

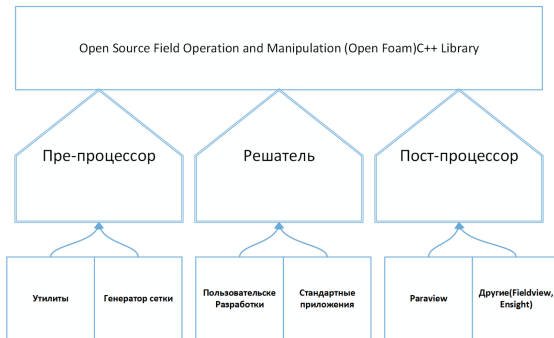
Статический анализ кода пакетов прикладных программ для поиска протяженных участков для замены на параллельный эквивалент

М.В. Чернова¹, А.Н. Сальников^{1,2}

ВМК МГУ имени М.В. Ломоносова¹,
ФИЦ ИУ РАН²,

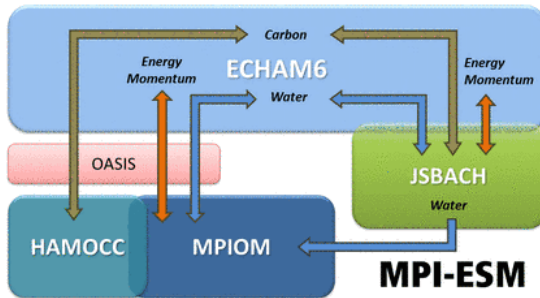
- В научной среде существует множество пакетов прикладных программ (ППП), активно использующихся в решении различных задач.
- Ежегодная разработка новых программных продуктов вызывает устаревание некоторых программных решений, в частности отдельных частей в пакетах прикладных программ.
- Возникает вопрос о повышении эффективности работы такого кода, а также о различных его модификациях с минимизацией затрат человеческого труда.

- OpenFOAM — открытое программное обеспечение для численного моделирования задач механики сплошных сред. В основе кода лежит набор библиотек, предоставляющих инструменты для решения систем дифференциальных уравнений в частных производных.



- Gromacs — (groningen machine for chemical simulations — гронингенская машина для химического моделирования) — пакет программ для моделирования физико-химических процессов в молекулярной динамике, специализирующийся на моделировании биомолекул (например, молекул белков и липидов), имеющих много связанных взаимодействий между атомами.

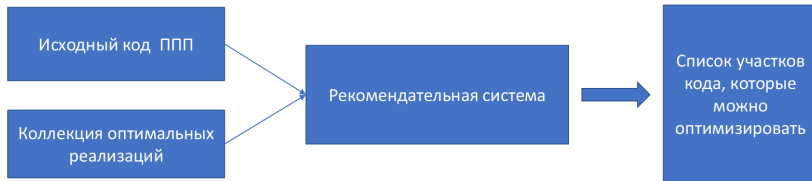
- MPIESM(ECHAM6/MPIOM/JSBACH/HAMOCC) — современная крупномасштабная совместная модель «океан-земля-атмосфера» европейского метеорологического сообщества.



Характеристики пакетов прикладных программ

	OpenFOAM	Gromacs	MPI-ESM
Общее количество файлов	15,639	4,563	3,990
Общее количество строк	4,517,786	3,006,701	4,034,054
Количество строк кода в файлах типа .c, .f*	1,626,725	2,246,960	1,018,752

- Создание системы выдачи рекомендаций пользователю о возможности замены некоторых фрагментов в пакете прикладных программ с целью дальнейшей оптимизации



- Создание коллекции реализаций часто используемых в ППП численных методов
- Создание метода поиска в исходном коде фрагментов, функционально эквивалентных заданному
- Создание коллекции целевых фрагментов с учетом различных характеристик «оптимальности». Под «оптимальностью» будем иметь в виду реализацию в соответствии с запросами пользователя (например, для конкретного окружения)

Цели и трудности решения

Цели:

- Найти в исходном коде ППП расположение фрагментов кода, функционально эквивалентных данному.
- Иметь возможность находить не только фрагменты, идентичные искомому, но также имеющие расхождения по различному ряду признаков, таких как: имена переменных/функций, порядок расположения операторов, наличие или отсутствие некоторых операторов.

Трудности решения:

- Большой объем исходного кода, в котором должен производиться поиск (следовательно, трудность построения графа потока управления).
- Нет однозначности при определении функциональной эквивалентности двух участков кода.

В качестве искомых мотивов будут использованы реализации математических методов, однако сам алгоритм может быть использован и для других методов.

В качестве реализаций математических методов было принято решение использовать следующие библиотеки:

- ATLAS(Automatically Tuned Linear Algebra Software) – реализация BLAS для языков C и Fortran. Средний объем методов 268 строк.
- OpenCV (Open Source Computer Vision Library) – библиотека алгоритмов компьютерного зрения, обработки изображений и численных алгоритмов общего назначения для языка C. Средний объем методов 297 строк.

Построение решения для поиска схожих фрагментов

- Создание списка зависимостей, требуемых мотивом.
- Обнаружение рекурсий.
- Генерация нового файла, включающего все необходимые вызовы процедур.
- Удаление незначащих фрагментов (комментарии, невызываемые процедуры).

Предобработка исходного кода пакета прикладных программ

- Анализ пакета: обнаружение и составление списка файлов реализации на конкретном языке программирования.
- Кластеризация файлов: отобранные файлы разбиваются на множество кластеров по наличию зависимостей для поэтапного сравнения с шаблоном.
- Удаление комментариев: до кластеризации или после, отдельно для каждого кластера.

Создание комплексной метрики для определения близости фрагментов

- 1 Фрагментом кода в заданном файле f будем считать последовательность лексем L_i , однозначно определяемую парой $\langle l_{begin}, l_{end} \rangle$, где l_{begin} и l_{end} – номера строк начала и конца фрагмента соответственно.
- 2 $\mathcal{P} = \{\mathcal{P}_i | 1 \leq i \leq N_p\}$ – множество исходного кода в пакете прикладной программы,
- 3 $M = \{M_j | 1 \leq j \leq N_m\}$ – множество реализаций из коллекции искомых фрагментов.
- 4 Тогда для двух фрагментов $L_{\mathcal{P}_i}$ и L_{M_j} зададим функцию $\rho(L_{\mathcal{P}_i}, L_{M_j})$ – расстояние между ними.
- 5 Фрагмент $L_{\mathcal{P}_i}$ будем называть дубликатом фрагмента L_{M_j} в отношении φ , если $\rho(L_{\mathcal{P}_i}, L_{M_j}) \leq \varphi \leq 1$.

Реализация алгоритма с использованием выбранных метрик

- На данном этапе был выбран способ вычисления метрики схожести на основании абстрактных синтаксических деревьев¹:

$$S(L_{\mathcal{P}_i}, M_j) = \frac{2N}{2N + L + R}, \quad (1)$$

где N — число совпадающих узлов, L — число различных узлов в первом поддереве, R — число различных узлов во втором поддереве.

- Построение деревьев производилось с помощью генератора парсеров ANTLR.

¹Baxter I. D. et al. Clone detection using abstract syntax trees

Результаты работы предобработчика

- При кластеризации использовалась глубина рекурсии, равная 3.

	OpenFOAM	Gromacs	MPI-ESM
Время анализа файлов	3.06	0.815	0.862
Количество обработанных файлов	8336	3206	1884
Время кластеризации	203.475	93.467	53.48
Количество кластеров	3643	2484	1167
Время удаления комментариев	257.103	121.467	85.225
Количество удаленных строк	15975	5775	5280

Таблица: Предобработка ППП (время в секундах)

- Описанный алгоритм поиска фрагментов был применен на пакете MPI-ESM. В качестве базы, на основании которой производился поиск, была создана коллекция реализаций численных методов, наиболее часто встречающихся в пакете
- Далее представлены результаты поиска при пороговом значении метрики 0.7. При увеличении метрики увеличиваются требования к идентичности фрагментов, и количество найденных фрагментов стремится к нулю.

Количество найденных функций в пакете MPI-ESM

	Метод Ньютона	Метод Гаусса	Метод релаксации	Метод простой итерации	Метод Ричардсона
Вызовы функций	5	11	2	6	5
Схожие реализации	12	28	15	8	3

- Создание базы шаблонов параллельных реализаций некоторых методов, которые будут предложены пользователю в качестве замены
- Реализация алгоритма выдачи информации о всех найденных схожих шаблонах и рекомендаций о замене и подстановке оптимальных эквивалентов
- Рассмотреть возможность интеграции такой системы в системы автоматического распараллеливания кода

Спасибо за внимание!

- *Дубликаты* – два функционально эквивалентных фрагмента кода.
- Под *методом* будет подразумеваться один из алгоритмических способов решения конкретной задачи.
- Каждый метод может иметь несколько *мотивов* – конкретных вариантов реализации, которые разрабатываются в зависимости от возможностей вычислительной системы и алгоритма, используемого разработчиком.
- *Код пакета прикладных программ* (код ППП) – реализация пакета прикладных программ, в которой необходимо обнаружить дубликат искомого мотива.
- *Шаблон* – искомый мотив, по которому непосредственно будет осуществляться поиск в коде ППП.