

Сравнение вариантов алгоритма Борувки в задаче построения минимального остовного дерева*

Зайцев Вадим, Новосибирский государственный университет, zaic101@gmail.com

Калгин Константин, ИВМиМГ СО РАН, kalginkv@gmail.com

*Работа поддержана Грантом Президента РФ для молодых учёных МК-3644.2014.9

Введение

Задача обработки и анализа больших данных стоит повсеместно и возникает во множестве областей: известные поисковики, провайдеры, социальные сети, сотовые операторы, такие области науки как биоинформатика (обработка генов и генных сетей) [1] и многие другие. Объёмы обрабатываемых данных составляют от десятков гигабайт до терабайт и выше в зависимости от области. Такие данные зачастую представимы в виде графов, где вся необходимая информация хранится в вершинах и рёбрах или дугах графов и анализ таких данных сводится к классическим алгоритмам на графах. Большие объёмы данных неизбежно ведут к тому, что размеры получаемых графов так же достаточно велики. В ряде работ исследуется проблематика получения эффективных параллельных реализаций алгоритмов на графах [2, 3, 4].

В данной работе рассматривается задача построения минимального остовного дерева (minimum spanning tree — MST) на мультипроцессоре с ccNUMA архитектурой. В теории графов остовным деревом называется ациклический связный подграф данного связного неориентированного графа, в который входят все его вершины. Минимальным остовным деревом называется остовное дерево во взвешенном графе с минимальной суммой весов рёбер. В случае, когда данный граф не является связным, для каждой компоненты будет построено своё остовное дерево и полученные деревья называются остовным лесом.

Минимальные остовные деревья имеют широкое практическое применение. Они используются при построении различных сетей: коммуникационных линий, компьютерных сетей, транспортных сетей [5]. Так же они используются в задачах кластеризации [6], сегментации при обработке изображений [7], распознавании рукописного ввода [8]. Минимальные остовные деревья используются как приближённое решение переборной задачи: например, задачи Штейнера — построение кратчайшей сети, соединяющей заданный конечный набор точек на плоскости [9].

1 Обзор алгоритмов построения MST

Существуют различные последовательные алгоритмы построения MST, наиболее распространёнными из которых являются алгоритмы Крускала, Прима и Борувки [10]. Не каждый из алгоритмов в его оригинальном варианте возможно эффективно распараллелить, однако, для некоторых из них существуют модификации, за счёт которых использование нескольких ядер одного узла кластера и даже использование нескольких узлов кластера может оказаться оправданным [11, 12].

Алгоритм Крускала строит остовное дерево в два этапа: сначала все рёбра графа сортируются в порядке неубывания веса (рёбра одинакового веса могут идти в произвольном порядке), затем рёбра последовательно обходятся и в остовное дерево добавляются те из них, которые не создают цикла.

Сортировка рёбер в общем случае может быть выполнена за время $O(E \cdot \log(E))$, проверка образования цикла для каждого ребра может быть реализована с помощью системы непересекающихся множеств и выполняться за время $O(\alpha(E))$, где $\alpha(E)$ — обратная функция Аккермана, значения которой малы и, например, для тестируемых RMat и SSCA2 графов, в которых количество рёбер составляет от 10 до 100 миллионов, не превосходят 5. Итоговая оценка времени работы получается $O(E \cdot \log(E) + E \cdot \alpha(E))$.

Алгоритм Прима начинает построение остовного дерева с произвольно выбранной вершины. Затем, пока дерево не будет построено полностью, выбирается ребро минимального веса, соединяющее вершину уже построенной части дерева с вершиной не из дерева, и выбранное ребро добавляется в строящееся остовное дерево.

Число шагов в алгоритме Прима равно числу вершин в графе. Время поиска ребра минимального веса на каждом шаге зависит от выбранной структуры данных: при использовании матрицы смежности для поиска минимального ребра потребуется просмотреть всю строку матрицы за линейное время, а при использовании кучи добавление и удаление каждого ребра будет происходить за логарифмическое время. Таким образом, алгоритм Прима возможно реализовать со временем работы $O(V^2)$, что является лучшей оценкой для полных графов, либо со временем $O(E \cdot \log(V))$, что более предпочтительно для разреженных графов.

Алгоритм Борувки. Это итерационный алгоритм: изначально каждая вершина графа считается отдельной компонентой, затем на каждой итерации для каждой компоненты находится инцидентное ей ребро, соединяющее её с другой компонентой, минимального веса и по найденным рёбрам компоненты объединяются. Алгоритм работает до тех пор, пока не останется одна компонента и пока между имеющимися компонентами есть рёбра.

На каждой итерации алгоритма множество компонент уменьшается не менее, чем в два раза, следовательно количество итераций можно оценить сверху как $\lceil \log_2(V) \rceil$. На каждой итерации в худшем случае необходимо просмотреть все рёбра, при этом процедура обработки одного ребра достаточно простая (определение каким компонентам принадлежат инцидентные вершины и, возможно, обновление рекорда для данных

компонент). Итоговая временная сложность получается $O(E \cdot \log(V))$.

1.1 Параллельные алгоритмы

С точки зрения параллелизма про алгоритмы Крускала и Прима можно сказать следующее:

- они состоят из большого количества шагов, которые зависимы между собой и не могут исполняться параллельно
- вычислительная сложность одного шага достаточно маленькая и её распараллеливание будет неэффективно

Это делает распараллеливание данных алгоритмов в оригинальном варианте бесперспективным.

Исключением является реализация алгоритма Прима для плотных графов с обновлением матрицы смежности вместо использования кучи: в этом случае сложность одного шага составляет $O(V)$ и его возможно распараллелить, разделив матрицу между потоками. Однако, не для плотных графов такой подход асимптотически хуже других алгоритмов и при распараллеливании даже на большое число вычислителей будет медленнее последовательных реализаций других алгоритмов на больших графах.

В то же время алгоритм Борувки обладает противоположными свойствами:

- количество итераций небольшое
- вычислительная сложность одной итерации большая и может быть распределена между потоками

Это делает алгоритм Борувки наиболее пригодным для распараллеливания и в перспективе может дать хорошую масштабируемость.

В статьях [13, 11, 14] представлены различные варианты параллельных реализаций алгоритма Борувки. Разные реализации используют разные структуры данных для хранения рёбер (хранение рёбер одним большим списком или хранение для каждой вершины списка инцидентных ей рёбер) и разные способы их объединения (слияние и копирование списков, объединение списков за константное время и т. д.).

В статье [11] представлен параллельный алгоритм построения MST на основе последовательного алгоритма Прима. Суть данного алгоритма заключается в том, что каждый поток выбирает случайную вершину в графе и начинает “растить” дерево из выбранной вершины, пока не попытается присоединить вершину, уже принадлежащую другому дереву, после чего два дерева встретившихся потоков объединяются, один из них продолжает “растить” объединённое дерево, а второй заново выбирает вершину графа. К преимуществам данного подхода относится отсутствие барьерной синхронизации, которая требуется в алгоритме Борувки после каждой итерации и, в зависимости от реализации, между разными шагами одной итерации. Другим преимуществом перед алгоритмом Борувки является то, что каждое ребро будет просмотрено ровно один раз.

2 Варианты алгоритма Борувки

В данных работах [13, 11, 14] результаты производительности показаны либо для старых архитектур [13, 14], либо масштабируемость была ограничена шестью одноядерными процессорами (SMP UMA, UltraSPARC II) [11].

Настоящая работа направлена на адаптацию существующих и реализацию новых алгоритмов, ориентированных на высокую эффективность на больших графах на современных вычислителях с общей памятью с ccNUMA архитектурой. Исследование началось с реализаций, предложенных в статье [11], а также их модификаций. Затем, в разделе 2.3 на основе сделанных выводов будет предложен собственный подход к хранению графа в алгоритме Борувки.

В частности, далее будут представлены результаты тестирования описанных реализаций на двухсокетной системе $2 \times$ Intel Xeon CPU E5-2690 (8 ядер 2.9 GHz, 32KB L1 cache, 256KB L2 cache, 20MB L3 cache). Производительность полученных реализаций будет показана в MTEPS — величине, обратно пропорциональной времени. Более подробно об используемом окружении и замерах времени описано в разделе 4.

2.1 Списки смежности (Adjacency Lists — AL)

Первый вариант алгоритма Борувки для хранения графа использует списки смежности для каждой компоненты. Изначально каждая компонента представлена одной вершиной и её список смежности совпадает со списком смежности вершины, затем после каждой итерации для каждой образовавшейся компоненты явно строится список всех рёбер, инцидентных данной компоненте. В качестве списков может быть использована любая структура данных, позволяющая последовательно обходить элементы. В данном подходе использовались динамически выделяемые массивы.

Изначально множество вершин графа разделяется между потоками равномерно по количеству инцидентных вершинам рёбер. Далее на каждой итерации алгоритма Борувки выполняются следующие шаги, пока не останется одна компонента или между несколькими оставшимися компонентами не будет рёбер:

1. **Минимальное инцидентное ребро.** Каждый поток для каждой своей компоненты находит инцидентное ей ребро минимального веса. Это делается обходом всех инцидентных компоненте рёбер.
2. **Объединение деревьев.** Компоненты объединяются по найденным рёбрам: выделяются образовавшиеся компоненты, для каждой такой компоненты определяется номер на следующей итерации. В качестве номера компоненты используется номер одной из вершин, входящих в данную компоненту. Так же решается проблема с возможными циклами (когда более двух компонент по кругу выберут следующую компоненту в качестве ближайшей), которые возможны в алгоритме Борувки при наличии в графе рёбер одинакового веса.

3. **Перенумерация компонент.** Осуществляется перенумерация компонент: для каждой вершины вычисляется номер компоненты, в которую она входит, используя параллельный алгоритм Pointer Jumping [15].
4. **Слияние списков.** Происходит объединение компонент путём слияния списка рёбер. В модификации с неотсортированными списками рёбер слияние делается последовательным копированием одного списка в другой. При использовании отсортированных по весу списков вместо копирования и последующей сортировки рёбер происходит слияние с использованием приёма, который используется в сортировке слиянием: из голов копируемых списков выбирается минимальный элемент, извлекается и добавляется в создаваемый список.

Были опробованы два подхода относительно сортированности списков смежности:

- **Отсортированные по весу списки смежности.** Отсортированные списки позволяют осуществлять шаг **минимальное инцидентное ребро** за константное время, просмотрев лишь голову списка, однако поддержание таких списков сортированными на шаге **слияния списков** требует дополнительного времени.
- **Неотсортированные списки смежности.** Когда поддержание списков отсортированными не требуется, их объединение возможно осуществлять простым копированием памяти, но поиск по неотсортированному списку возможен только за линейное время.

На практике подход с сортировкой списков по весу оказался заметно медленнее. На рисунке 1 показано сравнение производительности обоих подходов на RMat и SSCA2 графах: реализация со слиянием отсортированных списков обозначена как AL-Merge, реализация с копированием неотсортированных списков — AL-Copy. Подробное описание графов, на которых проводилось тестирование, приведено в разделе 4.

В подходе со списками смежности имеется две существенные проблемы. Первая из них заключается в том, что на шаге **слияния списков** происходит копирования большого объёма данных и этот шаг выполняется долго. Вторая проблема лежит в том, что при распределении компонент между потоками на шагах **минимальное инцидентное ребро** и **слияния списков** происходит дисбаланс нагрузки на последних итерациях, когда количество оставшихся компонент становится меньше числа потоков — в этом случае некоторые потоки просто простаивают без работы, что негативно сказывается на масштабируемости алгоритма.

2.2 Списки вершин (Vertexes lists — VL)

В качестве решения первой проблемы предыдущего алгоритма — большого объёма копируемых данных после каждой итерации — списки смежности каждой компоненты

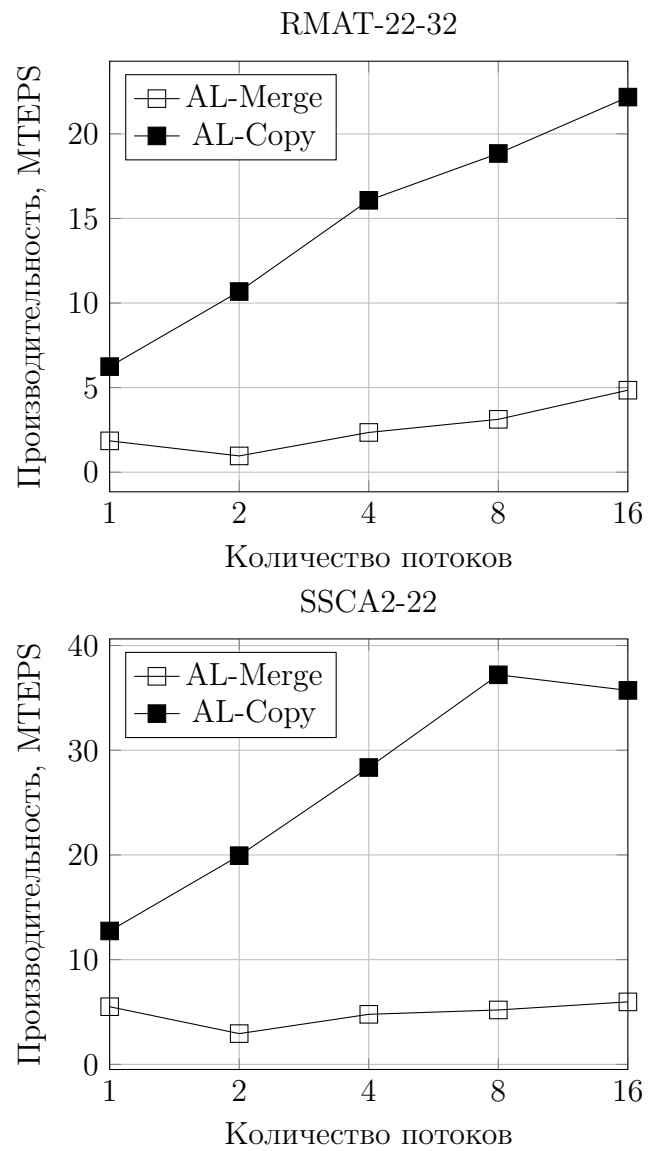


Рис. 1: Сравнение производительности реализаций алгоритма AL

были заменены на список входящих в неё вершин. Во-первых, такая модификация минимально влияет на другие шаги алгоритма. Во-вторых, на шаге **слияния списков** это позволяет в несколько раз сократить размер объединяемых списков.

Общая схема алгоритма осталась прежней со следующими изменениями:

1. На шаге **минимальное инцидентное ребро** теперь для обхода рёбер требуется обойти не список рёбер, а список входящих в компоненту вершин и уже для каждой вершины обойти список инцидентных ей рёбер.
2. На шаге **слияние списков** теперь копируются не списки рёбер, а списки вершин, что значительно сокращает объём перемещаемых данных.

Для хранения списка вершин были опробованы две структуры данных:

- **Односвязные списки.** Их преимущество заключается в возможности объединения за константное время, однако, обход таких списков связан со случайными обращениями в память.
- **Динамические массивы.** Объединение массивов возможно только за линейное время путём копирования одного из них в другой (на практике используется копирование меньшего массива в больший). Но во время обхода массива хорошо работает аппаратная предвыборка, что позволяет в значительной степени избежать кэш-промахов.

В качестве решения второй проблемы — дисбаланс на последних итерациях — шаг **минимальное инцидентное ребро** был разделён на два: сначала выделяются компоненты с небольшим количеством входящих в них вершин и обрабатываются как и раньше. Затем остаётся несколько компонент, состоящих из большого числа вершин, и работа по таким компонентам уже разделяется между потоками: каждый поток обрабатывает свою часть вершин, находит среди них минимальное ребро и в конце по всеми потокам находится итоговое ребро минимального веса.

На рисунке 2 отображено сравнение производительности обоих описанных подходов: с использованием односвязных списков (график VL-List) и с использованием динамических массивов (график VL-Vec). Для реализации на векторах так же показано её сравнение с оптимизацией (график VL-Vec-Bal), устраняющей дисбаланс и без неё: как видно из графиков, это заметно увеличило производительность и сделало подход более масштабируемым. Для реализации на связных списках применить описанную оптимизацию в простом варианте не удастся, поскольку при делении списка между потоками требуется обращение к произвольному элементу списка, чего односвязные списки в простой реализации дать не могут.

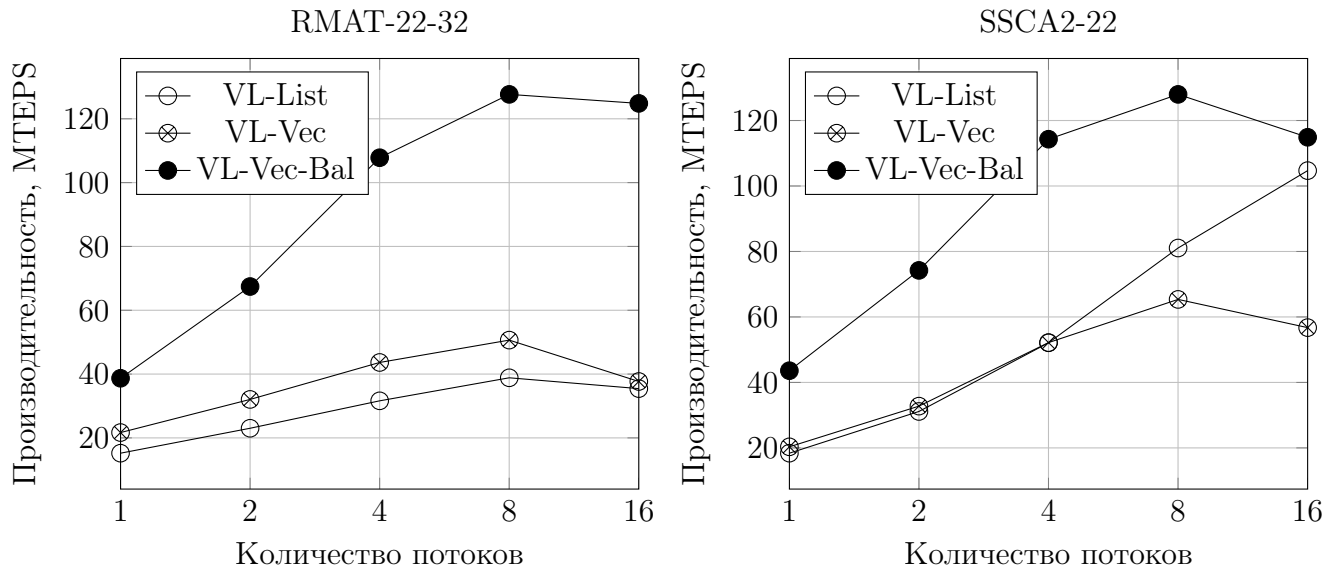


Рис. 2: Сравнение производительности реализаций алгоритма VL

2.3 Массив рёбер (Edges Array — EA)

Поскольку за счёт предыдущей модификации полностью решить проблему копирования больших объёмов данных не удалось, то следующий алгоритм разрабатывался с целью полностью избежать перемещения списков рёбер или вершин после каждой итерации. Если ранее использовался подход, когда для каждой компоненты хранилась информация, позволявшая восстановить список инцидентных рёбер, то теперь будет использоваться информация для каждой вершины о том, какой компоненте она принадлежит.

Изначально все вершины распределяются между потоками равномерно по количеству рёбер. Далее на каждой итерации алгоритма Борувки выполняются следующие шаги:

1. **Минимальное инцидентное ребро.** Каждый поток обходит все свои вершины и находит минимальное ребро для всех компонент, которые представлены в данном потоке хотя бы одной вершиной. В результате каждый поток будет иметь массив, содержащий минимальное ребро для каждой компоненты среди просмотренного подмножества.
2. **Редукция.** Происходит редукция полученного на предыдущем шаге массива: для каждой компоненты находится ребро минимального веса по всем потокам. Таким образом, для каждой компоненты определяется минимальное ребро уже во всём графе.
3. **Объединение деревьев.** Для каждой компоненты определяется её номер на следующей итерации. Так же решается проблема с возможными циклами (когда более двух компонент по кругу выберут следующую в качестве ближайшей), которые возможны в алгоритме Борувки при наличии рёбер одинакового веса.

4. **Перенумерация.** Осуществляется перенумерация компонент: для каждой вершины вычисляется номер компоненты, в которую она входит, используя параллельный алгоритм Pointer Jumping [15].

Основная проблема данного подхода заключается в шаге **редукции**: объём редуцируемого массива растёт линейно с ростом числа потоков, таким образом, среднее время выполнения данного шага почти не уменьшается с ростом количества потоков. Однако, даже с описанной проблемой подход с уменьшением объёма копируемых данных себя оправдал и данный подход показала лучшие результаты производительности среди описанных алгоритмов. Далее будут подробно рассмотрены оптимизации, применённые к данному подходу хранения графа.

3 Оптимизации алгоритма ЕА

3.1 Алгоритмические оптимизации

Предварительная сортировка рёбер. Алгоритм Борушки предполагает на каждой итерации просмотр всех рёбер графа, что и создаёт наибольшую вычислительную сложность первого шага. Объём “лишней” работы при этом достаточно велик: во-первых, постоянно будет просматриваться большое количество петель (рёбер, которые на данной итерации соединяют компоненту саму с собой), которых с каждой итерацией становится всё больше и больше, во-вторых, при обходе рёбер на шаге **поиска минимального ребра** хотелось бы уменьшить количество просматриваемых рёбер за счёт пропуска тех, вес которых уже больше, чем рекорд для текущей компоненты.

Проблему с петлями возможно попытаться решить их удалением, однако, это повлечёт за собой частичное перестроение внутреннего представления графа, что, как было выявлено в ходе исследования, негативно сказывается на производительности.

Решить обе обозначенные проблемы возможно за счёт предварительной сортировки списков смежности каждой вершины по возрастанию веса рёбер. Затем для каждой вершины поддерживается индекс последнего просмотренного ребра: изначально данный индекс указывает на первое ребро в списке, затем, по мере отбрасывания петель, он сдвигается и на следующей итерации поиск минимального ребра начинается не с начала списка, а с позиции, на которую указывает индекс. С другой стороны, если мы берём очередное ребро и его вес больше, чем текущий рекорд для компоненты, то поиск можно прервать, так как вес последующих рёбер данной вершины заведомо больше текущего рекорда.

Данная оптимизация хорошо масштабируется и на большом количестве потоков время выполнения сортировки невелико. В среднем производительность от данной оптимизации увеличивается на 20%.

Перенумерация вершин. В работе [16] предлагается алгоритм перенумерации вершин, который повышает локальность данных при последующих обходах графа, за счёт чего уменьшается количество кэш-промахов. Предложенный алгоритм перенумеровывает вершины в порядке обхода в ширину с предварительной отсортировкой списков смежности каждой вершины по возрастанию степени вершины, в которую ведёт ребро.

Балансировка больших вершин. Одной из особенностей RМАТ-графов является наличие вершин, степень которых значительно отличается от средней степени всех вершин графа. Например, в RМАТ-графе с 2^{20} вершинами есть 1351 вершина степени выше 1000, а максимальная степень вершины достигает 15078, при этом средняя степень вершины равна 32. Такие вершины создают дисбаланс на последних итерациях, когда просматривается большое количество рёбер и отбрасываются петли.

С целью устранения данного дисбаланса на стадии предварительно обработки графа выделяются вершины степени более 1000 и равномерно распределяются между пото-

ками. Значение степени 1000 было подобрано эмпирически и хорошо показало себя на практике.

В среднем оптимизации **перенумерация вершин** и **балансировки больших вершин** дают увеличение производительности в 11%.

Отдельный цикл пропуска петель. Шаг **поиск минимального ребра** с учётом оптимизации **предварительной сортировки рёбер** выглядит следующим образом:

Data: v — индекс обрабатываемой вершины

```

1  $myComp \leftarrow component[v]$ ;
2 foreach  $edgeId \leftarrow edgesIdsBegin[v]$  to  $edgesIdsEnd[v]$  do
3    $destVertex \leftarrow edges[edgeId].destination$ ;
4    $destComp \leftarrow component[destVertex]$ ;
5   if  $destComp = myComp$  then перейти к следующей итерации ;
6   if  $edges[edgeId] < currentRecord[myComp]$  then
7      $currentRecord[myComp] \leftarrow edges[edgeId]$  ;
7   закончить цикл;
8 end
```

Большая часть времени исполнения цикла находится в строках 2-5 и заключается в пропуске петель до тех пор, пока не найдётся ребро, ведущее в другую компоненту. Явное вынесение этих строк в отдельный цикл и разбиение цикла на две части ведёт к тому, что более простой цикл исполняется быстрее:

Data: v — индекс обрабатываемой вершины

```

1  $myComp \leftarrow component[v]$ ;
2  $edgeId \leftarrow edgesIdsBegin[v]$ ;
3 while  $edgeId \neq edgesIdsEnd[v]$  and  $component[edges[edgeId].destination] =$ 
    $myComp$  do
4    $edgeId \leftarrow edgeId + 1$ ;
5 end
6 if  $edges[edgeId] < currentRecord[myComp]$  then
7    $currentRecord[myComp] \leftarrow edges[edgeId]$ 
8 end
```

В среднем данная оптимизация ведёт к увеличению производительности на 10%.

3.2 Оптимизация редукции

Поскольку основные проблемы масштабируемости лежат именно в шаге **редукции**, то было сделано несколько оптимизаций, непосредственно затрагивающих данный шаг. Приведённые оценки увеличения производительности получены при запуске 32 потоков на всех 16 ядрах процессоров (исполнялось по 2 потока на одно ядро, поскольку в процессоре была включена технология Hyper Threading, использование которой описано далее в разделе 3.3) и показывают ускорение не одного только шага **редукции**, а исполнения всего алгоритма.

Отсутствие редукции на первой итерации. Поскольку каждая компонента изначально представлена одной вершиной, то от редукции на первой итерации можно отказаться. Данная оптимизация выглядит простой и очевидной, однако, объём редуцируемого массива на первой итерации больше, чем на всех последующих, и данная модификация даёт ускорение в 20%.

Перенумерация и сжатие массива компонент по ходу исполнения алгоритма. В начальной реализации при объединении компонент в одну в качестве её номера выбирался один из номеров вершин, которые в неё входят. Таким образом, несмотря на сокращение числа компонент, диапазон используемых номеров не изменялся (минимальный и максимальный номер были приблизительно равны соответственно первой и последней вершине). Такой подход затруднял обход всех компонент на шаге редукции: требовалось либо обходить весь диапазон номеров и совершать лишнюю работу, просматривая “пустые” номера, не являющиеся компонентами, либо отдельно поддерживать массив, в котором хранить индексы всех “живых” на данный момент компонент, что усложняло реализацию и приводило бы к увеличению кэш-промахов из-за обхода массива с непостоянным шагом.

Явная перенумерация компонент на каждой итерации на последовательные номера позволила решить обозначенную проблему. В среднем производительность увеличилась на 20%.

Иерархическая редукция. Особенностью NUMA архитектуры являются высокие задержки при обращении в память другого узла NUMA. При стандартном подходе к редукции, когда для каждой компоненты последовательно обходятся все потоки и выбирается лучший результат, таких обращений при двух NUMA узлах будет как минимум половина, а при большем количестве узлов NUMA — ещё больше.

Данную проблему можно решить за счёт иерархического подхода к редукции. При таком подходе шаг редукции разбивается на два этапа: на первом происходит стандартная редукция внутри NUMA узла, а на втором — между узлами. Количество обращений к чужому узлу NUMA в этом случае уменьшается в количество раз, равное количеству используемых потоков в одном узле NUMA. Иерархическая редукция повысила производительность на 5%.

3.3 Архитектурные оптимизации

Программная предвыборка. Обход рёбер в общем случае неизбежно связан со случайными обращениями в память: невозможно отсортировать рёбра так, чтобы индексы обеих инцидентных вершин шли последовательно. Аппаратная предвыборка современных процессоров ориентирована только на последовательные обращения в память (на обход с постоянным, возможно отрицательным, шагом, размер которого не превосходит размер страницы памяти) и подобный обход рёбер вызывает простои, связанные с ожиданием данных.

Однако, адреса случайных обращений в память при таких обходах в большинстве

случаев можно предсказать: обходя рёбра последовательно, нам известны индексы обеих инцидентных вершин для следующих итераций цикла. Как раз в таком случае и помогает программная предвыборка, давая возможность заранее запросить нужные данные для будущих вершин (например, какой компоненте они принадлежат) и сократить время ожидания.

В среднем программная предвыборка позволила повысить производительность на 30%, а шаг **перенумерация вершин**, где используется алгоритм Pointer Jumping, ускоряется до 2-3 раз.

Иногда в одном цикле целесообразно использование программной предвыборки для двух массивов с разным шагом. Так, в коде шага **поиска минимального ребра**:

```

1 for vertexId  $\leftarrow$  1 to vertexCount do
2   | edgeId  $\leftarrow$  startEdge[vertexId];
3   | destComp  $\leftarrow$  component[edges[edgeId].destination];
4   | ...
5 end
```

предвыборка для переменной *edgeId* делается достаточно легко, а вот переменная *destComp* зависит от значения *edgeId* и прежде, чем будет возможно сделать предвыборку для *destComp*, необходим получить значение *edgeId*. Предвыборки для обеих переменных без лишних простоев возможно реализовать, делая их с разным шагом: если предвыбирать *edgeId* с удвоенным шагом относительно *destComp*, то обе предвыборки не будут зависеть друг от друга:

```

1 for vertexId  $\leftarrow$  1 to vertexCount do
2   | prefetch(&edge[startEdge[vertexId + prefetchStep  $\times$  2]]);
3   | prefetch(&component[edge[startEdge[vertexId + prefetchStep].dest]]);
4   | edgeId  $\leftarrow$  startEdge[vertexId];
5   | destComp  $\leftarrow$  component[edges[edgeid].destination];
6   | ...
7 end
```

Выделение памяти с учётом NUMA архитектуры Как уже отмечалось ранее, обращение в память другого узла NUMA связано с высокими задержками. Выделение памяти на тех ядрах, где она используется, является более правильным на NUMA архитектуре.

Все используемые в реализации массивы ограничиваются двумя размерностями. Более того, во всех двумерных массивах первая размерность соответствует индексу потока, в котором он используется (например, таким является массив, над которым происходит редукция: первая размерность соответствует номеру потока, вторая — компоненте, а значением является ребро минимального веса, которое данный поток нашёл для данной компоненты), и с точки зрения NUMA архитектуры такие массивы проблем не создают.

Когда разные части одномерного массива одновременно используют несколько по-

токов (например, таким является массив рёбер), то возникает необходимость выделить непрерывную область в памяти на разных узлах NUMA. Ядро Linux предоставляет разработчику возможность явно управлять выделением страниц памяти на NUMA системах с использованием библиотеки `libnuma`. Однако, в реализациях использовалась особенность ядра Linux известная как `memory overcommit`: для программы выделение страниц памяти происходит не в момент вызова функции `malloc` (или аналогичной, запрашивающей у ядра операционной системы страницы памяти, но не инициализирующей их), а в момент первого обращения к странице. Таким образом, для разделения массива по узлам NUMA сначала с использованием функции `malloc` выделялся весь массив нужного размера, далее каждый поток инициализирует те страницы, которые он будет использовать. При таком подходе, конечно, может возникнуть пересечение, когда два потока используют одну страницу, но размеры таких пересечений невелики по сравнению с размером всего массива, а их количество ограничено количеством потоков.

Был также опробован подход с дублированием массива, который используют одновременно несколько узлов NUMA. Одним из примеров таких данных является массив, где для каждой вершины хранится индекс компоненты, в которую она входит на текущей итерации. Однако, из-за необходимости после каждой итерации дополнительно обновлять массив на каждом узле NUMA данный подход себя не оправдал и производительность в лучшем случае не уменьшалась.

Аппаратные потоки процессора. Современные процессоры позволяют на одном ядре исполнять одновременно несколько аппаратных потоков. Так в процессорах Intel и AMD данные технологии называются соответственно `Hyper-Threading` и `Modules` и позволяют запускать по 2 аппаратных потока на ядро, а ускоритель Intel Xeon Phi на данный момент имеет 4 аппаратных потока на каждое ядро. Использование такого вида многозадачности позволяет не только более эффективно использовать процессорные элементы, но и скрыть задержки в памяти за счёт переключения на другой поток.

Проблема использования аппаратных потоков для алгоритма ЕА заключается в увеличении объёма данных на шаге **редукции**. В начальных реализациях использование `Hyper-Threading` давало небольшое ускорение на малом количестве потоков, но ухудшало время на 16 потоках. После применения описанных выше оптимизаций редукции `Hyper-Threading` стал давать выигрыш в среднем 5%.

Linux Huge Pages. Объём обрабатываемых данных, находящихся в памяти, составляет гигабайты, что во много раз больше размера одной страницы памяти в процессорах архитектуры x86, равного 4 килобайтам. Использование большого числа страниц виртуальной памяти приводит к частым промахам в TLB кэше, что, конечно же, отрицательно сказывается на производительности. Размер страницы памяти зависит от архитектуры и во многих архитектурах, в том числе и в x86, реализована поддержка “больших страниц” — страниц, размер которых равен нескольким мегабайтам или даже гигабайтам.

Увеличение объёма страниц памяти ведёт к уменьшению количества используемых

страниц, что в свою очередь уменьшает число промахов в TLB кэше. К сожалению, на данный момент хоть и huge pages включены в ядре Linux, однако, настройки системы по умолчанию не позволяют их использовать без дополнительных действий со стороны администратора системы, потому данная оптимизация не была протестирована на центральном процессоре Intel Xeon.

Лучше дела обстоят на ускорителе Intel Xeon Phi, где Huge Pages разрешены “из коробки” и не требуют дополнительных настроек. Использование Huge Pages на данном ускорителе дало выигрыш от 5% до 15%.

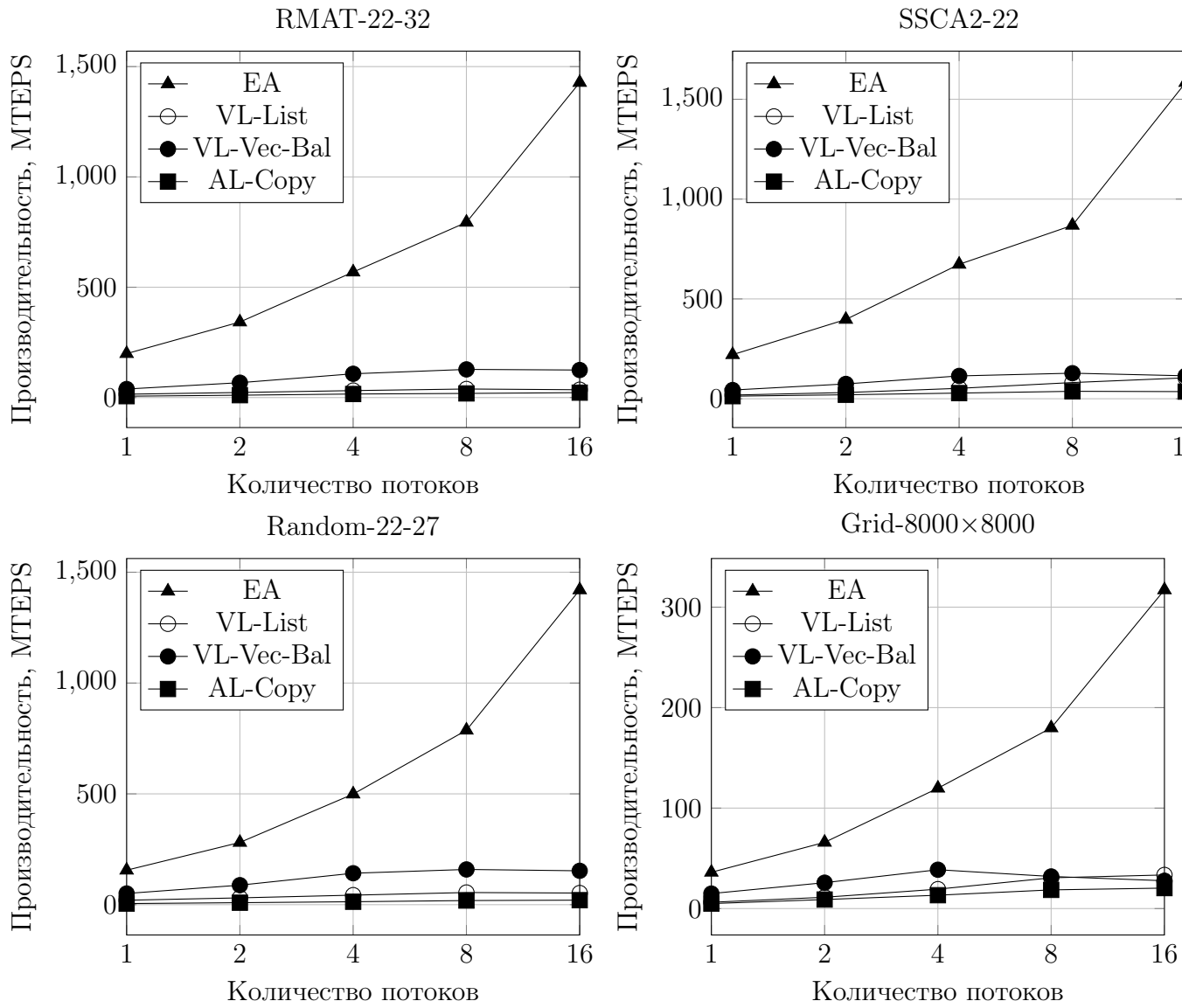


Рис. 3: Производительности алгоритмов на разных графах

4 Результаты

Тестирование проводилось на следующих графах:

- **RMat-графы.** Структура RMat-графа, процесс генерации которого подробно описан в [17], напоминает структуру социальной сети. Особенностью таких графов является наличие одной большой компоненты, в которую входят почти все вершины графа, а так же множество небольших компонент, в том числе и отдельно висячие вершины. Степень вершин в таких графах различна и может очень сильно отличаться от средней степени вершин во всём графе: присутствуют как множество вершин небольшой степени, так и отдельные вершины, степень которых на 2-3 порядка выше средней степени.
- **SSCA2-графы.** SSCA2 графы представляют из себя множество небольших клик,

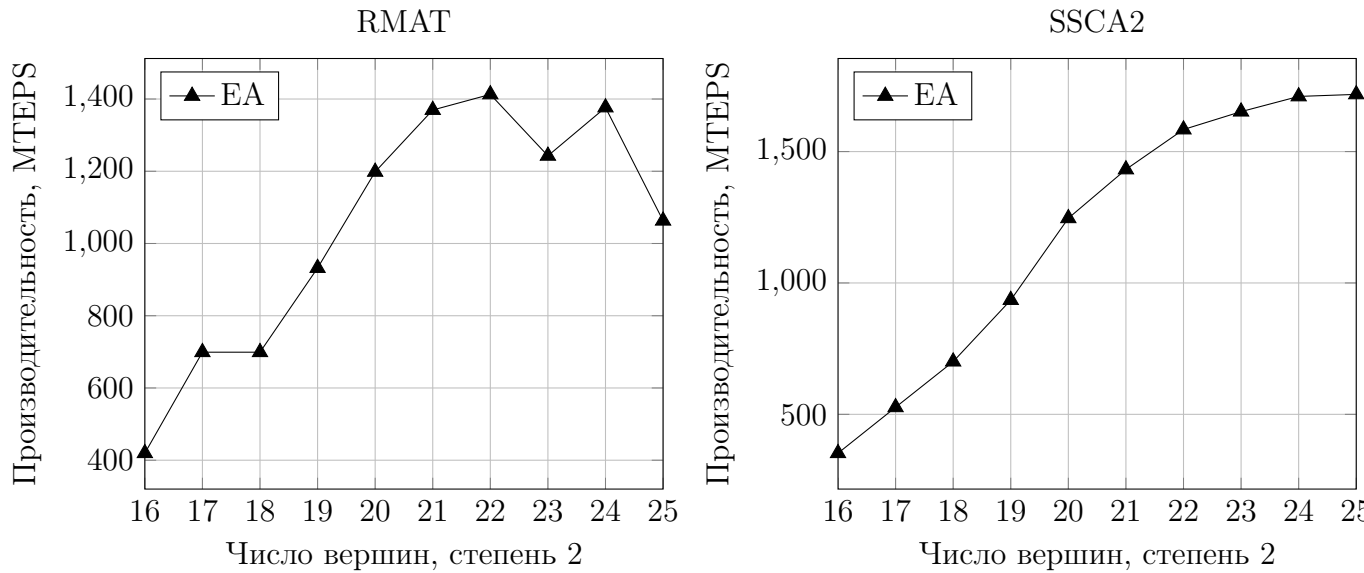


Рис. 4: Сравнение производительности алгоритма EA на графах разного размера

между которыми случайно сгенерировано некоторое число рёбер [18].

- **Двумерные и трёхмерные решётки.** Графы, структура которых является N -мерной решёткой: каждая вершина графа, за исключением крайних, соединена рёбрами с 2^{N+1} соседними вершинами.

Данный граф можно схематично представить как решётку в N -мерном пространстве, где каждая вершина v имеет целочисленные координаты $(v_0, v_1, \dots, v_{N-1})$, $0 \leq v_i < L_i$, где L_i — длина i -ой размерности решётки. Расстояние между двумя вершинами v и u в такой решётки полагается равным сумме модулей разностей по каждой координате: $\sum_{i=0}^{N-1} |v_i - u_i|$. Всего в решётке будет $L_0 \cdot L_1 \cdot \dots \cdot L_{N-1}$ вершин и смежными являются те вершины, расстояние между которыми равно 1.

- **Случайные графы.** При генерации случайных графов задаётся количество вершин графа и требуемое количество рёбер. Затем генерируются рёбра графа: для каждого ребра случайно выбираются обе инцидентные ему вершины, при этом каждая из вершин может быть выбрана с равной вероятностью.

Вес каждого ребра графа генерировался независимо и является случайной величиной на интервале от 0 до 1, имеющей равномерное распределение. Во всех графах, кроме решёток, допускаются петли и кратные рёбра.

В качестве формата хранения графа использовался сжатый формат хранения разреженных матриц CSR. В данном формате все ненулевые значения матрицы располагаются в памяти последовательно таким образом, что все элементы одной строки идут подряд слева направо, а строки матрицы идут последовательно сверху-вниз. Для каждого ненулевого элемента помимо его значения хранится индекс столбца, в котором он находится. Также для каждого элемента необходимо уметь определять и номер строки,

на которой он расположен, но за счёт того, что элементы одной строки идут подряд, достаточно для каждой строки исходной матрицы хранить индекс первого и последнего элемента в массиве элементов, а за счёт того, что строки идут последовательно, индекс последнего элемента строки i будет на единицу меньше, чем индекс первого элемента строки $i + 1$ (на практике обычно для хранения индексов обычно используют полуоткрытые интервалы — в этом случае данные индексы будут совпадать).

Запуски проводились на двухsocketной машине с процессорами Intel Xeon CPU E5-2690. Каждый процессор состоит из 8 ядер с тактовой частотой 2.9 GHz, каждое ядро имеет кэш первого уровня 32 килобайта и кэш второго уровня 256 килобайт. Так же каждый процессор имеет общий для всех ядер кэш третьего уровня 20 мегабайт. На тестируемой машине была включена технология Hyper Threading, позволяющая исполнять на каждом ядре процессора по 2 аппаратных потока.

Время работы программы измерялось с помощью ассемблерной инструкции `rdtsc`. Каждая реализация была запущена на каждом из графов 5 раз и в качестве итогового значения времени выбирался минимальный результат.

На всех графиках результаты представлены в MTEPS (TEPS, traversed edges per second — количество обработанных рёбер в секунду, приставка М означает mega) — величина, обратно пропорциональная времени и равная:

$$MTEPS = \frac{\text{количество рёбер в графе}}{(\text{время работы программы в секундах}) \cdot 10^6}$$

На рисунке 3 представлено сравнение описанных выше алгоритмов на графах следующих размеров:

- RMat-22-32 — RMat-граф с 2^{22} вершинами и средней степенью вершины 32.
- SSCA2-22 — SSCA2-граф, в котором 2^{22} вершин и примерно 2^{27} рёбер.
- Random-22-27 — случайный граф в котором, соответственно, 2^{22} вершин и 2^{27} рёбер.
- Grid-8000×8000 — двумерная решётка размера 8000 на 8000. В такой решётке находится 2^{26} вершин и приблизительно 2^{27} рёбер.

Размеры графов выбирались таким образом, чтобы количество рёбер было примерно одинаковым. Количество вершин в разных графах, за исключением решётки, так же примерно одинаково.

На графиках отображена зависимость производительности от числа используемых потоков. Результаты алгоритма ЕА заметно лучше остальных. Огромная разница в производительности обосновывается тем, что не все оптимизации, описанные в разделе 3 были применены и для остальных алгоритмов, но даже неоптимизированный вариант показывал лучшие результаты. Среди остальных алгоритмов относительно неплохую

производительность показывает VL-алгоритм на динамических массивах с дополнительной балансировкой нагрузки, однако, имеются проблемы с масштабируемостью на два сокета.

Решётка же получается более разреженной и имеет при таком же количестве рёбер в 16 раз больше вершин. Это ведёт к большему количеству итераций, который должен выполнить алгоритм Борувки прежде, чем остовное дерево будет построено, и как следствие увеличивается время работы. Так, для построения остовного дерева в графе RМAT-22-32 требуется 6 итераций, а в решётке 8000×8000 — 13 итераций, чем и объясняется более низкая производительность при равном количестве рёбер.

На рисунке 4 представлено изменение производительности алгоритма ЕА в зависимости от размера графа. Сравнение проводилось на RМAT графах и на SSCA2 графах со средней степенью вершины 32.

На малых размерах графа производительность низкая, поскольку в реализации присутствуют барьерные синхронизации между шагами алгоритма, а такие шаги, как **перенумерация вершин**, дополнительно содержат барьерные синхронизации и внутри себя. На больших графах объём вычислений между синхронизациями достаточно большой и время, затраченное на такие синхронизации, очень мало. На графах же небольшого размера время, которое потоки проверили на синхронизации, уже негативно сказывается на производительности.

Список литературы

- [1] Y. Xu, V. Olman, and D. Xu. Clustering gene expression data using a graph-theoretic approach: an application of minimum spanning trees. *Bioinformatics* 18 (4), 2002.
- [2] Kamesh Madduri, David A. Bader, Jonathan W. Berry, Joseph R. Crobak, and Bruce A. Hendrickson. Multithreaded algorithms for processing massive graphs. *Petascale Computing: Algorithms and Applications*, 2007.
- [3] D. Migov and A. Rodionov. Parallel implementation of the factoring method for network reliability calculation. *Springer Lecture Notes in Computer Science (in ICCSA 2014)*. Vol. 8584, issue 6, pages 654–664, 2014.
- [4] Родионов А. С., Сакерин А. В., and Мигов Д. А. Некоторые вопросы реализации полного перебора (на примере задач надёжности сетей). *Вестник СубГУТИ. №4*, pages 79–84, 2014.
- [5] R. L. Graham and Pavol Hell. *On the history of the minimum spanning tree problem*. 1985.
- [6] T. Asano, B. Bhattacharya, M. Keil, and F. Yao. *Clustering algorithms based on minimum and maximum spanning trees*. Fourth Annual Symposium on Computational Geometry, 1988.
- [7] Ma, A. Hero, J. Gorman, and O. Michel. *Image registration with minimum spanning tree algorithm*. 2000.
- [8] Ernesto Tapia and Raúl Rojas. *Recognition of On-line Handwritten Mathematical Expressions Using a Minimum Spanning Tree Construction and Symbol Dominance*. 2004.
- [9] Белоусов А. И. and Ткачев С. Б. *Дискретная математика*. МГТУ, 2006.
- [10] Т. Кормен. *Алгоритмы: построение и анализ*. Вильямс, 2005.
- [11] David A. Bader and Guojing Cong. Fast shared-memory algorithms for computing the minimum spanning forest of sparse graphs. <http://www.cc.gatech.edu/~bader/papers/MST-JPDC.pdf>, 2006.
- [12] Vladimir Loncar, Srdjan Skrbic, and Antun Balaz. Parallelization of minimum spanning tree algorithms using distributed memory architectures. 2014.
- [13] Frank Dehne and Silvia Gotz. Practical parallel algorithms for minimum spanning trees. <http://people.scs.carleton.ca/~dehne/publications/2-47.pdf>, 1998.

- [14] Sun Chung and Anne Condon. Parallel implementation of boruvka minimum spanning tree algorithm. <http://minds.wisconsin.edu/bitstream/handle/1793/60020/TR1297.pdf>, 1996.
- [15] Guy E. Blelloch and Bruce M. Maggs. Parallel algorithms. <http://homes.cs.washington.edu/~arvind/cs424/readings/BlellochM96b.pdf>, 2004.
- [16] Elizabeth Cuthill and J. McKee. Reducing the bandwidth of sparse symmetric matrices. *ACM '69 Proceedings of the 1969 24th national conference*, pages 157–172, 1969.
- [17] Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. R-mat: A recursive model for graph mining. 2004.
- [18] David A. Bader and Kamesh Madduri. Design and implementation of the hpcs graph analysis benchmark on symmetric multiprocessors. 2005.