

**I.I. Levin**

## **Intelligent Multiprocessor Systems: Creation and Application Experience**

### **Abstract**

"Super-Computer and Neurocomputer Research Centre" (Taganrog) develops and carries out mass production of a range of intelligent multiprocessor problem-oriented computing systems based on the field-programmable gate arrays (FPGAs). The latest development is the universal reconfigurable computing system (RCS) "Arcturus" designed to solve the problems in various subject areas: mathematical physics, digital signal processing and artificial intelligence technologies.

Currently, the computational complexity of the tasks implying the use of neural network technologies is constantly growing, since the volume of analysed data is constantly increasing, and requirements to the speed and quality of its processing are becoming more stringent. While processing the images, texts or sounds, the neural networks demonstrate high accuracy in case the input data correlates with the training data, otherwise quality of solving a problem significantly reduces. At the same time, the rate of emergence of the new information, unprocessed during the training, can be quite high, therefore there arises the need for constant retraining of neural networks. Thus, the neural networks training under the standard conditions is carried out over the several months, whereas retraining must be completed within the acceptable timeframe (several days or even hours) regardless the volume of data, which is feasible only in case of multiple intensification of the intelligent system performance. At the same time, in the systems based on the graphic processing units (GPUs) the hardware costs scale-up linearly, and although they provide a linear growth of declared peak performance, their machine learning rate demonstrates just a logarithmic growth at the most.

To ensure the required quality of neural network retraining within an acceptable timeframe, regardless the volume of data, the FPGA-based RCS can be used as an alternative solution. In this regard, the reconfigurable computing device (RCD) "Arcturus" is deemed to be the most advanced hardware platform. These RCDs implement a number of breakthrough engineering solutions that make it possible to fulfil the complex tasks in a single computing circuit, and perform calculations without interruptions due to the powerful information exchange system with the necessary level of power supply and cooling.

The RCD "Arcturus 3U19" is designed of 16 calculating blades (CBs) vertically installed on a cross-board, each with six XCVU37P FPGAs by Xilinx, UltraScale+ family. The RCD "Arcturus" has a unique layout density of the computing resource – 96 FPGAs of the high integration degree (more than 270 million logic cells in total) in one device. The calculating blades are located on a cross-board, which provides information and control data transfer between the blades. A unique cross-board ensures transit between the media (liquid-air). The external area of the cross-board includes various interfaces for connecting peripheral devices and 144 optical information channels used for information exchange between the devices.

The RCDs use Xilinx XCVU37P FPGA of UltraScale+ family of "HBM" (High Bandwidth Memory) series, each of them having two DDR memories of total capacity of 8 GB ensuring the multi-channel access with the peak capacity of up to 460 GB/s. The implementation of the present technology will significantly expand the opportunities for rational use of the memory in application programs for solving the various problems.

The RCD “Arcturus” ensures the unique real performance due to the powerful information exchange subsystem, which consists of the multiple communication channels between FPGAs within the cross-board, as well as between FPGAs of the neighbouring cross-boards. 24 differential lines with data transfer rate of up to 24 Gbit/s within a CB and 24 differential lines with data transfer rate of up to 12 Gbit/s between the CBs are implemented between a pair of FPGAs. The total capacity of CM communication channels is 15.6 Tbit/s, including between the CMs – 9 Tbit/s.

When building the computer systems, it is possible to organize the information exchange between the RCDs through the optical channels. The interdevice optical transceivers are installed on the cross-board and provide the capacity of up to 4.5 Tbit/s.

The maximum power consumption of the RCD “Arcturus” is 25 kW. To ensure cooling of the RCD “Arcturus” loaded electronic components, an immersion cooling system (immersion cooling) is used, which provides direct immersion of all the electronic modules of the device in the dielectric coolant with high electrical strength and thermal conductivity and the highest possible heat capacity upon low viscosity.

The RCD “Arcturus” application prospectivity for solving various computational-time-consuming problems within the structural paradigm of calculations has been experimentally confirmed. The structural paradigm allows organization of the pipeline data processing for problem solving, including such problems as machine or deep learning. The pipeline computing allows avoiding the overhead costs on organization of the computational process management, since in addition to direct data processing, it can execute the key procedures of analysing the results and generating a decision on continuing or completion of the work.

The conducted research revealed that when solving the problem of training a neural network image classifier, the performance of a single “Arcturus” calculating blade (CB) corresponds to the performance of the GPU Nvidia Tesla A100, and increase of the number of CBs in the system ensures almost linear increase in the performance of the RCS. Besides, upon equal performance value, the energy efficiency of the RCS created on the basis of the RCD “Arcturus” will be 2 times higher than that of the Nvidia DGX system created on the basis of the Tesla A100.