# COGNITIVE COMPUTING FOR PERSONALIZED MEDICINE

**Project Supervisor:**
Alfonso Valencia
alfonso.valencia@bsc.es
https://www.bsc.es/discover-bsc/organisation/scientific-structure/life-sciences

**Summary:**
The Computational Biology group led by ICREA Professor Alfonso Valencia, within the Life Sciences Department at Barcelona Supercomputing Center (BSC) is seeking a highly motivated post-doctoral researcher to develop machine learning and cognitive computing solutions to aggregate and analyse complex phenotypic and genotypic data in the context of Personalized Medicine. The project is intended to be carried out in collaboration with IBM Academic Initiative (https://developer.ibm.com/academic/) in the framework of the existing agreement between BSC and IBM. The selected candidate will work in a highly sophisticated HPC environment, will have access to systems and computational infrastructures, and will establish collaborations with experts in different areas.

The development of Personalized Medicine critically depends on advance data analytic methods able to extract value from massive and complex data sets. In 2016 the group has been awarded with the BBVA Foundation Grant for the project "PerMed: Precision Medicine from Big Data to Cognitive Computing" that proves Precision Medicine as the optimal environment for application of advanced machine learning techniques and innovative artificial intelligence technologies.

The overall goal of the current project is to predict disease consequences and potential therapeutic interventions using both machine learning techniques, such as Support Vector Machines (SVM), and cognitive computing systems, namely IBM Watson (https://www.ibm.com/watson/). Although SVM-based algorithms have proved to be well suited for biomedical language processing and text analytics (Krallinger et al. Ann N Y Acad Sci. 2009; Krallinger et al. Methods Mol Biol. 2010, Cañada et al. Nucleic Acods Res. 2017), novel approaches such as cognitive computing are generating competitive results for text and data mining and integration.

IBM Watson is a cognitive system that is able to leverage big data to accelerate discovery in biomedical research (Chen et al. Clim Ther. 2016; Ahmed et al. IEEE Pulse. 2017). IBM Watson achieves accurate reasoning and decision making by combining capabilities in natural language processing, dynamic learning, and hypothesis generation and evaluation. In particular, IBM Watson is able to tease apart unstructured data to dynamically identify and evaluate inferences between text passages by using contextual understanding and supporting evidence (High, Rob. REDP-4955-00. 2012). IBM Watson has been recently applied to Oncology studies (Marchevsky et al. Hum Pathol. 2017; Lim and Lee. J Gynecol Oncol. 2017) and many large research centers are participating in IBM programs such as Watson Genomics, a service that combine cognitive computing and tumour sequencing data.

The proposed activity will be centered on acquisition, processing, integration and analysis of large volumes of data extending across genotypic characteristics (mutations and genomic alterations, epigenomic features, chromatin structure) and phenotypic characteristics (gene expression levels, disease and symptoms descriptions at ontological level and from medical records). The candidate

will develop methods to compare and improve performances of SVM- and IBM Watson-based analytical solutions for genotype-phenotype data integration to assist personalized medical practices.

**References:**

- Krallinger M, Rojas AM, Valencia A. Creating reference datasets for systems biology applications using text mining. Ann N Y Acad Sci. 2009 Mar;1158:14-28.
- Krallinger M, Leitner F, Valencia A. Analysis of biological processes and diseases using text mining approaches. Methods Mol Biol. 2010;593:341-82.
- Cañada A, Capella-Gutierrez S, Rabal O, Oyarzabal J, Valencia A, Krallinger M. LimTox: a web tool for applied text mining of adverse event and toxicity associations of compounds, drugs and genes. Nucleic Acids Res. 2017 May 22.
- Chen Y, Elenee Argentinis JD, Weber G. IBM Watson: How Cognitive Computing Can Be Applied to Big Data Challenges in Life Sciences Research. Clin Ther. 2016 Apr;38(4):688-701.
- Ahmed MN, Toor AS, O'Neil K, Friedland D. Cognitive Computing and the Future of Health Care Cognitive Computing and the Future of Healthcare: The Cognitive Power of IBM Watson Has the Potential to Transform Global Personalized Medicine. IEEE Pulse. 2017 May-Jun;8(3):4-9.
- Marchevsky AM, Walts AE, Wick MR. Evidence-based pathology in its second decade: toward probabilistic cognitive computing. Hum Pathol. 2017 Mar;61:1-8.
- Lim S, Lee KB. Use of a cognitive computing system for treatment of cervical cancer. J Gynecol Oncol. 2017 May 29.
- High, Rob. The Era of Cognitive Systems: An Inside Look at IBM Watson and How it Works. REDP-4955-00. December 2012. http://www.redbooks.ibm.com/abstracts/redp4955.html?Open

**Main Responsibilities:**

- Development of his/her research and career
- High standard of research performance
- Cooperation with supervisor and co-workers

**Expected skills:**

- Strong computational background
- Knowledge of principles and techniques of the subject discipline
- Organizational and project management skills
- Written and oral communication skills
- Leading, coordinating and supervising abilities

**Qualifications & Experience**

Interesting candidates must have Ph.D. training in Computational Biology, Bioinformatics, or related disciplines, must be able to work independently and must have critical thinking.

# *In silico* Engineering of Drug Efficacy in Precision Medicine

**Project Supervisor:**
Víctor Guallar
victor.guallar@bsc.es
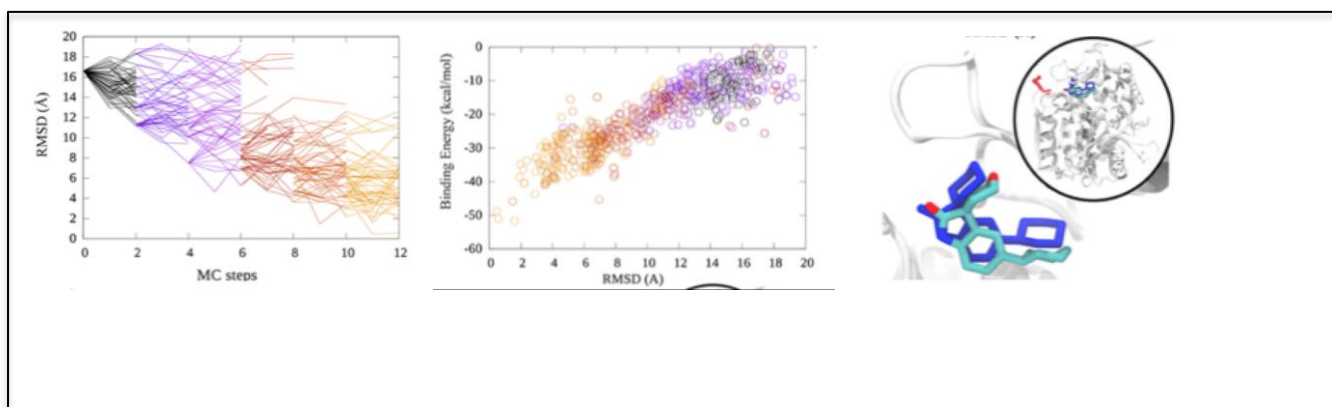https://www.bsc.es/life-sciences/electronic-and-atomic-protein-modelling

**Summary:**
Developments in software and hardware have revamped the use of molecular modeling in mapping protein-drug interaction mechanisms. As a consequence, we see a clear trend among industry to prioritize early in silico efforts in design pipelines. Sanofi, for example, signed in 2015 a $120M deal with Schrödinger, the largest molecular modeling software company. Related to these developments, our research has centered on two main aspects: i) developing state of the art software for mapping protein-ligand interactions, ii) usage of this (an others) software in applied studies focusing on solving biophysical and biochemical problems.

The software developed, Protein Energy Landscape Exploration (PELE), is today one of the best techniques for studying the mechanism of protein-ligand association, having been described as "an impressive accomplishment" by the latest international blind assessment of docking techniques (http://www.csardock.org) (*Carlson*, 2016). Application studies involved multiple pharmacological projects, for example in collaboration with AstraZeneca (*Edman*, 2015), and several enzyme engineering studies, for example with the world leader industrial partner Novozymes (*Acebes,* 2016). Importantly, with the development of Nostrum Biodiscovery, the first spin off from BSC, we have transferred our technology into the industrial world.

These advances have allowed the involvement of molecular predictions in personalized medicine (a critical step in the coming precision medicine). Using (PELE), our lab has created a protocol to predict resistance to antiretroviral drugs in the protease of HIV-1. Together with the IrsiCaixa Aids Institute, we tested the protocol with real patients sequences, providing the first tool that accurately can predict resistance using molecular modeling techniques (Hosseini, 2016).

Still, an accurate description of the receptor-drug interaction mechanism is a daunting task, requiring several hours/days of heavy computation. Adding an adaptive machine learning procedure to our Monte Carlo sampling technique, we have introduced a breakthrough in the sampling procedures: Adaptive-PELE. In a recent study (*Lecina, 2017*) we show remarkable performance in mapping the protein-ligand energy landscape, being able to reproduce the binding mechanism in complex systems (from solvent to the buried active site) in less than half an hour, or the active site induced fit in less than 5 minutes.

**Figure.** Enzyme-substrate induced fit sampling for epoxide hydrolase using the adaptive technique. Top: RMSD to the bound crystal. Middle: all-atom binding energy identifying the bound state (<2Å RMSD). Bottom: agreement between the predicted (blue) and the bound crystal (atom-type colorized); initial exploration position (red molecule in the inset). Notice that the bound state is found in ~10 MC steps (~3 minutes).

In this project, we aim at two goals: 1) Combine Adaptive-PELE with free energy procedures, such as + Markov State Models (MSM) and Free Energy perturbation (FEP) methods; 2) Test the developed combined techniques into lead optimization of drugs (bypassing drug resistance) and into new personalized medicine cases.

The first goal will developed fast and accurate binding free energy predictors, capable of rationalizing the resistance introduced by mutations, changes in affinity in congeneric drug series, etc. The second goal will expand our current tests cases into new applications and real drug design efforts. Current applications include antibodies development with the IrsiCaixa Aids institute, predicting Hepatitis C drug resistance with the Vall the Hebron hospital, etc.

**References:**

- Carlson, HA et al. CSAR 2014: A Benchmark Exercise Using Unpublished Data from Pharma. *J. Chem. Inf. Model*, **2016**, *56*:1063.
- Edman, K et al. Ligand binding mechanism in steroid receptors: From conserved plasticity to differential evolutionary constraints. *Structure*, **2015**, 12:2280
- Acebes, S et al. Rational enzyme engineering through biophysical and biochemical modeling. *ACS Catalysis*, **2016**, 6:1624
- Hosseini, A et al. Computational prediction of HIV-1 resistance to protease inhibitor. *J. Chem. Inf. Model*, **2016,** 56: 915-923
- Lecina, D et al. Adaptive simulations, towards interactive protein-ligand modeling. *Scientific Reports*, **2017**, In press

**Main Responsibilities:**
To develop (program) the (external to PELE) free energy module. Test the develop software in benchmark test and real drug design applications

**Expected skills:**
Programming in python. Knowledge in molecular modelling and viewing packages

**Qualifications & Experience:**
It is required a doctor degree in computational chemistry, biophysics or related. Experience in some of the following: i) molecular modelling, ii) i*n silico* drug design, iii) free energy methods.

## Development of a computational platform for personalized medicine in oncology

**David Torrents**
david.torrents@bsc.es
http://cg.bsc.es/cg/

Health care systems around the world will soon include the sequencing and analysis of the genome as routine, for prognosis, diagnosis and treatment of a growing number of diseases, reaching a point where a large fraction of the human population will have their genome sequenced. The essential genomic information that needs to be captured from each individual are all the differences and variations that their genome presents compared to the others in the population, and from those, identify those that are related to different traits, and in particular, to disease. Therefore, while the knowledge about the genomic basis of disease is growing fast, the complete reconstruction of genomes and its interpretation will soon be required for health. But the identification and analysis of this variation is a challenging task that requires complex algorithmic and computational approaches, often integrating different programing models and computing architectures. But the constant improvements on sequencing techniques (e.g. the new technologies for the sequence of single DNA molecules providing large sequence reads instead of the current methods based on short reads), not only keep increasing the possibilities of detection of mutations and clinical analysis, but also demands constant adaptations to different data types and formats, adding a changing layer of complexity to the development and adaptation to solutions. For this reason, it is essential to build on the BSC strong expertise in biomedical genomics, bioinformatics and computing sciences, to build new methods able to cope with the complexity of the mix environment of sequencing technologies.

The Computational Genomics group within the Life Science Department at the Barcelona Supercomputing Centre offers a position for scientists with a PhD in bioinformatics or similar disciplines to contribute to the generation of a computational platform for the analysis of tumour genomes within a clinical context. Our group is focused in the analysis of whole genome data to uncover the mechanisms and functional consequences of somatic chromosome rearrangements that are often associated with a more aggressive tumour progression. The proposed activity will be centered on the development and application of computational tools in the context of a local initiative to start applying personalized medicine protocols to oncological patients. This activity will be carried out within a collaborative frame that involves sequencing centers, hospitals, and our center, where we will be focusing in the identification of mutations. These specific tasks are part of the activity of the group, which is currently involved in several worldwide initiatives and collaborations that aim to understand the genomic basis of Cancer. Our aim is, ultimately, to generate and develop translational protocols for a more efficient and personalized treatments of tumours. This has positioned our group and center among the top institutions involved in the analysis of disease genomes (Puente et al, Nature 2011; Moncunill et al, Nature Biotech. 2014, Puente et al, Nature 2015).

Selected candidates will be dealing with the analysis of BIG biological data, i.e. with thousands of whole genome sequencing, in a highly competitive computational environment, at the level of research, and at the level of resources. This research frame offers postdoctoral scientist a unique opportunity to develop and improve their bioinformatics skills in the growing and expanding field of biomedical genomics. This activity will be carried out within a dynamic group that combines biological and computer sciences to seed light into the genetic basis of disease.

More information about the group and the activity can be found at: http://cg.bsc.es/cg/

**Main Responsibilities:**

- Coordinating other scientist and the research related to the generation of bioinformatics tools for genome analyses in the context of cancer research and within an HPC environment.

- Coordination of different projects related to the analysis of tumour genomes.

**Skill Specifications**

- We seek for highly motivated postdocs that want to contribute to the understanding of the genetic basis of cancer.

- Candidates must be ready to push this project with certain independence and, at the same time, to work within a cooperative environment.

**Qualifications & Experience**

- Candidates with a PhD in bioinformatics will be prioritized.
- Strong experience in computational environments related to biomedical genomics.
- Experience with Next Generation Sequencing data.

*Proposals for postdoctoral projects by:*
*Prof. Vassil Alexandrov and Dr. Isaac Rudomin*
*Extreme Computing Group, Computer Science Department, BSC*

**1.Monte Carlo and Deep Learning Methods for Enhancing Crowd Simulation:** Simulating realistic large-scale crowd behaviours is a complex endeavour. The use of real data and realistic perception models is required. And once the behaviour is established, one must generate and animate varied characters for realistic visualization without consuming too much memory and computing resources. At the Extreme Computing group at BSC, we have been working on the development of methods for simulating, generating, animating and rendering crowds of varied aspect and a diversity of behaviours. The focus is on efficient simulations of large crowds that can be run on low cost systems because we use of modern programmable GPUs and to scale up for even larger crowds: We subdivide the simulation into different regions and distribute the work to different nodes by using MPI and to the different CPUs and GPUs in each node by using OmpSS and within each GPU we use CUDA, striving to use all the computational resources available in these heterogeneous clusters. The ultimate goal is to simulate and visualize very large crowds (of over a million or several million characters) using a variety of advanced architectures in real time. Further given our groups experience in Monte Carlo methods we would like to explore advantages of the methods and combine those with other machine learning approaches, in particular Deep Learning, and develop new highly parallel Monte Carlo and Deep Learning approaches that use real data and simulation for more realistic large scale crowd simulation.

Experience in machine learning and HPC (High Performance Computing) as well as knowledge of Deep Learning frameworks such as Tensorflow and Keras and HPC programming paradigms such as MPI-OpenMP and CUDA are desirable. The postdoc will be expected to develop further the existing code and implement a Monte Carlo Deep Learning algorithm for crowd simulation. It is expected to test it on an NVIDIA P100 GPU based cluster and perform further comparative studies. Contact Dr Vassil Alexandrov and Dr Isaac Rudomin. ([vassil.alexandrov@bsc.es]() )

**2.Use of Machine Learning techniques for the generation of realistic Appearance and Behaviors of Characters in Crowds:** Simulating realistic large-scale crowd behaviours is a complex endeavour. The use of real data and realistic perception models is required. And once the behaviour is established, one must generate and animate varied characters for realistic visualization without consuming too much memory and computing resources. At the Extreme Computing group at BSC, we have been working on the development of methods for simulating, generating, animating and rendering crowds of varied aspect and a diversity of behaviours. The focus is on efficient simulations of large crowds that can be run on low cost systems because we use of modern programmable GPUs and to scale up for even larger crowds. Further, given our groups experience in generating and visualizing crowds, we would like to enhance the realism of the character's appearance and behavior by using machine learning techniques.

Expertise in Machine learning and Computer graphics, and knowledge of OpenGL/GLSL, WebGL and Deep Learning frameworks such as Tensorflow and Keras are desirable. The postdoc is expected to generate authoring tools that enhance the realism in appearance and behavior of characters in crowds. Contact: Dr Isaac Rudomin. ([Isaac.rudomin@bsc.es]() )

**Proyectos Biomecánica de CASE**

**Implementación de la cuantificación de la incertidumbre de modelos FSI, HPC para investigar biomarcadores no invasivos de la enfermedad coronaria.**

Las simulaciones con valor predictivo son deseables en prácticamente todos los escenarios de aplicación para modelos numéricos. Para ese fin, es necesario tomar en cuenta información incompleta e inexacta de las mediciones clínicas para obtener distribuciones o barras de error en lugar de resultados determinísticos.

Este trabajo se centra en el desarrollo de un marco de cuantificación de la incertidumbre basado en un modelo de sustitución eficiente basado en el adjunto junto con una simulación de Monte Carlo. El marco desarrollado se aplicará a modelos computacionales de enfermedad arterial coronaria (CAD) y el análisis de biomarcadores clínicos, en particular la reserva de flujo fraccional no invasivo. El objetivo de este trabajo es el desarrollo y aplicación de un nuevo enfoque para cuantificar el impacto de parámetros de entrada de modelo inciertos y para permitir simulaciones predictivas proporcionando distribuciones de probabilidad, barras de error o estimaciones de peor caso.

Con el fin de modelar el sistema cardiovascular se utiliza un modelo de interacción fluido-estructura (FSI) de las arterias coronarias. La razón para emplear un enfoque de FSI es que el estado del arte de la hipótesis de arterias coronarias rígidas y estáticas es una simplificación extrema de la circulación coronaria y su función. Estas simplificaciones pueden resultar útiles en casos simples sin grandes estenosis, pero no han podido tener un valor predictivo para las arterias coronarias tortuosas y altamente enfermas, que es el caso de la mayoría de los pacientes que sufren de CAD y que están en riesgo real cuando son sometidos a mediciones invasivas, hiperémicas inducidas por fármacos.

Project Leaders:  Jazmín Aguado-Sierra y Mohammad Kouhi (jazmin.aguado@bsc.es )

*Estudio y caracterización de emisiones en motores dual-fuel utilizando el método large-eddy simulation (LES)*

El proyecto está dedicado el desarrollo e implementación de un modelo de atomización y evaporación que permita estudiar las emisiones de NOx y otros contaminantes en motores de combustión interna alternativos. Este proyecto se enmarca en el contexto de investigación aplicada asociado a la mejora de la eficiencia y reducción de emisiones de los sistemas de transporte y en concreto, de los motores *dual-fuel*. Los nuevos diseños y tecnologías en motores de automoción han supuesto una mejora considerable de su eficiencia y prestaciones, aunque todavía dependen mayoritariamente del proceso de atomización y quemado del combustible principal. En este contexto, el proyecto desarrolla una metodología numérica basada en simulaciones numéricas de alta fidelidad mediante el método de las grandes escalas o *large-eddy simulation* (LES) en combinación con el uso de la supercomputación (HPC) para modelar el proceso de atomización, evaporación y combustión en condiciones de motor. Mediante el uso de métodos numéricos y algoritmos de HPC con tecnologías de *exascale*, se hace posible la resolución de estos problemas y se presentan iniciativas avanzadas para el desarrollo de tecnologías de simulación en ingeniería. El proyecto incluye el desarrollo e implementación del modelo de atomización Σ-Y en el contexto LES acoplado con un modelo de combustión avanzado basado en *unsteady flamelets*. La propuesta de trabajo se caracteriza por una alta interdisciplinaridad ya que involucra actividades relacionadas con el modelado fisicoquímico, matemática avanzada, ingeniería mecánica y de computación en aplicaciones de alta relevancia a nivel industrial y tecnológico. El proyecto está centrado en un tema relevante en el contexto internacional y posiciona la investigación española en el nivel más alto en temas de modelado y simulación numérica avanzada. El proyecto incluye colaboraciones con la industria y otros centros tecnológicos de alta relevancia a nivel internacional como Cerfacs, Siemens AG, CMT Motores Térmicos, Idiada y Technical University of Berlin.

**Responsable:**
Daniel Mira, daniel.mira@bsc.es
Senior Researcher, Computer Applications of Science and Engineering (CASE) Department.
Barcelona Supercomputing Centre (BSC)

**Wind Resource Assessment for wind farms over complex terrain including thermal effects.**

The Postdoc proposal will be in the framework of the development of Computational Fluid Dynamic (CFD) RANS models, and strategies for coupling them with meteorological models to improve the wind resource assessment (WRA) over complex terrain wind farms. For WRA the wind industry is increasingly relying on CFD models. The use of CFD-RANS simulations provides an accurate prediction of wind velocity and turbulent kinetic energy fields for the entire wind farm. These fields are necessary to know where is possible to place the wind turbines in the wind farms, and to accurately predict the potential energy production of the wind farm as function of the wind turbine positions.

In emplacements where strong convective or nocturnal boundary layers are developed, neutrally stratified CFD models cannot predict the time averaged wind velocity distribution correctly. In order to decrease the uncertainty of WRA, the effect of buoyancy forces and heat transport over the atmospheric boundary layer (ABL) should be included in CFD models. Thermal effects in the atmosphere are dynamic and must be simulated as transient. Moreover, Coriolis forces need to be included in CFD models when simulating over hilly terrains or in presence of strong nocturnal boundary layers. This further complicates the problem and makes transient thermal simulation the only viable choice.

A transient daily cycle simulation is needed to characterize the wind resources and power distribution over a wind farm with strong thermal stratification in the boundary layer. For such simulations the amplitude of the thermal variations of a representative onsite daily cycle needs to be imposed as boundary condition. The imposition of the boundary conditions for the daily cycle simulation remains an open question, it is not known a priori which is the best methodology to find the proper boundary conditions and forcings that drive the microscale flow. In the present proposal the boundary conditions will be inferred from onsite mesoscale simulations of at least one-year duration, together with onsite velocity and temperature measurements.

To run a daily cycle simulation in a given wind farm, a precursor single column model that simulates a daily cycle over flat terrain is used to impose inlet boundary conditions along time on the wind farm. This single column simulation needs to be driven from a mesoscale model (wall temperature and geostrophic velocities).

The Postdoc will work on the development and investigation of new downscaling methodologies from mesoscale to microscale simulations for WRA using HPC applied to the wind energy industry through the imposition of asynchronous boundary condition and forcing mechanisms. Moreover, he will be involved in the implementation of improved physical and numerical CFD models for an accurate prediction of the wind distribution over complex terrain with thermal coupling and/or Coriolis forces.

The developed methodologies will be tested in the frame of NEWA-ERANET project, where several complex sites are being instrumented. Wind farm measurements disposed by Iberdrola Renovables will also be available for testing.

**Semantic Technologies:**

One of the most difficult problems concerning big data is the integration of many heterogeneous data sources of different granularity and quality. Semantic technologies are the most promising approach to deal with this problem, but there are still many issues concerning scalability of querying, enabling types of operations that are typical for graph structures that were non-existent in the classical database community, and making this technology useable for non-IT experts. A smart city is the singular example where huge amounts of data are produced, and usually can be made publicly available for integration and analysis. This is not only a domain of great potential impact, but also the ideal testbed for other types of application domains such as, for instance, finance, fraud, or social services.

Our focus topic is scalability and usability for Semantic Web queries. In terms of scalability we are especially interested in graph algorithms for efficient distribution and parallelization, graph query languages, and/or query optimization. In terms of usability (and possibly also scalability) we are interested in improving data quality, query answering in the presence of uncertain/conflicting information, and/or approximate query answering.

**Profile:**
- Strong background in graph algorithms, **or alternatively** in probabilities, reasoning with uncertainty, or fuzzy logic.

- Ideally familiar with (1) language semantics / logic theory, or (2) parallelization.
- The candidate should have a good programming background and a good level of English.

- Responsible research staff: Maria-Cristina Marinescu (maria.marinescu@bsc.es )

**Análisis de incertidumbre en geofísica de exploración a gran escala**
Responsable: **Josep de la Puente** (Geophysical Applications, CASE)

Los estudios de exploración geofísica tienen como objetivo determinar la existencia, y en su caso cuantía, de reservorios subterráneos de recursos naturales, principalmente mineros y energéticos. Para dichos estudios se requiere de un complejo y computacionalmente costoso análisis de datos que resulta en el mapa 3D "más probable" de las propiedades físicas del subsuelo. Sin embargo existe mucha incertidumbre asociada a las limitaciones de los datos empleados (numero y disposición geométrica de receptores, banda de frecuencia empleada, ruido, …) así como a la información a priori empleada (modelos previos del subsuelo). La manera en que esa incertidumbre se traduce en confiabilidad en el modelo final es costosa de obtener y difícil de interpretar.

Dado el enorme coste asociado a la explotación de los recursos naturales del subsuelo, un preciso análisis de riesgo permite clarificar lo arriesgado de la apuesta en términos económicos, medioambientales o sociales.

Proponemos establecer un sistema de análisis de incertidumbre para exploración geofísica, aplicado a problemas 3D sísmicos y electromagnéticos. Si bien su espectro de aplicación será variado, focalizaremos el esfuerzo en el hallazgo y cuantificación de reservorios de hidrocarburos. Para ello basaremos el trabajo en las herramientas de altas prestaciones desarrolladas en BSC (www.bsc.es/bsit ) que garantizan una alta precisión y eficiencia HPC.

**Multipartitioning for the solution of PDEs**

**Keywords:** *Finite element, iterative solvers, direct solvers, Fortran, MPI, runtime, OpenMP, OmpSs, load balance, task parallelism*.

Simulations of multiphysics problems driven by PDEs involve two main kernels: the *algebraic system assembly* (ASA), and the *algebraic system solution* (ASS) using iterative or direct methods. On the one hand, the assembly of the algebraic system involves a loop over the elements of the mesh. On the other hand, iterative methods consist mainly of sparse matrix vector product (SpMV) and dot product operations, thus involving a loop over the nodes. If hybrid solvers are considered (DD solvers), additional operations like LU or incomplete LU factorizations are also required.

The parallel solution on distributed memory supercomputers of such problems is based on the partitioning of the computational domain into subdomains. The algebraic system assembly does not require communication. On the one hand, iterative solvers involve mainly point-to-point communications to compute SpMV, and reduction operations to compute dot products. On the other hand, direct solvers involve very complex communication patterns.

Traditionally, the partitioning has been carried out in such a way to:

> 1) Balance the number of elements across the subdomains in order to have a balanced matrix assembly;
>
> 2) Minimize the size of the interfaces between these subdomains in order to minimize the communication in SpMV. Therefore, no control on the iterative solver operations is possible.

This is an important issue when dealing with hybrid meshes, where the number of elements can greatly differ from the number of nodes and the size of the graph of each subdomains, thus leading to very unbalanced iterative and direct solvers.

This projects aims developing, validating and testing an *adaptive and dynamic multipartitioning method* to provide an optimum load balance to each kernel: ASA and ASS.

**Tasks**:

> 1) To develop an *adaptive and dynamic mesh partitioning* strategies for both assembly and different solvers/preconditioners: DD, Krylov and direct solvers.
>
> 2) To develop a *coupling strategy* between both partitions: sending of the matrix from ASA to ASS and sending of the solution from ASS to ASA.
>
> 3) To develop strategies for the *optimum placement* of the different MPI processes to minimize the cost of the coupling communications.
>
> 4) To implement *task parallelism* (OmpSs) and *dynamic load balance* (DLB) strategies at the shared memory level to enhance the parallelism and the load balance.
>
> 5) To test the methodology for different physical problems, involving different iterative solvers and preconditioners as well as different direct solvers, on different platforms, taking advantage of the heterogeneous clusters of BSC.

**Place**: The project will involve computational mechanics and computer science aspects and will be carried out in collaboration with the CASE and CS department of BSC.

**Responsible**: Guillaume Houzeaux (CASE) and Marta Garcia-Gasulla (CS) (guillaume.houzeaux@bsc.es )

**Mejora de la escalabilidad del entrenamiento distribuido de redes neuronales (Deep Learning)**

Antonio J. Peña, Sr. Researcher, Activity Lead, Manager del NVIDIA GPU Center of Excellence
http://www.bsc.es/pena-  – antonio.pena@bsc.es
Accelerators and Communications for HPC, Grupo Programming Models, Dep. Computer Sciences

**Motivación**
Actualmente la escalabilidad del entrenamiento de redes neuronales está severamente limitada por varios factores. El principal motivo es la gran cantidad de comunicaciones colectivas periódicas que se necesitan para mantener actualizado el modelo en los diferentes nodos de cómputo, aunque también intervienen factores como la entrada/salida de datos distribuida, los *kernels* de cálculo en GPUs, los algoritmos de entrenamiento y los propios diseños de las redes.

**Propuesta de Trabajo**
Para este puesto postdoctoral se propone estudiar el modo de mejorar el entrenamiento distribuido de redes neuronales tipo "Deep Learning" utilizando  GPUs, atacando los limitadores mencionados y otros que se puedan detectar. El investigador postdoctoral se unirá a los esfuerzos que ya se han comenzado dentro del equipo de trabajo. Para ello contará con el apoyo del amplio *know-how* interno del equipo en modelos de programación, *runtime systems*, GPUs y comunicaciones.

# Performance Analysis of Large Scale Graph Algorithms in a High Performance Computing Context

**Supervisor:** Dario Garcia-Gasulla (dario.garcia@bsc.es)
**Research Group:** High Performance Artificial Intelligence

**Summary:**
Graph analytics have become an essential tool for data mining in multiple fields of interest, both for academia and industry. The particularities of graph-based algorithms however, imply that most traditional solutions for scaling those algorithms are inappropriate for large-scale graphs. At the HPAI research group we are interested in understanding, which are the bottlenecks of large-scale graph computation, with two main goals in mind.
First, to obtain experimental results by applying graph algorithms of interest to the growing set of large-scale graphs available today.
Second, to learn about the limitations of current computer architectures, and how can these are properly adapted to large-scale graph analytics.
We intend to tackle both those goals using state-of-the-art technologies, thanks to the privileged high-performance computing resources available at BSC.

**Main responsibilities:**
The candidate will identify and manage graph-processing frameworks that are appropriate for running in a HPC environment. If needed significant software contributions will be made.
The candidate will implement graph analytics algorithms of interest (new or from the state-of-the-art) on the chosen platforms. He/she will also identify use cases of potential scientific interest derived from available large-scale graphs.
The candidate will measure the performance in terms of computational efficiency using HPC tools specific for that purpose.
The candidate will evaluate the scientific relevance of the results obtained by the algorithms, and publish and present works to the community.

**Expected skills:**
- Strong programming skills
- Teamwork skills
- Experience on the analysis of software scalability, parallelism and computational efficiency