion

## How to distribute an HPC application on these heterogeneous resources?

- · Asynchronous Phase Execution with correct Overlapping
- · Multi-Phase Distributions is not a individual Task

#### Correctaly using Node Heterogeneity improves performance







This study was financed in part by the "Coordenação de Aperfeiçoamento de Pessoal de Nível Superior" - Brasil (CAPES) - Finance Code 001, the National Council for Scientific and Technological Development (CNPq), grant no 141971/2020-7 to the first author, and the projects:

- FAPERGS (Data Science 19/711-6, MultiGPU 16/354-8, and GreenCloud 16/488-9)
- CNPq (447311/2014-0)
- CAPES (Brafitec 182/15, and Cofecub 899/18)

Experiments were carried out using Grid'5000, supported by a scientific interest group hosted by Inria and including CNRS, RENATER and several Universities as well as other organizations (see https://www.grid5000.fr).

## **Questions?**

Contact: lucas.nesi@inf.ufrgs.br

## Lucas Leandro Nesi, Arnaud Legrand, Lucas Mello Schnorr August 11, 2021

Institute of Informatics, Federal University of Rio Grande do Sul (UFRGS), Brazil University Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, France











