# CONVERGENCE OF BIG DATA
# WITH HIGH PERFORMANCE COMPUTING
# AND ARTIFICIAL INTELLIGENCE

Interdisciplinary Centre for Mathematical
and Computational Modelling (ICM),
University of Warsaw
as an exemplary case for
Centre of Excellence in Big Data

Discussion paper to recommend ICM
as a CoE in Big Data

Marek Michalewicz, PhD (ANU)
Director (Acting) ICM

# 1. Introducing ICM

[Interdisciplinary Centre for Mathematical and Computational Modelling](#) (ICM) was established in 1993 as a basic research and computational services unit of the University of Warsaw.

ICM is involved in interdisciplinary scientific research based on mathematical modeling, computer simulations and modeling, multi-scale and large-scale calculations, and teaching in the above areas. In addition, for the last 23 years ICM has been involved in Numerical Weather Prediction. ICM computed weather forecasts have been published for over twenty years on a very well known [meteo.pl](#) weather portal visited close to 200 million times annually. ICM researchers study problems related to civil aviation, modeling of social processes and scientometry – all based on exclusive access to specific Big Data resources. ICM is involved in securing access of Polish scientists to entire body of scientific literature, including over 8,000 journal titles and hundreds of thousands scientific books, by maintaining the [Virtual Library of Science](#), including the entire content and the rights to text mining.

ICM manages two data centers in Warsaw. The ICM Technology Center (CT-ICM) in Białołęka commissioned in 2016 has approx. 10,000 m2 of technical space with two world-class supercomputers: Petaflop Cray XC-40 (Okeanos) for traditional intensive numerical calculations, and the Huawei cluster (Enigma) for large data analytics (Hadoop, Spark) and cloud computing. The CT-ICM server room also has data storage equipment for about 20 PetaBytes of data in variety of file systems: high performance Lustre to object storage. The newly built lecture hall for 70 listeners has unique visualization equipment with 16 monitors and software that allows for displaying of huge datasets and the transmission of lectures or images from around the world.

In the domain of Big Data, High Performance Computing (HPC) and cloud services, ICM supports approximately 1,000 users from all over Poland using our supercomputers and computing, network and storage infrastructure for very large data.

ICM currently employs 110 staff members.

Separate units of ICM are devoted to different areas of Big Data techniques, including scientific visualisation, data science, real time data processing and critical decision support systems (in aviation industry), weather forecasting or data management. Besides research, ICM employs its own software development team, allowing production of quality implementations of research concepts.

ICM team has pioneered in a number of cutting edge networking solutions, both for high throughput and low latency requirements. Recently, ICM engineers have established a production 100Gbps connection over 12,375

miles CAE-1 (Collaboration Asia Europe-1) line between Warsaw and Singapore.

## 2. Infrastructure

### 2.1 Computer Centres

ICM manages two data centers in Warsaw. The ICM Technology Center (CT-ICM) in Białołęka, opened in 2016, has approx. 10,000 m² of technical space and is one of the most modern and best equipped scientific computing centers in Poland. The center is secured with continuous power supply systems at the level of 5.5 MW, computer cooling system with ice water and glycol, three independent fire safety systems and access protection systems. The second, older server room is in the IBIB Polish Aacademy of Science building at ul. Pawińskiego. It houses older computing equipment, cloud resources and storage.



*Figure 1. New ICM  Technology Centre in Białołęka suburb of Warsaw.*

### 2.2 Compute equipment

At the ICM Technology Centre in Białołęka there are two large computers: a Petaflop Cray XC-40 supercomputers (Okeanos) with 1084 compute nodes for traditional intensive numerical calculations, and Huawei (Enigma) cluster with

340 compute nodes for large data analytics (Hadoop, Spark) and cloud computing. In the server room at Pawiński street there is a Huawei (Topola) cluster with 223 compute nodes designed for smaller scale calculations.



*Figure 2. One of the most powerful research computers in Poland: Cray-XC40 - Okeanos at CT ICM in Białołęka.*

Last year, ICM built a cluster consisting of six servers equipped with the latest generation GPU accelerators (Nvidia Tesla V100) and the unique NEC Aurora Tsubasa computer with eight vector cards. This year, the cluster will be expanded with a node with a large amount of RAM (1.5TB) and further V100 Nvidia GPU accelerators. In total, ICM has computing resources of about two PFLOPS.

**2.3 Big Data storage**
In addition to computing resources, ICM has very substantial storage resources: 10 PB available on the Okeanos computer, a total of 4 PB on local drives in the Enigma node, and 6 PB of object storage intended for long-term storage of valuable data.

**2.4 Global, National and Metropolitan Networks**
ICM UW is a member of the Pionier consortium affiliating 21 national municipal networks and four other academic supercomputer centers. Pionier, a nationwide partnership to build and maintain exclusive academic and research fiber optic network has enabled access to a modern network with significant bandwidth and numerous international connections.

Two ICM server rooms are connected with each other and with other computing centers and the Pionier network with 100Gbps connections. Physical connections to network users and between backbone nodes are made as part of own network to ensure a redundant connection system while maintaining the geographical separation of fiber optic routes.

Thanks to two server rooms, ICM is able to significantly reduce the risk of failure or minimize the effects in the event of their occurrence by appropriate deployment of computing and network infrastructure components.
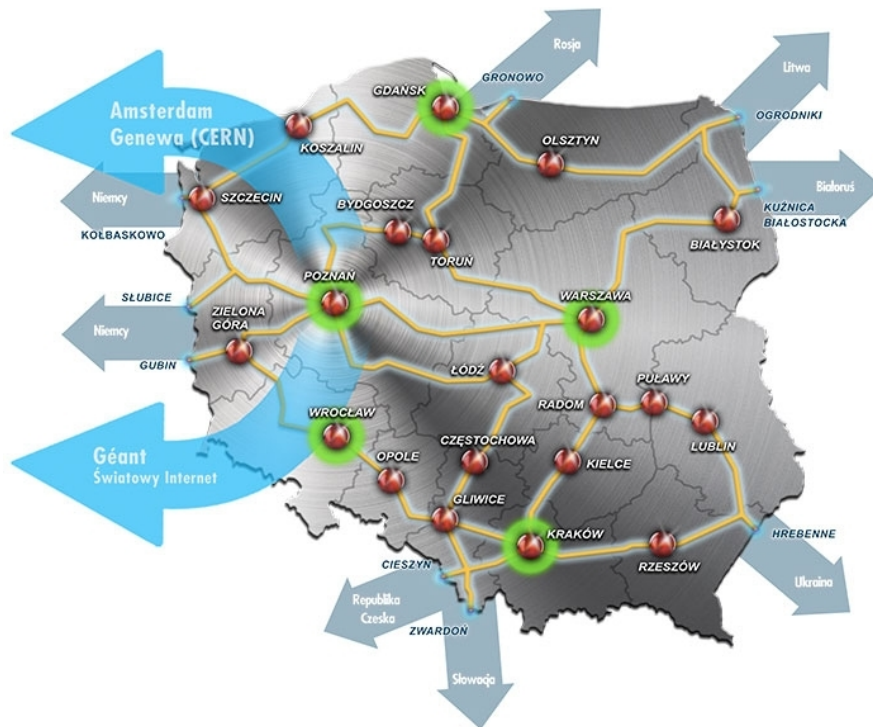


*Figure 3. Overview of the topology of the PIONIER academic fiber optic network. The green points indicate the location of five academic supercomputer centers..*

## 2.5 Visualisation
The newly built lecture auditorium at the ICM Technology Center for 70 students has unique visualization equipment. Video wall with a resolution of 8K and a diagonal of 220" (16 monitors) is fed with image content directly from a dedicated server equipped with 4 GPU cards, 48 CPU cores and 1 TB of operational memory. The software anables of interactive visual work with large data (including VisNow, ICM created visualisation software), as well as remote joint work (including SAGE2) or the transmission of lectures or images from around the world.

# 3. Professional Expertise at ICM

ICM is involved in research, exploration and development of new computational techniques and the latest technological solutions. In 2009-2016, the current ICM director, Dr. Marek Michalewicz, held the position of CEO (Chief Executive Officer) at the A*STAR Computational Resource Center in Singapore. It is the largest scientific computing center in the region of Southeast Asia belonging to the scientific organization A*STAR (Agency for Science, Technology and Research) in Singapore. A*STAR employs over 5,500 employees at 14 scientific institutes. Dr. Michalewicz was responsible for creating the National Supercomputing Center in Singapore (NSCC). NSCC was an investment of 100 million SGD and started its operations at the beginning of 2015.

The unique solutions created and introduced under the leadership of Dr. Marek Michalewicz in Singapore include:

1. Special network architecture for big data transfers: Designing a unique network connection architecture between two A*STAR campuses separated by a distance of about 1.2 km: Biopolis, and Fusionopolis 1 (headquarters of NSCC). Total inter-campus network capacity is unprecedented in the world and amounted to 1.18 Tbps in 2016.
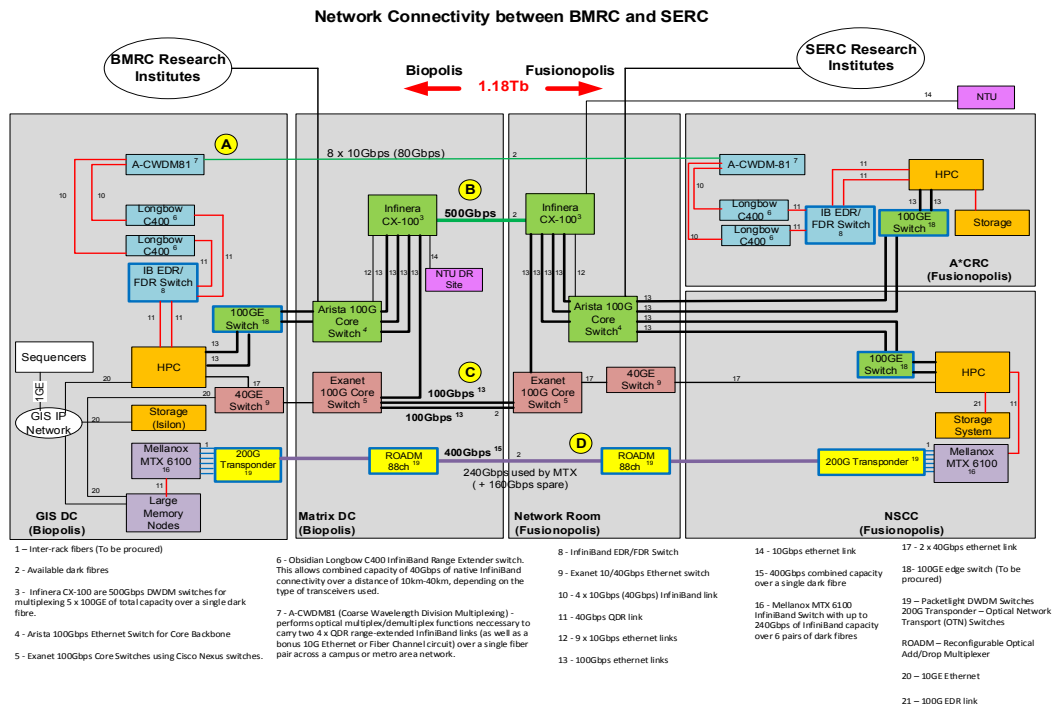


*Figure 4. Network connection between Biopolis and Fusionopolis in Singapore. 1.18 Tbps realized through five connections, including 500 Gbps CloudXpress-1 (Infinera, tcp/ip) and direct InfiniBand connection (MetroX, Mellanox).*

One of the most interesting applications of the connection between Genome Institute of Singapore (GIS) and NSCC was the direct connection of Illumina sequencers with the NSCC storage 1.2 km away using the InfiniBand protocol.
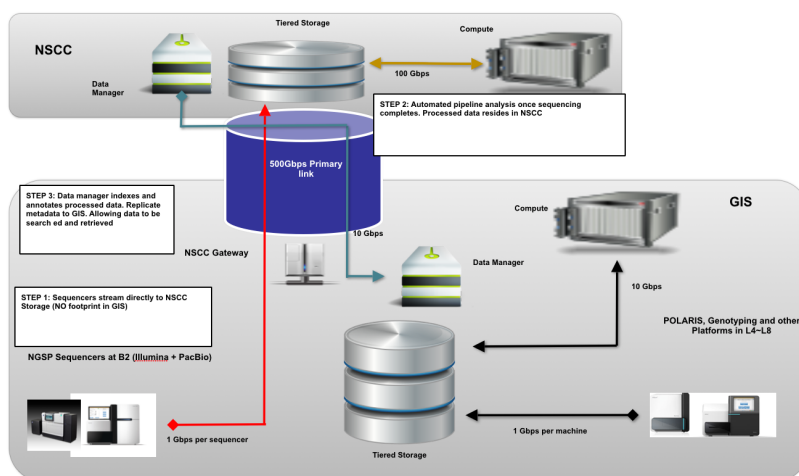


Figure 5. The Illumina sequencers in GIS directly transmit experimental data to the NSCC storage and then for analysis on a supercomputer - without any intermediate steps

1.  InfiniCortex

Between 2014 and 2016, Dr. Michalewicz was the initiator and the leader of the InfiniCortex project - creating a concurrent global computer, located on four continents and in seven countries (Singapore, Australia, Japan, Canada, USA, France, Poland) connected by a global InfiniBand network with bandwidth 100Gbps.
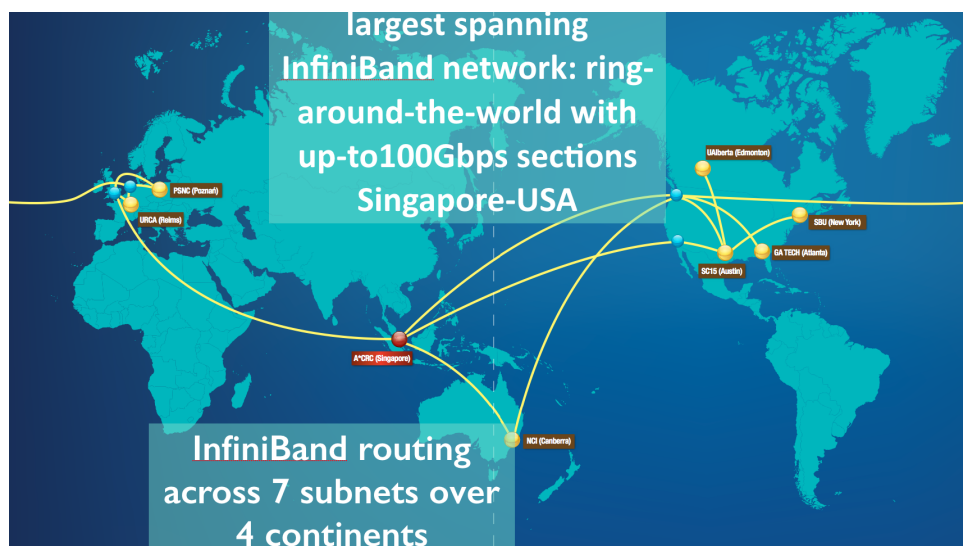


*Figure 6. InfiniCortex in 2015, presented at the Supercomputing 15 conference in New Orleans. Global InfiniBand connection with 100gbps bandwidth on four*

*continents. The creation of an InfiniBand routing and global sub-networks was also shown for the first time.*

This achievement is unprecedented and not repeated to this day. InfiniCortex is not only a unique infrastructure for distributed computing and data transmission worldwide, but also a platform on which many biological and medical applications have been demonstrated:

i.  transmission of genomic data between Australia and Singapore, via Seattle. A 1.14 Tb genomic file was sent using an InfiniCortex connection in 24 minutes. A normal transfer using a standard tcp/ip connection requires 12.5 hour.
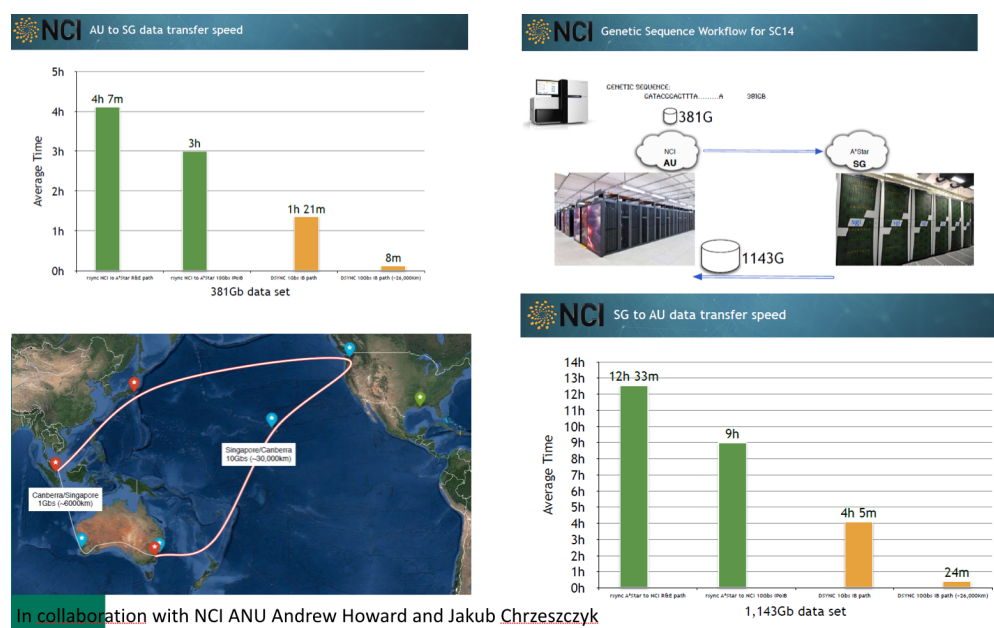


*Figure 7. Transfer of genomic data between National Computing Infrastructure at the Australian National University (NCI ANU) in Canberra and Singapore, using InfiniCortex (InfiniBand protocol at global distances).*

ii. InfiniCloud.
    On the basis of InfiniCortex infrastructure, a global InfiniCloud virtual computer was created, which was configured at the push of a button and was located in four places in the world: Australia, Singapore, France and the USA. InfiniCloud served to demonstrate the use of cancer genomics: the "cancer mutation calling pipeline".
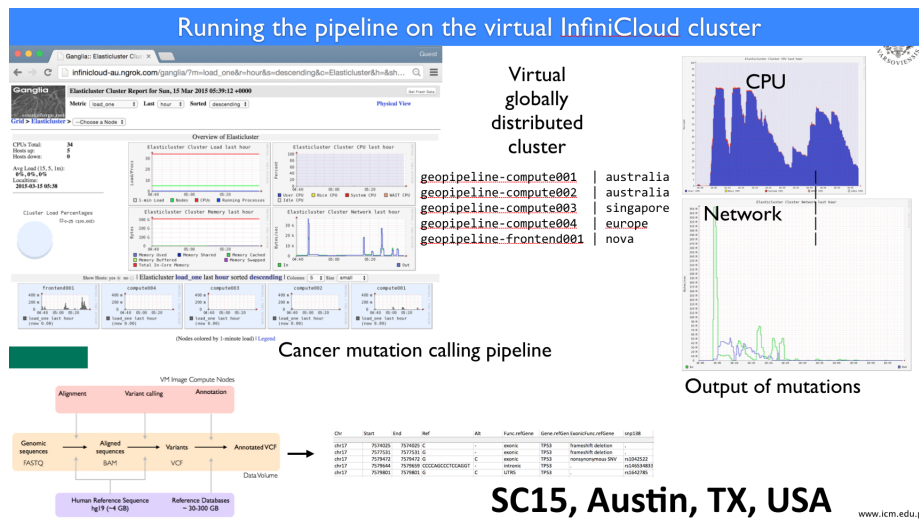
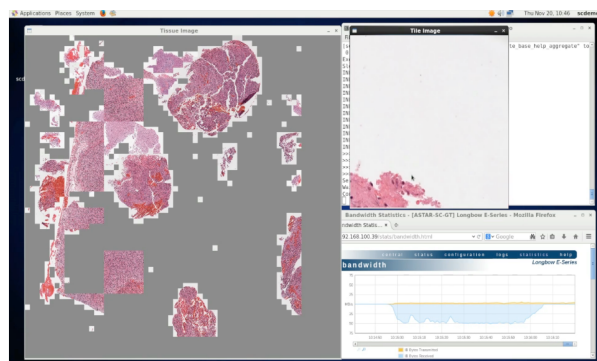*Figure 8. InfiniCloud -a global virtual computer used for calculations in cancer mutation genomics*

iii. On-the-fly segmentation of large pathological tissue images (Streaming Segmentation of Large Pathology Tissue Images). This project was carried out with partners from Stony Brook Univesrity, Oak Ridge National Laboratory and Georgia Institute of Techchnology (GT). Huge pathological tissue images up to 120,000 x 120,000 pixels stored in Singapore were sent in real time using the InfiniCortex connection to GT for analysis. GT results were transferred to a tablet (read: doctors next to the patient) in New Orleans.

## Streaming Segmentation of Large Pathology Tissue Images

**Abstract:**

Demo showed nuclear segmentation on a high resolution whole slide tissue color image (53,000x36,000 pixels, RGB) stored on a cluster at ACRC in Singapore.

The image is partitioned into tiles. The tiles that contain enough tissue data are streamed from multiple cluster nodes via the ORNL ADIOS system over Longbow connections to a cluster at Georgia Tech (GT) in the US. The tiles are processed on the GT cluster using the SBU RT middleware as they are received to segment nuclei. Segmented nuclei and the tiles are assembled into an image.



**Team:**

Tahsin Kurc, SBU and ORNL     Scott Klasky, ORNL     Jong Choi, ORNL     Joel Saltz, SBU

*Figure 9. Tissue pathology images segmentation "on the fly" at trans-continental distances.*

3. The architecture of the Aspire-1 supercomputer and the storage resources at NSCC in Singapore, designed by Dr. Michalewicz and colleagues, was created to meet the diverse users needs in HPC, Big Data and AI. The unique features of this solution are:

i.    Direct InfiniBand connection between GIS and NSCC
ii.   unprecedented fast I/O (Input/Output) between Aspire-1 supercomputer and storage: 500 Gbps using 10x Infinite Memory Engine units burst buffer.
iii.  ten servers with very large shared memory designed specifically for real-time genomics and biology (outside the queue system). One server with 6TB RAM, 4x with 2TB and 5x 1TB.
iv.   Hierarchical management of storage (HSM) and special solutions for distributed meta-data management (Arcitecta).
v.    Special storage partitions for genomic data.
vi.   A separate part of 500TB storage with hardware secure storage, especially for sensitive genomic and medical data for long-term storage.
vii.  Direct login nodes connected by InfiniBand with NSCC in three universities (NUS, NTU, SUTD) and in the future on campuses of many Singapore hospitals and institutes - the solution unique on a global scale.

For several years ICM has been involved in unique technological projects that places us in the league of leaders in the application of network and computing solutions on the global scale. Here are some examples:

1. Connection of two ICM data centers (Ochota - Białołęka) at a distance of approx. 20 km with a throughput of 1.2 Tbps using the latest CloudXpress-2 demonstration equipment from Infinera (ICM as the first to test this technology in Europe, right after Amazon, Facebook and Google conducted their tests).
2. Data Transfer Nodes (DTN) - combined with data transfer and computations on a global scale. A demonstration at Supercomputing 2018 conference in the US of connection of ICM in Warsaw with the Pawsey Center in Perth, Australia and launching containerized programs alternatively either on resources in Australia or in Warsaw.
3. Establishing InfiniBand connections between ICM and TASK Gdańsk (about 900 km light path), and ICM and NCBJ, Świerk (about 40km light path) and building highly distributed concurrent computer system
4. Data transfer between Warsaw and Singapore on the new 100Gbps CEA-1 (Collaboration Asia Europe 1) connection at 100Gbps in cooperation with US based company zettar.

# 4. Educational Activities at ICM

## 4.1 Computational Engineering

Computational engineering is a practical full-time Masters level program launched by the ICM UW in the winter semester 2016/2017. The offer is primarily addressed to graduates of engineering studies or master's degrees holders. The studies last three semesters and lead to a master's degree in Computational Engineering.

During their studies, students gain basic knowledge in the field of large-scale systems, their architecture, as well as their management and use in selected fields. They are trained in the methods of processing, analyzing and visualizing Big Data, including machine learning methods and artificial intelligence. They learn methods and paradigms of programming large-scale systems, with particular emphasis on parallel and distributed programming, including cloud computing. They gain knowledge about the selected field of applications of large-scale calculations, algorithms and computational methods, and also take part in computational projects using supercomputers. Students may elect topics in visualization of multimodal data, including medical data, or lectures on computational chemistry, molecular modeling or bioinformatics. Students have at their disposal a computer laboratory, they use Cray XC40 supercomputer (Okeanos), a computing cluster (Topola) and a large data analytics system (Enigma).

## 4.2 OMICS Data Science: Bioinformatics and Analysis of large-scale biomedical data

The Omics Data Science Course Program - Bioinformatics and Analysis of large-scale biomedical data was developed by the program council and introduced in 2019 as part of the project "Genetically conditioned diseases - education and diagnostics (EDUGEN)" co-financed by the European Union under the European Social Fund, Operational Program Knowledge Education Development. Priority axis: IV. Social innovation and transnational cooperation Actions: 4.3 in cooperation between the ICM UW and the Institute of Mother and Child in Warsaw.

The head of the program is the ICM employee dr Katarzyna Suski-Grabowski. Dr Katarzyna Suski-Grabowski is a graduate of the University of Paris-Sud, Department of molecular genetics, where she defended her doctoral dissertation on DNA replication. In 2003-2010 she conducted research at the Memorial Sloan Kettering Cancer Center in New York and then at the Weatherall Institute of Molecular Medicine Cancer Research Oxford University. In the years 2011-2017 dr Katarzyna Suski-Grabowski managed a

research group at the Institute of Biochemistry and Biophysics of the Polish Academy of Sciences. Currently, dr Katarzyna Suski-Grabowski is the head of the Omics Data Science bioinformatics course led by the Interdisciplinary Center for Mathematical and Computational Modeling at the University of Warsaw and the Flagship 1 coordinator at the European University 4EU +.

The OMICS educational project aims to build a new community of scientists, physicians and medical technicians who will be skilled in new computational techniques and data analysis in modern fields of systemic biology, medicine and genetics. The program prepares participants to tackle specific OMICS issues in the context of medical genetics.

Introductory classes for high-throughput data analysis:
- High-throughput tests in medicine - introduction
- Infrastructure in high-throughput research
- Basic IT tools (Python)
- Basic databases and tools for analyzing high-throughput data
- Basics of analysis using the R package
- The use of Big Data tools in omics analysis
- Deep Learning methods in biomedical research
- Ethical aspects of biomedical high-throughput research

This interdisciplinary program consists of all quantitative OMICS (OMICS data science) and is prepared for students from all fields of science. Students gain knowledge of high-throughput data analysis:
Practical use of high-pass analyzes - practical classes preceded by a theoretical introduction

- Genomics
- Transcriptomics
- Metagenomics/Microbiome
- Epigenomics
- Proteomics
- Metabolomics

## 4.3 *Warsaw Team* world Student Cluster Competitions

In December 2016 ICM staff members created *the Warsaw Team*, a team of Polish students who take part in the global Student Cluster Competition (SCC). Three independent competitions are held annually in China (April-May), Germany (June) and the USA (November). *The Warsaw Team* has already participated in seven world finals and is currently in the 23rd place in the world ranking, and 6th in Europe, the Middle East and Africa (EMEA). It is worth emphasizing in this context that students often have to run and accelerate the most advanced scientific software, often in the field of Big Data and AI. For example, in the Chinese competition ASC 2017 one of the programs was FALCON: the experimental PacBio diploid assembler, at Supercomputing 2017 (SC17) students accelerated the operation of MrBayes

program, Bayesian inference program and model selection for a wide selection of phylogenetic and evolutionary models. At the 2018 ASC competition in China, one of the competition tasks was the RELION program - the "gold standard" created for image analysis at Cryo-EM, and a year later at ASC 2019 students dealt with the Redbean program: A fuzzy Bruijn graph approach to long noisy reads assembly.

## 5. Conferences organised by ICM

Supercomputing Frontiers is a series of conferences initiated and organised in Singapore in 2015 by Dr Michalewicz. The first three editions (2015-2017) were held in Singapore, after which conference series has been moved to Warsaw under the name Supercomputing Frontiers Europe since 2018[1].

## 6. Other relevant developments at ICM

ICM plans to expand its involvement in graph calculations. Late last year we have bought a license for Urica XC software system for Okeanos supercomputer and we are finalising it's installation. This suite of programs is specifically designed for AI, Big Data and graph calculations. We are also in discussions to install a special suite of programs for property or attribute graph computations (from an undisclosed vendor).

For AI problems we also have a specialised NEC Tsubame Aurora vector computer with special SOL software.

## 7. ICM as exemplary Centre of Excellence in Big Data

For the last 27 years ICM has been building expertise and reputation in computational research, data analytics, Machine Learning, Artificial Intelligence, Medical Image Analysis in Medical Diagnosis, Virtual Library of Science, Numerical Weather Prediction, and many other important areas of mathematical and computational research. We have been involved in building and maintanance of substantial datacentres, network solutions of the global scale, Big Data projects and in users services in all above pursuits. ICM is a quintessential Center of Excellence in Big Data, despite not having anything of the kind in its name and despite not being directly and explicitly funded for specific "Big Data" activities.
ICM is looking forward to collaborations with other Polish, European and global players in Big Data, but also in HPC, broadly defined AI, and other pursuits of computer based scientific discovery as well as applied research.

---

[1] https://supercomputingfrontiers.eu