

УДК 004.7 + 007.3  
ББК 32.973.202-02

## **АРХИТЕКТУРА РАСПРЕДЕЛЁННОЙ ПАМЯТИ ДЛЯ ИЗМЕНЧИВОЙ КОМПЬЮТЕРНОЙ СРЕДЫ**

**Знаменский С. В.<sup>1</sup>**

*(ФГБУН Институт программных систем  
им. А.К. Айламазяна РАН, Переславль-Залесский)*

*Представлена методология использования ретроспективной общей памяти в качестве основы для безгранично перестраиваемого высокодоступного компьютерного сервиса.*

Ключевые слова: управление большими системами, разделяемая память, распределённые системы, системная архитектура.

### **Введение**

Целью работы является проработка архитектурных принципов высокоэффективных системы информационной поддержки сложного долговременного целенаправленного взаимодействия большого числа людей и вычислительных средств в будущем десятилетии.

Задачами такого рода систем будут разработка гладко обновляемой операционной системы и конструктивное разрешение содержательных социально-экономических и экологических проблем.

Имеется в виду, что мощность распределённой вычислительной системы в процессе бесперебойного обслуживания увеличивается путем замены устаревших кластеров на новые, которые могут принципиально отличаться по архитектуре и управлению. Речь о системе, которая не выключается и постоянно предоставляет внешний доступ к данным. Без остановок в ней же разраба-

---

<sup>1</sup> *Сергей Витальевич Знаменский, доктор физико-математических наук, доцент, (svz@latex.pereslavl.ru).*

тывается уникальное системное и прикладное программное обеспечение, отлавливаются и исправляются ошибки, сопоставляются результаты, полученные разными путями, оптимизируется выполнение задач, устраняются любые изолированные поломки.

Такая организация высокопроизводительного бесперебойного сервиса в долговременной перспективе является проблематичной: обновление оборудования не сводится к простой замене, а требует использования новых парадигм, качественно изменяются запросы и интерфейсы пользователей и средства разработки программного обеспечения.

Прогноз погоды даёт не единственный возможный пример будущей прикладной длительно функционирующей системы, в которой количественно и качественно растущие данные должны параллельно в тесном сотрудничестве обрабатываться в разных взаимосвязанных целях.

Повышению качества прогноза послужат реорганизация ввода информации от новых источников, наращивание вычислительных мощностей, совершенствование алгоритмов обработки и интерфейсов пользователя.

Единая информационно-вычислительная система должна обеспечивать не только вычисления по фиксированному алгоритму на основе заранее определённых структур входных данных, но и включение в обработку новых видов входной информации и новых алгоритмов обработки и совершенствование старых алгоритмов. Основой всего должна стать развёрнутая организация параллельных работ с общими данными на расширяющейся вычислительной основе. Единой организацией доступа к данным должны быть поддержаны работы по сопоставительному анализу разных подходов и алгоритмов обработки, выявлению возможностей ситуативного выбора наиболее адекватных алгоритмов обработки и архитектур локального оборудования. На этой же основе должны попутно развиваться и исследоваться алгоритмы интервальных оценок будущих погодных явлений.

На этом примере отчётливо видна потребность повторного использования полной информации о входных данных, результа-

тах их обработки, использованных алгоритмах и конфигурации системы при вычислениях.

Видна на нём и потребность в длительно функционирующей системе с гибко изменяемой конфигурацией оборудования и программного обеспечения и организации работ. Длительное существование информационной системы связано с рисками принципиальной несовместимости оборудования.

Исследования, связанные с организацией высокопроизводительного бесперебойного сервиса в долговременной перспективе [17, 18, 20], не подсказывают ясного пути к решению поставленной задачи.

Особенности оборудования, программных интерфейсов или общей структуры управления не могут претендовать на архитектурную основу такой системы. Этой архитектурной основой становится ретроспективность: доступность в неискажённом виде всей ранее доступной информации, включая программные и пользовательские интерфейсы и сопроводительную документацию по указанию времени прежнего доступа.

## **1. Идея ретроспективности**

Ретроспективная парадигма означает непризнание абстракции текущего состояния. Неправомерно требовать от распределённой системы, чтобы она мгновенно корректно изменила состояние: она может устанавливать новое логически непротиворечивое состояние в течение некоторого времени. Обычно это время мало, но даже в самых надёжных системах нередко ситуации, когда оно увеличивается по причинам неизбежных сбоев или мелких поломок.

Мы приходим к требованию *ретроспективной согласованности данных информационной системы*:

1) в любой момент времени  $t$  в системе доступен прошедший момент истины  $\mu(t) < t$ , которым завершилась доступная история согласованных изменений данных системы;

2) для любого доступного в истории предыдущего момента  $\tau < \mu(t)$  система возвращает корректно согласованные реплики.

Разность  $t - \mu(t)$  назовём *длительностью согласования*. При удачном разделении длительность согласования для части системы оказывается меньше, чем для целой системы.

*Ретроспективная* информационная система может не иметь единой концептуальной модели. Более того, её концептуальная модель может изменяться не мгновенно и не одновременно во всех частях. Этим она принципиально отличается от *последовательных (linearizable)* информационных систем, основанных на моделировании последовательности изменений глобального согласованного состояния системы.

Общепринятое с прошлого века [16] моделирование распределённой системы конечным автоматом (*linearizability*) корректно лишь при отсутствии сбоев, поломок и вызванных ими перезапусков процессов. Никакая реальная система без нарушений последовательности не обходится. Удивительно, что эта привычная мотивированная простотой модель, несмотря на очевидную неадекватность, остаётся общепризнанно эталонной для критических приложений, в том числе для банковских систем. Даже порой неизбежная необратимость составляющих сложной транзакции почему-то принимается как несовершенство мира, а не как принципиальный недостаток логической модели, базирующейся на идее отката транзакций.

Практически все существующие информационные системы *последовательные*. Это косвенно подтверждается растущей популярностью теоремы CAP Брюера [19]. Возможно, это объясняется двумя причинами: (1) неадекватная модель проста; (2) именно она канонизирована в стандартах передачи и обработки данных.

Заметим, что простота здесь тоже сомнительная: хорошо известно, что задача верификации разделяемой памяти в этой модели является *NP*-полной [12, 14, 21].

Ретроспективная парадигма разрешает системе иметь согласованное состояние не в каждый момент времени. Это означает, что пользователь на свой запрос может получить ответ о старом (уже неактуальном) согласованном состоянии с указанием момента времени, на который информация действительна:

Неверно полагать, что отказ пользователю в возможности доступа к текущей актуальной информации ущемляет его. В аналогичной ситуации стандартная система действительно покажет информацию о текущем состоянии, но ответ будет получен через значительное время, и к какому именно промежуточному моменту в этом промежутке будет относиться информация, так и останется неизвестным. Для ретроспективной системы легко сделать пользовательский интерфейс, эффективно эмулирующий подробное поведение: интерфейс должен не спешить с посылкой запроса, повторять запрос, пока он не будет датирован более новым моментом и не отразит время актуальности. Сомнительно, что пользователю всегда нужно дополнительное ожидание и неинтересен момент актуальности при задержке ответа.

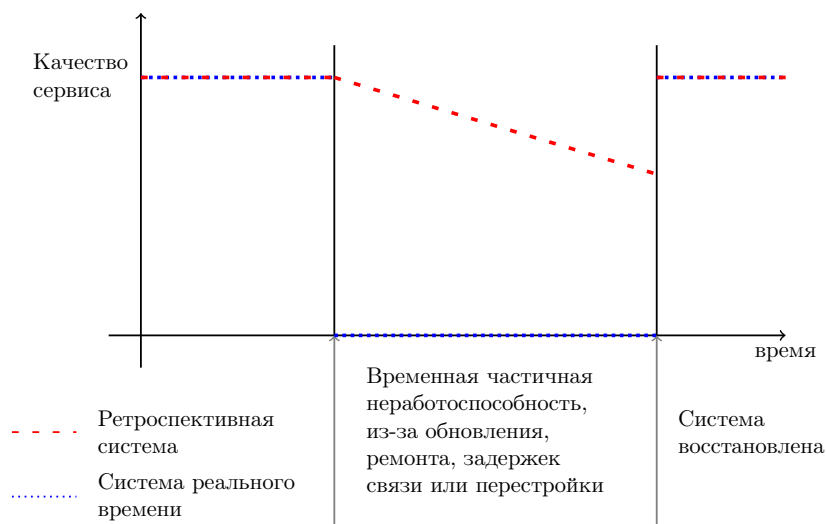


Рис. 1. Сравнение ретроспективной системы с системой реального времени

На примере банковских систем ретроспективная парадигма означает, что сделки совершаются не мгновенно, поэтому значение денег на счёте в момент незавершённой транзакции есть абстракция, не имеющая практического смысла. Поэтому не состо-

яния счетов, а записи в истории изменений являются в ретроспективной парадигме базовыми объектами для таких систем. Точное значение на счету реально существует лишь для достаточно отдаленного прошлого, все транзакции из которого можно считать так или иначе завершёнными.

Записи об изменении счетов только создаются и никогда не изменяются. Можно считать, что они никогда не исчезают. Поэтому на прикладном уровне в такой основанной на истории изменений банковской системе нет никакой конкуренции доступа и нужды в блокировках и очередях (записи ведь все разные и у каждой свой автор), да и по всей видимости в откатах транзакций.

То, что система обрабатывает изменения не по запросу, а по факту появления, создаёт эффект, заметный для пользователя, получающего результат сложной обработки: он, как и в современных интернет-сервисах, получает практически мгновенно результат обработки, начатой до того, как он отправил запрос.

Учитывая техническую сложность согласования времени с высокой точностью, ретроспективный подход требует, чтобы погрешность согласования времени не была существенно больше, чем время передачи данных. Однако трудно вообразить ситуацию, в которой это ограничение было бы обременительным.

В отличие от популярной парадигмы «согласованности в конечном счёте» (Eventual consistency [22]) ретроспективная парадигма учитывает, что обрабатывать меняющиеся данные корректно лишь по их состоянию на *фиксированный* момент времени.

Покажем, как ретроспективность позволяет резко повысить устойчивость распределённых систем к частичным поломкам и задержкам связи.

## **2. Устойчивость ретроспективных систем к задержкам и разделением**

Известная [8] своими небезупречными популяризациями *теорема CAP* [10] утверждает невозможность сочетания трёх свойств:

*Consistency (Согласованность):*

ответы на запросы к системе логически непротиворечивы.

*Availability (Доступность):*

ответ на любой запрос может быть получен незамедлительно.

*Partition-tolerance (Устойчивость к разделению):*

при восстановлении потерянной на время связи функциональность и внутренняя целостность восстанавливаются.

Формальное доказательство этой теоремы в [15] основано на прозрачной идее:

Как только в частях системы появились различия, то их доступность выявила противоречие

Практически после разрыва связи с филиалом компании работа *должна* быть продолжена с адекватными ситуациями ограничениями и после восстановления связи система *должна* полностью восстановиться без потерь информации, то логика доказательства очевидно противоречит потребностям практики.

Неприемлемость строго доказанного результата в науке встречается. Достаточно вспомнить классические безупречные доказательства невозможности рационального корня из числа 2, корней из отрицательных чисел и недифференцируемости функции Хевисайда. Их неприемлемость привела в своё время к развитию и широкому применению новых теорий, ставших основой современной математики — теории вещественного числа, теории комплексных чисел и теории обобщённых функций.

Рассмотрим возможность и целесообразность построения системы, сочетающей ретроспективность с *адекватностью* — безупречностью реакции на запросы и соразмерностью последствий временных поломок и сбоев.

## 2.1. ПРОБНАЯ ЗАДАЧА

Банк имеет тысячу клиентов в двух городах  $A$  и  $B$ . Расчётное максимальное время задержки передачи данных между городами  $T = 0,25$  секунды.

Города расположены так, что устранить эту задержку технически невозможно<sup>2</sup>

Для поддержки переводов денег со счёта на счёт банк нуждается в качественной информационной системе с по возможности быстрым временем отклика  $\tau$ , рассчитанной на обработку не более десяти запросов в секунду с одновременной поддержкой сессий:

- 1) результаты любых запросов гарантированно должны быть корректны и согласованы;
- 2) при любом обращении к сервису задержка обработки не должна превосходить  $\tau$  для любого из следующих запросов:
  - а) просмотр состояния счёта и истории его изменений;
  - б) перевод имеющейся суммы на другой счёт;
  - в) получение отчётов.

Для предельного упрощения задачи полагаем, что банк с наличностью не работает, кредитованием не занимается, вне городов никого не обслуживает, пользователи заведены с неизменными личными данными и для простоты не переезжают, а потерями и задержками внутри городов можно пренебречь.

## 2.2. ПОСЛЕДОВАТЕЛЬНЫЙ ПОДХОД

Информация о текущем состоянии каждого счёта должна где-то храниться. Запрос о переводе на этот счёт из другого города должен прийти и после этого ответ должен вернуться обратно, чтобы пользователь смог знать, что запрос принят. Поэтому реакция системы не может быть гарантирована быстрее, чем в течение  $\tau = 2T = 0,5$  секунды.

## 2.3. РЕТРОСПЕКТИВНЫЙ ПОДХОД

В мегаполисах размещаются идентичные серверы. Каждый из них ведёт счета пользователей из своего мегаполиса.

---

<sup>2</sup> При передаче со скоростью света через геостационарный спутник сигнал задерживается на  $1/4$  секунды. При недоступности геостационарного решения неизбежны дополнительные задержки.



Сервер отказывается выполнить перевод, если денег на счету отправителя недостаточно. Если на счету достаточно денег, то сервер немедленно изменяет состояние локальных счетов и надёжно запоминает своё обязательство по возможности быстро выполнить перевод. Это должно делаться за ощутимо меньшее время, чем ожидаемое время поступления следующего запроса. Если повторный запрос от пользователя на перевод денег следует через доли миллисекунды, то это скорее повторное срабатывание и уж никак не осмысленные действия пользователя, поэтому гораздо безопаснее оставить лишь один из серии таких практически одновременных запросов.

После обработки серии запросов (например каждые полсекунды) серверы передают друг другу полные списки накопившихся запросов и записей изменений в счетах и извещения о получении списков. Извещение о получении посылается только в том случае, когда пакет получен полностью и содержит момент времени, вплоть до которого все данные получены. Пакетная передача помогает гарантировать единый порядок и отсутствие пропусков в общей части информации серверов и корректно снизу оценить момент истины системы  $\mu(t)$ , вся история изменений вплоть до которого получена на обоих серверах.

Пользователь получает ответ на свой запрос с учётом реакции только ближайшего сервера. Реакция удалённого (например, уведомление о получении перевода в другой город) может быть учтена в отдельном последующем через секунду запросе (который может делаться автоматически либо клиентским приложением либо на основе технологии comet [11], позволяющей локальному серверу самостоятельно передать такое уведомление в момент получения).

#### **2.4. СРАВНЕНИЕ ПО СКОРОСТИ И ИНТЕРФЕЙСУ**

Условия задачи не влекут никаких технических препятствий к получению быстрого (в течение десяти миллисекунд) ответа на запрос на основе имеющейся на ближайшем сервере информации. Повышение реактивности системы в сто раз здесь достигнуто по сути включением в реакцию системы на запрос пользо-

вателя только немедленно доступной информации.

Это позволяет обеспечить пользователя удобным интерфейсом, дающим возможность выбора, наблюдать за процессом обработки либо спокойно выключить клиентское приложение либо заняться операциями с другими счетами.

Взаимодействие с внешней системой может породить сложную совместную транзакцию (например, перевод с одного счёта на другой в другой город, оттуда на третий и снова назад) и потери времени станут неизбежны. Для клиента, ожидающего поступлений на счёт, может быть создано клиентское приложение, ожидающее изменения состояния счёта для повторения запроса, отвергнутого из-за недостатка денег на счету.

Разумеется, при переводе денег упомянутая секунда задержки не составляет проблемы для клиента. Сюжет задачи выбран для ясности рассмотрения, а ничем не компенсированные задержки в других ситуациях могут оказаться критическими.

## 2.5. СРАВНЕНИЕ ПО ЛОГИЧНОСТИ И НАДЁЖНОСТИ

В конкретной задаче общая сумма денег на счетах в реальности не константа. Ясно и привычно, что отправленные с одного счёта деньги поступают на другой не мгновенно. Стандартная реализация на основе транзакций прячет время транзакции в неповоротливость пользовательского интерфейса, искусственно создавая иллюзию мгновенности сделки и ущемляя пользователя.

Ретроспективное решение делает находящиеся в пути деньги доступными для точного последующего анализа.

История изменений счёта, полученная пользователем, полна до момента  $\mu(t)$ , а после  $\mu(t)$  в ней могут временно отсутствовать деньги, «находящиеся в пути».

Запросы сводной информации на момент, предшествующий  $\mu(t)$ , не зависят от сервера. Для запроса сводных данных на последний момент у пользователя есть две возможности:

1) посмотреть сводную по всем доступным (неполным) данным с предупреждением о неполноте;

2) посмотреть *полную* сводную по всем данным на последний момент  $\mu(t)$ , для которого они доступны.

Если запросы следуют не чаще, чем может обработать традиционная система, то выдаётся в точности тот ответ, каким он был бы в традиционной системе на запрос, посланный из оптимально выбранного момента прошлого (неважно, был ли такой запрос в действительности).

Поломка связи между городами при стандартном подходе полностью блокирует операции в одном из них. При ретроспективном подходе блокируются лишь операции, непосредственно нуждающиеся в этой связи, остальное продолжает работать, и сервис немедленно полностью восстанавливается при восстановлении связи.

Кроме упомянутых, ретроспективная организация имеет и другие полезные особенности:

Поломка или неожиданная потеря одного сервера базы данных (попадание метеорита) при стандартном подходе должно быть заранее предупреждено специальными мерами, например, чёткой организацией резервного копирования или репликации базы. Иначе потери могут оказаться фатальными.

При ретроспективном подходе без каких-либо дополнительных мер исчезновение сервера приводит к потере лишь части обновлений за последние доли секунды, что оставляет шансы на спасение.

## **2.6. ИТОГИ СРАВНЕНИЯ**

**Теорема 1.** *Существует возможность реализации прикладной информационной системы с ретроспективным доступом, вполне адекватной функциональностью при разделении, быстрым автоматическим восстановлением после восстановления связи и с временем отклика в сто раз ниже, чем у любой реализации на стандартной основе.*

Эта теорема разумеется противоречит не самой теореме CAP, а её популярным интерпретациям.

Рассмотренный пример вырос из неоднократно обсуждавшихся на конференции [1, 3, 6] ретроспективных систем. Он ставит ряд вопросов:

1. Ведёт ли он к общей технологии разработки ретроспективных информационных систем?

2. Не чрезмерны ли требования ретроспективной парадигмы к объёмам памяти и вычислительных ресурсов?

3. Сможет ли ретроспективная технология кардинально упростить исправление концептуальных ошибок и создание эволюционирующих систем?

4. Сможет ли ретроспективная технология резко уменьшить потери и неудобства пользователей, связанные с обновлением систем?

5. Сможет ли положительный эффект ретроспективного подхода проявиться при организации многопроцессорных распределённых систем и высокопроизводительных вычислений с большими гранулами параллелизма либо с задержками переконфигурирования *FPGA* [9]?

### **3. Ретроспективная парадигма программной инженерии**

Проектирование системной архитектуры предполагает разделение системы на наиболее крупные составные части и принятие конструктивных решений, которые после их принятия с трудом поддаются изменению.

В рассматриваемой общей постановке неограниченно изменяются концептуальная модель, используемое оборудование и приложения (сервисы). Поэтому они не могут быть общей основой системной архитектуры для ретроспективной системы. Такой основой может быть только стабильная подсистема идентификации структур процессов и данных, обеспечивающая безусловную доступность неподдельной истории при произвольных изменениях.

Прежде чем её описывать, следует разобраться с понятиями модульности и информационных объектов, которые в системе с историей и изменчивостью имеют свои особенности.

### 3.1. СОХРАНЕНИЕ ПОЛНОЙ ИСТОРИИ

История изменений необходима

- для эффективного согласования распределённых данных;
- для быстрого восстановления при поломках;
- для ускорения отладки нового кода;
- для упрощения расследования вторжений и других происшествий;
- для повышения ответственности персонала.

Последние две цели требуют гарантий неподдельности истории. Эти гарантии нуждаются в комплексной поддержке истории изменений *на более низком уровне*, чем механизмы разрешение конфликтов доступа и механизмы обеспечения прикладной функциональности.

Экономное сохранение истории достигается разделением данных на *небольшие независимо изменяемые части*. Например, для веб-сайта организации, управляемого *CMS*, должны сохраняться история изменений шаблонов фрагментов страниц, версии фрагментов заполнений, версии исполняемых процедур, обеспечивающих формирование страницы на этой основе (об истории активности пользователей разговор особый).

**Пример.** *За два года существования сайта каждый из этих элементов изменяется в среднем пять раз (редкие 100, многие 1), а средний размер элемента превышает длину штампа времени и текущий размер информации сайта 100 Мб.*

Общий объём всех версий исходных файлов, помеченных моментом сохранения, составляет

$$(100 + 100) \times 5 = 1000 \text{ Мб.}$$

Объём выданных пользователям веб-страничек, которые получают разными сочетаниями исходных файлов и шаблонов, может на порядки превосходить эту цифру.

Для быстрого доступа к любой страничке из прошлого нужно лишь добыть актуальное на этот момент сочетание файлов и

шаблонов, что несложно организовать с логарифмически растущими с объёмом издержками [3]. Много резервных копий системы не потребуется, поскольку вся история остаётся нетронутой и сможет пропасть только вместе с используемыми данными.

Пример показывает возможность хранения высоко доступной полной истории изменений при вполне умеренном росте ресурсозатрат при двух условиях:

- 1) выделение независимо изменяющихся «первичных» фрагментов информации в отдельные информационные объекты;
- 2) несохранение объектов, просто однозначно реконструируемых из первичных.

Если ненужные версии данных занимают слишком много места, часть давно неактуальной информации может автоматически удаляться [4], создавая *белые пятна* истории на оси времени — промежутки, состоящие из моментов времени, информация на которые недоступна. Этот процесс мог бы проводиться автоматически субоптимально [5], оставляя относительный размер каждого белого пятна близким к минимально возможному при заданном ограничении общего объёма хранилища.

### 3.2. МОДУЛЬНОСТЬ

Модульность освобождает разработчиков прикладных программ от сложностей организации исполнения. Она традиционно понимается как разделение на изолированные модули со своим кодом и данными, обменивающиеся сообщениями. Рассмотрим, как такое понимание может при структурных изменениях приводить к неконтролируемому росту сложности организации обмена сообщений.

Сложность эта порождается двумя принципиальными дефектами неуправляемого обмена сообщениями:

- 1) новому исполнителю недоступна информация, разосланная участникам совместно выполняющейся работы;
- 2) сообщения теряются по разным причинам;
- 3) сообщение читается с неизвестным запозданием, за это время оно может потерять адекватность, а отменить непрочитанное сообщение автор не может.

Отсюда опасность использования изоляции модулей с организацией обмена сообщениями в коллаборативном приложении: после успешного тестирования малейшее изменение структуры взаимодействия или возникновение задержек может неожиданно выявить необходимость корректирующих сообщений. В итоге безупречная работа обеспечивается лишь до замены модулей и откатов неправильных действий, а запуск новых исполнителей на замененном оборудовании и (или) с новыми алгоритмами исполнения требует серьёзного рефакторинга системы.

Альтернативой обмена сообщениями является публикация данных с подпиской на изменения, и, в частности, стандартное решение *DDS (Data Distribution Service)* [13].

Механизм обработки запросов на подписку и на публикацию авторы не регламентируют. Во избежание коллизий запросы, вероятно, должны обрабатываться строго, т.е. проходить через общий вычислительный узел, что принципиально ограничивает масштабируемость. В любом случае это не снижает ценности *DDS* как перспективного новейшего технического решения для общей памяти распределённой системы, которое, по-видимому, может эффективно масштабироваться при устойчивых структурах данных с достаточно редко меняющимися настройками.

Подписка производится с настройками, устанавливающими количество хранимых версий, время жизни данных, количество реплицируемых копий, планомерно допустимые задержки, учитываемыми период изменения данных и приоритетность передачи. *DDS* ничем не ограничивает свободы программиста по выбору этих настроек.

Получая полный доступ к гибкому распределению вычислительных ресурсов, прикладной программист, к сожалению, практически лишён разумных ориентиров. Это создаёт сложнейшую проблему избыточной гибкости настроек, которые для долгоживущей системы должны быть ещё и динамически управляемыми.

Решение этой проблемы требует системной организации для гибкого регулирования минимальным количеством прозрачных настроек.

Построив дерево всех объектов, настройки будем считать вектор-функцией от узла. Если предположить, что настройки чаще должны совпадать для смежных узлов, то задание настроек сводится к разрезанию дерева и заданию настроек для каждой связной части. Поскольку каждая связная часть имеет ближайший к корню элемент, то именно в нём можно задавать или менять настройки. Тогда настройки, сделанные в любом узле, наследуются во всей ветке, за исключением подветок со своими настройками.

Это означает, что проблема настроек разделяемой памяти сводится к разумному представлению всех данных системы в виде дерева объектов.

Примерами таких веток в рассмотренном примере являются данные о изменениях в счёте отдельного пользователя и данные по городу. В общем случае такое разделение означает иерархию управления изменениями в системе.

### 3.3. КОНТЕКСТНАЯ АВТОНОМНАЯ АРХИТЕКТУРА

Архитектура контекстно-автономной системы [1] формируется на основе иерархии управления изменениями. Используется идея процессного подхода, рекомендованного стандартами ИСО серии 9000. Однако вместо термина *процесс* мы будем использовать более общий термин *активность*, не предполагающий обязательности входных данных и не ассоциирующийся с процессами операционной системы. Так, множество всех ИСО9000-процессов является подмножеством всех активностей системы.

#### **Другие примеры активностей:**

- нажатие пользователем кнопки, переход по ссылке или ввод текста;
- обнаружение события, нуждающегося в обработке;
- обработка данных, полученных от разных пользователей, при совместном редактировании.

Мониторинг сложных активностей осуществляют активности управления качеством, в частности, корректирующие взаимные приоритеты дочерних активностей.



Системная активность управляет качеством системы в целом через такие дочерние активности, заведомо неполный список которых включает:

- мониторинг наличия и состояния физических устройств и каналов связи;
- мониторинг информационного обмена между активностями и между физическими узлами и перемещение активностей.
- мониторинг взаимных приоритетов активностей с общим родителем;
- мониторинг качества обновления информации;
- документирование кода и данных системы;
- уточнение приоритетов информационного обслуживания активностей;

Большинство задач системной активности в начальной фазе существования системы может осуществляться человеком, а затем переводиться на автоматическое управление. Автоматизация управления требует постановок задач поиска эффективных алгоритмов. Эти алгоритмы должны будут разумным образом перестраивать размещения активностей по физическим устройствам и оптимизировать коммуникации, используя в качестве входных данных изменения физической конфигурации системы, изменения приоритетов и статистику мониторинга активностей.

### **3.4. ЛОГИЧЕСКАЯ ИДЕНТИФИКАЦИЯ ДАННЫХ**

Элемент данных (например, атрибут объекта) может сохраняться и реплицироваться в разных узлах системы. Изменение данного создаёт его новую версию. Прежде чем идентифицировать версии в различных местах системы, мы должны идентифицировать то, версии чего в ней сохраняются. Идентификатор самого элемента (например, поле конкретной формы, заполненное конкретным пользователем, содержащее код исполняемого в системе скрипта), не учитывающий версии и места хранения, будем называть *логическим идентификатором элемента данных*. Прозрачность идентификации данных крайне важна для доработок системы.

Для каждой активности  $A$  определим её *контекст*  $\mathcal{D}(A) =$

$[d_0(A), d_1(A)] \subset U$  как сегмент в линейно упорядоченном универсуме  $U$  всех возможных логических идентификаторов данных [4].

Универсум  $U$  будем считать *неисчерпаемым* в следующем смысле:

$$\forall x, y \in U \exists z \in U : x < z < y.$$

Неисчерпаемость обеспечивает возможность размещения в любом контексте множества данных любой конечной мощности. Например, универсум строк символов неограниченной длины с отношением лексикографической упорядоченности неисчерпаем.

Для вложенных контекстов активность с более широким контекстом будем называть родительской, а с более узким — дочерней.

Принцип разделения ответственности и принцип единоначалия приводят к следующим требованиям:

1) *контексты любых двух активностей либо вложены, либо не пересекаются;*

2) *активность с более широким контекстом не модифицирует контексты дочерних активностей;*

3) *перемещение активности в новую область резервирует старый контекст для возможного возврата и изоморфно переносит все дочерние активности.*

Добавляя совокупную активность, поддержанную системой, получаем в каждый момент времени дерево активностей. Оно изменяется: активности появляются, замирают, глубина вложенности активности может изменяться. Это обеспечивает произвольные структурные перестройки и гарантирует доступ к неискажённому прошлому.

### 3.5. ОРГАНИЗАЦИЯ ИСПОЛНЕНИЯ

Иерархия активностей упрощает проблему регулирования разделения ресурсов. Приоритеты допустимости задержек, уровня защищённости и ценности истории устанавливаются (и могут в любой момент быть пересмотрены) для общих активностей верхнего уровня:

1) состояние изменяющегося объекта однозначно определено лишь в достаточно давнем прошлом;

2) безупречный ответ может быть дан лишь на запросы об устоявшемся (линейно упорядоченном) прошлом, предшествующим моменту истины множества данных, нужных для обработки запроса;

3) *момент истины* (ранее которого прошлое устоялось) быстро определяется по времени и идентификатору;

4) обновление информации в автономном контексте данных происходит с ограниченной частотой при появлении изменений, необработанных на момент, когда все входные данные устоялись;

5) изменения помечаются комбинированным временем обработки, интерпретируемым и как частично упорядоченное логическое (линейно упорядоченное до момента истины), и как приближение к физическому времени.

Ограниченность частоты изменений обеспечивает локализацию последствий ошибок и открывает возможность оперативного автоматического выявления проявившихся в работающей системе противоречий. Локализация последствий ошибок вместе с автоматической диагностикой особенно важна для ретроспективной системы, поскольку перестройку «на ходу» невозможно осуществить без временных рассогласований и исправляемых ошибок.

Поясним сказанное.

**Пример.** Один процесс обеспечивает  $a = -b$ , а другой  $b = a + 1$ .

Подобное нарушение логики, проскользнувшее в обычную работающую систему, либо останется незамеченным, либо заставит систему бесконечно пересчитывать  $a$  и  $b$ .

Фиксация частоты изменений ограничивает любые негативные эффекты, в том числе и этот.

При появлении нескольких кандидатов для значения в близком времени возникает конфликт, требующий разрешения. Иногда для разрешения могут использоваться либо выбор значения, поступившего последним, либо выбор наибольшего (рекордного) из поступивших значений.

Остальные значения игнорируются или удаляются при записи. Если конфликты сложнее, то могут вводиться дополнительные активности для разных источников значений с целью разрешения конфликтов.

Если, например, от пользователя поступает более десяти форм в секунду, то это возможная попытка взлома. Классическое эффективное средство борьбы — это игнорирование промежуточных запросов от одного пользователя в малом промежутке времени. При этом запросы от всех разных пользователей должны быть сохранены и корректно обработаны.

Уменьшение частоты изменений контекстов ресурсозатратных активностей, таких как профиля пользователя или объёмная статистика, экономит ресурсы.

### 3.6. РЕТРОСПЕКТИВНАЯ СУБД КАК ОСНОВА ФИЗИЧЕСКОЙ РЕАЛИЗАЦИИ

Доступ к физическим ресурсам (процессоры, память) обеспечивается ретроспективной СУБД [6]. На каждом вычислительном узле она хранит лишь часть общей истории изменений объектов, а именно:

- 1) данные активностей, осуществляющихся на узле;
- 2) наличие и давность необработанных изменений во входных данных;
- 3) данные активностей, используемых в качестве входных данных, в частности:
  - исполняемый код осуществляющихся в узле активностей;
  - относительные приоритеты исполнения, шаги дискретизации и моменты истины;
  - списки активностей-подписчиков<sup>3</sup> с их шагами дискретизации;
  - список узлов, на которых осуществляются активности-подписчики;

---

<sup>3</sup> *Активности-подписчики используют результаты осуществляющихся в узле активностей в качестве входных данных.*

4) состояния процессов репликации данных на узлы активностей-подписчиков, давность непереданных изменений.

Ретроспективная СУБД реализует на узле следующую функциональность:

1) быстрое сохранение изменённой версии объекта с штампом времени изменения;

2) очень быстрое чтение версии объекта на заданный момент времени;

3) передача изменений на узлы с активностями-подписчиками;

4) сохранение изменений входных данных, полученных по подписке;

5) поддержка ссылок на согласованные версии данных активности с изменениями в отдельных объектах.

Подробности конструкции ретроспективной СУБД описаны в [4] и [6].

#### **4. Выводы:**

Рассмотрен базирующийся на сохранении полной истории способ организации информационной системы.

Показано, как он обеспечивает резкое повышение доступности и защищённости системы.

В общих чертах описана архитектура подобной системы, основанная на общей ретроспективной памяти, поддерживаемой на узлах ретроспективными СУБД. Она позволяет гибко расширять и заменять оборудование, обновлять структуры, алгоритмы и пользовательские интерфейсы без приостановок обслуживания.

Описаны решения, позволяющие полностью контролировать накладные расходы на сохранение полной истории.

#### **Литература**

1. АБРАМОВ С.М., ЗНАМЕНСКИЙ С.В., ЖИВЧИКОВА Н.С., КОТОМИН А.В., ТИТОВА Е.В. *Информационная система для разработки технологий организа-*

- ции сложной совместной деятельности // RCDL2009. – С. 186–192.
2. АНДРЕЕВ С.С., ДАВЫДОВ А.А., ДБАР С.А., ЛАЦИС А.О., ПЛОТКИНА Е.А. *О моделях и технологиях программирования суперкомпьютеров с нетрадиционной архитектурой* // Научный сервис в сети Интернет: суперкомпьютерные центры и задачи. – М.: Изд-во МГУ, 2010. – С. 186–187.
  3. ЗНАМЕНСКИЙ С.В. *Гибкая основа информационной системы для обучения* // RCDL2010. – С. 451–460.
  4. ЗНАМЕНСКИЙ С.В. *Глобальная идентификация данных в долговременной перспективе* // Прогр. сист. теор. прил. – 2012. – Т. 3, №2(11). – С. 77–88.
  5. ЗНАМЕНСКИЙ С.В. *Показатели эффективности распараллеливания резервного копирования* // Прогр. сист. теор. прил. – 2012. – Т. 3, №2(11). – С. 51–60.
  6. ЗНАМЕНСКИЙ С.В. *Ретроспективная основа совместной реорганизации сложных информационных ресурсов* // RCDL2011. – С. 93–101.
  7. КЛИМОВ Ю.А., ОРЛОВ А.Ю., ШВОРИН А.Б. *Программный инструментарий для трафаретных вычислений на гибридных суперкомпьютерах* // Прогр. сист. теор. прил. – 2012. – Т. 3, №2(11). – С. 23–49.
  8. КУЗНЕЦОВ С.Д. *Транзакционные параллельные СУБД: новая волна* // Труды Института системного программирования РАН. – 2011. – Т. 20. [Электронный ресурс]. – Режим доступа: <http://cyberleninka.ru/article/n/tranzaktsionnye-parallelnye-subd-novaya-volna>.
  9. BAUER L., BRAUN C., IMHOF M.E., KOCHTE M.A., ZHANG H., WUNDERLICH H.-J., HENKEL J. *OTERA: Online Test Strategies for Reliable Reconfigurable Architectures* // AHS12. – 2012. – P. 38–45.
  10. BREWER E. *Towards Robust Distributed Systems* // Proc. XIX Ann. ACM Symposium on Principles of Distributed Computing. – 2000. – P. 7.

11. BOZDAG E., MESBAH A., VAN DEURSEN A. *A comparison of push and pull techniques for AJAX* // Web Site Evolution, WSE 2007. – 2007. – P. 15–22.
12. CANTIN J.F., LIPASTI M.H., SMITH J.E. *The complexity of verifying memory coherence and consistency* // IEEE Transactions on Parallel and Distributed Systems. – 2005. – Vol. 16, №7. – P. 663–671.
13. CORSARO A., SCHMIDT D.C. *The Data Distribution Service – The Communication Middleware Fabric for Scalable and Extensible* // System of Systems. – 2012.
14. GIBBONS P., KORACH E. *Testing shared memories* // SIAM Journal on Computing. – 1997. – №26. – P. 1208–1244.
15. GILBERT S., LYNCH N. *Perspectives on the CAP Theorem* // Computer-IEEE Computer Magazine. – 2012. – №2. – P. 30–36.
16. HERLIHY M.P., WING J.M. *Linearizability: A correctness condition for concurrent objects* // ACM Transactions on Programming Languages and Systems. – 1990. – Vol. 12, №3. – P. 463–492.
17. JI HONG YAN, CHUN HUA FENG *Sustainability-Oriented Product Modular Design Using Design Structure Matrix (DSM)* // Method. Appl. Mechan. and Mater. – 2011. – P. 1468–1471.
18. METZGER A., POHL K., PAPAZOGLU M., NITTO E. DI, MARCONI A., KARASTOYANOVA D. *Research challenges on adaptive software and services in the future internet: Towards an scube research roadmap* // ICSE 2012.
19. SHIM S.S.Y. *The CAP Theorem's Growing Impact* // Computer. – 2012. – Vol. 45, №2. – P. 21–22
20. *Software Engineering for Self-Adaptive Systems: A Research Roadmap* / Eds. B.H.C. Cheng et al. Eds. // Self-Adaptive Systems. – 2009. – LNCS 5525. – P. 1–26.
21. TAYLOR R.N. *Complexity of analyzing the synchronization structure of concurrent programs* // Acta Informatica. – 1983. – Vol. 19, №1. – P. 57–84.

22. VOGELS W. *Eventually Consistent* // Communications of the ACM - Rural engineering development. – 2009. – Vol. 52, №1. – P. 40–44.

## **DISTRIBUTED MEMORY ARCHITECTURE FOR CHANGING COMPUTING ENVIRONMENT**

**Sergej Znamenskij**, Program Systems Institute of RAS, head of laboratory (svz@latex.pereslavl.ru).

*Abstract: The new methodology (retrospective shared memory) is presented as a basis for highly available and infinitely adaptive computing services.*

Keywords: large systems control, shared memory, distributed systems, system architecture.

*Статья рекомендована к публикации программным комитетом международной конференции «Параллельные вычисления и задачи управления» (РАСО), Россия, Москва, 24-26 октября 2012 г.  
Поступила в редакцию 21.04.2013.  
Опубликована 31.05.2013.*