

# Решение больших СЛАУ на гетерогенных СнК серии «Мультикор»

Е.С. Янакова, Г.Т. Мачарадзе, Н.В. Костулин, А.А. Тюрин

АО НПЦ «ЭЛВИС», г. Москва, helen@elvees.com, gmacharadze@elvees.com, nkostulin@elvees.com, alex.tyurin1420@gmail.com

**Аннотация** — В работе исследуется эффективность решения больших систем линейных алгебраических уравнений (СЛАУ) на российских гетерогенных СнК серии «Мультикор» с использованием шестнадцати ядер ELCore, разработанных АО НПЦ «ЭЛВИС», с целью использования вычислительного устройства для решения задач моделирования сложных систем. Описана архитектура гетерогенной СнК и модель программирования такого класса устройств для эффективной реализации локально-параллельных интенсивно вычислительных алгоритмов. Представлены результаты сравнительного анализа тестирования российских и зарубежных процессоров общего назначения на задаче решения СЛАУ с плотными матрицами.

**Ключевые слова** — решение СЛАУ, ELCore, моделирование сложных систем, гетерогенные СнК, модель программирования.

## I. ВВЕДЕНИЕ

Исследование и моделирование сложных физических систем приводит к математическим моделям в виде систем линейных алгебраических уравнений (СЛАУ). Иногда СЛАУ является промежуточным этапом для решения более сложных задач, однако значительное количество научно-технических задач моделирования нелинейных систем посредством дискретизации или линеаризации сводится к решению СЛАУ. Такие вычислительные задачи широко используются в различных областях науки и промышленности, в том числе в авиационно-космической и атомной отраслях [2]. Решение задач моделирования в указанных отраслях требует вычислительных устройств петафлопсного и эксафлопсного класса. Например, моделирование поведения отдельных элементов летательного аппарата потребляет  $10^{17} - 10^{19}$  FLOP, где FLOP - floating point operation [3], более фундаментальные задачи моделирования - более сотни эксафлоп. Поэтому задача решения СЛАУ большой размерности является актуальной задачей и начальным этапом исследования применения российских гетерогенных СнК на основе процессоров ELCore для моделирования сложных систем.

Прямые или точные методы решения СЛАУ позволяют получить решение за определенное число шагов и используются для уравнений с  $N$  порядка  $10^4$ -

$10^5$ . На основе прямых методов СЛАУ построен тест HPL (High-Performance Linpack) [4], который ориентирован на работу с плотными матрицами в формате плавающей точки двойной точности (fp64) и направлен на получение наивысшей производительности на вычислительных системах. Использование в качестве эталона стандартного теста, на основе которого строится международный рейтинг суперкомпьютеров ТОП-500 [5], позволяет проводить исследование и достоверный сравнительный анализ вычислительных устройств и систем мирового уровня.

В качестве объекта исследования выбран российский процессор серии «Мультикор» [1], состоящий из восьми процессоров общего назначения с архитектурой MIPS и шестнадцати проприетарных ядер ELCore, предназначенных для решения интенсивно вычислительных задач, объединённых высокоскоростной сетью на кристалле NoC (Network on Chip). Основная проблема в проводимом исследовании заключается в эффективности использования вычислительных ядер разной архитектуры для решения СЛАУ, синхронизации вычислительных ядер, балансировке загрузки с учетом пропускной способности сети NoC. Немаловажной задачей является выбор модели программирования такого класса устройств для внедрения в существующие вычислительные экосистемы и использование существующего программного обеспечения (ПО).

Поэтому основная цель данной работы заключается в исследовании эффективности решения больших СЛАУ на российских гетерогенных СнК серии «Мультикор» с использованием шестнадцати ядер ELCore. Использование гетерогенных СнК для решения СЛАУ является относительно новой задачей, поэтому для достижения указанной цели необходимо выполнить ряд дополнительных исследований, связанных с глубоким анализом архитектуры гетерогенной СнК, с анализом методов и алгоритмов многоуровневой параллельной обработки информации, с выбором модели программирования для такого класса систем и со способами синхронизации вычислительных ядер, а также с особенностями реализации. В завершение приведены результаты экспериментальных исследований и сравнительный анализ с процессорами российских и зарубежных компаний.

## II. ФОРМАЛИЗАЦИЯ ЗАДАЧИ РЕШЕНИЯ БОЛЬШИХ СЛАУ НА ГЕТЕРОГЕННЫХ СнК

Задача решения больших СЛАУ с плотными матрицами сводится к оптимизации набора алгебраических функций под выбранную платформу. Методика программной реализации и повышения эффективности состоит из трех этапов:

1) реализация требуемых базовых подпрограмм линейной алгебры (Basic Linear Algebra Subprograms, BLAS);

2) адаптация технологии межпроцессорного обмена MPI (Message Passing Interface) под аппаратную платформу;

3) подбор параметров теста HPL, при которых получаются наилучшие результаты.

Базовые подпрограммы линейной алгебры, необходимые для решения СЛАУ методом Гаусса, с результатами профилирования перечислены в табл. 1.

Таблица 1

*Базовые подпрограммы линейной алгебры, используемые для решения СЛАУ*

Уровень функции	Описание функций BLAS для теста HPL	Доля времени для N = 8192, %
1	DAXPY: выполняет расчет по формуле $y = a*x + y$	Менее 1
	DCOPY: выполняет копирование $x$ в $y$	Менее 1
	DSCAL: выполняет расчет по формуле $x = a*x$	Менее 1
	DSWAP: выполняет - swap $x$ and $y$	Менее 1
	IDAMAX: выполняет - index of max abs value	Менее 1
2	DGEMV: выполняет - matrix vector multiply	Менее 1
	DGER: выполняет расчет по формуле $A = \alpha*x*y' + A$	Менее 1
	DTRSV: выполняет решение треугольных матричных задач	Менее 1
3	DGEMM: выполняет умножение матриц	80
	DTRSM: решение СЛАУ с треугольной матрицей	6

Наиболее ресурсоемким при адаптации и реализации теста HPL на СнК серии «Мультикор» является первый этап, так как на нём должны нивелироваться проблемы, связанные с балансировкой загрузки и синхронизации мелкозернистых локально-параллельных алгоритмов функций BLAS. Мелкозернистые локально-параллельные алгоритмы предполагают параллелизм на уровне входных данных, разбивая задачи на множество мелких однотипных подзадач, которые будут исполняться параллельно на однотипных вычислительных узлах, которым являются процессорные ядра ELcore. При этом обязательным условием является локальность взаимодействия, когда обмен данными происходит только в пределах ограниченного физического и структурного радиуса, независимо от размеров задачи и системы. Таким образом задача реализации библиотечных функций BLAS сводится к следующему классу задач:

- выбор модели программирования гетерогенных СнК для мелкозернистых локально-параллельных алгоритмов;
- способ синхронизации однотипных вычислительных ядер;
- способ масштабирования вычислительных ядер.

После выполнения трех этапов разработанной методики реализации и адаптации алгоритмов задача решения больших СЛАУ с плотными матрицами сводится к экспериментальному исследованию и

сравнительному анализу вычислительных устройств мирового уровня.

## III. АРХИТЕКТУРА ГЕТЕРОГЕННОЙ СнК СЕРИИ «МУЛЬТИКОР» ДЛЯ РЕШЕНИЯ ВЫСОКОПРОИЗВОДИТЕЛЬНЫХ ЗАДАЧ

В качестве вычислительных устройств для решения СЛАУ используется российский СнК серии «Мультикор» с 8 процессорными ядрами общего назначения CPU с MIPS архитектурой и 16 проприетарными IP-ядрами ELcore, которые предназначены для высокоэффективной параллельной векторной обработки [1, 6]. Их функциональность учитывает современные тенденции в данной области и специфику решаемых прикладных задач. Целевыми задачами для ядра ELcore-50 являются высокоинтенсивная сигнальная обработка и приложения мультиспектрального (радио-, инфракрасного, оптический каналы) компьютерного зрения.

Процессорные ядра объединены NoC (network on chip, сеть на кристалле), предназначенной для высокоскоростного обмена информацией внутри СнК с наименьшей латентностью, обеспечивая тем самым высокоскоростной интерфейс к внешней памяти (рис. 1). При использовании NoC каждое ядро соединено с маршрутизатором, через который происходит его общение с другими блоками. Сами маршрутизаторы объединены в сеть, по которой пакеты

данных передаются от одного блока к другому, таким способом снижаются ограничения по масштабированию.

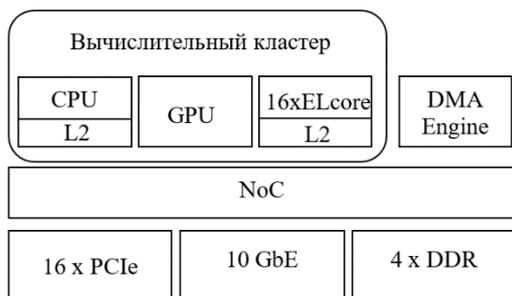


Рис. 1. Схема вычислительного кластера и периферии СпК серии «Мультикор»

Эффективная реализация функций BLAS (табл. 1) достигается благодаря использованию архитектурных особенностей и хорошо продуманной системы инструкций процессорного IP-ядра ELcore. В составе одного ядра объединены два тесно связанных между собой сопроцессора – скалярный ELcore-50S и векторный ELcore-50V с использованием одного реконфигурируемого мультиформатного регистрового файла (рис. 2). Скалярный сопроцессор ELcore-50S представляет собой RISC-ядро и образует скалярный канал обработки данных. Векторный сопроцессор ELcore-50V (или EVX - ELcore Vector eXtension) предназначен для выполнения векторных высоко параллельных вычислений, включая тензорные инструкции умножения матриц и фильтрации. Реконфигурируемый мультиформатный регистровый файл состоит из регистров общего назначения (RF) и регистрового файла для векторных инструкций (VF). Реконфигурируемость обеспечивается за счет возможности использования регистров VF для скалярных инструкций, мультиформатность достигается за счет использования регистров RF и VF в зависимости от инструкции для вычислений в разных форматах. Программы для скалярного и векторного каналов ELcore-50 записываются в виде единого потока инструкций и отлаживаются в рамках единой среды программирования и отладки. Тем самым обеспечивается единство программного кода и автоматическая синхронизация потоков обработки в скалярном и векторном каналах.

Программы для скалярного и векторного каналов процессорного ядра ELcore-50 выполняются одновременно. Это достигается с помощью VLIW-распараллеливания, которое предполагает одновременное выполнение на каждом такте работы процессора нескольких инструкций, которые записываются в программном коде в виде длинного командного слова (VLIW - very long instruction word). Архитектура процессорного ядра ELcore поддерживает одновременное исполнение до четырёх скалярных и до четырёх векторных инструкций – всего до восьми инструкций на каждом такте. Каждая инструкция кодируется 32-разрядным командным словом. Таким образом, размер длинного командного слова,

извлекаемого на каждом такте из программной памяти ядра ELcore-50, кратен 32 и составляет от 32 до 256 разрядов. В каждом 32-разрядном командном слове имеется 1-битовое кодовое поле, указывающее на то, является ли данное командное слово последним во VLIW-пакете. Применение такого VLIW-распараллеливания, с одной стороны, значительно повышает производительность обработки данных, с другой стороны, обеспечивает строгую синхронизацию работы скалярного и векторного каналов.

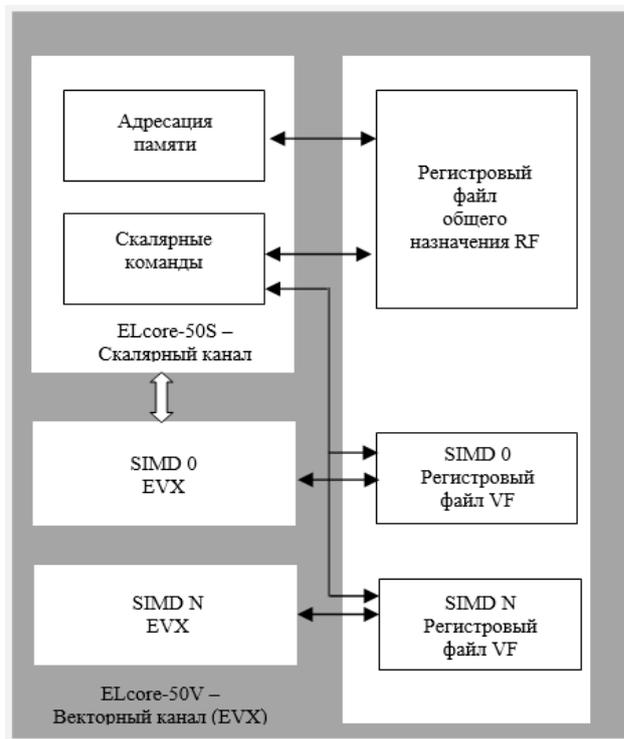


Рис. 2. Структура ядра ELcore-50

Таким образом, используя параллелизм всех типов, можно достичь очень высокой производительности вычислений – пиковая производительность 16 IP-ядер ELcore на частоте 1 ГГц для чисел в формате fp64 равна 1 TFLOPS.

#### IV. ВЫБОР МОДЕЛИ ПРОГРАММИРОВАНИЯ ГЕТЕРОГЕННЫХ СпК

В рамках реализации алгоритмов решения СЛАУ рассмотрены следующие технологии программирования: OpenCL [6], OpenMP [7], MPI [8] и OpenAMP [9]. Основными критериями выбора являются: поддерживаемые типы алгоритма; быстрый запуск существующего программного обеспечения; поддержка механизма управления памятью; поддержка единого виртуального адресного пространства; возможность выбора модели параллелизма; наличие аппаратной поддержки со стороны вычислительных устройств; универсальность модели для различных конфигураций гетерогенной СпК. Сравнительный анализ перечисленных технологий параллельного программирования для гетерогенных СпК показал

необходимость использования трехуровневого параллелизма для решения поставленной задачи:

1. *Параллелизм на уровне СнК*: использование технологии MPI для межпроцессорного взаимодействия.

2. *Параллелизм на уровне гетерогенных ядер*: использование технологии OpenCL, преимущества которой заключаются в поддержке механизмов управления памятью и единого виртуального адресного пространства, а также в возможности выбора типа параллелизма на уровне однотипных процессорных ядер и использования алгоритмов крупноблочного, блочного и мелкозернистого типа. Основные недостатки связаны с необходимостью модификации исходных кодов и наличием накладных расходов, связанных с запуском задач и синхронизации IP ядер ELcore.

3. *Параллелизм на уровне IP ядер ELcore*: использование технологии OpenMP.

Проведен анализ влияния недостатков технологии OpenCL при использовании на гетерогенной СнК, результаты которого представлены в таблице 2. Суммарное время запуска задачи не превышает одной миллисекунды при частоте СнК 1 ГГц, что является незначительным при решении больших СЛАУ.

Для синхронизации IP-ядер ELcore в СнК предусмотрены следующие средства:

- атомарные операции для синхронизации аргументов между IP-ядрами ELcore;
- аппаратный блок спин-блокировок для барьерной синхронизации между IP-ядрами ELcore;
- механизм передачи сообщений MailBox для синхронизации промежуточных этапов алгоритмов;

- использование внешней памяти для синхронизации доступа к большим блокам данных.

Таблица 2

*Результаты профилирования запуска задачи с использованием реализации технологии OpenCL*

Название этапа	Время выполнения этапа, мкс
Подсчёт числа аргументов	2
Настройка аргументов для задания	72
Настройка ELF-секций задания	632
Создание объекта для виртуальной памяти	60
Настройка кэш-памяти	< 1
Установка регистров	1
Копирование программы	< 1
Дополнительные расходы	15
Итого:	783

Сравнительный анализ средств синхронизации (см. табл. 3) показал, что наиболее эффективным способом для синхронизации ядер ELcore является совместное использование атомарных операций и блока аппаратной поддержки спин-блокировки. Под высокой скоростью работы понимается, что время синхронизации не превышает десятки микросекунд.

Таблица 3

*Сравнительный анализ средств синхронизации*

Способ синхронизации	Преимущества	Недостатки
Атомарные операции	Высокая скорость работы	Ограниченный набор операций
Спин-блокировка	Относительно высокая скорость работы. Неограниченный функционал	Существует возможность взаимной блокировки
Механизм передачи сообщений	Возможность настройки прерываний	Низкая скорость из-за накладных расходов
Внешняя память		Крайне низкая скорость за счет применения программных решений

Таким образом, гетерогенным СнК при использовании гибридных вычислений свойственен трехуровневый параллелизм с применением известных стандартных технологий: MPI - на уровне СнК, OpenCL - на уровне гетерогенных ядер и OpenMP – на уровне гомогенных ядер. Такой подход позволяет внедрить достаточно сложную СнК в существующую экосистему вычислительных устройств и использовать существующее открытое программное обеспечение.

Модификация программного кода необходима в случае отсутствия его адаптации под технологии параллельного и распределенного алгоритма.

V. ОСОБЕННОСТИ РЕАЛИЗАЦИИ И РЕЗУЛЬТАТЫ ЭКСПЕРИМЕНТАЛЬНЫХ ИССЛЕДОВАНИЙ

Для решения СЛАУ реализована библиотека FastBLAS, содержащая необходимые

оптимизированные функции (табл. 1) для гетерогенных СнК серии «Мультикор». Внутренняя структура реализации каждой функции представлена на рис. 3. Вызов функции BLAS происходит из приложения, запущенного на CPU. Каждая функция определяет число свободных или разрешенных для использования ядер ELcore и запускает задачу непосредственно на доступном числе ядер. После окончания расчетов управление передается основной программе на CPU.

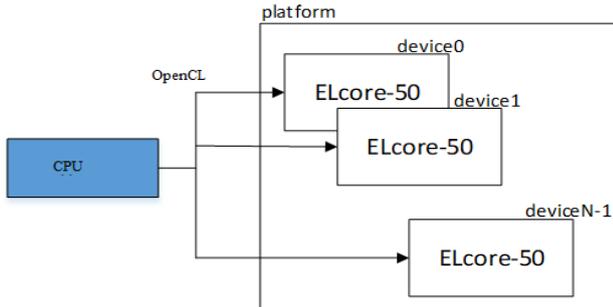


Рис. 3. Структура запуска функций FastBLAS с использованием технологии OpenCL

Время работы функций BLAS уровня 3 на разном числе ядер ELcore представлено в таблицах 4 и 5. Следует отметить, что чем выше N, тем выше процент загрузки процессорных ядер. При N = 5120 для двух ядер он приближается к 90%. При 16 ядрах значительные потери наблюдаются из-за конфликтов доступа к внешней памяти.

Результаты теста HPL при решении СЛАУ показали высокую эффективность при больших N из-за значительной временной доли функции DGEMM (вторая строчка таблицы 6), которая составляет 80% при N = 8192 в неоптимизированной версии функции.

Таким образом, полученные результаты производительности на задаче решения СЛАУ на плотных матрицах методом Гаусса соответствуют мировому уровню (таблица 7). Используемые технологии параллельного и распределенного программирования позволяют эффективно использовать аппаратные возможности гетерогенных СнК серии «Мультикор».

Таблица 4

Производительность функции DGEMM уровня 3 библиотеки FastBLAS на 16 ядрах ELcore

N	512	1024	2048	4096	5120
1 ядро					
Время, мс	5	38	304	2422	4721
Загрузка, %	87	88	88	89	89
2 ядра					
Время, мс	2	19	152	1212	2361
Загрузка, %	87	88	88	89	89
16 ядер					
Время, мс	<1	3	22	178	347
Загрузка, %	73	74	74	75	75

Таблица 5

Производительность функции DTRSM уровня 3 библиотеки FastBLAS на 16 ядрах ELcore

N	512	1024	2048	4096	5120
1 ядро					
Время, мс	3	20	153	1214	2414
Загрузка, %	80	85	87	88	87
2 ядра					
Время, мс	1	10	77	616	1208
Загрузка, %	76	84	87	87	87
16 ядер					
Время, мс	<1	2	16	128	336
Загрузка, %	46	50	52	52	52

Таблица 6

Производительность СнК с 16 ядрами ELcore на решении СЛАУ с плотными матрицами (тест HPL)

N	1024	2048	4096	8192
Доля времени функции DGEMM, %	11	37	60	80
Доля времени функции DTRSM, %	4	6	6	5
Результаты теста HPL, GFLOPS	120	240	410	600
Загрузка СнК, %	12	24	41	60

Таблица 7

Сравнительный анализ производительности российских процессоров и процессора Intel 8180 на тесте HPL

Наименование процессора	Solaris	Эльбрус -8СВ	Xeon 8180
Фирма-производитель	ЭЛВИС	МЦСТ [10]	Intel [11]
Пик. производительность SP/DP, вект. инстр., GFLOPS	512/256	576/288	>3 500/ >1700
Пик. производительность HP/SP/DP, матричн. инстр., TFLOPS	16/4/1	-	-
Результаты теста HPL, GFLOPS	600	244	1350

## VI. Выводы

В работе рассмотрена задача использования гетерогенных СнК серии «Мультикор» для решения СЛАУ большой размерности. Исследования показали, что архитектура СнК с использованием ядер ELcore позволяет эффективно выполнять высокопараллельные

вычисления, включая тензорные операции матричного умножения и фильтрации.

Новыми результатами проведенных исследований являются способ управления вычислительным процессом с использованием технологии MPI, OpenCL и OpenMP, а также модель трехуровневого параллелизма для решения СЛАУ.

Полученные численные результаты по производительности на задачах решения СЛАУ с плотными матрицами (табл. 7) показывают высокие характеристики гетерогенных SnK серии «Мультикор» с шестнадцатью ядрами ELcore.

#### ЛИТЕРАТУРА

- [1] URL: <https://www.multicore.ru> (дата обращения: 25.03.2020)
- [2] URL: <http://www.logos.vniief.ru> (дата обращения: 25.03.2020)
- [3] А.В. Горобец [и др.] Производительность процессора Эльбрус-8С в суперкомпьютерных приложениях

вычислительной газовой динамики // Препринты ИПМ им. М.В.Келдыша. 2018. № 152. 20 с.

- [4] URL: <https://www.netlib.org/benchmark/hpl> (дата обращения: 25.03.2020)
- [5] URL: <https://www.top500.org> (дата обращения: 25.03.2020)
- [6] Elena Yanakova, Andrey Belyaev, Georgij Macharadze. Efficient software and hardware platform «MULTICORE» for cloud video analytics IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering 2020.
- [7] URL: <https://www.khronos.org/opencl> (дата обращения: 25.03.2020)
- [8] URL: <https://www.openmp.org> (дата обращения: 25.03.2020)
- [9] URL: <https://www.open-mpi.org> (дата обращения: 25.03.2020)
- [10] URL: [https://www.multicore-association.org/workgroup/OAMP\\_Brochure.pdf](https://www.multicore-association.org/workgroup/OAMP_Brochure.pdf) (дата обращения: 25.03.2020)
- [11] URL: <http://mcst.ru/> (дата обращения: 25.03.2020)
- [12] URL: <https://www.intel.ru> (дата обращения: 25.03.2020)

## Solving Large SLE on Heterogeneous SoC of «Multicore» Series

E.S.Yanakova, G.T.Macharadze, N.V.Kostulin, A.A.Tiurin

JSC R&D Center «ELVEES», Moscow, [helen@elvees.com](mailto:helen@elvees.com), [gmacharadze@elvees.com](mailto:gmacharadze@elvees.com),  
[nkostulin@elvees.com](mailto:nkostulin@elvees.com), [alex.tyurin1420@gmail.com](mailto:alex.tyurin1420@gmail.com)

**Abstract** — A significant number of scientific and technical problems of nonlinear system simulation through discretization or linearization are reduced to solving systems of linear algebraic equations (SLE). Solving simulation problems in these industries requires computing devices of the petaflops and exaflops classes, so the problem of solving large-dimensional SLE is an actual problem. In this paper we research the efficiency of solving large SLE on Russian heterogeneous SoC of «Multicore» series using sixteen ELcore cores developed by JSC ELVEES research center [1] in order to use a computing device for solving problems of complex systems simulation. The architecture of heterogeneous SoC of «Multicore» series is described, which allows to perform vector highly parallel computations. In the double float format at a core clock of 1 GHz the peak performance is one TFLOPS. Based on the result of a comparative analysis of programming technologies OpenCL [6], OpenMP [7], MPI [8] and OpenCL [9], as well as synchronization tools, we proposed a three-level programming model for the effective implementation of locally parallel intensive computing algorithms. This programming model includes MPI for SoC level parallelism, OpenCL for heterogeneous cores level parallelism and OpenMP for IP cores ELcore level parallelism. The results of comparative analysis of Russian and foreign general-purpose processors testing, such as Elbrus-8SV [10] and Intel Xeon 8180 [11], to the problem of solving SLE with dense matrices, are presented. According to the results, the performance of the Russian SoC series "Multicore" for the problem of solving SLE slows on dense matrices corresponds to the world level.

**Keywords** — solving SLE, ELcore IP core, complex systems, simulation, heterogeneous SoC, programming model.

#### REFERENCES

- [1] URL: <https://www.multicore.ru> (access date: 25.03.2020)
- [2] URL: <http://www.logos.vniief.ru> (access date: 25.03.2020)
- [3] A.V.Gorobec [i dr.] Proizvoditel'nost' processora El'brus-8S v superkomp'yuternyh prilozheniyah vychislitel'noj gazovoj dinamiki (Elbrus-8S performance in supercomputer applications of computational gas dynamics) // Preprinty IPM im. M.V.Keldysha. 2018. № 152. 20 s.
- [4] URL: <https://www.netlib.org/benchmark/hpl> (access date: 25.03.2020)
- [5] URL: <https://www.top500.org> (access date: 25.03.2020)
- [6] Elena Yanakova, Andrey Belyaev, Georgij Macharadze. Efficient software and hardware platform «MULTICORE» for cloud video analytics IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering 2020.
- [7] URL: <https://www.khronos.org/opencl> (access date: 25.03.2020)
- [8] URL: <https://www.openmp.org> (access date: 25.03.2020)
- [9] URL: <https://www.open-mpi.org> (access date: 25.03.2020)
- [10] URL: [https://www.multicore-association.org/workgroup/OAMP\\_Brochure.pdf](https://www.multicore-association.org/workgroup/OAMP_Brochure.pdf) (access date: 25.03.2020)
- [11] URL: <http://mcst.ru/> (access date: 25.03.2020)
- [12] URL: <https://www.intel.ru> (access date: 25.03.2020)