

# **Двумерная виртуализация процессора: новый подход к повышению реальной эффективности суперкомпьютеров**

Андрей Ефимов ([eai@kbsp.by](mailto:eai@kbsp.by), [eai\\_andr\\_kb@mail.ru](mailto:eai_andr_kb@mail.ru))

"Конструкторское бюро системного программирования"  
(г.Гомель, Беларусь).

**Реальная эффективность и скорость: вызовы высокопроизводительных параллельных вычислений (ВПВ)**

**60-е годы – наиболее актуальны ВПВ научных расчетов (матфизика)**

**Вызов Сеймур Крэя:** “Каждый : может построить быстрый процессор. Фокус в том, чтобы построить быструю систему”

**Его ответ - иновационный гетерогенный суперкомпьютер CDC 6600 (1960е) :**

- векторный центральный процессор (ЦП);

- мультитрединговый сервисный процессор: 10 тредов в сервисном процессоре обеспечили толерантность доступа быстрой логики в медленную ферритовую память.

**Новое время и актуальности - ВПВ общего назначения (СУБД, облака, реальное время [1-4] )**

**Вызовы Бартона Смита:**

– переизобрести компьютер и вывести ИТ индустрию из спирали специализации, при которой “высокопроизводительными считаются вычисления, хорошо выполняемые существующими высокопроизводительными системами”;

- универсальные микропроцессоры (killer micros) также плохи для параллелизма общего назначения, как и динозавры, которые на них охотятся. ...Что не в порядке с микропроцессорами? ...Ответ очевиден – отсутствие толерантности к тонко-гранулированной латентности всех источников” **(далее – обобщенной латентности).**

**Его ответ:**

- 30-летнее развитие экстремального мультитрединга (ХМТ) в системах HEP Denecollor, HORIZON и Tera MTA, Cray MTA/ХМТ [4] (до 128 тредов) обеспечило толерантность для NUMA и неограниченное масштабирование NUMA-кластеров.

# Обобщенная латентность как основа метрик эффективности и реальной скорости

“Четыре всадника Апокалипсиса” (Томас Стерлинг) вызывают состояния обобщенной латентности (ОЛ), снижающей эффективность и скорость выполнения реальной работы (КПД и  $V_r$ ):

**накладные расходы** -- дополнительные работы аппаратуры на управление параллельностью;

**латентность** - простои при доступе к памяти или другим частям системы;

**конкуренция** - бесполезная активность при соперничестве за разделяемые ресурсы;

**голодание** - простои из-за слабого параллелизма и дисбаланса загрузки (в частности, из-за стены ILP в суперскалярных и VLIW процессорах).

**На основе ОЛ** в дополнении к экономическому показателю эффективность/стоимость, определим **чисто технический показатель**

**КПД =  $1 - C_{ол} / C$** , где  $C$  – стоимость аппаратуры при выполнении работы, а  $C_{ол}$  - стоимость ее части, находившей в состоянии обобщенной латентности.

**Определим реальную скорость  $V_r$**  (пропорциональную КПД) как

**$V_r = V_b * VM$** , где  $V_b$  (байт в секунду) – скорость обработки операндов в приведенном к байту формате, а  $VM$  (байт) – объем памяти системы, на которой обеспечивается данная скорость  $V_b$ .

# Повышение эффективности и реальной скорости на основе 2х-мерной виртуализации

**Рост скорости ВПВ исторически обеспечивался аппаратной виртуализацией процессоров в 2х ортогональных направлениях - памяти и тредов, как схематически показано на следующем слайде.**

**По горизонтали** – типы памяти, их размеры и оценки времен доступ в процессорных циклах (clk).

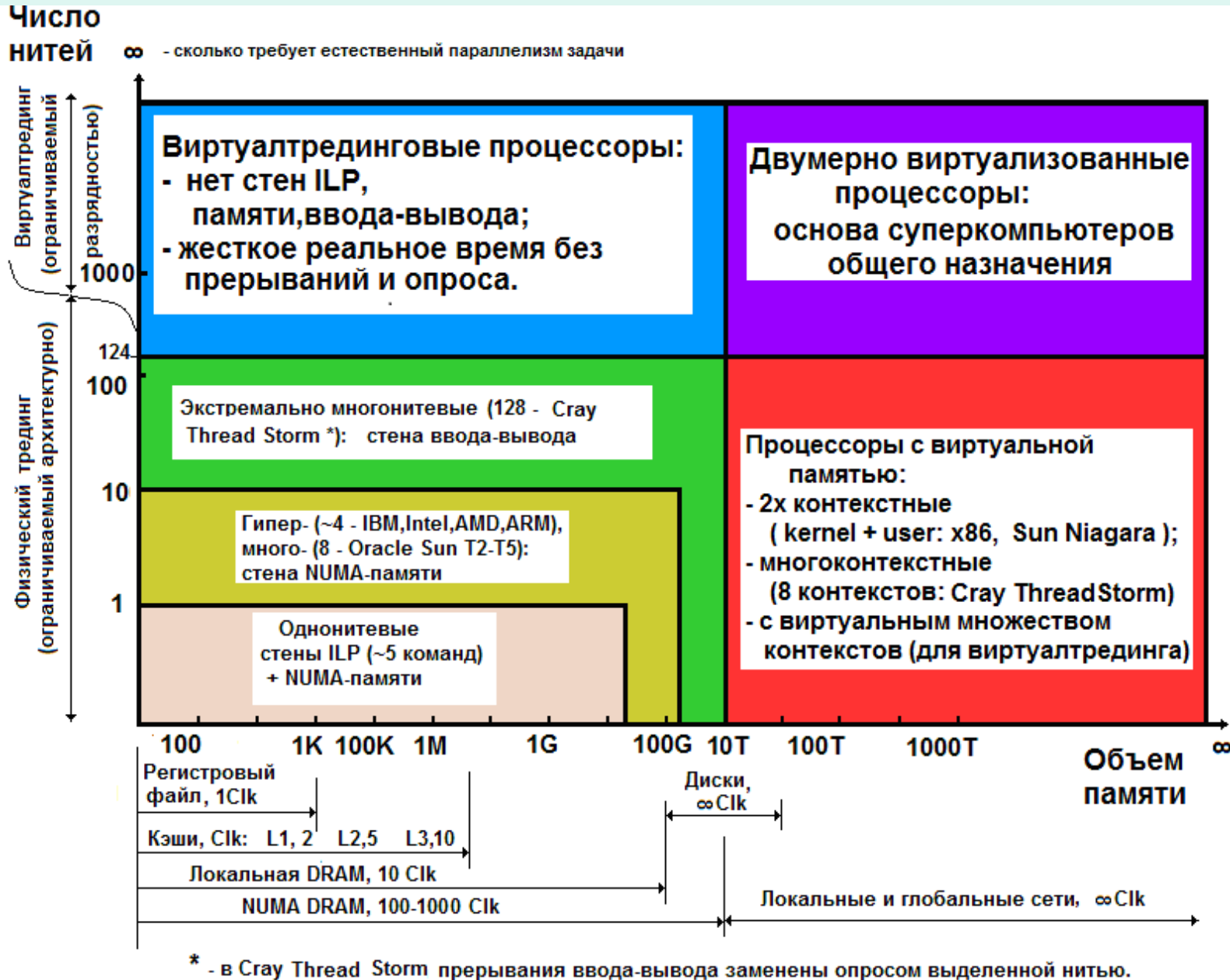
**По вертикали** – типы систем и количество тредов в них.

**Выделенные цветом прямоугольники связывают типы систем и типы памяти, к которой эти системы обеспечивают толерантный доступ:**

- простые 1-тредовые - к данным в регистровом файле;
- 1-тредовые с ILP (суперскалярный и VLIW) – к кэшам L1, L2 и частично к локальной памяти (совмещения обращения к памяти с вычислениями, но ILP wall позволяет запускать в такте около 5 команд);
- гипертрединговые, используя ILP и TLP – к локальной памяти;
- ХМТ-процессоры (Cray Thread Storm, [4]) (ILP + экстремальный TLP), исполняя в общем конвейере потоки 128 тредов, толерантны к NUMA.

**Но ХМТ-процессоры не обеспечивают толерантность доступа к виртуальной памяти, отображаемой на внешние устройства (локальные или сетевые диски) – они сталкиваются со стеной, которую естественно назвать стеной ввода-вывода (IO wall).**

# Двумерная виртуализация процессоров



## От к ХМТ к виртуалтредингу

Увеличение числа тредов в ХМТ до некоторого предела снижает стоимость оборудования, повышает КПД и реальную скорость  $V_r$ . Но это увеличение приводит к пропорциональному росту размера микроархитектурного регистрового файла (МРФ) и квадратичному росту его связей с конвейером, что приводит к снижению КПД из-за роста объема латентных элементов МРФ и связей.

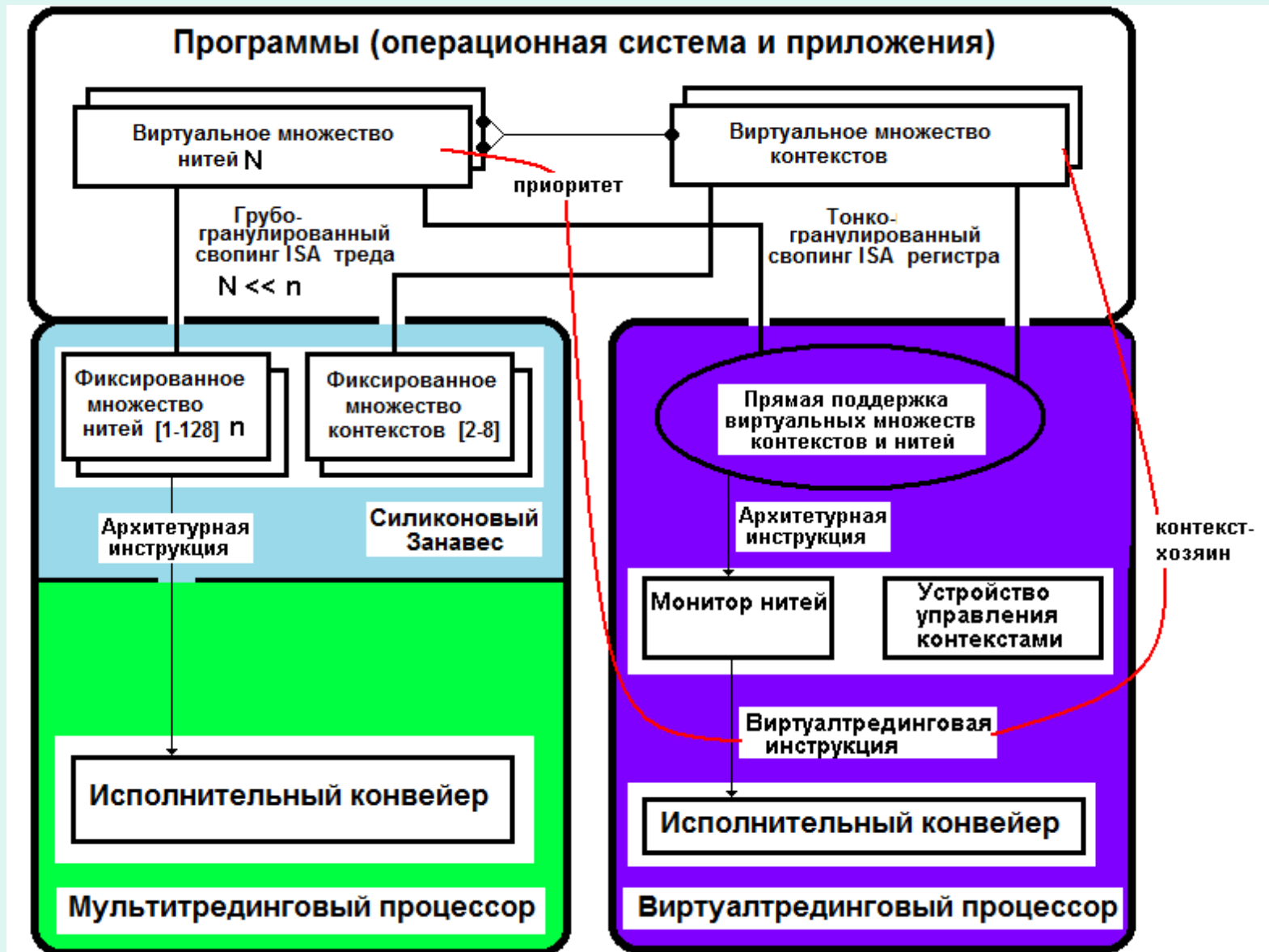
Способ уменьшения размера МРФ за счет тонко-гранулированной виртуализации предложен в архитектуре виртуального контекста (VCA) [6], которая наиболее близка к предлагаемой виртуалтрединговой. В VCA впервые тонко-гранулированным элементом свопинга стал отдельный архитектурных регистр в отличие от крупного гранулированного свопинга полного набора регистров треда в ХМТ. Однако VCA оптимизирует только регистровый файл рабочего множества нитей, а поддержка полного множества оставляется операционной системе.

Предлагаемая в докладе виртуалтрединговая архитектура (ВТА) является развитием архитектур ХМТ и VCA и представляется их логическим завершением.

# Основные концепции виртуалтрединговой системы

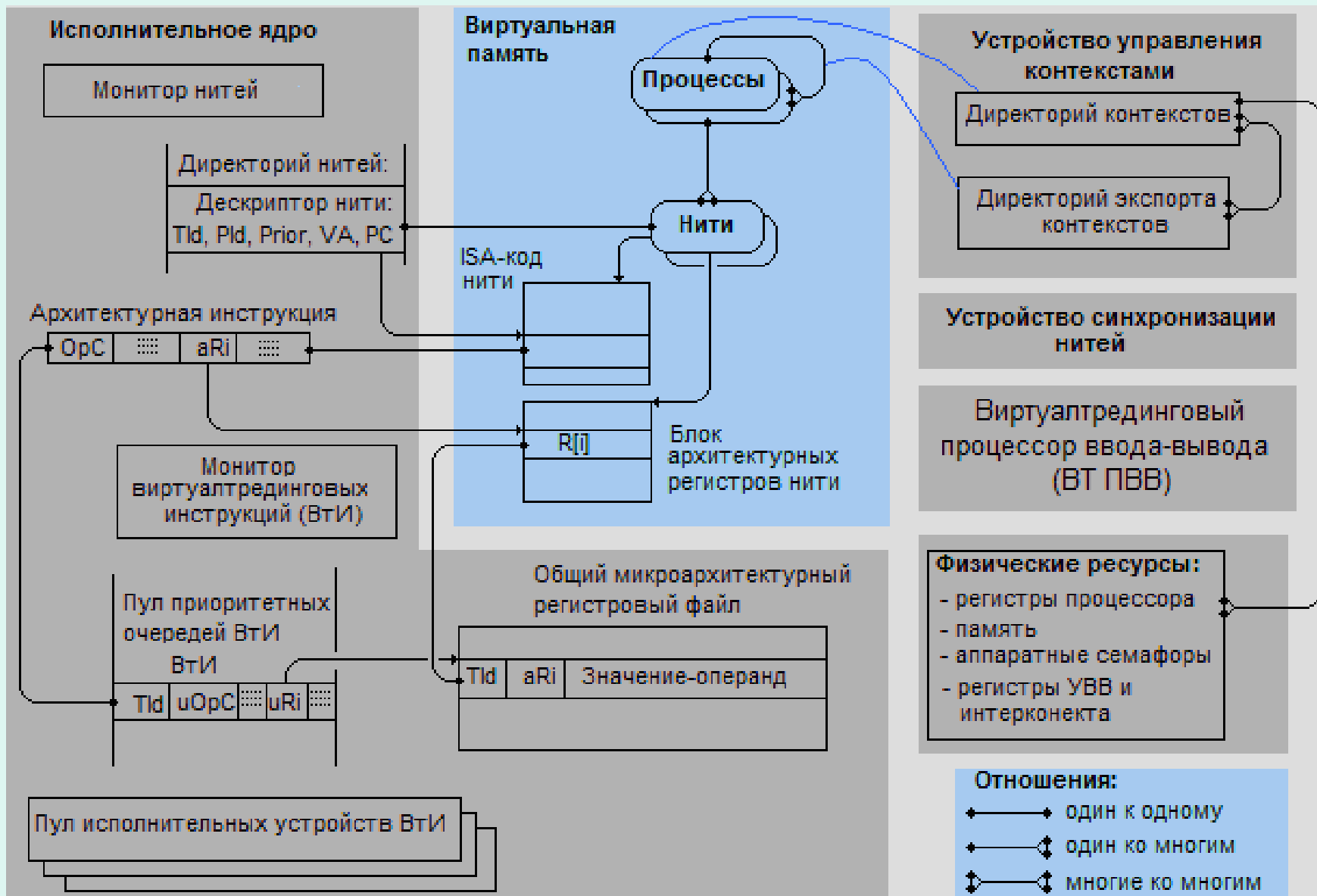
- 1. Устранена концептуальная причина низкой эффективности существующих реализаций ВПВ - Силиконовый Занавес в виде фиксированного множества тредов и аппаратных контекстов** (схема на следующем слайде). Виртуалтрединговая ISA и транзакционная микроархитектура поддерживают, в развитие известной VCA, виртуальные множества программных нитей и процессов (контекстов).
- 2. Транзакционная микроархитектура реализует прямую тонко-гранулированную приоритетную виртуализацию активностей полного множества нитей** за счет замены тредов как аппаратной единицей активности их своих тонко-гранулированными элементами – виртуалтрединговыми инструкциями (ВТИ), которые аппаратура генерирует динамически для полного множества нитей. При этом каждая ВТИ наследует все атрибуты, представляющие порождающую нить как объект управления в ОС, в том числе принадлежность к процессу-хозяину и приоритет.
- 3. Устройство выборки команд и исполнительный конвейер в виртуалтрединговой архитектуре обеспечивают параллельное приоритетное исполнение полного множества нитей** (вычислительных и ввода-вывода) за счет исполнения ВТИ без участия программного обеспечения.

# Различия мультитредингового и виртуалтредингового процессоров





# Эталонная модель виртуалтрединговой СнК



# Двухуровневая организация исполнения инструкций

**Монитор нитей реализует базовую ISA** для обеспечения полной преемственности унаследованного ПО и **виртуалтрединговую ISA**, в которой системные вызовы ОС, реализующие управление процессами, разделяемой памятью, нитями, их синхронизацией и вводом-выводом заменены виртуалтрединговыми инструкциями (ВТИ).

**Исполнение реализуется по двухуровневой схеме:**

**уровень 1** - монитор нитей потактово анализирует готовности к исполнению и приоритет полного множества нитей и формирует на этой основе из первичных ISA-инструкций пул приоритетных очередей ВТИ .

**уровень 2** - монитор ВТИ ведет пул приоритетных очередей ВТИ в микроархитектурной памяти, через который он взаимодействует с монитором нитей и пулом исполнительных устройств.

Взаимодействуя с монитором нитей, монитор ВТИ принимает ISA инструкции, переводит их в формат ВТИ и помещает в пул ВТИ. После отработки ВТИ пулом исполнительных устройств монитор ВТИ выдает результаты в монитор нитей.

# Организация управление доступом к разделяемым данным, синхронизация и ввод-вывод

**Устройство управления контекстами** совмещает функции устройств управления памятью и памятью ввода/вывода и **обеспечивает аппаратную поддержку полного множества контекстов** - ВТИ любой нити и контекста могут обращаться к данным других контекстов за счет прямой аппаратной поддержки экспорта разделяемых регионов с использованием директория экспорта контекстов.

**Устройство синхронизации нитей** реализует работу в критических секциях на основе аппаратных семафоров.

**Виртуалтрединговый процессор ввода-вывода** позволяет организовать обмены с внешними устройствами без прерываний.

Логику совместной работы этих устройств иллюстрируют следующие слайды.

# Транзакционная микроархитектура

**Монитор ВТИ – ключевой элемент транзакционной микроархитектуры, поддерживающий ВТИ на разных стадиях исполнения в конвейерах исполнительных устройств, конкурируя за ее ресурсы – обрабатывающую логику, связи и микроархитектурную память. Размещение всех операндов завершает этап подготовки операндов ВТИ как транзакции, которая далее может быть отработана в нескольких фазах исполнения и ожидания. Важнейшим свойством виртуалтрединга является поддержка перехода как ВТИ, так и породившей нити, в состояние ожидания произвольной длительности без участия ОС.**

Транзакционная микроархитектура дает ответ практически на все приведенные в [3] требования к “переизобретенному” универсальному суперкомпьютеру:

- приоритетный пул ВТИ можно рассматривать как пул транзакций, для которого клиентами являются монитор и исполнительные устройства ВТИ;
- в терминах требований [3], серверами для ВТИ, являются:
- транзакционная память, используемая в существующих системах;
- устройство синхронизации нитей средствами аппаратных семафоров;
- виртуалтрерединговый процессор ввода-вывода, обеспечивающий организацию обменов без прерываний.

# Синхронизация без инверсии приоритетов

**Прямая аппаратная поддержка монитором ВТИ смены состояний “ожидание-активность” ВТИ** (и, как следствие, их нитей-хозяев) позволяет заменить в системах на основе виртуалтрединга программные функции ОС аппаратными инструкциями синхронизации.

**ВТИ синхронизации обрабатывает устройство синхронизации нитей как транзакции, на основе аппаратных семафоров -- активных объектов, реализующих работу таймера и ведение очередей ВТИ.** Это устройство автономно, локально для каждого семафора, реализует в расширенном виде объединенную функциональность имеющихся в современных ОС мутексов (от MUTual EXclusion) и условных переменных (condition variables) (следующий слайд).

**Аппаратное таймирование и повышение приоритета нити при работе в критической секции** очень просто и эффективно **решают проблему инверсии приоритетов**, особенно актуальную при работе в жестком реальном времени.

Кроме того, замена системных вызовов аппаратными инструкциями синхронизации, делает ненужными фьютексы (Fast User space muTEXes).

# Инструкции синхронизации

**Виртуалтрединговые средства синхронизации** представлены следующими инструкциями, реализующими наиболее важные функции синхронизации стандарта POSIX:

- **SGet** с операндами **tw** –(тайм-аут ожидания входа нити в критическую секцию – (КС) ) и **prs** (приоритет при работе в критической секции) объединяет функций **pthread\_mutex\_init** и **pthread\_cond\_init**, выделяет программе аппаратный семафор (из пула свободных) и возвращает его указатель **ps**;
- **SLock(ps, ts)** заменяет функцию **pthread\_mutex\_lock** и обеспечивает либо вход нити в КС с тайм-аутом пребывания в ней **ts**, приоритетом **prs** (если эта секция свободна), либо ожидание входа в нее в течении **ts** до выдачи команды **SPas(ps)**;
- **SWait(ps, ts)** (аналог **pthread\_cond\_timedwait**) реализует неактивное ожидание входа в КС до выдачи команды **SPas(ps)**, позволяя устранить ресурсоемкий опрос разделяемых переменных в критической секции;
- **SPas(ps)** объединяя действия функций **pthread\_mutex\_unlock** и **pthread\_cond\_broadcast**, выводит нить из критической секций, охраняемой семафором **ps** и освобождает ее для других нитей.

# Ввод-вывод: Interruptions considered harmful

**Ввод-вывод в виртуалтрединговой СнК использует однородное управление вычислительными нитями и нитями ввода-вывода**, причем вычислительные нити управляют нитями ввода-вывода как ведомыми. Виртуалтрединговый процессор ввода-вывода (ПВВ) организован аналогично исполнительному ядру и поддерживает полное множество нитей ввода-вывода для каждого канала – независимой активности ввода-вывода.

**Выделение пары нитей каждому каналу реализует прерывания как сообщения**, отображающие завершение обращений в удаленную память с разницей, что при обращениях в память агентами являются регистр и память, а при вводе-выводе – оперативная память и память устройства.

Поддерживая виртуальное множество контекстов и нитей обоих типов, **виртуалтрединг устраняет понятия “аппаратное ядро”** как средство поддержки треда и ядра как средств реализации физического параллелизма в связанном контексте. **Это свойство виртуалтрединга позволяет утверждать, что прерывания в современных СнК еще более вредны, чем оператор goto в языках программирования.** Отсутствие прерываний существенно упрощает как организацию фон-Неймановских машин, так и их программирование.

# Вместо заключения: Время собирать в камни

**Мир ИТ становится экстремально параллельным**– запустив диспетчер задач Windows на своем ПК, каждый может увидеть стократное превышение числа программных нитей над числом тредов.

**Виртуалтрединговая архитектура (ВТА) позволяет устранить Силиконовый Занавес аппаратных тредов и “собрать в камни” прямое тонко-гранулированное приоритетное мультипрограммирование (= виртуализацию активностей) и на этой основе объединить программистов и инженеров-схемотехников для переноса Силиконового Занавеса в самое подходящее место – перед всадниками Апокалипсиса. Это позволит объединившимся комплексным специалистам, как и пользователям их разработок, избежать судьбы динозавров и обеспечить длительную эволюции в таком параллельном мире.**

**Введенные КПД и метрика реальной скорости** могут стать основой чисто технической системы оценки суперкомпьютеров.

**При конвергенции усилий** сообщества российских специалистов **может быть достаточно быстро создана виртуалтрединговая СнК как фундамент для российских суперкомпьютеров общего назначения.**

Одним из первых камней в этом фундаменте (сделанным, возможно силами студентов и аспирантов технических университетов ) должна стать СнК с виртуалтрединговой реинкарнацией выдающейся русскоязычной системы Эльбрус-2. Это будет хорошей подстраховкой переводу ряда важнейших систем на совершенно новые аппаратные и программные платформы.



# Список цитируемых работ

1. Smith, J.E. et al, Future general purpose supercomputer architectures. Proceedings of the 1990 ACM/IEEE conference on Supercomputing.
2. Burton Smith, The Quest for General Purpose Parallel Computing 1994, Developing a computer science agenda for high-performance computing, ACM New York, NY, 1994
3. Burton Smith, Reinventing Computing Microsoft Research Faculty Summit 2007, <http://www.cct.lsu.edu/~estrabd/LACSI2006/Smith.pdf>
4. Burton Smith, The end of architecture, ACM SIGARCH Computer Architecture News Homepage archive, Vol. 18, № 4, Dec. 1990.
5. Meeting the Demands of Unstructured Data with an extreme multithreaded machine. Introducing the Cray XMT Supercomputer, <http://www.cray.com/Assets/PDF/products/xmt/CrayXMTOverviewWhitepaper.pdf>
6. How to Fake 1000 Registers, Proceedings of the 38th annual IEEE/ACM International Symposium on Microarchitecture, <http://aggregate.ee.engr.uky.edu/LAR/micro05.pdf>

# Спасибо за внимание

Андрей Ефимов

[eai@kbsp.by](mailto:eai@kbsp.by)

[eai\\_andr\\_kb@mail.ru](mailto:eai_andr_kb@mail.ru)

"Конструкторское бюро системного программирования"  
(г.Гомель, Беларусь).