

**ОРДЕНА ЛЕНИНА**  
**ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ**  
**ИМЕНИ М.В.КЕЛДЫША**  
**РОССИЙСКОЙ АКАДЕМИИ НАУК**

**В.Н.КОВАЛЕНКО, Е.И.КОВАЛЕНКО,**  
**Д.А.КОРЯГИН, Э.З.ЛЮБИМСКИЙ**

**МЕТОД ОПЕРЕЖАЮЩЕГО**  
**ПЛАНИРОВАНИЯ ДЛЯ ГРИД**

**Москва**  
**2005 г.**

УДК 519.68

***В.Н.Коваленко, Е.И.Коваленко, Д.А.Корягин, Э.З.Любимский. Метод опережающего планирования для грид.***

В работе предлагается метод планирования распределения ресурсов между вычислительными заданиями для практически важной формы организации грид, когда ресурсы не отдаются в грид полностью, а используются в разделяемом с их владельцами режиме. Особенностью этого метода, названного опережающим планированием, является то, что ресурсы распределяются не в момент их освобождения, а строится план их использования на некоторый период времени. Благодаря этому, планирование в условиях разделения ресурсов с владельцами обладает свойством детерминированности – выполнение плана обеспечивается путем предварительного резервирования ресурсов. Свойство детерминированности делает метод опережающего планирования пригодным для обслуживания многопроцессорных заданий.

***Ключевые слова:*** грид, распределенный компьютеринг, планирование, резервирование ресурсов, кластер, коаллокация.

***V.N. Kovalenko, E.I. Kovalenko, D.A.Koryagin, E.Z. Ljubimsky. Method of Lookahead Scheduling for Grid.***

In this paper we propose a new method of resource distribution between computational jobs for a practically important form of a grid, where resources are used not only by grid, but are shared with their owners. The characteristic feature of the method, called Lookahead scheduling, is it's ability to build up a plan of resource allocations for some period of a time, unlike more common scheme, in which a scheduling is limited to forwarding jobs to free resources. Such "predictive" scheduling even in the conditions of shared resources has the property of a determinacy – realization of the plan is provided by means of advance reservation. The property of the determinacy makes the Lookahead method suitable for a multiple-processor job management.

***Key words:*** grid, distributed computing, scheduling, resource reservation, cluster, coallocation

---

***Работа выполнена при поддержке Российского фонда фундаментальных исследований (проекты 05-01-00626-а, 04-07-90299-в).***

**ОГЛАВЛЕНИЕ**

1. ВВЕДЕНИЕ.....	4
2. МЕСТО ПЛАНИРОВАНИЯ В АРХИТЕКТУРЕ ГРИД.....	6
3. СПЕЦИФИКА ПЛАНИРОВАНИЯ В ГРИД.....	8
3.1. Гетерогенность ресурсов и их отбор для заданий.....	9
3.2. Способы организации грид и современные интерфейсы доступа к ресурсам.....	10
3.3. Режимы использования ресурсов.....	12
3.4. Проблемы планирования при использованием стандартных интерфейсов менеджеров ресурсов.....	13
4. ОСНОВНЫЕ ПОЛОЖЕНИЯ ПРЕДЛАГАЕМОГО МЕТОДА ПЛАНИРОВАНИЯ.....	14
4.1. Разделение ресурсов между пользователями.....	16
4.2. Разделение ресурсов с владельцами.....	17
5. ОБЕСПЕЧЕНИЕ ДЕТЕРМИНИРОВАННОСТИ ПЛАНИРОВАНИЯ.....	18
5.1. Прогноз загрузки кластерных ресурсов.....	18
5.2. Предварительное резервирование ресурсов.....	21
6. СХЕМА ОПЕРЕЖАЮЩЕГО ПЛАНИРОВАНИЯ.....	22
6.1. Масштабируемость системы планирования.....	25
7. СРАВНЕНИЕ С СОВРЕМЕННЫМИ ПЛАНИРОВОЩИКАМИ.....	26
8. ЗАКЛЮЧЕНИЕ.....	29
9. ЛИТЕРАТУРА.....	30

## 1. ВВЕДЕНИЕ

Хотя единой точки зрения на то, что такое грид, и, тем более, общепринятого определения этого понятия до сих пор не существует, дискуссии [1] ведутся главным образом по поводу принципов организации, стандартов архитектуры и конкретных технологий. В то же время, расхождений по поводу базовой идеи – интеграции распределенных в сети ресурсов и обеспечения к ним эффективного удаленного доступа – существенно меньше, что позволяет сформулировать следующее общее определение грид.

**Грид** – это программно-аппаратная среда, которая построена из вычислительных установок, находящихся в различных административных доменах телекоммуникационной сети и которая позволяет дистанционно использовать любое количество ресурсов этих установок – процессорных, оперативной и постоянной памяти, программ и данных.

Усилия специалистов, разрабатывающих технологии грид, направлены на то, чтобы обеспечить качественно новый уровень интеграции распределенных ресурсов, при котором можно было бы рассматривать грид как единую операционную среду обработки запросов (заданий), а работа с грид была не сложнее, чем с более привычными компьютерными архитектурами: современными компьютерами, многопроцессорными и кластерными системами.

С помощью разработанного к настоящему времени программного обеспечения грид [2] созданы функционирующие в производственном режиме компьютерные комплексы сверхвысокой производительности и пропускной способности с агрегированной мощностью порядка нескольких терафлопс и петабайтными объемами хранимых данных EGEE ([www.eu-egee.org](http://www.eu-egee.org)), NorduGrid ([www.nordugrid.org](http://www.nordugrid.org)), Grid2003 ([www.ivdgl.org/grid2003](http://www.ivdgl.org/grid2003)), TeraGrid ([www.teragrid.org](http://www.teragrid.org)). Построение подобных систем дает возможность решать совершенно разные прикладные задачи беспрецедентного уровня сложности как в научно-технической, так и в производственной сфере [3].

Грид принадлежит к архитектурам распределенного компьютеринга с коллективной формой обслуживания пользователей, и у него, конечно, есть аналоги - в качестве ближайшего можно указать кластерные системы. Поэтому многие из задач, решаемых программным обеспечением грид, известны, однако в условиях грид постановка этих задач значительно усложняется, а для их решения требуются новых подходы.

Среди наиболее значимых отличительных свойств грид отметим:

- Глобальную распределенность. Составляющие грид вычислительные установки распределены в открытой глобальной сетевой среде, в которой велики коммуникационные издержки и предъявляются повышенные требования к безопасности.

- Автономность ресурсов. Ресурсная база грид формируется из независимых друг от друга вычислительных центров или отдельных компьютеров. Ресурсы обслуживаются и администрируются владельцами, которые имеют право проводить собственную политику доступа к ним.

- Коллективный режим работы. Ресурсы грид используются в коллективном режиме, поэтому должно обеспечиваться гибкое и скоординированное их распределение между пользователями, решающими разные задачи.

- Динамичность среды. Как состав ресурсов, так и состав пользователей может меняться, причем это является правилом, а не исключением.

Набор функций, которые требуются в этих условиях от программного обеспечения, в главных чертах определен – это средства обеспечения безопасности, надежности, мониторинга заданий и устройств, учета и протоколирования заданий. Ключевое место в этом ряду занимает функция диспетчеризации заданий: она обеспечивает распределение ресурсов из общего ресурсного пула грид между заданиями, доставку программ и данных, реализуя одну из важнейших концепций грид – виртуализацию ресурсов.

Уже сейчас существуют программные разработки, в которых достигнут практически приемлемый уровень реализации перечисленных функций. С наибольшей полнотой круг вопросов, связанных с поддержкой функционирования грид, решен в программных платформах LCG [4] и gLite [5], в которых интегрированы и взаимосвязаны компоненты программного обеспечения грид, разработанные разными исследовательскими коллективами. Тем не менее, едва ли можно утверждать, что получены окончательные решения.

Тема данной работы – планирование распределения заданий по ресурсам грид. Планирование происходит в контексте системы управления заданиями, обеспечивающей их дистанционное выполнение. Объекты планирования – задания и ресурсы, но задания запускаются пользователями, а ресурсы имеют владельцев, и мы рассматриваем такое планирование, которое учитывает интересы тех и других (раздел 2). В разделе 3 вводится понятие детерминированного планирования: такого, что построенный планировщиком план действительно реализуется. Оказывается, что в условиях грид (разделы 3.1 – 3.3) сделать это непросто. В первую очередь это обусловлено тем, что архитектура грид имеет два уровня управления: на верхнем действует планировщик грид, а на нижнем – локальный менеджер ресурсов. В большинстве применяемых на практике систем планирования в качестве аппарата взаимодействия этих двух уровней используются традиционные для локальных менеджеров ресурсов интерфейсы, изначально предназначенные для обслуживания пользователей систем пакетной обработки (СПО). Использование этого аппарата в грид приводит к существенным дефектам планирования (раздел 3.4), например к невозможности обеспечить точное

время старта заданий, а также накладывает жесткие ограничения на способы организации грид, что не позволяет в полной мере удовлетворить такие необходимые требования, как возможность включения в грид разнотипных компьютеров или возможность использования ресурсов грид совместно с их владельцами.

В разделе 4 излагаются исходные положения метода опережающего планирования, определяющие стратегию планирования, основанную на соглашениях о разделении ресурсов между пользователями и владельцами ресурсов. В разделе 5 показывается, каким образом можно добиться детерминированности планирования путем расширения СПО двумя новыми функциями: прогнозирования загрузки ресурсов и предварительного резервирования. В разделе 6 приводится программная схема опережающего планирования, использующего эти функции. Дана также оценка масштабируемости предлагаемого метода (раздел 6.1). Раздел 7 посвящен обзору современных систем планирования в грид и их сравнению с предлагаемым подходом.

## **2. МЕСТО ПЛАНИРОВАНИЯ В АРХИТЕКТУРЕ ГРИД**

Получивший общее признание способ организации программного обеспечения грид основан на открытой архитектуре служб OGSA [6], в соответствие с которой грид квалифицируется как программная система, состоящая из распределенных компонентов – служб, взаимодействующих между собой посредством стандартных, открытых и универсальных протоколов и интерфейсов [7]. Имея в виду главным образом грид вычислительного типа, можно рассматривать его функционирование как процесс обслуживания стандартизированных запросов на выполнение вычислений, оформленных в виде заданий для общераспространенных операционных систем, причем выполнение этих заданий производится на ресурсах, которые выбираются из общего пула. Перечислим основные этапы обработки задания (рис. 1):

- Планирование ресурсов. Специальная компонента программного обеспечения – планировщик выделяет из общего пула исполнительные ресурсы - те, на которых задание будет выполняться.
- Доставка исполняемых файлов и входных файлов на исполнительные ресурсы.
- Выполнение задания.
- По окончании задания доставка результирующих файлы на серверы хранения (в частности, на рабочее место пользователя).

Все перечисленные этапы обработки задания выполняются автоматически, без участия субъекта, выдавшего запрос (в частности, пользователя, хотя это может быть и программа), так что грид действительно представляет собой единую операционную среду.

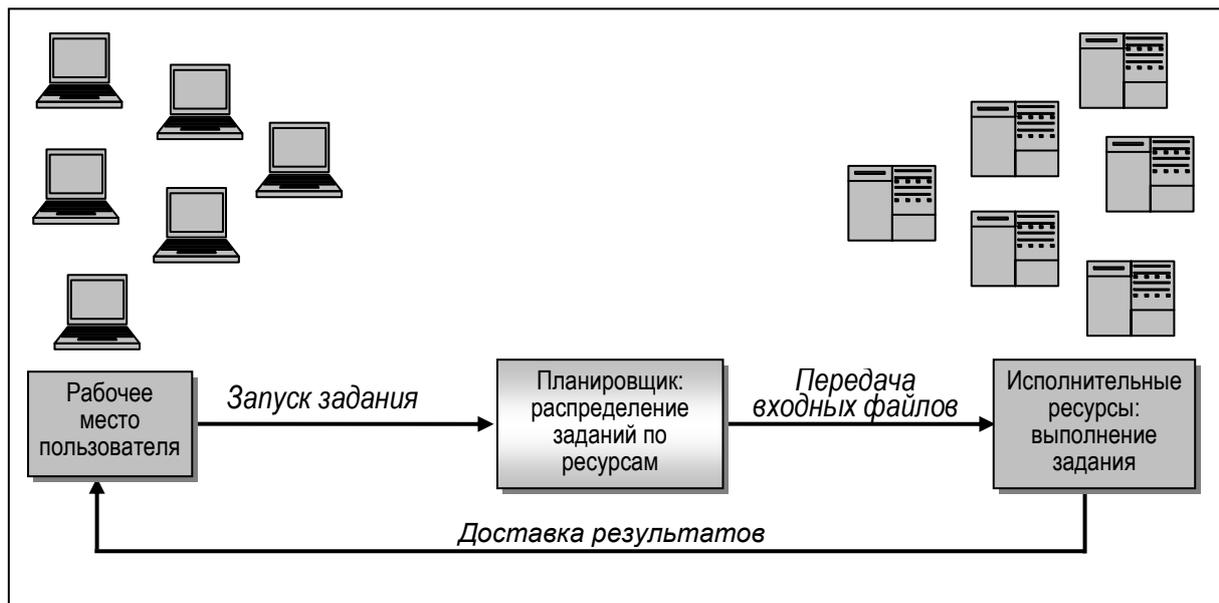


Рис. 1. Схема обработки запросов и виртуализация ресурсов.

Поскольку грид - распределенная среда, то, как показала практика, необходимы специальные службы, выполняющие:

- мониторинг и протоколирование процесса обработки запроса, что позволяет в каждый момент получать информацию о местонахождении и состоянии задания, а также управлять его обработкой (удаление, модификация);
- обеспечение надежности - перезапуск задания при сбоях программных или аппаратных средств на ресурсах и телекоммуникациях;
- мониторинг и учет потребления ресурсов, на основе чего производится планирование и диагностика сбоев оборудования.

Таким образом, планирование ресурсов – основная тема данной работы – происходит в контексте сложной многокомпонентной системы обработки заданий, примером которой является WMS (Workload Management System) [8], входящая в состав упоминавшейся выше платформы LCG. В этом контексте назначение планирования – распределение ресурсов грид между заданиями. Планирование имеет два аспекта: во-первых, оно определяет исполнительные ресурсы, на которых будет выполняться каждое задание, и, во-вторых, время, когда исполнительные ресурсы отводятся тому или иному заданию. Цель планирования – обеспечить скоординированное разделение ресурсов, учитывающее как интересы пользователей, так и владельцев ресурсов. Эти интересы выражаются правилами и соглашениями, которые устанавливаются объединениями пользователей и владельцев ресурсов – виртуальными организациями. Задача планирования в грид представляет интерес как самостоятельное исследовательское направление, которому посвящено большое количество публикаций (см., например, [9]).

### 3. СПЕЦИФИКА ПЛАНИРОВАНИЯ В ГРИД

Планирование в грид относится к широкому классу задач, общая постановка которых характеризуется наличием множества ресурсов, множеством их потребителей и определенными условиями выделения ресурсов потребителям [10]. Известно, что способы планирования в зависимости от этих условий, целей планирования и других факторов, могут быть совершенно разными.

Что касается постановки задачи в данной работе, то в качестве ресурсов будут рассматриваться вычислительные ресурсы компьютеров  $\mathbf{RS}=\{RS_i, i=1,\dots,N\}$  – процессоры, вместе со связанными с ними оперативной и дисковой памятью. Компьютеры глобально распределены, то есть могут находиться в произвольных точках глобальной сети. В систему планирования (планировщик) от работающих независимо друг от друга пользователей поступает неограниченный во времени поток заданий  $\mathbf{Z}(T)=\{z_i(t_i), i=1,2,\dots\}$ , где  $t_i$  – время поступления задания. Мы предполагаем, что все компьютеры однопроцессорные и выделяются заданиям в эксклюзивном режиме: в любой момент времени на компьютере выполняется не более одного задания.

Результатом планирования является последовательность **аллокаций** ресурсов  $\mathbf{A}=\{A_i, i=1,2,\dots\}$ . Нас интересуют такие методы планирования, способные строить **точные аллокации**  $A \in \mathbf{A}$ ,  $A=[z, \mathbf{IRS}, t_b, t_e]$ , в которых определяется множество исполнительных ресурсов  $\mathbf{IRS} \subseteq \mathbf{RS}$  для выполнения задания  $z$  и временной интервал  $[t_b, t_e]$ , на который они отводятся заданию  $z$ . На этом интервале ресурсы считаются занятыми и недоступными другим заданиям.

Назовем планирование **детерминированным**, если оно строит точные аллокации, и запуск заданий на исполнительных ресурсах осуществляется в соответствии с этими аллокациями. Для этого необходимо, чтобы:

- характеристики среды исполнительных ресурсов отвечали требованиям задания;
- задание могло быть доставлено на исполнительные ресурсы к моменту начала аллокации;
- исполнительные ресурсы были к этому моменту свободны, и политика управления ресурсами не препятствовала их выделению для задания.

Как будет показано в следующих разделах, особенности грид делают возможность построения точных аллокаций проблематичной, и вместо этого во многих алгоритмах и практических реализациях планировщиков результатом являются аллокации, в которых некоторые элементы (исполнительные ресурсы, время начала аллокации) определены примерно, что делает планирование недетерминированным в той или иной степени. Проявляется это в серьезных дефектах обработки заданий: например, если не известно точно время начала аллокаций, то в некоторых ситуациях запуск задания вообще может не состояться – есть шанс, что задание “зависнет”.

Рассмотрим далее обстоятельства, в которых происходит планирование в грид и которыми в большой степени определяется возможность построения метода планирования, обладающего свойством детерминированности: гетерогенность ресурсов; организация ресурсов; режимы использования ресурсов.

### *3.1. Гетерогенность ресурсов и их отбор для заданий*

Задание грид представляет собой обычный исполняемый файл (скрипт, программный код). Службами грид оно доставляется на исполнительные ресурсы, а собственно выполнение происходит в среде операционной системы (ОС) этих ресурсов. Как правило, программа, подготовленная на определенном компьютере, не требует каких-либо модификаций для использования в грид. Однако, любая программа рассчитана на определенную среду выполнения – ОС, архитектуру компьютера, объемы и характеристики его ресурсов.

Строя аллокации, планировщик должен учитывать, что грид является гетерогенной инфраструктурой, в состав которой могут включаться компьютеры с разной архитектурой и комплектацией. Гетерогенность проявляется в том, что разные классы ресурсов (процессор, основная память, кэш-память, дисковая память) различаются типом (процессоры – архитектурой) и характеристиками (процессоры - производительностью, память – объемом). В связи с этим, выбираемые исполнительные ресурсы не могут быть произвольными, а должны соответствовать требованиям задания. Эти требования не обязательно ограничиваются приведенными примерами: любая особенность компьютера, существенная для выполнения задания, например, операционная система, может рассматриваться как класс ресурсов.

При запуске через систему планирования требования задания оформляются в виде **ресурсного запроса**. Используемые на практике формализмы записи ресурсного запроса, хотя и разнообразны, но ориентированы на то, чтобы рассматривать ресурсы унифицированным способом, предполагая, что набор классов является стандартизированным (с возможностью расширения). Простейший вид ресурсного запроса можно представить таким образом:

`<ресурсный запрос>:={<класс>={<тип>|<характеристики>}...} [<время использования ресурсов>]`

Например, по ресурсному запросу: CPU=1 ГГц, OS=Linux должны выделяться машины с процессорной частотой не меньше 1 ГГц и операционной системой Linux.

Обратим внимание на параметр ресурсного запроса, определяющий время использования ресурсов, то есть время выполнения задания. Известно, что пользовательская оценка этого времени редко бывает точной, однако, наличие этого параметра представляется весьма важным по двум причинам. Во-первых, даже приблизительная оценка времени выполнения позволяет использовать более эффективные алгоритмы планирования ([11]). Во-вторых, существует

общепринятая практика работы в системах с разделяемыми ресурсами, согласно которой время исполнения (как и ресурсный запрос в целом) служит защитой от программных ошибок, представляя собой ограничения (по объемам и по времени) на потребляемые ресурсы: при превышении указанных в запросе лимитов, задание принудительно завершается. Наличие этого параметра не препятствуют тому, чтобы в грид могли обрабатываться задания с неограниченным временем выполнения, однако, они, по-видимому, должны рассматриваться, как особый случай.

Отметим, что параметр времени исполнения заказывается в расчете на определенную производительность ресурсов. Как в процессе планирования, так и при запуске задания должен производиться пересчет времени в соответствии с конкретными исполнительными ресурсами. Этот вопрос обсуждается в [12].

На практике получили распространение несколько разных языков ресурсных запросов. Язык RSL [13], применяемый в системе Globus Toolkit [14], ориентирован на запуск как однопроцессорных, так и многопроцессорных MPI-заданий. Языки ClassAd [15] системы Condor [16] и JDL [17] (WMS) позволяют определять альтернативные варианты ресурсного запроса и специфицировать пользовательские предпочтения при выборе ресурсов. В ряде работ предлагаются существенные расширения формализма языка запросов, направленные на спецификацию связанных заданий и цепочек заданий [18].

При планировании отбор ресурсов производится по информационной базе, содержащей сведения о составе и характеристиках ресурсов грид. Поставка этих данных в информационную базу осуществляется в оперативном режиме специализированными распределенными системами мониторинга ресурсов, из которых наибольшую распространенность получили MDS [19] и R-GMA [20]. Ресурсный запрос выступает в качестве формального критерия отбора, однако, он лишь сужает множество возможных исполнительных ресурсов. Для их однозначного определения должны приниматься во внимание дополнительные соображения.

### ***3.2. Способы организации грид и современные интерфейсы доступа к ресурсам***

Важной особенностью грид является то, что процесс обработки заданий в нем происходит на двух уровнях: глобальном и локальном. На глобальном уровне управление заданиями осуществляется программным обеспечением грид, выполняющим планирование, дистанционную передачу заданий и другие перечисленные в п. 2 функции. После того, как задание доставляется на исполнительные ресурсы, оно передается программному обеспечению локального уровня - **менеджеру ресурсов**. Менеджер ресурсов является неотъемлемой составляющей архитектуры грид, так как именно эта компонента фактически реализует поступающие дистанционные запросы,

такие как запуск и управление заданиями в локальной среде их обработки. Отметим также, что менеджер ресурсов регулирует с помощью очередей загрузку ресурсов в соответствии с политикой, которая определяется их владельцем.

Такая двухуровневая организация управления приводит к тому, что принадлежащий глобальному уровню планировщик не имеет непосредственного контроля над ресурсами. Все операции над заданиями, а также операции по получению информации о состоянии ресурсов, могут производиться только через менеджер ресурсов. Из этого следует важный вывод: возможности планирования в грид определяются теми интерфейсами, которыми располагают менеджеры ресурсов. В связи с этим рассмотрим, как устроен локальный уровень грид.

На современном этапе развития практически полностью доминируют менеджеры, управляющие кластеризованными ресурсами. При организации ресурсов в виде кластеров множество компьютеров, принадлежащих одному административному домену, включаются в грид как один узел, а в качестве менеджера выступает система пакетной обработки (СПО), обслуживающая всю совокупность компьютеров узла. Кластер в совокупности с СПО во многом аналогичен грид, однако он функционирует в замкнутой корпоративной среде и администрируется централизованно. СПО имеет интерфейсы управления заданиями и выполняет их обработку, реализуя распределение заданий по входящим в состав кластера компьютерам в соответствии с ресурсными запросами, аналогичными описанным выше. СПО ведет собственные очереди, буферизуя поступающие задания, если требуемые для их выполнения ресурсы временно заняты.

Большинство используемых на практике СПО (PBS [21], SGE [22], LSF [23], Condor [16]) были созданы задолго до появления грид и проектировались, как самодостаточные системы. Поэтому возможности взаимодействия программного обеспечения грид с такими менеджерами ограничены пользовательскими (и эквивалентными им программными) интерфейсами двух типов:

- запуска и управления заданиями;
- информирования о состоянии отдельных заданий и очередей.

Схема взаимодействия системы управления заданиями грид с СПО приведена на рис. 2. Поступающие дистанционно запросы на обработку заданий грид передаются в СПО через пользовательский интерфейс, где попадают в одну из очередей (которых в СПО как правило бывает несколько). После этого порядок выделения заданиям ресурсов определяет планировщик СПО в соответствии с устанавливаемой владельцем политикой управления.

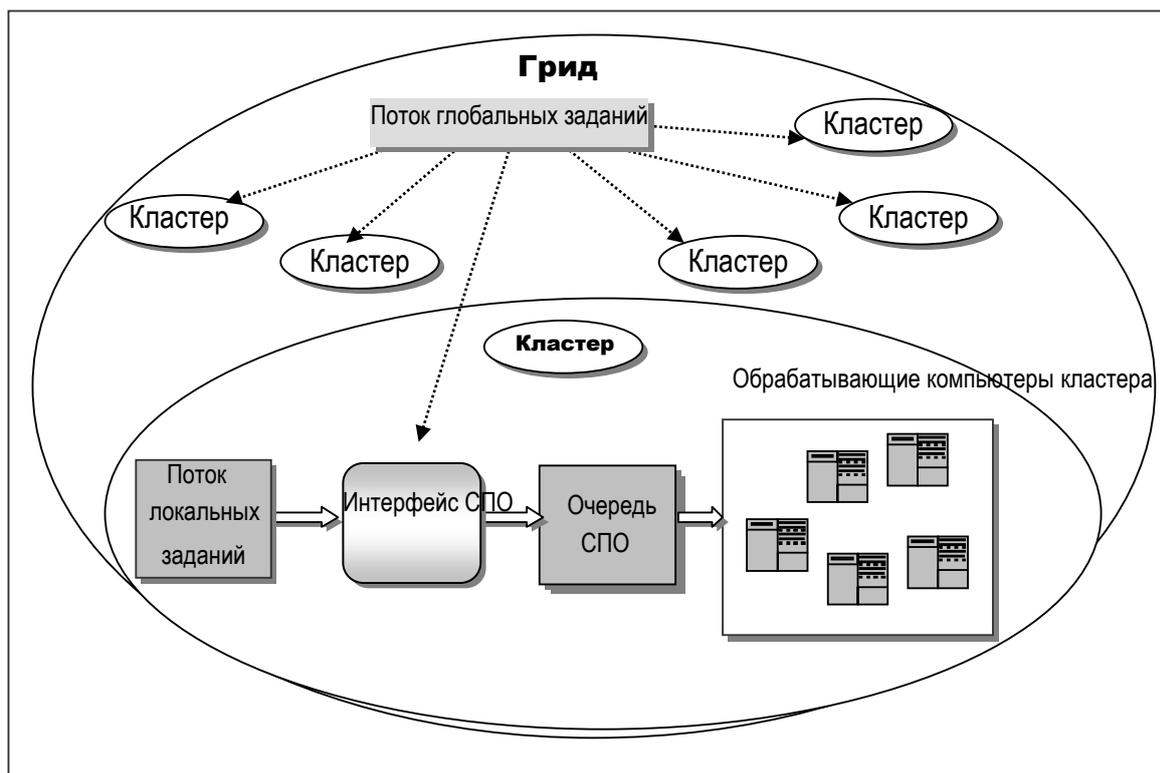


Рис. 2. Схема обработки заданий в системе пакетной обработки (СПО).

### 3.3. Режимы использования ресурсов

Сейчас, особенно в крупных проектах для поддержки научных приложений, наиболее распространен способ создания грид, когда владельцы выделяют в него свои ресурсы и не используют их сами, то есть выступают в роли профессиональных провайдеров. Этот способ имеет свою область существования, но представляет интерес также грид, в котором организации или отдельные лиц объединяют имеющиеся ресурсы, но не передают их в грид полностью – они должны совместно использоваться и в грид, и владельцами. Ресурсы, используемые в режиме разделения между владельцами и пользователями грид, мы будем называть **неотчуждаемыми**. Достоинство грид с неотчуждаемыми ресурсами в том, что, в отличие от первого способа, здесь не требуется затрат на формирование специальной ресурсной базы, и грид может создаваться динамически на ограниченный период времени с целью решению какой-либо крупной проблемы, для которой ресурсов отдельных членов кооперации не хватает.

С точки зрения планирования неотчуждаемость ресурсов усложняет ситуацию, поскольку в этом случае на них поступает два потока заданий: 1) поток из грид (глобальный), управляемый планировщиком, и 2) локальный поток заданий, которые запускаются средствами, отличными от пользовательских интерфейсов грид. Например, в условиях кластеризованных ресурсов локальные задания вводятся непосредственно через интерфейсы

СПО. Таким образом, локальные задания неподконтрольны планировщику, хотя создаваемая ими загрузка ресурсов должна учитываться планировщиком при распределении глобальных заданий.

### ***3.4. Проблемы планирования при использовании стандартных интерфейсов менеджеров ресурсов***

Стандартный набор интерфейсов менеджеров ресурсов является слишком ограниченным, для того, чтобы на их базе можно было построить детерминированное планирование. Покажем, какие проблемы порождает их использование.

**Определение исполнительных ресурсов.** В рамках интерфейсов СПО планировщик может получить только интегральную информацию о состоянии кластеров: количество заданий в очередях, общее число процессоров, число загруженных процессоров и т.п. Однако информация о характеристиках обрабатывающих компьютеров – платформа, характеристики их ресурсов – остается недоступной. Информация общего плана, которой располагает планировщик, позволяет построить аллокацию, в которой в качестве исполнительных ресурсов определен кластер, однако, не гарантируется, что в нем есть необходимые заданию компьютеры с нужной архитектурой, достаточным объемом памяти, и оно действительно может быть в нем запущено.

Планирование, основанное на такой информационной базе, может применяться только в грид с однородными по ресурсам кластерами. Если же это не так, то, по-видимому, единственно возможный способ планирования – спекулятивный запуск заданий во все кластеры грид.

**Оценка времени начала аллокации.** При построении аллокации планировщик может точно оценить время ее начала только при условии, что в кластере есть достаточное для планируемого задания количество свободных компьютеров - тогда оно может начать выполняться сразу. В противном случае в результате использования аллокации задание будет поставлено в очередь, но время его запуска на счет будет неопределенным. Говоря точнее, оно зависит от всех заданий, находящихся в кластере – в состоянии выполнения или ожидания в очереди, а также от того, какие задания будут поступать в кластер в дальнейшем.

**Возможность получения аллоцированных ресурсов.** Даже в ситуации, когда в кластер не поступают локальные задания, неопределенность времени начала аллокации приводит к задержкам в обработке глобальных заданий (и простою ресурсов), однако, при разделении ресурсов с владельцами последствия будут более тяжелыми. Если в кластер поступает неконтролируемый планировщиком локальный поток, задание грид может оставаться в очереди кластера неограниченно долго – оно может “зависнуть”.

Для полноты картины рассмотрим положение дел с планированием при ином способе организации ресурсов, который не предполагает их

кластеризации. Хотя грид с кластеризованными ресурсами получил наибольшее распространение, развивается и альтернативный, весьма перспективный по нашему мнению подход, в котором грид строится из отдельных компьютеров [24], [25]. Сравнению этих двух способов посвящена работа [26].

При включении ресурсов в грид без объединения их в кластеры менеджер ресурсов устанавливается на каждый компьютер и обслуживает только его, осуществляя локальное управление заданиями путем непосредственного взаимодействия с ОС. При таком “индивидуальном” подходе информационные интерфейсы менеджера позволяют предоставлять планировщику полный набор данных (объемы и характеристики ресурсов), необходимых для того, чтобы отбирать исполнительные компьютеры в точном соответствии с ресурсным запросом. Решается также проблема гарантированного запуска заданий. Если в варианте с кластеризованными ресурсами время начала аллокации определяется путем опроса состояния кластера через интерфейсы СПО, что приводит к дефектам планирования, то менеджер компьютера сам определяет возможность выполнения задания грид и выдает планировщику запрос на его получение.

Основная проблема грид с некластеризованными ресурсами – непредсказуемость времени окончания задания. Такой грид наиболее выгоден, если ресурсы используются в режиме разделения с владельцами. Разделение ресурсов происходит здесь непосредственно во время выполнения: ресурсы, в том числе и процессорные, делятся между заданием грид и приложениями, выполняемыми владельцем компьютера. Владелец компьютера устанавливает политику разделения ресурсов таким образом, чтобы обработка задания грид не оказывала существенного влияния на его деятельность, то есть, фактически, разрешает занимать компьютер в периоды простоя. Помимо того, компьютеры могут выключаться в произвольные моменты времени. Эти обстоятельства приводят к тому, что реальное время обработки задания на таких ресурсах будет существенно отличаться от процессорного времени, заказанного в ресурсном запросе. Проблема планирования в этом случае - сделать реальное время выполнения приемлемым для пользователя.

#### **4. ОСНОВНЫЕ ПОЛОЖЕНИЯ ПРЕДЛАГАЕМОГО МЕТОДА ПЛАНИРОВАНИЯ**

Как показано в предыдущем разделе, возможности планирования в условиях кластеризованных неотчуждаемых ресурсов на основе интерфейсов современных СПО весьма ограничены. Используемые на практике брокеры ресурсов, например, GRB [8] системы WMS и брокер платформы ARC/NorduGrid [27], имеют довольно жесткие ограничения по применению и не способны исключить такие нежелательные эффекты, как непредсказуемость времени обработки заданий, задержка обработки в ситуациях, когда имеются простаивающие ресурсы. Существенный недостаток этих систем -

невозможность планирования многопроцессорных параллельных заданий с синхронным выделением ресурсов в нескольких кластерах.

Предлагая новый подход к планированию, мы исходим из того, что грид, по-видимому, допускает многообразные формы воплощения в зависимости от мотивов его создания, организации ресурсов, предметной ориентации, способа использования и иных обстоятельств. Многообразие форм грид предполагает и различие в методах планирования, поэтому, ориентируясь главным образом на кооперативный грид некоммерческих организаций, определим условия применения и основные принципы, которые характеризуют предлагаемый подход.

**Ресурсы и задания.** Ресурсы грид организованы в кластеры, управляемые СПО, и используются совместно с владельцами, то есть рассматривается грид с кластеризованными неотчуждаемыми ресурсами. Объектами планирования являются однопроцессорные и многопроцессорные задания, причем выполнение последних может происходить как в пределах одного кластера, так и на ресурсах нескольких кластеров.

**Стратегия планирования.** Отбор ресурсов для заданий на основе ресурсных запросов не в полной мере определяют способ планирования, полностью оставляя в стороне вопрос о согласовании распределения ресурсов между разными заданиями. Необходимые дополнительные условия, которые можно назвать стратегией планирования, часто удается формализовать в виде целевой функции, оценивающей качество планирования для различных вариантов распределения заданий на некотором временном интервале. Планирование в этом случае сводится к оптимизационной задаче. Такой способ применим и к грид, где в качестве оптимизационного критерия могут выступать максимизация загрузки ресурсов, дохода владельцев или минимизация расходов на оплату ресурсов со стороны пользователей ([28]).

В любом случае выбор стратегии должен исходить из определенных соглашений по разделению ресурсов. В грид это соглашения двух типов: 1) между пользователями и 2) между пользователями и владельцами ресурсов. В качестве основы для выработки этих соглашений наиболее перспективным представляется экономический подход [29], в котором делается попытка формализовать процесс распределения ресурсов с помощью модели рынка: на рынке ресурсов грид в качестве продавцов выступают владельцы ресурсов, а в качестве покупателей – пользователи. В рамках этой модели рассматриваются проблемы ценообразования, обеспечения надежности в расчетах, представления политик переговоров, а также применения соответствующих программных разработок [30] в различных операционных компонентах программного обеспечения грид, в том числе в брокерах ресурсов [31]. Основываясь на основных идеях экономического подхода, рассмотрим далее конкретный вид предлагаемых соглашений по разделению ресурсов и способ их поддержки при планировании.

#### ***4.1. Разделение ресурсов между пользователями***

Наш выбор стратегии разделения ресурсов между пользователями направлен на то, чтобы каждый пользователь мог контролировать полное время обработки своих заданий и, в том числе, управлять скоростью получения ресурсов на этапе планирования. Конечно, прямой заказ времени получения ресурсов невозможен, поскольку за ресурсы конкурирует множество заданий, однако, если задействовать механизм приоритетов и позволить пользователю устанавливать приоритеты своих заданий, то тем самым у него появится средство для изменения их положения в очереди на получение ресурсов.

Существенен вопрос, как совместить присвоение приоритетов заданиям самим пользователем с принципом “справедливого” разделения ресурсов между пользователями. В локальных системах расстановка приоритетов осуществляется административно, однако для грид такое простое решение не годится. В рамках экономического подхода эквивалентом приоритета может выступать плата за задание, вносимая пользователем. Изменяя плату, пользователь регулирует время старта задания. Отметим, что платность не исключает возможности обслуживания и “бесплатных” заданий, но они будут получать ресурсы в последнюю очередь.

Важно, ведутся ли расчеты реальными (в коммерческих виртуальных организациях) или условными деньгами (в кооперативных виртуальных организациях). В последнем варианте возможность оперативного управления отдельными заданиями должна быть дополнена интегральными ограничениями на общее количество получаемых пользователем ресурсов. Проблема “справедливого” (fair-share) разделения ресурсов возникает и в корпоративных системах коллективного пользования, а некоторые решения [32] применимы в условиях грид. При использовании экономического подхода виртуальная организация выделяет каждому пользователю ограниченный бюджет на определенный период времени, из которого он вносит плату при запуске задания. Определение величин бюджетов для множества пользователей грид производится иерархическим способом, аналогично, например, определению прав доступа к ресурсам. Совокупность двух механизмов – оплаты заданий и ограничения бюджета, - позволяя управлять скоростью обработки каждого отдельного задания, не дает монополизировать ресурсы отдельному пользователю.

Подводя итог, скажем, что в качестве стратегии планирования в аспекте разделения ресурсов между пользователями мы выбираем приоритетное обслуживание заданий. Этот выбор определяет **архитектуру** системы планирования: централизованная архитектура (в противоположность распределенной [33]) с общей очередью глобальных заданий.

#### 4.2. Разделение ресурсов с владельцами

Соглашение о разделении ресурсов с владельцами является ключевым принципом предлагаемого метода планирования. В рассматриваемой ситуации с неотчуждаемыми ресурсами имеется два потока заданий – глобальный и локальный, и соглашение о разделении ресурсов определяет, какие из распределяемых ресурсов и в какое время могут быть отданы тому или иному глобальному заданию. Соответствующие правила должны быть специфицированы таким образом, чтобы планировщик грид был в состоянии установить возможность занятия ресурсов заблаговременно, в момент построения аллокации.

На практике применяются два варианта разделения ресурсов между локальным и глобальным потоками.

1. Ресурсы кластера разделяются на две непересекающиеся части: для глобальных заданий и для локальных заданий. Система очередей СПО конфигурируется таким образом, чтобы задания каждого сорта распределялись только на свою часть ресурсов.

2. Ресурсы кластера образуют одно множество. Задания грид запускаются в СПО с некоторым фиксированным приоритетом.

И первый, и второй вариант имеют недостатки и не позволяют эффективно управлять разделением ресурсов. В первом варианте одна часть ресурсов будет простаивать, если нет распределяемых на эту часть заданий, в то время как ресурсов для обслуживания заданий другого сорта может не хватать. Имеющаяся статистика показывает [34], что при таком статическом разделении средний уровень использования ресурсов не превышает 25%. Во втором варианте время получения ресурсов глобальными заданиями будет неопределенным, так как оно зависит от приоритетов локальных заданий, находящихся в СПО. Если же глобальным заданиям присваивается максимально возможный приоритет, то они всегда будут иметь преимущественное право получения ресурсов, и при интенсивном их потоке обработка локальных заданий может быть заблокирована полностью.

Мы предлагаем соглашение о разделении ресурсов, которое осуществляет динамическую балансировку потоков заданий. Соглашение состоит в том, что ресурсы, которые получает глобальное задание, должны оплачиваться, и глобальное задание может занять ресурсы при условии, что назначенная пользователем при его запуске плата не меньше их стоимости. Принципиальной особенностью этого соглашения является то, что стоимость ресурсов меняется во времени и определяется не только их объемами, а вычисляется с учетом приоритетов локальных заданий, которые могут занять те же ресурсы. Цена использования ресурсов в единицу времени задается формулой  $PR(t) = C * P(t) * F(\mathbf{R})$ , где  $F(\mathbf{R})$  – неубывающая функция, зависящая от объема ресурсов  $\mathbf{R}$  каждого типа, которыми располагает компьютер, а  $P(t)$  – приоритет локального задания, которое должно выполняться на компьютере во время  $t$ . Это дает гибкий механизм балансировки потоков, позволяющий, с

одной стороны, администратору кластера путем изменения цены ресурсов (коэффициент  $C$ ) ограничивать интенсивность потока внешних заданий (при достаточно больших значениях коэффициента  $C$  глобальные задания вообще не будут поступать в кластер), и, с другой стороны, дающий возможность пользователю грид влиять на скорость получения ресурсов для своих заданий, повышая плату.

Поскольку стоимость ресурсов определяется тем, какие локальные задания претендуют на их получение, планировщик грид должен заблаговременно снабжаться информацией, описывающей будущее состояние каждого компьютера грид.

## 5. ОБЕСПЕЧЕНИЕ ДЕТЕРМИНИРОВАННОСТИ ПЛАНИРОВАНИЯ

Предлагаемый метод планирования, базирующийся на рассмотренных принципах, мы называем **опережающим** планированием, так как ключевой механизм в нем - опережающее события в кластерах построение точных аллокаций на основании прогноза состояния ресурсов в будущем. Цель разработки нового метода - сделать планирование детерминированным. Для этого требуется, чтобы построенные планировщиком точные аллокации могли быть реализованы, то есть чтобы по ним мог быть осуществлен запуск заданий на исполнительных ресурсах. Обеспечить это возможно, но, как показано в п. 3.4, необходимо расширение механизмов взаимодействия системы планирования с менеджерами ресурсов. Вопрос о развитии функциональности СПО и соответствующем расширении их интерфейсов назрел, и есть достаточно большое количество предложений на этот счет. Их систематизации посвящена работа [35]. Два новых механизма: прогноз загрузки и предварительное резервирование ресурсов, на основе которых строится метод опережающего планирования, лежат в русле этих предложений.

### 5.1. Прогноз загрузки кластерных ресурсов

Опережающее планирование опирается на прогноз распределения локальных заданий по ресурсам кластеров – расписание загрузки компьютеров. Это расписание описывает, какие задания будут занимать компьютеры в любой момент времени из некоторого временного диапазона. Расписание состоит из множества слотов вида:

Слот := [Компьютер], [Локальное задание], [Начало, Конец]

Каждый слот определяет интервал времени [Начало, Конец], в течение которого компьютер выделяется локальному заданию. Совокупность интервалов, относящихся к одному компьютеру, может покрывать временной диапазон расписания лишь частично: промежутки между слотами соответствуют периодам, когда компьютер не занят.

Прогноз, содержащийся в расписании, зависит от алгоритма планирования СПО, параметров политики управления, а для его получения требуется учет характеристик (ресурсных запросов, включая заказанное время исполнения, а также приоритетов) всех заданий, находящихся в кластере. По-видимому, единственный способ построения прогноза такого рода – моделирование процесса обработки заданий в кластере.

Хотя ни одна из известных СПО не имеет готовых средств для построения расписаний, возможность реализации такой функциональности была нами показана [36] на примере кластерного планировщика Maui [37]. Реализация выполнена на базе входящей в состав Maui подсистемы моделирования. Исходное назначение этой подсистемы – оценка эффективности загрузки ресурсов: получая на вход модельный поток заданий и используя текущие настройки кластерного планировщика, она выдает повременное распределение заданий по ресурсам. Реализация построения расписаний свелась, в основном, к внедрению в цикл планирования кода, фиксирующего текущее состояние кластера, запускающего процесс моделирования и передающего ему исходные данные. Опыт этой разработки позволяет утверждать, что создание компоненты построения расписаний в СПО любого типа, в том числе и с отличными от Maui алгоритмами планирования, возможно: задача сводится к эмуляции управления очередями и выполнения заданий на ресурсах.

Для того, чтобы планировщик грид мог вычислять стоимость ресурсов, требуется несколько иная форма расписаний. Поэтому перед передачей планировщику грид слоты расписания преобразуются так, что в них вместо заданий указывается вычисляемая с учетом приоритетов тех же заданий цена ресурсов в единицу времени:

Слот := [Компьютер], [Цена ресурсов], [Начало, Конец]

Построенные в разных кластерах расписания, которые мы будем называть **локальными расписаниями**, поставляются в информационную базу планирования, причем поставка ведется в оперативном режиме, отражающем происходящие изменения. Напомним, что в ту же базу заносятся данные о характеристиках компьютеров, но эта информация более статична. По информационной базе планировщик грид производит размещение глобальных заданий, строя для них аллокации ресурсов. В планировщике может использоваться любой алгоритм (важно только, чтобы он был модифицирован с учетом соглашения о разделении ресурсов), а само планирование может вестись на любую глубину по времени, согласованную с временным диапазоном расписаний. Так, например, если планировщик обслуживает только однопроцессорные задания, в качестве алгоритма может выступать, например, FCFS [38], и распределять достаточно только те компьютеры, на которых освобождение ресурсов, то есть смена заданий, прогнозируется на ближайшее время. Однако, в полной мере достоинства опережающего

планирования проявляются в том, что его можно применить для многопроцессорных заданий.

Одна из наиболее трудных проблем планирования в грид – **коаллокация**: синхронное выделение ресурсов в нескольких кластерах для многопроцессорных заданий. Известно, что простые алгоритмы приоритетного планирования не способны обеспечить запуск заданий, требующих большого количества ресурсов, независимо от того, насколько высок их приоритет. Причина заключается в фрагментации ресурсного пула из-за того, что ресурсы постоянно выделяются более мелким заданиям. Для предотвращения зависания многопроцессорных заданий необходимо каким-то образом накапливать ресурсы для них. Лучший на сегодня алгоритм для планирования многопроцессорных заданий – алгоритм “обратного заполнения” BackFill [39], разработанный для систем типа MPP и достаточно широко применяющийся на практике. Идея этого алгоритма состоит в том, что ресурсы выделяются заданиям не непосредственно в момент освобождения, а заблаговременно: планировщик строит план распределения ресурсов – расписание запусков заданий. При построении расписания ресурсы распределяются заданиям в порядке их приоритетов, причем задание может получить некоторые ресурсы только при условии, что они уже не отведены более приоритетным заданиям. Алгоритм гарантирует получение ресурсов высокоприоритетными заданиями в минимально возможное время и, в то же время, допускает нарушение порядка очереди, позволяя низкоприоритетным заданиям занимать оставшиеся неиспользованными ресурсы, что способствует повышению коэффициента общей загрузки ресурсов.

Адаптация алгоритмов планирования, разработанных для локальных систем, в том числе и алгоритма BackFill, к условиям грид опирается на прогноз загрузки кластерных ресурсов. Планировщик грид должен размещать глобальные задания с учетом распределения ресурсов между локальными заданиями кластеров, используя локальные расписания. Главный пункт адаптации – применение критерия, соответствующего соглашению о разделении ресурсов, при определении возможности выделения ресурсов глобальному заданию.

Несмотря на перспективность использования расписаний, реализация основанного на них планирования имеет свои трудности. Содержащийся в локальном расписании прогноз получается в результате работы реального планировщика СПО, соответствует параметрам политики управления, и в этом отношении является точным. Тем не менее, расписание представляет собой лишь ориентировочный (tentative) прогноз, поскольку ход реальной обработки будет совпадать с расписанием лишь при условии, что в кластер не будут поступать новые задания, и не будет меняться состав ресурсов. Есть и еще одна причина возможных отклонений: расписание строится, исходя из оценки времени выполнения, заданной в ресурсных запросах заданий, которая редко бывает точной.

В реальности, непредсказуемые события, такие как преждевременное окончание заданий, будут влиять на расписание. Для того, чтобы планирование в грид соответствовало состоянию ресурсов, локальные расписания должны регулярно обновляться и, соответственно, должен перестраиваться план распределения глобальных заданий. В связи с этим возникает вопрос о способе поставки расписаний. Хотя в современных информационных службах грид в основном используется метод **опроса** (pull) источников данных, применительно к интерфейсу поставки расписаний более предпочтительным может быть метод **нотификации** (push) об изменениях в расписаниях кластеров или даже комбинация этих методов.

Укажем еще одну форму прогноза о состоянии ресурсов, которая в какой-то мере альтернативна расписаниям, – это список возможных аллокаций для отдельных глобальных заданий. В этом варианте планировщик грид опрашивает расширенные соответствующей функцией менеджеры кластеров, посылая им ресурсный запрос глобального задания и получая в ответ временные слоты, в которые оно может получить ресурсы. Списки возможных аллокаций могут получаться таким же образом, как и расписания – путем моделирования работы менеджера кластера. Отметим также, что если возможные аллокации основаны на бинарном прогнозе (слот занят/свободен локальными заданиями), как в [40], то их едва ли можно использовать в грид с неотчуждаемыми ресурсами: по нашему мнению разделение ресурсов между локальными и глобальными заданиями должно быть более гибким.

## ***5.2. Предварительное резервирование ресурсов***

Предварительное резервирование ресурсов – один из немногих новых механизмов, реализованных в практических системах, хотя примеры реализаций немногочисленны: кластерный планировщик Maui и коммерческий вариант СПО PBS - PBSPro [41]. Соответствующее этому механизму расширение интерфейса СПО позволяет заказывать какое-то количество ресурсов на определенный в запросе интервал времени, отводя их исключительно для выполнения заданий некоторого пользователя или группы пользователей [42]. Заметим, что выполнение резервирований обеспечивается планировщиком кластера, который модифицирует распределение ресурсов таким образом, чтобы к моменту наступления времени резервирования, заказанные ресурсы были свободны.

В связи с применением в грид предприняты усилия по разработке и стандартизации дистанционных интерфейсов предварительного резервирования, которые надстраиваются над локальными интерфейсами СПО. В этой связи важными представляются предложения, выработанные в специализированном комитете Global Grid Forum [43]. Стандартизация направлена на обеспечение независимости от интерфейсов конкретных СПО и, помимо того, унифицирует работу с различными видами ресурсов: помимо вычислительных, рассматриваются сетевые ресурсы и ресурсы хранения.

Предложенный аппарат позволяет описывать разнообразные сценарии резервирования и включает такие, например, команды как отказ и подтверждение резервирования, уточнение состава резервируемых ресурсов.

Очевидное применение предварительного резервирования для планирования в грид –гарантировать выделение ресурсов под построенные планировщиком аллокации. Однако, сам по себе механизм резервирования не решает вопросы разделения ресурсов и не определяет, какие из запросов на резервирование должны удовлетворяться. Основная цель использования прогноза загрузки ресурсов как раз состоит в том, чтобы создать условия, при которых планировщик грид строит аллокации, удовлетворяющие соглашениям о разделении ресурсов. Запросы на ресурсы по таким аллокациям должны, естественно, удовлетворяться, и предварительное резервирование позволяет это сделать.

Сценарий использования предварительного резервирования состоит в следующем. Когда подходит время доставки какого-то задания, планировщик посылает в кластер запрос на резервирование, который в штатной ситуации должен безусловно выполняться. Затем происходит передача задания в кластер, и СПО обеспечивает его запуск в запланированное время и на запланированных ресурсах. При определенных условиях возможны отказы при резервировании. Это, однако, обнаруживается планировщиком, и задание никуда не перемещается, а поступает на перепланирование. В грид, как распределенной среде, важно обеспечить надежность доверительных отношений между планировщиком грид и кластером. Для этого при обработке в кластере запроса на резервирование может производиться проверка его корректности в смысле соответствия соглашению о разделении ресурсов.

## 6. СХЕМА ОПЕРЕЖАЮЩЕГО ПЛАНИРОВАНИЯ

Представим схему опережающего планирования следующим псевдокодом.

1. Цикл планирования:
  - {
  - 2. Получение стоящих в очереди заданий грид;
  - 3. Получение расписаний всех кластеров;
  - 4. Цикл размещения заданий:
    - {
    - 5. Построение аллокации для задания;
    - 6. Если время начала аллокации близко:
    - 7. { резервирование; доставка задания в кластер;}
    - }
  - }

Планирование представляет собой циклический процесс, в каждом цикле которого строится ориентировочное распределение заданий по ресурсам и определяется подмножество заданий, для которых начало выполнения достаточно близко. Для этих заданий выполняется резервирование ресурсов, гарантирующие их выделение в соответствии с построенной планировщиком аллокацией, и инициируется доставка заданий в кластеры.

Первые два оператора цикла планирования получают из информационной базы стоящие в глобальной очереди задания (оператор 2) и множество локальных расписаний (3). Заполнение информационной базы осуществляется отдельным процессом, который может работать параллельно с процессом планирования. Каждый цикл планирования, однако, происходит при фиксированном состоянии информационной базы, а происходящие за время одного цикла изменения состояния учитываются на следующем. Локальные расписания объединяются в общий массив слотов, по которому ведется планирование глобальных заданий.

Собственно планирование заключается в построении аллокаций для заданий в порядке их приоритетов. Обозначим  $N$  требуемое заданию количество процессорных узлов, и  $T$  - длительность выполнения задания. Задача построения аллокации (5) состоит в том, чтобы найти такие  $N$  процессорных узлов и такой интервал времени  $[t_0, t_0+T]$ , на котором ресурсы этих процессорных узлов могут быть выделены данному заданию. Возможность выделения определяется следующими условиями:

- характеристики ресурсов каждого процессорного узла соответствуют ресурсному запросу задания,
- суммарная стоимость ресурсов, вычисляемая по слотам, попадающим в интервал действия аллокации, не превышает платы за ресурсы, назначенной пользователем.

После построения аллокации могут быть инициированы действия (7) по обеспечению запуска задания: предварительное резервирование ресурсов и передача задания в кластеры. Эти действия выполняются только для заданий, время начала аллокаций которых достаточно близко к моменту планирования. Для остальных заданий эффект построения аллокации сводится к тому, что слоты, по которым она построена, становятся полностью или частично недоступными на протяжении того же шага планирования для других глобальных заданий. Заметим, что на последующих шагах такие задания могут получить совсем другие ресурсы.

Способ отбора аллокаций (6), для которых необходимо инициировать запуск заданий, иллюстрирует рис. 3.

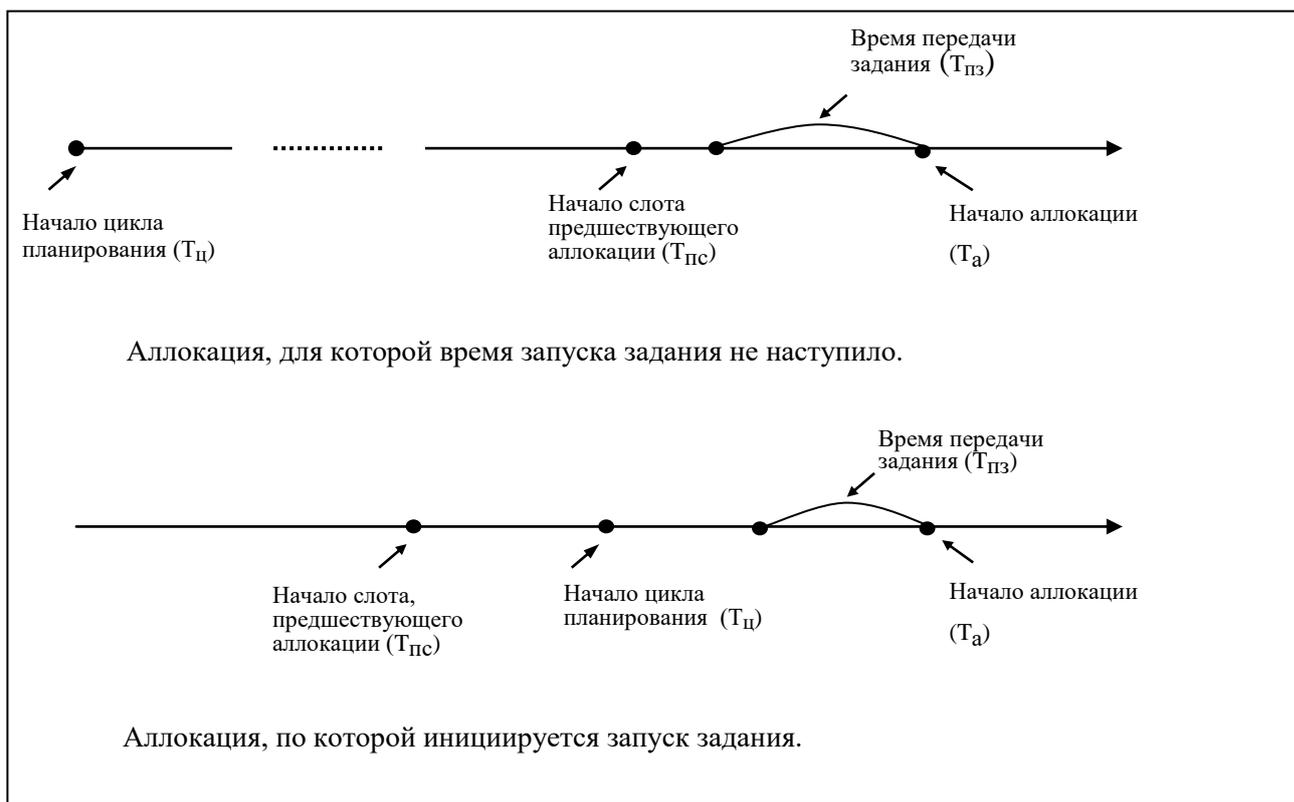


Рис. 3. Выбор времени запуска задания.

Выбор момента резервирования отражает компромисс между стремлением основывать реальное выделение ресурсов на оперативной информации о состоянии грид и необходимостью гарантировать запуск заданий. Если резервирования делаются рано, например, непосредственно в момент поступления заданий в глобальную очередь, то изменения состояния (появление новых заданий, отклонение хода обработки заданий от первоначального прогноза) за все время их нахождения в очереди никак не будут отражаться, например, на порядке получения ресурсов, то есть не будет работать механизм приоритетов.

Что же касается минимального упреждения, необходимого для обеспечения запуска задания, то оно определяется двумя условиями. Первое условие: упреждение должно быть таким, чтобы ресурсы аллокации не оказались занятыми локальными заданиями. При этом надо учесть, что эти ресурсы могут освободиться раньше прогнозируемого в расписании времени. Если происходит такое “преждевременное” освобождение ресурсов, на них может начать выполняться локальное задание, окончание которого происходит после начала аллокации. В связи с этим, первое условие начала доставки задания задается неравенством:  $T_{пс} > T_{ц}$ , где  $T_{пс}$  – начало слота, непосредственно предшествующего аллокации на тех же ресурсах,  $T_{ц}$  – время начала цикла планирования.

Второе условие заключается в том, что должно хватать времени для передачи задания на исполнительные ресурсы. Процесс передачи включает

пересылку всех его входных файлов, так что время передачи может быть значительным. Это условие можно записать в виде:  $(T_a - T_{ц}) > T_{пз}$ , где  $T_a$  – время начала аллокации,  $T_{пз}$  – время передачи задания.

Первое из сформулированных условий позволяет с высокой вероятностью рассчитывать, что ресурсы для построенных планировщиком аллокаций будут выделены в кластере. Тем не менее, отказы в резервировании возможны: ресурсы все же могут оказаться занятыми локальным заданием. Такое, однако, может произойти лишь в том случае, когда в промежуток между построением аллокации и резервированием аварийно (в момент старта) завершается несколько заданий, предшествующих началу аллокации. Эта ситуация рассматривается как исключительная, а задание поступает на перепланирование.

### ***6.1. Масштабируемость системы планирования***

Способность системы планирования обслуживать большое количество компьютеров – масштабируемость – является ограниченной, как и у всякой системы с централизованной архитектурой. В связи с этим важными представляются вопросы оценки возможностей рассматриваемой схемы планирования, диагностики адекватности функционирования, а также разработка способов повышения масштабируемости.

Штатный режим работы планировщика – когда он, обслуживая поток поступающих в глобальную очередь заданий, успевает своевременно реагировать на изменения, происходящие в кластерах. Заметим, что в отличие от локальных СПО, планировщик не обрабатывает каждое отдельное событие. Вместо этого он периодически получает от кластеров обновление расписаний, и на каждом цикле планирования перестраивает заново аллокации для глобальных заданий. Нарушение штатного режима происходит, когда отставание  $(T_{ир} - T_{пр})$ , где  $T_{пр}$  – время производства расписания в кластере,  $T_{ир}$  – время использования расписания при планировании, становится большим. Как меру допустимого отставания введем понятие **времени жизни расписания** (TTL – Time To Live): планировщик не должен использовать расписание, время жизни которого истекает ко времени конца цикла планирования.

Для каждого кластера значение TTL может устанавливаться индивидуально, исходя из смысла этого параметра: значение TTL определяет, с каким предельным отставанием планировщик должен реагировать на изменения состояния кластера. При большом TTL получение ресурсов для вновь поступающих кластерных заданий будет иметь тенденцию к задержке: полагающиеся им ресурсы могут занять глобальные задания. Однако, задержка всегда ограничена, и, если, например, TTL устанавливается равным минимально возможной длительности задания, то локальное задание пропустит не более одного глобального задания вне очереди, причем задерживаться будут только такие локальные задания, которые имеют самый высокий приоритет в локальной очереди и могут быть сразу запущены.

Необходимое условие правильного функционирования заключается в том, что длительность цикла планирования должна быть меньше TTL всех расписаний. Исходя из этого условия, дадим оценку максимальной производительности планировщика, отмечая, что анализ всей системы планирования требует учета ряда факторов: скорости генерации локальных расписаний, производительности сетевые операции по доставке расписаний и управлению заданиями.

В цикле планирования при построении аллокаций выполняется сопоставление глобальных заданий, находящихся в очереди, со слотами агрегированного по всем кластерам расписания. Обозначая количество глобальных заданий  $N_z$  и число слотов  $N_s$ , представим длительность цикла планирования в виде  $CP = A_1 * N_z * N_s$ . Здесь  $A_1$  – время одной операции сопоставления, которое мы оцениваем  $A_1 = 10^{-2}$  мск. Чтобы получить оценку зависимости  $CP$  от числа компьютеров в грид  $N$ , будем исходить из того, что  $N_z = C_2 * N$  и  $N_s = C_3 * N$ , где константы  $C_2, C_3 < 5$  (см. ниже). Тогда  $CP = A_1 * C_2 * C_3 * N^2$ . При  $N = 1000$   $CP < 10^{-2} * 25 * 10^6$  мск =  $25 * 10^4$  мск = 4 мин.

С учетом вышеприведенное определение TTL, должно выполняться неравенство  $TTL > 2 * CP$ , и полагая  $TTL = L$ , где  $L$  минимальная длина задания, получим ограничение масштабируемости:  $N < \sqrt{L / (2 * A_1 * C_2 * C_3)}$ . При  $L = 30$  мин. планировщик способен обслуживать несколько тысяч компьютеров.

В этих выкладках нуждается в разъяснениях предположения о количестве глобальных заданий и количестве слотов. Смысл первого предположения ( $N_z = C_2 * N$ ) – в том, что количество запускаемых в грид заданий должно быть сбалансировано с имеющимися ресурсами: как предполагается в оценке, в пересчете на один компьютер длина очереди не превышает 5 заданий.

Второе предположение ( $N_s = C_3 * N$ ) исходит из того, что поставляются не полные, а ограниченные расписания, в которых прогнозируется распределение только тех локальных заданий, которые находятся в голове очереди своего кластера. Если планирование выполняется только для однопроцессорных заданий, длина расписаний может быть ограничена двумя слотами на каждом компьютере – такого прогноза хватит, чтобы построить аллокацию и заранее произвести резервирование ресурсов. В случае многопроцессорных заданий теоретически нужно иметь расписания полного объема. Тем не менее, при выборе  $C_3$  равным небольшому целому числу (4-5), будет во-первых, гарантировано выделение ресурсов самым приоритетным глобальным заданиям, и, во-вторых, снижение эффективности алгоритма Backfill, выражающееся в недополучении ресурсов менее приоритетными заданиями, как мы рассчитываем, будет незначительным.

## 7. СРАВНЕНИЕ С СОВРЕМЕННЫМИ ПЛАНИРОВЩИКАМИ

Не претендуя на полный обзор довольно многочисленных реализаций систем управления заданиями в грид, появившихся в последнее время, рассмотрим их возможности планирования.

Большинство разработок предназначены для обслуживания кластеризованных ресурсов, в том числе система WMS [8] наиболее крупного на сегодняшний день грид проекта EGEE [44], имеющего мировой масштаб. Брокер системы WMS реализует централизованную схему планирования, работающую в двух режимах распределения заданий: жадного и ленивого (ленивый режим был впервые реализован в проекте Alien [45]). В жадном режиме распределение заданий происходит без глобальной очереди: поступившее в брокер задание сразу же отправляется в некоторый кластер. Для выбора кластера используется информационная база, содержащая данные о текущем состоянии всех ресурсов грид. Задания, запускаемые пользователями в ленивом режиме поступают в глобальную очередь. Кластеры, имеющие свободные ресурсы, информируют брокер о характеристиках свободных ресурсов, а тот, выбирая из очереди подходящее задание, передает его в кластер.

В обоих режимах происходит сопоставление характеристик ресурсов и заданий на предмет выяснения возможности выполнения. Информационная база жадного режима, построенная на информационной службе MDS [19] системы Globus Toolkit [14], недостаточно полна: она не содержит данных, описывающих процессорные узлы, так что этот режим может использоваться только в грид с однородными ресурсами. Сопоставление характеристик ресурсов и заданий в ленивом режиме производится на основе метода *matchmaking* [15] и не имеет такого недостатка.

Основное решение планирования – выбор исполнительных ресурсов – в жадном режиме делается по интегральной характеристике состояния кластера: выбирается тот, у которого в локальной очереди меньше всего заданий. Как было показано, это потенциально ведет к зависанию заданий, тем более, что глобальная очередь в этом режиме не поддерживается и распределение происходит в момент поступления задания в брокер. Подход ленивого режима более аккуратный – задания грид находятся в очереди, а кластер запрашивает у брокера новое задание в моменты времени, определяемые политикой, заданной администратором. Однако, вопрос о том, какой должна быть эта политика остается открытым. Простейший подход – запрашивать задание при пустой локальной очереди – не разрешает проблем неоднородного кластера и неотчуждаемых ресурсов.

Система управления заданиями ARC [27] проекта NorduGrid представляет интерес как пример децентрализованной архитектуры планирования. В целом, схема планирования в ARC аналогична жадному режиму брокера WMS, однако планировщик здесь входит в состав программного обеспечения каждого рабочего места пользователя, а планирование и передача задания в кластер, выполняется непосредственно с места выдачи запроса. Преимущества децентрализованной архитектуры известны – это отсутствие “критической точки”, живучесть и масштабируемость. Отметим также, что в данной

архитектуре планирование производится отдельно для каждого задания и принимаемые решения могут оказаться не согласованными друг с другом.

В некоторых системах находят применение два новых механизма – резервирование и миграция. Сценарий Moab Grid Scheduler [40] выглядит следующим образом:

- задания грид поступают в глобальную очередь,
- планировщик посылает запросы в кластеры и получает от них время возможных аллокаций,
- исходя из полученной информации, выбирается кластер, производится резервирование ресурсов, задание отправляется в кластер.

Moab, также как и CSF (Community Scheduler Framework) [46], который также поддерживает резервирование в службе планирования, не дает все-таки полного решения, предполагая, что кластерное обеспечение имеет средства для принятия решений о том, когда могут быть отведены ресурсы для резервирования. В стандартных СПО таких средств нет, и их создание составляет отдельную задачу.

Системы GridWay [47] и GridLab Resource Management System (GRMS) [48] используют механизм перераспределения заданий - миграцию. В GridWay после того как глобальное задание отправлено в кластер специальная служба осуществляет его мониторинг. В случае, если задание не запущено на счет в локальном менеджере ресурсов за некий интервал времени, происходит перепланирование данного глобального задания, и оно перемещается на другие ресурсы.

Миграция представляет собой мощный механизм, который способен расширить возможности планирования. Однако, его применение сопряжено с большими временными затратами на перемещение файлов заданий, так что миграция, по-видимому, может рассматриваться как дополнение к качественному способу планирования.

Приведем также ряд исследовательских работ, в которых рассматриваются вопросы, близкие к теме данной работы. Идея использования прогноза для выбора ресурсов грид не является новой. Одна из первых работ такого рода – система Gallop [49], планирование в которой осуществляется на основе предсказания свободной вычислительной мощности по статистическим данным, для чего на каждом кластере Gallop устанавливаются специальная компонента (Prophet), вырабатывающая локальный прогноз. Значительное количество публикаций посвящено прогнозированию загрузки ресурсов на основе статистической информации (см., например, [50]).

Детерминированное планирование рассматривается в работе [12], где исследуется эффективность использования предварительного резервирования для многопроцессорных заданий. В предложенной модели планирования прогноз ограничен одним шагом - временем освобождения ресурсов, которое оценивается по заказанному времени выполнения заданий.

В [51] обосновывается необходимость расширения функциональности современных систем планирования функциями предварительного резервирования ресурсов и составления расписаний запуска заданий. Авторы проводят различие между системами планирования двух типов: традиционными - *Queuing systems* и системами, основанными на расписаниях - *Planning systems*. В то время как первые распределяют задания, исходя из текущего состояния ресурсов, вторые основываются на полноценном плане на будущее. По мнению авторов такой подход, реализованный в системе CCS [52], способствует применению аппарата предварительного резервирования и открывает перспективу для решения задачи управления параллельными заданиями. Разделяя эти общие положения, отметим, что в архитектуре CCS только один уровень управления: все задания распределяются через CCS.

Наиболее близко к нашей работе стоит описанный в [34] планировщик *Ursala*. Для него предложен способ размещения заданий, в том числе параллельных, основанный на механизмах предварительного резервирования и списках возможных аллокаций. Получая такие списки от кластеров, *Ursala* подбирает такое временное окно, в котором свободных ресурсов достаточно для запуска задания, и резервирует их. В работе отмечается способность локального планировщика *Maui* строить расписание загрузки ресурсов кластера. Сравнивая нашу работу с *Ursala*, отметим, что мы не ограничиваемся размещением отдельных заданий: предложена динамическая схема приоритетного планирования очереди заданий и предложен конкретный способ балансировки потоков локальных и глобальных заданий, основанный на плате за ресурсы.

В работе [33] представлена распределенная архитектура планирования с множественными экземплярами планировщиков, управляющих непересекающимися множествами ресурсов. Подбор ресурсов поступающим в планировщик заданиям осуществляется в два этапа: сначала они ищутся внутри одного множества, а если там они не обнаружены, то поиск продолжается путем передачи задания соседним планировщикам. При выборе ресурсов используется экономическая модель: ресурсы подбираются не только по ресурсному запросу, но и с учетом целевых функций пользователя и владельца ресурсов. Для определения возможного времени старта задания используются расписания, полученные от планировщиков. Как и в описанной выше системе CCS в этой схеме подразумевается наличие только одного потока заданий, поступающего на каждый планировщик.

## 8. ЗАКЛЮЧЕНИЕ

В работе рассмотрена совокупность механизмов: резервирование, предсказание загрузки ресурсов, приоритетное обслуживание заданий, разделение ресурсов между грид и владельцами – на базе которых разработан метод опережающего планирования. Основным достоинством этого метода является, по-видимому, способность гарантировать точное время запуска

одно- и многопроцессорных заданий в условиях кластеризованных ресурсов, используемых совместно с их владельцами.

Настоящая публикация является развитием работ [53,54]. Программная реализация системы планирования находится в следующем состоянии. Завершен действующий прототип для управления однопроцессорными заданиями, что позволяет приступить к исследованию характеристик системы: производительности, качества планирования и масштабируемости. В стадии разработки находятся расширения базовых алгоритмов механизмами, адаптирующими параметры планирования к меняющейся нагрузке и масштабам грид. Второе направление - обслуживание многопроцессорных заданий - получило развитие в модификации алгоритма Backfill применительно к двухуровневому планированию [55].

## 9. ЛИТЕРАТУРА

- [1]. Ian Foster. What Is The GRID? A Three Point Checklist. GRID Today, July 22, 2002: Vol. 1 No. 6, <http://www.gridtoday.com/02/0722/100136.html>. Русский перевод: <http://www.gridclub.ru/library/publication.2004-11-29.5830756248>
- [2]. Коваленко В.Н., Корягин Д.А. Оценка возможностей программных платформ грид, Труды международной конференции "Распределенные вычисления и Грид-технологии в науке и образовании" (Дубна, 29 июня-2 июля 2004 г.).- Дубна: 11-2004-205, ОИЯИ, 2004, сс. 128-133. <http://www.gridclub.ru/library/publication.2005-03-17.0023427070>
- [3]. Приложения, работающие в грид EGEE. <http://public.eu-egee.org/files/Application%20running%20on%20EGEERusnew.pdf>
- [4]. LCG middleware documentation. <http://grid-deployment.web.cern.ch/grid-deployment/cgi-bin/index.cgi?var=documentation>
- [5]. gLite – Lightweight Middleware for Grid Computing. <http://www.glite.org>
- [6]. Foster I., Kesselman C., J. Nick, Tuecke S. The Physiology of the Grid: An Open Grid Services Architecture for Distributed Systems Integration. <http://www.globus.org/research/papers/ogsa.pdf>. Русский перевод: <http://www.gridclub.ru/library/publication.2004-11-29.8307957187>
- [7]. Modeling Stateful Resources with Web Services. Ian Foster, Jeffrey Frey, Steve Graham, Steve Tuecke, Karl Czajkowski, Don Ferguson, Frank Leymann, Martin Nally, Igor Sedukhin, David Snellin, Tony Storey, William Vambenepe, Sanjiva Weerawana. <http://www-106.ibm.com/developerworks/library/ws-resource/ws-modelingresources.pdf>. Русский перевод: <http://www.gridclub.ru/library/publication.2005-04-25.6831747572/view>
- [8]. Workload Management System User and Reference Guide, <https://edms.cern.ch/file/572489/1/WMS-guide.pdf>
- [9]. Grid Resource Management: State of the Art and Future Trends. Edited by J.Nabrzycki, J.M.Schopf, J.Weglarz. Kluwer Academic Publishers, 2003.
- [10]. E.G. Coffman, P.J. Denning, Operating Systems Theory, 1973, Prentice-Hall, Englewood Cliffs.

- [11]. D.A.Lifka. The ANL/IBM SP Scheduling System. <http://www.tc.cornell.edu/~lifka/>.
- [12]. Jon B.Weissman, P. Srinivasan. Ensemble Scheduling: Resource Co-Allocation on the ComputationalGrid. [http://www-users.cs.umn.edu/~jon/papers/grid2001\\_ens.ps](http://www-users.cs.umn.edu/~jon/papers/grid2001_ens.ps)
- [13]. The Globus Resource Specification Language RSL v1.0. [http://www.globus.org/toolkit/docs/2.4/gram/rsl\\_spec1.html](http://www.globus.org/toolkit/docs/2.4/gram/rsl_spec1.html)
- [14]. Globus Toolkit. <http://www.globus.org>
- [15]. R. Raman, M. Livny, and M. Solomon. Matchmaking: An extensible framework for distributed resource management. *Cluster Computing*, 2(2), 1999.
- [16]. Condor. <http://www.cs.wisc.edu/condor/>
- [17]. Job Description Language HowTo, [http://server11.infn.it/workload-grid/docs/DataGrid-01-TEN-0102-0\\_2-Document.doc](http://server11.infn.it/workload-grid/docs/DataGrid-01-TEN-0102-0_2-Document.doc).
- [18]. J. Frey, T. Tannenbaum, I. Foster, M. Livny, and S. Tuecke. Condor-G: A computation management agent for multiinstitutional Grids. *Cluster Computing*, 5(3):237–246, 2002.
- [19]. K. Czajkowski, S. Fitzgerald, I. Foster, and C. Kesselman. Grid information services for distributed resource sharing. In *Proceedings of the Tenth IEEE International Symposium on High-Performance Distributed Computing (HPDC-10)*, August 2001.
- [20]. Relational Grid Monitoring Architecture (R-GMA). A.Cooke et al, UK e-Science All Hands Conference, Nottingham, September 2003. [http://www.r-gma.org/pub/ah03\\_148.pdf](http://www.r-gma.org/pub/ah03_148.pdf)
- [21]. Portable Batch System. <http://www.openpbs.org>
- [22]. Sun Grid Engine. <http://www.sun.com/software/gridware/sge.html>
- [23]. Platform LSF.  
<http://www.platform.com/Products/Platform.LSF.Family/Platform.LSF/>
- [24]. Calder B., Chien A., Wang Ju, Yang Don. *The Entropia Virtual Machine for Desktop Grids*, 2003.
- [25]. Scheduling in Bag-of-Task Grids: The PAUÁ Case. W.C. Francisco Brasileiro Lauro Costa Daniel Paranhos Elizeu Santos-Neto Nazareno Andrade. <http://walfredo.lsd.ufcg.edu.br/papers/PAUA.pdf>.
- [26]. Коваленко В.Н., Корягин Д.А. Организация ресурсов грид. Препринт ИПМ РАН, № 63, стр. 1-25, Москва, 2004
- [27]. The NorduGrid architecture and tools. P. Eerola et al. *Computing in High Energy and Nuclear Physics*, 24-28 March 2003, La Jolla, California. <http://www.slac.stanford.edu/econf/C0303241/proc/papers/MOAT003.PDF>
- [28]. A.Abraham, R.Buyya, B.Nath. Nature's Heuristics for Scheduling Jobs on Computational Grids, *The 8th IEEE International Conference on Advanced Computing and Communications (ADCOM 2000)*, December 14-16, 2000, Cochin, India. <http://www.buyya.com/papers/nhsjcg.pdf>

- [29]. R. Buyya, J. Giddy, and D. Abramson. A case for economy Grid architecture for service-oriented Grid computing. In Proceedings of the Heterogeneous Computing Workshop, April 2001.
- [30]. R. Buyya, J. Giddy, and D. Abramson. An Economy Grid Architecture for Service-Oriented Grid Computing, 10th IEEE International Heterogeneous Computing Workshop (HCW 2001), San Francisco, USA, April 2001.
- [31]. C. Ernemann, V. Hamscher, and R. Yahyapour. Economic scheduling in Grid computing. In D. Feitelson and L. Rudolph, editors, Job Scheduling Strategies for Parallel Processing (Proceedings of the Eighth International JSSPP Workshop; LNCS #2537), pages 129–152. Springer-Verlag, 2002.
- [32]. How Sun Grid Engine, Enterprise Edition 5.3 Works.  
<http://www.sun.com/software/gridware/sgeee53/wp-sgeee/wp-sgeee.pdf>
- [33]. V. Hamscher, U. Schwiegelshohn, A. Streit, and R. Yahyapour. Evaluation of job-scheduling strategies for grid computing. Lecture Notes in Computer Science, 1971:191–202, 2000.
- [34]. Quinn Snell, Mark Clement, David Jackson, Chad Gregory. The Performance Impact of Advance Reservation Meta-scheduling. Computer Science Department Brigham Young University Provo, Utah 84602-6576, 2000, <http://supercluster.org/research/papers/ipdps2000.pdf>
- [35]. U.Schwiegelshohn and R.Yahyapour. Attributes for Communication Between Grid Scheduling Instances. In “Grid Resource Management. State of the Art and Future Trends”. Edited by J.Nabrzyski, J.M.Schopf, J.Weglarz. Kluwer Academic Publishers, pp. 41-52.
- [36]. Коваленко Е.И., Шорин О.Н. Расширение возможностей кластерных систем управления для информационного обслуживания Грид-диспетчера. Тезисы докладов международной конференции «Распределенные вычисления и Грид-технологии в науке и образовании» Дубна, 29 июня - 2 июля 2004 г., стр. 56.
- [37]. Maui scheduler. <http://www.supercluster.org/maui>
- [38]. J. Krallmann, U. Schwiegelshohn, and R. Yahyapour. On the design and evaluation of job scheduling systems. In D. Feitelson and L. Rudolph, editors, Job Scheduling Strategies for Parallel Processing (Proceedings of the Fifth International JSSPP Workshop; LNCS #1659), pages 17–42. Springer-Verlag, 1999.
- [39]. D.G. Feitelson and A.M. Weil. Utilization and Predictability in Scheduling the IBM SP2 with Backfilling. In Proceedings of IPPS/SPDP 1998, pages 542–546. IEEE Computer Society, 1998.
- [40]. Moab Grid Scheduler.  
<http://www.clusterresources.com/products/mgs/specoverview.shtml>
- [41]. PBS Professional. <http://www.pbspro.com/>
- [42]. В.Н.Коваленко, А.В.Орлов, “Управление заданиями в распределенной среде и протокол резервирования ресурсов”. Препринт ИПИМ РАН, № 1, стр. 1-25, Москва, 2002.

- [43]. A. Roy and V. Sander. Advance Reservation API. Technical Report GFD-E.5, Global Grid Forum (GGF), 2002.
- [44]. Проект Enabling Grids for E-Science (EGEE), <http://www.eu-egee.org>
- [45]. P. Buncic, A. J. Peters, P.Saiz. The AliEn system, status and perspectives. Computing in High Energy and Nuclear Physics, 24-28 March 2003, La Jolla, California.  
<http://www.slac.stanford.edu/econf/C0303241/proc/papers/MOAT004.PDF>
- [46]. Open source metascheduling for Virtual Organizations with the Community Scheduler Framework (CSF). [http://www.cs.virginia.edu/~grimshaw/CS851-2004/Platform/CSF\\_architecture.pdf](http://www.cs.virginia.edu/~grimshaw/CS851-2004/Platform/CSF_architecture.pdf)
- [47]. GridWay. [http://www.gridway.org/int\\_arch.php](http://www.gridway.org/int_arch.php)
- [48]. GridLab Resource Management System.  
<http://www.gridlab.org/WorkPackages/wp-9>
- [49]. J.B. Weissman, "Gallop: The Benefits of Wide-Area Computing for Parallel Processing," Journal of Parallel and Distributed Computing, Vol 54(2), November 1998.
- [50]. L. Yang, I. Foster, and J. M. Schopf. Homeostatic and tendency-based CPU load predictions. In Proceedings of International Parallel and Distributed Processing Symposium (IPDPS), 2003.
- [51]. Matthias Hovestadt, Odej Kao, Axel Keller, and Achim Streit. Scheduling in HPC Resource Management Systems: Queuing vs. Planning. <http://www.uni-paderborn.de/pc2/papers/files/409.pdf>
- [52]. A. Keller and A. Reinefeld. Anatomy of a Resource Management System for HPC Clusters. In Annual Review of Scalable Computing, vol. 3, Singapore University Press, pages 1–31, 2001.
- [53]. С.А.Богданов, В.Н.Коваленко, Е.В.Хухлаев, О.Н.Шорин, "Метадиспетчер: реализация средствами метакомпьютерной системы Globus". Препринт ИПМ РАН, № 30, стр. 1-23, Москва, 2001.
- [54]. V.N.Kovalenko, E.I.Kovalenko, D.A.Koryagin, E.Z.Ljubimskii, A.V.Orlov, E.V.Huhlaev. Resource manager for GRID with global job queue and with planning based on local schedules. Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment. Volume 502, Issues 2-3 , 21 April 2003 , pp 411 – 414.
- [55]. Коваленко В.Н., Семячкин Д.А. Использование алгоритма Backfill в грид, Труды международной конференции "Распределенные вычисления и Грид-технологии в науке и образовании" (Дубна, 29 июня-2 июля 2004 г.).- Дубна: 11-2004-205, ОИЯИ, 2004, сс. 139-144.  
<http://www.gridclub.ru/library/publication.2004-12-27.6965428621>