

Масштабируемая массивно-параллельная реализация схемы переноса примесей MPDATA на распределенных GPU/SMP системах

Левченко Олег – Петров Артем – Микушин Дмитрий

Геофизический Центр РАН – СибНИГМИ РАН – ИВМ РАН

октябрь, 2010

- Физическая задача
- Выбор численного метода
- Численная схема MPDATA
- Проект параллельной реализации MPDATA для гибридных суперкомпьютеров
 - Общая схема вычислений
 - Архитектура GPU GT200
 - Анализ реализации для GPU
 - Метод декомпозиции прямоугольной сетки
- Заключение

Перенос примеси как физическая задача



- Тип источника (точечный, распределённый)
- Тип аэрозоля (пепел, песок, соль, ...)
- Распределение размеров частиц (мониторинг концентрации опасной мелкодисперсной пыли)
- Динамика: эмиссия (+ сальтация), перенос (+ химические процессы), осаждение

Модели и численные схемы

- Физические уравнения модели прогноза погоды WRF v.3 [1]
- Численные схемы адвекции
 - Ориентированы на пассивные примеси
 - Сохранение массы
 - Эйлеровы модели
 - Ограничения на временной шаг
 - Полу-Лагранжевые модели:
 - Большие временные шаги без потери стабильности
 - Порядок точности
 - Сохранение массы/положительно-определенность
 - Overshooting/undershooting

$$\partial_t U + (\nabla \cdot \mathbf{V}u) - \partial_x(p\partial_\eta\phi) + \partial_\eta(p\partial_x\phi) = F_U \quad (2.3)$$

$$\partial_t V + (\nabla \cdot \mathbf{V}v) - \partial_y(p\partial_\eta\phi) + \partial_\eta(p\partial_y\phi) = F_V \quad (2.4)$$

$$\partial_t W + (\nabla \cdot \mathbf{V}w) - g(\partial_\eta p - \mu) = F_W \quad (2.5)$$

$$\partial_t \Theta + (\nabla \cdot \mathbf{V}\theta) = F_\Theta \quad (2.6)$$

$$\partial_t \mu + (\nabla \cdot \mathbf{V}) = 0 \quad (2.7)$$

$$\partial_t \phi + \mu^{-1}[(\mathbf{V} \cdot \nabla \phi) - gW] = 0 \quad (2.8)$$

Характерные особенности MPDATA

Специально разработана под метеорологические нужды:

- Положительная определенность
- Сохранение массы
- Вычислительная простота, сравнимая с донорной схемой

Расширения MPDATA

Численная модель вычислений MPDATA хорошо теоретически проработана и включает следующие расширения [6]:

- Возможность увеличить пространственную точность до третьего порядка, временную точность до пятого
- Возможность эмпирической подстройки для увеличения точности без увеличения вычислительной сложности
- Физическая диффузия в составе адвективного потока
- Транспорт скалярного поля с переменным знаком
- Моделирование дивергентных полей

Почему MPDATA?

- Численная модель вычислений MPDATA потенциально может быть оптимизирована под CUDA:
 - Соотношение доступов в память к вычислениям
 - Явная модель вычислений
- Является частью численных гидродинамических моделей погоды EULAG и NH3D

Одномерный случай - формулировка адвекции[5]

1 $\frac{\partial \Psi}{\partial t} + \frac{\partial(\Psi \cdot v)}{\partial x} = 0$, где Ψ - концентрация, v - скорость ветра

2
$$\Psi_i^{n+1} = \Psi_i^n + \overbrace{[F(\Psi_{i-1}^n, \Psi_i^n, U_{i-1/2}^{n+1/2})]}^{\text{входной поток}} - \overbrace{[F(\Psi_i^n, \Psi_{i+1}^n, U_{i+1/2}^{n+1/2})]}^{\text{выходной поток}}$$

3 $F(\Psi_L, \Psi_R, U) = 0.5 \cdot (U + |U|) \cdot \Psi_L + 0.5 \cdot (U - |U|) \cdot \Psi_R$, где $U = v \cdot \frac{\Delta t}{\Delta x}$

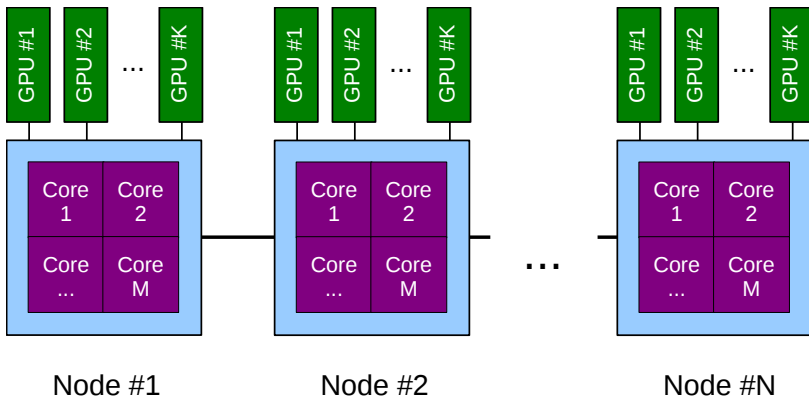
4 $\max_{i,n} |v_{i \pm 1/2}| \cdot \frac{\Delta t}{\Delta x} \leq 1$ - условие стабильности схемы

5 $O(\Delta x)$ - недостаточно для практических применений

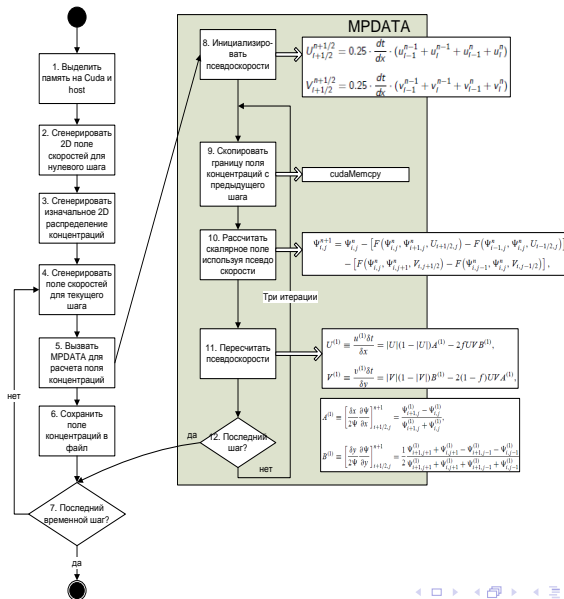
Одномерный случай - MPDATA[5]

- 1 $\frac{\partial \Psi_i^n}{\partial t} + \frac{\partial(\Psi_i^n \cdot v)}{\partial x} = \frac{\partial[\frac{1}{2} \cdot (|v| \cdot \Delta x - v^2 \cdot \Delta t) \cdot \frac{\partial \Psi_i^n}{\partial x}]}{\partial x}$ вместо $\frac{\partial \Psi_i^n}{\partial t} + \frac{\partial(\Psi_i^n \cdot v)}{\partial x} = 0$
- 2 $\frac{\partial \Psi}{\partial t} + \frac{\partial(\Psi \cdot v)}{\partial x} = \frac{\partial(K \cdot \frac{\partial \Psi}{\partial x})}{\partial x}$ где $K = \frac{1}{2} \cdot (|v| \cdot \Delta x - \Delta t \cdot v^2)$ - коэффициент диффузии
- 3 $\frac{\partial \Psi}{\partial t} - \frac{\partial(K \cdot \frac{\partial \Psi}{\partial x})}{\partial x} = 0$ - "скрытое" уравнение адвекции
- 4 $\frac{\partial \Psi}{\partial t} + \frac{\partial(\Psi \cdot v_{diff})}{\partial x} = 0$ где $v_{diff} = -\frac{K}{\Psi} \cdot \frac{\partial \Psi}{\partial x}$
- 5 $u_{antidiff} = -u_{diff}$ - "проигрыш" диффузии назад во времени
- 6 $\Psi_i^* = \Psi_i^n + [F(\Psi_{i-1}^n, \Psi_i^n, U_{i-1/2}^{n+1/2}) - F(\Psi_i^n, \Psi_{i+1}^n, U_{i+1/2}^{n+1/2})]$
- 7 $\Psi_i^{n+1} = \Psi_i^* + [F(\Psi_{i-1}^*, \Psi_i^*, v_{i-1/2}^{antidiff}) - F(\Psi_i^*, \Psi_{i+1}^*, v_{i+1/2}^{antidiff})]$ где $v_{i+1/2}^{antidiff} = \frac{(|v_{i+1/2}| \cdot \Delta x - \Delta t \cdot v_{i+1/2}^2) \cdot (\Psi_{i+1}^* - \Psi_i^*)}{(\Psi_i^* + \Psi_{i+1}^* + \epsilon) \cdot \Delta x}$
- 8 $O(\Delta x^2)$ - достаточно для практических применений

Целевая гибридная вычислительная система

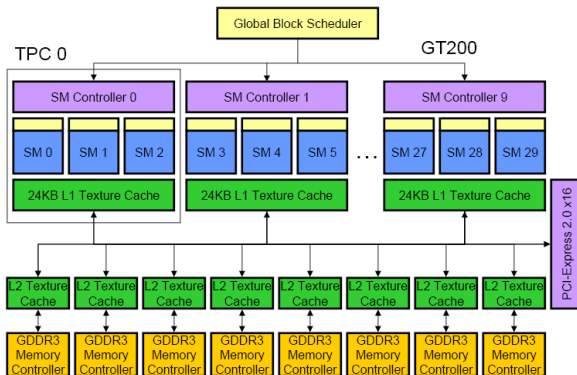


Структура GPU-CPU программного комплекса[4]



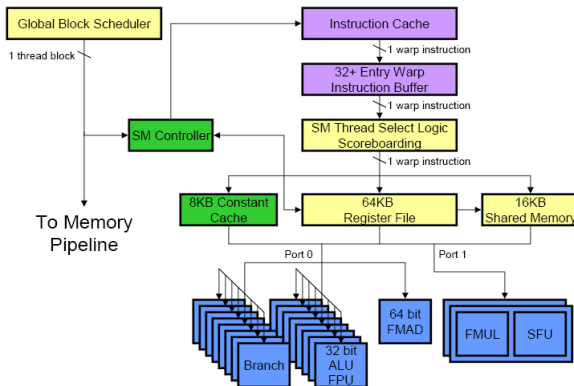
GPGPU - Кластер из мультипроцессоров[2]

- Асинхронные вычислительные мультипроцессоры
- Общая многопортовая память
- Конвейеры текстурной памяти



Мультипроцессор GT 200[2]

- WARP и `__syncthreads()`
- Треугольник "регистры - разделяемая память - потоки"
- Специфичные оптимизации по паттернам доступа к памяти

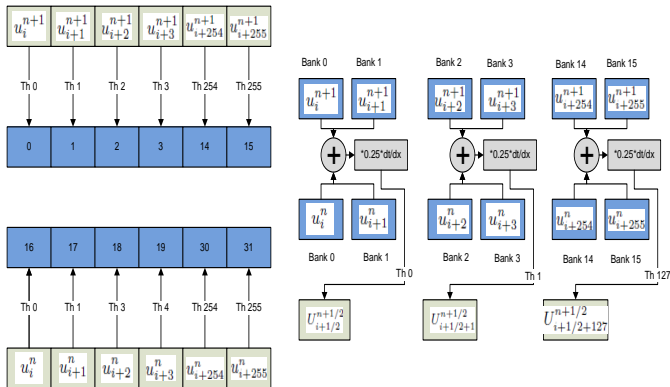


Dual Issue[2]

- Двойное опережение
- Зависимости между 64-bit fpu и 32-bit fpu



Инициализация псевдоскоростей на CUDA



$$U_{i+1/2}^{n+1/2} = 0.25 \cdot \frac{dt}{dx} \cdot (u_{i-1}^{n-1} + u_i^{n-1} + u_{i-1}^n + u_i^n)$$

- Склейка по доступу в глобальную память (Coalescing)
- Эффективная блокировка варпов (Warps Stall)
- Возможен двухкратный конфликт по банкам разд. памяти

Анализ производительности мультипроцессора[3]

1 Загрузка $v_{i+0+..255}^{n+1}$

- $N_{threads} = 256 \Rightarrow N_{warps} = \frac{256}{32} = 8$
- $ShMemSize_1 = N_{threads} \cdot 2 \cdot 4 = 2048$ Bytes

2 Загрузка $v_{i+0+..255}^n$

- $N_{threads} = 256 \Rightarrow N_{warps} = \frac{256}{32} = 8$
- $ShMemSize_2 = N_{threads} \cdot 2 \cdot 4 = 2048$ Bytes

3 Вычисление $U_{i+1/2+..127}^{n+1/2}$

- $N_{threads} = \frac{256}{2} = 128 \Rightarrow N_{warps} = \frac{128}{32} = 4$

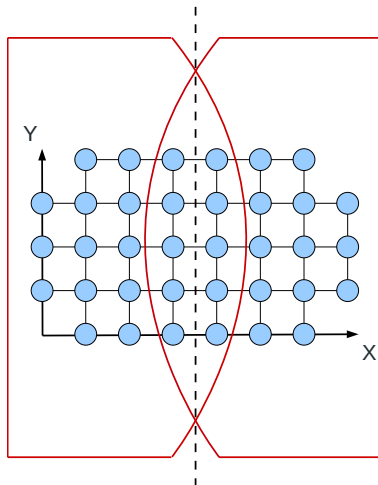
4 Анализ загрузки мультипроцессора

- Максимальное количество блоков на мультипроцессоре по потокам $MaxBlocks_{threads} = \frac{1024}{256} = 4$
- Максимальное количество блоков на мультипроцессоре по разделяемой памяти

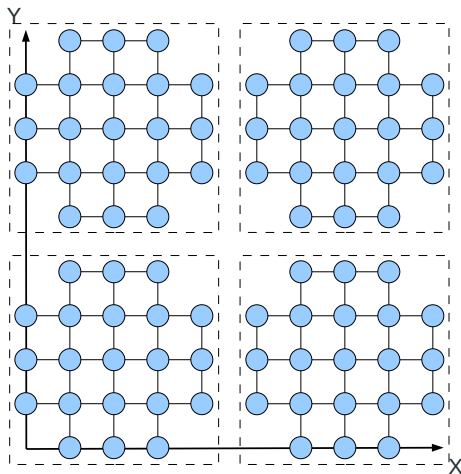
$$MaxBlocks_{shmem} = \frac{MaxShmem}{ShMemSize_1 + ShMemSize_2} = \frac{32768}{2048 + 2048} = 8$$

5 Вывод - так как алгоритм memory-bound, то фактором, сдерживающим производительность мультипроцессора, может стать нехватка регистров, если $N_{regs} > \frac{16384}{1024} = 16$

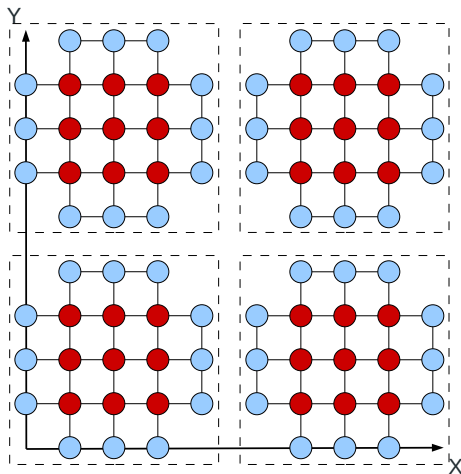
Дублирование границ при декомпозиции



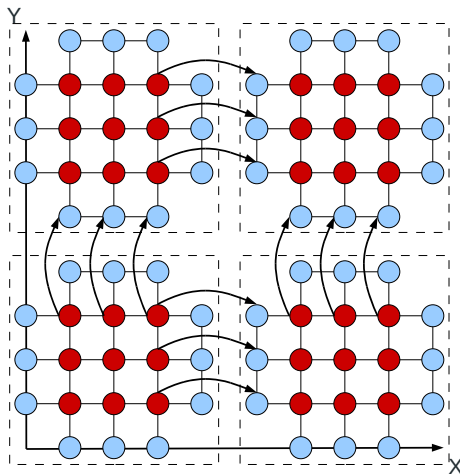
Данные в распределённой памяти



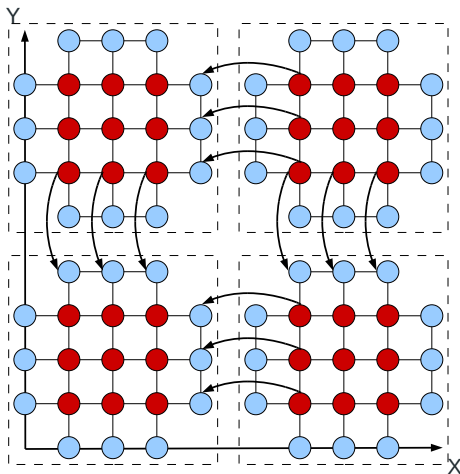
Параллельный расчёт внутренних узлов



Синхронизация граничных узлов



Синхронизация граничных узлов



Заключение

- Изучена предметная область численных решений уравнения адвекции
- Протестирована последовательная версия алгоритма MPDATA
- Начата разработка CUDA реализации алгоритма MPDATA
- Реализована проблемно-независимая распределенная среда обмена границами



A Description of the Advanced Research WRF Version 3.

http://www.mmm.ucar.edu/wrf/users/docs/arw_v3.pdf.



D. Kanter.

Nvidia's gt200: Inside a parallel processor.

<http://www.realworldtech.com/page.cfm?ArticleID=RWT090808195242&p=8>.



D. B. Kirk and W. mei W. Hwu.

Programming Massively Parallel Processors.

Morgan Kaufmann, 2010.



D. N. Mikushin.

Численное моделирование мезомасштабного переноса примеси над гидрологически неоднородной поверхностью.



P. K. Smolarkiewicz.

A simple positive definite advection scheme with small implicit diffusion.



P. K. Smolarkiewicz and L. G. Margolin.

Mpdata: A finite-difference solver for geophysical flows.