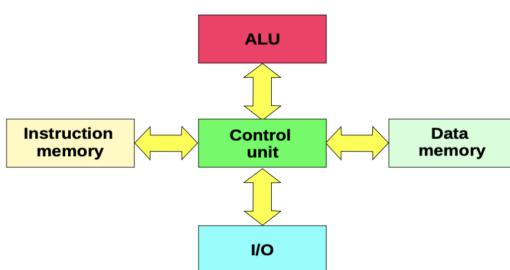
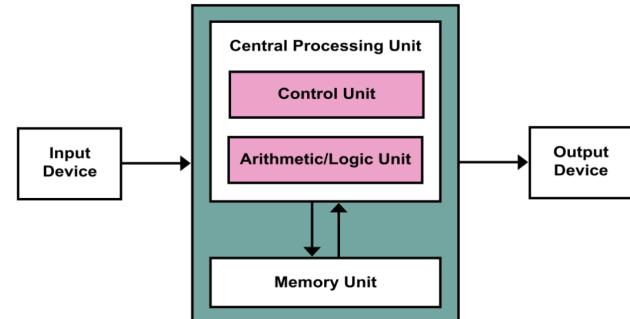


1. Архитектура компьютерных систем. Архитектура Фон-Ндская архитектура. Принципы архитектуры Фон-Неймана и Гарвардской. Архитектуры NUMA и UMA.

Гарвардская архитектура



Архитектура фон Неймана



Существует два основных вида архитектуры: фон Неймана и Гарвардская.

В то фон Неймана данные и команды хранятся в одной и той же памяти

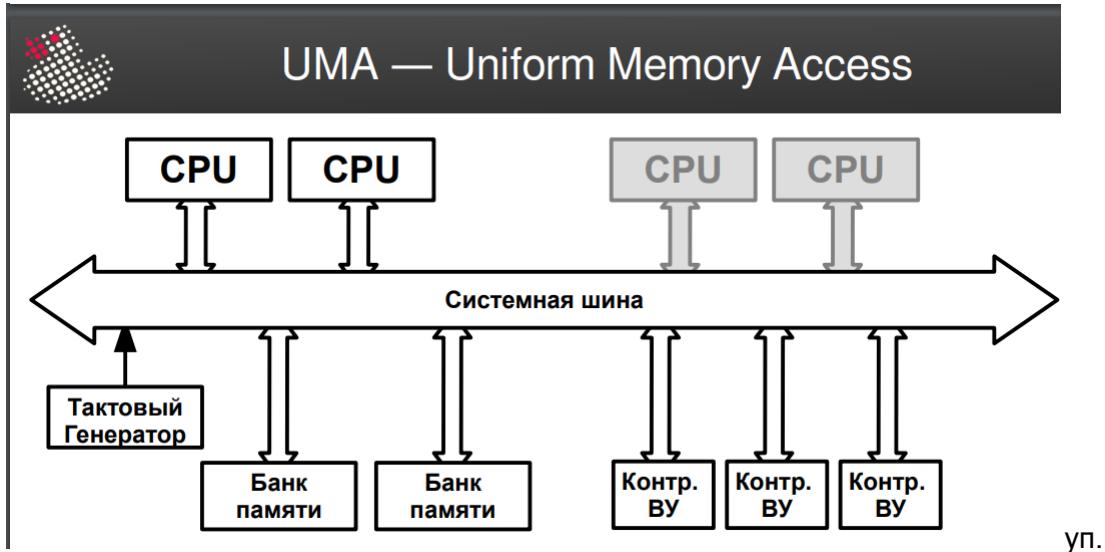
В Гарвардской - отдельная память для команд, отдельная для хранения данных

- Принципы Фон-Неймана:
 - однородность памяти (команды и данные хранятся вместе, внешне неразличимы)
 - принцип адресности(у каждой ячейки памяти есть свой адрес, к которому имеет доступ процессор)
 - принцип программного управления(вычисления представлены в виде программы, состоящей из последовательности команд)
 - принцип двоичного кодирования(данные и команды кодируются нулями и единицами)
- Принципы Гарвардской
 - отдельно данные, отдельно инструкции
 - канал инструкций и канал данных также отличаются

С точки зрения пользовательских штук две штуки UMA и NUMA:

- UMA (Uniform Memory Access / однородный доступ к памяти): Все устройства одного ранга и все устройства имеют одинаковый доступ к памяти, а значит

время доступа для всех одинаково. Задержка одинакова и тд



уп.

Плюсы:

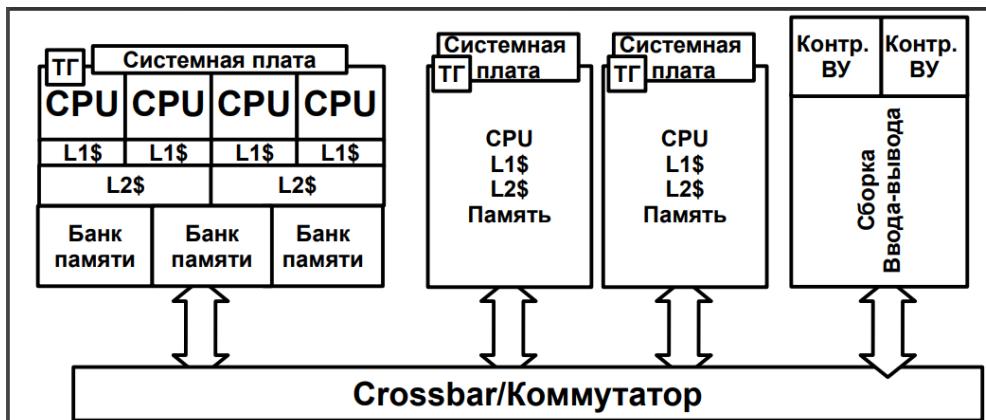
- одинаковое время доступа к памяти
- простота реализации

Минусы:

- Плохая масштабируемость (с увеличением количества процессов возникают конфликты и конкуренция за доступ к памяти, а это замедляет работу системы)

Работает нормально только до двух процессоров.

- NUMA (Non Uniform Memory Access / неоднородный доступ к памяти)



- все системные ресурсы расположены на отдельных системных платах
- имеется локальная память (расположенная "ближе" к процессору)
- общая память и устройства ввода-вывода доступны через коммутатор, координирующий взаимодействие компонентов
- Используется в серверах
- адресное пространство общее для всех процессоров!!!

Плюсы:

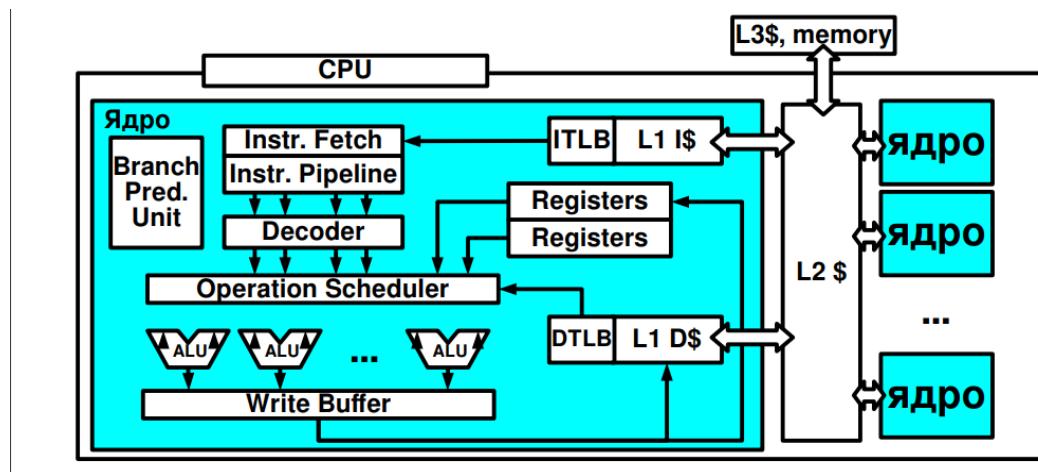
- возможность замены практически любого компонента во время работы системы (hot swap / горячая замена)
- Хорошая масштабируемость
- Надежность и отказоустойчивость

Минусы:

- требовательна к операционной системе(она должна поддерживать все эти фишки)
- важно учитывать время доступа и расположение данных в памяти(так как время доступа зависит от того, где данные расположены)

2. Общая организация процессора, памяти, организация вычислений.

Организация процессора



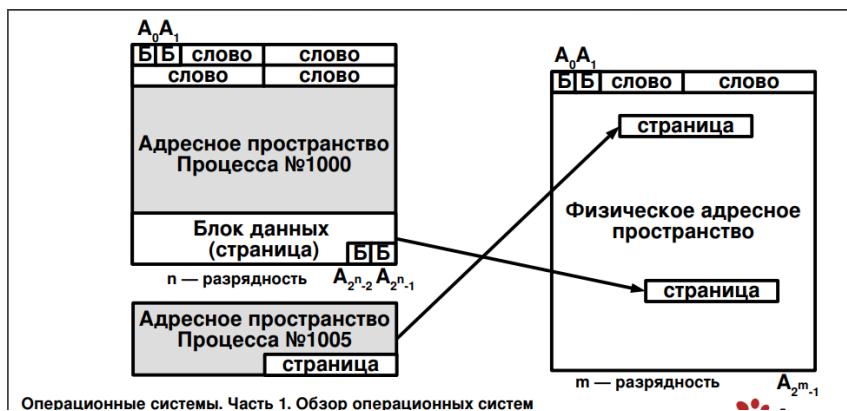
- Кэш второго уровня (L2 \$): кэш общий для всех ядер
- Кэш первого уровня (L1\$): разделен на кэш с данными (L1 D\$) и кэш с инструкциями (L1 I\$), оба имеют буфер TLB, это помогает параллельно читать данные и инструкции
- Буфер ассоциативной трансляции ([TLB](#)): хранит записи, которые помогают быстрее найти записи внутри кэша или понять есть они в кэше этого уровня или нет (ускоряет трансляцию адреса виртуальной памяти в адрес физической памяти)
- Hyper-threading (два блока Registers), суть технологии [hyper-threading](#): у процессора есть два набора регистров, что позволяет быстрее переключать контекст двух потоков. Для операционной системы это выглядит как два логических процессора.
- предсказатель переходов (branch prediction unit): в случае когда выполняется условие в условном переходе, загруженные команды оказываются ненужными, а значит конвейер встанет, чтобы избежать этого создан предсказатель переходов, который дает понять какие команды грузить.
- [Вычислительный конвейер](#) (Упомянуто в лекции криво, лучше почитать вики или [тут](#))

- параллельность циклов исполнения и выборки команды за счет множества АЛУ, позволяющих распараллелить очередь операций
 - конвейер - получает команду с устройства выборки команды, декодирует и отдает планировщику операций
 - планировщик - решает пустить пришедшие команды параллельно или нет и пустить ли вообще команду на выполнение, в зависимости от наличия данных
 - branch pred unit - предсказатель перехода, для оптимизации загрузки команд (чтобы конвейер не останавливался)
 - write buffer нужен т. к. скорость записи в память может не успевать за ходом исполнения команд
 - есть обычные алу а есть специфичные (для плавающей точки, для векторных операций и т.п.)
-

Есть еще картинка из Столлингса попроще: СК,РК,АЛУ, РА,РБ,регистр адреса ввода-вывода, регистр буфера ввода-вывода



Организация памяти

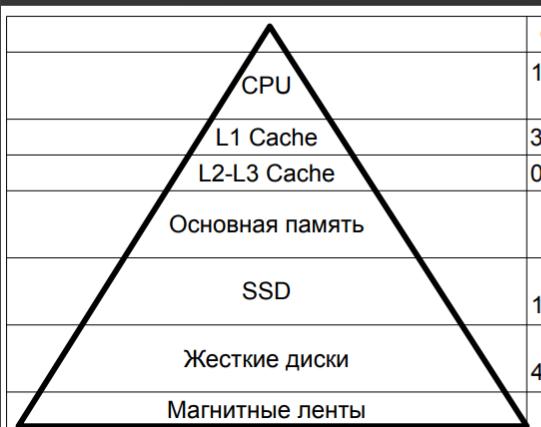


Операционные системы. Часть 1. Обзор операционных систем

- Виртуальная память – способ организации памяти, при котором каждый процесс имеет собственное адресное пространство, думая, что у него есть своя независимая память.
- Адреса виртуальной памяти преобразуются в физические при помощи таблиц трансляций.

- Единица доступа к информации – байт. Данные объединены в страницы (размер обычно 4 Кб).
- Large Pages – страницы с большим объемом данных (к примеру 1 Мб)

Преимущество: уменьшение области трансляции при выделении адресного пространства, релокация, swapping, защита



	Объем	Тд	*	Тип	Управл.
CPU	100-1000 б.	<1нс	1с	Регистр	компилятор
L1 Cache	32-128Кб	1-4нс	2с	Ассоц.	аппаратура
L2-L3 Cache	0.5-32Мб	8-20нс	19с	Ассоц.	аппаратура
Основная память	0.5Гб-4Тб	60-200нс	50-300с	Адресная	программно
SSD	128Гб-1Тб/drive	25-250мкс	5д	Блочн.	программно
Жесткие диски	0.5Тб-4Тб/drive	5-20мс	4м	Блочн.	программно
Магнитные ленты	1-6Тб/к	1-240с	200л	Последов.	программно

Память внутри компа организована иерархически. Самая быстрая и дорогая память наверху, снизу наибольший объем. Сделано так, чтобы к более медленной памяти было меньше обращений.

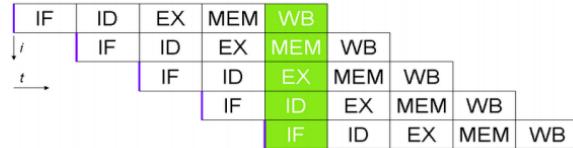
Тд(время доступа) – интервал между первым запросом доступа к памяти и первыми данными оттуда

Виды управления:

- Компилятор: управляет при помощи ключей компиляции (уровни оптимизации)
- Аппаратура: поиск осуществляется при помощи схем в ассоциативной памяти (TLB)
- Программно

Организация вычислений

- Ядро выполняет каждую команду последовательно. Цикл команды:
 - Выборка команды (Instruction Fetch)
 - Декодирование инструкций (Instruction Decode)
 - Исполнение (Execution)
 - Чтение памяти (MEM)
 - Запись (Write Back)



Конвейерная обработка команд. Следующая команда начинает выполняться, пока не закончила выполнение предыдущая.

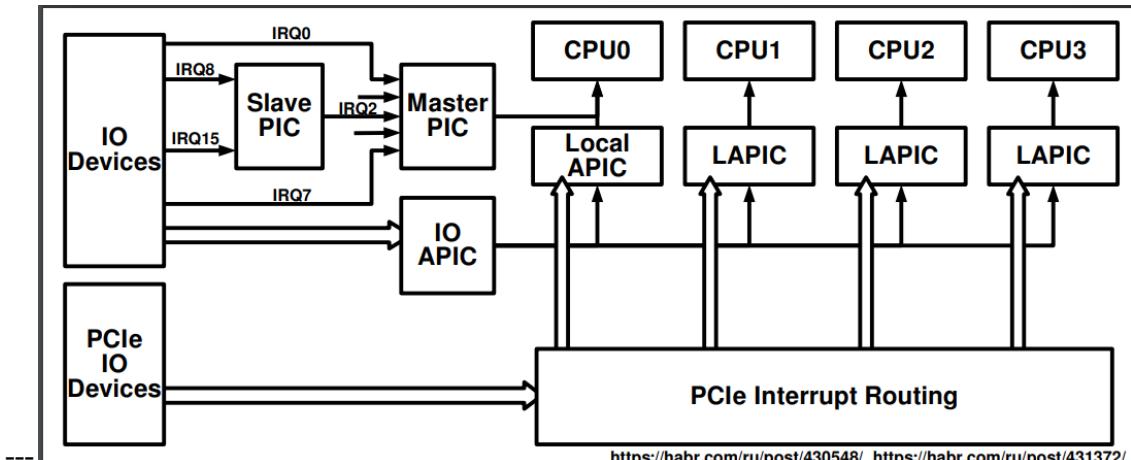
3. Организация прерываний, типы прерываний, контроллер прерываний.

Прерывание (англ. interrupt) — сигнал от программного или аппаратного обеспечения, сообщающий процессору о наступлении какого-либо события, требующего немедленного внимания. Прерывание извещает процессор о наступлении высокоприоритетного события, требующего прерывания текущего кода, выполняемого процессором. Процессор отвечает приостановкой своей текущей активности, сохраняя свое состояние и выполняя функцию, называемую обработчиком прерывания (или программой обработки прерывания), которая реагирует на событие и обслуживает его, после чего возвращает управление в прерванный код.

- Могут быть вызваны самой программой или сигналом с устройств IO.
- Более высокий приоритет прерываний могут прерывать более низкие.
- Выполняются в конце цикла команды
- Могут быть вложенными. Есть ограничения по вложенности обработчика в зависимости от архитектуры.

КОНТРОЛЁР ПРЕРЫВАНИЙ

Контроллер прерываний (англ. Programmable Interrupt Controller, PIC) — микросхема или встроенный блок процессора, отвечающий за возможность последовательной обработки запросов на прерывание от разных устройств



<https://habr.com/ru/post/430548/> <https://habr.com/ru/post/431372/>

Надо читать всю статью с [хабра](#)

Слайд интерпретируется так:

1. slave-master PIC - был придуман для увеличения линий прерываний
 2. IO APIC, LAPIC - был придуман для поддержки многопроцессорности
 3. PCIe_device -> {PCIe_Interrupt_Routing} -> Local_APIC -> CPU - пришло с появлением PCIe, информация о прерывании пишется (почти) напрямую в память LAPIC. {PCIe_Interrupt_Routing} - обобщение представляющее роутинг.
-

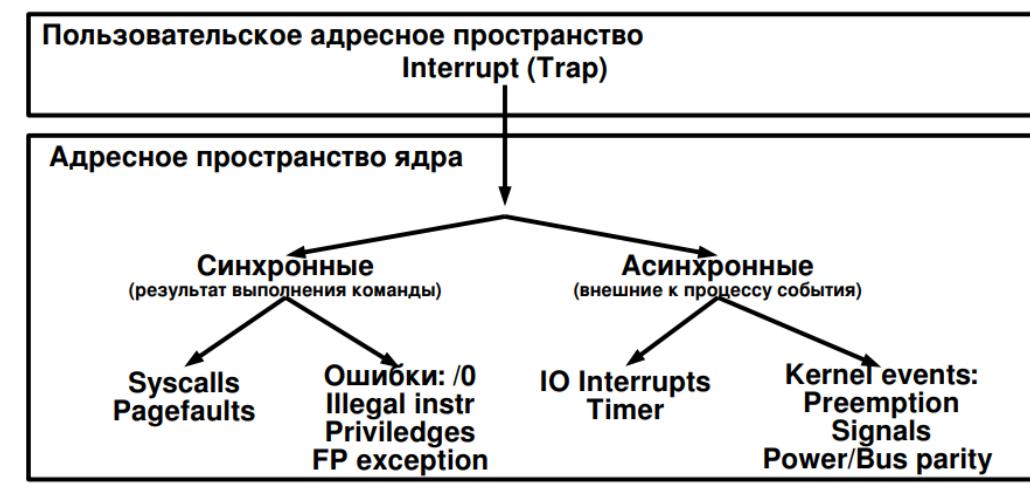
Разные подходы:

- запрет повторных прерываний, пока обрабатываются предыдущие (последовательная обработка прерываний).
- обработка прерываний по уровню приоритета

(схема для x86 проца)

В контроллере прерываний есть несколько процессоров. У каждого процессора свой локальный маршрутизатор/контроллер прерываний, обработчик прерываний. Есть устройства IO, посылающие запрос на прерывание, логика, занимающаяся обработкой, отвечающая за приоритеты и прочее, локальный контроллер посылает процессору запрос на прерывание.

ТИПЫ ПРЕРЫВАНИЙ:



асинхронные, или внешние (аппаратные) — события, которые исходят от внешних аппаратных устройств (например, периферийных устройств) и могут произойти в любой произвольный момент: сигнал от таймера, сетевой карты или дискового накопителя, нажатие клавиш клавиатуры, движение мыши. Факт возникновения в системе такого прерывания трактуется как запрос на прерывание (англ. Interrupt request, IRQ) — устройства сообщают, что они требуют внимания со стороны ОС;

синхронные, или внутренние — события в самом процессоре как результат нарушения каких-то условий при исполнении машинного кода: деление на ноль или переполнение стека, обращение к недопустимым адресам памяти или недопустимый код операции;

4. Типичные функции операционной системы. Интерфейсы ОС. Работа ОС как замена оператора ЭВМ.

Наиболее важной из системных программ является операционная система, которая скрывает от программиста детали аппаратного обеспечения и предоставляет ему удобный интерфейс для использования системы. Операционная система выступает в роли посредника, облегчая программисту и программным приложениям доступ к различным службам и возможностям.

Функции:

- Разработка программ
- Выполнение программ
- Доступ к устройствам ввода-вывода
- Контролируемый доступ к файлам
- Доступ к системе и системным ресурсам
- Обнаружение и обработка ошибок
- Учет использования и диспетчеризация ресурсов
- Предоставление ключевых интерфейсов ОС:
 - ISA (Instruction Set Architecture) — Набор команд
 - ABI (Application Binary Interface) — Бинарный интерфейс приложения
 - API (Application Programming Interface) — Интерфейс прикладных программ

Если кратко, то это менеджмент доступа программ к ресурсам, а также упрощение разработки программ (использование API, как библиотеки низкого уровня)

Существуют редакторы, отладчики(программы-утилиты), которые не входят в состав ядра ОС, но поддерживаются ею, и помогают в разработке программ.

Контролируемый доступ к файлам состоит в том, что ОС определяет природу устройств и структуры их данных.

Структура системы команд (instruction set architecture — ISA). Определяет набор команд машинного языка, которые может выполнять компьютер. Этот интерфейс является границей между аппаратным и программным обеспечением. Обратите внимание!

Бинарный интерфейс приложения (application binary interface — ABI). ABI определяет стандарт бинарной переносимости между программами. ABI определяет интерфейс системных вызовов операционной системы и аппаратных ресурсов и служб, доступных в системе через пользовательскую ISA.

Интерфейс прикладного программирования (application programming interface — API). API обеспечивает программе доступ к аппаратным ресурсам и службам, доступным в системе через пользовательскую ISA с библиотечными вызовами на языке высокого уровня. Обычно любые системные вызовы выполняются через библиотеки. Применение API обеспечивает легкую переносимость прикладного программного обеспечения на другие системы, поддерживающие тот же API, путем перекомпиляции.

До оператора ученые = программисты сами набирали программу, запускали и уходили на какое-то время, пока она не выполнит свою работу. Было расписание, по которому пользователи приходили со своими программами.

Потом компьютеры стали быстрее, стало сложно соблюдать такое расписание, появились операторы, которым отдавались программы и которые уже сами, по мере освобождения компьютера, запускали новую программу.

- В первых ЭВМ был только пульт управления
- Оператор должен был:
 - получить программу с данными от программиста;
 - подготовить программу к загрузке (н-р, с перфокарт);
 - загрузить программу и компилятор;
 - запустить программу на вычисление;
 - распечатку с результатами передать программисту.
- Минусы:
 - Наличие расписания машинного времени
 - Долгое время подготовки к работе

5. Пакетная обработка. Системный монитор.

Потом появился системный монитор

По сути весь доп. софт, который был нужен для работы программы, заранее загрузили в ЭВМ, плюс добавили прерывания, которые нужны, чтобы сигнализировать о том, что программа выполнилась, и планировщик, чтобы по этому сигналу запускать следующую программу. Программы все еще выполняются последовательно.

Центральная идея, лежащая в основе простых пакетных схем обработки, состоит в использовании программы, известной под названием **монитор** (monitor). Используя операционную систему такого типа, пользователь не имеет непосредственного доступа к машине. Вместо этого он передает свое задание на перфокартах или магнитной ленте оператору компьютера, который собирает разные задания в пакеты и помещает их в устройство ввода данных. Так они передаются монитору. Каждая программа составлена таким образом, что при завершении ее работы управление переходит к монитору, который автоматически загружает следующую программу.

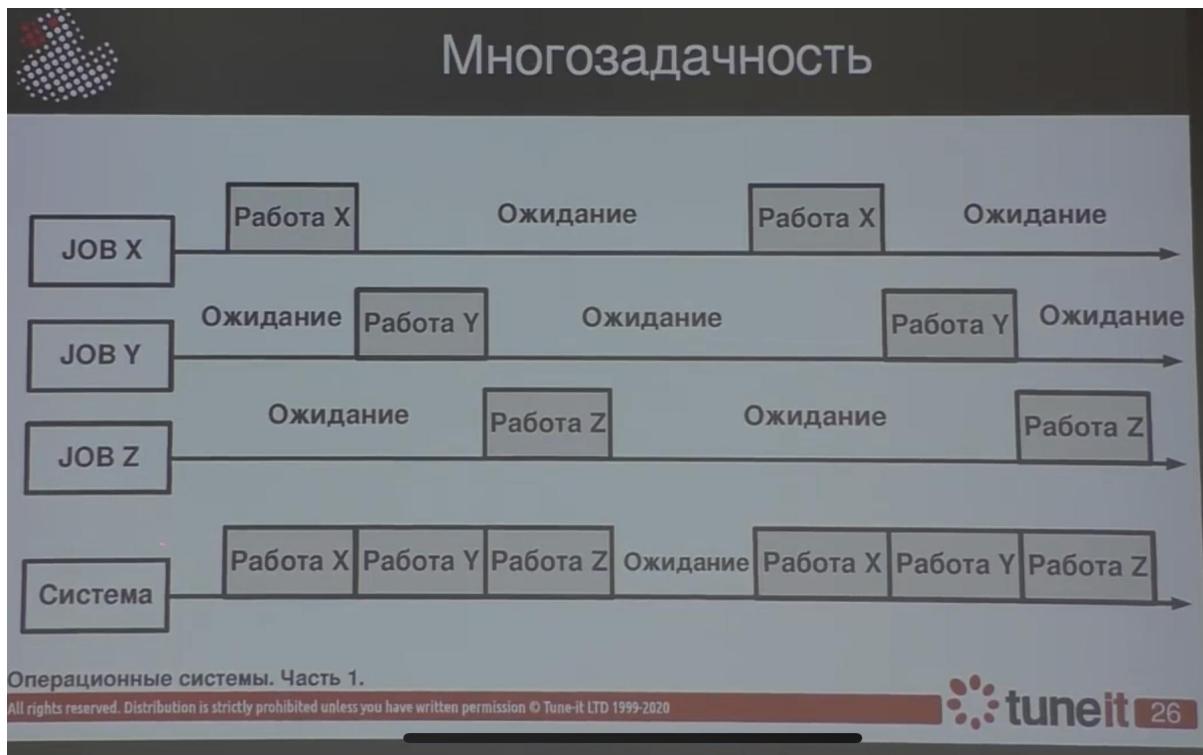
- Машинное время дорогое, его простои необходимо минимизировать
- 1950 г., General Motors, IBM 701
- Наборы программ и данных передавались оператору и запускались



Драйверы устройств – осуществляют работу с устройствами (например, ввода-вывода)

6. Анализ общесистемной эффективности, как предусловие многозадачности. Многозадачность, как способ повышения системной эффективности. Системы разделения времени.

- Одно задание плохо загружает CPU:
 - Read (15 мкс) → Compute 100 instr (1 мкс) → Write (15 мкс)
 - Общая загрука CPU ~ 3.2%
- Давайте запустим много задач, и пока одни занимаются вводом-выводом, другие будут производить вычисления



- Хорошо бы исключить оператора и добавить пользователей!
 - Посадим юзеров за терминалы, пусть сами работают
 - Будем выдавать им часть времени процессора с использованием квантования времени (time slices)
- CTSS (Compatible Time-Sharing System), MIT 1961, IBM 709
 - Выгрузка и загрузка задач
 - 32 пользователя
- Появились проблемы разделения ресурсов и защиты одних программ от других

Пользователей много, всем нужен компьютер, и заметили, что кучу времени процессор простояивает, потому что программе нужно получить данные, например, с диска.

Решили сделать так, чтобы в это время на процессоре выполнялась другая программа, которая ничего не ждет. Концептуально ОС появилась уже тогда, в виде систем разделения времени, но это был по сути только скелет.

7. Процессы, проблемы современных процессов.

Планирование выполнения процессов и управление ресурсами.

Процесс – экземпляр программы во время ее исполнения.

Процесс – это единица активности ОС, в которой существуют последовательные действия, текущее состояние и набор связанных ресурсов. Данное определение с точки зрения ОС.

[статья про процессы](#)

Структура процесса

- Исполняемая программа (процесс представляет собой программу)
- Набор потоков исполнения (внутри программы включено множество потоков исполнения, это по сути единица потребления процессора)
- Связанные структуры ядра (когда процесс создается внутри ОС, ядро должно построить структуры, области памяти, которые содержат описания ресурсов которые процесс потребляет)
- Адресное пространство -- совокупность всех допустимых адресов каких-либо объектов вычислительной системы — ячеек памяти, секторов диска, узлов сети и т. п., которые могут быть использованы для доступа к этим в определенном процессе
- Контекст исполнения (например, регистры) (состояние процесса - то, что нужно восстановить когда мы возвращаемся допустим из другого процесса к этому)
- Контекст безопасности (всякие `suid/sgid` в linux)
- Ресурсы (файлы и прочее) + структуры, связанные с ними структуры
- Динамические библиотеки

Проблемы современных процессов

- Защита памяти процесса (недетерминированное поведение программы, то есть непредсказуемое)
- Взаимные блокировки (`deadlocks`, `starvation`, `livelocks`)
- Проблемы синхронизации
- Взаимное исключение доступа к ресурсам

Динамическая взаимоблокировка (livelock) – Состояние, при котором два или более процессов постоянно изменяют свое состояние в ответ на изменения в другом процессе (или процессах) без какой-либо полезной работы. Это похоже на взаимоблокировку тем, что при этом отсутствует какой-либо прогресс, но отличается тем, что ни один процесс не заблокирован и не ждет чего-либо

Взаимное блокирование (**deadlock**¹) можно определить как *перманентное* блокирование множества процессов, которые либо конкурируют в борьбе за системные ресурсы, либо сообщаются один с другим. Множество процессов оказывается взаимно блокиро-

Голодание (starvation)

Состояние, при котором процесс откладывается на неопределенное время, потому что предпочтение постоянно отдается другим процессам

Планирование выполнение процессов и управление ресурсами

- Равноправие (Пользователи должны получать ресурсы равноправно)
- Дифференциация отклика (необходимость снижения времени отклика)
- Общесистемная эффективность
- Планировщики процессов, дисков и прочее (Разные классы диспетчеризации (Time Sharing, interactives, real time, system, fair share (типа каждому пропорционально приоритету~количество акций), fixed..)) [I/O scheduling](#)

8. Управление памятью, виртуальная память. Защита информации и безопасность ОС.

Управление памятью

1. Изоляция процессов (отслеживание, что ни один из процессов не смог изменить чужую память)
2. Управление выделением и освобождением памяти (программы должны динамически подключаться к памяти - т.е. надо в том числе уметь передать ей управление)
3. Поддержка модулей (возможность определять, создавать, уничтожать и менять размер)
4. Защита и контроль доступа (ОС должна следить каким образом различные пользователи могут осуществлять доступ к различным областям памяти)
5. Долгосрочное хранение - приложениям требуются средства, с помощью которых можно хранить информацию после выключения компьютера.
6. Страницочный обмен ([Подкачка страниц](#) (англ. *paging*; иногда используется термин *swapping* от *swap*, [/swэр/](#)) — один из механизмов [виртуальной памяти](#), при котором отдельные фрагменты памяти перемещаются из оперативной памяти во вторичное хранилище ([жёсткий диск](#) или другой внешний накопитель), освобождая оперативку для загрузки других активных фрагментов памяти. Такими фрагментами в современных [ЭВМ](#) являются страницы памяти.). В винде: page file sys.

Виртуальная память обслуживается MMU и TLB, делает отдельное адресное пространство для каждого процесса и ядра.

Виртуальная память (англ. *virtual memory*) — метод [управления памятью компьютера](#), позволяющий выполнять программы, требующие больше [оперативной памяти](#), чем имеется в компьютере, путём автоматического перемещения частей программы между основной памятью и вторичным хранилищем

- + Существуют невыгружаемые страницы.

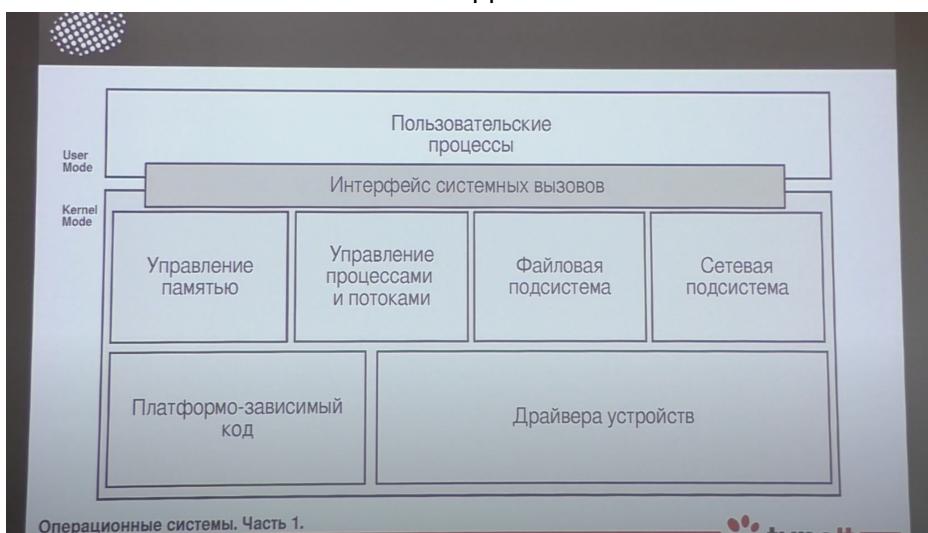
Защита информации и безопасность ОС. Все, что дается по безопасности: регламенты. Не факт, что будет работать, но есть процедуры (протоколы), соблюдение которых может помочь.

- Доступ к системе (защита от несанкционированного доступа)
- Конфиденциальность (невозможность получить доступ без авторизации)
- Целостность данных (защита данных от неавторизованного и нецелостного изменения)

- Аутентификация и авторизация (**Аутентификация** — процедура проверки подлинности, например проверка подлинности пользователя путем сравнения введенного им пароля с паролем, сохраненным в базе данных. **Авторизация** — предоставление определенному лицу или группе лиц прав на выполнение определенных действий.)

9. Структура ядра операционной системы. Архитектуры монолитного ядра, ядра динамически загружаемыми модулями и микроядра.

ЯДРО ОС



Есть пользовательский режим, есть режим ядра, есть интерфейс вызовов, где данные передаются между пользовательскими процессами и самой системой.

1) Монолитное ядро

Да

**Монолитное ядро
(monolithic kernel)**

Большое ядро, содержащее практически всю операционную систему, включая планирование, файловую систему, драйверы устройств и управление памятью. Все функциональные компоненты ядра имеют доступ ко всем его внутренним структурам данных и процедурам. Как правило, монолитное ядро реализовано как единый процесс со всеми элементами, разделяющими одно и то же адресное пространство

Монолитное ядро (аналог мониторов); есть служебная программа, информация загружается в память и там работает; ядро нельзя изменить.

Сейчас в чистом виде не используется нигде, кроме как встроенные системы.

2)Модульное ядро — современная, усовершенствованная модификация архитектуры монолитных ядер операционных систем.

В отличие от «классических» монолитных ядер, модульные ядра, как правило, не требуют полной перекомпиляции ядра при изменении состава аппаратного обеспечения компьютера. Вместо этого модульные ядра предоставляют тот или иной механизм подгрузки модулей ядра, поддерживающих то или иное аппаратное обеспечение (например, драйверов). При этом подгрузка модулей может быть как динамической (выполняемой «на лету», без перезагрузки ОС, в работающей системе), так и статической (выполняемой при перезагрузке ОС после переконфигурирования системы на загрузку тех или иных модулей).

Команды для загрузки и удаления в линукс: `insmod`, `modprobe`, `rmmod`. Удобно включить/выключить не перегружая всю машину.

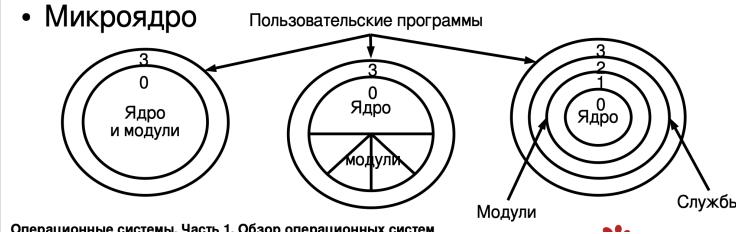
3)Микроядро

Микроядро (microkernel)

Небольшое привилегированное ядро операционной системы, обеспечивающее планирование процессов, управление памятью и службы связи и опирающееся на другие процессы для выполнения некоторых функций, традиционно связываемых с ядром операционной системы

Ядро выполняет базовые функции, а остальные функции осуществляются в виде сервисов и находятся на своих уровнях безопасности. Работает медленно за счет частого переключения контекстов. (пример: Windows NT)

- Монолитное ядро
- Ядро с динамически загружаемыми модулями
- Микроядро



10. Потоки исполнения, многопоточность, модели многопоточности.

Процесс с точки зрения ОС:

- 1)владение процессом(ресурсы)
- 2)планирование/исполнение (как раз наши потоки)

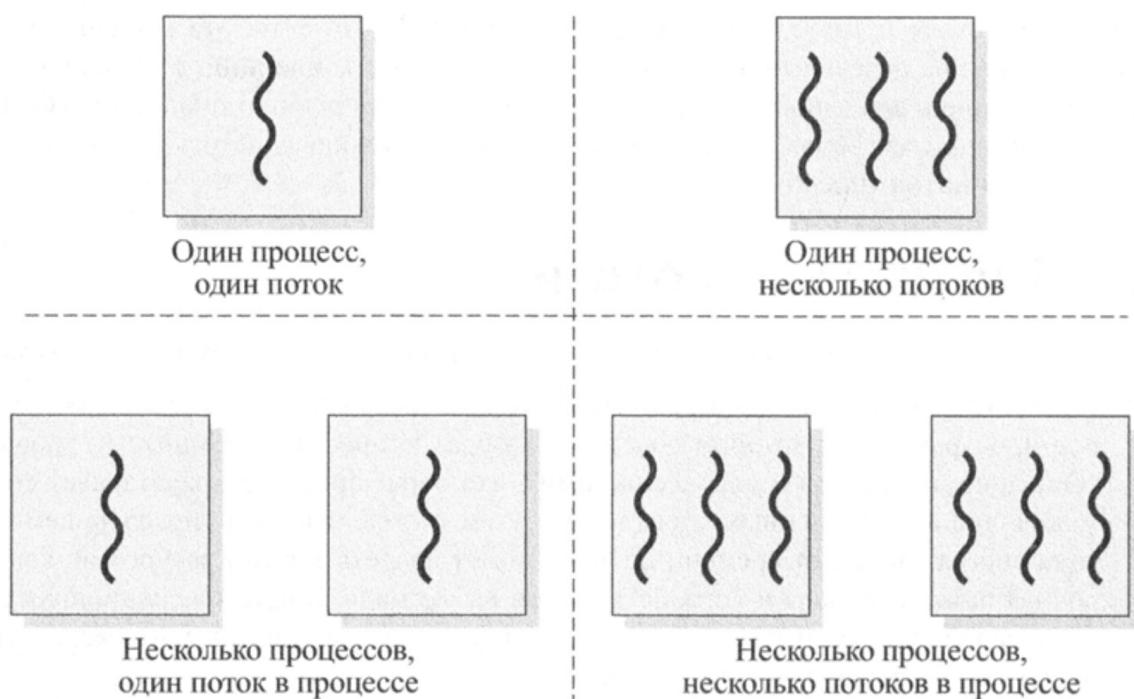
Поток (нить выполнения, thread) - отдельный набор регистров и единица диспетчеризации и выполнения ОС (по сути это элемент выполнения процесса).

Многопоточностью (multithreading) называется способность операционной системы поддерживать в рамках одного процесса несколько параллельных путей выполнения.

Это удобно, потому что создавать весь процесс с нуля - дорого и долго.

Множественные нити исполнения в одном процессе называют *потоками* и это базовая единица загрузки CPU, состоящая из идентификатора потока, счетчика, регистров и стека.

Многопоточные модели:



Потоки POSIX

В конце 1980-х и начале 1990-х было несколько разных API, но в 1995 г. *POSIX.1c* стандартизовал потоки POSIX, позже это стало частью спецификаций *SUSv3*.

В программировании **зелёные потоки** ([англ. green threads](#)) — это **потоки выполнения**, управление которыми вместо **операционной системы** производит **виртуальная машина** (ВМ). Green threads эмулируют многопоточную среду, не полагаясь на возможности ОС по реализации легковесных потоков. Управление ими происходит в [пользовательском пространстве](#), а не пространстве [ядра](#), что позволяет им работать в условиях отсутствия поддержки встроенных потоков.

Преимущество заключается в том, что вы получаете функциональность, подобную обычному треду. Недостатком является то, что зеленые потоки фактически не могут использовать несколько ядер.

- Green threads значительно превосходят встроенные потоки Linux-системы по времени активации потоков и [синхронизации](#).
- Встроенные потоки Linux имеют несколько более высокую производительность операций [ввода-вывода](#) и [переключения контекста](#).

В [Java 1.1](#) green threads являлись единственной потоковой моделью (моделью распараллеливания потоков), используемой в JVM, по крайней мере, в [Solaris](#). Ввиду того, что green threads обладают ограничениями в сравнении с native threads, в последующих версиях Java основной упор сделан на native threads.

11. Симметричная и асимметричная многопроцессорная обработка.



SMP vs ASMP

- Asymmetric Multiprocessing — есть Master CPU, он управляет Slaves CPU
 - CPU — GPU
- Symmetric Multiprocessing – процессоры равны, процесс выполняется на нескольких процессорах одновременно
 - «Простота» разработки и производительность
 - Более высокая надежность. При отказе одного выполнять процессы могут другие
 - Масштабируемость приложений
 - Динамическое добавление ресурсов процессора
- Многопоточность \neq Многопроцессорность

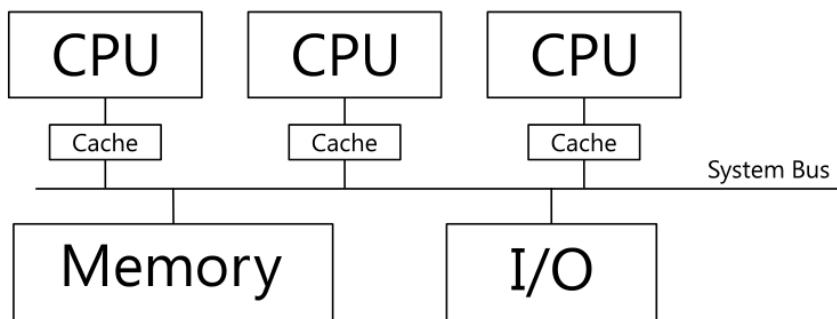
AMP или **ASMP** (от англ.: **Asymmetric multiprocessing**, рус.: **Асимметричная многопроцессорная обработка** или **Асимметричное мультипроцессирование**) — тип многопроцессорной архитектуры компьютерной системы, который использовался до того, как была создана технология [симметричного мультипроцессирования](#) (SMP). Также использовался как более дешевая альтернатива в системах, которые поддерживали SMP.

Асимметричное мультипроцессирование – наиболее простой способ организации вычислительного процесса в системах с несколькими процессорами. Этот способ часто называют также “ведущий-ведомый”.

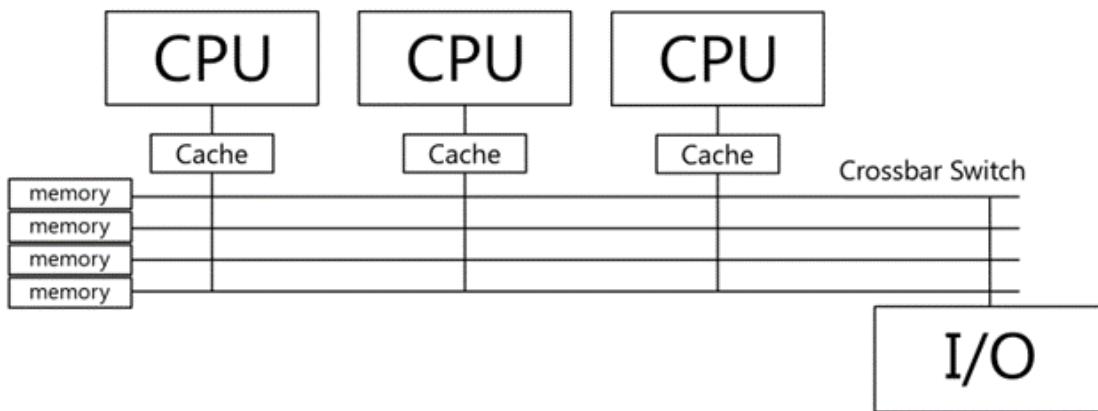
GPU – графический процессор (тут пример: что ЦПУ - главный и распределяет задачи, а ГПУ - просто процесцит переданные задачи).

Симметричная многопроцессорность (англ. *symmetric multiprocessing*, сокращённо **SMP**) — архитектура [многопроцессорных компьютеров](#), в которой два или более одинаковых [процессора](#) сравнимой производительности подключаются единообразно к общей памяти (и периферийным устройствам) и выполняют одни и те же функции (почему, собственно, система и называется *симметричной*). Просто, потому что все одинаково, но надо бороться со всякими локами и захватами одного и того же участка. Каждый процесс управляет собой.

Разные SMP-системы соединяют процессоры с общей памятью по-разному. Самый простой и дешевый подход — это соединение по **общей шине** (system bus/UMA). В этом случае только один процессор может обращаться к памяти в каждый данный момент, что накладывает существенное ограничение на количество процессоров, поддерживаемых в таких системах. Чем больше процессоров, тем больше нагрузка на общую шину, тем дольше должен ждать каждый процессор, пока освободится шина, чтобы обратиться к памяти. Снижение общей производительности такой системы с ростом количества процессоров происходит очень быстро, поэтому обычно в таких системах количество процессоров не превышает 2-4. Примером SMP-машин с таким способом соединения процессоров являются любые многопроцессорные серверы начального уровня.



Второй способ соединения процессоров — через **коммутируемое соединение** (crossbar switch/NUMA). При таком соединении вся общая память делится на банки памяти, каждый банк памяти имеет свою собственную шину, и процессоры соединены со всеми шинами, имея доступ по ним к любому из банков памяти. Такое соединение схемотехнически более сложное, но оно позволяет процессорам обращаться к общей памяти одновременно. Это позволяет увеличить количество процессоров в системе до 8-16 без заметного снижения общей производительности. Примером таких SMP-машин являются многопроцессорные рабочие станции RS/6000.



12. Виртуализация. Типы виртуализации.

Виртуализация — это создание виртуальной (логической) версии каких-то ресурсов на одной вычислительной машине . Т.е. это некий интерпретатор, работающий под управлением другой программы.

- **Виртуальные машины (интерпретаторы)**
 - Java VM, JavaScript в браузере, Python
- **Контейнеры приложений**
 - Docker, Solaris containers, Linux Containers (LXC), ...
- **Аппаратная виртуализация**
 - KVM, Hyper-V, VMWare, Virtual Box, ...
 - Виртуализация аппаратных устройств
- **Облачные технологии**
 - Построены на базе аппаратной виртуализации
 - Дополнительно включают provisioning и общий мониторинг

В целом, виртуальная машина – совокупность ресурсов, которые симулируют поведение реальной машины. При этом процессы идущие на host-машине и гостевой платформе изолированы.

Контейнеры приложений - поменять корень и чтобы все оттуда выполнялось, можно положить только нужные либы и ничего лишнего. (Суть изоляции приложения заключается в том, чтобы оно «не знало» о других программах, работающих в той же операционной системе одновременно с ним, чтобы оно «считало» себя единственным работающим.) Пример: запуск в песочнице

Для каждого приложения во время его запуска создается, так называемая, среда выполнения.

Грубо говоря, виртуализация - эмуляция аппаратного окружения (несколько ОС на одной машине), контейнеризация - выделение изолированного окружения (контейнера) в рамках одной ОС. Все контейнеры используют одно ядро ОС, в виртуализации у каждого окружения свое ядро.

Облачные технологии позволяют использовать с устройства ресурсы удаленных устройств.

13. Сбои и отказоустойчивость ОС. Причины появления отказов в ОС и способы борьбы с ними..



Отказоустойчивость

- Способность системы продолжать работу при аппаратных или программных ошибках
 - Избыточность аппаратуры (двойное, тройное резервирование)
 - Аппаратная «горячая» замена компонентов (диски, контроллеры, процессоры, системные платы)
 - Программная поддержка ОС выведения компонентов из системы и их подключения
 - Организация уровней хранения RAID (Redundant Array of Inexpensive Disks) в дисковой подсистеме

Двойное резервирование: два устройства вместо одного (например, 2 драйвера). По дефолту можно размазать работу на них обоих, а если 1 выйдет из строя - ничего страшного, продолжаем работать: медленнее, но все же - система не рухнет.

Тройное резервирование: смотрится сигнал с всех 3х и берется “большинство” голосов.

Сложно ли менять в горячую? Надо уметь отключить банку (изолировать пассивно и активно), и затем ее заменить.

При наличии множества дисков различные запросы ввода-вывода могут обрабатываться параллельно, если блок данных, к которому производится обращение, распределен по множеству дисков. Кроме того, даже единственный запрос ввода-вывода может быть выполнен параллельно, если блок данных, к которому осуществляется доступ, распределен по нескольким дискам.

В случае применения нескольких дисков имеется большое количество вариантов организации данных и добавления избыточности для повышения надежности (что может

RAID - см. далее. Это дисковые массивы: дорого, зато эффективно(избыточный массив независимых дисков, рассматривается операционной системой как единый логический диск). Избыточность необходима для возможности восстановления данных в случае отказа одного из дисков ([подробнее про уровни](#)).

Таблица 11.4. Уровни RAID

Категория	Уровень	Описание	Требуемое количество дисков ¹	Доступность данных	Пропускная способность передачи больших данных	Скорость запросов малого ввода-вывода
Расщепление	0	Без избыточности	N	Меньше, чем у одного диска	Очень высокая	Очень высокая для чтения и записи
Отражение ²	1	Отражение	$2N$	Больше, чем у RAID 2, 3, 4 и 5, но меньше, чем у RAID 6	Для чтения больше, чем у одного диска; для записи сравнимо с одним диском	Для чтения почти вдвое больше, чем у одного диска; для записи сравнимо с одним диском
Параллельный доступ	2	Избыточность с кодами Хэмминга	$N+m$	Гораздо выше, чем у одного диска, сравнимо с RAID 3, 4 и 5	Наибольшая среди всех перечисленных альтернатив	Почти вдвое больше, чем у одного диска
	3	Четность с чередующимися битами	$N+1$	Гораздо выше, чем у одного диска, сравнимо с RAID 2, 4 и 5	Наибольшая среди всех перечисленных альтернатив	Почти вдвое больше, чем у одного диска
Независимый доступ	4	Четность с чередующимися блоками	$N+1$	Гораздо выше, чем у одного диска, сравнимо с RAID 2, 3 и 5	Для чтения аналогична RAID 0; для записи значительно меньше, чем у одного диска	Для чтения аналогична RAID 0; для записи значительно меньше, чем у одного диска
	5	Распределенная четность с чередующимися блоками	$N+1$	Гораздо выше, чем у одного диска, сравнимо с RAID 2, 3 и 4	Для чтения аналогична RAID 0; для записи меньше, чем у одного диска	Для чтения аналогична RAID 0; для записи в общем случае меньше, чем у одного диска
	6	Двойная распределенная четность с чередующимися блоками	$N+2$	Наибольшая среди всех перечисленных альтернатив	Для чтения аналогична RAID 0; для записи меньше, чем у RAID 5	Для чтения аналогична RAID 0; для записи значительно меньше, чем у RAID 5

====

ВСЕ эти моменты должны быть решены самой ОС.

Причины отказов:

- Ошибочное состояние аппаратуры или ПО в результате сбоя компонентов (может и из-за лунных волн, а может постоянно, тогда уже программа в аппаратуре)
- Ошибки оператора (ввел что-то не так:)
- Физические помехи окружающей среды (кабель влияет на кабель)
- Ошибки проектирования, программирования, структур данных и пр.

Могут быть: постоянные, временные (однократные или периодические)

14. Надежность. Среднее время восстановления.

Коэффициент доступности и время простоя.

Надежность (Reliability)

- $R(t)$ - Вероятность бесперебойной работы системы до времени t , при условии ее корректной работы в $t=0$
- Бесперебойная работа — корректная работа и защита данных
- Среднее время наработки на отказ (Mean Time To Failure)
$$MTTF = \int_0^x R(t)dt$$

Операционные системы. Часть 1.
All rights reserved. Distribution is strictly prohibited unless you have written permission © Yate-Ltd LTD 1999-2008

Среднее время восстановления (Mean Time To Recover)

- Обычно время для перезагрузки, ремонта или замены неисправного компонента, установки (или переустановки) ОС и ПО

График работы и сбоя (downtime). Ось времени t . Время работы обозначено как U , время сбоя как D . На графике изображены три периода работы: U_1 , U_2 , U_3 . Периоды сбоя: D_1 (между U_1 и U_2), D_2 (между U_2 и U_3). Время сбоя $D = D_1 + D_2$.

$$MTTF = \frac{U_1 + U_2 + U_3}{3}$$
$$MTTR = \frac{\text{Boot time} + \text{Reboot time} + \text{Repair time}}{3}$$

Операционные системы. Часть 1.
All rights reserved. Distribution is strictly prohibited unless you have written permission © Yate-Ltd LTD 1999-2008

Коэффициент доступности (Availability)

- Доля времени (%), когда система или служба доступна для запросов пользователей
- Простой (downtime) — время, в течении которого система недоступна
- Безотказная работа (uptime) — время, когда она находится в продуктивной работе

$$\text{Availability} = \frac{MTTF}{MTTF + MTTR}$$

Надежность (reliability) $R(t)$ системы определяется как вероятность ее **бессбойной** работы до времени t при условии ее корректной работы в момент времени $t=0$. Для операционных систем и компьютеров термин **бессбойная работа** означает правильное выполнение набора программ и защиту данных от случайного изменения. **Среднее время наработки на отказ** (mean time to failure — $MTTF$) определяется как

$$MTTF = \int_0^x R(t)dt$$

MTTR – среднее время восстановления

Доступность = вероятность, что система находится в рабочем состоянии

Ну тут вроде все понятно.

Есть еще классы доступности систем (процент коэффициента доступности), но это +
условная штука.

15. Резервирование и отказоустойчивость.

Методы резервирования

- Физическая избыточность (компонентов, серверов) - использовать два прибора/контроллера для одного и того же.
- Временная избыточность (повтор вычислений)
- Информационная избыточность (ECC, RAID - контрольные суммы)

Отказоустойчивость означает способность системы или компонента к продолжению нормальной работы, несмотря на наличие ошибок аппаратного или программного обеспечения. Обычно отказоустойчивость предполагает определенную степень избыточности. Отказоустойчивость предназначена для повышения степени надежности системы. Как правило, увеличение отказоустойчивости (и соответственно, повышение надежности) имеет определенную стоимость, либо финансовую, либо выражющуюся в падении производительности (либо и то, и другое одновременно). Таким образом, определение желаемой степени отказоустойчивости должно учитывать, что именно является критическим ресурсом.

Методы повышения отказоустойчивости ОС

- Изоляция процессов (сейчас с виртуальным пространством очень удобно, но в ядре все еще одно пространство).
- Разрешение блокировок при параллелизме.
- Виртуализация - полностью изолировать выполнение (свои ядра).
- Точки восстановления (копии ключевых файлов) и откаты.

16. История и развитие ОС GNU/Linux. Single UNIX Specification и POSIX.

История

Система Linux возникла как вариант операционной системы UNIX, предназначенный для персональных компьютеров с архитектурой IBM PC (Intel 80386). Первоначальная версия была написана Линусом Торвальдсом (Linus Torvalds), финским студентом, изучавшим теорию вычислительных машин. В 1991 году Торвальдс представил в Интернете первую версию системы Linux. С тех пор множество людей, сотрудничая посредством Интернета, развивают Linux под общим руководством ее создателя. Благодаря тому что система Linux является бесплатной и можно беспрепятственно получить ее исходный код, она стала первой альтернативой для рабочих станций UNIX, предлагавшихся фирмами Sun Microsystems и IBM. На сегодняшний день Linux является полнофункциональной системой семейства UNIX, способной работать почти на всех платформах.

2.10. Linux 135

Залогом успеха Linux является то, что она бесплатно распространяется при поддержке Фонда бесплатно распространяемых программ (Free Software Foundation — FSF). Целью этой организации является создание надежного аппаратно-независимого программного обеспечения, которое было бы бесплатным, обладало высоким качеством и пользовалось широкой популярностью среди пользователей. Проект GNU² фонда предоставляет инструменты для разработки программного обеспечения под эгидой общедоступной лицензии GNU (GNU Public License — GPL). Таким образом, система Linux в таком виде, в котором она существует сегодня, является продуктом, появившимся в результате усилий Торвальдса, а затем и многих других его единомышленников во всем мире, и распространяющимся в рамках проекта GNU.

Linux используется не только многими отдельными программистами; она проникла и в корпоративную среду. В основном это произошло благодаря высокому качеству ядра операционной системы Linux, а не из-за того, что эта система является бесплатной. В эту популярную версию внесли свой вклад многие талантливые программисты, в результате чего появился впечатляющий технический продукт. К достоинствам Linux можно отнести то, что она является модульной и легко настраиваемой. Благодаря этому можно достичь высокой производительности ее работы на самых разнообразных аппаратных платформах. К тому же, получая в свое распоряжение исходный код, производители программного обеспечения могут улучшать качество приложений и служебных программ, с тем чтобы они удовлетворяли конкретным требованиям их пользователей. Имеются также коммерческие компании, такие как Red Hat и Canonical, которые обеспечивают высокопрофессиональную и надежную поддержку своих дистрибутивов Linux. В этой книге подробности внутреннего устройства ядра Linux излагаются на основе ядра Linux 4.7, выпущенного в 2016 году.

Большая часть успеха операционной системы Linux связана с используемой ею моделью развития. Разработчики пользуются единым списком рассылки под названием "LKML" (Linux Kernel Mailing List — список рассылки ядра Linux). Кроме того, имеется множество других списков рассылки, каждый из которых посвящен той или иной подсистеме ядра Linux (список рассылки netdev для сети, linux-pci — для подсистемы PCI, linux-acpi — для подсистемы ACPI и др.). Обновления, отправляемые в эти списки рассылки, должны соответствовать строгим правилам (главным образом — соглашениям о кодировании ядра Linux) и изучаются разработчиками со всего мира, подписанными на эти списки рассылки. Любой пользователь может отправлять свои обновления в эти списки рассылки. Статистика (например, время от времени публикуемая на сайте lwn.net) показывает, что многие обновления предлагаются разработчиками из известных коммерческих компаний, таких как Intel, Red Hat, Google, Samsung и др. Кроме того, многие разработчики являются сотрудниками коммерческих компаний (как Дэвид

того, многие разработчики являются сотрудниками коммерческих компаний (как Дэвид Миллер (David Miller), поддерживающий сетевые функции и работающий в компании Red Hat). Такие обновления изучаются и обсуждаются в списке рассылки, после чего в них вносятся исправления, и цикл обсуждения начинается заново. В конце концов принимается решение, следует ли принять или отклонить эти исправления. Каждый руководитель подсистемы время от времени отправляет запрос о размещении исправлений его части в основном ядре, который обрабатывается Линусом Торвальдсом. Сам Линус

² GNU — рекурсивная аббревиатура для GNU's Not Unix (GNU — не Unix). Проект GNU представляет собой бесплатный набор программных пакетов и инструментов для разработки UNIX-подобной операционной системы.

136 ГЛАВА 2. Обзор операционных систем

выпускает новую версию ядра примерно каждые 7–10 недель, причем каждый такой выпуск имеет около 5–8 предварительных версий-кандидатов.

Интересно попытаться понять, почему другие операционные системы с открытым кодом, такие как различные версии BSD или OpenSolaris, не имеют таких успеха и популярности, которыми обладает Linux. Тому может быть много причин; конечно, открытость модели развития Linux способствовала популярности и успеху этой операционной системы. Но эта тема выходит за рамки нашей книги.



Single UNIX Specification. POSIX

- Unix SUS: The Open Group, Austin Group
 - Основные определения
 - Системные интерфейсы
 - Командная оболочка и утилиты
 - Пояснения
 - X/Open Curses
 - UNIX: AIX, HP-UX, IRIX, Mac OS X, SCO OpenServer, Solaris, Tru64 и z/OS.
 - UNIX-like: FreeBSD, OpenBSD, NetBSD, OpenSolaris, BeleniX, Nexenta OS) и Linux
- POSIX (Portable Operating System Interface) ISO/IEC 9945

SUS – спецификация, предъявляющая требования к UNIX системам. При помощи данного документа можно определить, относится ли ОС к Unix
POSIX - стандарт отвечает за интерфейс ОС UNIX
(ISA - instructions set architecture, ABI - application binary interface, API - application program interface).

UNIX-like – операционки, не удовлетворяющие стандартам

17. Понятие дистрибутива, дистрибутивы Linux.

Дистрибутив - это ядро + набор софта для прикладных задач, этот набор в каждом дистрибутиве может быть свой. В этот набор входит “системный” софт: пакетный

менеджер, файловый менеджер, графическая оболочка и набор прикладных программ (текстовый редактор, текстовый процессор, браузер, аудио/видеопроигрыватели и т.п.) (но при этом далеко не все это обязательно входит, например, графическая оболочка далеко не везде используется, а некоторые дистрибутивы представляют из себя почти чистое ядро, например Arch)

Часть дистрибутивов создается компаниями(RedHat, Ubuntu, Suse),: хотя это и опен-сурс проекты, компании продают дистрибутивы, предлагая поддержку (обновления, траблшутинг)

Многие дистрибутивы поддерживаются сообществом, которое использует их и заинтересовано в их развитии (Arch, Debian и их многочисленные потомки)

Поскольку код открытый, многие берут старый дистрибутив за основу и что-то меняют там, создавая новый - этим объясняется огромное дерево дистрибутивов Линукса, но большинство из них наследники Debian и RedHat. Android использует ядро Linux, но не GNU, поэтому обычно не причисляется к дистрибутивам Linux.

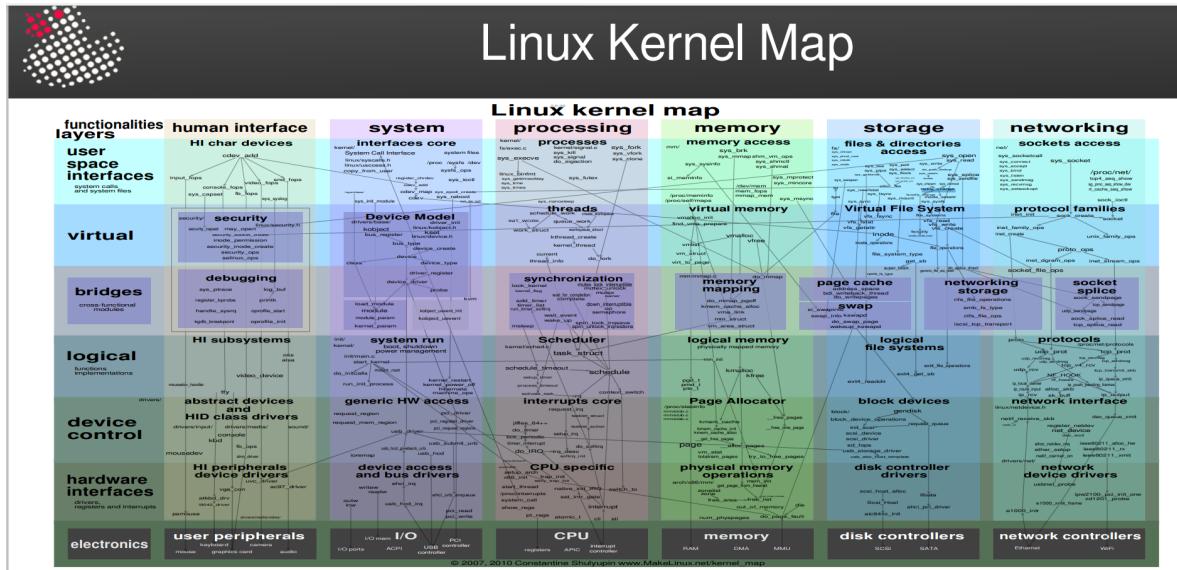
В Linux ядро обычно берется от Линуса Торвальдса(какой-либо версии) и включает в себя окружение(библиотеки, утилиты). Дистрибутивы отличаются различными менеджерами пакетов и обновлений. Кроме того есть графическая подсистема, есть какие-то прикладные программы, которые работают/не работают :) на Linux.
Поддержку дистрибутива можно купить за денежку.

18. Архитектура и основные подсистемы Linux. Linux Kernel Map.

Архитектура Linux



Рис 1 Компоненты Linux



<https://makelinux.github.io/kernel/map/>

Основные подсистемы Unix/Linux:

- Процессы и планировщик. – Создает, управляет и планирует процессы.
- Виртуальная память. – Выделяет виртуальную память для процессов и управляет ею.
- Физическая память. – Управляет пулом кадров страниц и выделяет страницы для виртуальной памяти.

(Физическая и виртуальная память связаны через MMU manager memory unit, много вкручено функций на уровне ядра, управляющих этим взаимодействием).

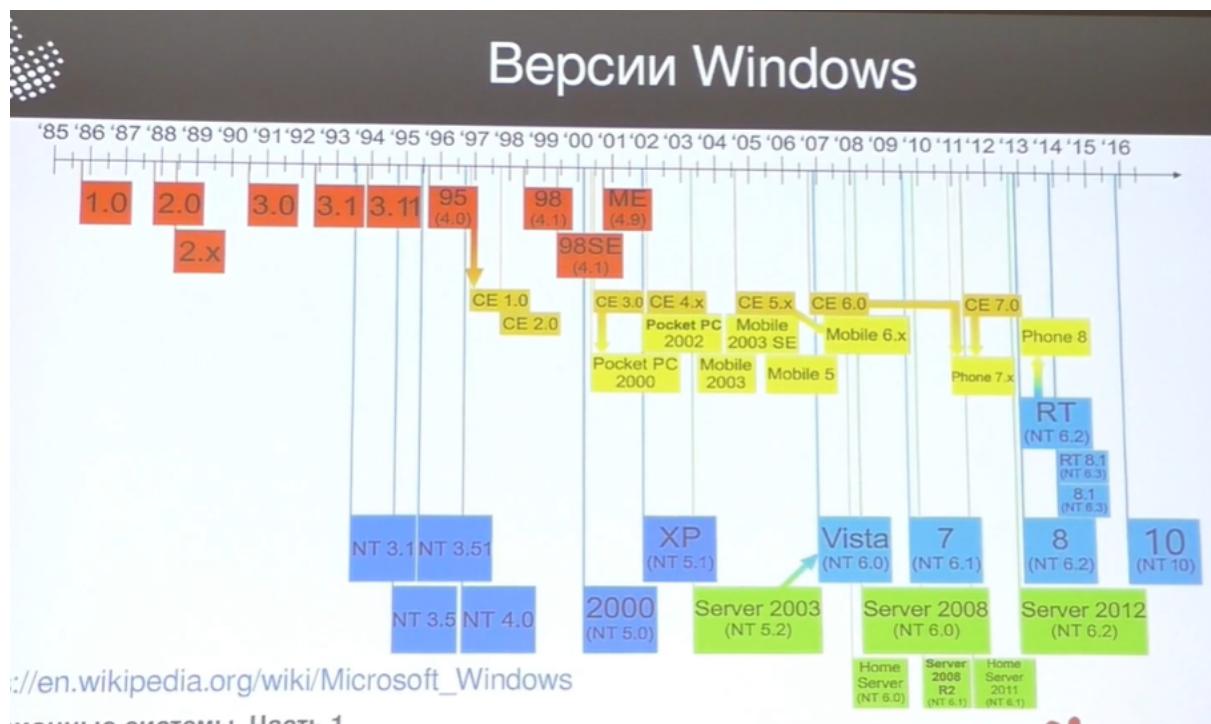
- Файловая система (бывают физические и виртуальные) – Предоставляет глобальное иерархическое пространство имен для файлов, каталогов и других объектов, связанных с файлами и функциями файловой системы.
В Unix все есть файл, поэтому все, с чем работаем - файловая система. Любая система обращается к драйверу vfs (virtual file system).
- Драйверы символьных устройств. (непосредственно передают данные на устройство) – Управление устройствами, которые требуют от ядра отправки или получения данных по одному байту, например терминалами, принтерами или модемами.
(ls -l: chain-oriented)
- Драйверы блочных устройств. (осуществляют буферизацию сначала) – Управление устройствами, которые читают и записывают данные блоками, как, например, различные виды вторичной памяти (магнитные диски, CD-ROM и т.п.).
(ls -l: block-oriented)
- Сетевые протоколы. TCP/IP. – Поддержка пользовательского интерфейса сокетов для набора протоколов (еще, например x25 - в банковской сфере). Их оч. много.
- Драйверы сетевых устройств. – Управление картами сетевых интерфейсов и коммуникационными портами, которые подключаются к сетевым устройствам, такими как мосты или роутеры.

- Ловушки и отказы. – Обработка генерируемых процессором прерываний, как, например, при сбое памяти. (+ реакция: паника, подгрузка страницек...)
- Прерывания. – Обработка прерываний от периферийных устройств.
- Сигналы и IPC – Управляет межпроцессным взаимодействием.

Утилита для посылки сигналов: kill, какие есть сигналы: kill -l.

IPC - средство межпроц. взаимодействия. Включает: разделяемая память, семафоры, месседжи.

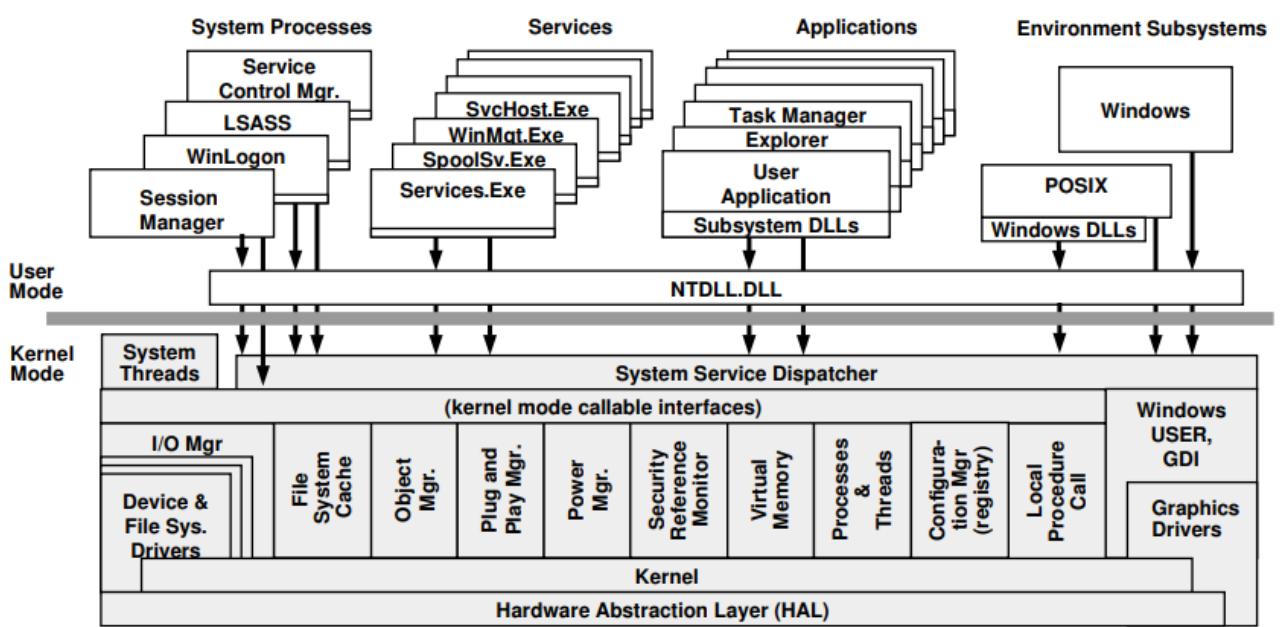
19. История и развитие Windows

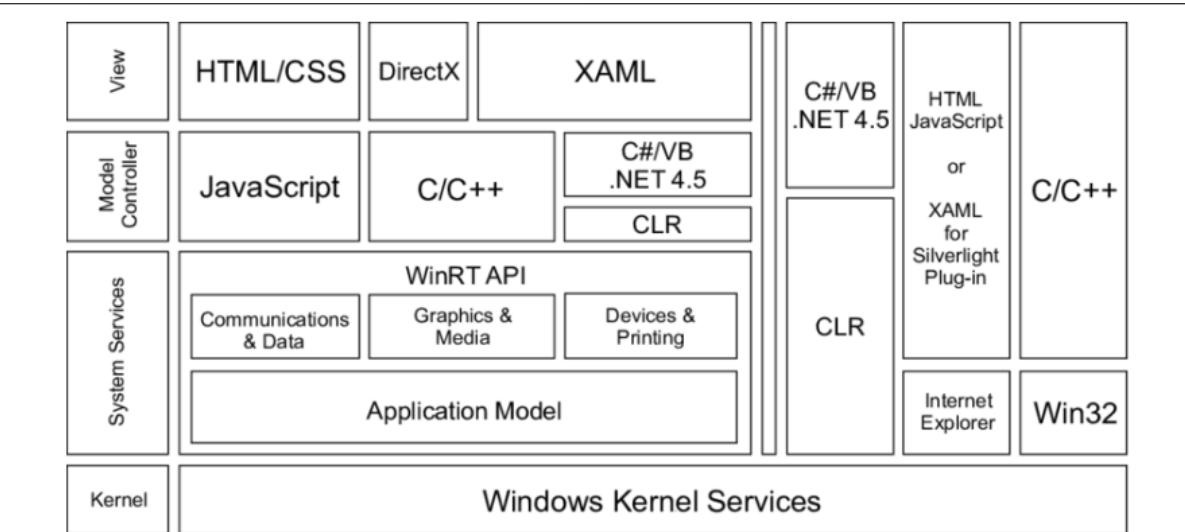


Основы

Впервые Microsoft использовала имя “Windows” в 1985 году для операционной среды, расширяющей возможности примитивной операционной системы MS-DOS, которая успешно использовалась на ранних персональных компьютерах. Такая комбинация Windows/MS-DOS в конечном итоге была заменена новой версией Windows, известной как Windows NT, впервые выпущенной в 1993 году и предназначеннной для ноутбуков и настольных систем. Хотя основная внутренняя архитектура остается примерно той же, что и в Windows NT, сама операционная система продолжает развиваться и дополняться новыми функциями и возможностями. Последняя версия на момент написания этой книги — Windows 10. Windows 10 включает в себя возможности предыдущей операционной системы для настольных и портативных компьютеров Windows 8.1, а также версий Windows, предназначенных для мобильных устройств для Интернета вещей (Internet of Things — IoT). Windows 10 также включает в себя программное обеспечение Xbox One. В результате единая унифицированная операционная система Windows 10 поддерживает настольные компьютеры, ноутбуки, смартфоны, планшеты и Xbox One.

20. Общая архитектура Windows. Windows API





Windows присуще четкое разделение на модули. Каждая функция системы управляет только одним компонентом операционной системы. Остальные ее части и все приложения обращаются к этой функции через стандартный интерфейс. Доступ к основным системным данным можно получить только через определенные функции. В принципе любой модуль можно удалить, обновить или заменить, не переписывая всю систему или стандартный интерфейс прикладного программирования (application program interface -API).

21. Сервисы, функции и важные компоненты Windows.



Сервисы и функции

- Windows API functions:
 - CreateProcess, CreateFile
- System calls (Native System Services)
 - NtCreateUserProcess
- Kernel support functions
 - ExAllocatePoolWithTag
- Windows services
 - Управляются Service Control Manager
- Dynamic link libraries (DLL)
 - msrvct.dll, kernel32.dll

Операционные системы. Часть 1.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Другие важные компоненты системы

- Гипервизор Hyper-V
 - Запуск гостевых ОС, Device Guard, Hyper Guard, Credentials Guard, Application Guard, ...
- Firmware (в том числе и для устройств)
- Terminal Servers
- Объекты и безопасность
- Registry (Реестр)
- Оснастки

Операционные системы. Часть 1.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



22. Процесс, характеристики процесса в момент выполнения. Состояние процесса. Разделение ресурсов.



Итак, процесс это:

- Выполняемая программа
 - Экземпляр программы, выполняющейся на компьютере
 - Сущность, которая может быть назначена процессору и выполнена на нем
 - Единица активности, характеризуемая выполнением последовательности команд, текущим состоянием и связанным с ней множеством системных ресурсов.
-
- **Идентификатор.** Уникальный идентификатор, связанный с этим процессом, чтобы отличать его от всех прочих процессов.
 - **Состояние.** Если процесс выполняется в настоящее время, он находится в **состоянии выполнения**.
 - **Приоритет.** Уровень приоритета по отношению к другим процессам.
 - **Программный счетчик.** Адрес очередной выполняемой команды программы.
 - **Указатели памяти.** Включают указатели на программный код и данные, связанные с этим процессом, а также на любые блоки памяти, совместно используемые с другими процессами.
 - **Данные контекста.** Это данные, присутствующие в регистрах процессора во время выполнения процесса.
 - **Информация о состоянии ввода-вывода.** Включает в себя внешние запросы ввода-вывода, устройства ввода-вывода, назначенные процессу, список файлов, используемых процессом, и т.д.
 - **Учетная информация.** Может включать количество процессорного времени и времени работы, учетные записи и т.д.

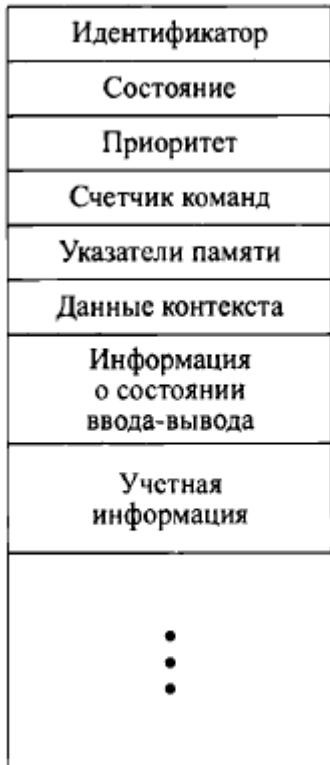
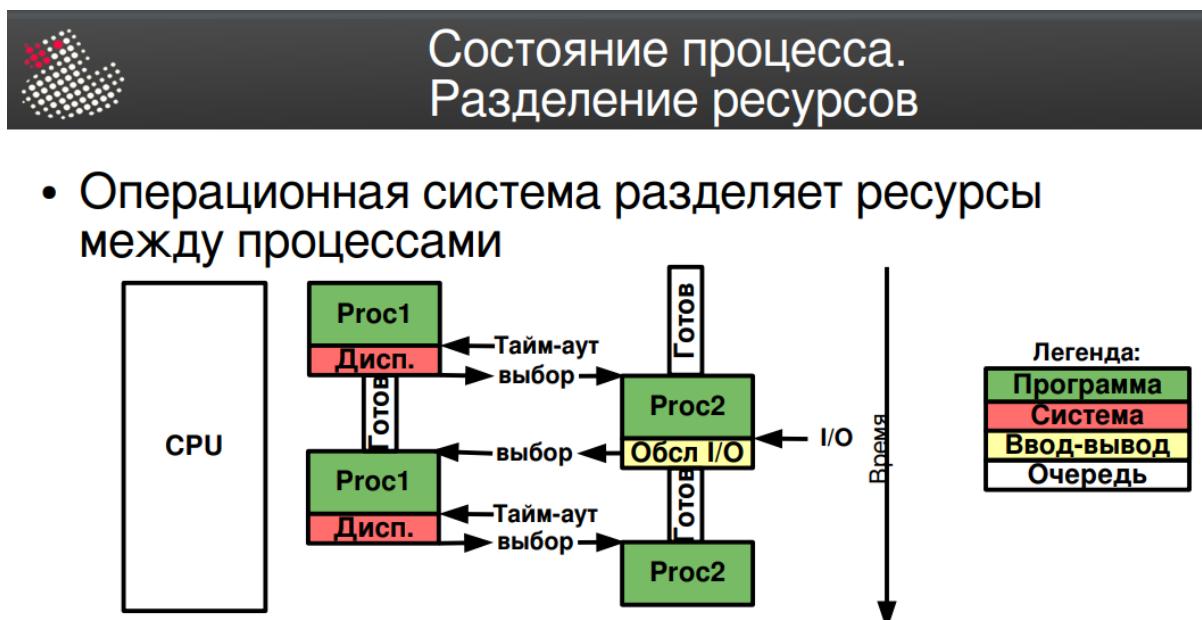


Рис. 3.1. Упрощенный управляющий блок процесса



Эту схему я честно говоря не особо понял, но судя по всему тут имеется ввиду многозадачность, что ОС распределяет процессорное время между разными процессами и у этих процессов разные состояния. 1 процесс вначале выполняется, затем завершается по тайм ауту и становится готовым к выполнению, в то время как 2 процесс начинается выполняться и заканчивает свое выполнение по прерыванию от ввода/вывода, после чего становится блокированным и так далее.

Здесь речь о системе разделения ресурсов: у нас есть очередь процессов и процессам отведено равное время на исполнение, как оно истекает, происходит переход к следующему процессу(в состоянии готовности), который исполнялся наименьшее время. Хотим таким образом достичь “честного” распределения ресурсов и не допустить starvation.

23. Модель процесса с пятью состояниями, назначение состояний.



Рис. 3.6. Модель с пятью состояниями

Опишем каждое из пяти состояний процессов, представленных на диаграмме.

1. **Выполняющийся.** Процесс, который выполняется в текущий момент времени. В настоящей главе предполагается, что на компьютере установлен только один процессор, поэтому в этом состоянии может находиться только один процесс.
2. **Готовый к выполнению.** Процесс, который может быть запущен, как только для этого представится возможность.
3. **Блокированный/Ожидаящий⁵.** Процесс, который не может выполняться до тех пор, пока не произойдет некоторое событие, например завершение операции ввода-вывода.
4. **Новый.** Только что созданный процесс, который еще не помещен операционной системой в пул выполнимых процессов. Обычно это новый процесс, который еще не загружен в основную память, хотя управляющий блок процесса уже создан.
5. **Завершающийся.** Процесс, удаленный операционной системой из пула выполнимых процессов из-за завершения его работы или аварийно прерванный по какой-либо иной причине.

24. Paging и Swapping. Модель процесса с семью состояниями.

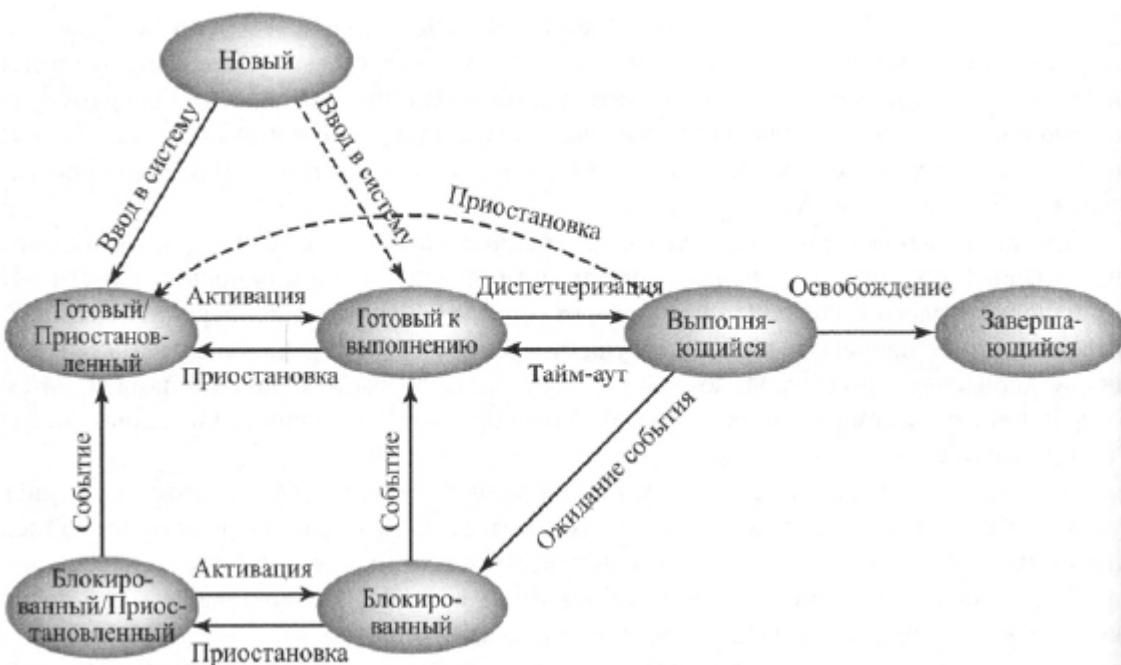


Paging/Swapping

- Основной памяти всегда мало =(
 - Программисты как только видят больше памяти, стараются ее максимально «использовать»
 - Большое количество запущенных процессов
- Давайте поместим блокированный процесс на диск и освободим основную память для других процессов!
 - Нужно организовать область подкачки процессов на диске
- Paging (пейджинг) — выгрузка (и загрузка) неиспользуемых страниц процесса на диск
- Swapping (свопинг) — выгрузка всего процесса, кроме критически важных для ядра структур управления

В принципе, Paging – это тот же Swapping, но осуществляемый «постстранично»: если вдруг оказывается, что какая-то страница не очень нужна, то она выгружается. Для этого в операционных системах существуют специальные алгоритмы, занимающиеся вытеснением из памяти не использующихся страниц. Такой алгоритм сводится к тому, что у страниц, находящихся в памяти «по кругу» постоянно проверяется – есть ли изменения информации на странице. И если изменений на странице нет в течение заданного времени, то страница считается неиспользуемой и специальным образом помечается. А далее, в зависимости от того, что это за страница, принимается решение, что с ней делать. Например, если на странице находится код программы, который всегда можно загрузить из файла заново, то эту страницу можно смело удалить; если же на странице находятся данные, то ее необходимо сохранить в swap, и так далее.

Модель процесса с 7 состояниями



Следует обратить внимание на то, что на схеме появились новые состояния процесса.

Одно из них - это состояние runnable/suspend, при котором процесс может выполнится, но - либо весь процесс, либо его существенная часть находится в области подкачки.

Для простоты дальнейшего описания сейчас будем говорить о Swapping, учитывая то, что при Paging все очень похоже.

Итак, процесс готов на выполнение, у него ввод-вывод закончился. И процесс хочет попасть на процессор. Но ему мешает то, что его часть находится в свопе. А оказалась она там потому, что до этого у процесса было другое состояние wait(suspend). Рассмотрим подробнее.

Вначале процесс находился на CPU. Потом процесс «перешел к вводу-выводу» или «повис» на какой-либо блокировке. В какой-то момент операционная система «решила», что процесс блокировки тянется слишком долго (или ОС решила, что памяти слишком мало), поэтому операционная система решает **освободить** от этого процесса память, потому что он все равно ничего не делает и находится на блокировке. После этого процесс **переходит** в состояние wait(suspend). После этого в какой-то момент «ввод-вывод» «пришел», и процесс необходимо «пробудить». Система «перебрасывает» процесс в очередь ожидания запуска. Но процесс выполнится не может, потому что его часть находится на диске. После этого запускается процедура, загружающая процесс с диска в память. И только после того, как все необходимые страницы загружены, процесс переходит в состояние **unnable** и может быть выполнен на процессоре.

В этой схеме есть еще ряд интересных вещей. На схеме явно видно, что мы разделили состояние wait на две части. И так же разделили на две части состояние runnable. Помимо этого, на схеме появились «стрелочки», которые могут вызывать вопросы. Дело в том, что процесс порождения процесса очень длителен и он занимает сотни миллисекунд (а иногда и больше). И все современные системы, учитывая наличие виртуальной памяти и наличие страницек, поступают следующим образом. Чтобы процесс создать как можно быстрее, он создается практически без создания маппингов памяти. И может получиться так, что при приходе процесса в состояние runnable/desperate memory conditions, в его сегменте кода не будет ни одной загруженной страницы. И поэтому процесс не сможет продолжить свое выполнение. При этом получается, что процесс уже создан, все структуры ядра созданы, но процесс выполнится еще не может. Поэтому необходимо задействовать уже механизмы пейджинга и свопинга для подкачки необходимых для выполнения страниц.

Также у операционной системы в зависимости от количества свободной памяти есть несколько режимов работы. И на одной из вновь появившихся стрелочек происходит **desperate memory conditions** – это состояние системы, когда процессы потребляют столько памяти, что ее начинает не хватать **самой** системе. В этот момент операционная система начинает массово отправлять процессы в своп, чтобы освободить память **для своих** нужд. Таким образом процесс оказывается полностью готов, но у него нет возможности выполнится, потому, что у операционной системы не хватит ресурсов, чтобы его выполнить. И очень часто, когда система очень загружена, случается так, что все процессы встают в «локи» на моменте выделения новой памяти.

25. Управляющие таблицы процесса. Образ процесса

Если смотреть на операционную систему укрупненно, то у нее есть четыре основных подсистемы ядра и некоторые внутренние структуры, которые их описывают.

В частности, существует подсистема с ее собственными внутренними структурами; существуют файлы и структуры файловой системы, обеспечивающие работу с файлами. Аналогичными структурами подсистем можно описать устройства и процессы.

Внутри операционной системы существуют списки процессов. Примерами таких списков могут являться «глобальный список» процессов и «список состояний», в которых процессы могут находиться.

У каждого процесса есть структура, которая отражает образ процесса в памяти машины. В эту структуру обычно входит «ядерная часть» и «часть на стороне пользователя», которая описывает и характеризует процесс, а также содержит данные процесса и некоторые другие характеристики. Таких структур внутри ядра большое количество, и они между собой сильно «перевязаны». При этом каждый разработчик драйвера **создает**, как правило, свои управляющие структуры.

«На стороне ядра» всегда должна находиться информация, которая показывает, что происходит с процессом, в каком он находится состоянии, что ему требуется для исполнения, сколько времени он потратил на CPU...

Если посмотреть на реальную структуру операционной системы (например, в Linux или в Solaris), можно заметить большое количество полей, для понимания которых необходимо **потребовать** достаточно количество времени. Они сложны по названиям, и перемещаться по их перекрестным ссылкам достаточно сложно, так как часто одни структуры «вкладываются» в другие структуры, что в результате приводит к очень объемным текстам Си-шного кода, которые трудны для понимания.

Детального понимания всех этих структур для **каждодневной** работы «обычным программистам», как правило, не требуется. Однако для формирования качественного представления о предмете, имеет смысл посмотреть, как в операционной системе реализованы основные структуры.

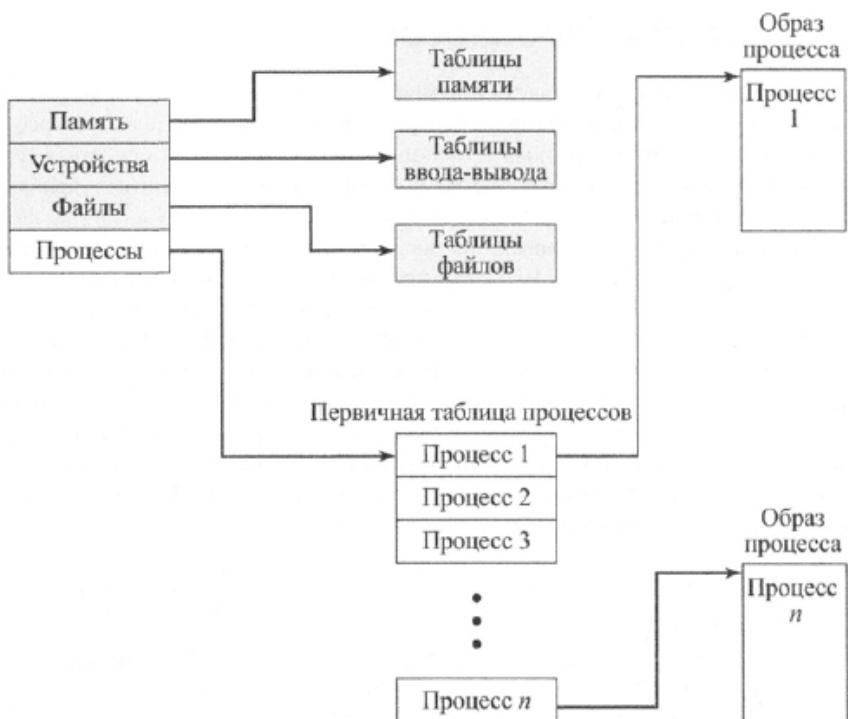
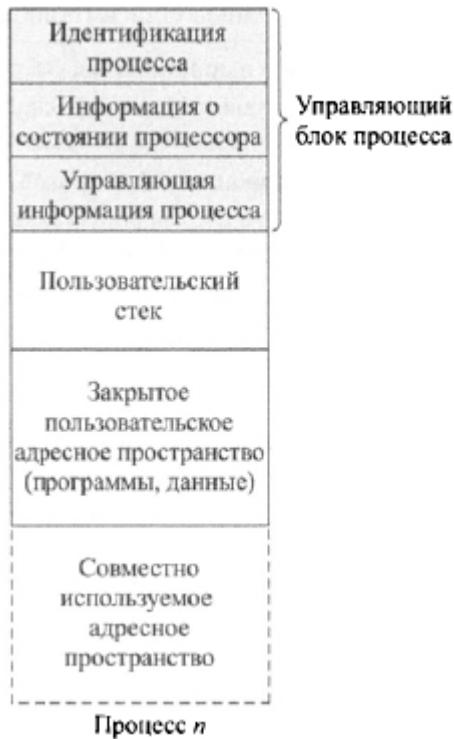


Рис. 3.11. Общая структура управляющих таблиц операционной системы

Таблица 3.4. Типичные элементы образа процесса

Данные пользователя	Допускающая изменения часть пользовательского адресного пространства. Сюда могут входить данные программы, пользовательский стек и модифицируемый код
Пользовательская программа	Программа, которую нужно выполнить
Системный стек	С каждым процессом связан один или несколько системных стеков. Стек используется для хранения параметров, адресов вызова процедур и системных служб
Управляющий блок процесса	Данные, необходимые операционной системе для управления процессом (см. табл. 3.5)



26. Управляющий блок процесса (PCB), состав PCB.

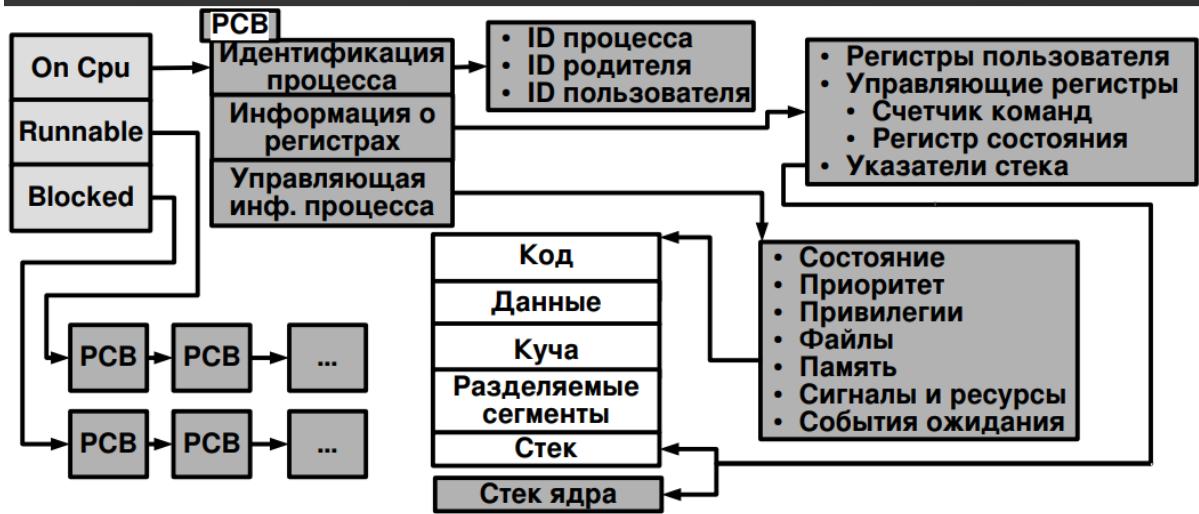
Роль управляющего блока процесса

Управляющий блок процесса — это самая важная структура данных из всех имеющихся в операционной системе. В управляющий блок каждого процесса входит вся необходимая операционной системе информация о нем. Информация в этих блоках считывается и/или модифицируется почти каждым модулем операционной системы, включая те, которые связаны с планированием, распределением ресурсов, обработкой прерываний, а также осуществлением контроля и анализа. Можно сказать, что состояние операционной системы задается совокупностью управляющих блоков процессов.

Ю



Управляющий блок процесса (PCB)



Операционные системы. Часть 2. Процессы и потоки

• • •

Те, как мы видим в PCB содержится вся информация о процессе, его айди, состояние регистров, код, данные и тп.

27. Функции ОС, связанные с процессами. Создание процесса, переключение процессов.



Функции ОС связанные с процессами

- Управление процессами
 - Создание и завершение процессов
 - Планирование и диспетчеризация процессов
 - Переключение процессов
 - Синхронизация и поддержка обмена информацией между процессами
 - Организация управляющих блоков процессов
- Управление памятью
 - Выделение адресного пространства процессам
 - Пейджинг и Свопинг
 - Управление страницами и сегментами
- Управление вводом-выводом
 - Управление буферами
 - Выделение процессам каналов и устройств ввода-вывода
- Функции поддержки
 - Обработка прерываний
 - Учет использования ресурсов
 - Текущий контроль системы



Создание процесса

- Присвоить процессу уникальный идентификатор
- Выделить память для процесса
- Инициализировать PCB
- Поставить процесс в очереди ядра
- Создать потоки ввода-вывода
- Создать другие управляющие структуры данных



Переключение процессов

- Процесс может работать в user и kernel mode
- Используется механизм прерываний:
 - Внешнее прерывание (IO)
 - Ловушка — обработка ошибки или исключительной ситуации
 - Вызов ОС
- Переход из user в kernel-mode ситуации:
 - Прерывания таймера
 - Прерывания ввода-вывода
 - Page fault — отсутствие блока памяти

Таблица 3.8. МЕХАНИЗМЫ ПРЕРЫВАНИЯ ВЫПОЛНЕНИЯ ПРОЦЕССА

Механизм	Причина	Что используется
Прерывание	Внешняя по отношению к выполнению текущей команды	Отклик на внешнее асинхронное событие
Ловушка	Связана с выполнением текущей команды	Обработку ошибки или исключительной ситуации
Вызов супервизора	Явный запрос	Вызов функции операционной системы

- **Прерывание таймера.** Операционная система определяет, что текущий процесс выполняется в течение максимально разрешенного промежутка времени, именуемого **квантом времени** (time slice). Квант времени представляет собой максимальное количество времени, которое процесс может выполнять без прерывания. Если это так, то данный процесс нужно переключить в состояние готовности и передать управление другому процессу.
- **Прерывание ввода-вывода.** Операционная система определяет, что именно произошло, и если это то событие, которого ожидают один или несколько процессов, операционная система переводит все соответствующие блокированные процессы в состояние готовности (соответственно, блокированные/приостановленные процессы она переводит в состояние готовых/приостановленных процессов). Затем операционная система должна принять решение: возобновить выполнение текущего процесса или передать управление готовому к выполнению процессу с более высоким приоритетом.
- **Ошибка отсутствия блока в памяти.** Допустим, что процессор должен обратиться к слову виртуальной памяти, которое в настоящий момент отсутствует в основной памяти. При этом операционная система должна загрузить в основную память блок (страницу или сегмент), в котором содержится адресованное слово. Сразу же после запроса на загрузку блока операционная система может передать управление другому процессу, а процесс, для продолжения выполнения которого нужно загрузить блок в основную память, переходит в блокированное состояние. После загрузки нужного блока этот процесс переходит в состояние готовности.

28. Процессы в ОС UNIX SVR4. Диаграмма состояний, основные структуры.

Состояния процессов

Всего в операционной системе UNIX SVR4 распознается девять состояний процессов, перечисленных в табл. 3.9; соответствующая диаграмма переходов состояний показана на рис. 3.17 (в ее основе — рисунок из [14]). Этот рисунок похож на рис. 3.9,б; нужно только принять во внимание, что два спящих состояния в системе UNIX соответствуют двум блокированным состояниям. Кратко перечислим основные различия между диаграммами.

Таблица 3.9. Состояния процессов в UNIX

Выполняющийся пользовательский	Выполняющийся в пользовательском режиме
Выполняющийся ядра	Выполняющийся в режиме ядра
Готовый к выполнению, загруженный	Готов к выполнению, как только ядро решит передать ему управление
Спящий, загруженный	Не может выполняться, пока не произойдет некоторое событие; процесс находится в основной памяти (блокированное состояние)
Готовый к выполнению, выгруженный	Процесс готов к выполнению, но прежде чем ядро сможет спланировать его запуск, процесс свопинга должен загрузить этот процесс в основную память
Спящий, выгруженный	Процесс ожидает некоторого события; он выгружен из основной памяти (блокированное состояние)
Вытесненный	В момент переключения процессора из режима ядра в пользовательский режим ядро решает передать управление другому процессу
Созданный	Процесс только что создан и еще не готов к выполнению
Зомби	Самого процесса больше не существует, но записи о нем остались с тем, чтобы ими мог воспользоваться родительский процесс

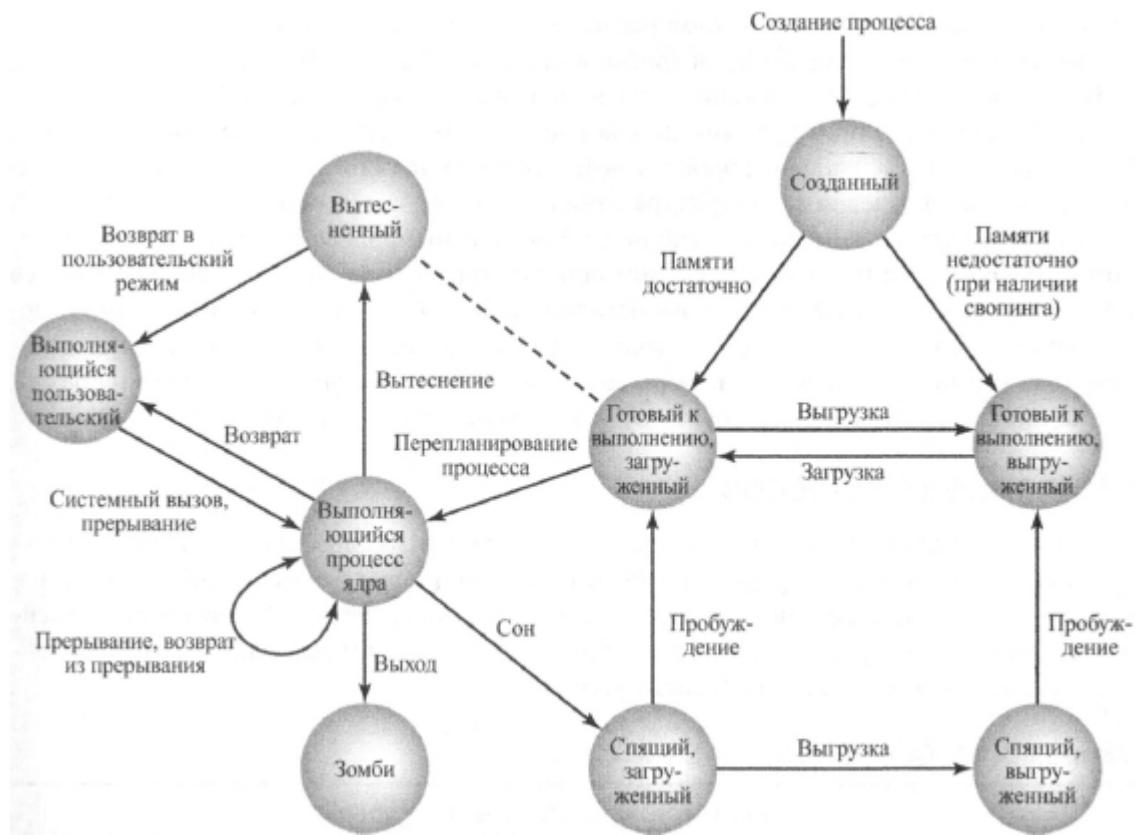
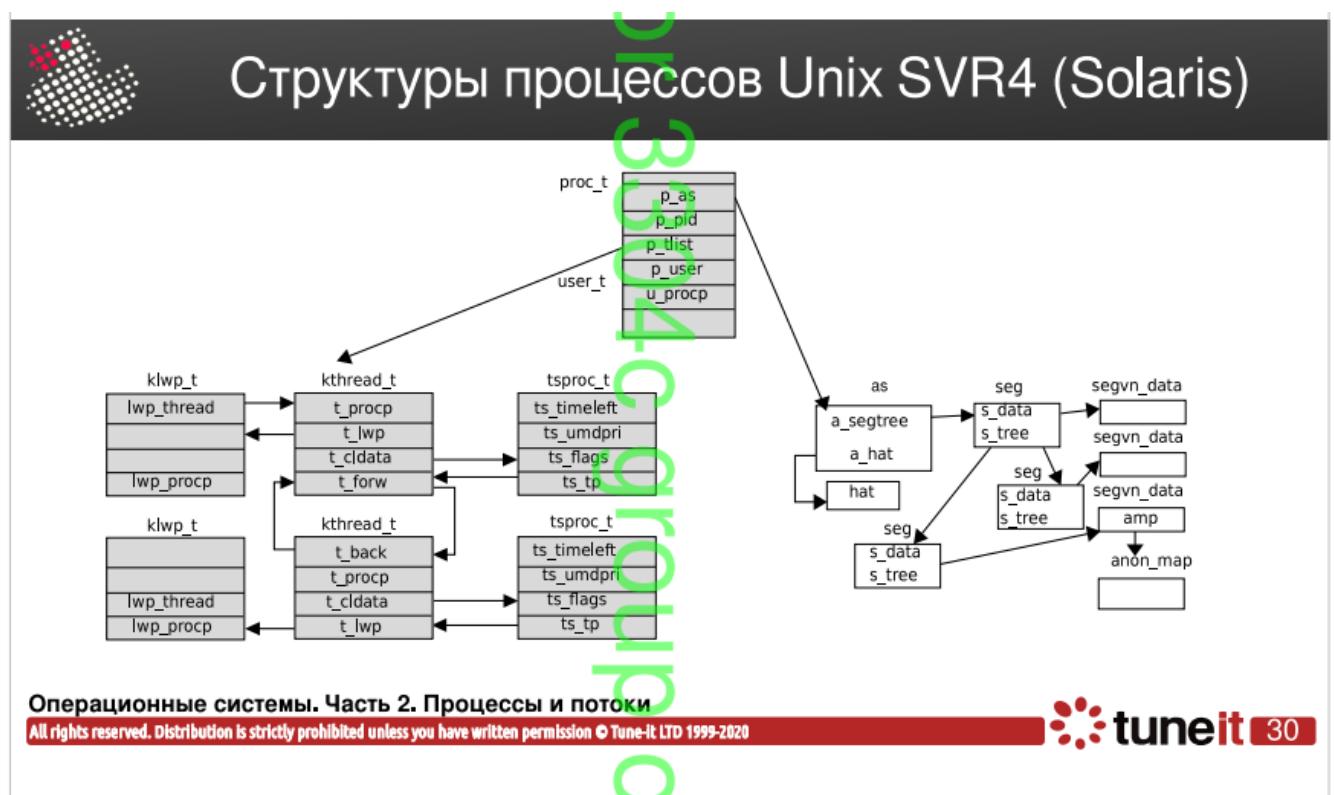


Рис. 3.17. Диаграмма переходов между состояниями процессов в системе UNIX



На слайде представлена небольшая схема, демонстрирующая то, как связаны между собой структуры PCB в Unix SVR4 (Solaris).

В Solaris PCB структурно состоит из двух частей: proc_t и user_t. Эти структуры как бы вставлены друг в друга. proc_t и user_t (т.е. PCB) занимают суммарно порядка 2Кб памяти.

p_as – это указатель на адресацию памяти.

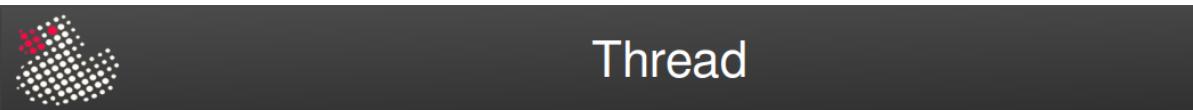
На схеме также показаны структуры, относящиеся к адресному пространству; структуры, относящиеся к конкретному сегменту адресного пространства; а также структуры, которые говорят о том именованное это адресное пространство, или анонимное. Признак именованности/неименованности означает, что адресное пространство связано с файлом или с ним не связано. Например: именованное пространство – это сегмент данных в файле, анонимное пространство – это, например «куча» или стек.

Внутри процесса есть такая абстракция как трэды. На схеме представлены трэды на уровне ядра и трэды на уровне пользователя. Обе эти категории трэдов отдельно диспетчеризуются и у них существуют специальные приоритеты. Все это связано между собой в некоторое количество структур.

На схеме приведены далеко не все существующие поля, так как различных полей, описывающих структуры процесса более сотни. Причем, чем глубже происходит изучение структуры, тем больше она начинает «разрастаться». Это объясняется тем, что ядро – это достаточно сложный объект для изучения.

Самое «неприятное» здесь в том, что большое количество переменных от версии к версии меняется. И если на уровне пользователя (через POSIX) операционная система представляется единым целым, то в действительности, всё содержание от версии к версии «дышит». Поэтому перед тем, как начинать заниматься отладкой необходимо внимательно загрузить именно необходимые версии всех используемых компонентов.

29. Понятие потока выполнения, связь потока и процесса. Преимущества потоков.



- Абстрактная модель процесса подразумевает:
 - Владение ресурсами (сегменты памяти, файлы, каналы ввода вывода, контекст безопасности, ...)
 - Планирование и диспетчеризацию (приоритеты, состояния) независимы
- Процесс — единица группировки общих ресурсов
- Thread (нить выполнения) — единица выполнения программного кода.
Содержит:
 - Состояние выполнения
 - Сохраненный контекст потока (регистры, ...)
 - Стек (ядра и пользователя)
 - Локальные переменные (thread_locals)
 - Доступ к памяти и другим ресурсам процесса-владельца

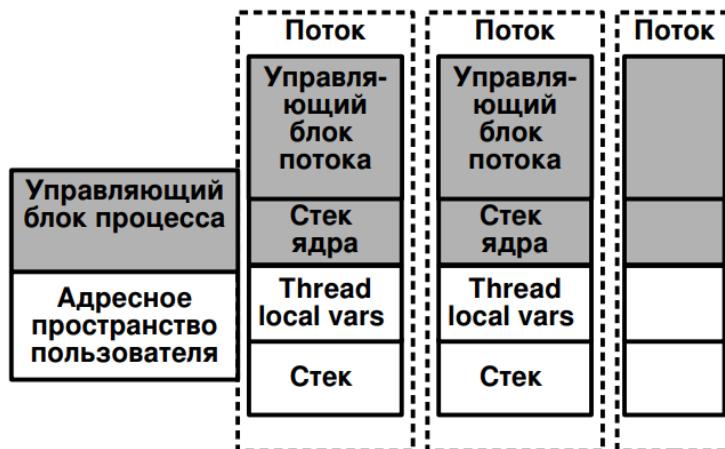


Связь структур ядра процесса и потока

Однопоточная модель



Многопоточная модель



Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



Преимущество потоков

- Общее преимущества в быстродействии:
 - Потоки создаются на порядок быстрее
 - Потоки переключаются быстрее
 - Потоки завершаются быстрее
 - Потоки могут обмениваться информацией быстрее (без участия ядра)
- Даже в однопроцессорной системе есть преимущества
 - Работа в приоритетном и фоновом режиме
 - Асинхронная обработка частей программы
 - Модульная структура программы

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

30. Состояния потока, User Level Threads vs Kernel Level Threads

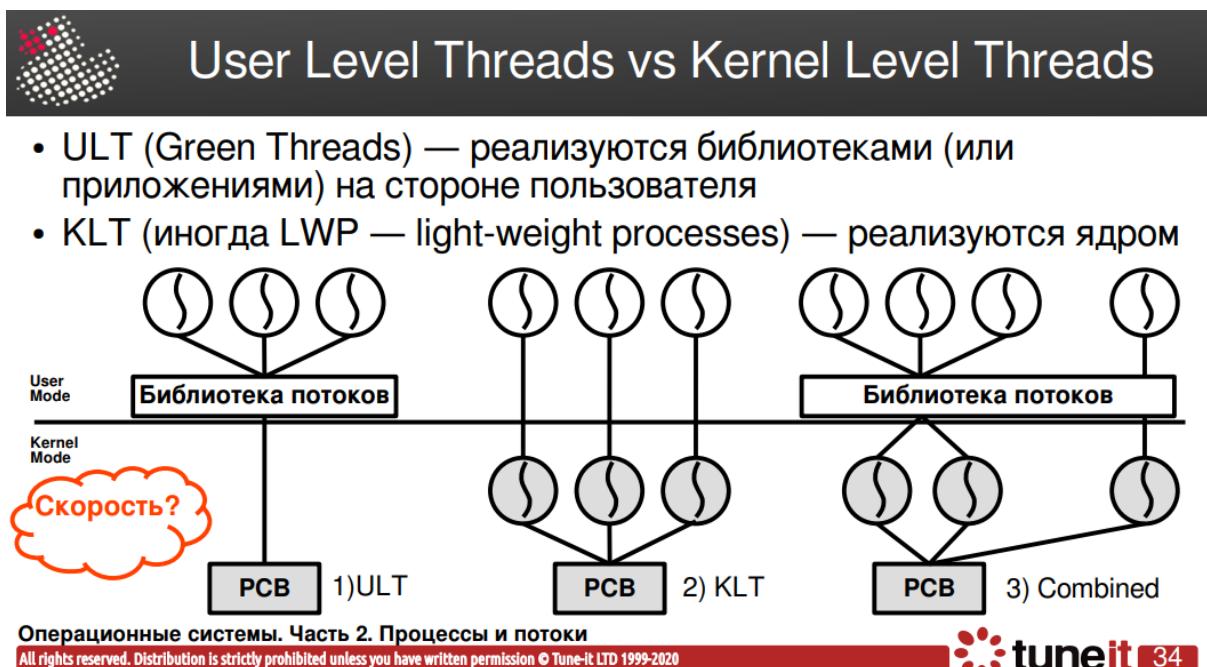


ТАБЛИЦА 4.1. ВРЕМЯ ЗАДЕРЖЕК ПОТОКОВ (ms)

Операция	Пользовательские потоки	Потоки на уровне ядра	Процессы
Нулевое ветвление	34	948	11300
Ожидание сигнала	37	441	1840

ТАБЛИЦА 4.2. Соотношение между потоками и процессами

Потоки: процессы	Описание	Примеры систем
1:1	Каждый поток реализован в виде отдельного процесса с собственным адресным пространством и со своими ресурсами	Традиционные реализации системы UNIX
M:1	Для процесса задаются адресное пространство и динамическое владение ресурсами. В рамках этого процесса может быть создано несколько потоков	OS/2, Windows NT, Solaris, Linux, OS/390, MACH
1:M	Поток может переходить из среды одного процесса в среду другого процесса. Это облегчает перенос потоков из одной системы в другую	Ra (Clouds), Emerald
M:N	Сочетает в себе подходы, основанные на соотношениях M:1 и 1:M	TRIX

31. Многопроцессорность и многопоточность. Закон Амдала.



Многопроцессорность и многопоточность

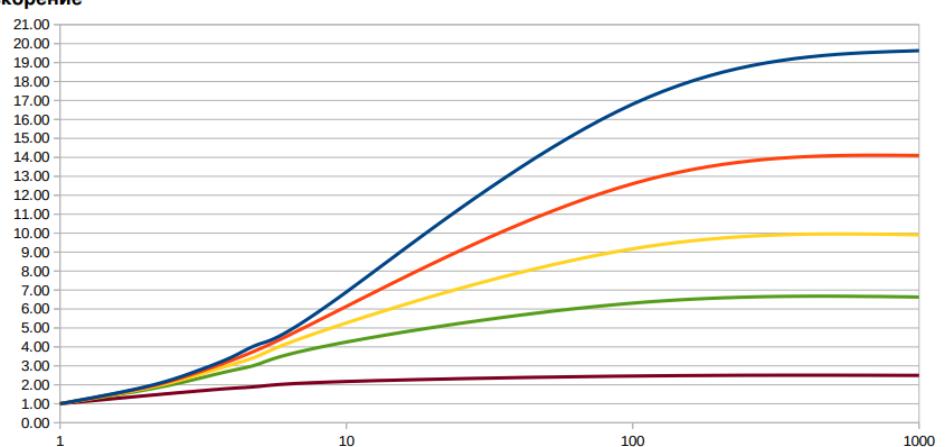
- Насколько можно ускорить программу на N процессорах с использованием потоков?
- Закон Амдала:

$$\text{Ускорение} = \frac{\text{Время работы на одном процессоре}}{\text{Время выполнения на } N \text{ процессорах}} = \frac{T \times (1-f) + T \times f}{T \times (1-f) + \frac{T \times f}{N}} = \frac{1}{(1-f) + \frac{f}{N}}$$

где, T - время работы; f — доля распараллеливания [0..1])

- Когда f мало, использование параллельного выполнения неэффективно
- Когда N → ∞, то ускорение ограничено 1/(1-f)

Ускорение



32. Механизм параллельных вычислений, функции ОС.



Механизм «параллельных» вычислений

- Однопроцессорные системы — процессы чередуются



- Многопроцессорные — чередуются и перекрываются.

- Конкуренция за общие ресурсы. Голодание.



Требуемые функции ОС

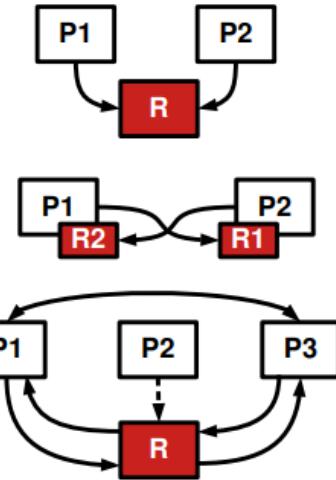
- Отслеживание ресурсов процесса/потока
- Распределение и освобождение ресурсов для каждого активного процесса/потока (CPU, файлы, память, ...)
- Защита ресурсов процесса/потока от непреднамеренного воздействия других процессов/потоков
- Независимость результата процесса/потока от скорости его выполнения и скорости других процессов/потоков

33. Проблемы параллельного выполнения:
взаимоисключений, взаимоблокировки, голодание.
Требования к взаимным исключениям. Уровни
взаимодействия процессов и потоков.



Проблемы параллельного выполнения

- Взаимоисключений (Mutual Exclusion)
— процессы/потоки не должны одновременно использовать критический ресурс
- Взаимоблокировки (DeadLocks, LiveLocks) — процессы/потоки не должны взаимозахватывать требуемые ресурсы.
- Голодание (Starvation) — конкуренция за ресурсы не должна порождать невозможность доступа к ресурсу



Требования к взаимным исключениям

- Взаимоисключений осуществляются в принудительном порядке.
 - В критическом участке кода должен находиться один процесс/поток
- Процесс/поток не должен влиять на другие процессы/потоки в некритическом участке
- Противодействие бесконечному ожиданию доступа к критическому участку
- Вход в свободный критический участок должен незамедлительно предоставляться
- Отсутствие предположений о количестве процессов или их скорости
- Ограничение времени нахождения в критических участках

- **Процессы не осведомлены о наличии друг друга.** Это независимые процессы, не предназначенные для совместной работы. Наилучшим примером такой ситуации может служить многозадачность множества независимых процессов. Это могут быть пакетные задания, интерактивные сессии или комбинация тех и других. Хотя эти процессы и не работают совместно, операционная система должна решать вопросы **конкурентного использования** ресурсов. Например, два независимых приложения могут требовать доступ к одному и тому же диску или к принтеру. Операционная система должна регулировать такие обращения.
- **Процессы косвенно осведомлены о наличии друг друга.** Эти процессы не обязательно должны быть осведомлены о наличии друг друга с точностью до идентификатора процесса, однако они совместно обращаются к некоторому объекту, например к буферу ввода-вывода. Такие процессы демонстрируют **сотрудничество** при совместном использовании общего объекта.
- **Процессы непосредственно осведомлены о наличии друг друга.** Такие процессы способны общаться один с другим с использованием идентификаторов процессов и изначально созданы для совместной работы. Также они демонстрируют **сотрудничество** при работе.

34. Примитивы синхронизации ОС. Предназначение примитивов синхронизации



Основные примитивы

- Семафоры (Semaphore) — захват и освобождение множественного ресурса
 - или одного (бинарные семафоры)
- Мьютексы (Mutex) — блокировка и освобождение ресурса единственным процессом/потоком
- Условные переменные (Conditional Variable) — Блокировка до выполнения какого-либо условия
- Блокировки чтения/записи (rw-lock) — отдельные блокировки на чтение и запись
- Мониторы (Monitor) — конструкции языков программирования, которые скрывают низкоуровневые примитивы синхронизации
- Флаги событий (Event Flags) — связывание условий продолжения выполнения с одним или несколькими флагами (битами блокирующей переменной)
- Почтовые ящики (Message Passing) — передача сообщений

Столлинг гл. 5.1 — Эволюция подхода к блокировке — самостоятельно прочитать!

ТАБЛИЦА 5.3. Основные механизмы параллельных вычислений

Семафор (Semaphore)	Целочисленное значение, используемое для передачи сигналов между процессами. Над семафором могут быть выполнены только три операции (все они являются атомарными): инициализация, уменьшение (декремент) и увеличение (инкремент) значения. Операция уменьшения может привести к блокировке процесса, а операция увеличения — к разблокированию. Известен также как семафор со счетчиком (counting semaphore) или обобщенный семафор (general semaphore)
Бинарный семафор (Binary semaphore)	Семафор, который может принимать только два значения — 0 и 1
Мьютекс (Mutex)	Аналогичен бинарному семафору. Ключевым отличием является то, что процесс, блокирующий мьютекс (устанавливающий его значение равным 0), должен и разблокировать его (установить его значение равным 1)

290 ГЛАВА 5. ПАРАЛЛЕЛЬНЫЕ ВЫЧИСЛЕНИЯ: ВЗАИМОИСКЛЮЧЕНИЯ И МНОГОЗАДАЧНОСТЬ

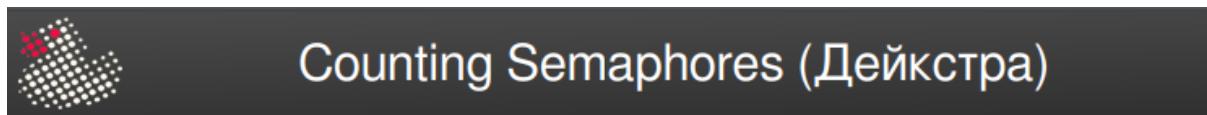
Окончание табл. 5.3

Условная переменная (Condition variable)	Тип данных, используемый для блокировки процесса или потока до тех пор, пока не станет истинным некоторое условие
Монитор (Monitor)	Конструкция языка программирования, инкапсулирующая переменные, процедуры доступа и код инициализации, в абстрактном типе данных. Переменные монитора могут быть доступны только через его процедуры доступа, и в любой момент времени только один процесс может активно работать с монитором. Процедуры доступа представляют собой <i>критические участки</i> . Монитор может иметь очередь процессов, ожидающих доступа к нему
Флаги событий (Event flags)	Слово памяти, используемое как механизм синхронизации. Код приложения может связать с каждым битом флага свое событие. Поток может ждать либо одного события, либо сочетания событий путем проверки одного или нескольких битов в соответствующем флаге. Поток блокируется до тех пор, пока все необходимые биты не будут установлены (И) или пока не будет установлен хотя бы один из битов (ИЛИ)
Почтовые ящики/ сообщения (Mailboxes/messages)	Средство обмена информацией между двумя процессами, которое может быть использовано для синхронизации
Спин-блокировки (Spinlocks)	Механизм взаимоисключения, в котором процесс выполняется в бесконечном цикле, ожидая, когда значение блокирующей переменной укажет доступность критического участка

35. Примитивы синхронизации ОС. Семафоры и мьютексы. Бинарный семафор

Семафоры (Semaphore) — захват и освобождение множественного ресурса – или одного (бинарные семафоры)

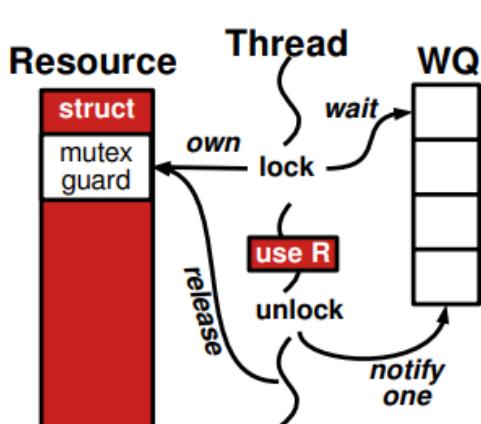
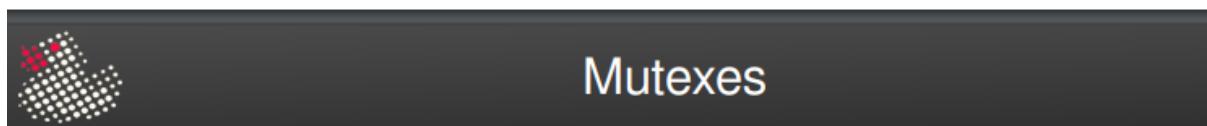
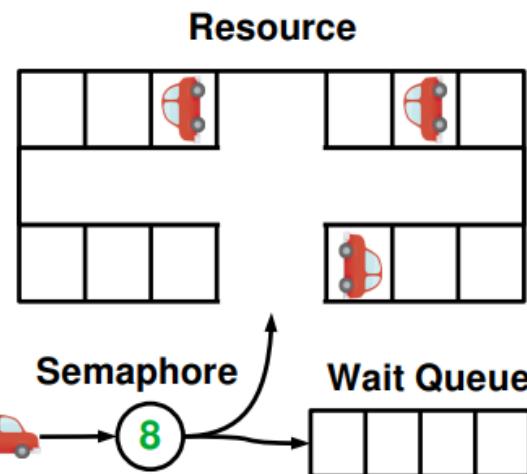
Мьютексы (Mutex) — блокировка и освобождение ресурса единственным процессом/потоком



- sema_p() «proberen», semWait

```
count--  
if count<0  
    thread blocks in WQ
```
- sema_v() «verhogen», semSignal

```
count++  
if count<=0  
    notify thread
```
- Политика выборки разная
– строгая, слабая, приоритет



- Множество различных реализаций:
 - Блокирующие, спин, адаптивные, фьютексы, ...
- Не- бинарный семафор!
 - Захвативший блокировку должен ее освободить
- Захват должен быть коротким
- Priority inversion

36. Примитивы синхронизации ОС. Условные переменные, rwlocks.

Условные переменные (Conditional Variable) — Блокировка до выполнения какого-либо условия

Блокировки чтения/записи (rw-lock) — отдельные блокировки на чтение и запись

Conditional Variables

- wait(condition, lock)
- signal(condition, lock)
- broadcast(condition, lock)

Producer/Consumer

```
void produce(char c) {
    lock_acquire(&lock);
    while (count == SIZE) {
        cond_wait(&spaceAvailable, &lock);
    }
    count++;
    buffer[putIndex] = c;
    putIndex++;
    if (putIndex == SIZE) {
        putIndex = 0;
    }
    cond_signal(&dataAvailable, &lock);
    lock_release(&lock);
}

char consume() {
    char c;
    lock_acquire(&lock);
    while (count == 0) {
        cond_wait(&dataAvailable, &lock);
    }
    count--;
    c = buffer[getIndex];
    getIndex++;
    if (getIndex == SIZE) getIndex = 0;
    cond_signal(&spaceAvailable, &lock);
    lock_release(&lock);
    return c;
}
```

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Multiple-reader, single-writer locks (rwlocks)

- Когда читатели читают, они захватывает readlock, количество одновременных читателей содержится в rwlock
- Когда писателю требуется записать — он устанавливает требование записи (want write) и ожидает на rwlock
- rwlock ждет освобождения readlock
- Оповещает писателей
 - Один захватывает writelock (читатели не могут читать)
- Оповещает читателей
 - Захватывается readlock (писатели должны требовать)



Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



37. Примитивы синхронизации ОС. Мониторы, флаги событий, передача сообщений.

Мониторы (Monitor) — конструкции языков программирования, которые скрывают низкоуровневые примитивы синхронизации

Флаги событий (Event Flags) — связывание условий продолжения выполнения с одним или несколькими флагами (битами блокирующей переменной)

Почтовые ящики (Message Passing) — передача сообщений



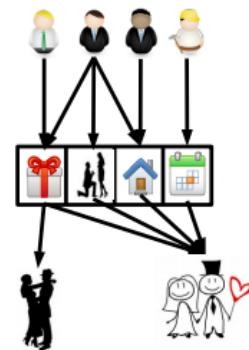
Monitors

- Набор процедур, выполняющих операции над общими данными
- Каждая процедура подразумевает захват блокировки общих данных
- Локальные переменные используются только внутри монитора
- Для ожидания и оповещения используются Conditional Variables
- Parallel Pascal, Java
- Реализуются высокоуровнево, программисту не требуется возится с захватом/освобождением или ожиданием/нотификацией
- ИМХО В операционных системах не используются



Event Flags

- Наборы (битовых) флагов ожидания событий
- Операции
 - Установка флага
 - Сброс флага
 - Ожидание флага
 - Ожидание какого-либо флага
 - Ожидание всех флагов
- Реализованы VMS, python



связывание условий продолжения выполнения с одним или несколькими флагами (битами блокирующей переменной)



- Две операции — send(получатель, сообщение) и receive(отправитель, сообщение)
- Необходим метод адресации получателя и отправителя
 - Прямая адресация (direct): явная и постфактум
 - Косвенная адресация (indirect) — номер «почтового ящика» - 1:1, 1:N, M:1, M:N
- Раздельная синхронизация посылки и получения
 - Блокирующая, неблокирующая, гарантия доставки
- Формат сообщения
 - Фиксированная, переменная длина, файл, ...
- Выборка из очереди — Приоритет, FIFO, ...

38. Примитивы синхронизации ОС. Неблокирующие примитивы синхронизации и неблокирующие структуры данных.

Блокировка вызывает переключение контекста – Дорого (долго)!

Преимущество неблокирующих алгоритмов — в лучшей масштабируемости по количеству процессоров. К тому же, если ОС прервёт один из потоков фоновой задачей, остальные, как минимум, выполнят свою работу, не простоявая. По максимуму — возьмут невыполненную работу на себя.

Сейчас модно пользоваться неблокирующими примитивами синхронизации и неблокирующими структурами – Wait-free (N-steps), lock-free (some N-steps, other retry), obstruction-free

Без препятствий (англ. obstruction-free)

Самая слабая из гарантий. Поток совершает прогресс, если не встречает препятствий со стороны других потоков. Алгоритм работает без препятствий, если поток, запущенный в любой момент (при условии, что выполнение всех препятствующих потоков приостановлено) завершит свою работу за детерминированное количество шагов. Синхронизация с помощью мьютексов не отвечает даже этому требованию: если поток остановится, захватив мьютекс, то остальные потоки, которым этот мьютекс нужен, будут простоявать.

Без блокировок (англ. *lock-free*)

Для алгоритмов без блокировок гарантируется системный прогресс по крайней мере одного потока. Например, поток, выполняющий операцию «[сравнение с обменом](#)» в цикле, теоретически может выполнять бесконечно, но каждая его итерация означает, что какой-то другой поток совершил прогресс, то есть система в целом совершает прогресс.

Без ожиданий (англ. *wait-free*)

Самая строгая гарантия прогресса. Алгоритм работает без ожиданий, если каждая операция выполняется за определённое количество шагов, не зависящее от других потоков

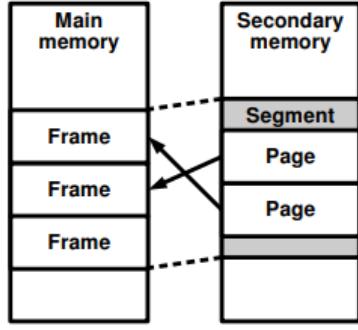
Примеры:

- Атомарные операции и примитивы синхронизации
- Неблокирующие структуры данных — списки, деревья, очереди, ... over 9000.

39. Управление памятью, основные определения и требования к организации.

Определения

- Основная память (main memory)
 - Область, где процессор может исполнять программы
- Вторичная память (secondary memory)
 - Область, где программы и данные могут храниться, в том числе и во время выполнения
- Кадр (frame) — область фиксированного размера основной памяти
- Страница (page) — область фиксированного размера во вторичной памяти, которая может быть скопирована в кадр
- Сегмент (segment) — область переменного размера в основной или вторичной памяти, может быть разделена на страницы



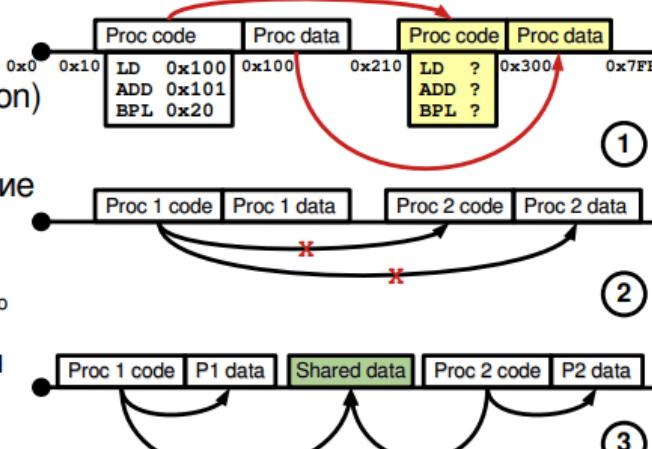
Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 3

Требования к организации памяти

Требования:

- Переместимость (relocation)
- Защита (protection)
- Совместное использование (sharing)
- Логическая организация
 - Управление основной и вторичной памятью
 - Модульная организация программ
- Физическая организация
 - Поддержка аппаратуры
 - Перекрытие (overlays) и виртуализация
 - Организация области подкачки



Операционные системы. Часть 3. Память и планирование

Перемещение

В многозадачной системе доступная основная память в общем случае разделяется среди множества процессов. Обычно программист не знает заранее, какие программы будут резидентно находиться в основной памяти во время работы разрабатываемой им программы. Кроме того, для максимизации загрузки процессора желательно иметь большой пул процессов, готовых к выполнению, для чего требуется возможность загрузки и выгрузки активных процессов из основной памяти. Требование, чтобы выгруженная из памяти программа была вновь загружена в то же самое место, где находилась и ранее, было бы слишком сильным ограничением. Крайне желательно,

чтобы программа могла быть перемещена (relocate) в другую область памяти. Таким образом, заранее неизвестно, где именно будет размещена программа, а кроме того, программа может быть перемещена из одной области памяти в другую при свопинге. Эти обстоятельства обусловливают наличие определенных технических требований к адресации, проиллюстрированных на рис. 7.1.

402 ГЛАВА 7. УПРАВЛЕНИЕ ПАМЯТЬЮ



Рис. 7.1. Требования к адресации процесса

На рисунке представлен образ процесса. Для простоты предположим, что образ процесса занимает одну непрерывную область основной памяти. Очевидно, что операционной системе необходимо знать местоположение управляющей информации процесса и стека выполнения, а также точки входа для начала выполнения процесса. Поскольку управлением памятью занимается операционная система и она же размещает процесс в основной памяти, соответствующие адреса она получает автоматически. Однако, помимо получения операционной системой указанной информации, процесс должен иметь возможность обращаться к памяти в самой программе. Так, команды ветвления содержат адреса, указывающие на команды, которые должны быть выполнены после них; команды обращения к данным - адреса байтов или слов, с которыми они работают. Так или иначе, но процессор и программное обеспечение операционной системы должны быть способны перевести ссылки в коде программы в реальные физические адреса, соответствующие текущему расположению программы в основной памяти.

Защита

Каждый процесс должен быть защищен от нежелательного воздействия других процессов, случайного или преднамеренного. Следовательно, код других процессов не

должен иметь возможности без разрешения обращаться к памяти данного процесса для чтения или записи. Однако удовлетворение требованию перемещаемости усложняет задачу защиты.

Поскольку расположение программы в основной памяти непредсказуемо, проверка абсолютных адресов во время компиляции невозможна. Кроме того, в большинстве языков программирования возможно динамическое вычисление адресов во время выполнения (например, вычисление адреса элемента массива или указателя на поле структуры данных). Следовательно, во время работы программы необходимо выполнять проверку всех обращений к памяти, генерируемых процессом, чтобы удостовериться, что все они - только к памяти, выделенной данному процессу. К счастью, как вы увидите позже, механизмы поддержки перемещений обеспечивают и поддержку защиты. Обычно пользовательский процесс не может получить доступ ни к какой части операционной системы - ни к коду, ни к данным. Код одного процесса не может выполнить команду ветвления, целевой код которой находится в другом процессе. Если не приняты специальные меры, код одного процесса не может получить доступ к данным другого процесса. Процессор должен быть способен прервать выполнение таких команд. Заметим, что требования защиты памяти должны быть удовлетворены на уровне процессора (аппаратного обеспечения), а не на уровне операционной системы (программного обеспечения), поскольку операционная система не в состоянии предвидеть все обращения к памяти, которые будут выполнены программой. Даже если бы такое было возможно, сканирование каждой программы в поиске предлагаемых нарушений защиты было бы слишком расточительным с точки зрения использования процессорного времени. Следовательно, соответствующие возможности аппаратного обеспечения - единственное средство определения допустимости обращения к памяти (данным или коду) во время работы программы

Совместное использование

Любой механизм защиты должен иметь достаточную гибкость для того, чтобы обеспечить возможность нескольким процессам обращаться к одной и той же области основной памяти. Например, если несколько процессов выполняют один и тот же машинный код, то будет выгодно позволить каждому процессу работать с одной и той же копией этого кода, а не создавать собственную. Процессам, сотрудничающим в работе над некоторой задачей, может потребоваться совместный доступ к одним и тем же структурам данных. Система управления памятью должна, таким образом, обеспечивать управляемый доступ к разделяемым областям памяти, при этом никоим образом не ослабляя защиту памяти. Как мы увидим позже, механизмы поддержки перемещений обеспечивают и поддержку совместного использования памяти.

Логическая организация

Практически всегда основная память в компьютерной системе организована как линейное (одномерное) адресное пространство, состоящее из последовательности байтов или слов. Аналогично организована и вторичная память на своем физическом уровне. Хотя такая организация и отражает особенности используемого аппаратного обеспечения, она не соответствует способу, которым обычно создаются программы. Большинство программ организованы в виде модулей, одни из которых неизменны (только для чтения, только для выполнения), а другие содержат данные, которые могут быть изменены. Если операционная система и аппаратное обеспечение компьютера

могут эффективно работать с пользовательскими программами и данными, представленными модулями, то это обеспечивает ряд преимуществ.

1. Модули могут быть созданы и скомпилированы независимо один от другого, при этом все ссылки из одного модуля во второй разрешаются системой во время работы программы.
2. Разные модули могут получать разные степени защиты (только для чтения, только для выполнения) за счет весьма умеренных накладных расходов.
3. Возможно применение механизма, обеспечивающего совместное использование модулей разными процессами.

Основное достоинство обеспечения совместного использования на уровне модулей заключается в том, что они соответствуют взгляду программиста на задачу и, следовательно, ему проще определить, требуется ли совместное использование того или иного модуля. Инструментом, наилучшим образом удовлетворяющим данным требованиям, является сегментация, которая будет рассмотрена в данной главе среди прочих методов управления памятью.

Физическая организация

Как указывалось в разделе 1.5, память компьютера разделяется как минимум на два уровня: основная и вторичная. Основная память обеспечивает быстрый доступ по относительно высокой цене; кроме того, она энергозависима, т.е. не обеспечивает долговременное хранение. Вторичная память медленнее и дешевле основной и обычно не энергонезависима. Следовательно, вторичная память большой емкости может служить для долговременного хранения программ и данных, а основная память меньшей емкости - для хранения программ и данных, использующихся в текущий момент. В такой двухуровневой структуре основной заботой системы становится организация потоков информации между основной и вторичной памятью.

Ответственность за эти потоки может быть возложена и на отдельного программиста, но это непрактично и нежелательно по следующим причинам.

1. Основной памяти может быть недостаточно для программы и ее данных. В этом случае программист вынужден прибегнуть к практике, известной как структуры с перекрытием - оверлеи (overlay), когда программа и данные организованы таким образом, что различные модули могут быть назначены одной и той же области памяти; основная программа при этом ответственна за перезагрузку модулей при необходимости. Даже при помощи соответствующего инструментария компиляции оверлеев разработка таких программ приводит к дополнительным затратам времени программиста.
2. Во многозадачной среде программист при разработке программы не знает, какой объем памяти будет доступен программе и где эта память будет располагаться. Таким образом, очевидно, что задача перемещения информации между двумя уровнями памяти должна возлагаться на операционную систему. Эта задача является сущностью управления памятью.

кратко по распределению памяти

Таблица 7.2. Технологии управления памятью

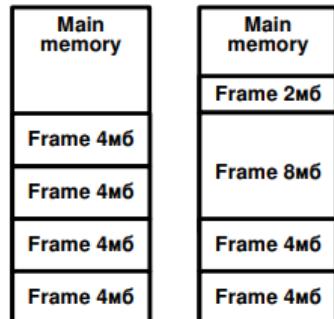
Технология	Описание	Сильные стороны	Слабые стороны
Фиксированное распределение	Основная память разделяется на ряд статических разделов во время генерации системы. Процесс может быть загружен в раздел равного или большего размера	Простота реализации, малые системные накладные расходы	Неэффективное использование памяти из-за внутренней фрагментации, фиксированное максимальное количество активных процессов
Динамическое распределение	Разделы создаются динамически; каждый процесс загружается в раздел строго необходимого размера	Отсутствует внутренняя фрагментация, более эффективное использование основной памяти	Неэффективное использование процесора из-за необходимости уплотнения для противодействия внешней фрагментации
Простая страницчная организация	Основная память разделена на ряд кадров равного размера. Каждый процесс разделен на некоторое количество страниц равного размера и той же длины, что и кадры. Процесс загружается путем загрузки всех его страниц в доступные, но не обязательно последовательные кадры	Отсутствует внешняя фрагментация	Наличие небольшой внутренней фрагментации
Простая сегментация	Каждый процесс разделен на ряд сегментов. Процесс загружается путем загрузки всех своих сегментов в динамические (не обязательно смежные) разделы	Отсутствует внутренняя фрагментация; по сравнению с динамическим распределением повышенная эффективность использования памяти и сниженные накладные расходы	Внешняя фрагментация
Страницчная организация виртуальной памяти	Все, как при простой страницочной организации, с тем исключением, что не требуется одновременно загружать все страницы процесса. Необходимые нерезидентные страницы автоматически загружаются в память	Нет внешней фрагментации; более высокая степень много задачности; большое виртуальное адресное пространство	Накладные расходы из-за сложности системы управления памятью
Сегментация виртуальной памяти	Все, как при простой сегментации, с тем исключением, что не требуется одновременно загружать все сегменты процесса. Необходимые нерезидентные сегменты автоматически загружаются в память	Нет внутренней фрагментации; более высокая степень многозадачности; большое виртуальное адресное пространство; поддержка защиты и совместного использования	Накладные расходы из-за сложности системы управления памятью

40. Фиксированное и динамическое размещение программ в памяти.



Fixed partitioning

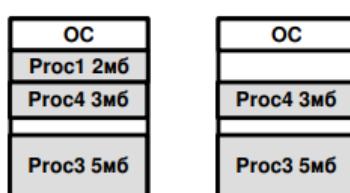
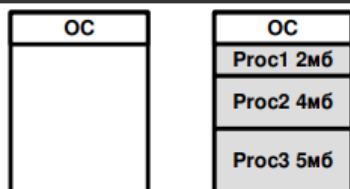
- У программ разная длина, их сложно упаковывать в память
- А что если оперировать кадрами заданного размера?
- Минусы
 - Память расходуется неэффективно для небольших процессов
 - Ограничение на количество одновременных процессов
 - Если процесс не помещается в раздел, он должен использовать перекрытия (overlays), самостоятельно загружая и выгружая свои части



Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Dynamic partitioning



- Блоки процессов разного размера, ОС встраивает их в свободные места
- Минусы
 - Фрагментация памяти
 - Необходимость упаковки памяти
 - Начинает за здравие, заканчивает ну вы понимаете...
 - Сложные и медленные алгоритмы размещения в памяти
- Подходит для загрузки драйверов
 - Поэтому и нужна была перезагрузка при добавлении драйверов
- Разумный компромисс - «Buddy System» - см. Столлингса.

Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



подробнее см столлингса страницы 405-411.

Система двойников

Как фиксированное, так и динамическое распределение памяти имеют свои недостатки. Фиксированное распределение ограничивает количество активных процессов и неэффективно использует память при несоответствии между размерами разделов и процессов. Динамическое распределение реализуется более сложно и включает накладные расходы по уплотнению памяти. Интересным компромиссом в этом плане является система двойников ([136], [191]).

В системе двойников память распределяется блоками размером 2^k , $L \leq k \leq U$, где

2^L — минимальный размер выделяемого блока памяти;

2^U — наибольший распределяемый блок; вообще говоря, 2^U представляет собой размер всей доступной для распределения памяти.

Вначале все доступное для распределения пространство рассматривается как единый блок размером 2^U . При запросе размером s , таким, что $2^{U-1} < s \leq 2^U$, выделяется весь блок. В противном случае блок разделяется на два эквивалентных двойника с размерами 2^{U-1} . Если $2^{U-2} < s \leq 2^{U-1}$, то по запросу выделяется один из двух двойников; в противном случае один из двойников вновь делится пополам. Этот процесс продолжается до тех пор, пока не будет сгенерирован наименьший блок, размер которого не меньше s . Система двойников постоянно ведет список “дыр” (доступных блоков) для каждого размера 2^i . Дыра может быть удалена из списка ($i+1$) разделением ее пополам и внесением двух новых дыр размером 2^i в список i . Когда пара двойников в списке i оказывается освобожденной, они удаляются из списка и объединяются в единый блок в списке ($i+1$). Ниже приведен рекурсивный алгоритм для удовлетворения запроса размером $2^{i-1} < k \leq 2^i$, в котором осуществляется поиск дыры размером 2^i .

```
void get_hole(int i)
{
    if (i == (U+1))  < Ошибка >;
    if (< Список i пуст >
    {
        get_hole(i+1);
        < Разделить дыру на двойники >;
        < Поместить двойники в список i >;
    }
    < Взять первую дыру из списка i >;
}
```

На рис. 7.6 приведен пример использования блока с начальным размером 1 Мбайт. Первый запрос А — на 100 Кбайт (для него требуется блок размером 128 Кбайт). Для этого начальный блок делится на два двойника по 512 Кбайт. Первый из них делится на двойники размером 256 Кбайт, и, в свою очередь, первый из получившихся при этом разделении двойников также делится пополам. Один из получившихся двойников размером 128 Кбайт выделяется запросу А. Следующий запрос В требует 256 Кбайт. Такой блок имеется в наличии и выделяется. Процесс продолжается с разделением и слиянием двойников при необходимости. Обратите внимание, что после освобождения блока Е происходит слияние двойников по 128 Кбайт в один блок размером 256 Кбайт, который, в свою очередь, тут же сливается со своим двойником.

На рис. 7.7 показано представление системы двойников в виде бинарного дерева, непосредственно после освобождения блока В. Листья представляют текущее распределение памяти. Если два двойника являются листьями, то по крайней мере один из них занят; в противном случае они должны слиться в блок большего размера.

Система двойников представляет собой разумный компромисс для преодоления недостатков схем фиксированного и динамического распределения, но в современных операционных системах ее превосходит виртуальная память, основанная на страничной организации и сегментации. Однако система двойников нашла применение в параллельных системах как эффективное средство распределения и освобождения параллельных программ (см., например, [118]). Модифицированная версия системы двойников используется для распределения памяти ядром UNIX (подробнее об этом вы узнаете в главе 8, “Виртуальная память”).

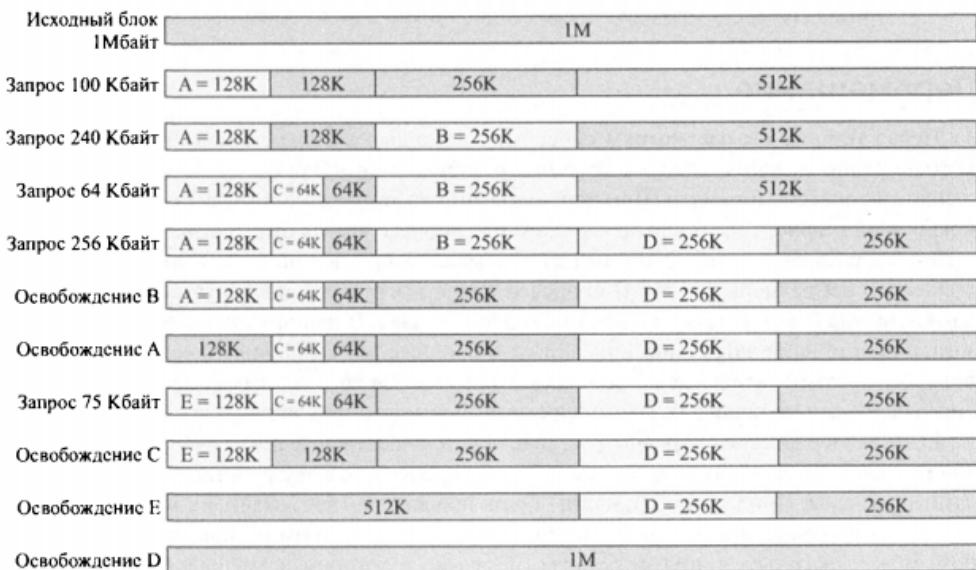


Рис. 7.6. Пример системы двойников

41. Модели аппаратного перемещения программ.

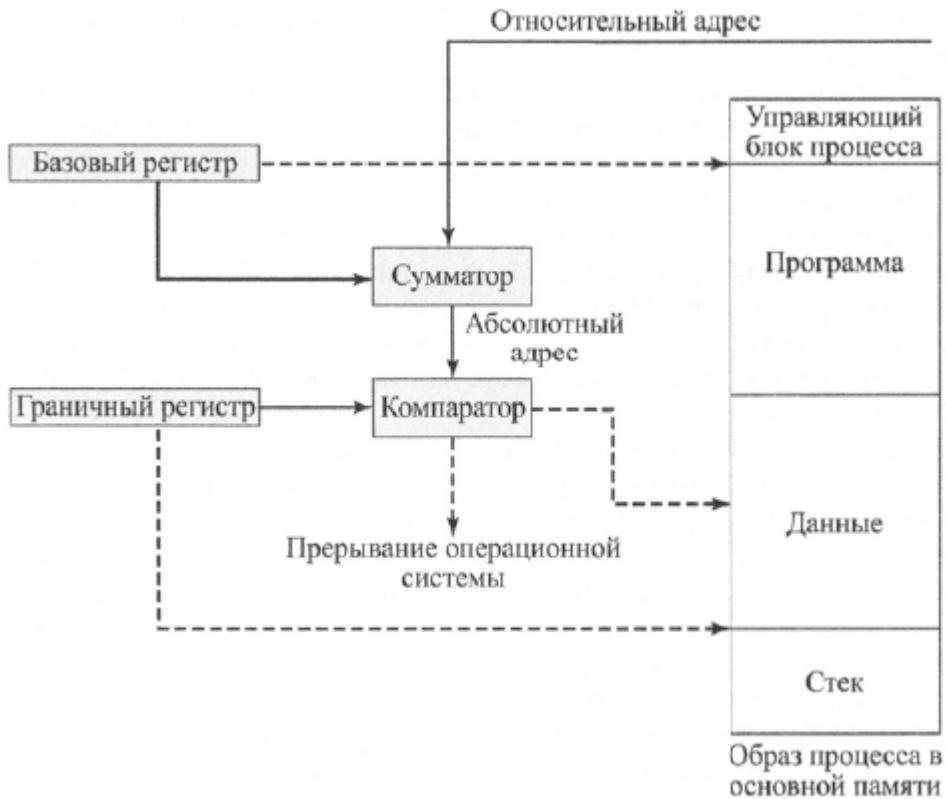
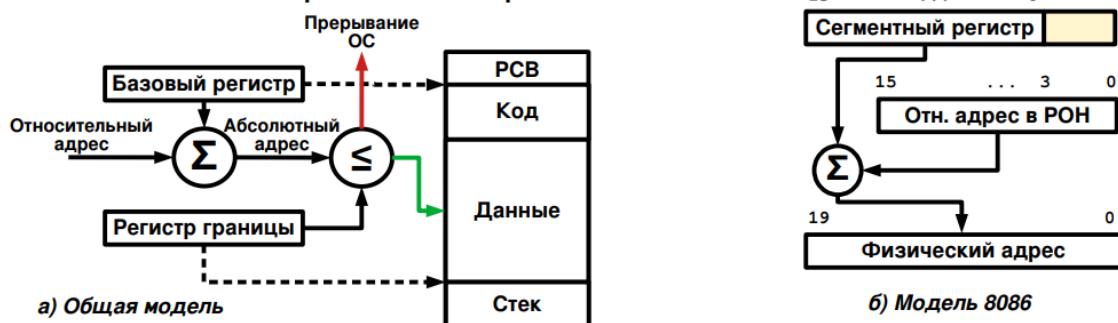


Рис. 7.8. Аппаратная поддержка перемещения

- Копировать программы средствами ОС — долго!
 - Нужна аппаратная поддержка, например, как в 8086
- Модели аппаратного перемещения:



42. Простой страничный поход и простая сегментная организация.



Simple paging

- Давайте разделим память на кадры одинакового (небольшого размера)
 - Пусть страница занимает целиком кадр
 - Для удобства размер кратен степени 2
- Уменьшится внутренняя фрагментация
 - Пустое место будет в последнем блоке каждого процесса
- Внешняя фрагментация исчезнет совсем
- Необходимы таблицы страниц
 - ОС должна знать, где ее герои

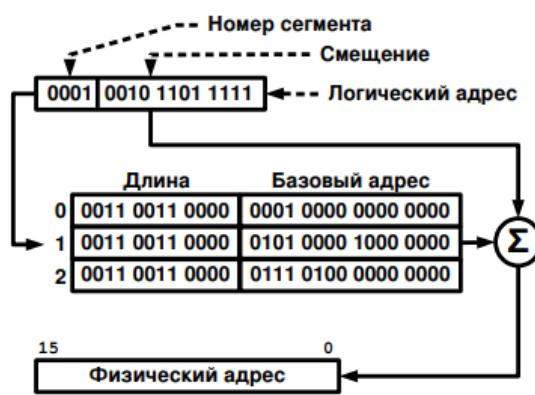
Таблицы страниц	
Proc A	1
	2
Proc B	6
	7
Proc C	8
	9
Proc A.1	3
	4
Proc C.1	5
	6
Proc B.1	7
	8
Proc C.2	9
	10
Proc C.3	11
	12
Proc C.4	13
	14
Proc B.2	15
	16
Proc B.3	17
	18
Proc C.4	19
	20
Proc A.2	21
	22

ОС	0
Proc A.1	1
Proc A.2	2
Proc C.1	3
Proc C.2	4
Proc C.3	5
Proc B.1	6
Proc B.2	7
Proc B.3	8
Proc C.4	9
Proc C.4	10
	11
	12
	13
	14
	15
	16
	17
	18
	19
	20
	21
	22

Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Simple Segmentation



- Таблица сегментов. Программист должен устанавливать:
 - Длину сегмента
 - Базовый адрес
- Внутренняя сегментация отсутствует
- Внешняя сегментация снижается

Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



43. Виртуальная память основные определения и принципы организации аппаратуры и управляемых программ.

Реальный (физический) адрес — адрес в основной памяти

Виртуальный адрес — логический адрес внутри процесса

Адресное пространство — диапазон адресов процесса

Виртуальное адресное пространство — область для одного процесса с виртуальными адресами

Виртуальная память — схема расположения процессов в памяти, в которой:

- Вторичная память адресуется так же как и основная
 - Виртуальные адреса транслируются в адреса в основной памяти
 - Основная память может быть расширена на вторичную
 - Размер памяти ограничен схемой адресации, но не фактическим количеством ячеек
- Организация виртуальной памяти это совокупность аппаратных и программных средств

Схемы реализации

- совокупность страничной и сегментной адресации
 - – Логические адреса динамически транслируются в физические
- Сегменты и страницы не должны располагаться последовательно в основной памяти
- Любая часть процесса в разные моменты выполнения может находиться во вторичной памяти и менять физические адреса в основной

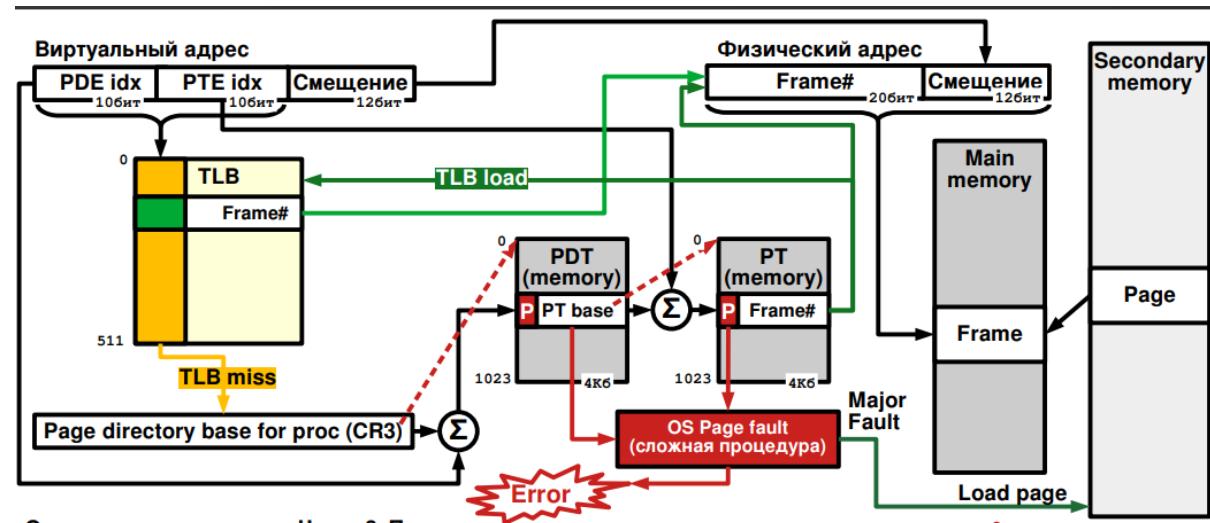
Resident set (резидентная часть) — часть процесса, находящаяся в основной памяти

Real memory — часть памяти, где процесс непосредственно выполняется

Virtual memory — часть памяти, которую процесс может занять

Следовательно, ОС может одновременно с бОльшим количеством процессов

44. Виртуальный страничный обмен. Двухуровневая организация MMU и TLB 80386. (для КОТ и ГТ — общие принципы)



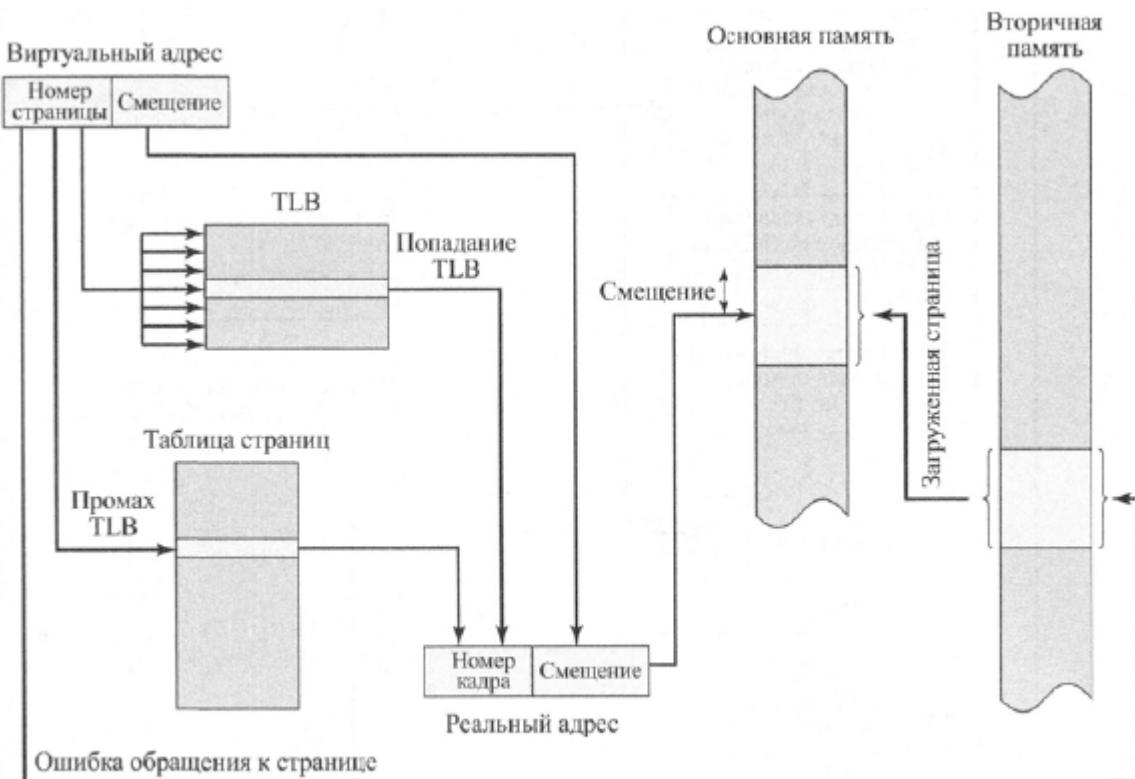
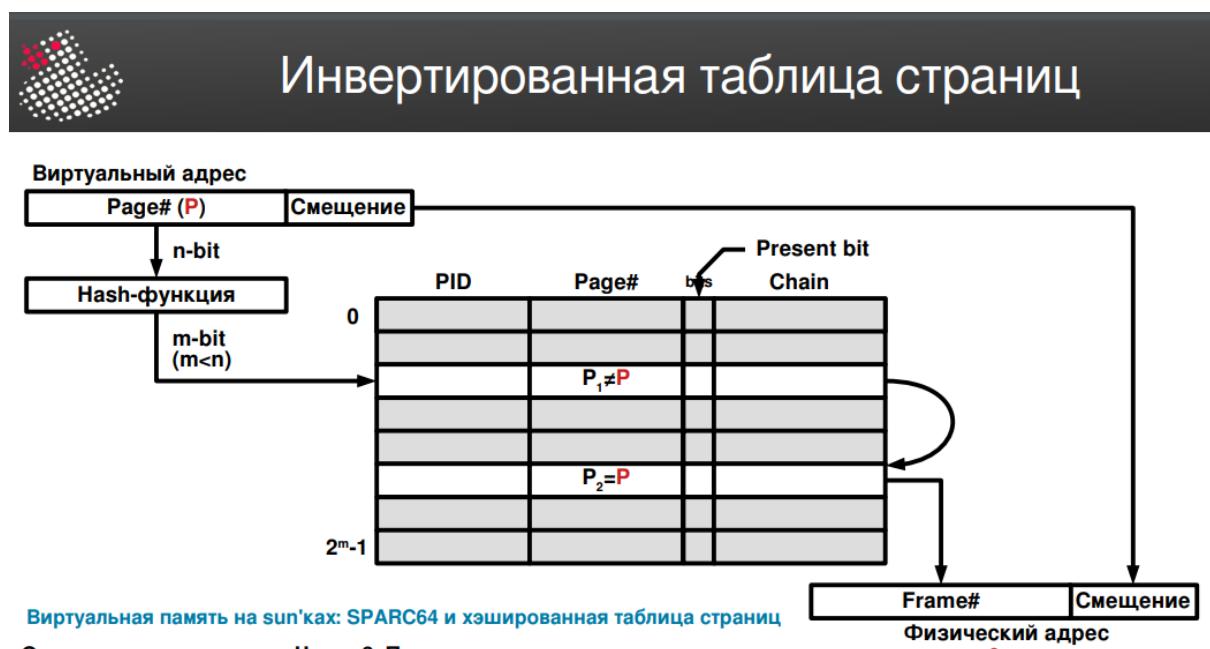


Рис. 8.6. Использование TLB

45. Инвертированная таблица страниц.



46. Сегментно-страничная виртуальная память.



454 ГЛАВА 8. Виртуальная память

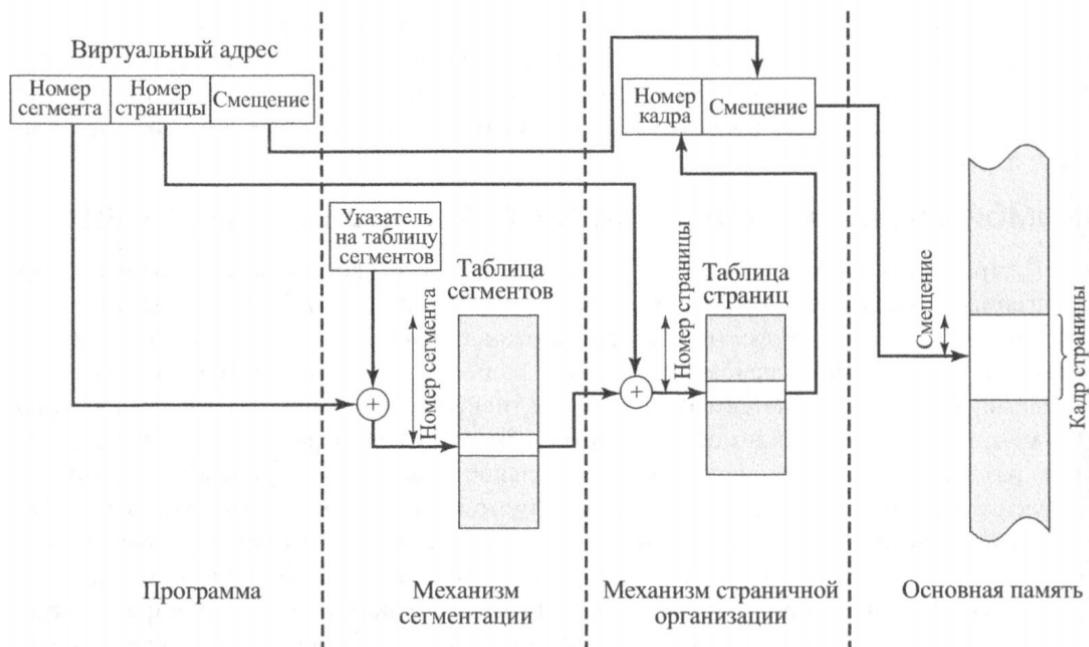


Рис. 8.12. Трансляция адреса при совместном использовании сегментации и страничной организации

47. Влияние размера страницы виртуальной памяти на ОС. Стратегии ОС по работе с виртуальной памятью.

Размер страницы

- При разработке ОС необходимо определится с размером страницы.
Важно учитывать:
 - Размеры таблиц страниц
 - Внутреннюю фрагментация
 - Количество page-fault при трансляции адреса
 - Скорость взаимодействия со вторичной памятью (и размер блока)
 - Локальность данных (в многопоточных приложениях ниже)
 - Количество промахов TLB, размер TLB, размер кэша L1, L2, L3, и др.
- Современные процессоры могут работать с разными размерами страниц (large, huge pages)
 - Intel 4Кб, 2Мб, 1Гб
 - Sparc64 VI 8Кб, 64Кб, 512Кб, 4Мб, 32Мб, 256Мб

Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 17

Стратегии ОС по работе с виртуальной памятью

- Стратегия выборки страниц из вторичной памяти:
 - По требованию, предварительная выборка
- Стратегия размещения
 - В современных NUMA системах скорость выше рядом с процессором процесса/потока
- Стратегия очистки (выгрузки во вторичную память)
 - По требованию и предварительная очистка
- Управление многозадачностью
 - Выгрузка процессов целиком
- Стратегии замещения (см. далее)
- Управление резидентной частью процессов (см. далее)

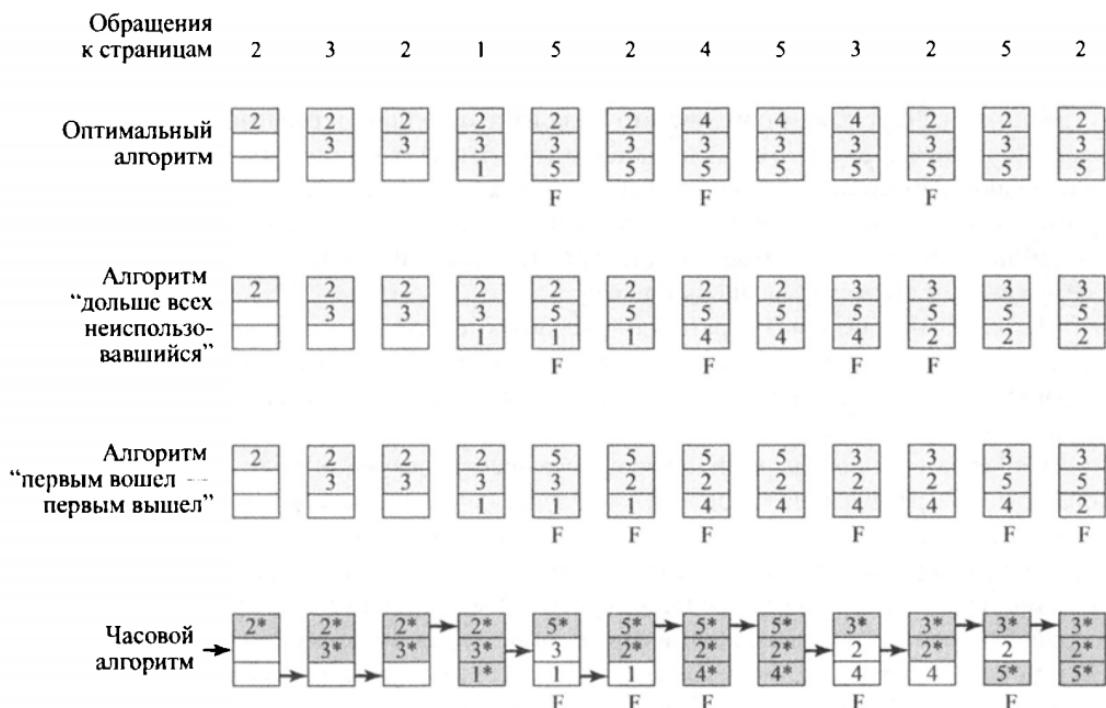
Операционные системы. Часть 3. Память и планирование

17

48. Стратегии замещения страниц ОС. Часовой Алгоритм. Управление резидентной частью процесса.

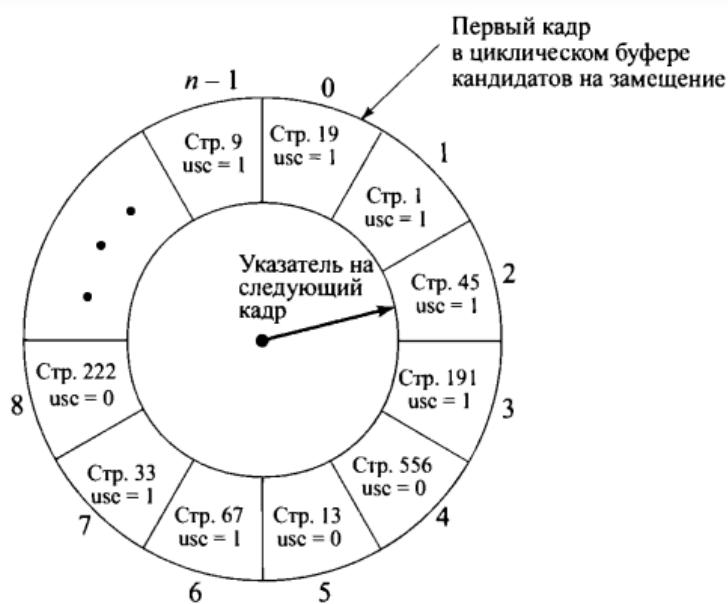
Какие именно кадры работающего процесса нужно выбрать для замещения и использования другими процессами?

- Выгрузим кадр, обращения к которому не последует в ближайшее время
- Используется статистика прошлого поведения
- Необходимо учитывать *locked frames*
- не все страницы можно выгрузить (часть ядра, буферы, ...)
- Если кадр совместно используется много раз (н-р код разделяемой библиотеки) то его выгрузка может повлиять на много процессов

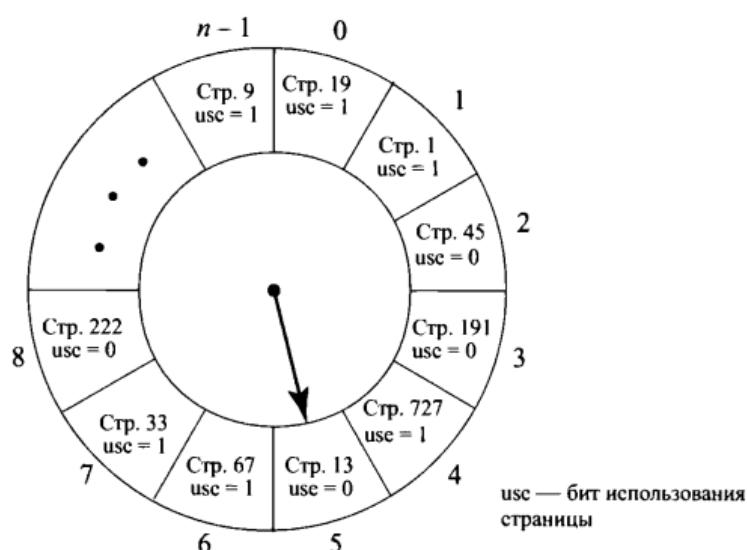


F --- прерывания обращения к странице после первоначального заполнения кадров

Рис. 8.14. Поведение четырех алгоритмов замещения страниц



а) Состояние буфера непосредственно перед замещением страницы

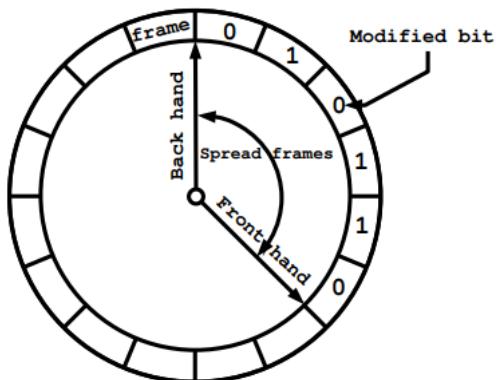


б) Состояние буфера непосредственно после замещения страницы

Рис. 8.15. Пример работы часового алгоритма



Модифицированный Clock Algorithm (Solaris 2.0 - 11)



- Алгоритм:

- FH — сбрасывает бит модификации
- Модификация фрейма устанавливается
- BH — Проверяет

- Расстояние между стрелками меняется динамически
- На современных размерах памяти занимает много времени
 - $16\text{Гб}/4\text{кб} = 4194304$ фреймов



Управление резидентной частью процессов

- Размер резидентной части:
 - При загрузке и работе процесса нет смысла держать в памяти все его страницы
 - Меньше памяти процессу → больше процессов в памяти
 - Меньше страниц процесса в памяти → больше pagefaults
 - После N-страниц дополнительное выделение памяти не приводит к кардинальному снижению pagefaults
 - Стратегии управления: Фиксированный и динамический размер
- Область видимости (scope)
 - Локальная — страница замещается у процесса вызвавшего pagefault
 - Глобальная — сканируются страницы всех процессов, кроме locked и сильно shared

Таблица 8.5. УПРАВЛЕНИЕ РЕЗИДЕНТНЫМ МНОЖЕСТВОМ

	Локальное замещение	Глобальное замещение
Фиксированное распределение	<ul style="list-style-type: none">• Количество кадров процесса фиксировано• Страница для замещения выбирается среди выделенных процессу кадров	<ul style="list-style-type: none">• Невозможно
Переменное распределение	<ul style="list-style-type: none">• Количество выделенных процессу кадров может время от времени изменяться• Страница для замещения выбирается среди выделенных процессу кадров	<ul style="list-style-type: none">• Страница для замещения выбирается среди всех доступных кадров в основной памяти; это приводит к изменению размера резидентного множества процесса

49. Виды планирования процессов. Критерии краткосрочного планирования. Приоритеты.

Таблица 9.1. Типы планирования

Долгосрочное планирование	Решение о добавлении процесса в пул выполняемых процессов
Среднесрочное планирование	Решение о добавлении процесса к числу процессов, полностью или частично размещенных в основной памяти
Краткосрочное планирование	Решение о том, какой из доступных процессов будет выполняться процессором
Планирование ввода-вывода	Решение о том, какой из запросов процессов на операции ввода-вывода будет обработан доступным устройством ввода-вывода

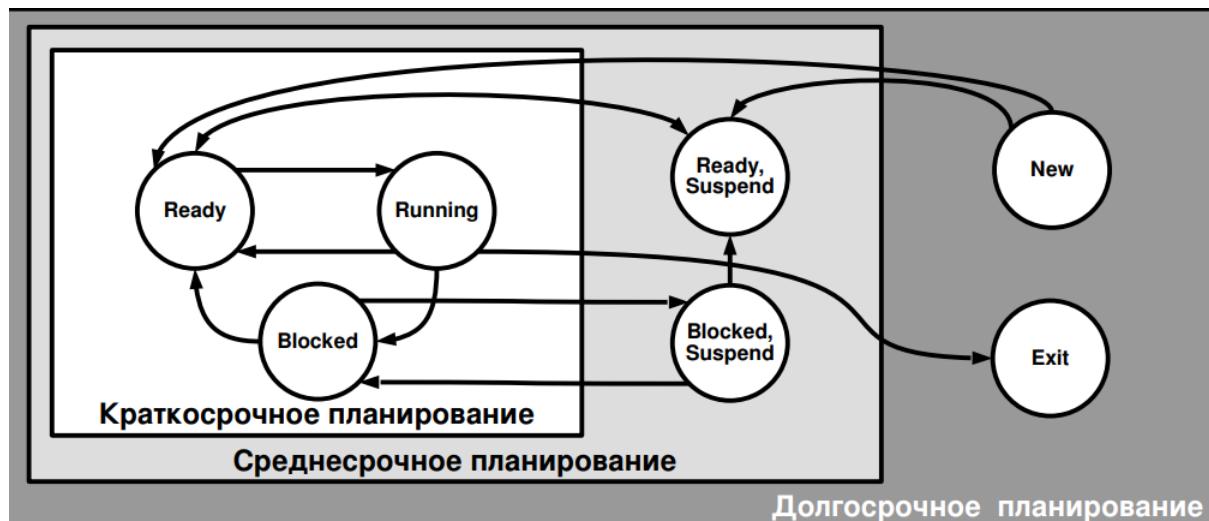


ТАБЛИЦА 9.2. КРИТЕРИИ ПЛАНИРОВАНИЯ

Пользовательские, связанные с производительностью	
Время оборота	Интервал времени между передачей процесса для выполнения и его завершением. Включает время выполнения, а также время, затраченное на ожидание ресурсов, в том числе процессора. Критерий вполне применим для пакетных заданий
Время отклика	В интерактивных процессах это время, истекшее между подачей запроса и началом получения ответа на него. Зачастую процесс может начать вывод информации пользователю, еще не окончив полной обработки запроса, так что описанный критерий — наиболее подходящий с точки зрения пользователя. Стратегия планирования должна пытаться сократить время получения ответа при максимизации количества интерактивных пользователей, время отклика для которых не выходит за заданные пределы
Предельный срок	При указании предельного срока завершения процесса планирование должно подчинить ему все прочие цели максимизации количества процессов, завершающихся в срок

Окончание табл. 9.2

Пользовательские, иные	
Предсказуемость	Данное задание должно выполняться примерно за одно и то же количество времени и с одной и той же стоимостью, независимо от загрузки системы. Большие вариации времени выполнения или времени отклика дезориентируют пользователей. Это явление может сигнализировать о больших колебаниях загрузки или о необходимости дополнительной настройки системы для устранения нестабильности ее работы
Системные, связанные с производительностью	
Пропускная способность	Стратегия планирования должна пытаться максимизировать количество процессов, завершающихся за единицу времени, что является мерой количества выполненной системой работы. Очевидно, что эта величина зависит от средней продолжительности процесса; однако на нее влияет и используемая стратегия планирования
Использование процессора	Этот показатель представляет собой процент времени, в течение которого процессор оказывается занятым. Для дорогих совместно используемых систем этот критерий достаточно важен; в однопользовательских же и некоторых других системах (типа систем реального времени) — менее важен по сравнению с рядом других
Системные, иные	
Беспристрастность	При отсутствии дополнительных указаний от пользователя или системы все процессы рассматриваются как равнозначные и ни один процесс не должен подвергнуться голоданию
Использование приоритетов	Если процессам назначены приоритеты, стратегия планирования должна отдавать предпочтение процессам с более высоким приоритетом
Баланс ресурсов	Стратегия планирования должна поддерживать занятость системных ресурсов. Предпочтение должно быть отдано процессу, который недостаточно использует важные ресурсы. Этот критерий включает использование долгосрочного и среднесрочного планирования

50. Использование приоритетов.



Использование приоритетов

- Процессам присваивается цифровой приоритет
 - В разных ОС - разные схемы назначения приоритетов
- Из очередей выбирается процесс с наивысшим приоритетом
- Если приоритеты процессов совпадают, то используется дополнительная стратегия
- Процессы с низким приоритетом могут голодать
- Приоритеты могут динамически изменяться

При выборе процесса планировщик начинает с очереди процессов с наивысшим приоритетом (RQ0). Если в очереди имеется один или несколько процессов, процесс для работы выбирается с использованием некоторой стратегии планирования. Если очередь RQ0 пуста, рассматривается очередь RQ1 и т.д.

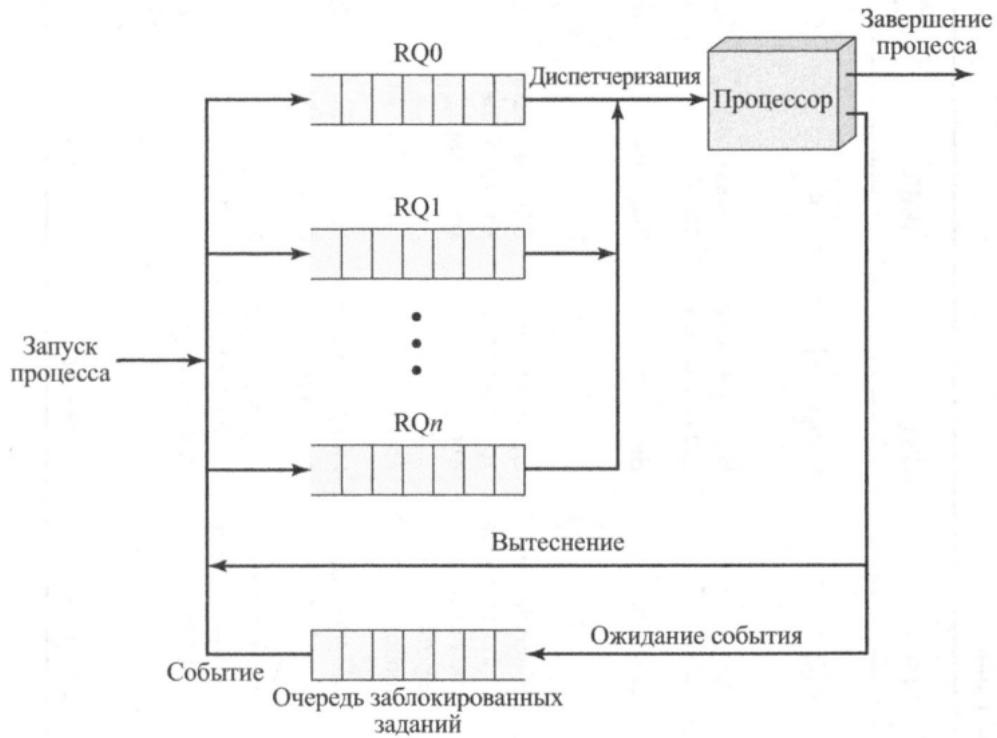


Рис. 9.4. Планирование с учетом приоритета процессов

Одна из основных проблем в такой чисто приоритетной схеме планирования состоит в том, что процессы с низким приоритетом могут оказаться в состоянии голодания. Это будет происходить при постоянном поступлении новых готовых к выполнению процессов с высоким приоритетом. Если такое поведение нежелательно, приоритет процесса может изменяться с его "возрастом" или историей выполнения (пример такой стратегии планирования будет приведен ниже).

51. Стратегии планирования FCFS, RR, SPN, SRT, HRRN, Feedback.



Стратегии планирования

- First Come First Served (FCFS) — аналог FIFO
- Round Robin — карусельное планирование
- Shortest Process Next — Короткие процессы вперед
- Shortest Remaining Time — Наименьшее время до завершения
- Highest Response Ratio Next — Наивысшее отношение отклика
- Feedback — снижения приоритета в зависимости от длительности времени исполнения
 - в случае blocked или preempted — feed it back

Операционные системы. Часть 3. Память и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Сравнение стратегий

	FCFS	RR	SPN	SRT	HRRN	Feedback
Функция выбора	max(w)	Const TQ	min(s)	min(s-e)	max($w+s$) s	$w \uparrow \rightarrow \text{prio} \uparrow$ $e \uparrow \rightarrow \text{prio} \downarrow$
Preemption	No	At Time Quantum	No	At arrival to ready queue	No	At Time quantum
Throughput	-	Can be low if TQ low	High	High	High	-
Response time	Can be high	Good for short	Good for short	Good	Good	-
Overhead	Minimum	Minimum	Can be high	Can be high	Can be high	Can be high
Effect	Bad for short and high I/O	fair	Bad for long	Bad for long	Balanced	Can prefer high IO proc.
Starvation	No	No	Possible	Possible	No	Possible

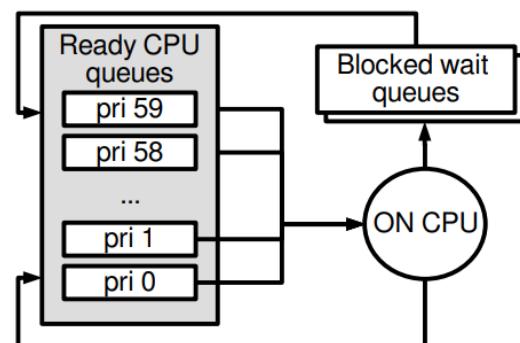
w — общее время ожидания процесса в очередях, e — общее время выполнения,
s — общее время, предполагаемое или заданное, для обслуживания процесса, включая e

Таблица 9.3. ХАРАКТЕРИСТИКИ РАЗЛИЧНЫХ СТРАТЕГИЙ ПЛАНРИРОВАНИЯ

	FSCF	Круговая	SPN	SRT	HRRN	Со снижением приоритета
Функция выбора	$\max[w]$	const	$\min[s]$	$\min[s - e]$	$\max\left(\frac{w+s}{s}\right)$	См. текст
Режим решения	Невытесняющий	Вытесняющий (по времени)	Невытесняющий	Вытесняющий (по решению)	Невытесняющий	Вытесняющий (по времени)
Пропускная способность	Не важна	Может быть низкой при малом кванте времени	Высокая	Высокая	Высокая	Не важна
Время отклика	Может быть большим, в особенности при больших отклонениях во времени выполнения процесса	Обеспечивает хорошее время отклика для коротких процессов	Обеспечивает хорошее время отклика для коротких процессов	Обеспечивает хорошее время отклика	Обеспечивает хорошее время отклика	Не важна
Накладные расходы	Минимальны	Минимальны	Могут быть высокими	Могут быть высокими	Могут быть высокими	Могут быть высокими
Влияние на процессы	Плохо оказывается на коротких процессах и процессах с интенсивным вводом-выводом	Беспристрастна	Плохо оказывается на длинных процессах	Плохо оказывается на длинных процессах	Хороший баланс	Может привести к предпочтению процессов с интенсивным вводом-выводом
Голодание	Отсутствует	Отсутствует	Возможно	Возможно	Отсутствует	Возможно

52. сро. Feedback планировщик и классы планирования ОС UNIX SVR4.

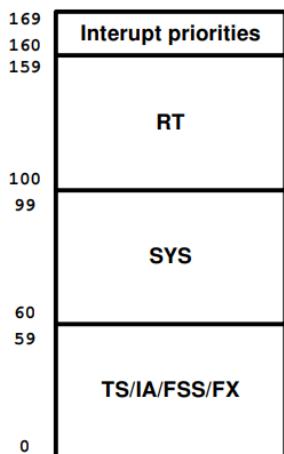
globpri	quantum	tqexp	slpriet	maxwait	lwait
59	2	49	59	32000000	59
58	4	48	58	0	59
57	4	47	58	0	59
		...			
40	4	30	55	0	55
		...			
30	8	20	53	0	53
29	12	19	52	0	52
		...			
21	12	11	52	0	52
20	12	10	52	0	52
19	16	9	51	0	51
		...			
12	16	2	51	0	51
11	16	1	51	0	51
10	16	0	51	0	51
9	20	0	50	0	50
		...			
1	20	0	50	0	50
0	20	0	50	0	50



globpri - глобальный приоритет,
 quantum - выделяемый квант времени (мс),
 tqexp - приоритет, если квант полностью выбран CPU,
 slpriet - приоритет, если ожидаем ресурса
 maxwait — граничное время длинного ожидания (мс)
 lwait - приоритет, в случае длинного ожидания ресурса



Классы планирования SVR4



- TimeSharing — разделение времени (Feedback)
- InterActive — интерактивный (boost к активному приложению)
- Fixed, System, RealTime — классы с фиксированными приоритетами
- Fair Share Scheduler — справедливый планировщик
- Отдельные приоритеты для Interrupt threads
- Учет афинити для NUMA

Класс приоритета	Глобальное значение	Последовательность планирования
Реальное время	159 • • • • 100	Первые
Ядро	99 • • 60	
Разделение времени	59 • • • • 0	Последние

Рис. 10.12. Очереди диспетчера SVR4

53. Справедливое планирование.

Планирование осуществляется исходя из приоритетов с учетом приоритета процесса, недавнего использования им процессора и недавнего использования процессора группой, к которой он принадлежит. Чем больше числовое значение приоритета, тем ниже сам приоритет. Для процесса j из группы k применимы следующие формулы:

$$CPU_j(i) = \frac{CPU_j(i-1)}{2}$$

$$GCPU_k(i) = \frac{GCPU_k(i-1)}{2}$$

$$P_j(i) = Base_j + \frac{CPU_j(i)}{2} + \frac{GCPU_k(i)}{4W_k}$$

где

$CPU_j(i)$ — мера загруженности процессора процессом j на интервале i ;

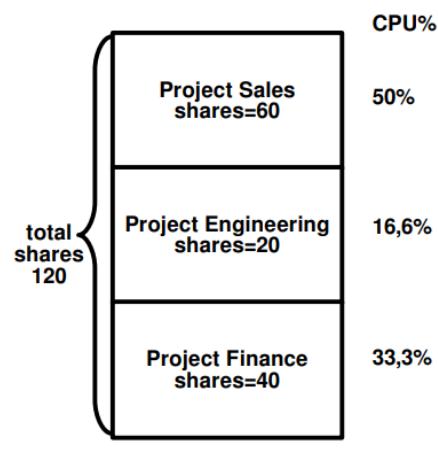
$GCPU_k(i)$ — мера загруженности процессора группой k на интервале i ;

$P_j(i)$ — приоритет процесса j в начале интервала i (меньшее значение соответствует большему приоритету);

$Base_j$ — базовый приоритет процесса j ;

W_k — вес, назначенный группе k ($0 < W_k \leq 1$ и $\sum_k W_k = 1$).

- Выделяет время процессора на основании заданных user shares
- Используется как замена TS/IA классов
- Учитывает все процессоры в системе



54. Планирование в многопроцессорных системах. Типы многопроцессорных систем с точки зрения организации планирования. Гранулярность и проектирование планировщиков процессов и потоков для многопроцессорных систем.



Типы многопроцессорных систем

- Слабосвязанные (распределенные, кластеры)
 - Набор систем со своей основной памятью и подсистемой ввода-вывода
 - Распределение заданий между системами
 - Распределение общих данных и результатов
- Функционально-специализированные
 - Ведущий процессор выполняет координацию
 - Ведомые процессоры выполняют вычисления
- Сильносвязанные процессоры
 - Имеют общую память и систему кэшей
 - Управляются одной ОС

- Гранулярность синхронизации — частота синхронизации между процессами в системе
- Fine (тонкая) — параллельные вычисления на уровне отдельных машинных команд
 - Менее 20 команд
- Medium (средняя) — на уровне одного приложения
 - 20-200 команд
- Coarse (грубая) — на уровне взаимодействующих процессов
 - 200-2000 команд
- Very Coarse (очень грубая) — на уровне распределенных систем
 - 2000-1 000 000
- Independent (независимая) — нет синхронизации, независимые процессы

Вопросы проектирования планировщиков в случае (сильно) многопроцессорных систем

- Назначение процессов процессорам
 - Статическое — выделяем процессор для процесса
 - Динамическое — общая очередь для всех процессов
 - Динамическая балансировка нагрузки
- Использование многозадачности на отдельных процессорах
 - Нужна ли для статического назначения многозадачность?
- Диспетчеризация процесса
 - Так ли плох будет FCFS? Влияние выбора алгоритма диспетчеризации снижается.

Подходы к планированию потоков

- Load Sharing — глобальная очередь потоков
 - Равномерное распределение нагрузки
 - Нет централизованного планировщика
 - Минусы — блокировки центральной очереди, промахи кэшей, неэффективность при тонкой средней гранулярности
- Gang scheduling — связанные потоки распределяются на связанные процессы по одному на процессор
 - Снижение накладных расходов на планирование при тонкой и средней гранулярности
- Dedicated processor assignment — назначается пул процессоров равный количеству потоков
 - В экстремальном случае увеличивает простой процессора
 - Полное устранение переключений повышает скорость работы
- Динамическое планирование
 - В отсутствии свободных CPU ресурсы изымаются из процесса, использующего несколько CPU

55. ОС реального времени и планировщики. Deadline-планирование.

Планирование реального времени

Планирование реального времени является одной из областей, в которых ведутся активные исследовательские работы. В этом подразделе мы приведем краткий обзор различных подходов к проблеме планирования реального времени и рассмотрим два распространенных класса алгоритмов планирования.

В обзоре алгоритмов планирования реального времени [196] замечено, что различные алгоритмы планирования зависят от того, 1) выполняет ли система анализ планируемости, и если выполняет, то 2) как именно, статически или динамически, и 3) к чему он приводит — к непосредственному выполнению функции планирования или же к построению расписания, согласно которому в процессе работы осуществляется диспетчeringизация заданий. На основе этого автор определяет следующие классы алгоритмов.

- **Статическое планирование с использованием таблиц.** При этом выполняется статический анализ осуществимости планирования; результатом анализа является план, который в процессе работы системы определяет, когда должно начаться выполнение заданий.
- **Статическое вытесняющее планирование на основе приоритетов.** В этом случае также выполняется статический анализ, но расписание не создается. Вместо этого на основе проведенного анализа заданиям назначаются приоритеты, с тем чтобы далее можно было использовать традиционный вытесняющий планировщик, работающий с учетом приоритетов заданий.
- **Динамическое планирование на основе расписания.** Осуществимость планирования определяется не статически, а динамически, в процессе выполнения.

Поступающее в систему задание принимается только в том случае, если определена возможность его выполнения с учетом всех временных требований. Одним из результатов анализа является расписание, используемое для принятия решения о диспетчеризации задания.

- **Динамическое планирование наилучшего результата.** При этом анализ осуществимости планирования не выполняется; система пытается удовлетворить все предельные сроки и снимает те выполняющиеся процессы, предельные сроки которых нарушены.

Статическое планирование с использованием таблиц применимо для периодических заданий. Входной информацией для анализа являются время поступления заданий в систему, время выполнения, предельные сроки выполнения и относительный приоритет каждого задания. Планировщик пытается разработать такой план работы, который удовлетворял бы всем временным требованиям заданий. Такой подход является предсказуемым, но абсолютно не гибким, поскольку любое изменение требований любого задания приводит к необходимости пересмотра всего расписания. Типичными представителями этой категории алгоритмов планирования являются планирование наиболее раннего предельного срока и другие, рассматриваемые далее, алгоритмы планирования периодических заданий.

Статическое вытесняющее планирование на основе приоритетов использует тот же вытесняющий механизм планирования с приоритетами, что и большинство обычных многозадачных операционных систем. В таких системах для определения приоритетов могут использоваться самые различные факторы. Например, в системе с разделением времени приоритет процесса изменяется в зависимости от того, на что ориентирован этот процесс — на вычисления или на операции ввода-вывода. В системах реального времени назначение приоритетов связано с временными ограничениями каждого из заданий. Примером такого подхода может служить частотно-монотонное планирование, рассматриваемое ниже в данной главе, которое назначает заданиям статические приоритеты на основе информации об их периодах.

В случае использования **динамического планирования на основе расписания** после поступления задания в систему (но до начала его выполнения) предпринимается попытка создать расписание, которое содержит как все имеющиеся в системе задания, так и вновь поступившее. Если удается создать такое расписание, что при выполнении удовлетворяются временные ограничения как нового задания, так и уже имеющихся в системе заданий, оно принимается планировщиком.

Динамическое планирование наилучшего результата используется во многих современных коммерческих системах реального времени. При поступлении нового задания система назначает ему приоритет на основе характеристик этого задания. При этом подходе обычно используется планирование, учитывающее предельные сроки, — наподобие планирования наиболее раннего предельного срока. Как правило, поступающие задания непериодические, а потому статический анализ планирования неприменим. При таком типе планирования мы не знаем, будут ли удовлетворены временные ограничения задания до тех пор, пока задание не будет полностью выполнено (или пока не будут нарушены временные ограничения). Именно это и является основным недостатком данной схемы; преимущество же динамического планирования наилучшего результата — в простоте реализации.

На рис. 10.8 показан простой пример RMS. Экземпляры заданий пронумерованы последовательно. Как можно видеть, для задания 3 второй экземпляр не выполняется, поскольку предельное время просрочено. Третий экземпляр испытывает вытеснение, но все равно остается в состоянии завершиться до предельного срока завершения.

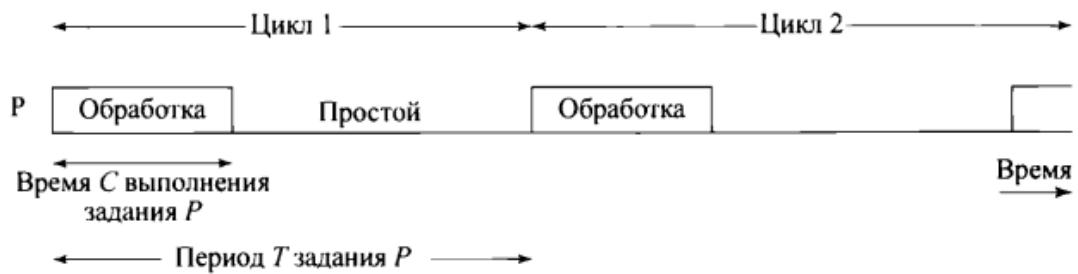


Рис. 10.7. Временная диаграмма периодического задания

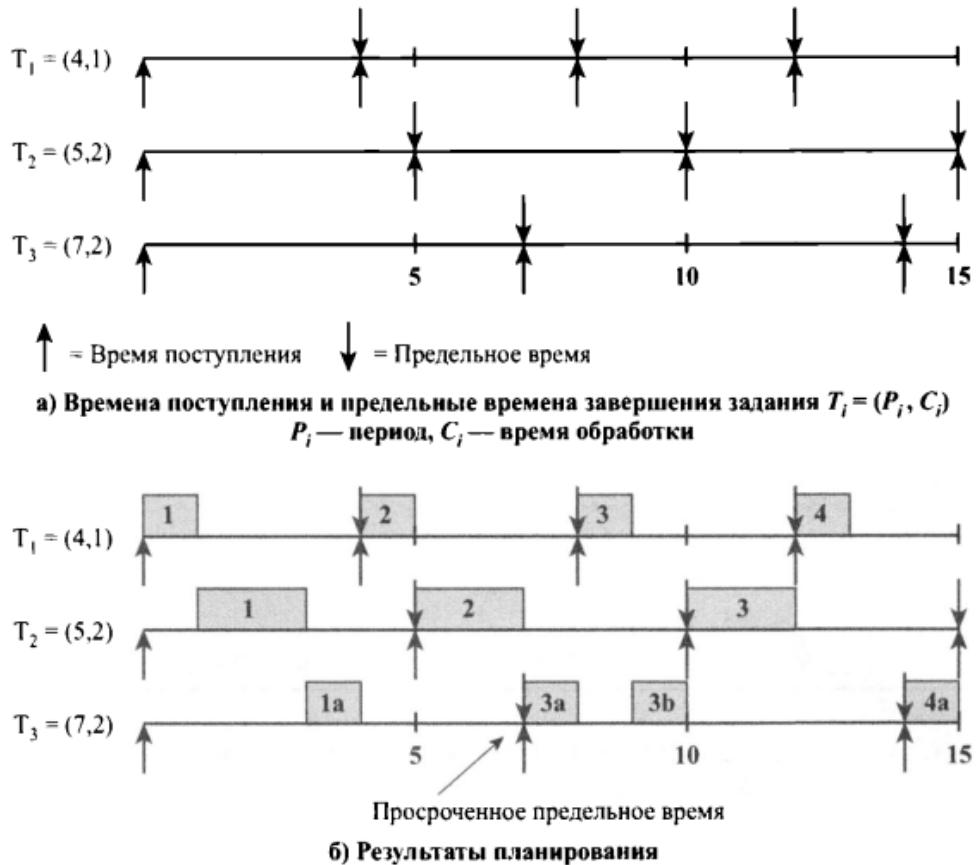


Рис. 10.8. Пример работы RMS

Одной из характеристик эффективности алгоритма периодического планирования является его гарантия соответствия всем жестким предельным срокам. Предположим, что у нас имеется n заданий, каждое из которых имеет свое фиксированное время выполнения и период. Тогда необходимым условием соответствия всем жестким предельным срокам является выполнение следующего неравенства:

$$\frac{C_1}{T_1} + \frac{C_2}{T_2} + \dots + \frac{C_n}{T_n} \leq 1 \quad (10.1)$$

Сумма загруженности процессора разными заданиями не может превышать 1, что соответствует полной загрузке процессора. Неравенство (10.1) определяет верхнюю границу количества заданий, которые может успешно обслуживать идеальный алгоритм планирования. Для конкретного реального алгоритма граница может оказаться ниже. Так, можно показать, что для алгоритма RMS справедливо следующее неравенство:

$$\frac{C_1}{T_1} + \frac{C_2}{T_2} + \dots + \frac{C_n}{T_n} \leq n(2^{1/n} - 1) \quad (10.2)$$

В табл. 10.5 приведены некоторые значения верхней границы для метода RMS. По мере возрастания количества заданий верхняя граница стремится к значению $\ln 2 \approx 0,693$.

Таблица 10.5. ЗНАЧЕНИЯ ВЕРХНЕЙ ГРАНИЦЫ ЗАГРУЖЕННОСТИ ДЛЯ МЕТОДА RMS

n	$n(2^{1/n} - 1)$
1	1,000
2	0,828
3	0,779
4	0,756
5	0,743
6	0,734
•	•
•	•
•	•
∞	$\ln 2 \approx 0,693$

В качестве примера рассмотрим три периодических задания (здесь $U_i = C_i/T_i$).

- Задание P_1 : $C_1 = 20$; $T_1 = 100$; $U_1 = 0,2$.
- Задание P_2 : $C_2 = 40$; $T_2 = 150$; $U_2 = 0,267$.
- Задание P_3 : $C_3 = 100$; $T_3 = 350$; $U_3 = 0,286$.

Общая загруженность процессора этими тремя заданиями составляет

$$0,2 + 0,267 + 0,286 = 0,753.$$

Верхняя граница загруженности этих трех задач при использовании метода RMS составляет

$$\frac{C_1}{T_1} + \frac{C_2}{T_2} + \frac{C_3}{T_3} \leq 3(2^{1/3} - 1) = 0,779.$$

Поскольку общая загруженность процессора по обработке приведенных заданий ниже верхней границы для метода RMS, можно сделать вывод,

что при RMS-планировании будут успешно выполнены все задания. Можно также показать, что верхняя граница из (10.1) справедлива для метода наиболее раннего предельного срока. Таким образом, при применении планирования с наиболее ранним предельным сроком можно достичь более высокой загрузки процессора и, соответственно, обработать большее количество заданий. Тем не менее метод RMS широко распространен и используется во многих промышленных приложениях. В работе [221] это поясняется следующими причинами.

1. На практике отличие производительности невелико. Кроме того, неравенство (10.2) консервативно, и на практике зачастую достигается 90%-ная загрузка процессора.

2. Большинство жестких систем реального времени содержат мягкие компоненты, такие как некритичный вывод на экран или встроенное самотестирование, выполняющееся с низким приоритетом. Эти компоненты используют процессорное время, которое остается после RMS-планирования жестких заданий.

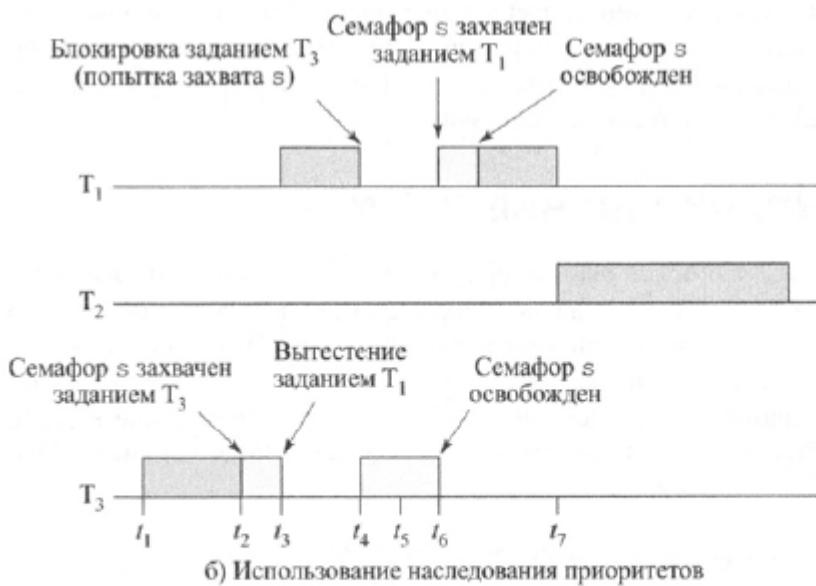
3. При использовании RMS проще обеспечить стабильность. Когда система не в состоянии обеспечить все предельные сроки в силу перегруженности или временных ошибок, необходимо гарантировать выполнение предельных сроков для подмножества обязательных заданий. При статическом назначении приоритетов гарантируется только корректность выполнения основных задач с относительно высокими приоритетами. Это может быть достигнуто и в RMS-планировании путем реструктурирования обязательных задач для повышения их частоты либо посредством изменения приоритетов RMS для подмножества обязательных задач. При планировании с наиболее ранним предельным сроком приоритеты периодических заданий изменяются от одного периода к другому, что усложняет обеспечение корректной работы обязательных заданий с жесткими предельными сроками.



Deadline планирование

- Важны своевременные завершение или начало выполнения задания
 - Конфликты, сбои и временные недостатки ресурсов не должны влиять
- Задания могут включать дополнительную информацию:
 - Ready time — время готовности задания к выполнению
 - Starting deadline — предельное время начала выполнения
 - Completion deadline — предельное время полного завершения задания
 - Processing time — время, необходимое для полного выполнения задания
 - Resource requirements — список ресурсов (не процессор) для задания
 - Priority — мера важности задания
 - Subtask structure — обязательные и необязательные задачи
- Rate Monotonic Scheduling — см. Столлингс гл. 10.2

56. Проблема инверсии приоритетов, типы инверсии и способы решения в планировщике.



■ Нормальное выполнение □ Выполнение критического участка кода

Рис. 10.9. Инверсия приоритетов

Основная идея **наследования приоритетов** (priority inheritance) заключается в том, что задача с более низким приоритетом наследует приоритет любой высокоприоритетной задачи, ожидающей совместно используемый ресурс. Это изменение приоритета происходит, как только более высокоприоритетное задание блокируется в ожидании ресурса; оно должно заканчиваться, когда ресурс освобождается задачей с более низким приоритетом. На рис. 10.9, б показано, что наследование приоритетов решает проблему неограниченной инверсии приоритетов, показанной на рис. 10.9, а. Соответствующий порядок событий оказывается следующим.

- t_1 : T_3 начинает выполняться.
- t_2 : T_3 захватывает семафор s и входит в критический участок кода.
- t_3 : T_1 , имеющий более высокий приоритет, чем приоритет T_3 , вытесняет T_3 и начинает выполнять.
- t_4 : T_1 пытается войти в критический участок кода, но блокируется, поскольку семафор уже захвачен T_3 . T_3 немедленно (временно) получает тот же приоритет, что и T_1 . T_3 продолжает выполнение критического участка.
- t_5 : T_2 готов к выполнению, но поскольку T_3 теперь имеет более высокий приоритет, T_2 не в состоянии вытеснить T_3 .
- t_6 : T_3 покидает критический участок и разблокирует семафор: уровень его приоритета опускается до бывшего у него ранее приоритета по умолчанию. T_1 вытесняет T_3 , блокирует семафор и входит в критический участок.
- t_7 : T_1 приостанавливается по причине, не связанной с T_2 , и T_2 начинает выполнение.

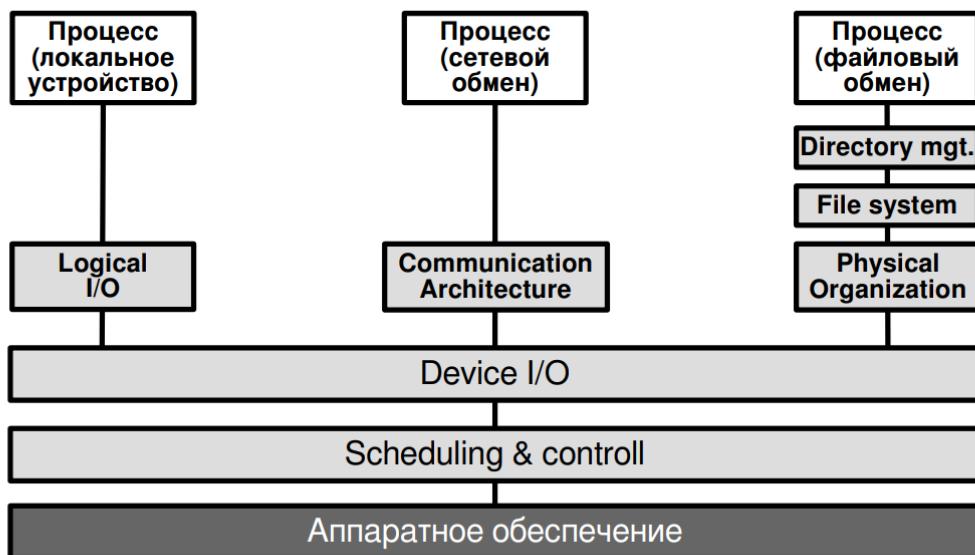
Именно этот подход был применен для решения проблемы Pathfinder.

В подходе с **потолком приоритета** с каждым ресурсом связан приоритет. Приоритет, назначенный ресурсу, на один уровень выше приоритета его наиболее приоритетного пользователя. Затем планировщик динамически назначает этот приоритет любой задаче, которая обращается к ресурсу. Когда задание заканчивает работу с ресурсом, его приоритет возвращается в нормальное состояние.

.57. Ввод-вывод. Современные устройства и скорости обмена, развитие способов ввода-вывода, логическая структура ввода-вывода.



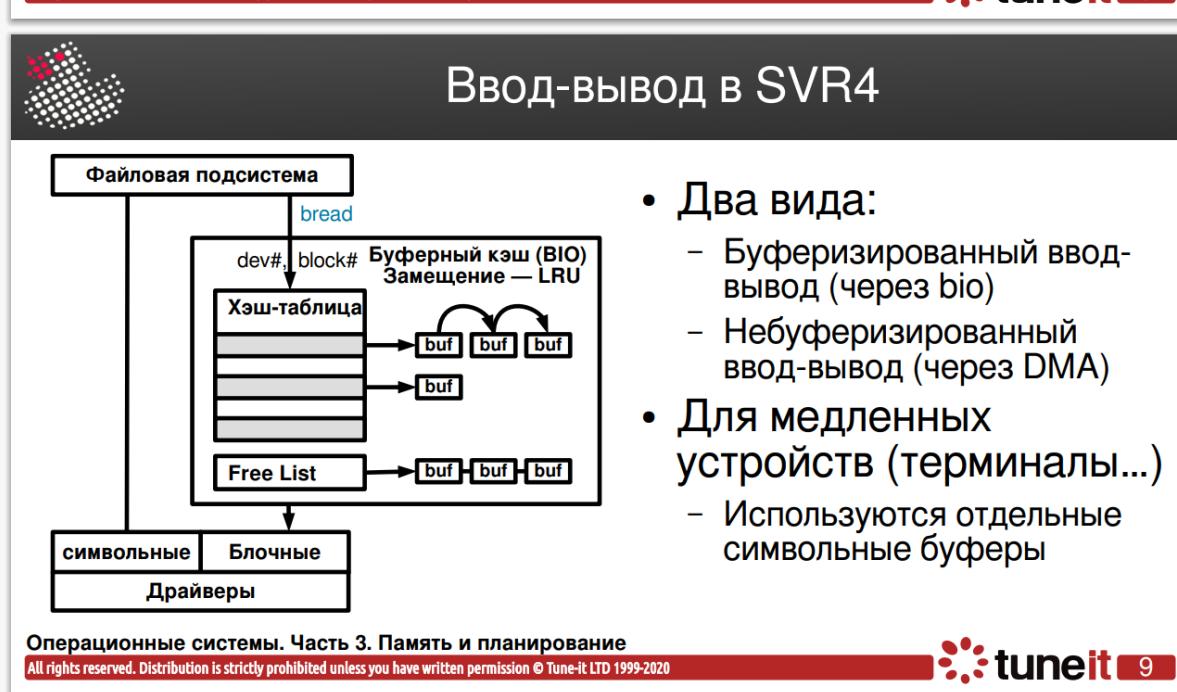
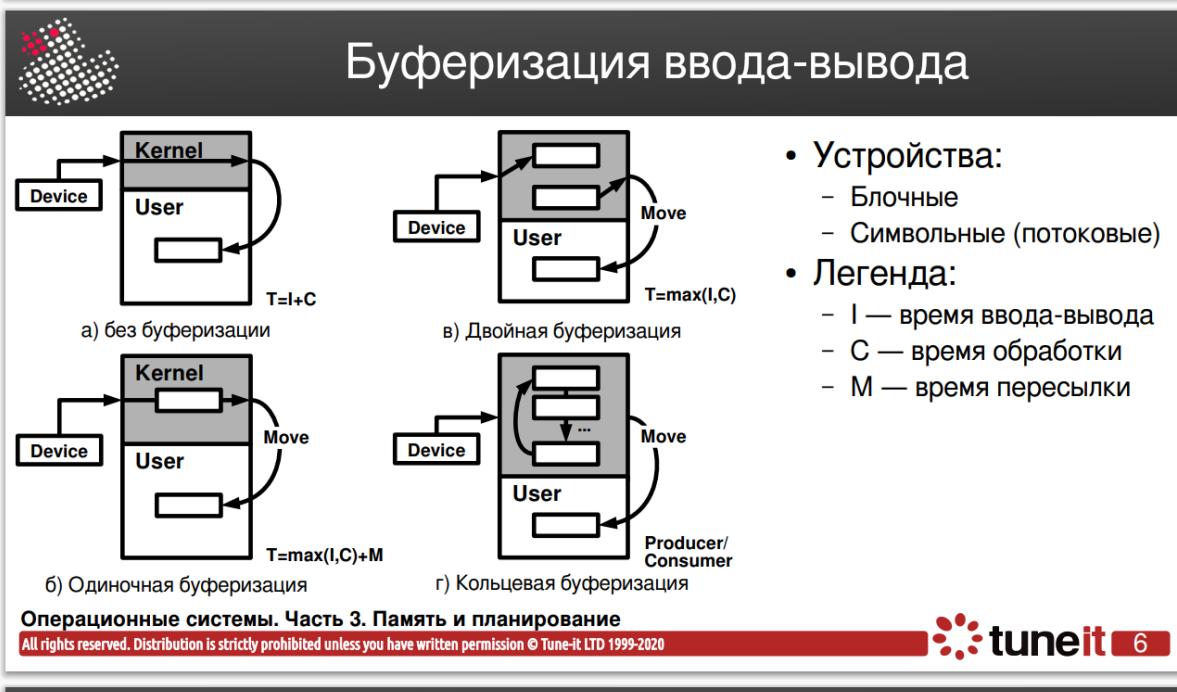
Логическая структура ввода-вывода



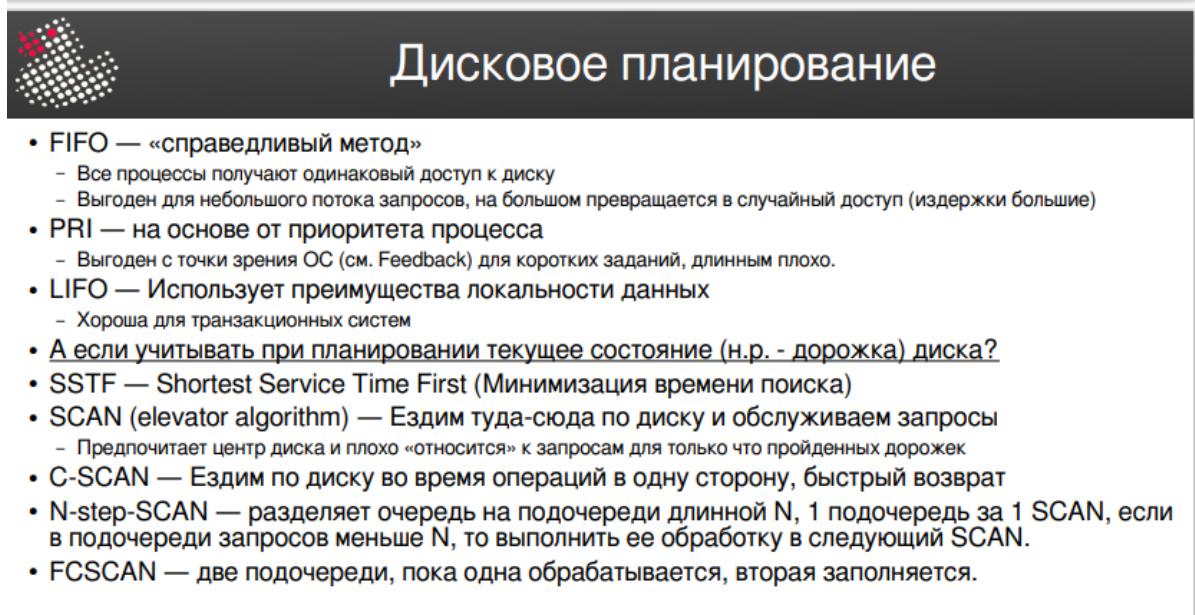
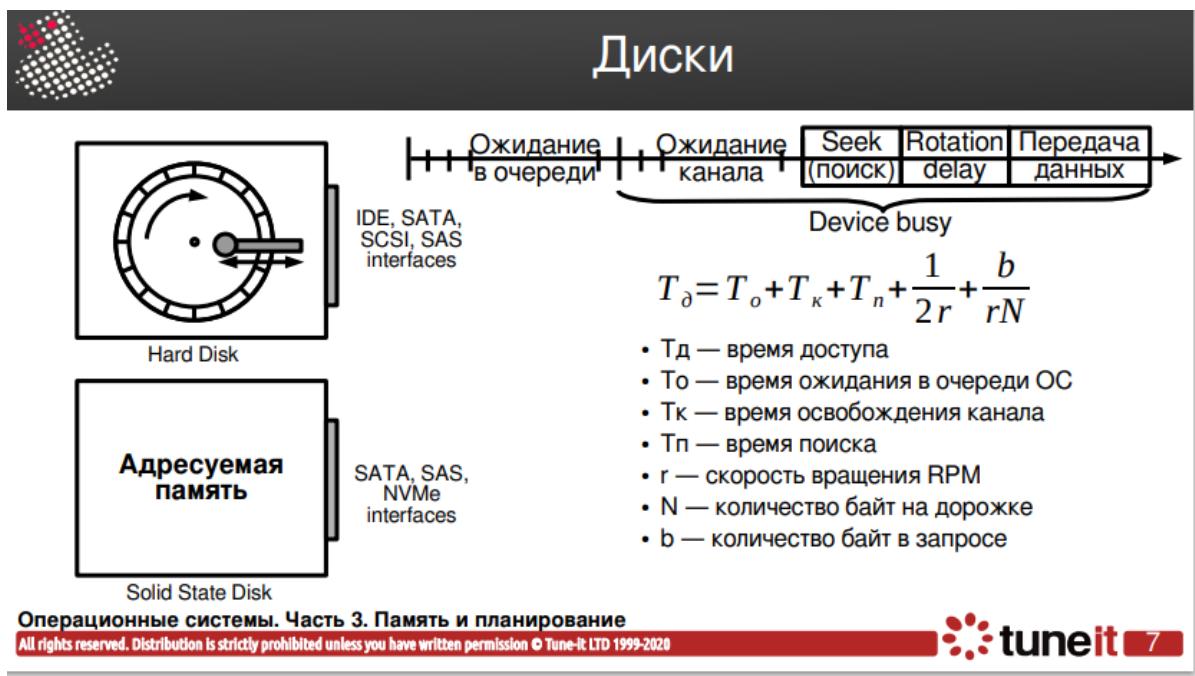
Ионные системы. Часть 3. Память и планирование

- **Логический ввод-вывод.** Модуль логического ввода-вывода обращается с устройством как с логическим ресурсом и не обращает внимания на детали фактического управления устройством. Логический модуль ввода-вывода работает посредником между пользовательскими процессами и устройством, позволяя им работать с последним с использованием идентификатора устройства и простых команд, таких как открытие, закрытие, чтение и запись.
- **Устройство ввода-вывода.** Запрошенные операции и данные (буферизированные символы, записи и т.п.) конвертируются в соответствующие последовательности инструкций ввода-вывода, команд управления каналом и команд контроллера. Для более эффективного использования устройства может быть применена буферизация.
- **Планирование и контроль.** На этом уровне происходят реальная организация очередей и планирование операций ввода-вывода, а также управление выполнением операций. Таким образом, на этом уровне осуществляются работа с прерываниями, получение и передача информации о состоянии устройства. Это уровень программного обеспечения, которое непосредственно взаимодействует с контроллером ввода-вывода, а следовательно, с аппаратным обеспечением устройства.

58. Буферизация ввода вывода. Ввод-вывод в UNIX SVR4.



59. Диски и дисковое планирование.



60. Концепции RAID.

(Redundant Array of Independent Disks - избыточный массив независимых дисков).

1. RAID - это набор физических дисков, рассматриваемых операционной системой как единый логический диск.
2. Данные распределены по физическим дискам массива.
3. Избыточная емкость дисков используется для хранения контрольной информации, гарантирующей восстановление данных в случае отказа одного из дисков.

Благодаря большому количеству дисков повышается производительность, но увеличивается вероятность сбоя. Именно поэтому RAID предусматривает хранение информации для восстановления данных.

Несколько дисков вместо одного используются потому, что производительность вторичных запоминающих устройств растет медленнее, чем других компонентов, а несколько дисков позволяют выполнять несколько операций ввода-вывода одновременно.

Из семи уровней RAID чаще всего используются 0, 1, 5 и 6.

RAID 2 и 3 не используются из-за того, что они призваны использоваться в среде с большим количеством ошибок, в то время как диски обладают высокой надежностью.

..61. RAID-0, 1, 10, 0+1.

0 не является настоящим RAID - уровнем, поскольку никак не пользуется избыточностью. Основным плюсом RAID- 0 является очень высокая скорость ввода-вывода. Она достигается за счет того, что данные расщепляются(striped) и записываются на все доступные диски сразу(по схеме 1 дорожка 1 диска,1 дорожка 2 диска,1 дорожка 3 диска,2 дорожка 1 диска,2 дорожка 2 диска и т.д.)



a) RAID0 (без избыточности)

Из-за такой записи повреждение одного диска приводит к повреждению всех данных.

1 зеркалирует данные (mirroring, создание зеркала).

Избыточность достигается не путем вычислений, а просто через дублирование данных. При этом используется такая же схема, как и при RAID-0.



б) RAID 1 (отражение)

В качестве плюсов можно выделить:

простоту восстановления данных,

возможность выбирать между дисками (при обращении выбирается тот, у которого минимальная задержка из-за вращения и минимальное время поиска).

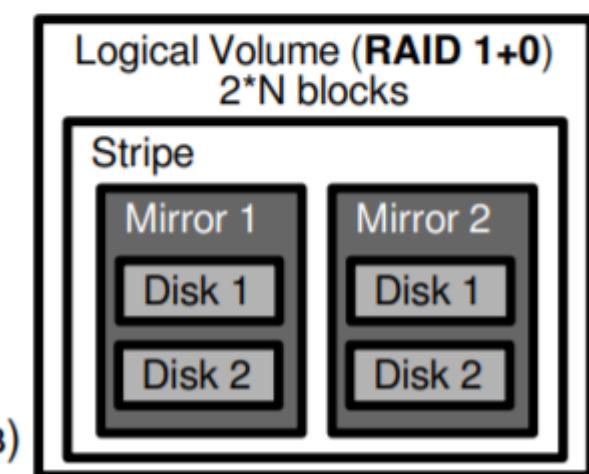
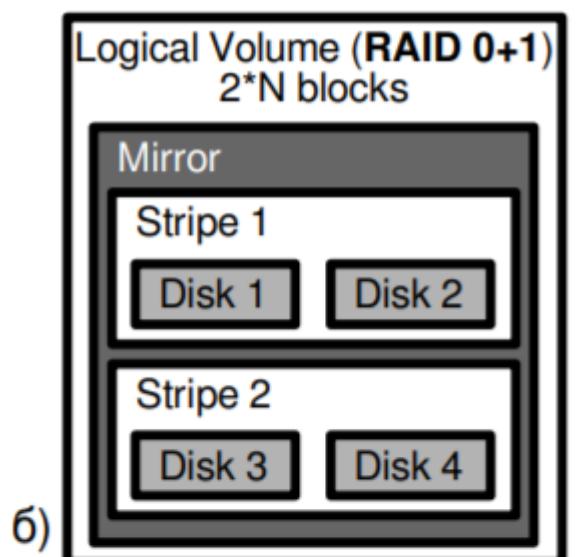
запись можно обновлять параллельно на обоих дисках(таким образом время определяется скоростью более медленной операции)

Минус - требуется в два раза больше дисков, поэтому RAID-1 используется только для самых важных файлов.

RAID-1 лучше подходит для чтения и может превышать скорость RAID-0 в два раза в среде, ориентированной на транзакции.

RAID - 1+0 =RAID 10

Массив RAID - 0, собранный из RAID-1



62. RAID 4,5,6. Аппаратные дисковые массивы.

RAID -4 также почти не используется.

- Подсчитывает четность для SU
 - $P(0,1) = SU(0) \oplus SU(1)$
 - Если сбой диска 2 SU(1)= $P(0,1) \oplus SU(0)$
- Характеристики (a):
 - Высокая надежность, избыточность Disk+1
 - Дополнительные вычисления при записи, и при чтении в случае сбоя диска
 - Пропускная способность на чтение как у RAID 0 и меньше, чем одинарная на запись
 - Скорость обработки запросов для чтения как у RAID 0 и меньше, чем одинарная на запись
- RAID 4 — выделенный диск с четностью
- RAID 6 — двойная размазанная четность, по разным алгоритмам, избыточность Disk+2

При уровне RAID 5 может быть восстановлен только один диск, при уровне RAID-6 - 2.



63. Файловый ввод-вывод, основные определения. Задачи ОС по управлению файлами. Совместное использование файлов.

Файл это коллекция данных со следующими свойствами:

- Наличие структуры (может быть сложной, может быть несколько в одном файле)
- Возможность долгосрочного существования
- Возможность совместного использования процессами
- Основные файловые операции
- Создание, удаление, открытие, закрытие, чтение, запись, выбор позиции, контрольные операции, блокирование

Поле — одиночный элемент записи, обычно одно значение.

Запись — набор полей

Файл — совокупность записей, относящихся к однородному набору данных.

База данных — файл со сложной взаимозависимой структурой полей и записей.

Операции, проводимые с файлом, могут влиять на его структуру

Файл может не иметь структуры (например, представлять собой поток символов)

Задачи ОС по управлению файлами:

- Возможность хранить файлы и выполнять пользователю над ними операции: не
 - Создание, удаление, чтение и изменения файлов
 - Управление доступом к файлам для себя и других пользователей
 - Перемещение данных между файлами
 - Выполнение резервного копирования и восстановления файлов
 - Обеспечение возможности работы с файлами по именам, удобным для пользователя
- Гарантия корректности данных
- Обеспечение приемлемой производительности
- Поддержка различных типов устройств хранения
- Минимизация или исключения потерь и повреждения данных
- Обеспечение базового набора функций ввода-вывода
- Обеспечение совместного использования файлов



Совместное использование файлов

- Владелец — пользователь, имеющий все права на файл, и распределяющий эти права:
 - Отсутствие прав — другие пользователи не знают о существовании файла
 - Знание — пользователи знают о существовании файла, и могут запросить права у владельца
 - Выполнение — пользователь может загрузить и выполнить программу
 - Чтение, добавление в конец файла, модификация, удаление
 - Изменение прав доступа
- Доступ к файлу может быть передан разным наборам пользователей:
 - Конкретному пользователю
 - Списку пользователей или группе пользователей
 - Всем пользователям
- Блокировка одновременного доступа
 - Всего файла
 - Отдельных записей

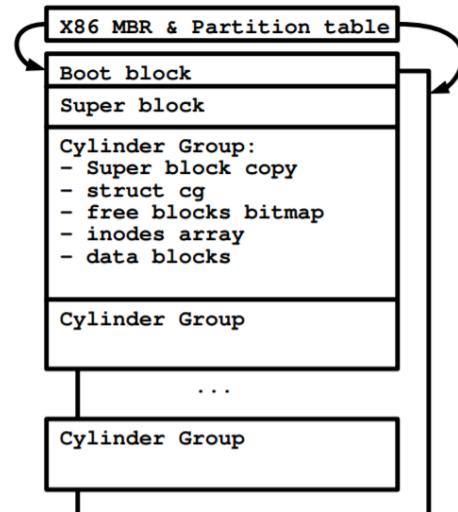
64. Управление файлами в UNIX SVR4



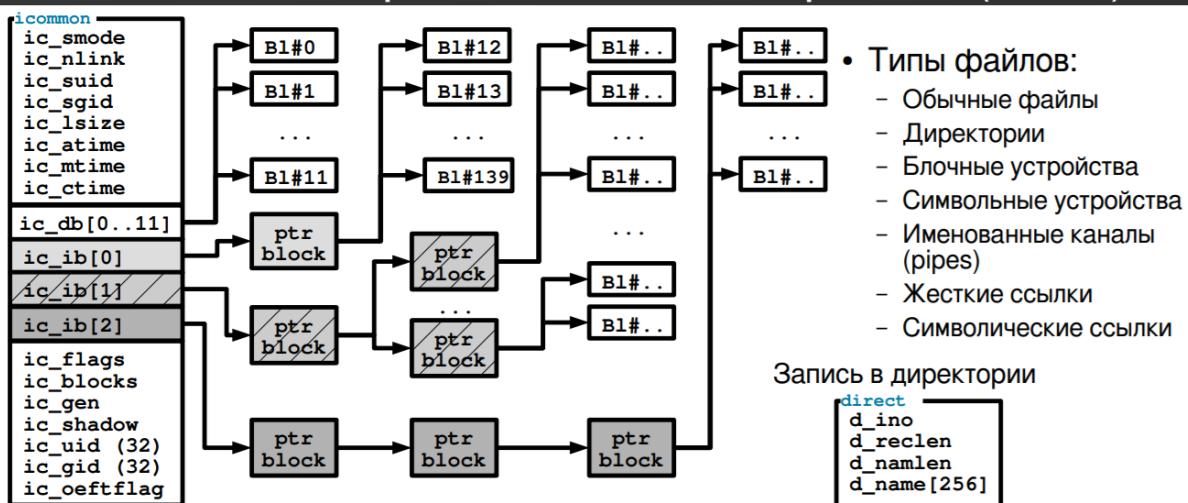
Структура файловой системы UNIX SVR4 UFS

- Файловая система содержит:

- Загрузчик ОС
- Суперблок — геометрия и служебные параметры файловой системы
- Цилиндровые группы — информация о свободных блоках, массив записей о файлах (inode), блоки данных — равномерно распределены по всему дисковому пространству



Файловая система UNIX SRV4. Индексированные записи файлов (inode)



65. Каталоги файлов. Элементы каталога, операции ОС.

Look



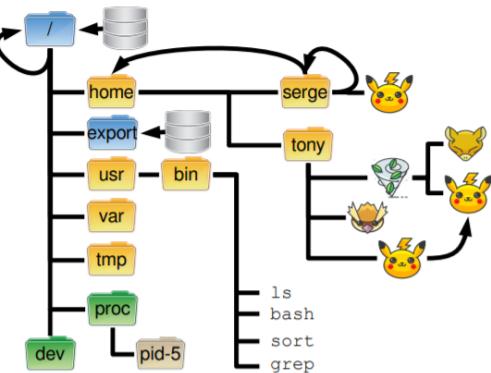
Каталоги файлов. Элементы каталога.

- Основные:
 - Имя файла, Тип файла (бинарный, текстовый, ...)
 - Организация файла — организация внутренней структуры
- Адресная информация
 - Том (носитель) — указатель на физическое устройство
 - Начальный адрес — номер первого блока файла
 - Занимаемый размер и зарезервированный размер— количество байт в файле и максимальное количество байт
- Информация об управлении доступом
 - Владелец — пользователь, который может управлять файлом
 - Разрешенные действия (чтение, запись, исполнение, поиск, и т.п.)
- Информация об использовании
 - Создатель и дата создания — создатель не обязательно может быть владельцем
 - Дата последнего чтения и последний пользователь, прочитавший файл
 - Дата последнего изменения и последний пользователь, изменивший файл
 - Дата последней резервной копии на другом устройстве
 - Текущее использование — информация о текущих действиях с файлом



Каталог, операции ОС

- Часть элементов каталога может находиться в записи о файле, уменьшая его размер
- Каталог может быть полностью или частично загружен в основную память
- Простейшая структура — список записей фиксированной длины
- Основные операции:
 - Поиск файловых записей,
 - Создание и удаление записей о файле
 - Получение списка записей
 - Обновление каталога, изменение записи
- Обычно структура каталогов — n-уровневая иерархия, дерево, связанный граф



• Именование:

- Имя файла
- Имя, включающее полный путь от корня
- Имя, относительно текущего каталога

Существует 2 типа фрагментации: Внутренняя: когда выделяется больше памяти, чем запрашивалось, избыток памяти не используется; Внешняя: свободная память в процессе выделения или освобождения разделяется на мелкие блоки и в результате не обслуживаются некоторые запросы на выделение памяти.

66. Размещение записей и файлов в блоках данных.

Сложность и типы организации размещения.



- Запись — логическая единица доступа к структурированному файлу
 - ОС осуществляет ввод-вывод блоками
- Какой должен быть размер блок относительно размера записи?
 - Блок больше → Больше записей передано за одну операцию ввода-вывода
 - При случайному доступе — передаем записи, которые не используются
 - Большим блокам нужны большие буферы ОС

Чем больше блок - тем больше объем данных, переданный за одну операцию

ввода-вывода => тем меньше нагрузка на ОС

Обычно размер блока кратен размеру страницы

Размещение файлов

Номер блока на дорожке							
Дорожка	0	1	2	3	4	5	6
0							
1							
2							
3							
4							
5							
6							
7							
N							

- Как на физической структуре диска, состоящего из блоков разместить файлы?

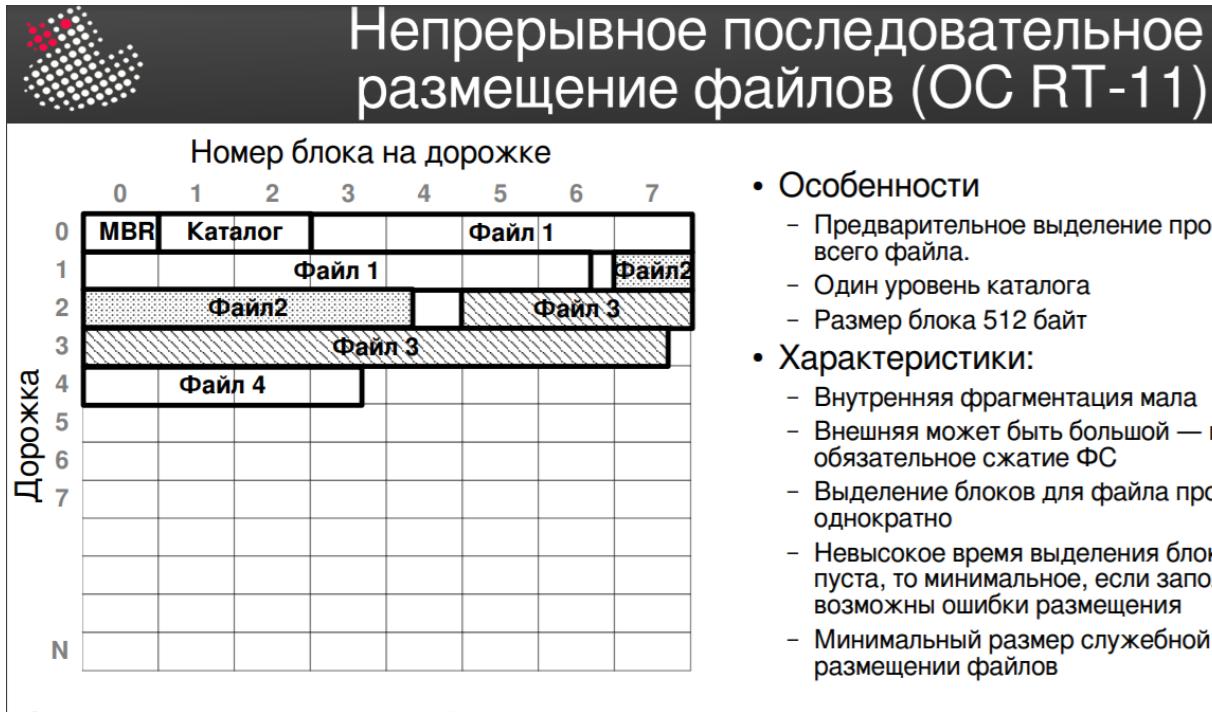
- Предварительное выделение пространства для всего файла?
- Размер порции файла?
- Фиксированный или переменный размер порции?

- Решения влияют на характеристики:

- Внутреннюю и внешнюю фрагментацию
- Частота выделения (allocation) блоков для всей файловой системы
- Время выделения блоков
- Размер служебной информации о размещении файлов

Больше про размещение в следующих билетах.

67. Непрерывное размещение файлов (на примере ОС RT-11)



- Особенности

- Предварительное выделение пространства для всего файла.
- Один уровень каталога
- Размер блока 512 байт

- Характеристики:

- Внутренняя фрагментация мала
- Внешняя может быть большой — необходимо обязательное сжатие ФС
- Выделение блоков для файла производится однократно
- Невысокое время выделения блоков — если ФС пуста, то минимальное, если заполнена, то возможны ошибки размещения
- Минимальный размер служебной информации о размещении файлов

Плюс - быстрое последовательное чтение файлов

Минус - дыры после удаления файлов

68. Цепочечное размещение файлов (на примере DOS FAT)



- Особенности

- Выделение пространства для файла по мере необходимости
- Размер порции файла — кластер (C) — от 512 (FAT12) до 32Кб (FAT 32)
- Запись каталога содержит имя (8+3), размер файла, ..., номер начального кластера: File1.txt
- FAT (File Allocation Table) — Содержит таблицу аллокации цепочки кластеров

CL#	10	16	17	27	38	39	39
Next CL#	FFFF	17	27	38	39	39	10

- Характеристики:

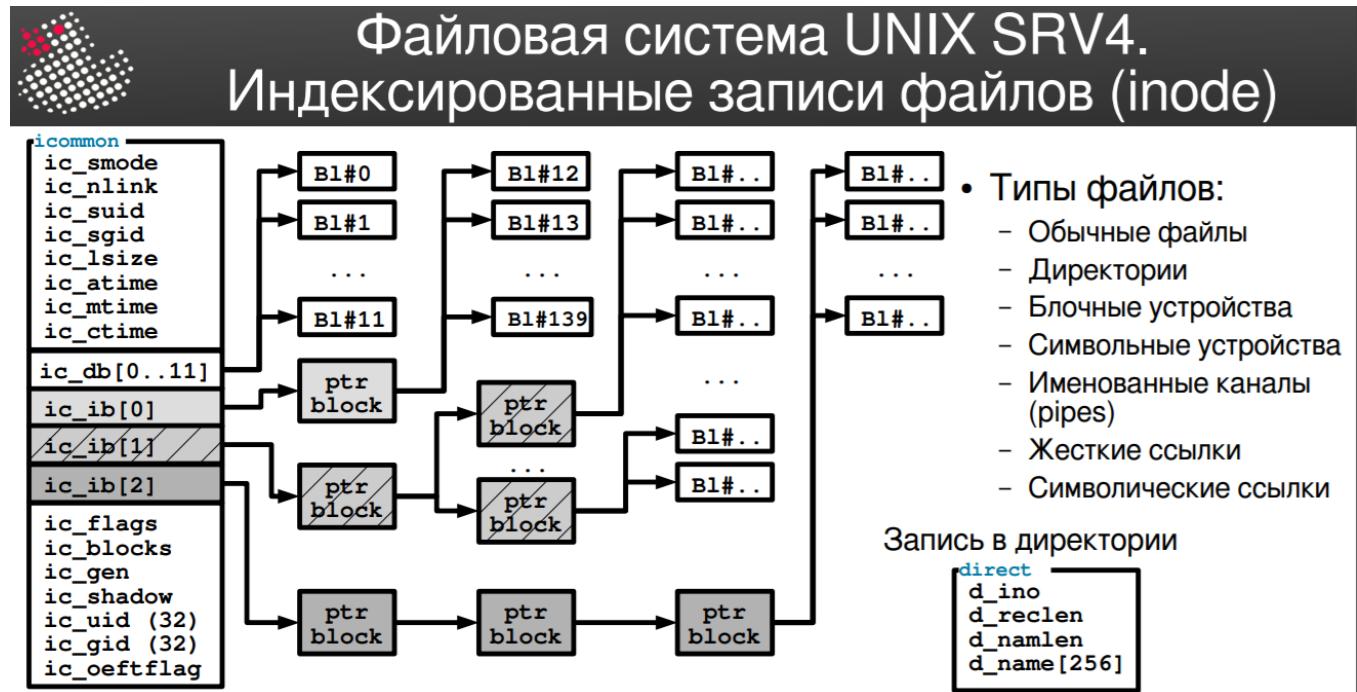
- Внутренняя фрагментация мала
- Внешняя может быть большой — файлы размещены непоследовательно, необходимо обязательное дефрагментирование ФС для увеличения скорости
- Высокое время выделения блоков файла
- Размер служебной информации о размещении файлов ограничен FAT

Файл разбивается на кластеры одинакового размера, которые затем размещаются в памяти.

Таблица указывает на следующий кластер для каждого кластера.

Удобно заменять файлы, поскольку можно заменять кластеры последовательно.

69. Индексированное размещение (на примере файловой системы UNIX UFS).



ic_db - где смотреть запись

ic_id[n] - указатели на блоки с разной косвенностью, итоговый пункт назначения-блок с данными.

Индексированная запись - потому что указывает индексом на блок, который нужно использовать.

Дополнительные вопросы для кафедры ВТ.

70. Linux: стандартные средства для наблюдения счетчиков ядра.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 69

Стандартные средства для наблюдения счетчиков ядра

- sar (system activity reporter): общесистемное средство, собирающая статистику по пейджингу (-B) и свопингу (-W), вводу-выводу (-b,-d), смонтированным системам (-F), прерываниям (-l), управлению питанием (-m), сети (-n), процессорам (-P,-u), очереди процессов и загрузке (-q,-w), памяти (-r), области подкачки (-s), терминалам (-y)
- Можно настроить сбор исторических результатов (crontab)
- Пример: sar -q 1 1 (одно измерение за одну секунду)

Linux 5.4.0-47-generic (ra) 01.10.2020 _x86_64_ (4 CPU)
14:35:36 runq-sz plist-sz ldavg-1 ldavg-5 ldavg-15 blocked
14:35:37 0 1034 1,44 1,70 1,75 0
Average: 0 1034 1,44 1,70 1,75 0

Операционные системы. Часть 1.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 70



Стандартные средства для наблюдения счетчиков ядра (2)

- Процессор: ps, top, tiptop, turbostat, rdmsr, numastat, uptime
- Виртуальная память: vmstat, slabtop, pidstat, free
- Дисковая подсистема: iostat, iotop, blktrace
- Сеть: netstat, tcpdump, iptraf, ethtool, nicstat, ip
- Интерактивные (типа top) или с указанием количества запуска и интервала (типа sar)
- Некоторые работают только с правами root!

Операционные системы. Часть 1.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 71

71. Linux: файловая система /proc.



/proc

- Виртуальная файловая система, содержащая файлы статистики и управляющая модулями ядра
- find /proc |wc -l
 - 398401(over 9000!)
- Для начала:
 - /proc/cpuinfo - информация о процессоре (модель, семейство, размер кэша и т.д.)
 - /proc/meminfo - информация о памяти, размере области подкачки и т.д.
 - /proc/mounts - список смонтированных файловых систем.
 - /proc/devices - список устройств.
 - /proc/filesystems - поддерживаемые файловые системы.
 - /proc/modules - список загружаемых модулей.
 - /proc/version - версия ядра.
 - /proc/cmdline - список параметров, передаваемых ядру при загрузке.

Более подробно смотри:

kernel.org -- ver-- Documentation/filesystems/proc.rst

Операционные системы. Часть 1.

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



/proc (2)

- Можно получить много полезной информации о выполняющихся процессах

```
serge@ra:~$ echo $$
13318
serge@ra:~$ cd /proc/13318
serge@ra:/proc/13318$ ls -F
arch_status      cpuset    loginuid    numa_maps      sched      status
attr/            cwd@     map_files/   oom_adj       schedstat   syscall
autogroup        environ   maps         oom_score    sessionid  task/
auxv             exe@     mem          oom_score_adj setgroups   timers
cgroup           fd/      mountinfo   pagemap      smaps      timersslack_ns
clear_refs       fdinfo/   mounts      patch_state  smaps_rollup uid_map
cmdline          gid_map   mountstats personality  stack      wchan
comm              io       net/        projid_map   stat
coredump_filter  limits   ns/         root@       statm
serge@ra:/proc/13318$ cat wchan
do_wait
```

Documentation/filesystems/proc.rst

72. Linux: трассировщики системных вызовов и библиотек.



Трассировщики

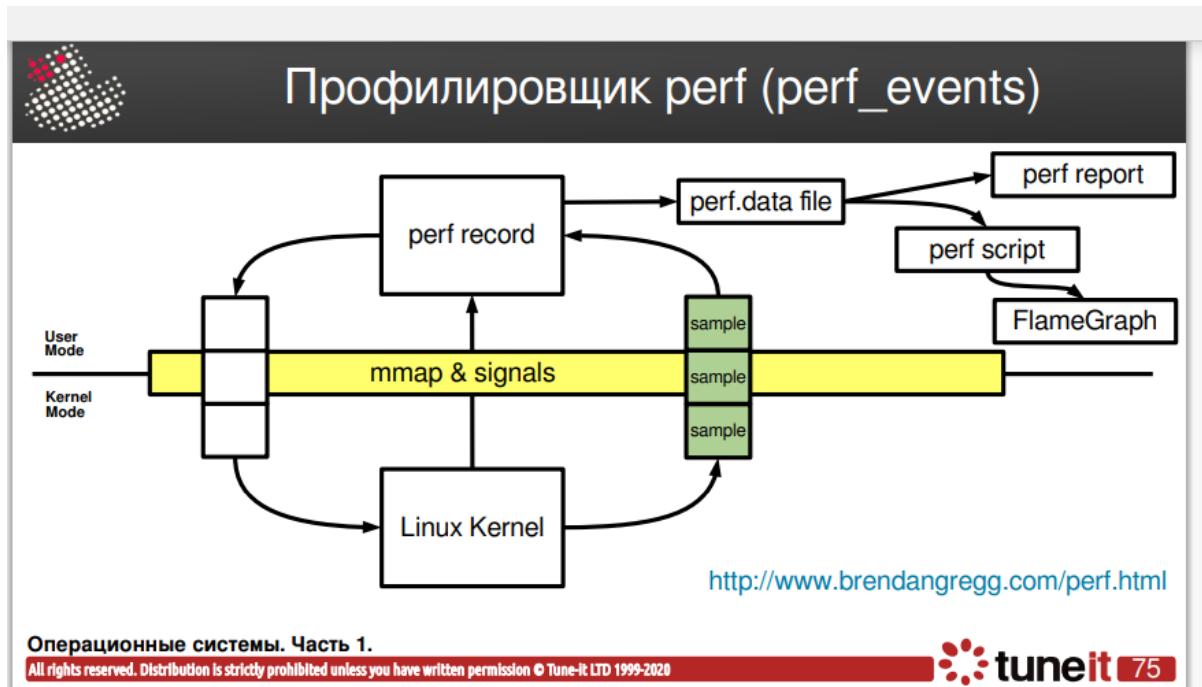
- Трассировка системных вызовов: strace
 - Трассировка вызовов библиотек: ltrace
 - Трассировка lock -офф(ф): bpftrace

Операционные системы. Часть 1.

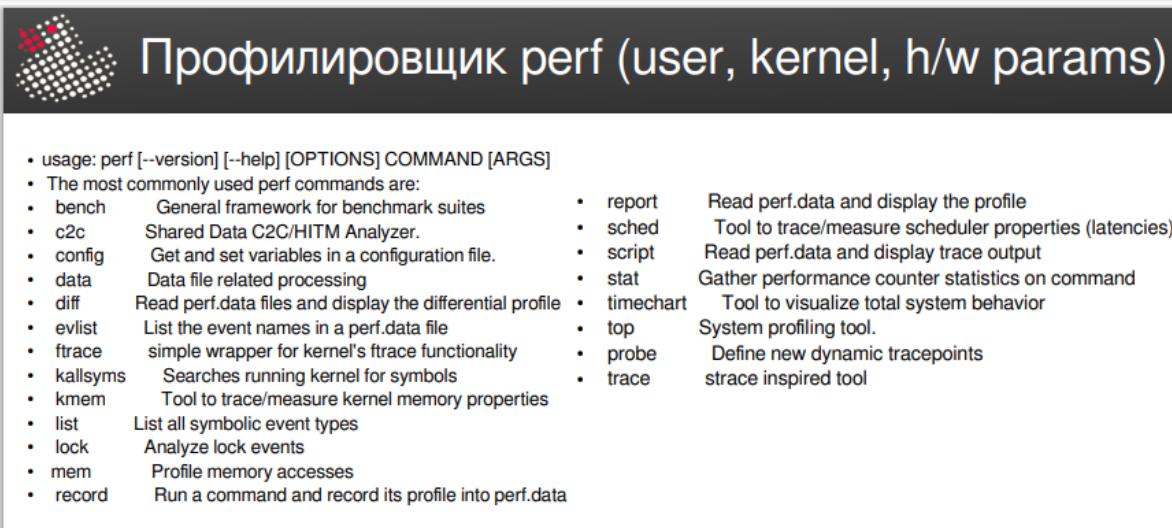
All rights reserved. Distribution is strictly prohibited unless you have written permission. © Type-it LTD 1999-2020



73. Linux: Профилировщик perf и FlameGraph.

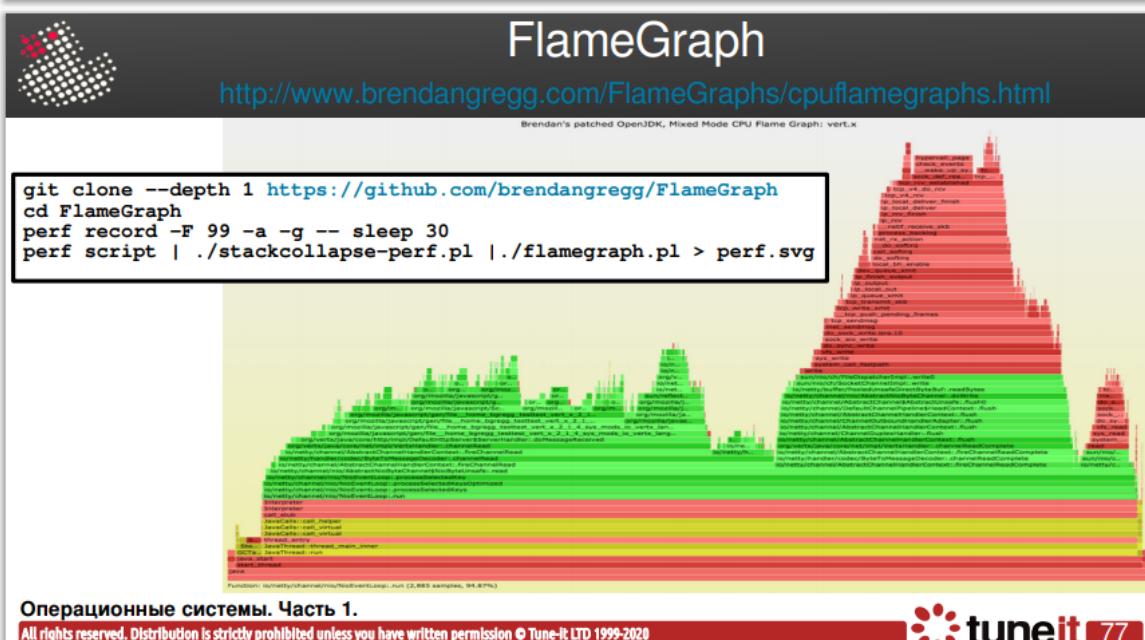


Операционные системы. Часть 1.
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



<http://www.brendangregg.com/perf.html>

[https://star-wiki.ru/wiki/Perf_\(Linux\)](https://star-wiki.ru/wiki/Perf_(Linux))



<http://www.brendangregg.com/FlameGraphs/cpuflamegraphs.html>

74. Linux: SystemTap.

<https://m.habr.com/ru/post/77502/>

75. Linux: Отладчик ядра.

- 1) Лучше этим не пользоваться чтобы не поломать свою ос, как говорил Миша Филиппов не играйте с ядром а то можно и проиграть
- 2) Есть команды ГО (запуск) и ? (хелп)
- 3) Это типа как гдб но для системных вызовов
- 4) Можно запускать удаленно
- 5) Это один из самых конченных билетов так что если он вам выпал вы плохо молились своему богу

(вот спасибо!)

<https://russianblogs.com/article/19411334724/>

.76. Windows: стандартные отладочные средства.

Resource Monitor

Имя	759%	17%
Память	Память	ЦП
Приложения (11)		
> IntelliJ IDEA (5)	395 622 427,7 МБ	19,4%

давай!! ДАВАААЙ УРААА
ДАВААЙ ДАААВАААЙ

Performance Monitor

Операционные системы. Часть 1. Обзор операционных систем
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 88

Операционные системы. Часть 1. Обзор операционных систем
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 89

<https://anydifferencebetween.com/performance-monitor-vs-resource-monitor/>

77. Windows: утилиты SysInternals



- Множество скриптов и программ для получения информации о системе
- Автор — Марк Руссинович, в настоящее время сотрудник Microsoft (соавтор книги Windows Internals)
- Отдельно загружается и устанавливается с сайта Microsoft

<https://docs.microsoft.com/ru-ru/sysinternals/>

Операционные системы. Часть 1. Обзор операционных систем
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



- PsList and PsKill — просмотр и остановка процессов (в том числе и удаленно)
- Process Explorer — просмотр ресурсов процесса, замена Task Manager
- Process Monitor — просмотр связанных с процессом ресурсов реестра
- Autoruns — поиск автозапускаемых программ
- Contig — дефрагментирует конкретный файл
- PSFile — позволяет показать открытые файлы, в том числе и удаленно
- MoveFile — перемещает заблокированные файлы во время перезагрузки.
- Sync — синхронизация файловой системы
- TCPview — информация о открытых сетевых соединениях
- SDelete — удалить файлы и папки без возможности восстановления

Операционные системы. Часть 1. Обзор операционных систем
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



<https://ru.wikipedia.org/wiki/Sysinternals>

78. Windows: отладчики WinDbg и KD

Абсолютно то же самое, что и отладчик ядра в линуке

.79. Аппаратная поддержка взаимных исключений.

- 1) Отключение прерываний

- 2) Атомарные операции на уровне ядра
(Столлингс стр. 285)
(задача о парикмахерской - проблема с очередью обслуживания)

.80. Эволюция подхода к блокировке (Столлингс, гл. 5.1)

?

81. Принципы взаимного блокирования (Столлингс, гл. 6.1)

История про машинки

Все взаимоблокировки предполагают наличие конфликта в борьбе за ресурсы между двумя или несколькими процессами. Наиболее ярким примером может служить транспортная взаимоблокировка. На рис. 6.1 показана ситуация, когда четыре автомобиля должны примерно одновременно пересечь перекресток. Четыре квадранта перекрестка представляют собой ресурсы, которые требуются процессам. В частности, для успешного пересечения перекрестка всеми четырьмя автомобилями необходимые ресурсы выглядят следующим образом.

- Автомобилю 1, движущемуся на север, нужны квадранты *а* и *б*.
- Автомобилю 2, движущемуся на запад, нужны квадранты *б* и *в*.
- Автомобилю 3, движущемуся на юг, нужны квадранты *в* и *г*.
- Автомобилю 4, движущемуся на восток, нужны квадранты *г* и *а*.

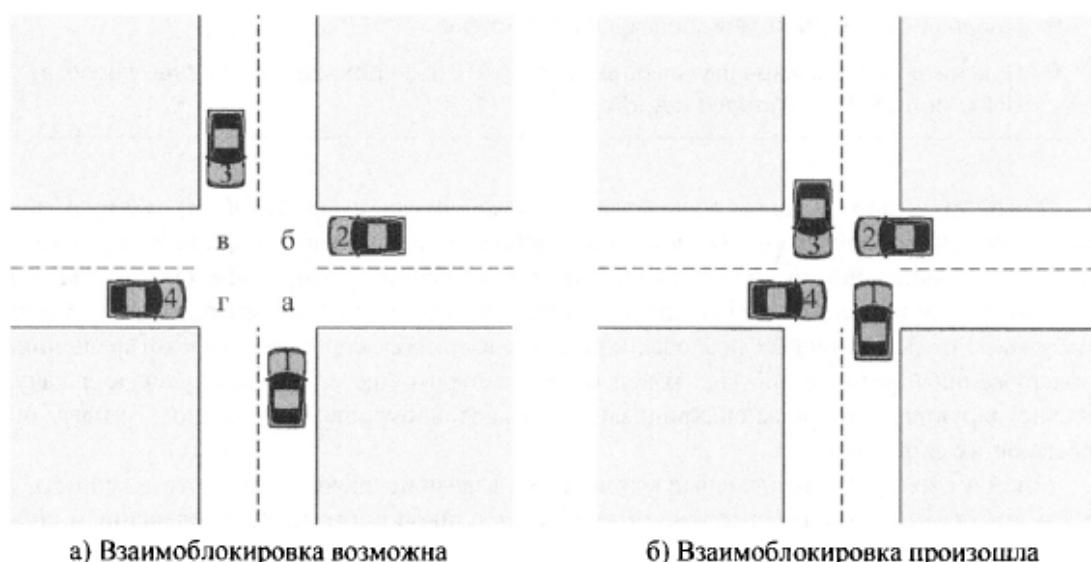


Рис. 6.1. Пример взаимоблокировки

A

82. Предотвращения взаимоблокировок, устранение взаимоблокировок, обнаружение блокировок. (Столлингс, гл. 6.2, 6.3, 6.4)

для того чтобы взаимоблокировка стала возможной, требуются три условия.

1. **Взаимные исключения.** Одновременно использовать ресурс может только один процесс.
2. **Удержание и ожидание.** Процесс может удерживать выделенные ресурсы во время ожидания других ресурсов.
3. **Отсутствие перераспределения.** Ресурс не может быть принудительно отобран у удерживающего его процесса.

Кроме перечисленных трех условий, необходимых, но не достаточных, для реального осуществления взаимоблокировки, требуется выполнение четвертого условия.

4. **Циклическое ожидание.** Существует замкнутая цепь процессов, каждый из которых удерживает как минимум один ресурс, необходимый процессу, следующему в цепи после данного (см. рис. 6.5, в и 6.6).

83. Задача об обедающих философах (Столлингс, гл. 6.6)

Теперь рассмотрим задачу об обедающих философах, представленную Дейкстрой в [65]. Итак, в некотором царстве, в некотором государстве жили вместе пять философов. Жизнь каждого из них проходила в основном в размышлениях, прерываемых приемом пищи. Философы давно сошлись во мнении, что только спагетти в состоянии восстанавливать их подточные непрерывными размышлениями силы.

Питались они за одним круглым столом (рис. 6.11), на который помещались большое блюдо со спагетти, пять тарелок, по одной для каждого философа, и пять вилок. Проголодавшийся философ садится на свое место за столом и, пользуясь двумя вилками, приступает к еде. Задача состоит в том, чтобы разработать ритуал (читай — алгоритм) обеда, который обеспечивает взаимоисключений (два философа не могут одновременно пользоваться одной вилкой) и не допускает взаимоблокировок и голодания (обратите внимание, насколько уместным оказался этот термин в данной задаче!).

Эта задача Дейкстры может показаться не очень важной, но она очень хорошо иллюстрирует проблемы взаимоблокировок и голодания. Кроме того, при решении данной задачи приходится сталкиваться со многими трудностями в организации параллельных вычислений (см., например, [89]). Задача об обедающих философах может рассматриваться как типичная задача, возникающая в многопоточных приложениях при работе с совместно используемыми ресурсами и, соответственно, может выступать в качестве тестовой при разработке новых подходов к проблеме синхронизации.

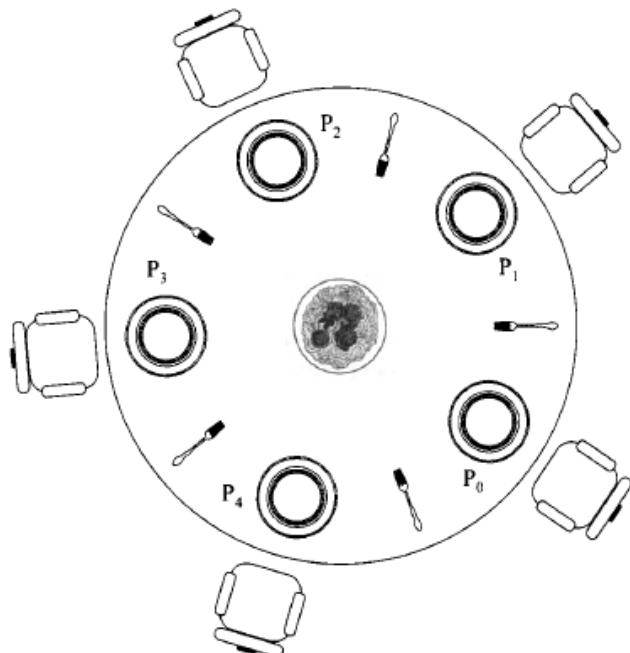
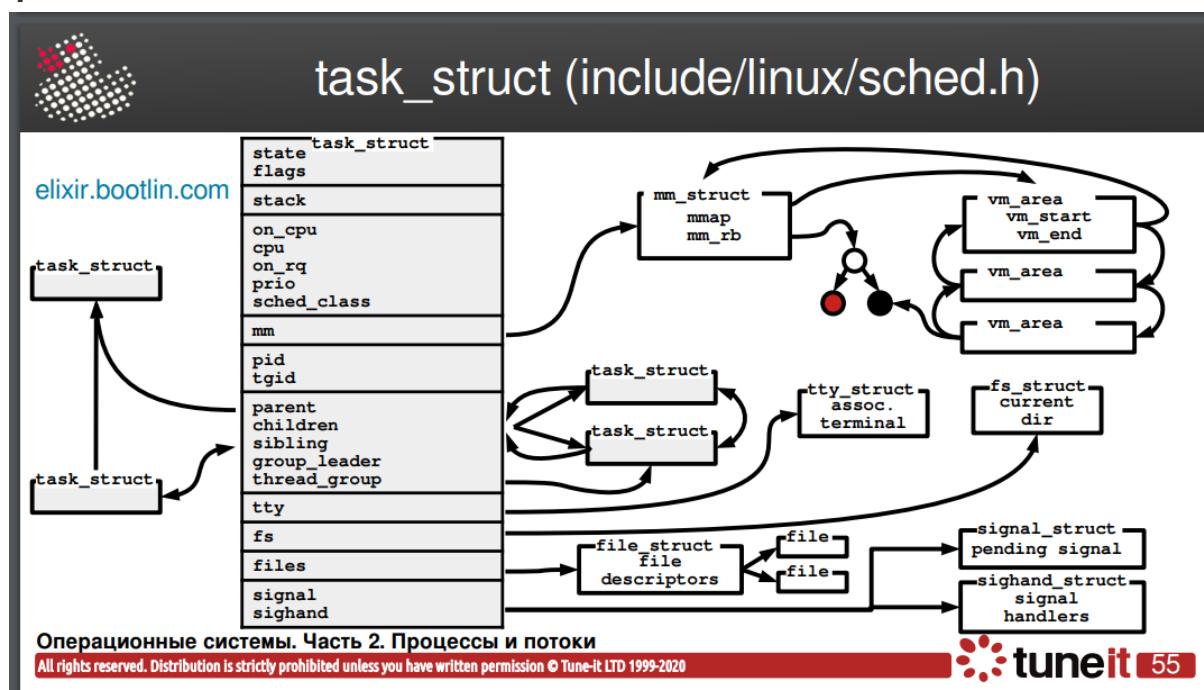


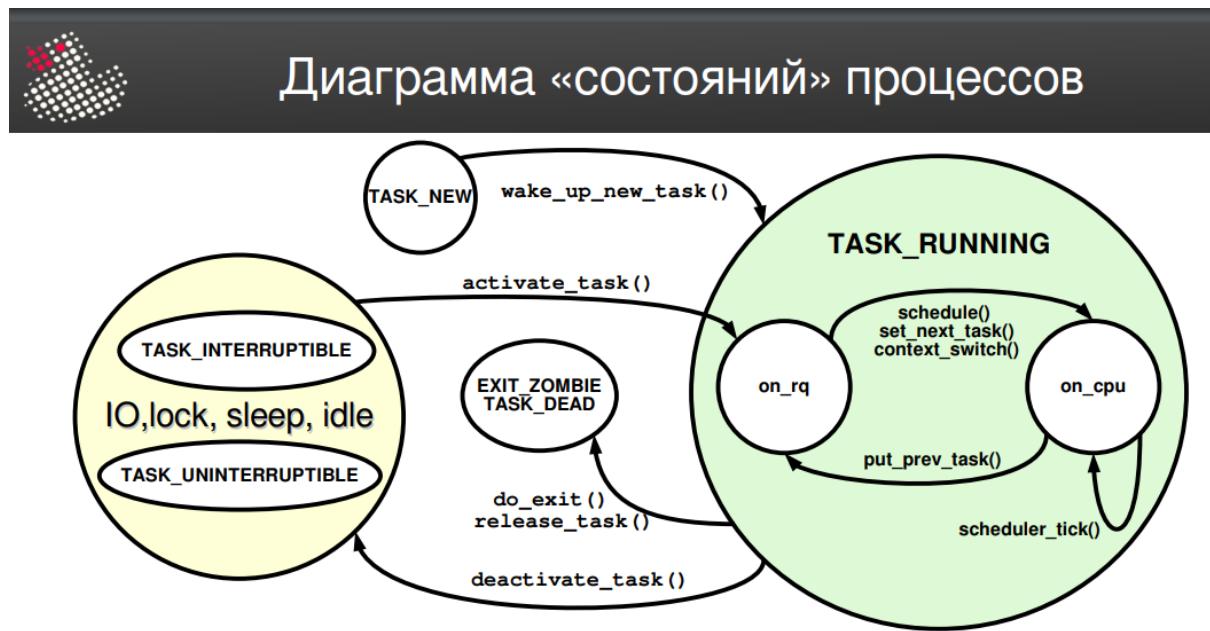
Рис. 6.11. Обеденный стол философов

84. Процессы в Linux: структура task_struct, поля структуры, связь с другими структурами ядра.

В ядре Linux процесс представлен довольно большой структурой `task_struct` (дескриптором процесса). Помимо самой необходимой для описания процесса информации эта структура содержит массу других данных, используемых для учета и связи с другими процессами (родительскими и порожденными). Полное описание структуры `task_struct` выходит за рамки статьи, однако, ее фрагмент, содержащий упомянутые в статье элементы, приведен в листинге 1. Стоит заметить, что структура `task_struct` объявлена в файле `./linux/include/linux/sched.h`.



85. Диаграмма состояния процесса Linux.



<https://ravesli.com/processes-v-linux/>

86. Создание процесса Linux на уровне пользовательского процесса.



Userland: Создание процесса

- `fork()` - создать процесс
 - Возвращает номер процесса родителю, 0 — созданному процессу, -1 в случае ошибки.
- `vfork()` - создать процесс с копией адресного пространства родителя
 - Родительский процесс блокирован до тех пор, пока дочерний не вызовет `exit()` или `execve()`
- `clone()` - создать процесс, управляя копированием выбранных частей процесса
 - `clone(CLONE_VM | CLONE_FS | CLONE_FILES | CLONE_SIGHAND, 0);`
- `execv(const char *path, const char *arg, ...)` - перекрытие образа процесса



Linux Threads

- В ядре Linux нет примитива Thread, соответствующего пользовательскому потоку
 - Соответственно нет планирования потоков
- Потоки создаются так же, как и процессы
 - Разделяют ресурсы вызвавшего `fork` процесса
 - Каждый поток — отдельная `task_struct`
 - **Создание:** `clone (CLONE_VM | CLONE_FILES | CLONE_FS | CLONE_SIGHAND | CLONE_THREAD | CLONE_SETTLS | CLONE_PARENT_SETTID | CLONE_CHILD_CLEARTID | CLONE_SYSVSEM, 0); (NPTL version)`

87. Создание и завершение процесса Linux на уровне ядра. Вызываемые функции.



88. Особенности реализации потоков в Linux. KThread. Tasklet.



KThread

- Предназначен для выполнения фоновых операций в ядре
 - Планировщик работает с ними так же, как с процессами
- Нет своего адресного пространства
 - mm в task_struct равен NULL
- kthread являются потомками kthreadd (pid==2)
 - Все фоновые потоки ядра ps --ppid=2
- Создание — kthread_create()
 - После создания нужно выполнить wake_up_process()
 - Одновременное создание и запуск макрос kthread_run()
- Для остановки другие потоки вызывают kthread_stop()
 - А сам поток должен проверять kthread_should_stop()

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Tasklets

```
include/linux/interrupt.h
struct tasklet_struct
{
    struct tasklet_struct *next;
    unsigned long state;
    atomic_t count;
    void (*func) (unsigned long);
    unsigned long data;
};
```



- Обслуживание «bottom half» для обработки программного прерывания (soft irq)
- Один tasklet одновременно на одном CPU, разные одновременно на разных CPU
- Два приоритета — normal и hi
 - tasklet_schedule и tasklet_hi_schedule
 - «Hi» выполняются раньше
- count - 1 - деактивирован

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



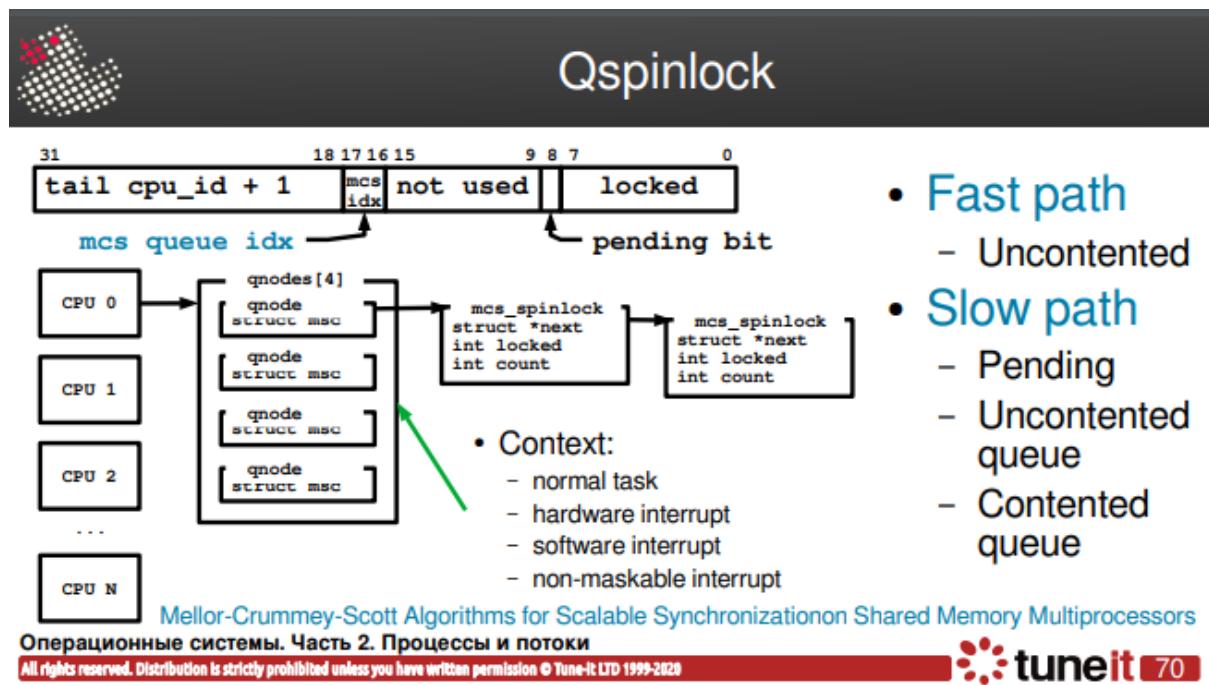
<https://habr.com/ru/company/embox/blog/244071/>

,89. Примитивы синхронизации Linux. Spinlock и qspinlock.

Сpin-блокировка или **спинлок** (англ. *spinlock* — циклическая блокировка) — низкоуровневый примитив **синхронизации**^[1], применяемый в многопроцессорных системах для реализации взаимного исключения исполнения **критических участков** кода с использованием цикла активного ожидания^[2]. Применяется в случаях, когда ожидание захвата блокировки предполагается недолгим^[2], либо если контекст выполнения не позволяет переходить в заблокированное состояние^[3].

Спин-блокировки являются аналогами **мьютексов**, позволяющими тратить меньше времени на процедуру блокировки потока, поскольку не требуется переводить поток в заблокированное состояние. В случае мьютексов может потребоваться задействование планировщика с переводом потока в другое состояние и добавлением его в список потоков, ожидающих разблокировки. Спин-блокировки не задействуют планировщик и используют цикл активного ожидания без изменения состояния потока, что приводит к трате процессорного времени на ожидание освобождения блокировки другим потоком. Типовой реализацией спин-блокировки является простая циклическая проверка переменной спин-блокировки на доступность^[1]

Qspinlock - spinlock с очередью



90. Примитивы синхронизации Linux. Semaphore и Mutex.



Semaphore

```
include/linux/semaphore.h
struct semaphore {
    raw_spinlock_t lock;
    unsigned int count;
    struct list_head wait_list;
};

#define __SEMAPHORE_INITIALIZER(name, n) \
{ \
    .lock = __RAW_SPIN_LOCK_UNLOCKED((name).lock), \
    .count = n, \
    .wait_list = LIST_HEAD_INIT((name).wait_list), \
}
static inline void sema_init(struct semaphore *sem, int val)
{
    static struct lock_class_key __key;
    *sem = (struct semaphore) __SEMAPHORE_INITIALIZER(*sem, val);
    lockdep_init_map(&sem->lock.dep_map, "semaphore->lock", &__key, 0);
}
```

- void **down**(struct semaphore *sem);
- void **up**(struct semaphore *sem);
- int **down_interruptible**(struct semaphore *sem);
- int **down_killable**(struct semaphore *sem);
- int **down_trylock**(struct semaphore *sem);
- int **down_timeout**(struct semaphore *sem, long jiffies);

Операционные системы. Часть 2. Процессы и потоки

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



Semaphores up/down

```
kernel/locking/semaphore.c
int down_interruptible(struct semaphore *sem)
{
    unsigned long flags;
    int result = 0;

    raw_spin_lock_irqsave(&sem->lock, flags);
    if (likely(sem->count > 0))
        sem->count--;
    else
        result = __down_interruptible(sem);
    raw_spin_unlock_irqrestore(&sem->lock, flags);
    return result;
}

__down_common(sem,
             TASK_INTERRUPTIBLE,
             MAX_SCHEDULE_TIMEOUT);

list_add_tail();
__set_current_state(state);

kernel/locking/semaphore.c
void up(struct semaphore *sem)
{
    unsigned long flags;

    raw_spin_lock_irqsave(&sem->lock, flags);
    if (likely(list_empty(&sem->wait_list)))
        sem->count++;
    else
        __up(sem);
    raw_spin_unlock_irqrestore(&sem->lock, flags);
}

static void __up(struct semaphore *sem)
{
    struct semaphore_waiter *waiter =
        list_first_entry(&sem->wait_list,
                        struct semaphore_waiter, list);
    list_del(&waiter->list);
    waiter->up = true;
    wake_up_process(waiter->task);
}
```

Операционные системы. Часть 2. Процессы и потоки

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

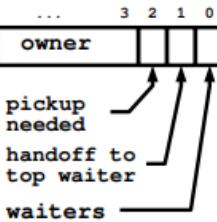




Mutex

```
include/linux/mutex.h
struct mutex {
    atomic_long_t owner;
    spinlock_t wait_lock;
#endif CONFIG_MUTEX_SPIN_ON_OWNER
/* Spinner MCS lock */
    struct optimistic_spin_queue osq;
#endif
    struct list_head wait_list;
    //DEBUG STUFF
};
```

- Only one task can hold the mutex at a time.
- Only the owner can unlock the mutex.
- Multiple unlocks are not permitted.
- Recursive locking/unlocking is not permitted.
- A mutex must only be initialized via the API.
- A task may not exit with a mutex held.
- Memory areas where held locks reside must not be freed.
- Held mutexes must not be reinitialized.
- Mutexes may not be used in hardware or software interrupt contexts such as tasklets and timers.



- ww_mutex**
 - deadlock free
- Wound-Wait**
 - wait-kill

- Fast path**
 - uncontended
- Mid path**
 - optimistic locking
- Slow path**
 - block waiters
- Unlock path**
 - wake up

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020

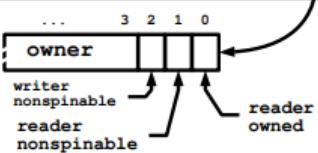
.91. Примитивы синхронизации Linux. rw_semaphore, seqloc



rw_semaphore

```
include/linux/rwsem.h
struct rw_semaphore {
    atomic_long_t count;
    atomic_long_t owner;
#endif CONFIG_RWSEM_SPIN_ON_OWNER
/* spinner MCS lock */
    struct optimistic_spin_queue osq;
#endif
    raw_spinlock_t wait_lock;
    struct list_head wait_list;
    //DEBUG STUFF
};
```

Bit 0	- writer locked bit
Bit 1	- waiters present bit
Bit 2	- lock handoff bit
Bits 3-7	- reserved
Bits 8-62	- 55-bit reader count
Bit 63	- read fail bit



- void down_read(struct rw_semaphore *sem);
- int down_read_killable(struct rw_semaphore *sem);
- int down_read_trylock(struct rw_semaphore *sem);
- void down_write(struct rw_semaphore *sem);
- int down_write_killable(struct rw_semaphore *sem);
- int down_write_trylock(struct rw_semaphore *sem);
- void up_read(struct rw_semaphore *sem);
- void up_write(struct rw_semaphore *sem);

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020



Sequence counters and seqlock

- Механизм для реализации неблокирующего чтения (с повторами) и записи без голодаания

- Несколько версий

- **seqcount_t** — запись должна быть синхронизирована внешними средствами
- **seqcount_LOCKTYPE_t** — запись синхронизирована средствами выбранного LOCKTYPE
- **seqcount_t с семантикой защелки** (две копии данных для четного/нечетного значения счетчика)
- **seqlock_t** — запись синхронизирована spinlock и «non preemptible»

```
do {
    seq = read_seqbegin(&foo_seqlock);
    /* [[read-side critical section]] */
} while (read_seqretry(&foo_seqlock, seq));
```

Sequence counters and sequential locks

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020

.92. Типы процессов и потоков Windows.

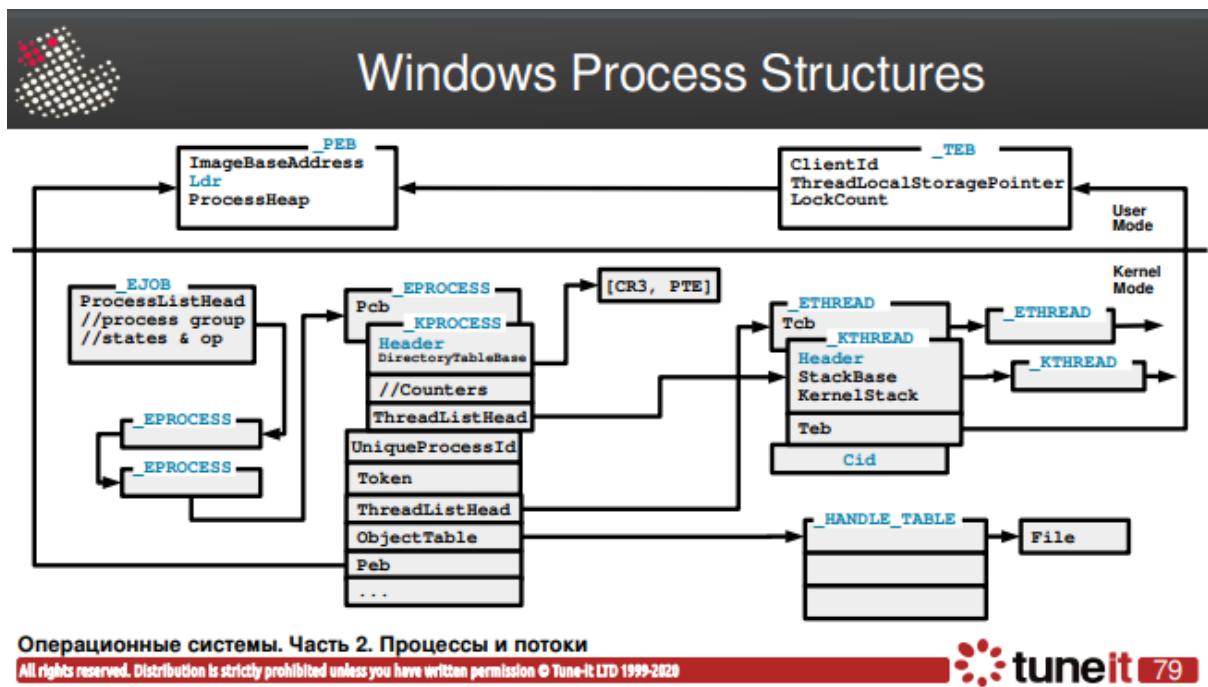
Типы процессов

- «Современные процессы»
 - Universal Windows Platform (UWP) processes (AKA Immersive processes) — Windows 8 (Windows Store)
 - Protected processes (Windows Media Certificate)
 - Protected Processes Light
- Minimal processes
 - Нет пользовательского адресного пространства
 - System process и Memory Compression process
- Pico-процессы
 - Drawbridge project (Pico providers, например Lxss.sys и LxCore.sys)
 - Ограничивают доступ к системным функциям, предоставляют набор callback
- Trustlets - Isolated User Mode (IUM) Processes
 - используют виртуализацию VTL1 (Virtual Trust Level 1), VTL0 — остальная система
- Windows on Windows (WOW) - 32 бита в 64 битном режиме
- JOBS — средство группировки процессов
- Dos, Win16, bat-файлы

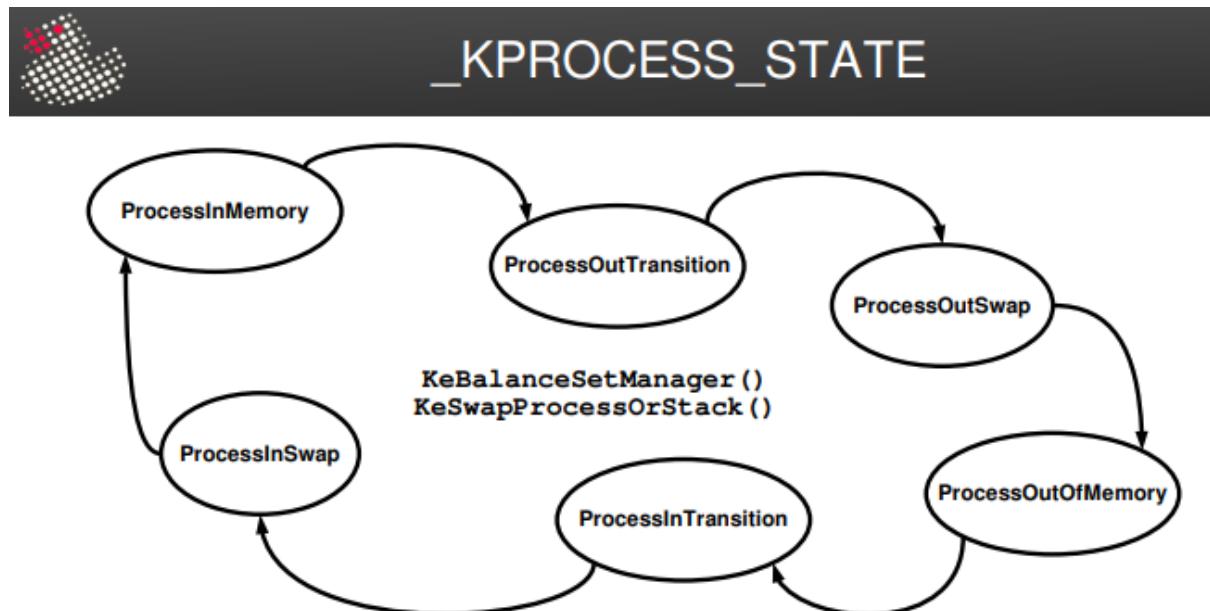
Операционные системы. Часть 2. Процессы, потоки и планирование
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020

tuneit 77

.93. Структура процесса и потока в Windows. Поля структур.

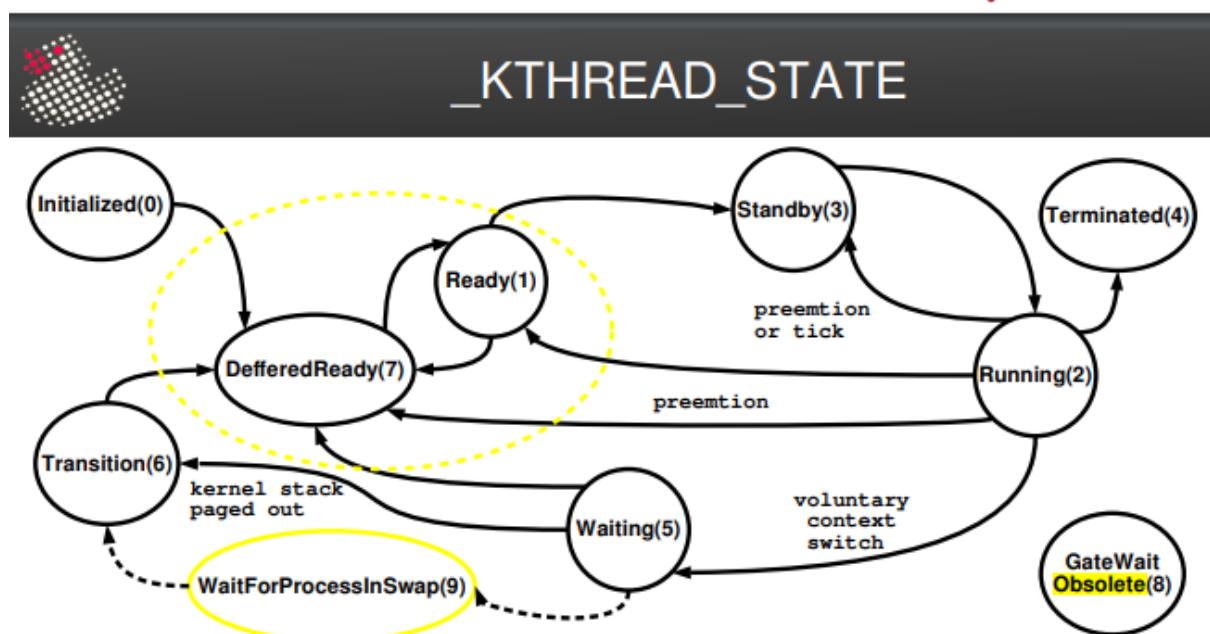


.94. Диаграммы состояний процесса и потока Windows



Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020

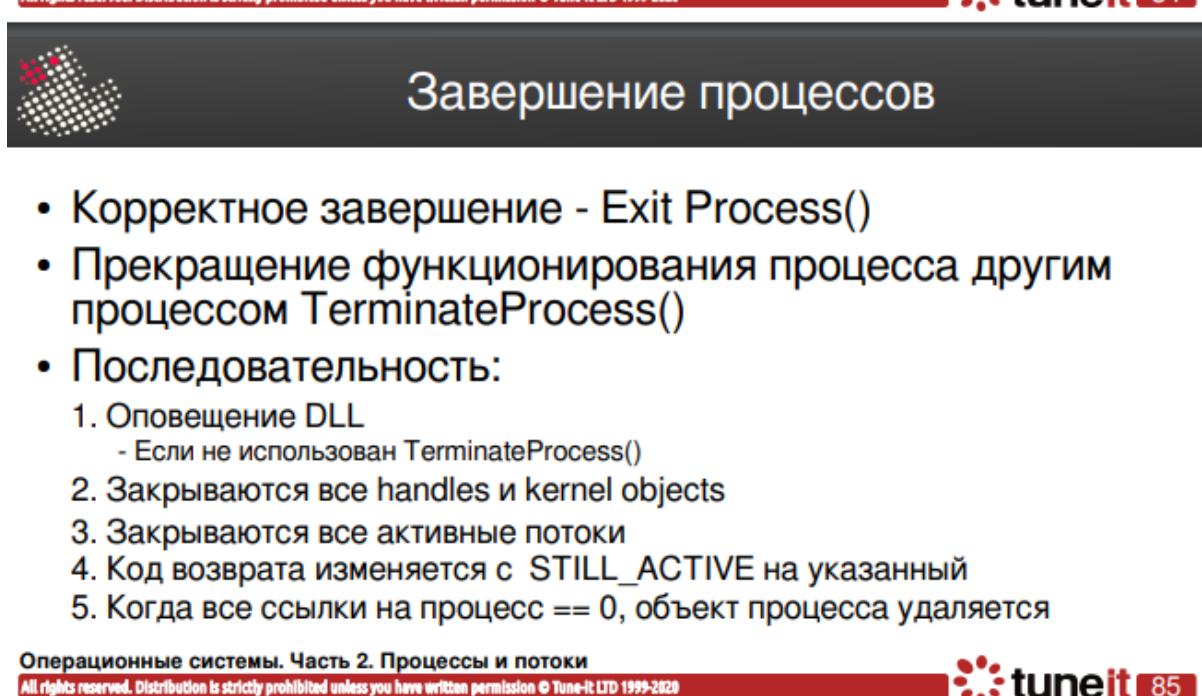
tuneit 80



Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-It LTD 1999-2020

tuneit 81

.95. Создание и завершение процесса Windows.



.96. Примитивы синхронизации Windows. Понятие Dispatcher Object. Ожидание наступление события, вызовы Wait.



Kernel Object = { Dispatcher Object
Control Object }

- Любой объект ядра, на котором можно ожидать события - «dispatcher object»
 - Некоторые предназначены только для синхронизации (events, mutexes, semaphores, queues, ...)
 - Другие для ожидания событий от процессов, потоков, файлов
- Общая структура данных
- Два состояния
 - «Signaled» - ожидание удовлетворено и «Not-signaled»
 - еще нет.
 - Различные объекты отличаются в том, что именно изменяет их состояние
 - Wait и unwait — операции общие.

```
ntos/inc/ntosdef.h
typedef struct _DISPATCHER_HEADER {
    union {
        struct {
            UCHAR Type;
            union {
                UCHAR Absolute;
                UCHAR NpxIrql;
            };
        };
        union {
            UCHAR Size;
            UCHAR Hand; В win10 больше!
        };
        union {
            UCHAR Inserted;
            BOOLEAN DebugActive;
        };
    };
    volatile LONG Lock;
} DISPATCHER_HEADER;
```

Операционные системы. Часть 2. Процессы и потоки

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



Вызовы Wait

- Гибкие вызовы ожидания — KeWaitForSingleObject, KeWaitForMultipleObjects
 - Ожидает одного или нескольких объектов ("any" или "all"). В случае всех — все объекты должны быть в состоянии «signaled»
 - Таймаут задается опционально
- Объекты, которые можно ожидать:
 - Events, Mutexes, Semaphores, Timers
 - Processes and Threads (exit or terminate)
 - Directories (change notification)
- Нет гарантированного порядка завершения ожидания
 - Потоки ожидающие события пробуждаются в неопределенном порядке, зависит от диспетчера
 - Обычный порядок - FIFO; Однако APC может изменить этот порядок

```
ntos/inc/ke.h
typedef enum _KWAIT_REASON {
    PowerUp,
    FreePage,
    PageIn,
    PoolAllocation,
    DelayExecution,
    Suspend,
    UserRequest,
    MrExecutive,
    MrFreePage,
    MrPageIn,
    MrPocllocation,
    MrUserRequest,
    MrSuspended,
    MrUserRequest,
    MrEventPair,
    MrQueue,
    MrIpcReceive,
    MrIpcSend,
    MrVirtualMemory,
    MrPageOut,
    MrRendezvous,
    Spare2,
    Spare3,
    Spare4,
    Spare5,
    Spare6,
    MrKernel,
    MrResource,
    MrPushLock,
    MrWait,
    MrQuantumEnd,
    MrDispatchInt,
    MrPreempted,
    MrTimeExecution,
    MrWaitStatus,
    MrGuardedMutex,
    MrRundown,
    MaximumWaitReason
} KWAIT_REASON, RE;
```

Операционные системы. Часть 2. Процессы и потоки

All rights reserved. Distribution is strictly prohibited unless you have written permission © Tune-it LTD 1999-2020



.97. Примитивы синхронизации Windows. EventObject, Mutex, Mutant.



EventObject

- Может сбрасываться вручную, автоматически и «pulsed»
- Используется для оповещения о наступлении события двух типов:
 - Синхронизация — один из ждунов продолжает выполнение когда наступает состояние Signaled. EO автоматически сбрасывается в Not-Signaled
 - Нотификация — пробуждает всех ждунов, требует ручного сброса состояния
- Операции:
 - KeInitializeEvent - Initialize an event object
 - KePulseEvent - Set/reset event object state atomically
 - KeReadStateEvent - Read state of event object
 - KeResetEvent - Set event object to Not-Signaled state
 - KeSetEvent - Set event object to Signaled state

```
ntos/inc/ntosdef.h
typedef struct _KEVENT {
    DISPATCHER_HEADER Header;
} KEVENT, *PKEVENT, *PRKEVENT;
```



Mutexes и Mutants

- Mutually exclusive, deadlock free доступ к разделяемым ресурсам
- В нормальном режиме создается в «signaled» состоянии
- Возможен рекурсивный захват (сколько захватов, столько освобождений)
- Захваченный mutex блокирует выход из Kernel Mode
- Снятие Lock — mutant: любой поток (abandoned state); mutex: владелец
- В mutex запрещены APC, в mutant — нет.
- Возможно использовать mutant в UserLand
- Increment — добавляются к приоритету потока если wait satisfied
- Wait — за KeRelease* сразу следует функция ожидания (освободить и занять в рамках «атомарной операции»)

```
ntos/inc/ke.h
typedef struct _KMUTANT {
    DISPATCHER_HEADER Header;
    LIST_ENTRY MutantListEntry;
    struct _KTHREAD *OwnerThread;
    BOOLEAN Abandoned;
    UCHAR ApcDisable;
} KMUTANT, KMUTEX;
```

- Mutant:
 - KeInitializeMutant
 - KeReadStateMutant
 - KeReleaseMutant
 - KPRIORITY Increment,
 - BOOLEAN Abandoned,
 - BOOLEAN Wait
- Mutex:
 - KeInitializeMutex
 - KeReadStateMutex
 - KeReleaseMutex
 - KeReleaseMutant(Mutex, 1, FALSE, Wait)

.98. Примитивы синхронизации Windows. Fast mutex, Guarded mutex.



Fast mutexts и Guarded mutexes

- Больше похожи на «нормальные» мьютексы
- Поле Count: 0-й бит — lock, 1-й — single waiter woken
- Нельзя захватывать рекурсивно
- Операции:
 - ExInitializeFastMutex
 - ExAcquireFastMutex
 - ExTryToAcquireFastMutex
 - ExReleaseFastMutex
- До Windows 8 — разные реализации, после - идентичны

```
ntos/inc/ex.h
typedef struct _FAST_MUTEX {
    LONG Count;
    PKTHREAD Owner;
    ULONG Contention;
    KEVENT Event;
    ULONG OldIrql;
} FAST_MUTEX, *PFAST_MUTEX;
```

.99. Примитивы синхронизации Windows. Semaphore, spinlock



Fast mutexes и Guarded mutexes

- Больше похожи на «нормальные» мьютексы
 - Поле Count: 0-й бит — lock, 1-й — single waiter woken
 - Нельзя захватывать рекурсивно
 - Операции:
 - ExInitializeFastMutex
 - ExAcquireFastMutex
 - ExTryToAcquireFastMutex
 - ExReleaseFastMutex
 - До Windows 8 — разные реализации, после - идентичны

```
ntos/inc/ex.h
```

```
typedef struct _FAST_MUTEX {  
    LONG Count;  
    PKTHREAD Owner;  
    ULONG Contention;  
    KEVENT Event;  
    ULONG OldIrql;  
} FAST_MUTEX, *PFAST_MUTEX;
```

Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless with written permission © Tuna-IT LTD 1998-2028

tuneit 92



Semaphores

- Управляют захватом разделяемого ресурса
 - Limit — максимальное количество такого ресурса
 - Count — начальное состояние, помещается в Header.SignaledState
 - Семафор открыт, когда SignaledState > 0
 - Если SignaledState+Adjustment > Semaphore→Limit
 - вызывается ExRaiseStatus(STATUS_SEMAPHORE_LIMIT_EXCEEDED)
 - Increment — добавляются к приоритету потока если wait satisfied
 - Wait — за KeReleaseSemaphore сразу следует функция ожидания (освободить и занять в рамках «атомарной операции»)

ntos/inc/ke.h

- KeInitializeSemaphore
 - LONG Count
 - LONG Limit
 - KeReadStateSemaphore
 - KeReleaseSemaphore
 - KPRORITY Increment,
 - LONG Adjustment,
 - BOOLEAN Wait

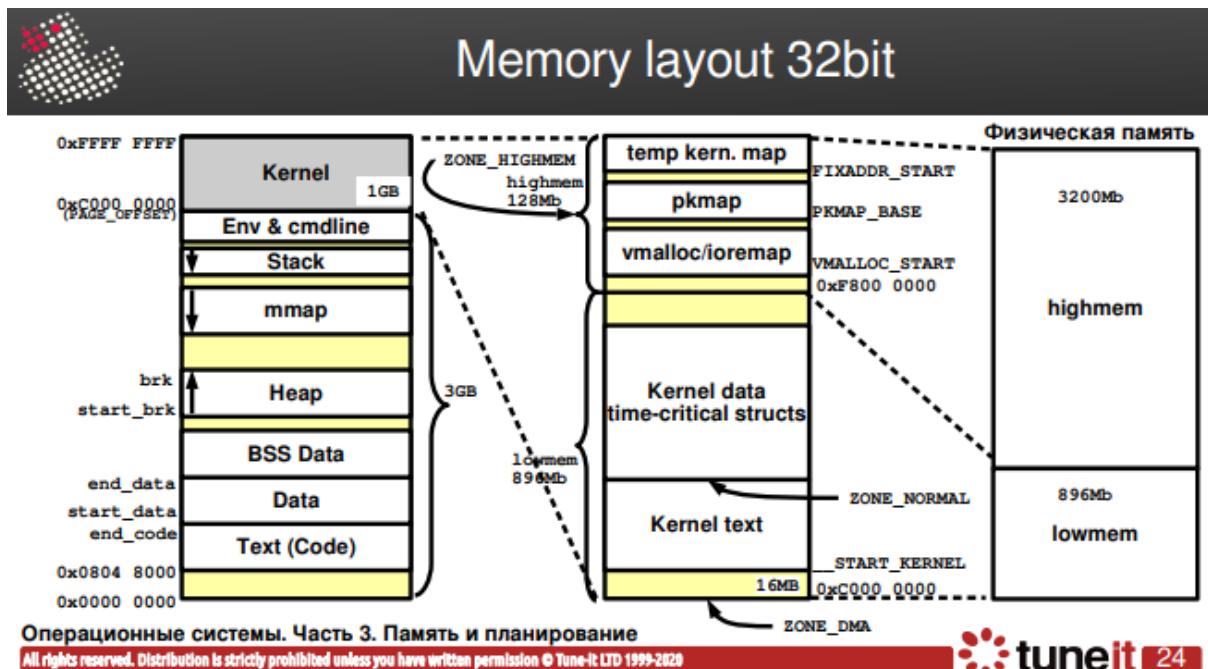
Операционные системы. Часть 2. Процессы и потоки
All rights reserved. Distribution is strictly prohibited unless with written permission © Tuna-IT LTD 1998-2028

tuneit 93

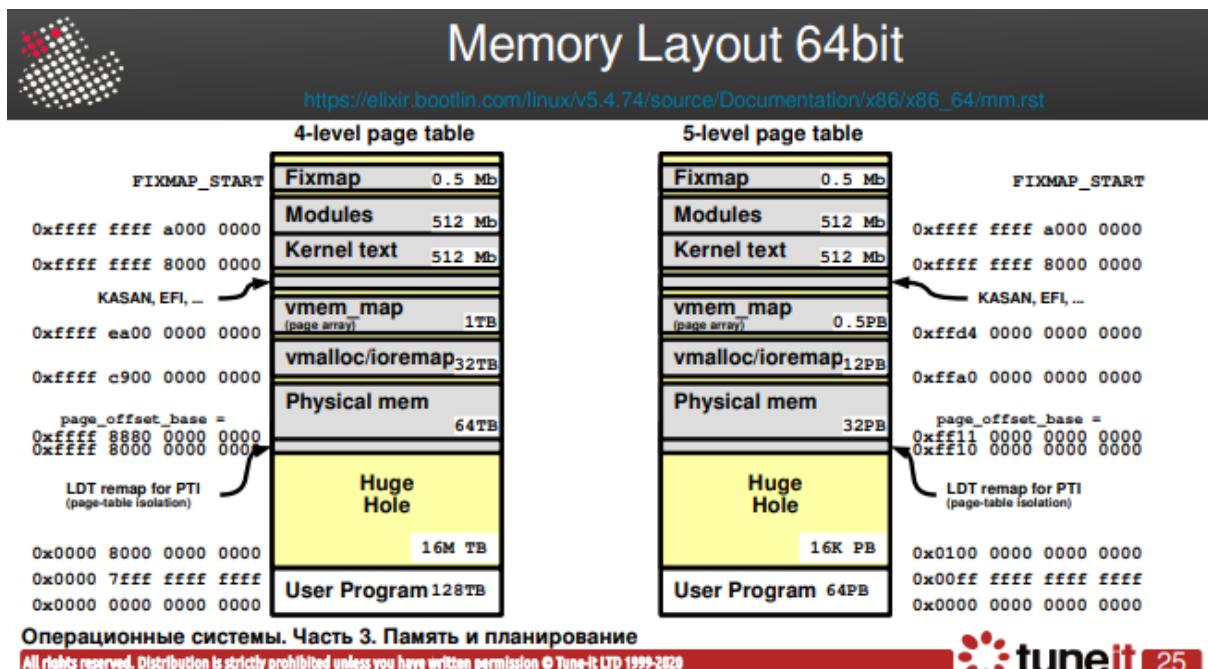
- .100. Хешированная таблица страниц SPARC64.

<https://www.tune-it.ru/web/il/home/-/blogs/%D0%B2%D0%B8%D1%80%D1%82%D1%83%D0%B0%D0%BB%D1%8C%D0%BD%D0%B0%D1%8F-%D0%BF%D0%B0%D0%BC%D1%8F%D1%82%D1%8C-%D0%BD%D0%B0-sun-%D0%BA%D0%B0%D1%85-sparc64-%D0%B8%D1%85%D1%8D%D1%88%D0%B8%D1%80%D0%BE%D0%B2%D0%B0%D0%BD%D0%BD%D0%B0%D1%8F-%D1%82%D0%B0%D0%B1%D0%BB%D0%B8%D1%86%D0%B0-%D1%81%D1%82%D1%80%D0%B0%D0%BD%D0%BD%D0%B8%D1%86>

.101. Виртуальная память Linux. 32-х разрядная модель.



.102. Виртуальная память Linux. 64-х разрядные модели.



.103. Виртуальная память Linux. Структуры памяти.



Ничего не понятно. Худший билет во всем экзамене)

На слайде представлен кусочек того, что имеет отношение к структурам внутри памяти.

Начнем с тех частей, которые нам хорошо известны.

В первую очередь это `task_struct` (в левом нижнем углу схемы) – место, где расположена информация, относящаяся к процессу и находящаяся в памяти ядра. И это `mm_struct`, где находится описание памяти для процесса ядра.

Здесь выделено несколько переменных (первая из них указывает на `mmap`. `Mmap` – эта такая штука, которая представляет связанный список всех сегментов текущего процесса (предпоследний «столбец» справа)).

Причем видно, что каждый сегмент описывается структурой `vm_area_struct`. И в этой структуре находятся описательные элементы этого сегмента (например, есть адрес, где сегмент начинается и где заканчивается: `vm_start`, `vm_end`; есть права, указывающие на то, что разрешено делать в этой самой страничке: `READ`, `EXEC`; и есть указатели на следующий и предыдущий сегмент, чтобы в случае необходимости их можно было подвернуть изучению).

В правом «столбце» в явном виде под каждым сегментом подписано, как бэкэндятся данные. Т.е. условно говоря, что «лежит под» этим самым сегментом. Имеется в виду, что сегменты «внизу себя» хранят либо своп-область и зарезервированное место в свопе, либо они хранят «внизу себя» файл.

И, соответственно, есть структура (`struct_file`), которая соответствует каждому сегменту, т. е. в явном виде есть ссылка на файл прямо из структуры `vm_area_struct`, где описано, что это за файл. И там есть поле `f_mapping`, которое указывает на маппинг файла.

Кроме того, есть анонимные странички `Stack anonymous` (показаны в верхнем правом углу схемы). У них есть отдельная структура, которая называется `anon_vma`, которая описывает анонимную область памяти.

Следует отметить, что сегментов в текущем процессе может быть очень много. Поэтому, для того чтобы найти сегмент по виртуальному адресу, существует красно-черное дерево поиска. Переменная в `mm_struct` (называемая `mm_rb`) ускоряет поиск по сегментам, необходимый для того, чтобы найти необходимый сегмент, которому принадлежит тот или иной адрес.

Таким образом, все сегменты, которые имеют «снизу» файл, имеют ссылку на `struct_file`. Все сегменты, которые анонимны – имеют структуру, соответствующую анонимной памяти.

Страницы «в регионе» `struct_file` могут или «поддерживаться» файлами, или быть анонимными. В каждой структуре памяти есть указатель на таблицу страниц. Таблица страниц «живет» относительно процесса: процесс при помощи регистра CR3 указывает на вход в нее; там есть специальные функции, которые осуществляют поиск по номеру страницы. Поэтому легко можно найти в `vmem_map` (верхний левый угол слайда) соответствующую страничку просто прибавив к ней frame number.

В `vmem_map` хранятся описания всех страниц. Внутри каждой таблицы есть «реверс-маппинг», показывающий, кому она принадлежит. Так же там существует некоторое количество указателей, по которым можно определить, что находится в таблице. Поэтому, например, пэйджер и свопер, будут работать с каждой страничкой в зависимости от ее типа.

В `address_space` есть специальная переменная `i_pages` которая указывает на `xarray` (это тоже разновидность дерева для быстрого поиска). И этот `xarray` содержит странички, которые временно скэнированы в памяти. Т. е. описание всех страниц у нас есть. А в `xarray` будут у нас находиться те страницы, которые в реальности находятся в памяти и не требуют загрузки с диска.

Если нужно установить все элементы (страницы) файла, то для этого есть другая переменная, которая называется `i_mmap`. При этом тоже используется красно-черное дерево, которое устанавливает связи с блоками данных, которые являются файлом.

.104. Способы выделения памяти для пользовательских процессов Linux



Выделение памяти user-space

- **Malloc**

- Несколько реализаций (buddy, glibc, ...)
- При нехватке памятидвигает program break — sbrk()

- **Mmap**

- Гибкий вызов создания областей памяти
- Создает «file backed» и анонимную память

```
void *malloc(size_t size);
void free(void *ptr);
void *calloc(size_t nmemb, size_t size);
void *realloc(void *ptr, size_t size);
void *reallocarray(void *ptr,
                  size_t nmemb, size_t size);
```

```
void *mmap(void *addr, size_t length,
           int prot, int flags,
           int fd, off_t offset);
int munmap(void *addr, size_t length);
```

PROT_NONE	MAP_SHARED	MAP_LOCKED
PROT_EXEC	MAP_PRIVATE	MAP_NONBLOCK
PROT_READ	MAP_32BIT	MAP_NORESERVE
PROT_WRITE	MAP_ANONYMOUS	MAP_POPULATE
	MAP_FIXED	MAP_STACK
	MAP_GROWSDOWN	MAP_UNINITIALIZED
	MAP_HUGETLB	
	MAP_HUGE_2MB	
	MAP_HUGE_1GB	

.105. Способы выделения памяти в пространстве ядра Linux.



Выделение памяти в пространстве ядра

- Выделение страниц

- page *alloc_pages(gfp_t flags,order) — выделяет 2^{order} последовательных страниц в памяти, возвращает struct page; free_pages() - освобождение; __get_free_pages(gfp_t flags ,order) - возвращает адрес

gfp_t — где и как создавать

- Динамическая аллокация

- vmalloc(unsigned long size),vzalloc() - выделение непрерывной области из нескольких страниц, достаточных для size; vmalloc_user() - выделение обнуленный страниц в пространстве пользователя
- ioremap(), iounmap() - мапинг области ввода-вывода в физической памяти в область виртуальной памяти
- page_frag_alloc(), page_frag_free() - выделение фрагментов страниц, оптимизация для сетевых драйверов
- kmalloc(size, gfp_t flags), kfree() — выделение памяти меньшей чем page в различных зонах ядра (использует SLAB-аллокатор)

- Отображение

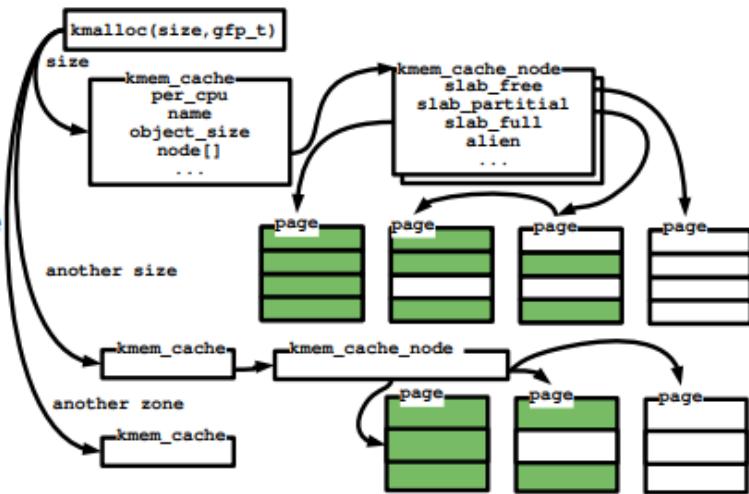
- kmap(page), kunmap(page) — отобразить страницу в highmem,
kmap_atomic(page), kunmap_atomic(kvaddr) — отобразить в FIXMAP, не блокируется!

.106. Слаб-аллокаторы SLAB/SLUB/SLOB.



SLAB/SLUB/SLOB allocators

- Несколько реализаций, в Linux осталось 2: SLOB и SLUB
- Быстрое управление небольшими объектами
- В качестве бекэнда используется alloc_page



.107. Copy on write и pagefault в Linux.



COW — Copy On Write

- Механизм использования данных в случае их записи
 - В реальной нагрузке процесс может умереть быстрее, чем изменит страницу данных
 - Давайте при fork создавать **мапинги**, а не сами страницы
 - В случае обращения на чтение — просто прочитаем
 - В случае записи — создадим копию страницы для порожденного процесса (и anonymous mapping), и потом запишем
- COW порождает множество совместно используемых для чтения страниц, и копирует их при записи



do_page_fault

- Обращение к странице, которой нет в требуемом сегменте памяти
 - Minor fault — на самом деле нужные данные в памяти есть, но по разным причинам недоступны для текущего процесса
 - Major fault — frame выгружен из памяти
 - * Очищен (выгружен) для страниц кода, статических данных, ...
 - * Находится в области подкачки для анонимных страниц, ...
 - Segmentation fault — если мы обращаемся в запрещенную область
 - Kernel panic — если мы обращаемся в запрещенную область из ядра
- Может быть вызвано дефектом кода, сбоем аппаратуры,
- Основная функция **do_page_fault**

.108. Замещение страниц в Linux. Kswapd.



Замещение страниц Linux. kswapd

- Для каждой зоны свой набор из двух списков page
- Кандидат на замещение — page у которой $PG_active=PG_referenced=0$
- kswapd запускается при нехватке памяти в заданной зоне
- Из списков исключаются страницы: SHM_LOCK, VM_LOCKED, ramfs

