

Языки высокого уровня для гибридных высокопроизводительных систем с ускорителями

**В.А. Бахтин, В.А. Крюков,
А.С. Колганов, Н.В. Поддерюгина, М.Н. Притула**

bakhtin@keldysh.ru

**Институт прикладной математики
им. М.В. Келдыша РАН**



План доклада

- ▶ DVM-модель параллельного программирования
- ▶ Принципы расширения DVM-модели
- ▶ Основные возможности DVMH (DVM for Heterogeneous systems)
- ▶ Распараллеливание тестов NAS NPB и реальных приложений
- ▶ Планы развития DVM-системы



DVM–модель параллельного программирования

- ▶ Объединяет достоинства модели параллелизма по данным и модели параллелизма по управлению (1993 г.)
- ▶ Базирующаяся на этих языках система разработки параллельных программ (DVM) создана в ИПМ им. М.В. Келдыша РАН
- ▶ Аббревиатура DVM (Distributed Virtual Memory, Distributed Virtual Machine) отражает поддержку виртуальной общей памяти на распределенных системах



Средства программирования

C-DVM = Язык Си + специальные прагмы

Fortran-DVM = Язык Фортран 95 + специальные комментарии

- ▶ Специальные комментарии и прагмы являются высокоуровневыми спецификациями параллелизма в терминах последовательной программы
- ▶ Отсутствуют низкоуровневые передачи данных и синхронизации
- ▶ Последовательный стиль программирования
- ▶ Спецификации параллелизма «невидимы» для стандартных компиляторов
- ▶ Существует только один экземпляр программы для последовательного и параллельного счета



Состав DVM-системы

DVM-система состоит из следующих компонент:

- ▶ Компилятор Fortran-DVM(H)
- ▶ Компилятор C-DVM(H)
- ▶ Библиотека поддержки LIB-DVM(H)
- ▶ DVM-отладчик
- ▶ Предсказатель производительности DVM-программ
- ▶ Анализатор производительности DVM-программ



Спецификации параллельного выполнения программы

- ▶ Распределение элементов массива между процессорами
- ▶ Распределение витков цикла между процессорами
- ▶ Спецификация параллельно выполняющихся секций программы (параллельных задач) и отображение их на процессоры
- ▶ Организация эффективного доступа к удаленным (расположенным на других процессорах) данным
- ▶ Организация эффективного выполнения редукционных операций - глобальных операций с расположенными на различных процессорах данными (таких, как их суммирование или нахождение их максимального или минимального значения)

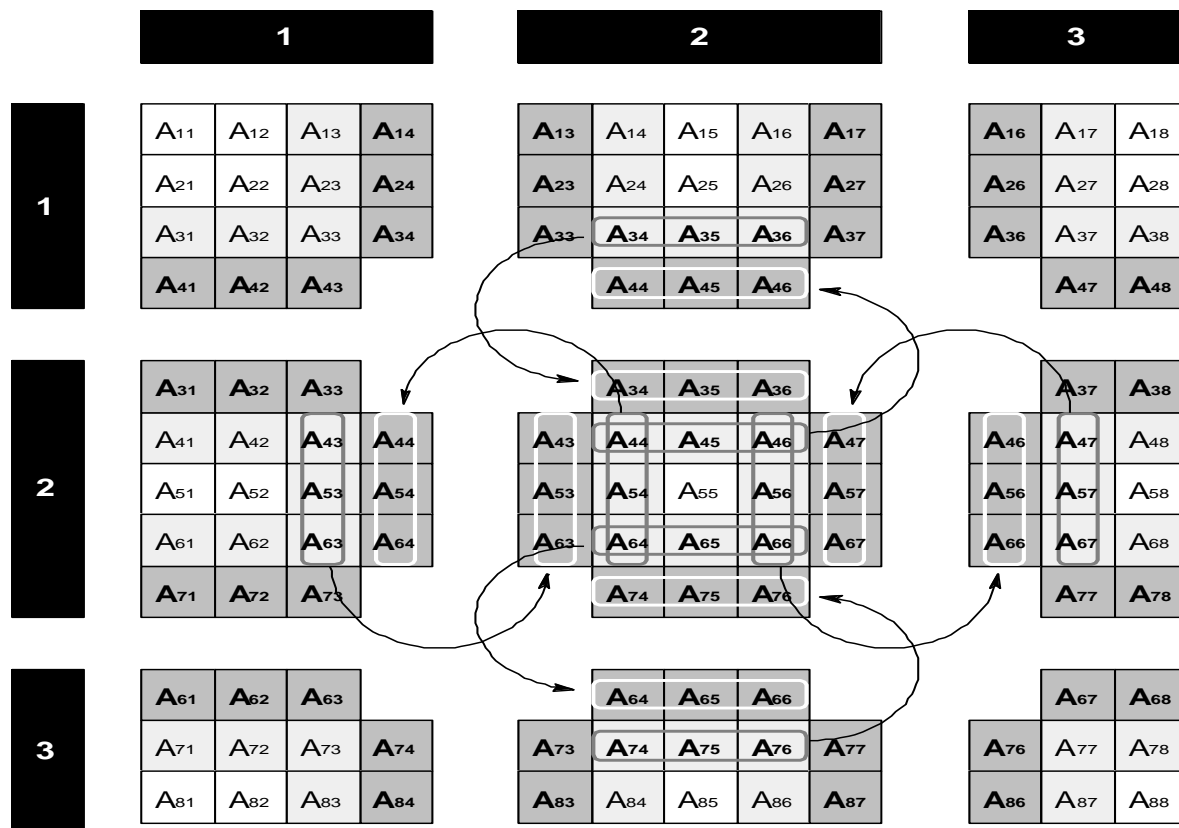


Алгоритм Якоби на языке Fortran

```
PROGRAM JACOB_SEQ
PARAMETER (L=4096, ITMAX=100)
REAL A(L,L), B(L,L)
PRINT *, '***** TEST_JACOBI *****'
DO IT = 1, ITMAX
    DO J = 2, L-1
        DO I = 2, L-1
            A(I, J) = B(I, J)
        ENDDO
    ENDDO
    DO J = 2, L-1
        DO I = 2, L-1
            B(I, J) = (A(I-1, J) + A(I, J-1) + A(I+1, J) +
*                      A(I, J+1)) / 4
        ENDDO
    ENDDO
ENDDO
PRINT *, B
END
```



Распределение массива



```

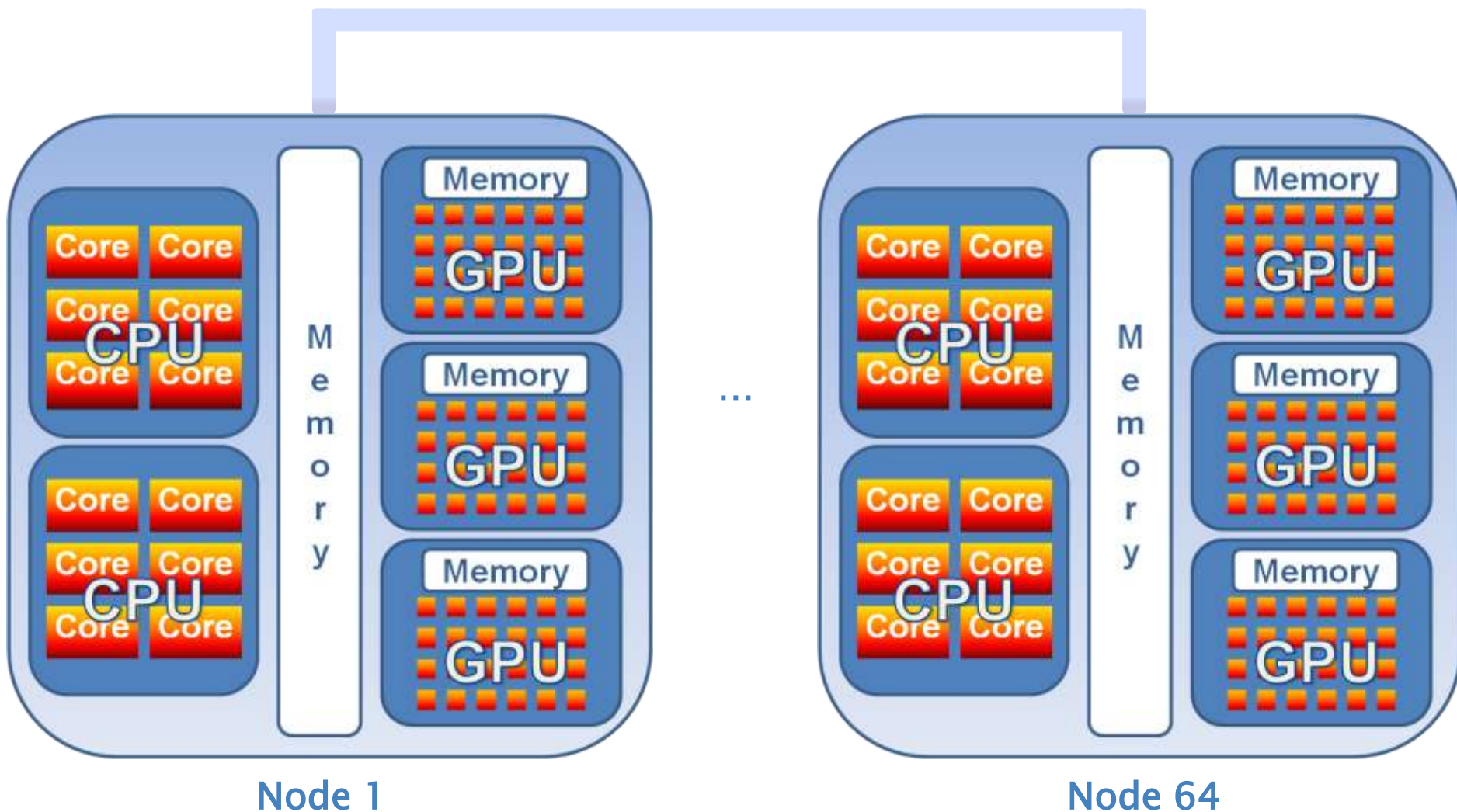
PROGRAM  JACOBY_DVM
PARAMETER  (L=4096, ITMAX=100)
REAL  A(L,L), B(L,L)
!DVM$  DISTRIBUTE  ( BLOCK, BLOCK) :: A
!DVM$  ALIGN B(I,J) WITH A(I,J)
PRINT *, '***** TEST_JACOBI *****'
DO IT = 1, ITMAX
!DVM$      PARALLEL (J,I) ON A(I, J)
            DO J = 2, L-1
                DO I = 2, L-1
                    A(I, J) = B(I, J)
                ENDDO
            ENDDO
!DVM$      PARALLEL (J,I) ON B(I, J), SHADOW_RENEW (A)
            DO J = 2, L-1
                DO I = 2, L-1
                    B(I, J) = (A(I-1, J) + A(I, J-1) + A(I+1, J) + A(I, J+1)) / 4
                ENDDO
            ENDDO
        ENDDO
PRINT *,B
END

```

Алгоритм Якоби
в модели DVM



Гибридная вычислительная система К-100



Принципы расширения DVM-модели

- ▶ Поддержка ускорителей различной архитектуры (GPU, Intel Xeon Phi(MIC)) в узлах кластера
- ▶ Гибкость в управлении распределением вычислений внутри узла кластера(между ускорителями и ядрами центрального процессора)
- ▶ (Полу-)автоматическое управление перемещением данных между оперативной памятью универсального процессора и памятьми ускорителей

Проблема перемещения данных

▶ Два подхода

- Ручное копирование
- Указание входных и выходных данных

▶ Недостатки ручного копирования

- Ориентированность на конкретный набор используемых вычислительных устройств

```
#pragma omp target device(acc0) map(A,B)
```

```
#pragma omp parallel for
```

```
for (i=0;i<N;i++) // OpenMP 4.0
```

```
    A[i] += A[i]*B[i];
```

▶ Проблема оптимизации копирований для разветвленных программ

- Излишние перемещения
- Сложно находимые ошибки



Основные возможности DVMH

- ▶ Определение *вычислительных регионов* (или просто регионов) - фрагментов программы, которые следует выполнять на том или ином ускорителе

!DVM\$ REGION [clause {, clause}]

<region inner>

!DVM\$ END REGION

Где <region inner>:

- ▶ Параллельный DVM-цикл
- ▶ Последовательная группа операторов
- ▶ Контрольная (хост) секция



Основные возможности DVMH

- Уточнение требуемых регионов данных и вида их использования (входные, выходные, локальные):

IN(subarray_or_scalar {, subarray_or_scalar})

OUT(subarray_or_scalar {, subarray_or_scalar})

INOUT(subarray_or_scalar {, subarray_or_scalar})

LOCAL(subarray_or_scalar {, subarray_or_scalar})

INLOCAL(subarray_or_scalar{,subarray_or_scalar})



Основные возможности DVMH

- ▶ Управление перемещением данных между оперативной памятью ЦПУ и памятью ускорителей при выполнении на ЦПУ фрагмента программы, не включенного в какой-либо регион

GET_ACTUAL[(subarray_or_scalar{,subarray_or_scalar})]

делает все необходимые обновления для того, чтобы на хост-памяти были самые новые данные в указанном подмассиве или скаляре.

ACTUAL[(subarray_or_scalar {, subarray_or_scalar})]

объявляет тот факт, что указанный подмассив или скаляр самую новую версию имеет в хост-памяти. При этом пересекающиеся части всех других представителей указанных переменных автоматически устаревают и перед использованием будут (по необходимости) обновлены.



```

PROGRAM  JACOBY_DVMH
PARAMETER  (L=4096, ITMAX=100)
REAL  A(L,L), B(L,L)
!DVM$  DISTRIBUTE  ( BLOCK,  BLOCK)  ::  A
!DVM$  ALIGN  B(I,J)  WITH  A(I,J)
PRINT *, '***** TEST_JACOBI *****'
DO IT = 1, ITMAX

!DVM$      REGION INOUT(A,B)
!DVM$      PARALLEL (J,I) ON A(I, J)
          DO J = 2, L-1
            DO I = 2, L-1
              A(I, J) = B(I, J)
            ENDDO
          ENDDO

!DVM$      PARALLEL (J,I) ON B(I, J), SHADOW_RENEW (A)
          DO J = 2, L-1
            DO I = 2, L-1
              B(I, J) = (A(I-1, J) + A(I, J-1) + A(I+1, J) + A(I, J+1)) / 4
            ENDDO
          ENDDO

!DVM$      END REGION
          ENDDO

!DVM$  GET_ACTUAL(B)
PRINT *,B
END

```

Алгоритм Якоби
в модели DVMH

```

#define L 4096
#define ITMAX 100
#pragma dvm array distribute[block][block] shadow[1:1][1:1]
double A[L][L];
#pragma dvm array (align([i][j] with A[i][j])
double B[L][L];
int main(int argc, char *argv[]) {
    for(int it = 0; it < ITMAX; it++) {
        #pragma dvm region
        {
            #pragma dvm parallel([i][j] on A[i][j])
            for (int i = 1; i < L - 1; i++)
                for (int j = 1; j < L-1; j++) A[i][j] = B[i][j];
            #pragma dvm parallel([i][j] on B[i][j]) shadow_renew(A)
            for (int i = 1; i < L - 1; i++)
                for (int j = 1; j < L - 1; j++)
                    B[i][j] = (A[i - 1][j] + A[i + 1][j] + A[i][j - 1] + A[i][j + 1]) / 4.;
        }
    }
    return 0;
}

```

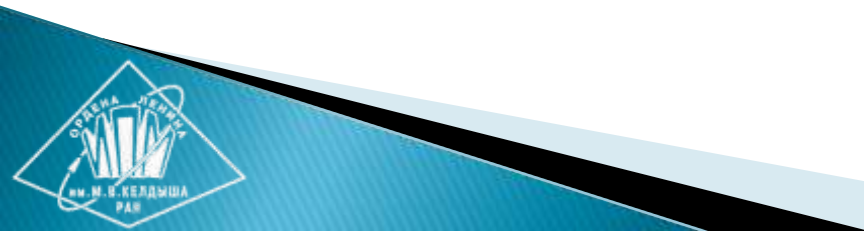
Алгоритм Якоби
в модели DVMH

Основные достоинства DVMH

- ▶ В качестве целевой архитектуры рассматривается кластер с гетерогенными узлами, а не отдельные узлы
- ▶ Перемещение информации между памятью ЦПУ и памятью ускорителей производится, в основном, не по директивам в программе, а автоматически в соответствии со спецификациями использования данных в регионах.

Это позволяет:

- ▶ динамически решать, где выгоднее выполнять тот или иной регион
- ▶ многократно выполнять регион для нахождения оптимального отображения вычислений на ГПУ
- ▶ сравнивать результаты выполнения региона на ЦПУ и ГПУ с целью обнаружения расхождений в результатах выполнения



Распараллеливание тестов NAS

- ▶ **EP** - генерация пар случайных чисел Гаусса
- ▶ **MG** - приближенное решение трехмерного уравнения Пуассона. Метод MultiGrid
- ▶ **BT** - 3D Навье-Стокс, блочная трехдиагональная схема. Метод переменных направлений
- ▶ **LU** - 3D Навье-Стокс. Метод последовательной верхней релаксации
- ▶ **SP** - 3D Навье-Стокс. Скалярная пятидиагональная схема. Beam-Warning approximate factorization



Тесты NAS для GPU

- ▶ OpenCL-версии

Center for Manycore Programming at
Seoul National University (SNU NPB Suite)

- ▶ CUDA-версии

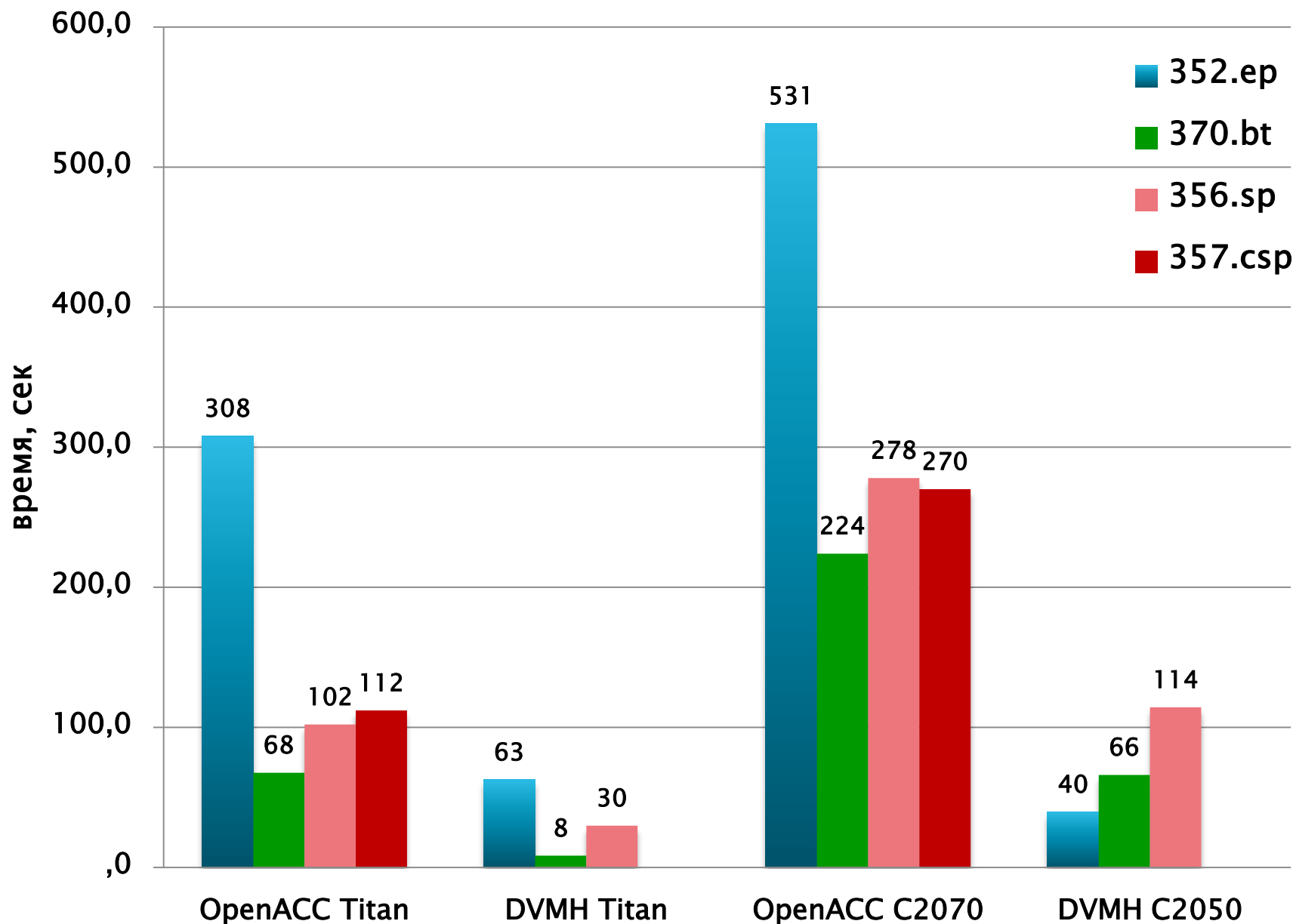
Chemnitz University of Technology (BT, LU, SP)
Laercio Lima Pilla from the Federal University of
Rio Grande do Sul (EP, CG, FT)

- ▶ OpenACC-версии

SPEC ACCEL V1.0 (BT, EP, CG, SP)



Времена выполнения OpenACC и DVMH-версий тестов NAS NPB



http://www.spec.org/accel/results/accel_acc.html

Оптимизация работы с памятью

- ▶ Вынесение часто используемых элементов массива в теле цикла в скалярные переменные
- ▶ Сокращение операций чтения из глобальной памяти GPU за счет избыточных вычислений
- ▶ Динамическое переупорядочивание массивов в памяти GPU



!DVM\$ PARALLEL (K,J,I) ON U(I,J,K,*), PRIVATE(M1,M2,M)

DO K = 2,NZ-1

DO J = 2,NY-1

DO I = 2,NX-1

M1 = 2

M2 = 3

DO M = 1,5

U(I,J,K,M) = U(I,J,K,M1) + U(I,J,K,M2)

ENDDO

DO M = 1,5

U(I,J,K,M) = U(I,J,K,M1+1) + U(I,J,K,M2+1)

ENDDO

DO M = 1,5

U(I,J,K,M) = U(I,J,K,M1-1) + U(I,J,K,M2-1)

ENDDO

ENDDO

ENDDO

ENDDO

Фрагмент программы LU



!DVM\$ PARALLEL (K,J,I) ON U(I,J,K,*), PRIVATE(M1,M2,M,U_)

DO K = 2,NZ-1

DO J = 2,NY-1

DO I = 2,NX-1

DO M=1,5

U_(M) = U(I,J,K,M)

ENDDO

M1 = 2; M2 = 3

DO M = 1,5

U_(M) = U_(M1) + U_(M2)

ENDDO

DO M = 1,5

U_(M) = U_(M1+1) + U_(M2+1)

ENDDO

DO M = 1,5

U_(M) = U_(M1-1) + U_(M2-1)

ENDDO

DO M=1,5

U(I,J,K,M) = U_(M)

ENDDO

ENDDO

ENDDO

ENDDO

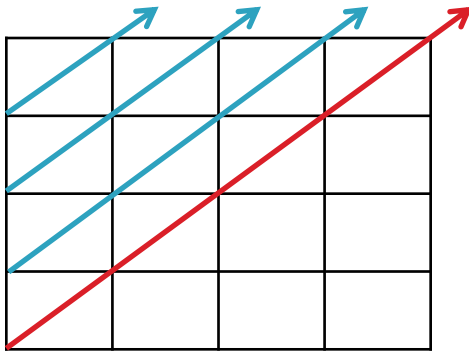
Использование приватных переменных для часто используемых элементов массивов, которые отображаются компилятором NVCC на регистры.

Алгоритм SOR на языке Fortran

```
PROGRAM SOR_DVMH
PARAMETER (L=1000, ITMAX=20, W = 0.5)
REAL A(L,L) , EPS, S
!DVM$ DISTRIBUTE A(BLOCK,BLOCK)
PRINT *, '***** TEST_SOR *****'
DO IT = 1, ITMAX
  EPS = 0.
!DVM$ ACTUAL(EPS)
!DVM$ REGION
!DVM$ PARALLEL(J, I) ON A(I, J), ACROSS(A(1:1,1:1)),
!DVM$& REDUCTION(MAX(EPS)), PRIVATE(S)
  DO J = 2, L-1
    DO I = 2, L-1
      S = A(I, J)
      A(I, J) = (W / 4) * (A(I-1, J) + A(I+1, J) + A(I, J-1) +
&      A(I, J+1)) + ( 1-W ) * A( I, J)
      EPS = MAX ( EPS, ABS( S - A( I, J )))
    ENDDO
  ENDDO
!DVM$ END_REGION
!DVM$ GET_ACTUAL(EPS)
PRINT 200, IT, EPS
200 FORMAT(' IT = ',I4, ' EPS = ', E14.7)
ENDDO
END
```

Выполнение гиперплоскостями

REAL A(4, 4)



**!\$DVM PARALLEL (K,I) ON A(I,K),
ACROSS(A(1:1, 1:1))**

DO K = 1,4

DO I = 1,4

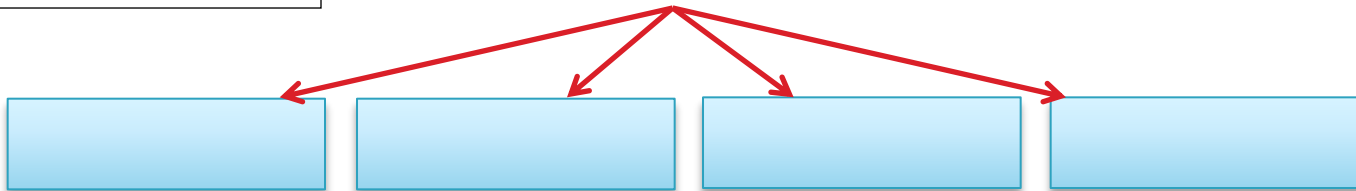
A(I, K) = A(I+1,K) + A(I-1,K) + A(I,K+1)+ ...

ENDDO

ENDDO

4 диагональ

warp load/store operations



Array A(4, 4)



Динамическое переупорядочивание массивов

1	5	9	13
2	6	10	14
3	7	11	15
4	8	12	16

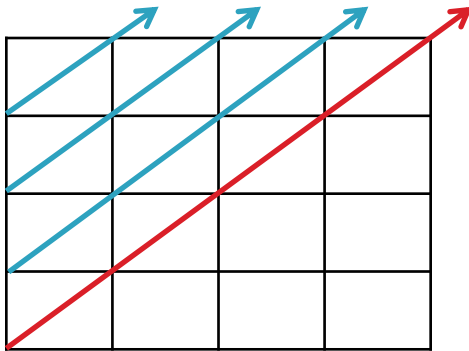


1	6	10	14
2	9	13	12
5	4	8	15
3	7	11	16

Выполняется поддиагональная трансформация матрицы - соседние элементы на диагоналях располагаются в соседних ячейках памяти

Выполнение гиперплоскостями (-autoTfm)

REAL A(4, 4)



**!\$DVM PARALLEL (K,I) ON A(I,K),
ACROSS(A(1:1, 1:1))**

DO K = 1,4

DO I = 1,4

A(I, K) = A(I+1,K) + A(I-1,K) + A(I,K+1)+ ...

ENDDO

ENDDO

4 диагональ

warp load/store operations **with DVMH DR** (-autoTfm)



Array A(4, 4)

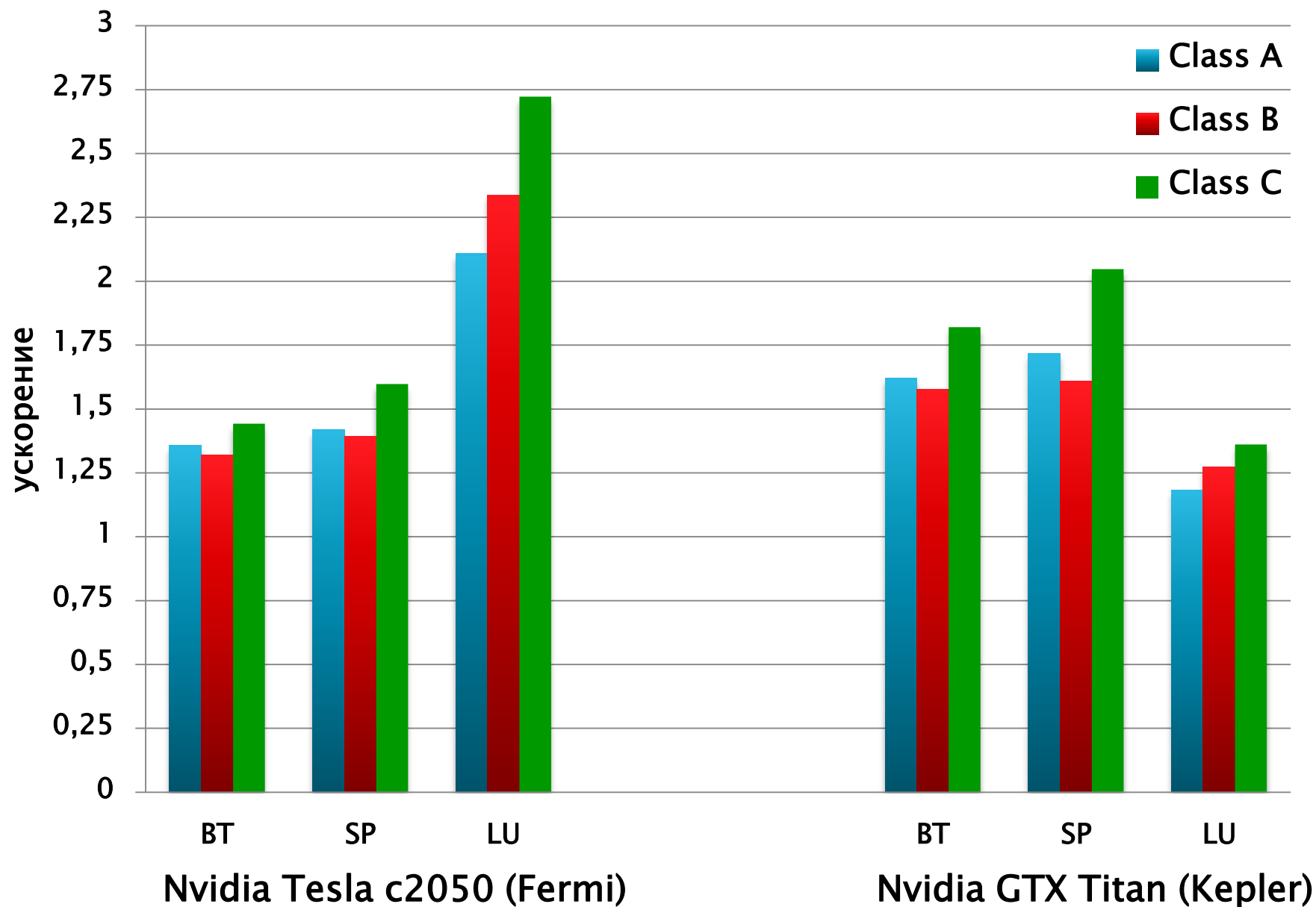


Переупорядочивание массивов

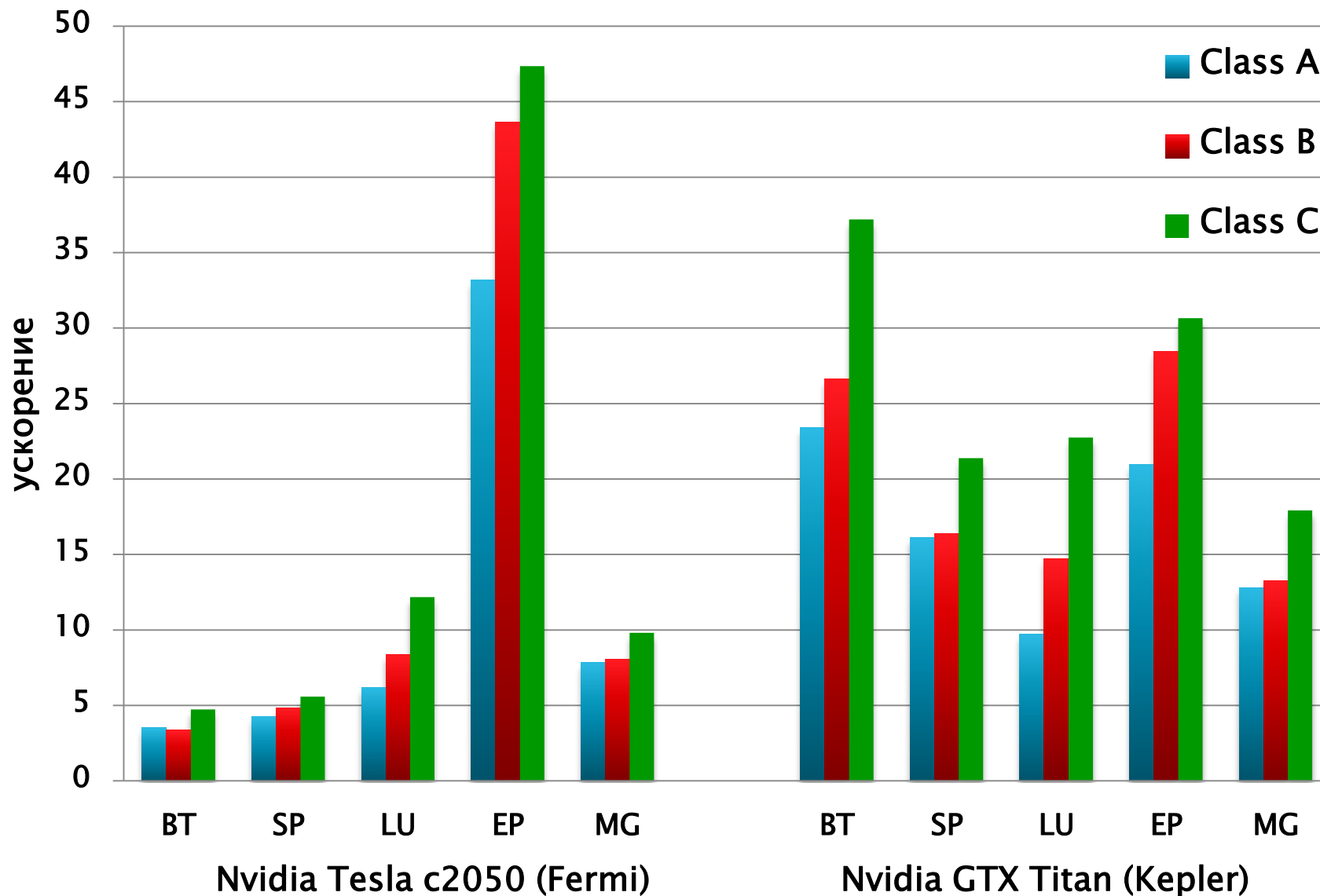
- ▶ Никаких дополнительных указаний в DVMH-программе
- ▶ Работает в динамике
- ▶ Для каждого цикла для каждого массива выбирает лучший порядок элементов
- ▶ Поддержка диагонализированных представлений
- ▶ Не происходит возврата состояния в конце цикла, только переход в требуемое



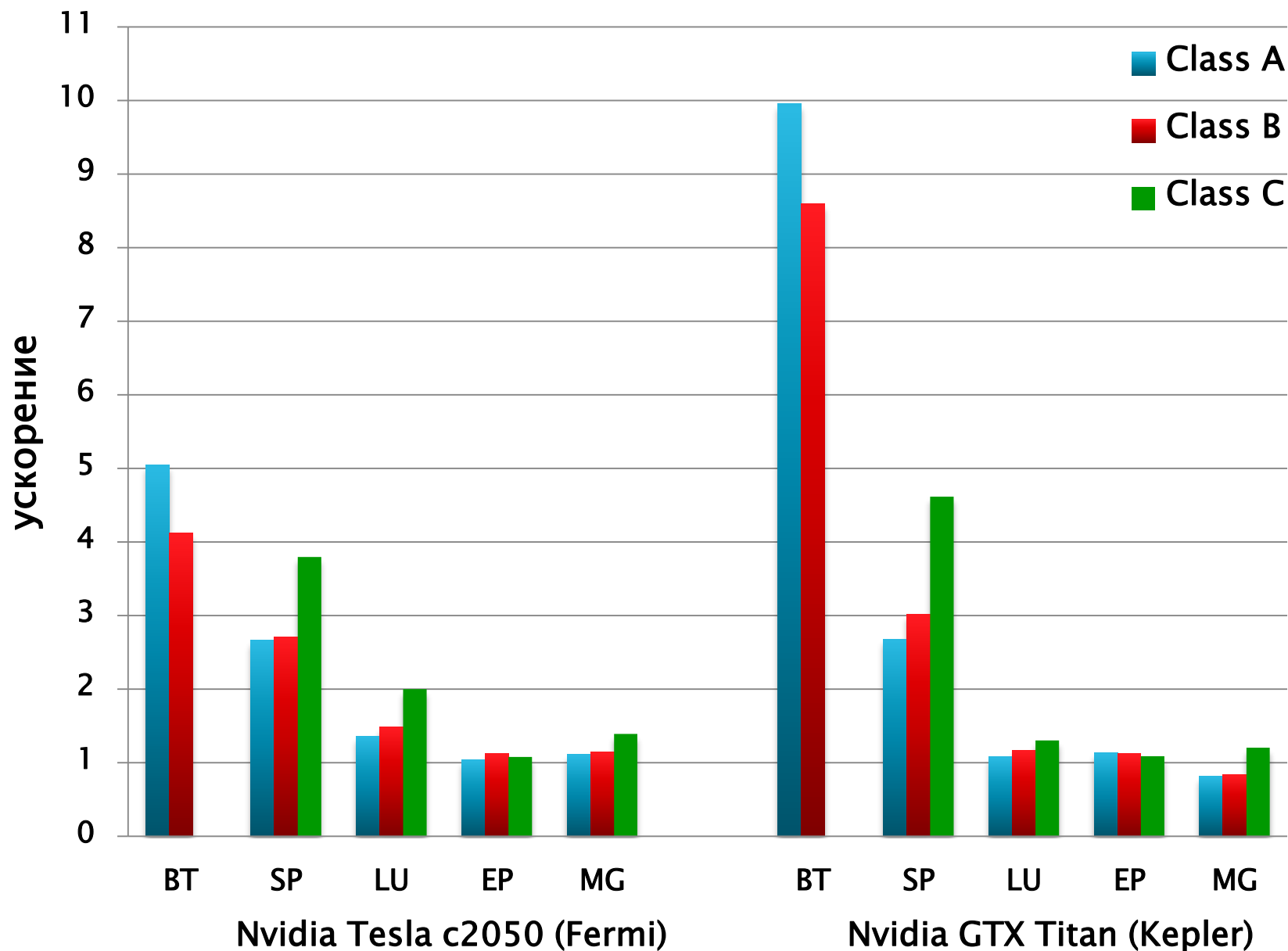
Ускорение выполнения DVMH-версий программ в результате выполнения динамического переупорядочивания массивов



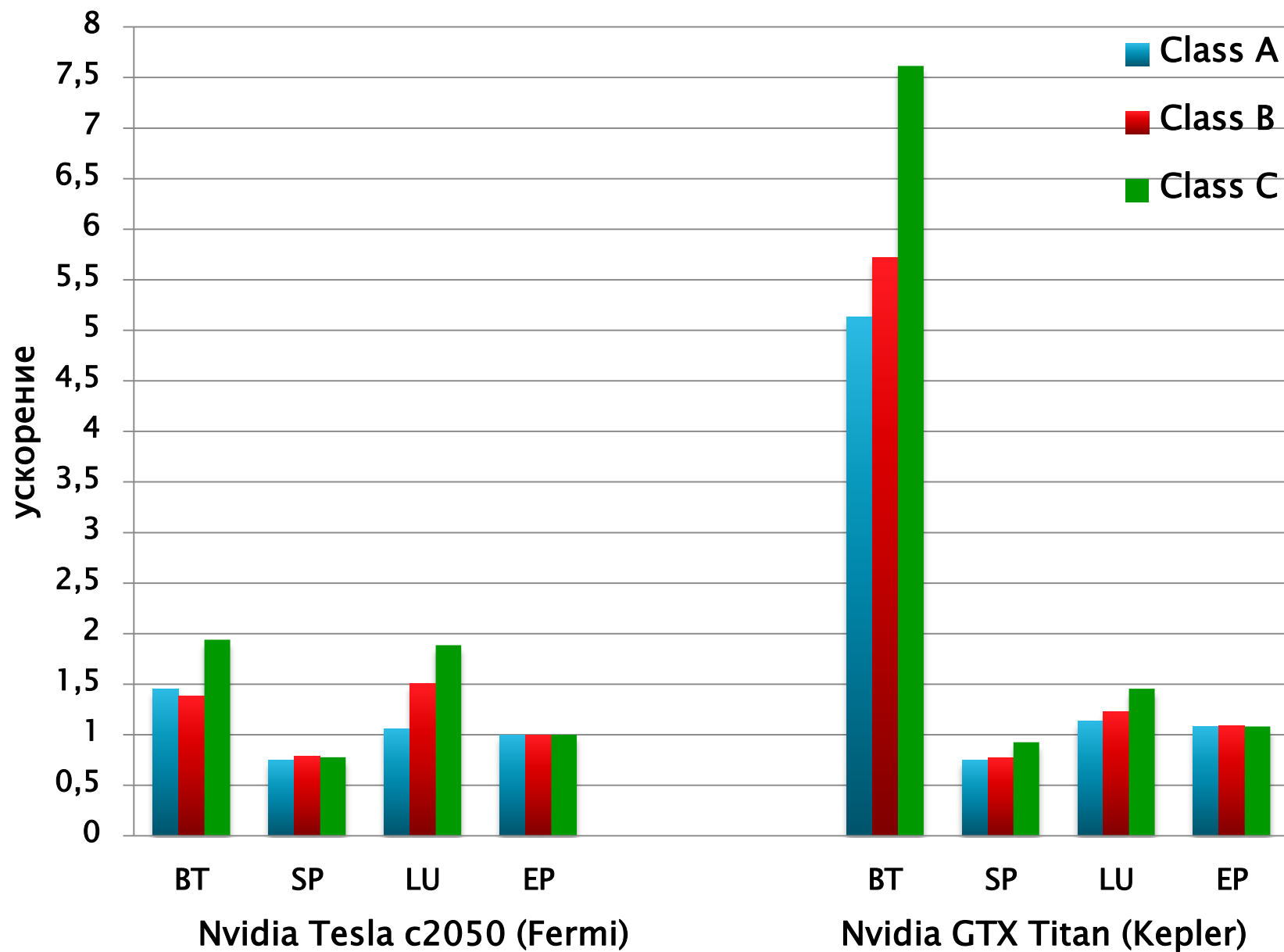
Ускорение выполнения DVMH-версий тестов NAS NPB по сравнению с последовательными версиями тестов на CPU Intel Xeon 5670



Ускорение выполнения DVMH-версий тестов NAS NPB по сравнению с OpenCL-версиями



Ускорение выполнения DVMH-версий тестов NAS NPB по сравнению с CUDA-версиями



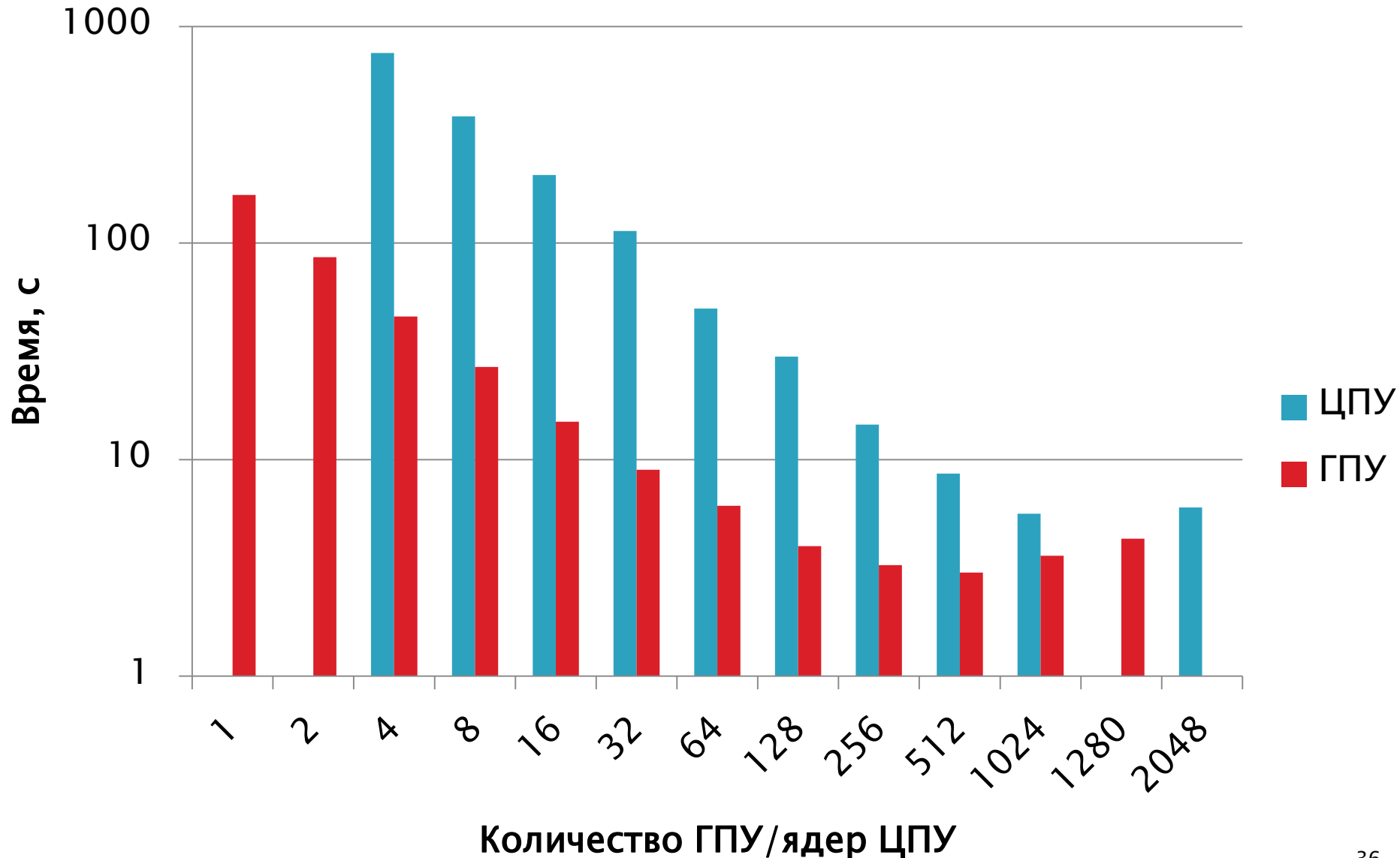
Применение DVMN для приложений

- ▶ Контейнер
- ▶ Каверна
- ▶ Состояния кубитов
- ▶ Кристаллизация3D
- ▶ Спекание2D
- ▶ Спекание3D

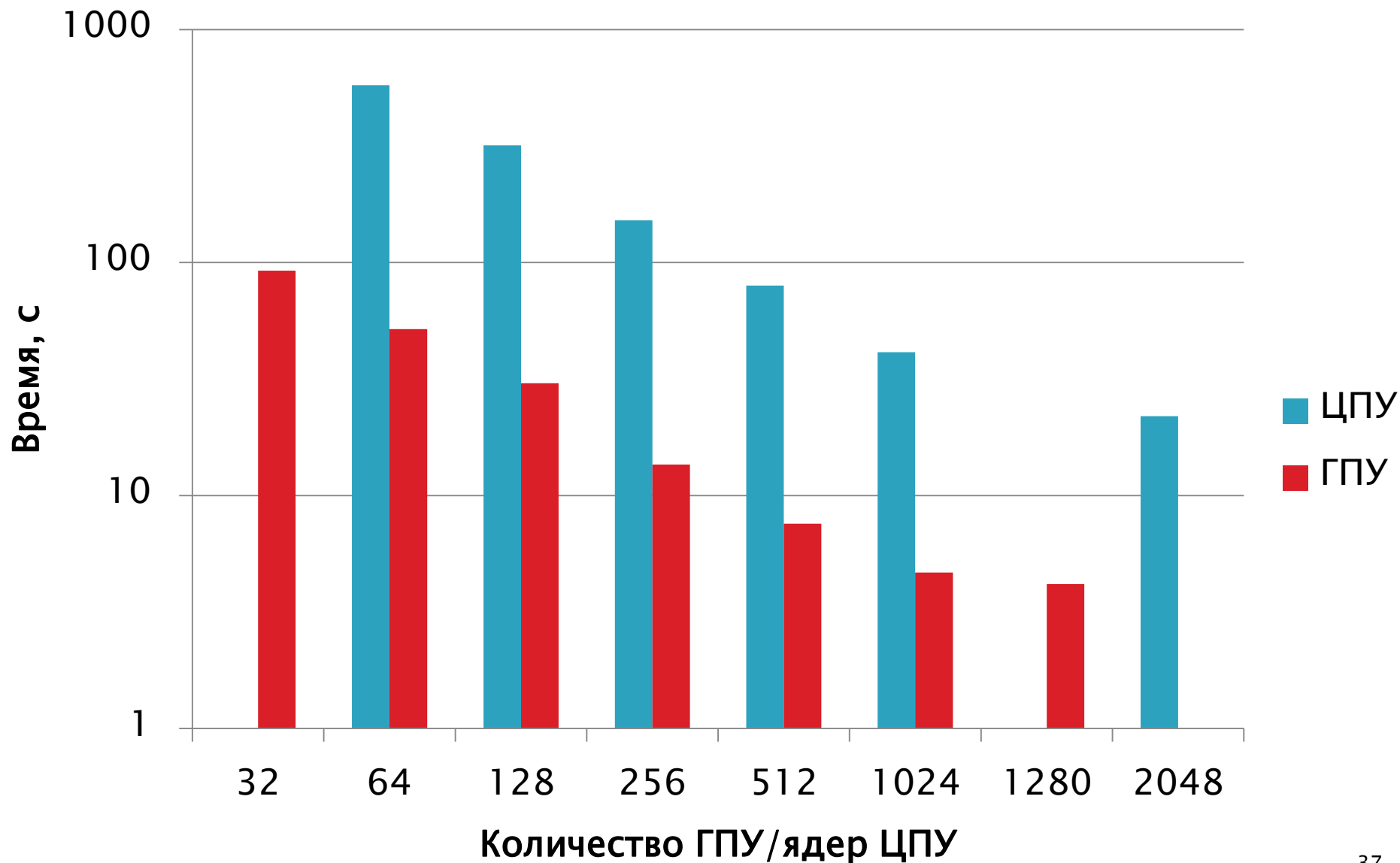
Программа Контейнер

- ▶ Моделирование течения вязкой тяжелой жидкости под действием силы тяжести в прямоугольном контейнере с открытой верхней стенкой и отверстием в одной из боковых стенок в трехмерной постановке
- ▶ Последовательная программа 828 строк
- ▶ Параллельная программа 942 строки
 - 26 распределенных массивов
 - 5 вычислительных регионов
 - 21 параллельный цикл
 - 7 директив `actual/get_actual`

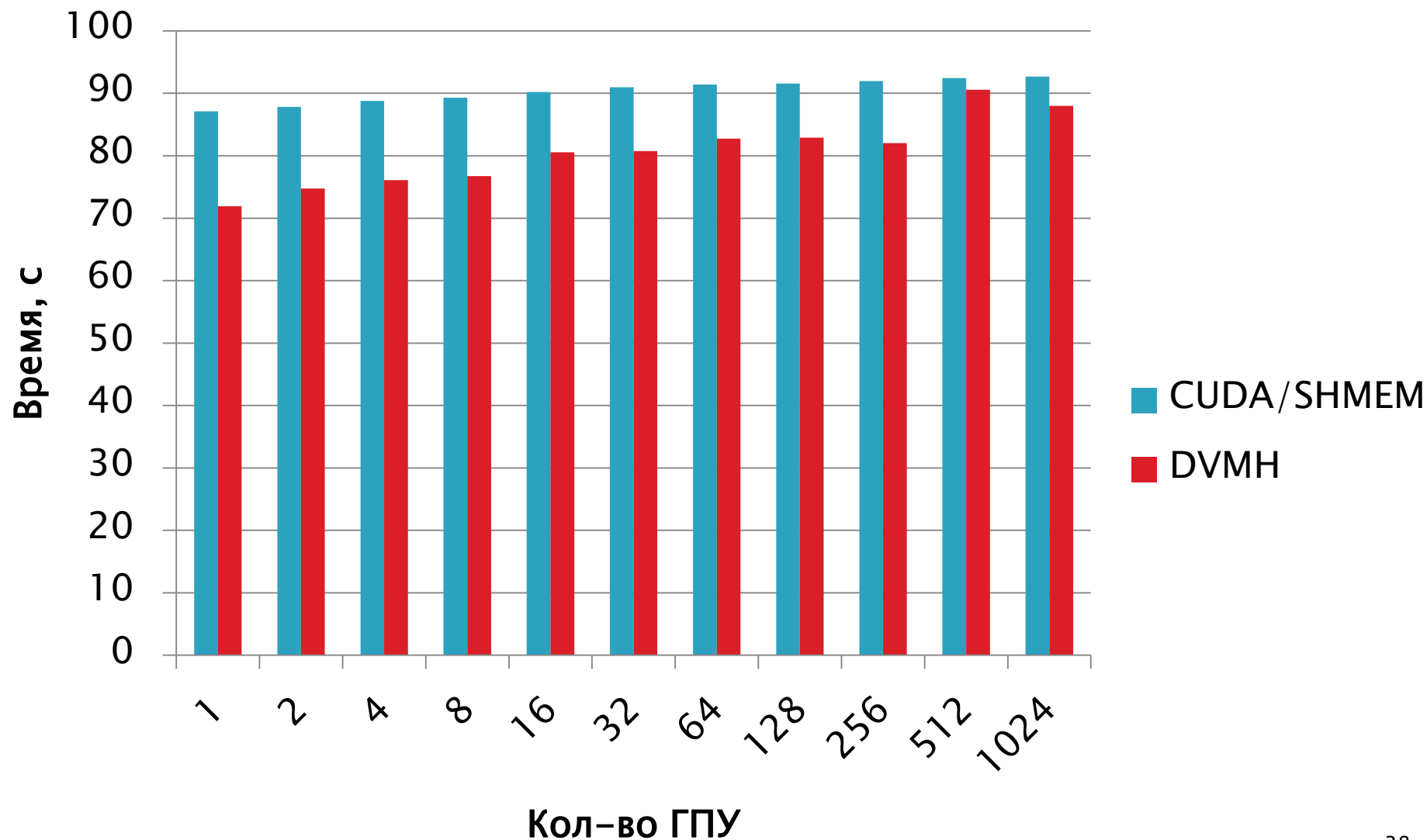
Время выполнения программы Контейнер 200x200x200, itmax=200



Время выполнения программы Контейнер 800x800x800, itmax=50



Программа Контейнер. Сравнение DVMH vs C+SHMEM+CUDA

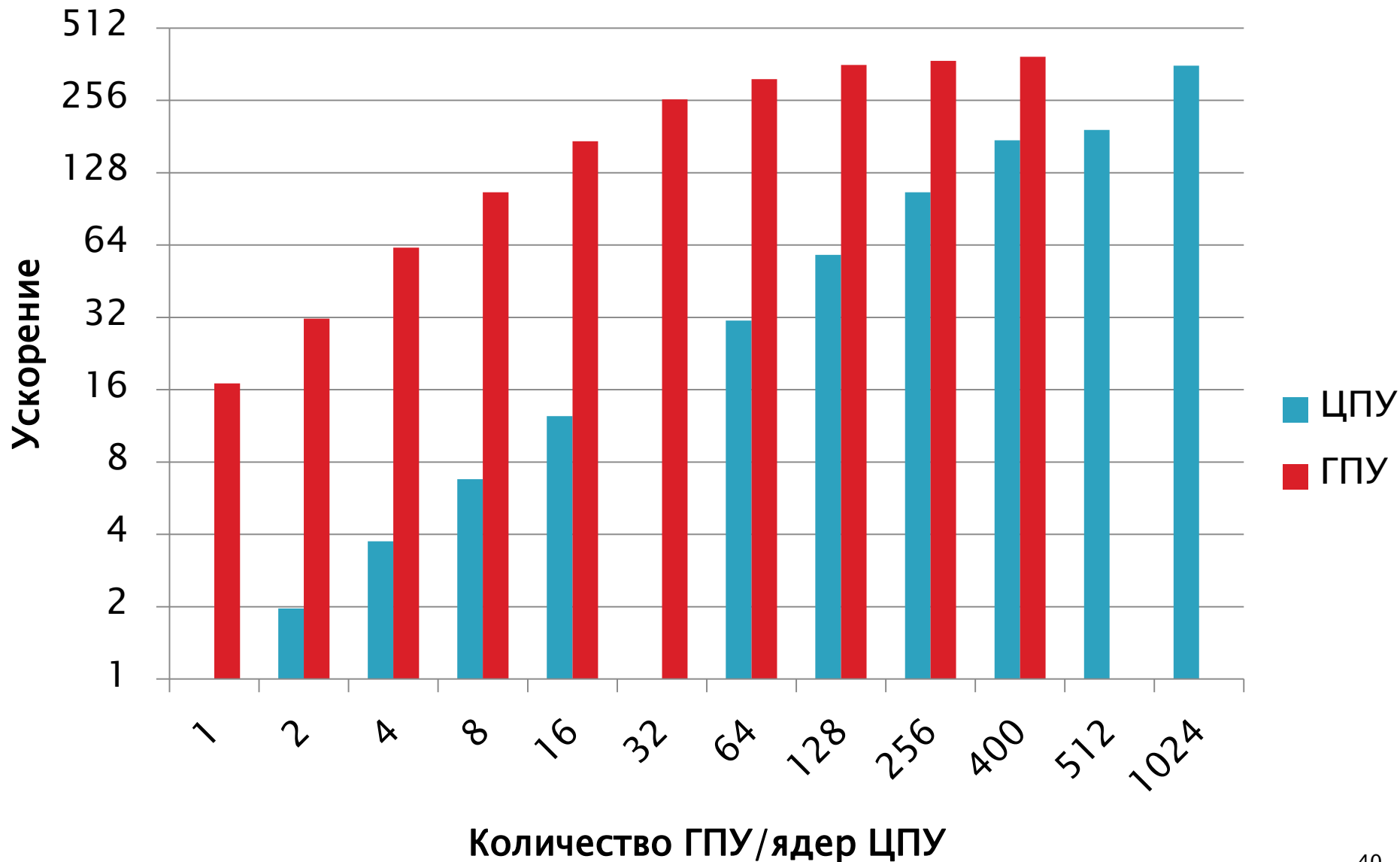


Программа Каверна

- ▶ Моделирование циркуляционного течения в плоской квадратной каверне с движущейся верхней крышкой в двумерной постановке
- ▶ Последовательная программа 496 строк
- ▶ Параллельная программа 613 строк
 - 18 распределенных массивов
 - 7 вычислительных регионов
 - 28 параллельных циклов
 - 11 директив `actual/get_actual`



Ускорение программы Каверна 3200x3200, itmax=200

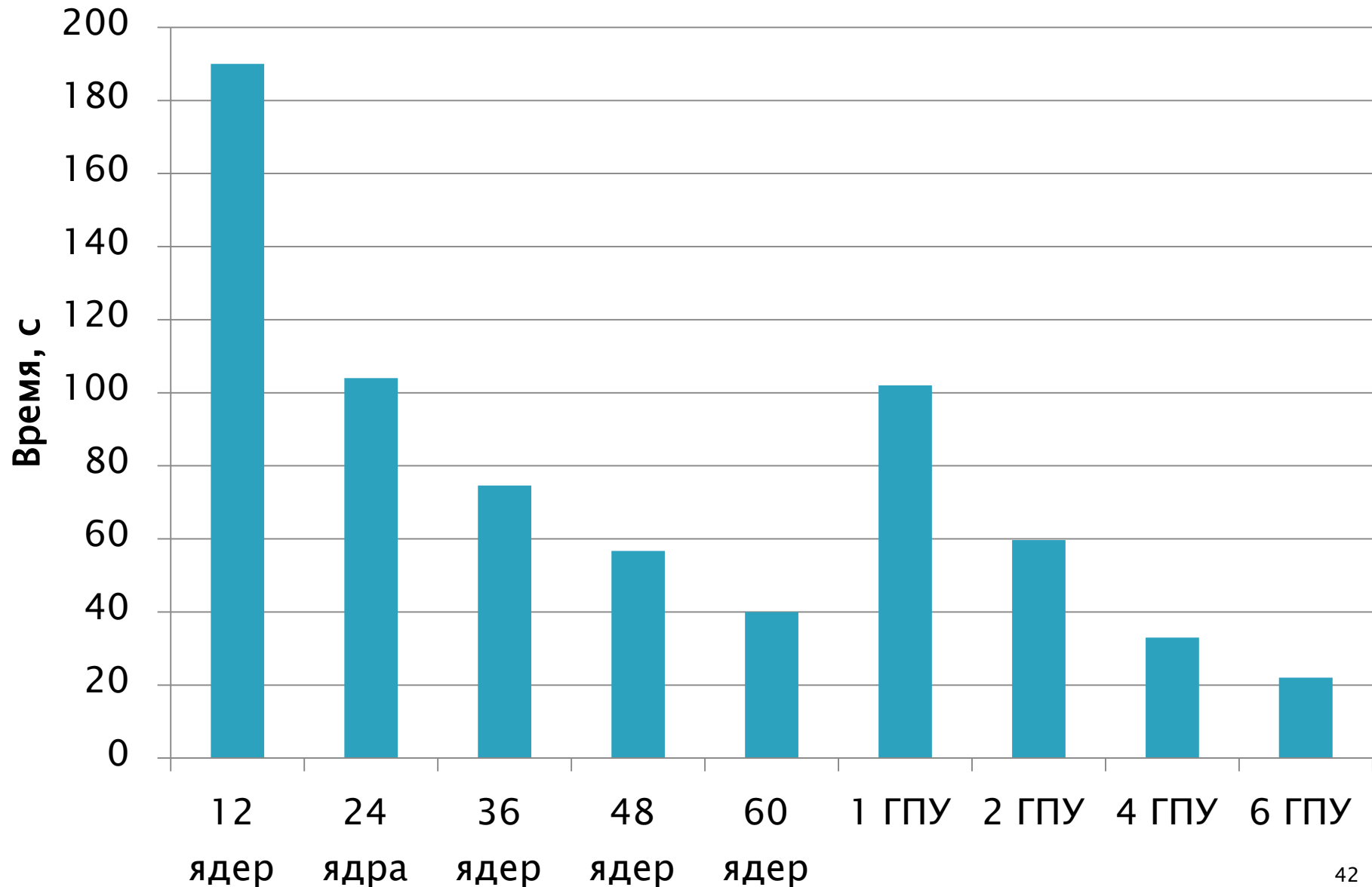


Программа «Состояния кубитов»

- ▶ Проведение трехмерных нестационарных расчетов состояния кубитов квантового компьютера на основе совместного решения трехмерного уравнения Пуассона и нестационарного уравнения Шредингера (683 строки).
- ▶ Параллельная программа 1011 строк
 - 23 распределенных массива
 - 5 вычислительных регионов
 - 61 параллельный цикл
 - 5 директив `actual/get_actual`



Программа «Состояния кубитов» 121x121x241, itmax=192

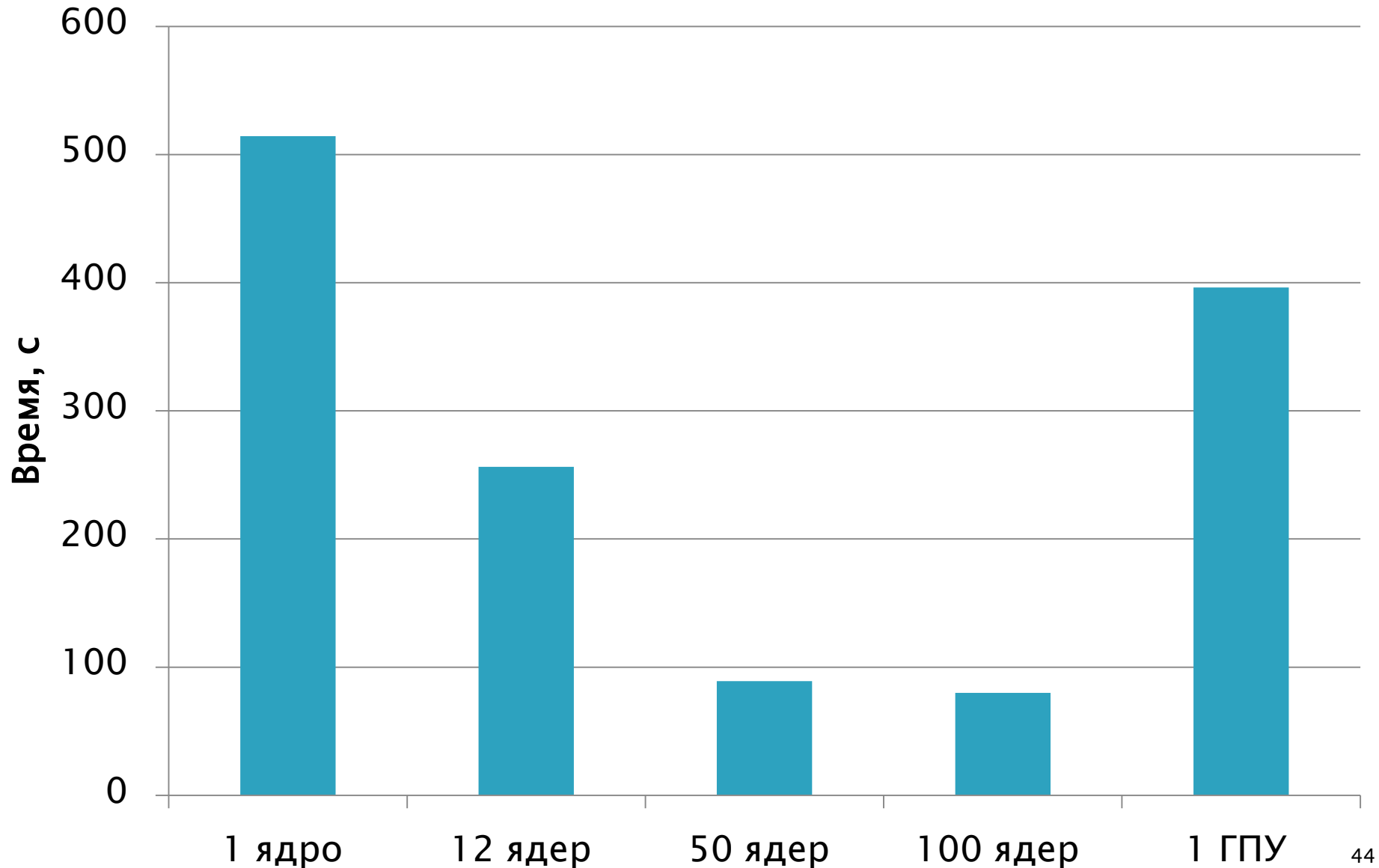


Программа Кристаллизация3D

- ▶ Трехмерное моделирование процессов объемной кристаллизации при воздействии на образец лазерного или электронного пучка на основе многокомпонентной и многофазной гидродинамической модели.
- ▶ Последовательная программа 530 строк
- ▶ Параллельная программа 865 строк
 - 21 распределенный массив
 - 6 вычислительных регионов
 - 27 параллельных циклов
 - 5 директив actual/get_actual



Время выполнения программы Кристаллизация 3D (103x5x103)

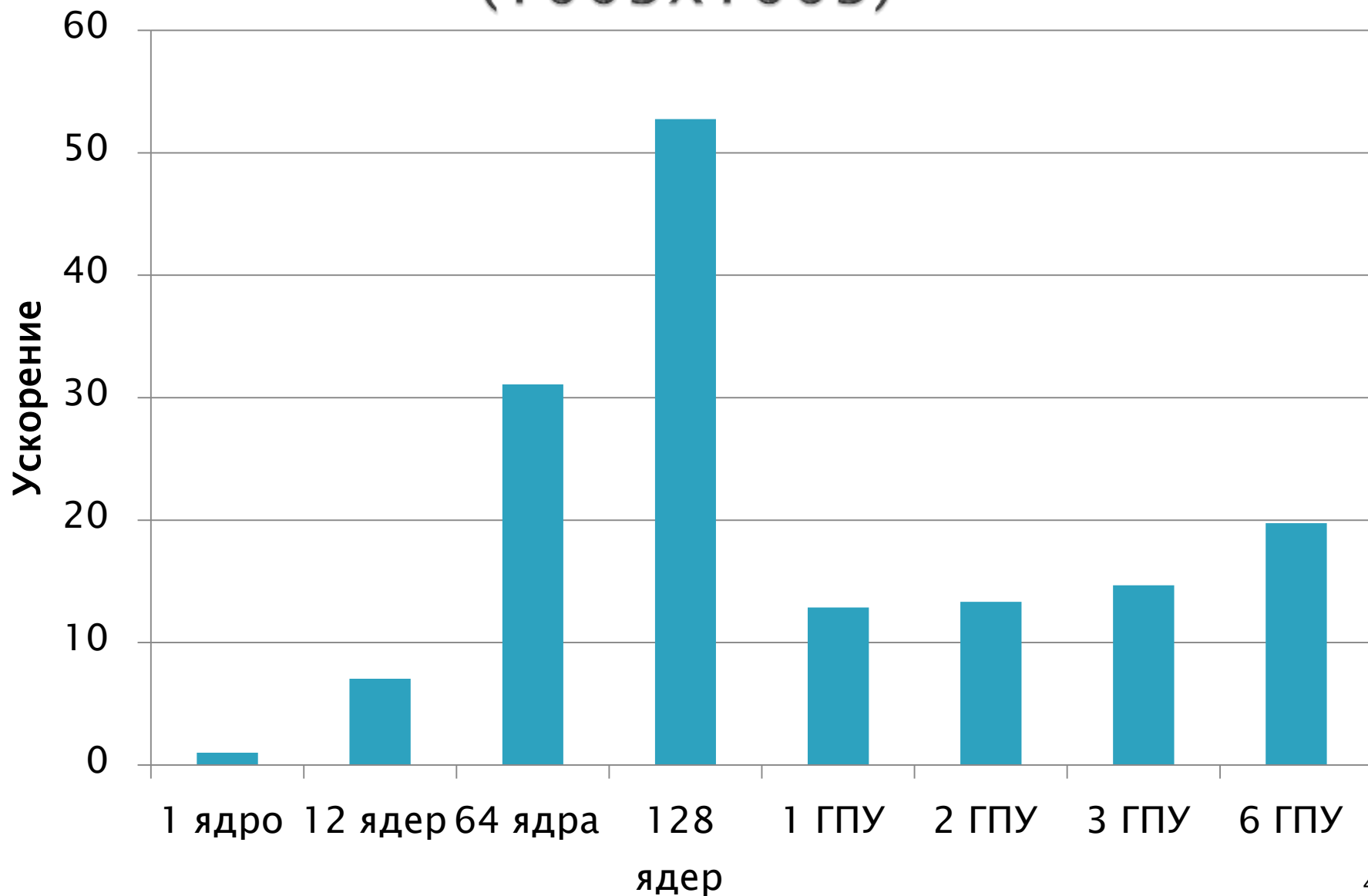


Программа Спекание2D

- ▶ Двухмерное моделирование процессов плавления многокомпонентных порошков при селективном лазерном спекании на основе многокомпонентной и многофазной гидродинамической модели
- ▶ Последовательная программа 831 строка
- ▶ Параллельная программа 1167 строк
 - 17 распределенных массивов
 - 8 вычислительных регионов
 - 21 параллельных циклов
 - 19 директив actual/get_actual



Ускорение программы Спекание2D (1003x1003)

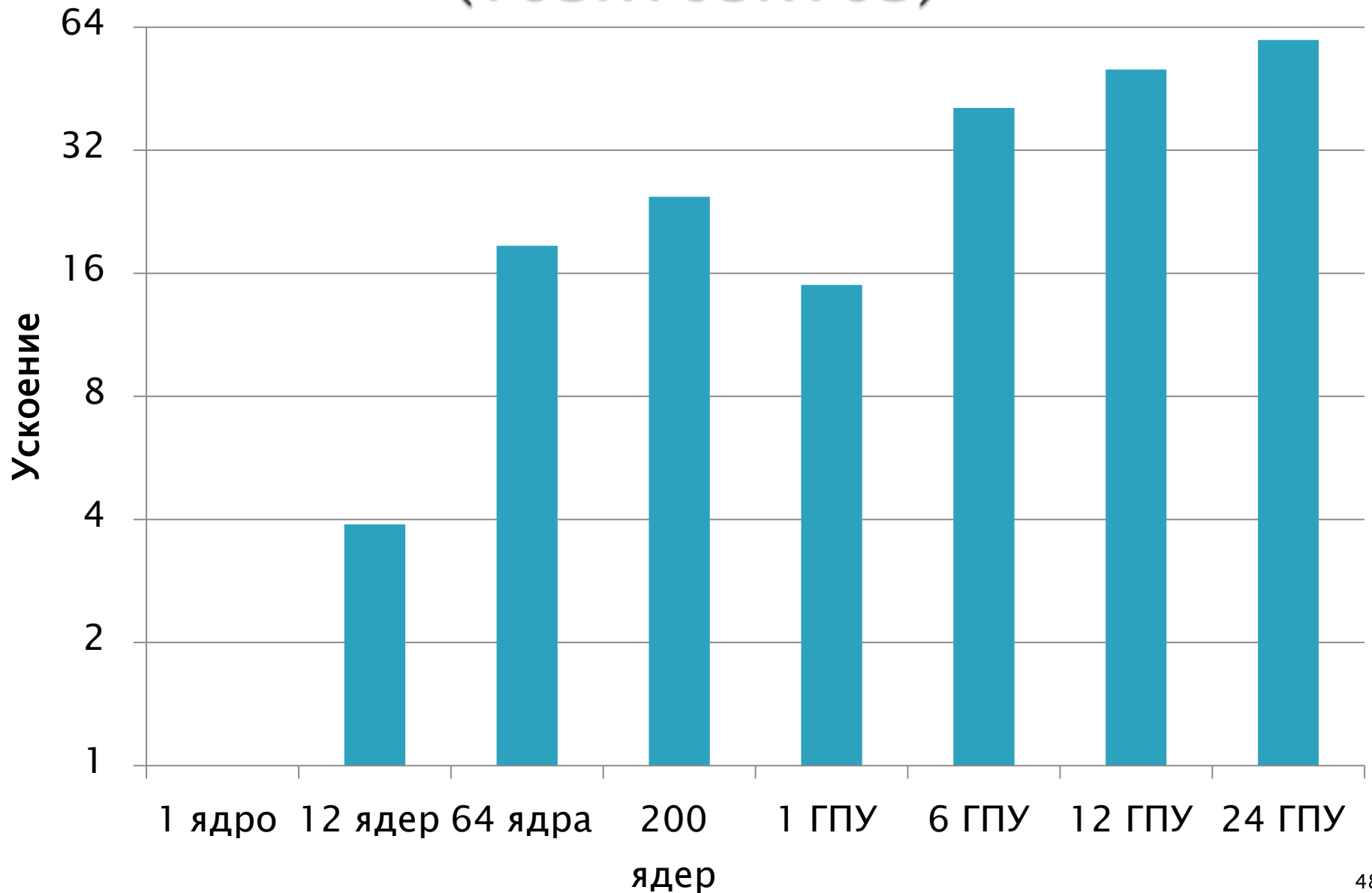


Программа Спекание3D

- ▶ Трехмерное моделирование процессов плавления многокомпонентных порошков при селективном лазерном спекании на основе многокомпонентной и многофазной гидродинамической модели
- ▶ Последовательная программа 677 строк
- ▶ Параллельная программа 1236 строк
 - 19 распределенных массивов
 - 4 вычислительных региона
 - 51 параллельных циклов
 - 22 директивы actual/get_actual



Ускорение программы Спекание3D (103x103x103)



Выводы

- ▶ Разработанный подход к созданию прикладного программного обеспечения (DVM-подход) существенно упрощает создание прикладных программ для суперкомпьютерных систем с ускорителями.
- ▶ Эффективность программ, написанных в высокоуровневой модели DVMH, близка к эффективности программ, написанных с использованием низкоуровневых моделей CUDA или OpenCL и превосходит эффективность программ в модели OpenACC.
- ▶ Высокоуровневые языки C-DVMH и Fortran-DVMH позволяют обеспечить эффективное отображение одной и той же DVMH-программы на вычислительные системы различной архитектуры.



Планы развития DVM–системы

- ▶ Поддержка неструктурных сеток и разреженных матриц
- ▶ Исследование эффективности отображения DVMH-программ на сопроцессоры Intel Xeon Phi
- ▶ Доработка системы отладки эффективности для DVMH-программ



Вопросы, замечания?

СПАСИБО !

<http://www.keldysh.ru/dvm>
dvm@keldysh.ru



Алгоритм Якоби

- ▶ Двумерный 5002×5002 , два цикла в регионе + редукция
 - 1 ГПУ – 57 с.
 - 1 ГПУ + 4 ядра ЦПУ – 49 с.
- ▶ Трехмерный $502 \times 252 \times 252$, один цикл в регионе
 - 1 ГПУ – 14,7 с.
 - 1 ГПУ + 4 ядра ЦПУ – 12,6 с.

