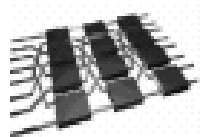# Research Activities of the GPPD
## Parallel and Distributed Processing Group

# Load Balancing Strategies
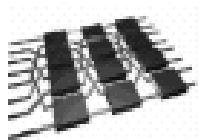
# Informatics Institute
# UFRGS

Philippe O. A. Navaux

Third Workshop of the CNPq-Inria HOSCAR Project

Bordeaux,  September 03, 2013

# We are located in Brazil

… in the State of

Rio Grande do Sul

# Porto Alegre, RS, Brasil

# **Porto Alegre**

Population: 1.44 Million

Climate: S

four well

# UFRGS
# Federal University of Rio Grande do Sul

- Created in 1.895
- One of the five principal universities in Brazil
- Approximately 2.278 faculty members
- Students: approximately 30.000 (undergraduate and graduate)
- Four campi



Porto Alegre
RS - Brasil

# The Institute in Numbers

74 Full-time Professors

2 Departments
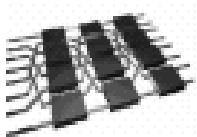- Theoretical and Applied Informatics

2 Undergraduate Programs
- Computer Science and Computer Engineering

3 Graduate Programs

Masters and PhD in Computer Science, Microelectronics, Informatics in Education

691 Undergraduate Students

250 Graduate Students

4 Buildings

Total Area: 7,600 m$^2$

37 Laboratories

Including Teaching and Research Labs

1 Center of Events

4 Auditoria

1 Library

Area: 633 m$^2$

1 Company Incubator (CEI)

# Faculty

- 74 full-time professors dedicated to research and education

- The largest group of researchers in Computer Science and Engineering in Brazil

Faculty members graduated from important institutions around the world

| | |
|---|---|
| Brazil (26) | Belgium (2) |
| France (16) | Portugal (2) |
| German (11) | Canada |
| United Kingdom (5) | Switzerland |
| USA (3) | |

# Research

- Covers all traditional areas of Computer Science and Engineering, from Microelectronics to Computing Theory

  Strong integration among hardware and software groups

| | | |
|---|---|---|
| Artificial Intelligence | Distance Learning | Microelectronics |
| Bioinformatics | Embedded systems | Multimedia |
| Computer Graphics | Fault tolerance | Parallel Processing |
| Computer Networks | Formal Methods | Robotics |
| Databases | Information Systems | Software Engineering |

- Continuous interaction with IT companies in the State
- Origin of many software and most important hardware companies of the region
- Origin of CEITEC the Microelectronic center of Rio Grande do Sul

| | |
|---|---|
| Digitel | CP Eletrônica |
| Altus | Perto |
| Digistar | Parks |
| Teracom | and many others |

# Presenting the GPPD - Parallel and Distributed Processing Group

# Main Research Areas

- Processors architectures
  - Thread Mapping, cache

- Power Consumption

- Cluster, Grid and Cloud computing
  - Virtualization

- Mobile and Ubiquitous Computing

- Parallel and distributed algorithms

- Performance evaluation

- Tools and environments for parallel and distributed programming

- Debugging and Monitoring

# People GPPD

**UFRGS**
UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL

## 2006

**Professors**
Alexandre Carissimi
Cláudio Fernando Resin Geyer
Nicolas Maillard
Philippe Olivier Alexandre Navaux
Tiarajú Asmuz Diverio

**Master Students**
Alexandre da Silveira Ilha
Caciano dos Santos Machado
Carlos Eduardo Benevides Bezerra
Clarissa Cassales Marquezan
Eder Stone Fontoura
Elton Nicoletti Mathias
Everton Hermann
Guilherme Peretti Pezzi
Gustavo Romano
Gustavo Cestari Frainer
João Vicente Lima
Luiz Sequeira Laurino
Marcelo Veiga Neves
Marco Antonio Zanata Alves
Rafael Keller Tesser
Rodrigo Virote Kassick
Rômulo Bandeira Rosinha
Sonia Andrea Lugo Vazquez

**Associated Researchers**
Patrícia Kayser Vargas Mangan
Rafael Bohrer Ávila
Roberto Pinto Souto
Tatiana Gadelha Serra dos Santos

**Undergraduate Students**
Danilo Fukuda Conrad
Eduardo Dias Camaratta
Francieli Zanon Boito
Laércio Lima Pilla
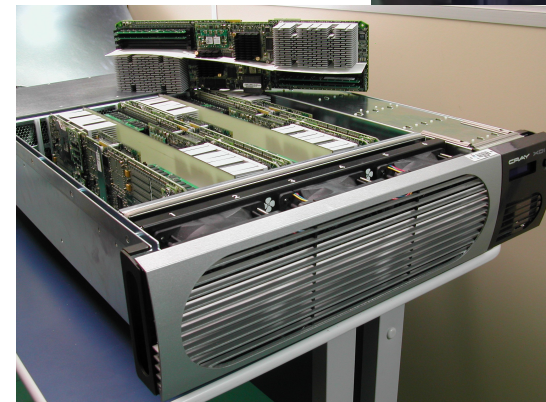Manuela Klanovicz Ferreira
Vicente Silva Cruz

**PhD Students**
Cristiano André da Costa
Débora Nice Ferrari Barbosa
Eduardo Rocha Rodrigues
Emerson André Fedechen
Fábio Reis Cecin
Henrique Cota de Freitas
Lucas Mello Schnorr
Luciano Cavalheiro da Silva
Marcos Ennes Barreto
Marko Petek
Mozart Lemos de Siqueira
Márcia Cristina Cera
Mônica Xavier Py
Rodrigo da Rosa Righi

# People 2011

# The group (12/2012)

# Clusters

- Clusters:
  - DELL Cluster 112 cores,
    - 14 nodes with 2 quad Xeon processors
  - CRAY XD1
    - 6 nodes, dual Opterons with FPGAs.
  - DELL Cluster - LABTEC
    - 20 nodes, PIII, fast-ethernet
  - Itautec Cluster
    - 16 nós, Myrinet
  - Numa Machines
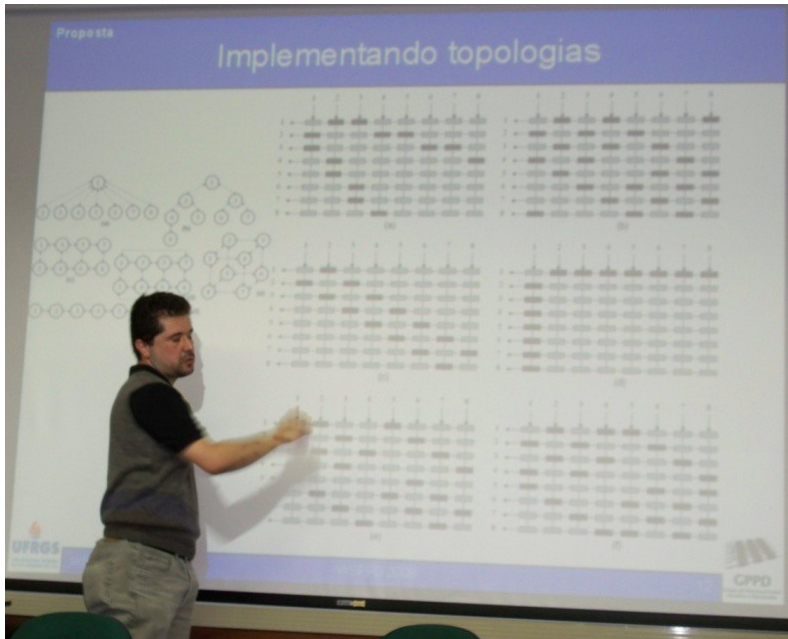  - Access to Grid 5000
  - Access to Cray Machines

# International Cooperation

# Some research in GPPD

# Architecture



Henrique Cota

Network on chip,
Memory Hierarchy,
Energy efficient parallel hardware…

Marco Zanatta

- # NoC

  ## Network on Chip

  – Henrique Costa Freitas



GPPD
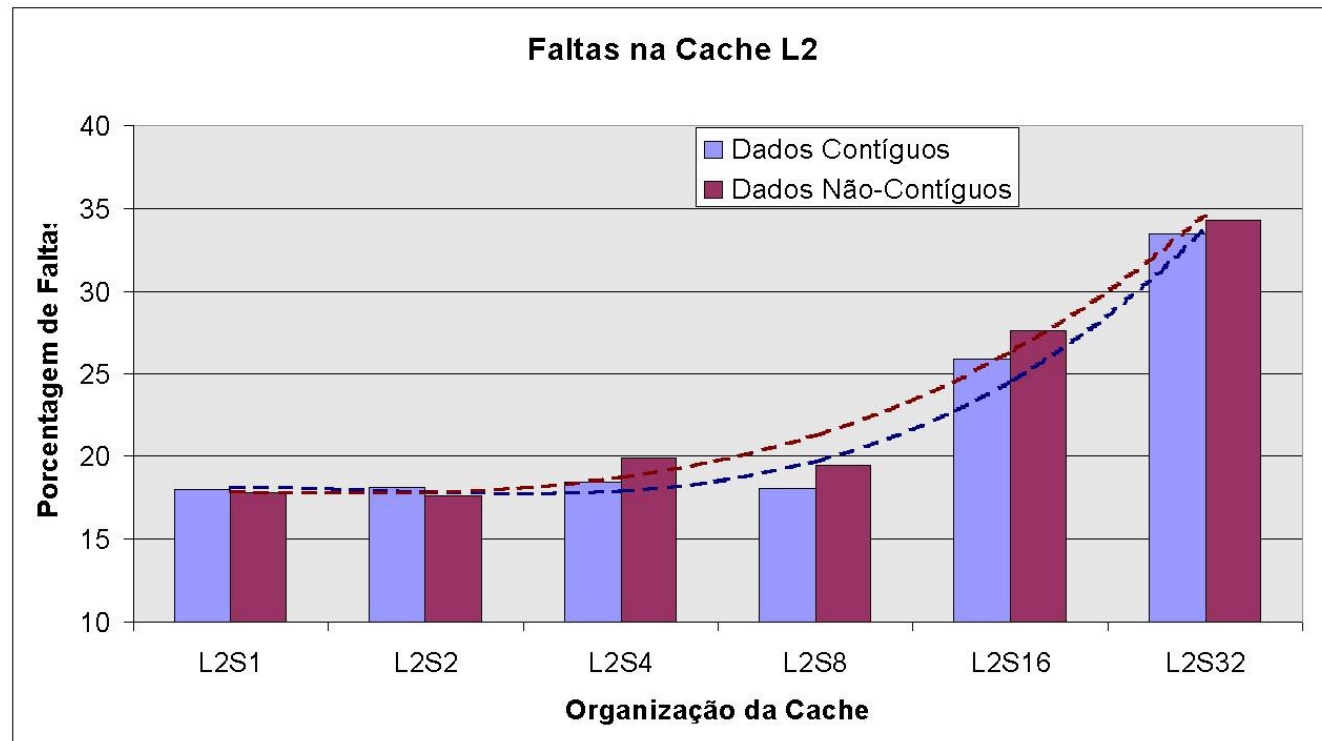Grupo de Processamento
Paralelo e Distribuído

.INf
INSTITUTO
DE INFORMÁTICA
UFRGS

- Multi Core Clusters on a Chip
  - Henrique Cota Freitas

# Architecture: Influence of Cache L2

- The influence of sharing L2 Cache in a Multiprocessor Chip
  - Marco Zanata Alves
  - Models for L2Cache:
    - Completely shared ?
    - Shared by groups/ clusters ?

- Actually using L2
  - Niagara – SUN
  - Power5 – IBM
  - Intel, AMD, etc

**Faltas na Cache L2**



GPPD
Grupo de Processamento
Paralelo e Distribuído

UFRGS

- Questions:

  - Which is impact on power consumption of a virtual machine instance?

  - How frequency reduction influences application performance and power consumption on virtualized system?

# NUMA – thread and data mapping

- NUMA challenges
  - Memory accesses to remote nodes
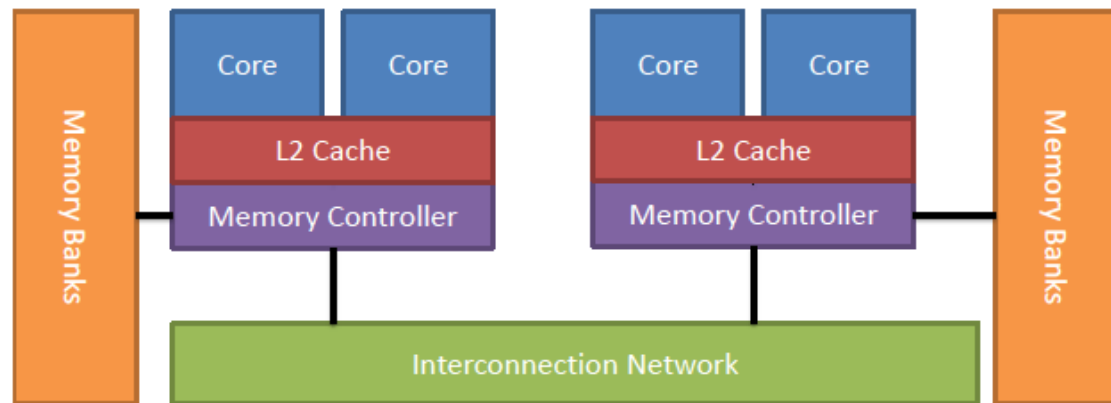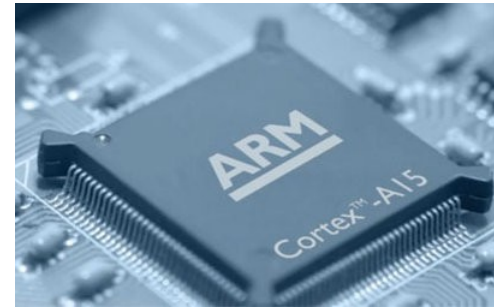    - Contention on interconnections
    - Communication across memories



Fig.1  NUMA Architecture

- Data distribution is important

- Solution: manage memory and thread affinity

- Goal: optimize memory access latency and bandwidth for the whole application

Matthias Diener, Eduardo Cruz

GPPD
Grupo de Processamento
Paralelo e Distribuído

.Inf
INSTITUTO
DE INFORMÁTICA
UFRGS

- **Energy challenges:**
  - Minimize energy consumption.
  - Attain performance restraints.

- **Solution: Include low power processors as an HPC resource.**

- **Goal: Efficient utilize low power and high performance processors to improve the energy efficiency.**



Cluster of ARM Processor

Edson Padoin, Daniel Oliveira

# Low-power processors for HPC

- Question:
  - How interesting is to employ ARM processors on HPC application? Which are the main limitations?

- Experimental platform:
  - 8 nodes clusters (pandaboards)   ARM 9 dual-core (1.2 GHz), 1 GB RAM (DDR 2)      Linux Kernel

- People:
  - Rodrigo K. Ferreira (undergrad student,

- Advisor: Alexandre Carissimi

- Evaluate network bandwidth
  - Netperf benchmark

- Evaluate MPI applications and their limitation on ARM processor clusters (memory)
  - NAS Benchmarks (MPI)

GPPD
Grupo de Processamento
Paralelo e Distribuído

.Inf
INSTITUTO
DE INFORMÁTICA
UFRGS

Laércio Pilla

- # Scheduling and data distribution

  - ## NUMA machines

    - NUMA factor
    - [Cache] memory latencies
    - Cache misses

  - ## Clusters of NUMA machines

  - ## Heterogeneous/hybrid machines

    - CPU + GPU

  - # ARM

- Case studies

  - SPECFEM3D

  - Ondes3D

- Tools

  - Charm++

  - Hwloc

  - SGPU_2

  - StarPU



Dimitri Komatitsch, http://komatitsch.free.fr/

# I/O Scheduling

The goal: **To investigate the scheduling of I/O operations on parallel file systems**

- The tool: aIOLi (Lebre et al., 2006)

# I/O Scheduling

- How to improve the scheduler's decisions?

  - Hypothesis: **include information about the applications' access patterns**

- Subject of **Francieli Zanon Boito**'s PhD

  - Advised by **Philippe Navaux** and **Yves Denneulin**

  - 2010 to 2013

# Tasks for Parallel Programming

**What we do**

# Kaapi

- A French project - http://kaapi.gforge.inria.fr/
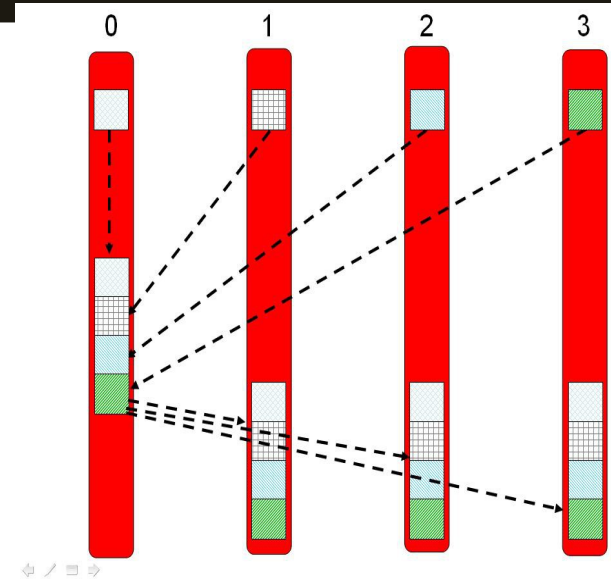- C++ library that provides an API for parallel programming based on tasks.
- A shared global address space
  - You create objects inside it with the keyword "shared"
- A task is a call to a function, prefixed by "fork"
  - Just like Unix / Cilk
  - Tasks communicate through objects *shared*.
  - Tasks specify the access mode to shared objects (read/write)
- The Kaapi runtime builds the data-flow graph and uses it to schedule the tasks.

# Tasks with MPI?

# The MPI task

- MPI defines tasks that:
  - Have their own address space,
  - Communicate with other tasks through messages.

  - Usually all are launched at the start of the program.
- The mapping "**MPI task" / O.S.** is not specified.
  - Usually, 1 task == 1 heavy process (O.S. view);
    - MPI-2 has somewhat reinforced this common understanding
  - Some MPI Distributions use threads (MPICH);
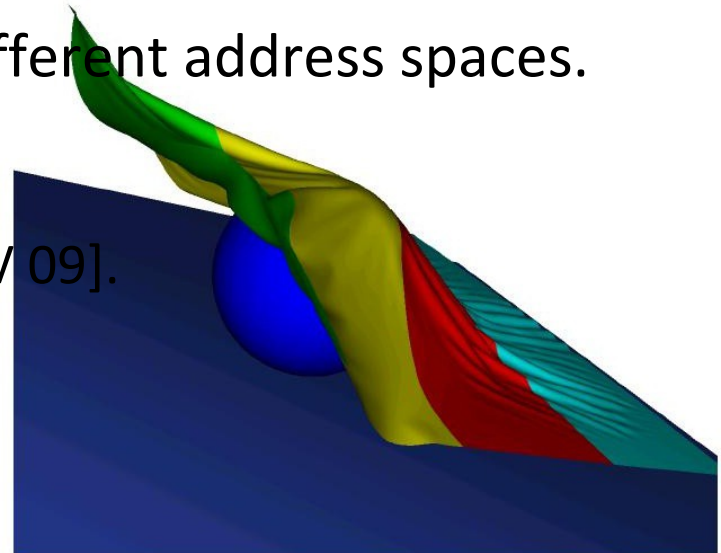    - A-MPI defines an abstract task (Urbana Champaign)

# Dynamic Tasks in MPI: D&C and Spawn

- Using MPI-2

- Program with Divide & Conquer techniques.
- Use MPI_Comm_spawn to (recursively) create new tasks.
  - 1 task == 1 (MPI) process.

- Make sure that the children tasks may communicate with their parent
  - Have the parent send the children input data, and then block.
  - Have the children send their results back to the parent.
- This implies very large-grained parallelism, but at least:
  - You can benefit from dynamic resources.
  - You can improve the load balance.

GPPD
Grupo de Processamento
Paralelo e Distribuído

.Inf
INSTITUTO
DE INFORMÁTICA
UFRGS

# Tasks & Heterogeneous Parallel Programming

- Adaptive Work Stealing already uses 2 algorithms
  - 1 sequential, 1 parallel.
- Why not using **2 different implementations** (methods) to run on a container, e.g. **one for CPU and one for GPU**?
  - The merge() method handles the different address spaces.
- This has been done partially

  by B. Raffin and E. Hermann [EGPGV 09].
- J. Lima´PhD
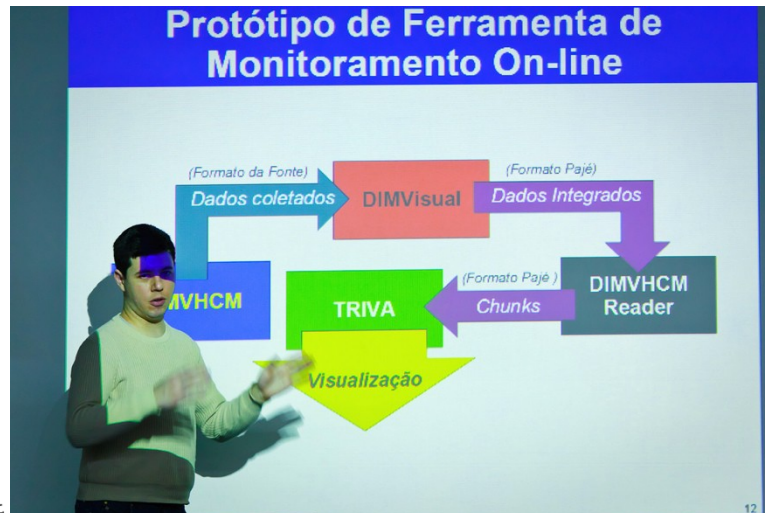  - V. Danjean, T. Gautier, N. Maillard

GPPD
Grupo de Processamento
Paralelo e Distribuído

INSTITUTO
DE INFORMÁTICA
UFRGS

# Monitoring parallel applications

- Using 3D representations to visualize the behavior of parallel programs
  - Resources (x, y) vs. Time (z)
  - Rotations, zooms…

Lucas Schnorr

Rafael Tesser

# Cloud Computing

- Question:
  - Is cloud computing appropriate to execute HPC applications?
- Experimental Platform: Microsoft Azure (PaaS and IaaS)

- People:
  - Eduardo Roloff (graduate student, MsC thesis,
  - Advisor: Philippe Navaux; co-advisor: Alex Carissimi

- Goals and methodology
  - Evaluate Azure PaaS and IaaS solutions on MPI and OpenMP applications
  - NAS benchmarks, Climate applications (BRAMs)
  - Analysis: execution time and system resources

GPPD
Grupo de Processamento
Paralelo e Distribuído

.InF
INSTITUTO
DE INFORMÁTICA
UFRGS

- Joint project with CPTEC, Brazilian National Weather Forecast

# Atmosfera Massiva Project

**Team:**
- Distributed and Parallel Processing Group – INF-UFRGS.
- Weather Forecasting and Climate Research Center – INPE.
- National Laboratory of Scientific Computing (LNCC).

**Study Topic:**
- Meteorological and environmental models.

**Objectives:**
- Study on the impact of multi-core architectures and multi-level parallelism for meteorological and environmental models.
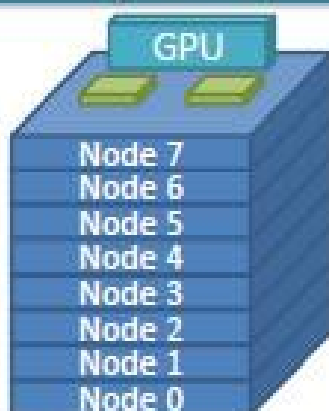
# Atmosfera Massiva Project

**Challenge:**
- Achieve better application's parallelism expression in order to use all the potential from new multi-core architectures.

Parallel file systems
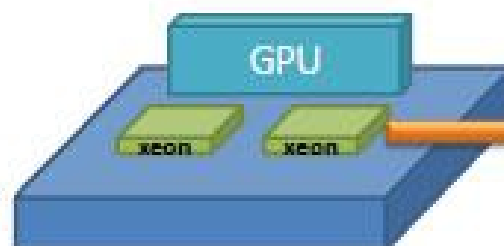
Message passing and shared memory communication

Cache memory architecture

Intra-chip interconnection

GPU

Node 7
Node 6
Node 5
Node 4
Node 3
Node 2
Node 1
Node 0

Multi-core parallelism level

CUDA programming for GPU

GPU

xeon    xeon

Core    L2    Core

Core    L2    Core

# High Performance Computing for Geophysics Applications

Partners: INRIA Grenoble, Bordeaux, Pau
BCAM, BRGM, UFRGS, UNAM, UJF

The HPC-GA project is gathering an international, pluridisciplinary consortium of leading European and South American researchers featuring complementary expertise to face the challenge of designing high performance geophysics simulations for parallel architectures.

# Load Balancing Strategies

# Motivation

## Climate and Weather Models

High computational demands

Typically run on large supercomputers

## Obstacles to Scalability

Lack of sufficient parallelism

Load imbalances

Static sources: day/night cycle, topography

Dynamic sources: atmospheric phenomena

Typical solution: changes in model's code

Requires intimate knowledge of the application

Needs to be redone for each source of imbalance

# Example of Imbalance

Weather Forecast, Feb.2010          Processor Load, P=64 (8x8)

# Dynamic Load Balancing for Weather Models via AMPI

## Eduardo R. Rodrigues

IBM Research – Brazil

edrodri@br.ibm.com

## Celso L. Mendes

University of Illinois – USA
cmendes@illinois.edu

## Philippe Navaux

Univ.Fed.RGS - Brazil

navaux@inf.ufrgs.br

## Laxmikant Kale

University of Illinois – USA

kale@illinois.edu

## Jairo Panetta

CPTEC/INPE - Brazil

panetta@cptec.inpe.br

# Approach on Charm++

## Leverage Charm++ Run-Time System

Support for MPI applications via Adaptive-MPI

Minimal changes required to original MPI code

Some of the changes can be automated

## Charm++ Load Balancing Framework

Explores migration capability in Charm++

Balancing policies based on observed load and/or communication traffic

Same balancers can be applied to different codes

Various balancing policies available

Easy to create/code new policies

Balancers are oblivious to application code details

# Case-Study: BRAMS Weather Model

## BRAMS Roots

RAMS, from Colorado State University

Model adapted for the tropics

Software structure modernized at CPTEC-Brazil

## Major BRAMS Features

Fortran90 + MPI parallelization

Open source, many thousands of lines of code

Research and production versions available

Support for multiple nested grids

Recently extended with a coupled aerosol and tracer transport model (CATT-BRAMS)

Daily production use, across Brazil and abroad

# Load Balance Approach

## Adaptive MPI (AMPI):

MPI implementation based on Charm++

Transforms MPI ranks into *user-level* threads

Processor Virtualization:

    Multiple threads (MPI ranks) per processor

    Mapping controlled by the Charm++ runtime

    Example: MPI code on 4 processors

**MPI:**



**AMPI:**

# Load Balance via AMPI

## Key AMPI Feature: Migration

May migrate threads across processors during execution to enable load balancing (LB)

Example: migration of thread 2 from P1 to P0



Dynamic, measurement-based load balancing

Measurements can be turned on/off in given phases

Balancing/migration points can be selected by user

## Load distribution after 600 tsteps

P=64, 1024 AMPI threads (=1024 MPI ranks)

# Mapping of Threads Before/After LB

(<u>Random colors</u> used here, for illustration only)

**initial mapping**
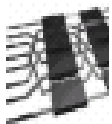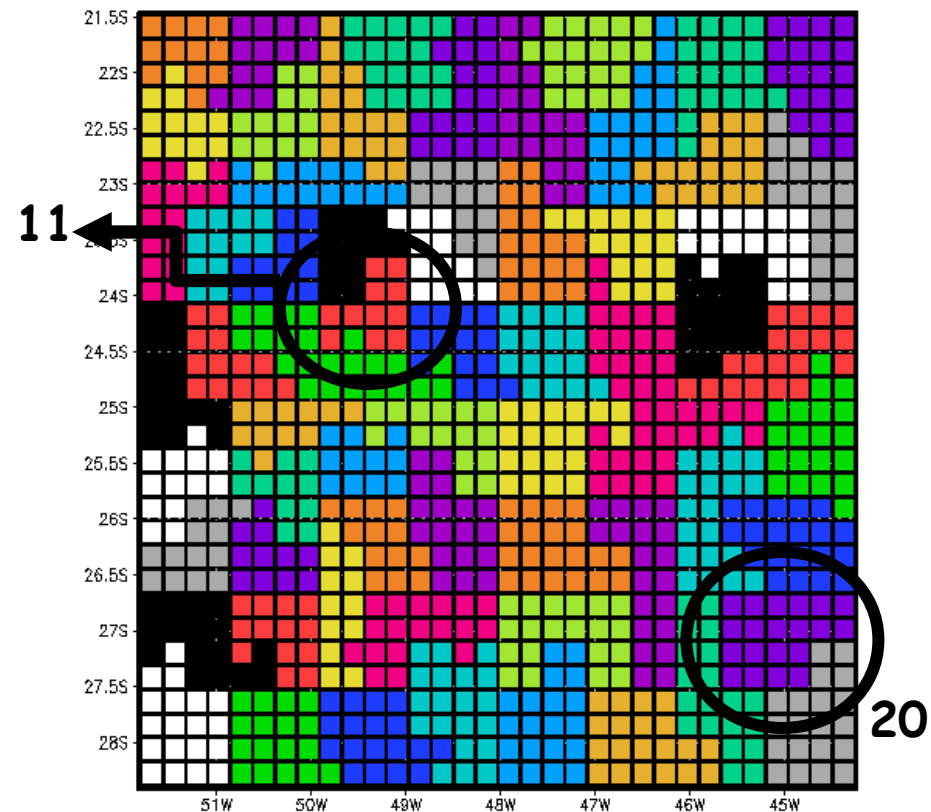
**mapping after LB**

# LB Leads to Variable #Threads/Proc

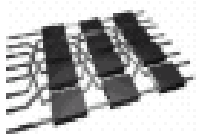Such that **load** is uniform across processors
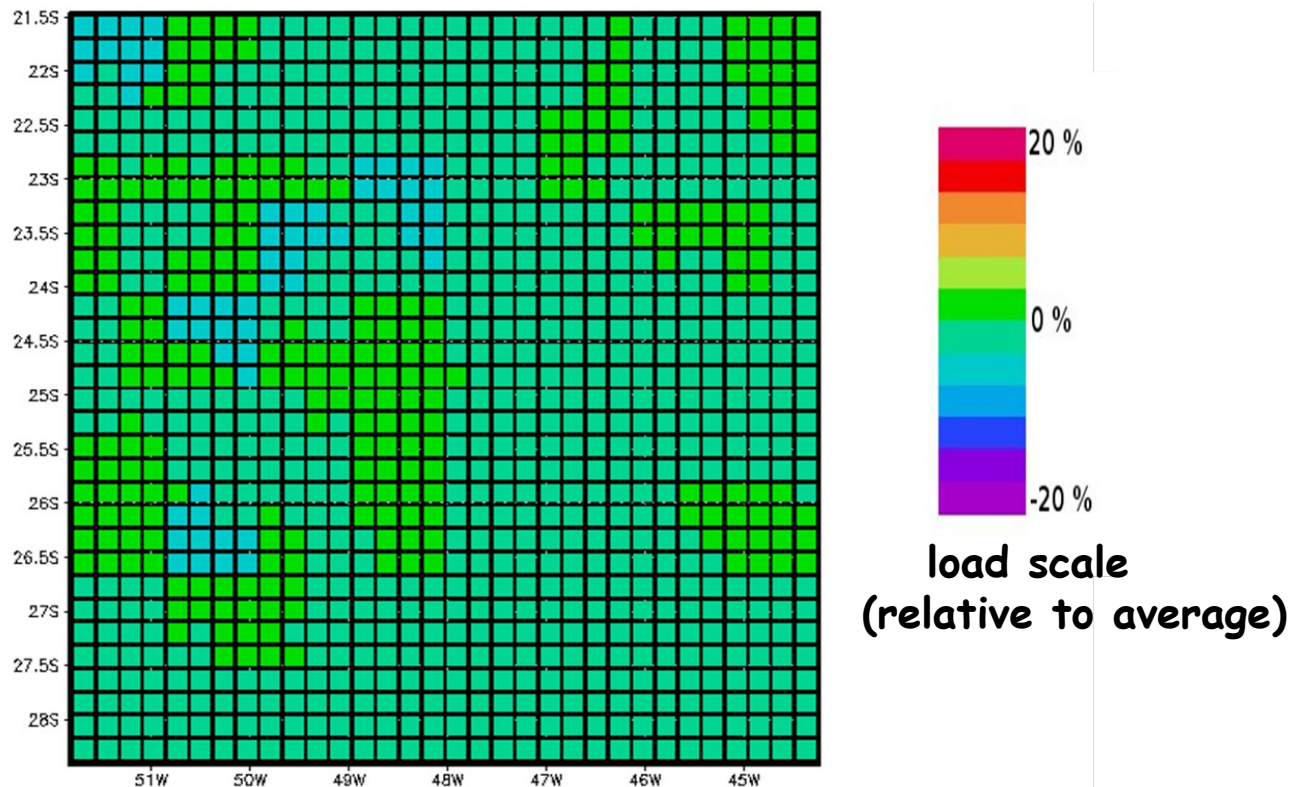


initial mapping    mapping *after* LB

# Load distribution after Load Balance

Hilbert-LB balancer used (good for 2D domains)



load scale
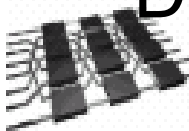(relative to average)

Virtual processors support:

Removal of global an static variables;

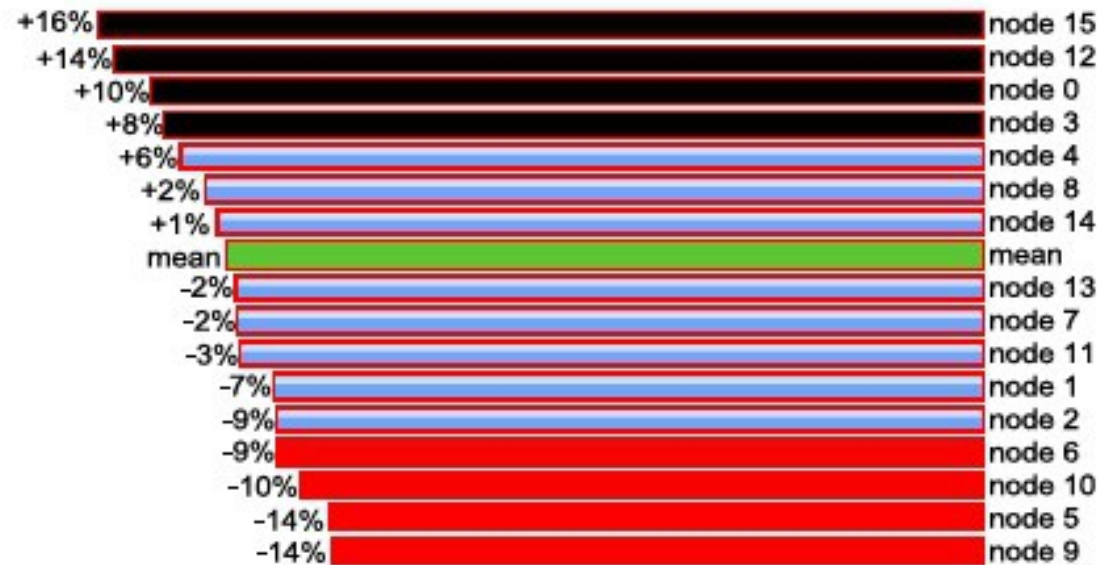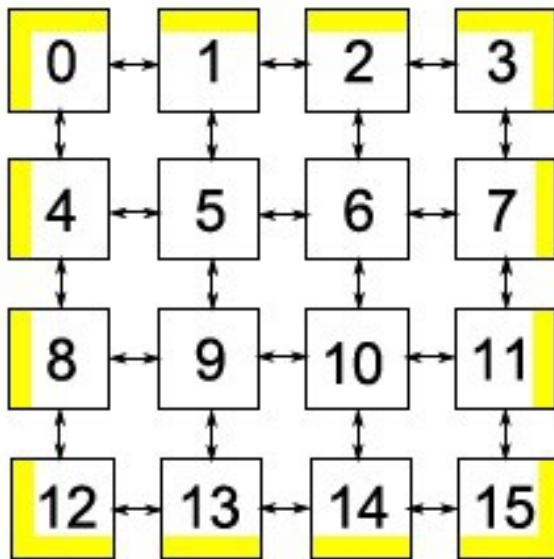- due to the use of user-level threads in place of processes;

Support to process migration:

- Implementation of functions for data serialization;
  - PUP functions: **P**acking and **U**n**p**acking;

Destruction and creation of MPI_Request variables.

GPPD
Grupo de Processamento
Paralelo e Distribuído

# Ondes 3D imbalance

# Thank you

navaux @ inf.ufrgs.br

http://gppd.inf.ufrgs.br/