

Факультет вычислительной математики и кибернетики  
Кафедра информатики и автоматизации научных исследований

Николай Старостин, к.т.н., доцент  
Маргарита Панкратова, аспирант

# Архитектурно-зависимая декомпозиция в методиках суперкомпьютерного моделирования

# Предметная область

## Инструменты инженерного анализа на супер-ЭВМ:

ЛОГОС, ЛЭГАК-ДК, ДАНКО+ГЕПАРД, НИМФА, APM WinMachine 2010, Autodesk Simulation CFD, Autodesk Simulation Mechanical, Autodesk Simulation MoldFlow

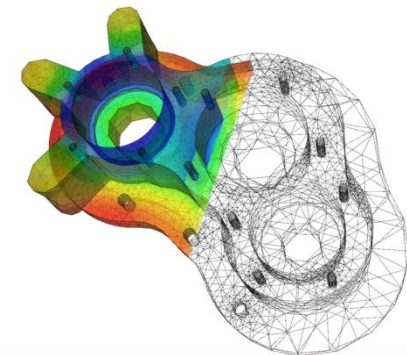
## Развитие пакетов



Многие методики численного моделирования физических процессов основываются на аппроксимационных сетках, которые моделируются с помощью графов. На практике сеточные структуры имеют большие порядки ( $10^6 - 10^9$  узлов), что приводит к необходимости их распределения по узлам параллельной вычислительной сети.

## Эффективность параллельных вычислений

в значительной мере зависит от объёма межпроцессорных коммуникаций и сбалансированности загрузки процессоров.

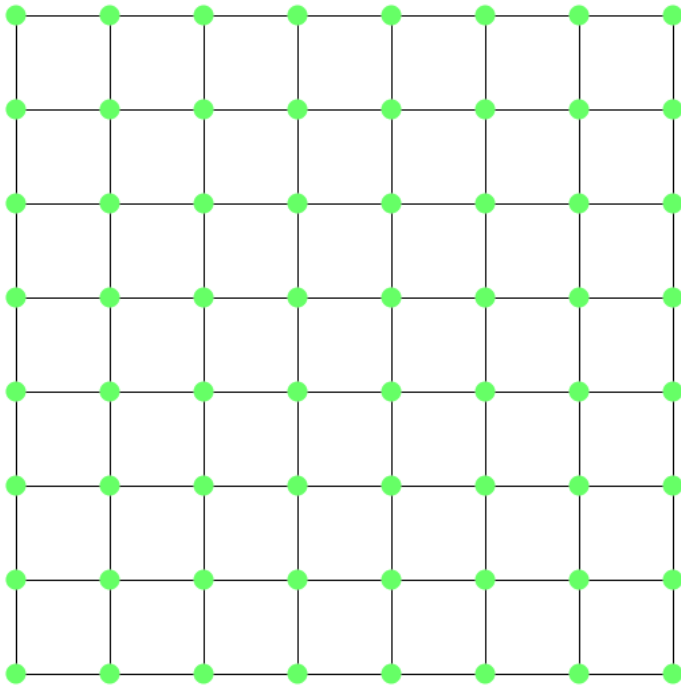


## Моделируемые процессы:

газодинамика,  
аэродинамика,  
гидродинамика,  
турбулентное  
перемешивание,  
прочность,  
разрушение,  
теплоперенос,  
многокомпонентная  
многофазная  
фильтрация и др.

# Общая постановка задачи

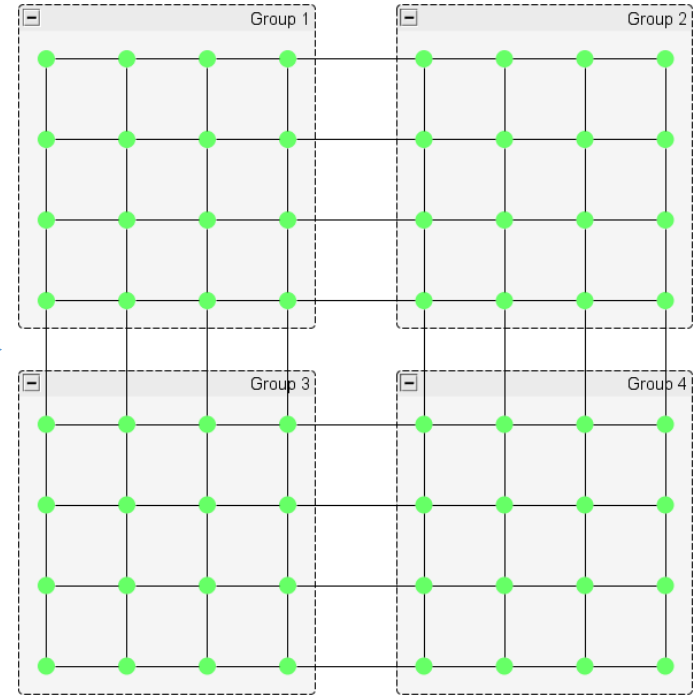
Расчетная сетка



Распределить



Вычислительная сеть,  
физическая топология



**Цель:** минимизация издержек на коммуникации.

? Как изначально разбить сетку и распределить ее по узлам вычислительной сети, чтобы минимизировать издержки на межпроцессорные коммуникации.

# Общая постановка задачи

## Исходные данные

$G(V, E, w, u)$  - расчетная сетка

$V = \{v_1, \dots, v_n\}$  – множество узлов сетки

$E \subseteq V^{(2)}$  – связи между узлами

Весы вершин  $w(v_i) \in N, i = \overline{1, n}$  - затраты на расчет

Весы ребер  $u(v_i, v_j) \in N, i, j = \overline{1, n}$  – интенсивность коммуникаций

Матрица  $T = (t_{ij})_{k \times k}$  - вычислительная сеть

$t_{ij} \in N$  - оценка затрат на передачу данных между парой вычислителей

# Общая постановка задачи

## Варьируемые параметры

Решение задачи - вектор  $x = \{x_1, \dots, x_n\}$ ,  $x_i \in \{1, \dots, k\}$ ,  $i = \overline{1, n}$

$i$ -ый элемент соответствует узлу расчетной сетки

Значение  $x_i$  - узел вычислительной сети, на который он распределен

# Общая постановка задачи

## Ограничения

$W = \sum_{i=1}^n w(v_i)$  – суммарный вес узлов сетки

$W_j = \sum_{\{i: x_i=j, i=\overline{1, n}\}} w(v_i), j = \overline{1, k}$  – загрузка (оценка относительного времени работы) вычислительного узла  $j$

$\tilde{W} = \frac{W}{k}$  – идеальная загрузка вычислительных узлов

$\varepsilon > 0$  – коэффициент отклонения от идеальной загрузки

$$\left| \frac{W_j}{\tilde{W}} - 1 \right| \leq \varepsilon, j = \overline{1, k} \quad (1)$$

- ограничения по балансировке нагрузки вычислительных узлов

# Общая постановка задачи

## Критерии

**Проблема:** очень сложно, а зачастую практически невозможно спрогнозировать: точное время каждой коммуникации; возникновение коллизий на уровне сетевых устройств; снижение пропускной способности каналов в связи с повышением интенсивности обменов; маршрутов, по которым происходит передача пакетов.

**Вывод:** нет возможности применения модели в терминах планирования коммуникационных обменов. С практической точки зрения удобнее оперировать оценками затрат на всю совокупность коммуникационных обменов.

$\beta(x, v_i, v_j) = u(v_i, v_j) \cdot t_{x_i x_j}$  - оценка затрат на коммуникационный обмен между парой узлов вычислительной сети

$$F_1(x) = \max_{(v_i, v_j) \in E} \beta(x, v_i, v_j) \rightarrow \min \quad (2)$$

$$F_2(x) = \sum_{(v_i, v_j) \in E} \beta(x, v_i, v_j) \rightarrow \min \quad (3)$$

# Проблемы организации исполнения параллельной программы

- Как правило, при организации исполнения параллельной программы заранее известна только структура расчетной сетки и общее число вычислительных узлов, которое требуется для расчета
- Структура выделенного участка вычислительной сети определяется только в момент начала расчета, когда планировщик выделяет очередной задаче требуемое число вычислителей
- Вывод: поэтому предварительная декомпозиция сеточной структуры (она необходима в силу слишком большого объема данных) осуществляется без учета информации о коммуникационной топологии
- Задача (1),(2) вырождается в задачу минимизации ширины ленты матрицы
- Задача (1),(3) – в задачу k-разбиения графа

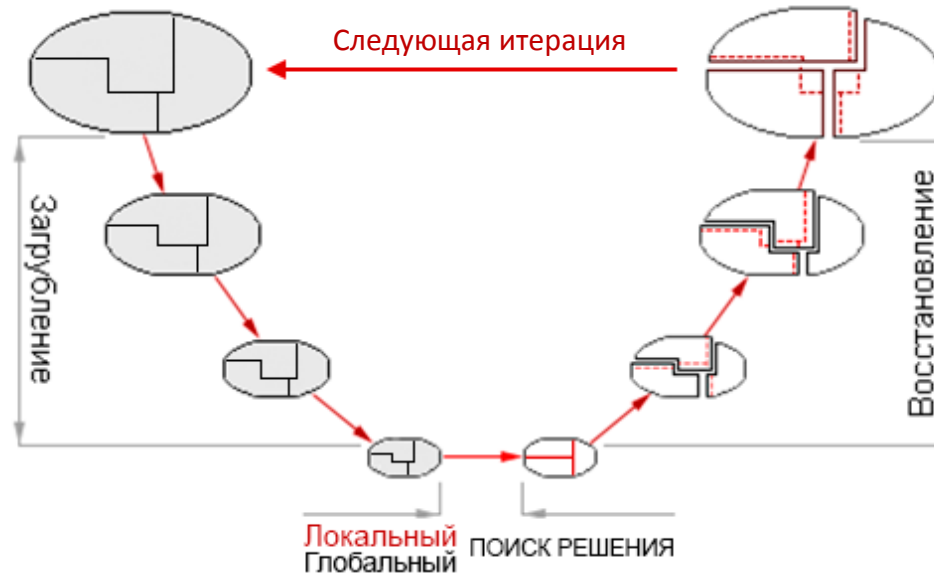
Пакеты, которые решают обозначенный класс задач:

**METIS, JOSTLE, SCOTCH, CHACO, MATRUZ**

**MATRUZ** – инструментарий для суперЭВМ, разработанных ННГУ совместно с РФЯЦ ВНИИЭФ.



# Многоуровневый итерационный метод



**Идея** использовать полученное решение на следующей итерации.

В традиционной схеме некоторое разбиение возникает только как результат работы фазы поиска. Возникают ВОПРОСЫ:

**1. Как изменится цель и принципы редукции?**

**2. Какие алгоритмы применимы на фазе поиска решений?**

Локальный поиск исследует окрестность текущего решения.

Алгоритмы глобального поиска пытаются найти новое решение.

## Фаза редукции

1. Эксплуатирует информацию о структуре найденного решения.
2. Реализует однозначную проекцию решения на грубый граф.
3. Обеспечивает принцип сравнения проекций на всех уровнях.

## Фаза поиска

1. Разные стратегии поиска решения.
2. Локальный поиск пытается улучшить проекцию решения.
3. Глобальный поиск пытается найти новое решение.

## Фаза восстановления

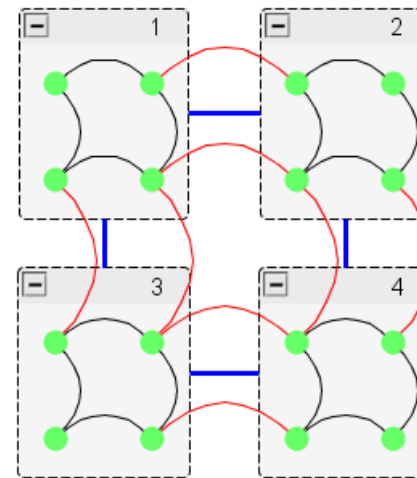
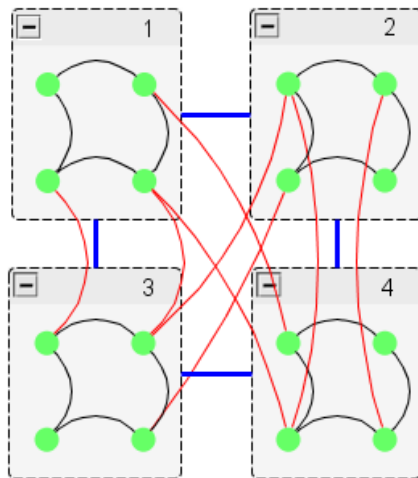
1. Проекция рекорда на исходный граф.
2. Пытается улучшить рекорд на разных уровнях редукции



# Процесс организации параллельного расчета

1. Подготовка данных
2. Декомпозиция расчетной сетки
3. Запуск параллельной программы численного расчета
  1. **Перенумерация вычислительных узлов с учетом внешних связей декомпозиции**
  2. Загрузка декомпозированной расчетной сетки
  3. Процесс численного расчета

В итоге даже для качественной декомпозиции в процессе работы параллельного расчета в выделенном сегменте физической топологии могут наблюдаться значительные потери при межпроцессорных обменах. **Подобные негативные эффекты можно минимизировать за счет удачного распределения расчетной сетки по узлам вычислительной сети.**



# Постановка задач перенумерации вычислительных узлов

## Расчетная сетка:

$G(V, E, w, u), V = \{v_1, \dots, v_k\}, E \subseteq V^{(2)}$

$w(v_i) \in N, i = \overline{1, k}$  - затраты на расчет,  $u(v_i, v_j) \in N, i, j = \overline{1, k}$   
интенсивность обменов

## Вычислительная сеть:

$T = (t_{ij})_{k \times k}, t_{ij} \in N$  - затраты на передачу данных между парой  
вычислителей

## Решение:

Перестановка  $x = \{x_1, \dots, x_k\}$

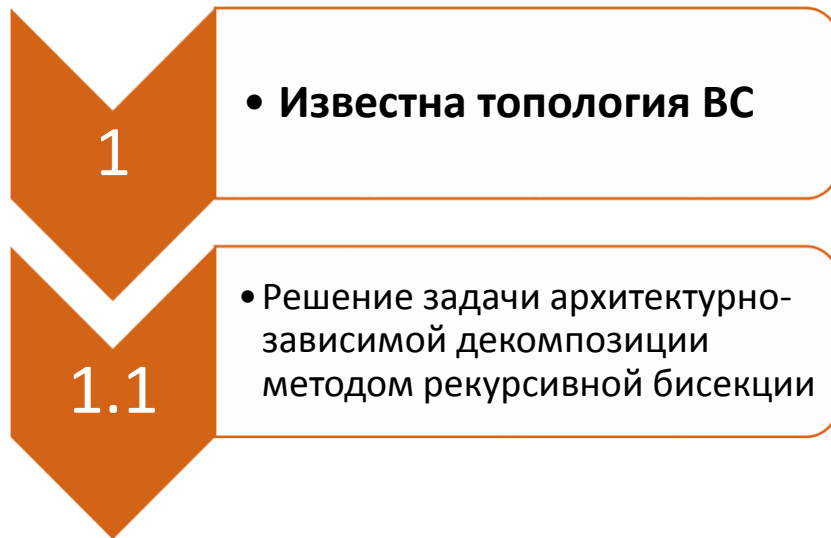
## Критерии:

$\beta(x, v_i, v_j) = u(v_i, v_j) \cdot t_{x_i x_j}$  - оценка затрат на коммуникационный обмен  
между парой узлов вычислительной сети

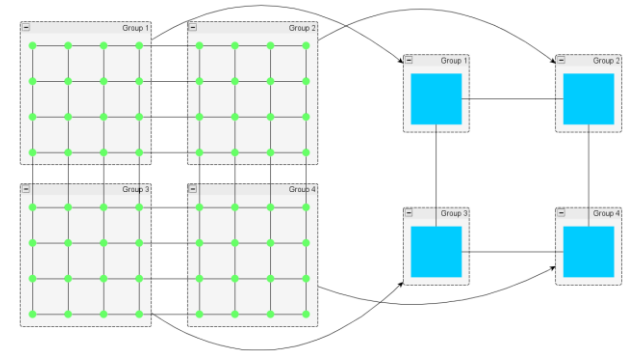
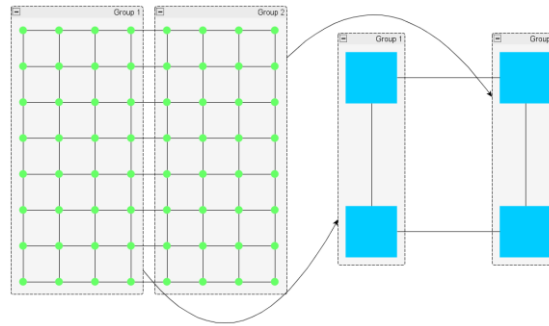
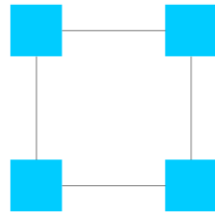
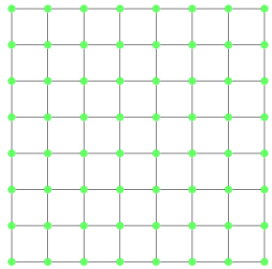
$$F_1(x) = \max_{(v_i, v_j) \in E} \beta(x, v_i, v_j) \rightarrow \min \quad (4)$$

$$F_2(x) = \sum_{(v_i, v_j) \in E} \beta(x, v_i, v_j) \rightarrow \min \quad (5)$$

# Пути решения задачи

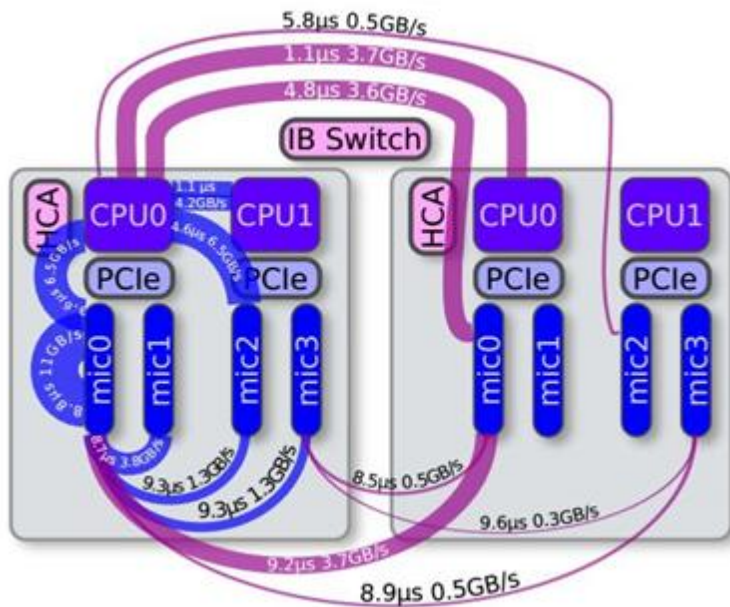


# Концепция рекурсивной бисекции

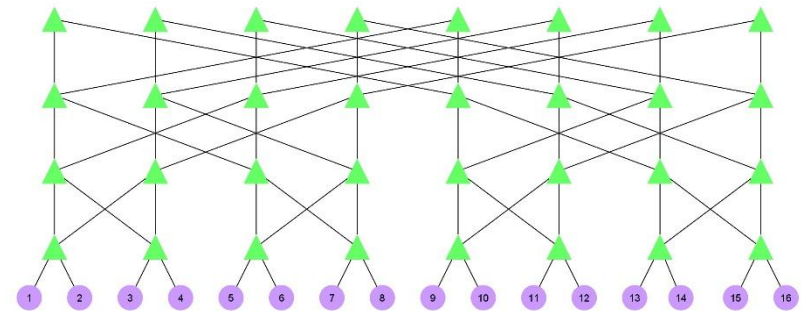


# Топология тестовых вычислительных сетей

Symmetric



Fat tree





# Вычислительный эксперимент

<http://staffweb.cms.gre.ac.uk/~wc06/partition>

<http://www.cise.ufl.edu/research/sparse/matrices/>

Расчетная сетка	Размер	BC	Размер	1-этапный		2-этапный	
				F <sub>1</sub>	F <sub>2</sub>	F <sub>1</sub>	F <sub>2</sub>
gr_30_30	900	Fattree	64	8	6734	8	6467
gr_30_30	900	Symmetric	288	430	619860	430	479564
ef_body	45087	Fat tree	64	8	29986	8	32671
ef_body	45087	Symmetric	288	430	3248324	430	3224784
fe_tooth	78136	Fat tree	64	8	218430	8	210407
fe_tooth	78136	Symmetric	288	430	18033219	430	16097627
fe_rotor	99617	Fat tree	64	8	290339	8	261618
fe_rotor	99617	Symmetric	288	430	24332380	430	22415356

**Вывод:** 2-этапный алгоритм показывает лучшие результаты.

1. Рекурсивная бисекция – эвристический алгоритм
2. Качественная предварительная декомпозиция расчетной сетки минимизирует внешние связи, меньше связей проявляется в виде межпроцессорных коммуникаций