



**Нижегородский государственный университет
им. Н.И.Лобачевского**

Факультет Вычислительной математики и кибернетики

Программирование для Intel Xeon Phi

**Оптимизация вычислений
в задаче матричного умножения. Оптимизация
работы с памятью**

Сиднев А.А.

Архангельск, 2014

GEMM

- GEneral Matrix Multiplication

$$C = \alpha \cdot A \cdot B + \beta \cdot C$$

α, β – скаляры

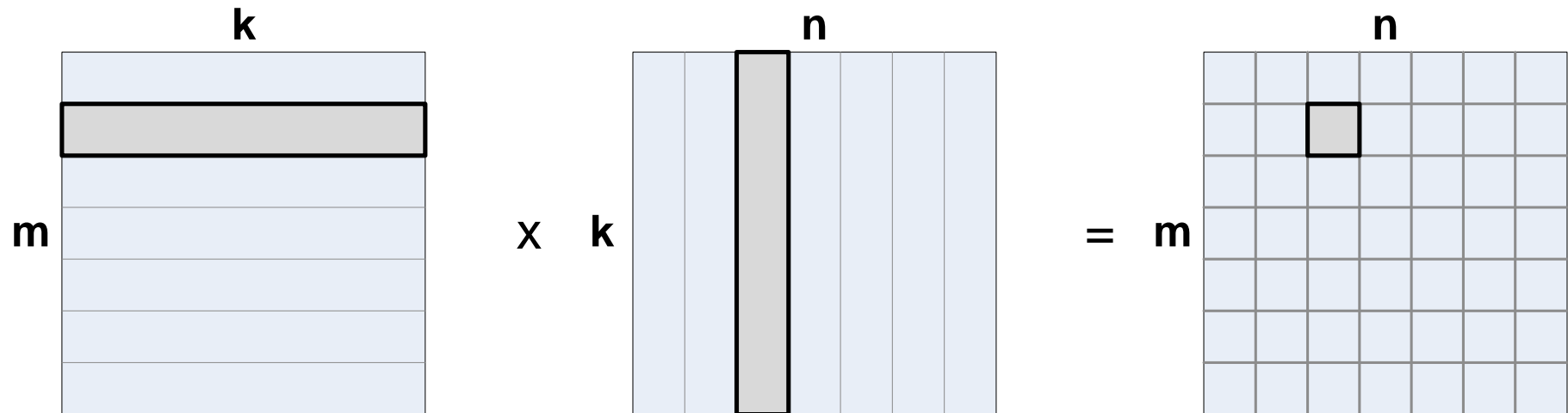
A, B, C – матрицы

$$A \in R^{m \times k} \quad B \in R^{k \times n} \quad C \in R^{m \times n}$$



Умножение матриц

- Основная операция GEMM



- Сложность алгоритма $O(n^3)$

Задача матричного умножения

- Известны последовательные алгоритмы умножения матриц, обладающие меньшей вычислительной сложностью:
 - алгоритм Штрассена $O(n^{2,81})$;
 - алгоритм Копперсмита-Винограда $O(n^{2,376})$;
 - алгоритм Вирджинии Вильямс $O(n^{2,373})$.
- На практике эти алгоритмы применяют редко.
 - Как правило, в оценке сложности присутствует большая константа.
 - Алгоритмы обладают лучшей производительностью только при слишком больших размерах матриц, не помещающихся в оперативную память современных компьютеров.



Эффективность матричного умножения

- Определяется эффективностью использования:
 - кэшей
 - TLB-кэшей
 - векторных команд
 - программной/аппаратной предвыборки данных
- Справедливо как для последовательной, так и параллельной версии



Компиляция и подготовка запуска на Tornado

- ❑ Скопировать директорию /tmp/xeonphilabs/lab2 в домашний каталог /home/fpkX/xeonphilabs/lab2

```
cp -R /tmp/xeonphilabs/lab2/ /home/fpk20/xeonphilabs/lab2/
```

- ❑ Перейти в каталог

```
cd ~/xeonphilabs/lab2/
```

- ❑ Скомпилировать программу под Xeon Phi

```
icc -openmp -mkl -mmic gemm.cpp -o gemm.mic
```

- ❑ Загрузить запускающий модуль для Xeon Phi

```
module load launcher/mic
```



Запуск на Tornado

❑ Установить переменную окружения

```
export
```

```
LD_LIBRARY_PATH=/opt/software/intel/composer_xe_2013.5.192/mk  
l/lib/mic:/opt/software/intel/impi/4.1.0.024/mic/lib:/opt/sof  
tware/intel/composer_xe_2013/lib/mic:/opt/software/intel/com  
poser_xe_2013.5.192/tbb/lib/mic
```

❑ Запустить программу на Xeon Phi

```
sbatch -N 1 --gres=mic:1 --reservation=scc_phi_training  
native_run.sh ./gemm
```



Умножение матриц

- Реализация через подматрицы
- Базовые элементы (форма матрицы)

□ - Блок (Block)



- Полоса (Panel)



- Матрица (Matrix)

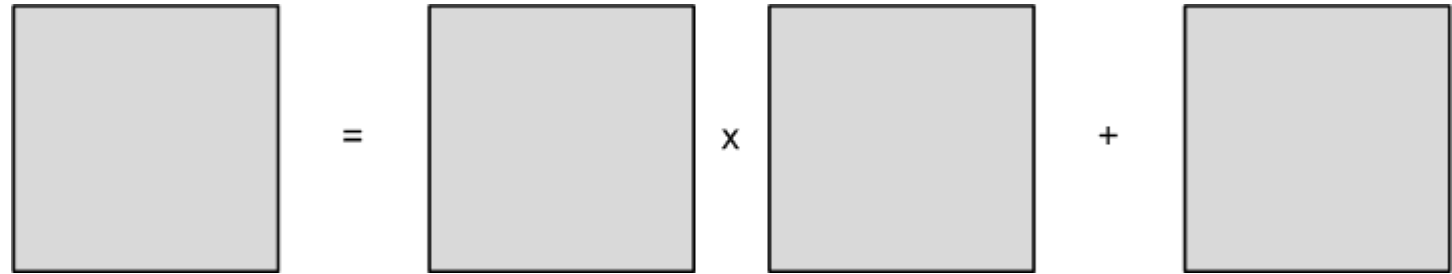
Варианты умножения матриц

- ❑ GEMM
- ❑ GEPP
- ❑ GEMP
- ❑ GEPM
- ❑ GEBP
- ❑ GEPB
- ❑ GEPDOT
- ❑ GEBB



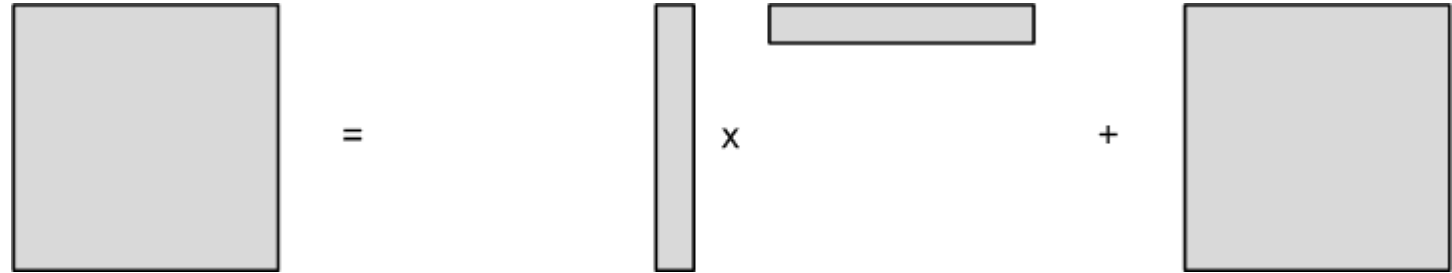
Варианты умножения матриц (GEMM)

- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEBB



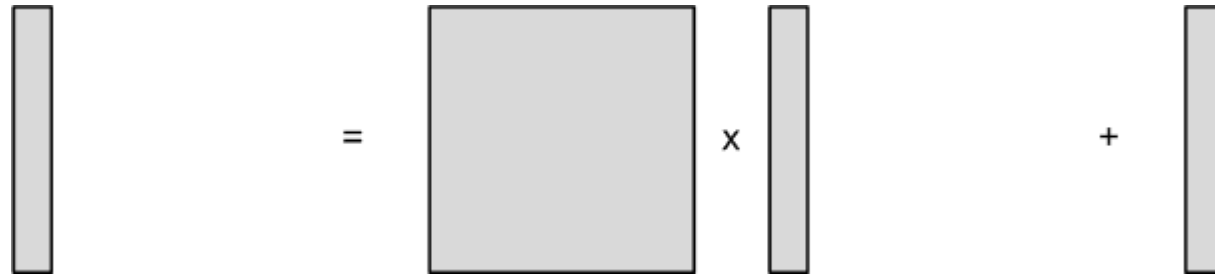
Варианты умножения матриц (GEMM)

- GEMM
- GEPB
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEBB



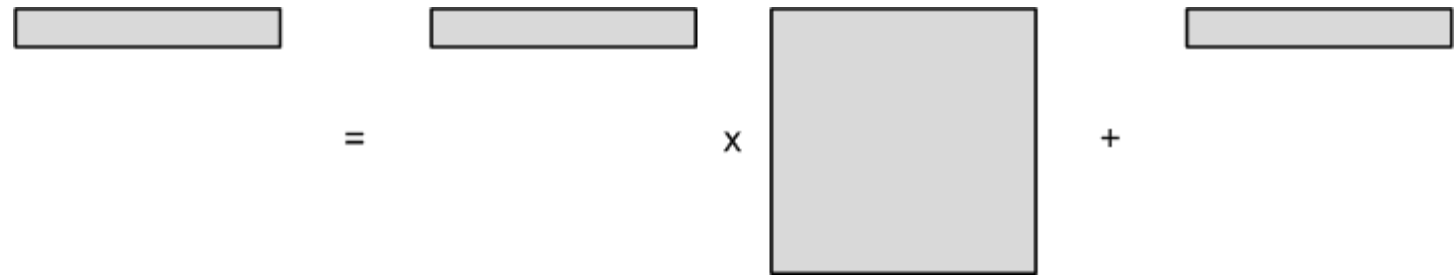
Варианты умножения матриц (GEMP)

- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEBB



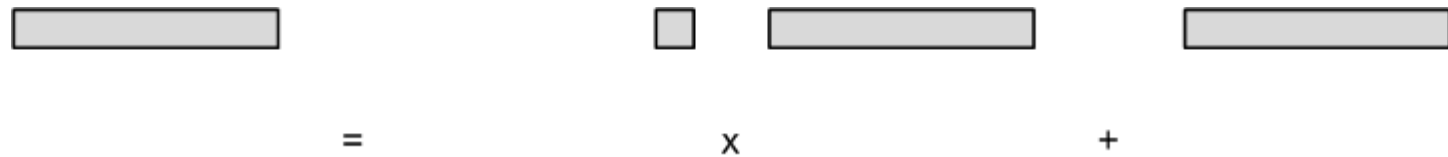
Варианты умножения матриц (GEMM)

- GEMM
- GEPP
- GEMP
- GERM
- GEBP
- GEPB
- GEPDOT
- GEBB



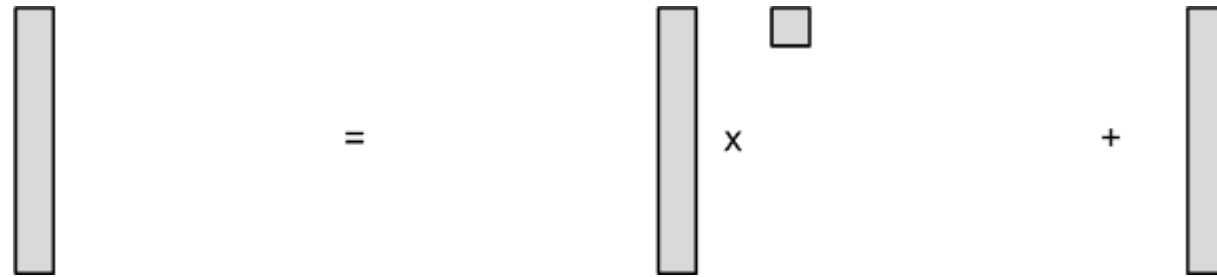
Варианты умножения матриц (GEBP)

- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEBB



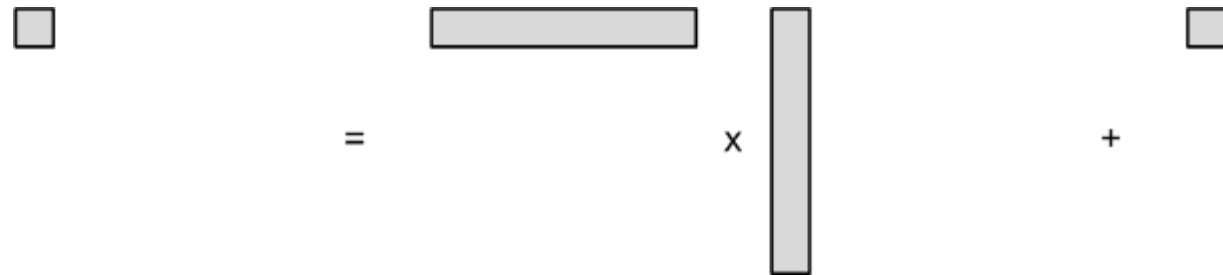
Варианты умножения матриц (GEPB)

- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB**
- GEPDOT
- GEBB



Варианты умножения матриц (GEPDOT)

- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEVB

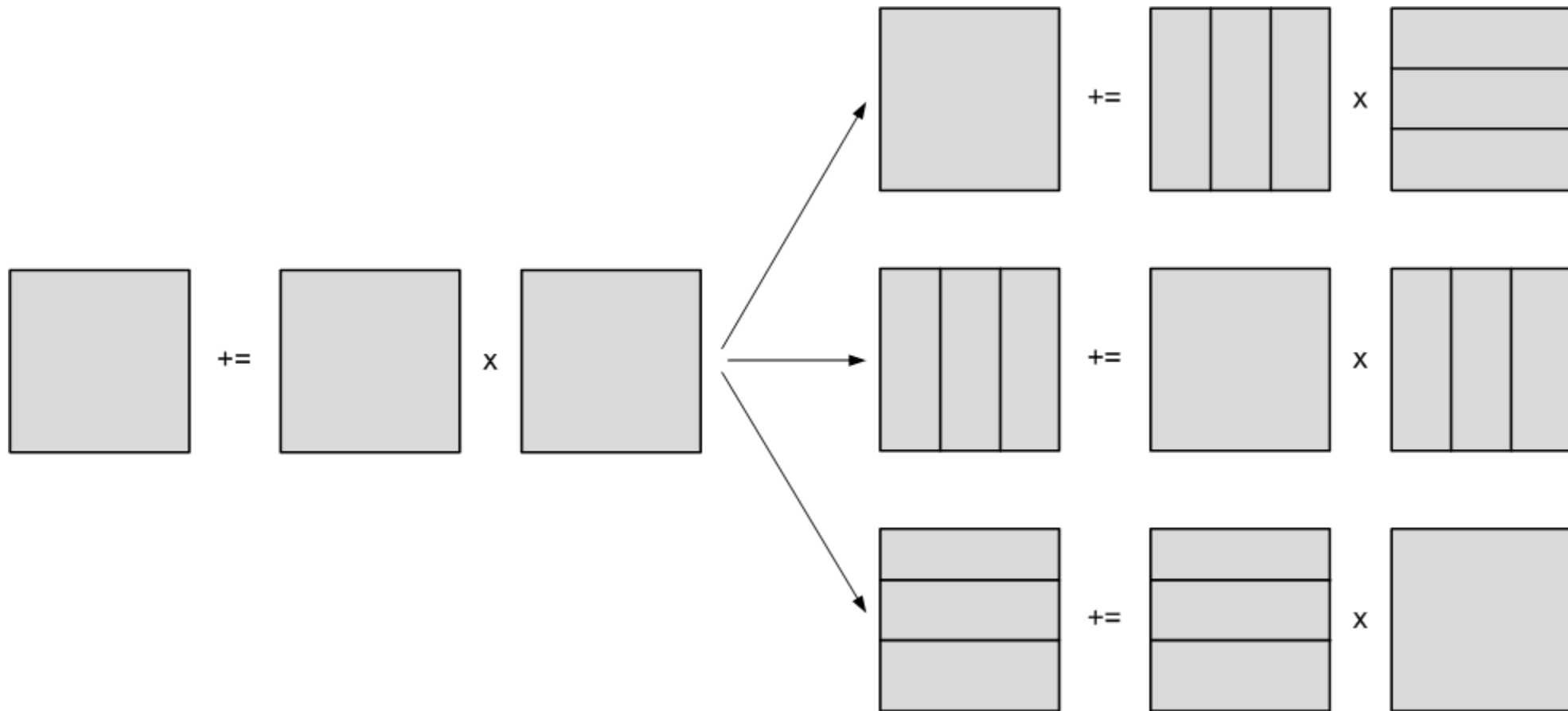


Варианты умножения матриц (GEBV)

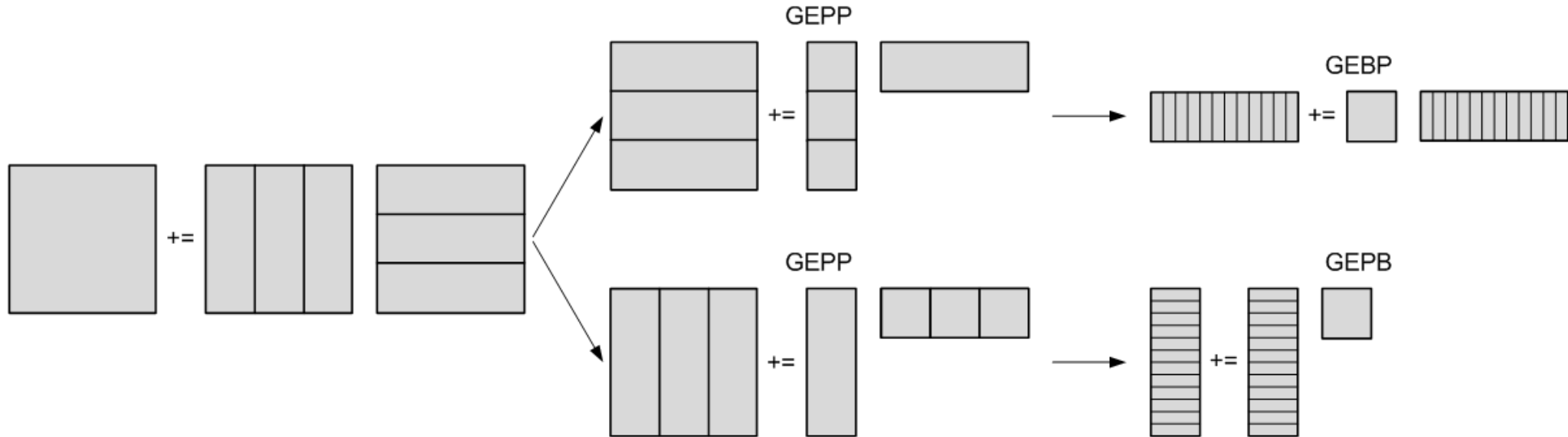
- GEMM
- GEPP
- GEMP
- GEPM
- GEBP
- GEPB
- GEPDOT
- GEBB



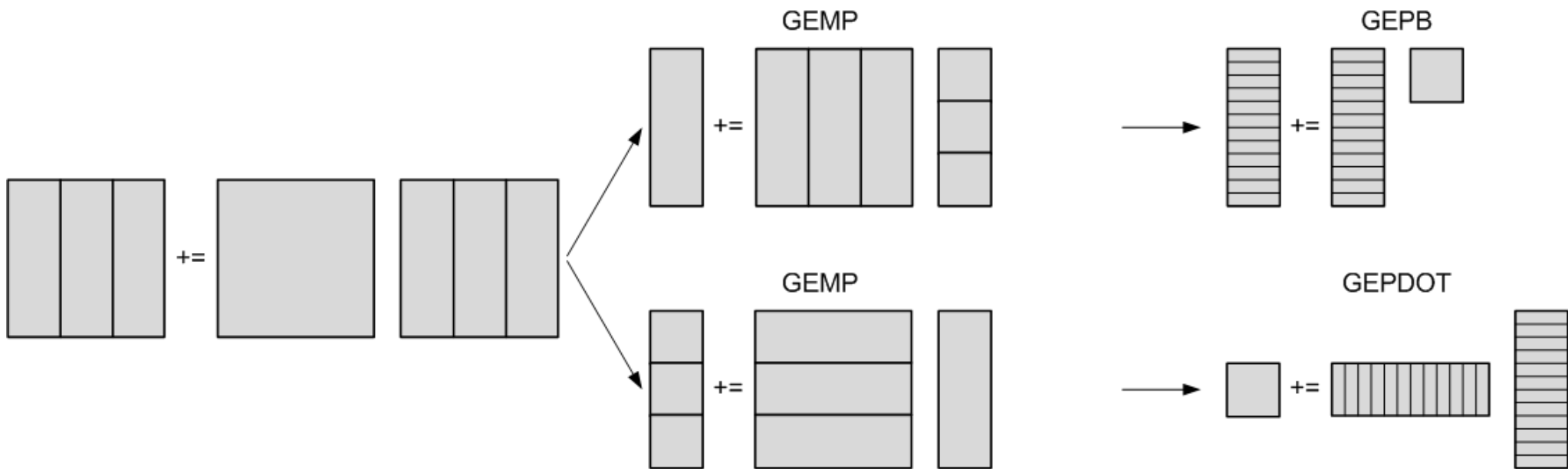
Верхний уровень



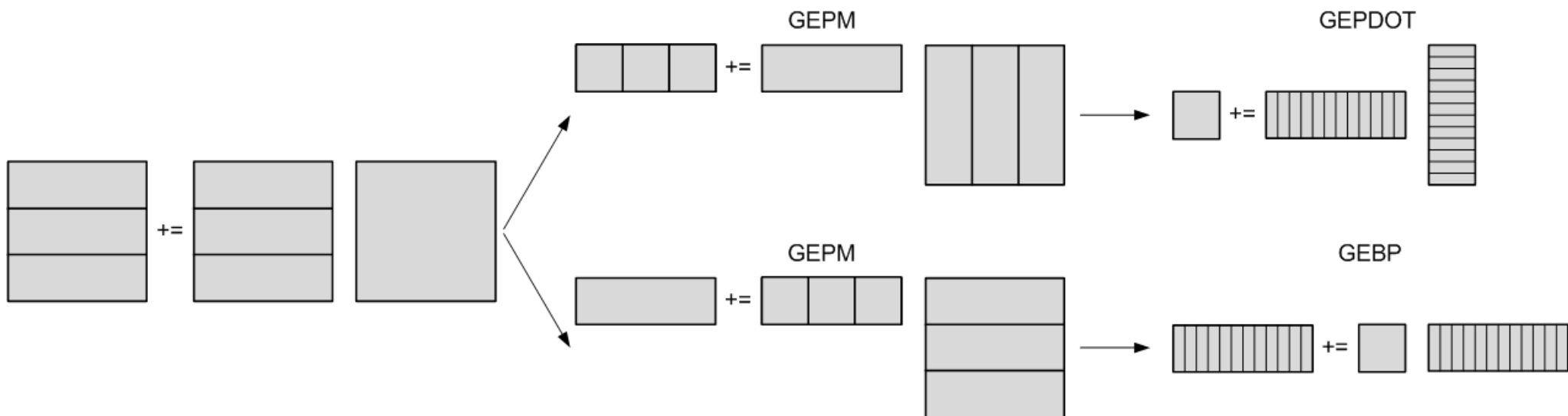
Средний и низкий уровни (1)



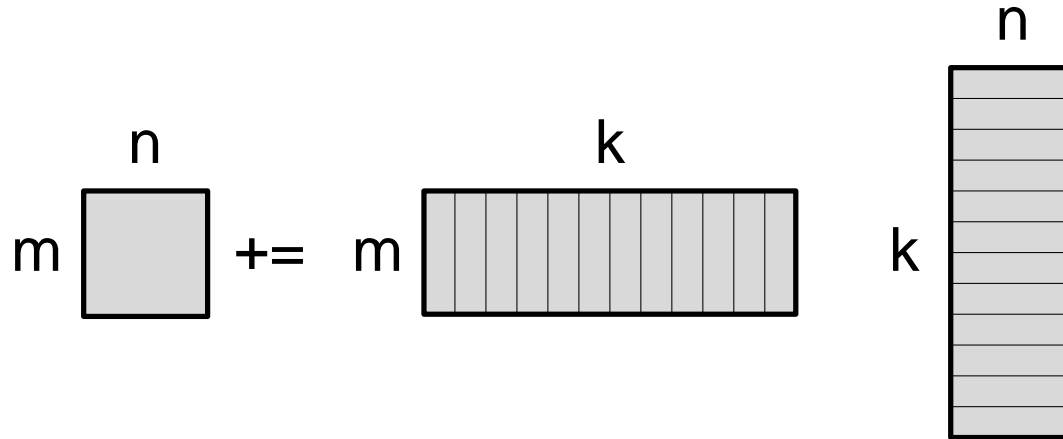
Средний и низкий уровни (2)



Средний и низкий уровни (3)



GEPDOT (1)



- ❑ Операций с памятью: $mk + kn + 2mn$
- ❑ Вычислительных операций: $2mkn$

GEPDOT (2)

- $m \ll k$
- $n \ll k$
- Количество вычислений на операцию:

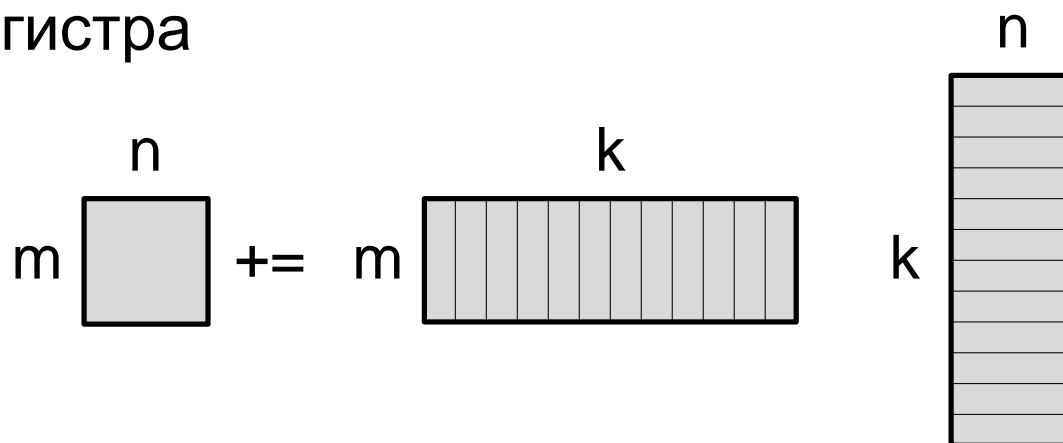
$$\frac{2mkn}{mk + kn + 2mn} \approx \frac{2mn}{m + n} \rightarrow \max$$

- Максимум достигается при $m=n$
- Ограничения:
 - Размер кэш-памяти
 - Количество регистров
 - Размер TLB-кэша



GEPDOT на Xeon Phi (1)

- Xeon Phi
 - Размер регистра 512 б
 - 32 векторных регистра



- $n=8$
 - 8 чисел двойной точности
- $m=31$
 - 31 регистр для хранения строк матрицы C
 - 1 регистр для хранения строки B

GEPDOT на Xeon Phi (2)

- Xeon Phi

- Размер L1 = 32 КБ

- Размер L2 = 512 КБ

- $k=400$

- $(400 \cdot 31 + 400 \cdot 8) \cdot 8 \approx 124 \text{ КБ}$

- Упаковка матриц A и B

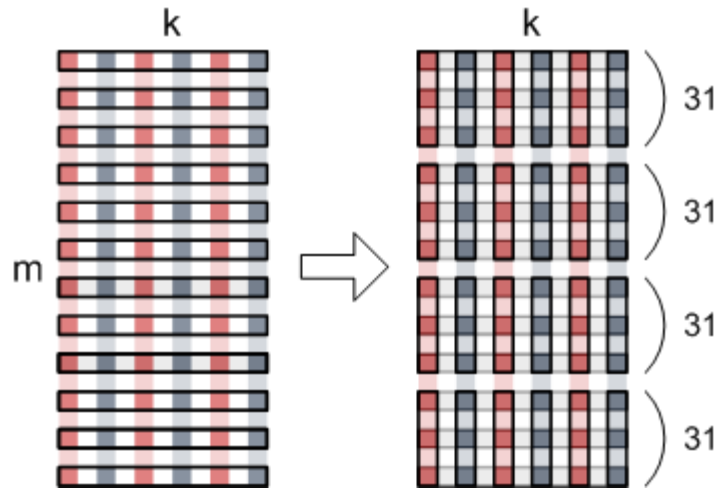
- Ограниченный размер кэша TLB

- Последовательный доступ (проще загружать в регистры)

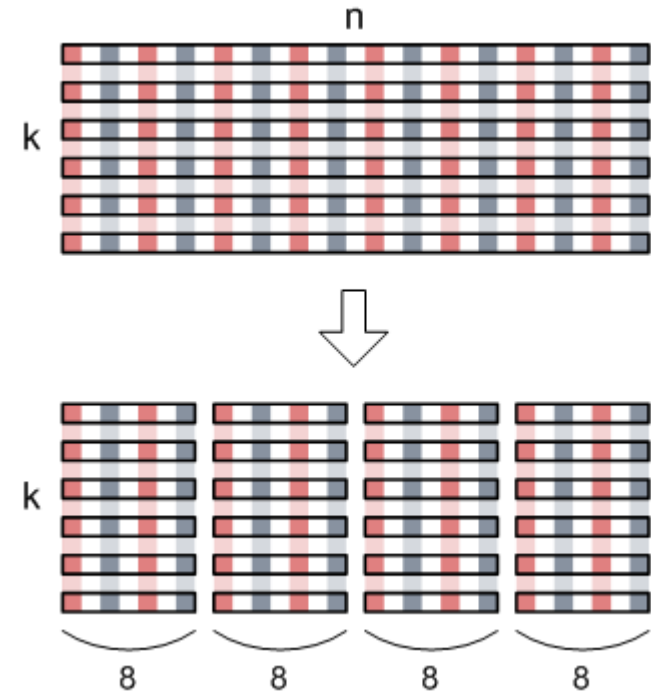


Упаковка матриц

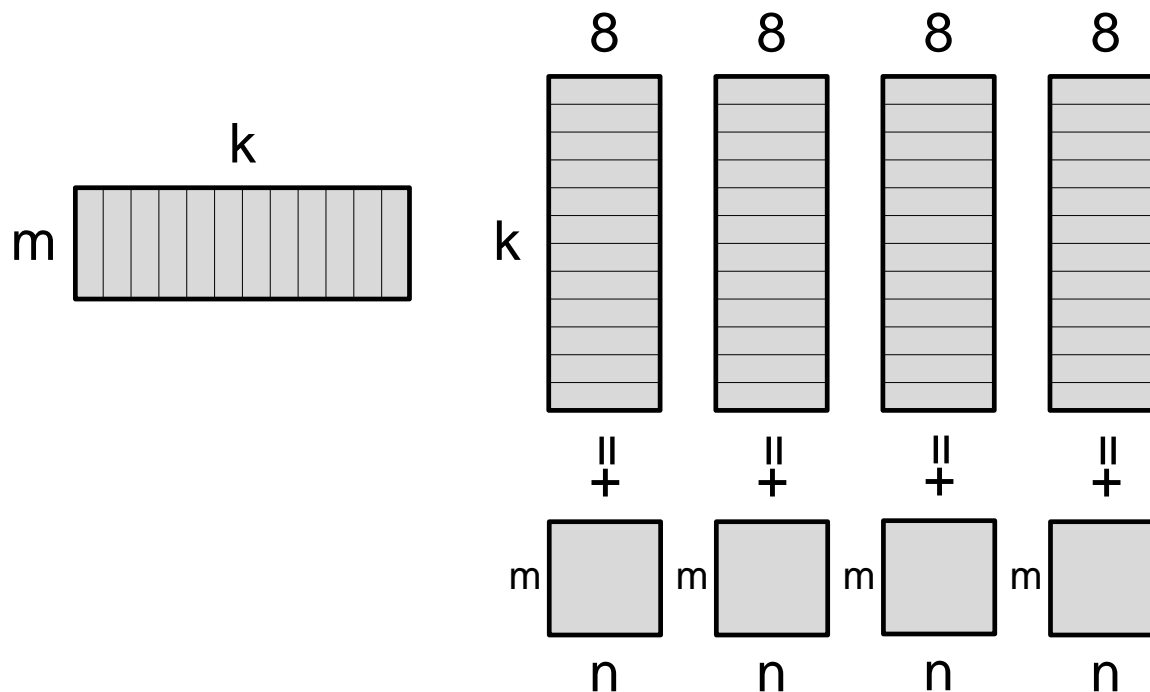
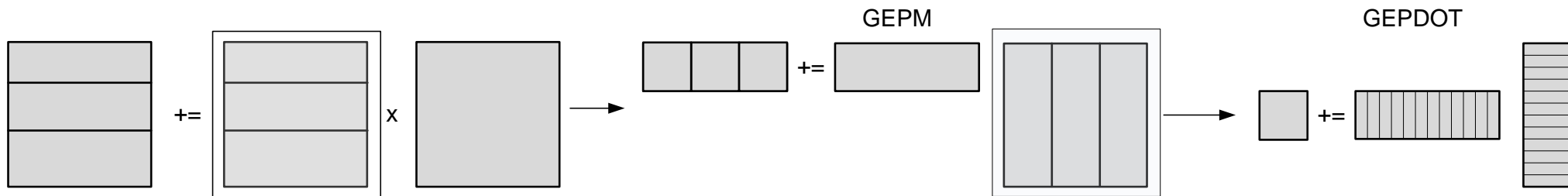
Матрица A



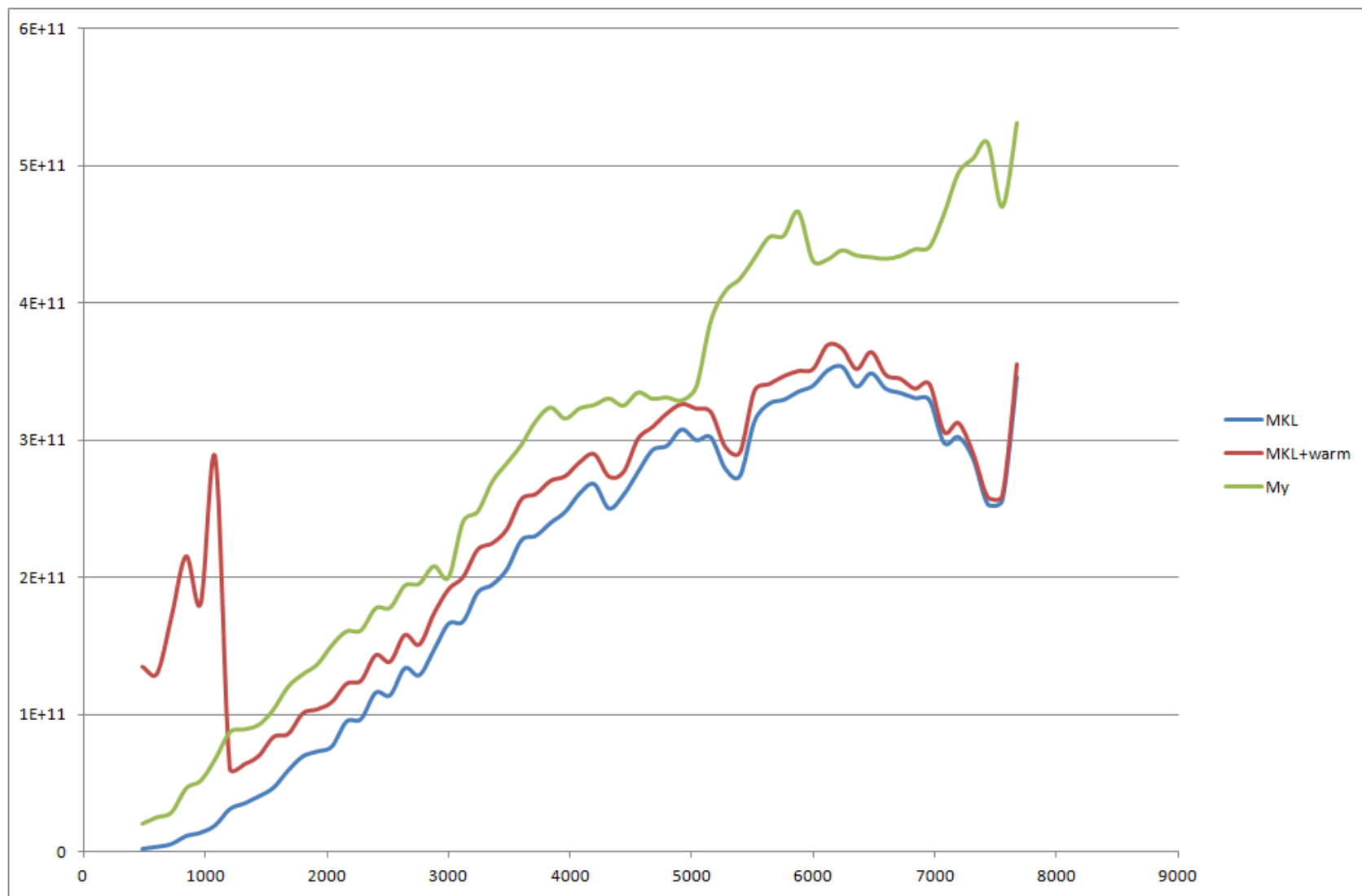
Матрица B



Параллельная реализация



Результаты



Литература

1. K. Goto and R. A. v. d. Geijn. Anatomy of highperformance matrix multiplication. ACM Trans. Math. Softw., 34(3):12:1–12:25, May 2008.
2. Design and Implementation of the Linpack Benchmark for Single and Multi-node Systems Based on Intel® Xeon Phi™ Coprocessor. 2013 IEEE 27th International Symposium on Parallel & Distributed Processing.
3. Opportunities for Parallelism in Matrix Multiplication. FLAME Working Note #71. October 23, 2013.



Авторский коллектив

- Сиднев Алексей Александрович,
ассистент кафедры
Математического обеспечения ЭВМ факультета ВМК ННГУ
sidnev@vmk.unn.ru

